

UNIVERSITY *of York*

  
UNIVERSITY OF **Hull**

# Filling the Digital Preservation Gap

*A Jisc Research Data Spring project*

*Phase Three report - October 2016*

Jenny Mitcham, Chris Awre, Julie Allinson,  
Richard Green, Simon Wilson

## Authors:

Jenny Mitcham ([jenny.mitcham@york.ac.uk](mailto:jenny.mitcham@york.ac.uk)) is Digital Archivist at the Borthwick Institute for Archives at the University of York

Chris Awre ([c.awre@hull.ac.uk](mailto:c.awre@hull.ac.uk)) is Head of Information Services, Library and Learning Innovation, at the University of Hull

Julie Allinson ([julie.allinson@york.ac.uk](mailto:julie.allinson@york.ac.uk)) is the manager of Digital York at the University of York

Richard Green ([r.green@hull.ac.uk](mailto:r.green@hull.ac.uk)) is an independent consultant working with the digital repository team at the University of Hull

Simon Wilson ([s.wilson@hull.ac.uk](mailto:s.wilson@hull.ac.uk)) is University Archivist at the University of Hull

## Acknowledgements

The authors of this report would like to thank our funders and several organisations and individuals who contributed information for it. In particular we would like to thank staff at Artefactual Systems for their advice and support throughout the project and The National Archives for their work on new research data file signatures in PRONOM and advice as we created our own. Also Fergus McGlynn and Cottage Labs for their hard work on our proof of concept implementations. Thanks also to other members of the digital preservation community that have engaged with our work on file formats; both those who have profiled the files in their digital archive using DROID (Max Eckart from Bentley Historical Library, Pawel Jaskulski a former trainee at Norfolk Record Office and Rachel MacGregor from Lancaster University) and those who have blogged about their attempts to create file format signatures for PRONOM (Andrea Byrne of Archives New Zealand and David Heelas a former trainee at Hull History Centre). Carl Wilson from the Open Preservation Foundation has also been very helpful in discussing the file formats problem with us.



This report was funded by Jisc as part of its Research Data Spring initiative.



This report is licensed under a Creative Commons CC-BY-NC-SA 2.0 UK licence.

## Contents

### [Contents](#)

### [Executive summary](#)

### [Introduction](#)

### [1. Implementation](#)

#### [Implementation at Hull](#)

##### [How it works](#)

##### [Future work](#)

#### [Implementation at York](#)

##### [How it works](#)

##### [Future work](#)

#### [Lessons Learned](#)

### [2. Research data file formats](#)

#### [Profiling research data](#)

#### [Identifying Research Data File Formats](#)

##### [Signature creation at The National Archives](#)

##### [Creating our own PRONOM signatures](#)

#### [Discussion](#)

##### [Why PRONOM?](#)

##### [What to accept?](#)

##### [Automation v. human intervention](#)

##### [Future work](#)

#### [Recommendations](#)

##### [For data curators](#)

##### [For TNA](#)

##### [For digital preservation tool providers](#)

##### [For educators](#)

##### [For funders](#)

##### [For digital preservation membership organisations](#)

##### [For researchers](#)

### [3. Outreach](#)

#### [Events](#)

[International Digital Curation Conference \(IDCC16\) - Amsterdam \(22-24 February 2016\)](#)

['Digital Preservation: Strategic Issues' - National Library of Wales \(25 February 2016\)](#)

[UK Archives Discovery Forum - Kew \(17 March 2016\)](#)

[UK Archivemata group meeting - York \(22 March 2016\)](#)

[Research Data, Records and Archives: Breaking the Boundaries - Edinburgh \(18 April 2016\)](#)

[Open Repositories \(OR16\) - Dublin \(13-16 June 2016\)](#)

[Jisc and CNI conference - Oxford \(6 July 2016\)](#)

[Hydra Virtual Connect \(7 July 2016\)](#)

[TNA Digital Transformation Day - Kew \(25 July 2016\)](#)

[Jisc Research Data Network meeting - Cambridge \(6 September 2016\)](#)

[UK Archivemata group meeting - Lancaster \(14 September 2016\)](#)

[iPRES conference - Bern \(3-6 October 2016\)](#)

[Hydra Connect - Boston \(3-6 October 2016\)](#)

[Research Data Spring Showcase - Birmingham \(20 October 2016\)](#)

[Other publications](#)

[The National Archives](#)

[Nestor](#)

[Blogs](#)

[Project website](#)

[Project reports](#)

[Glossary](#)

[Appendix 1: Draft instructions for use of Hull implementation](#)

[Depositing digital content for preservation and discovery using Box folders](#)

[Overview](#)

[Deposit options](#)

[One or more files to be represented in a single object](#)

[More than one file, each to be represented in its own object](#)

[A hierarchy of files to be represented in a single object](#)

[Appendix 2: A Draft PCDM-based Data Model for Datasets](#)

[PCDM Model](#)

[People and Organisations Model](#)

[Namespaces](#)

[Models](#)

[Dataset](#)

[CurrentPerson](#)

[CurrentOrganisation](#)

[Package](#)

## Executive summary

This report describes work carried out at the Universities of Hull and York on phase 3 of the Filling the Digital Preservation Gap project. The work described here built on the work that was carried out in phases 1 and 2 of the project. The report is in three parts.

The first section of the report describes proof of concept implementations of Archivemata for the preservation of research data at Hull and York. These implementations were integrated with other systems and services in use for managing research data at each institution. Though both implementations integrated Archivemata with a Fedora based repository, the implementations were not identical due to differing research information and storage systems as well as institutional requirements. Both institutions have successfully created largely automated systems for the longer term preservation of research data using Archivemata as a key element of the infrastructure.

The second section describes another strand of the project which looked specifically at the issue of research data file formats and the challenges involved in identifying these files automatically with current tools and registries. As well as exploring the nature of the problem (with the comparison of file format profiles and identification rates and methods) further research focused on how we might improve on this result by enhancing the available registries. A targeted piece of work in this area has increased the number of research data types within the PRONOM registry but there is a recognition that further community effort in this area is necessary. The report concludes with a set of recommendations for further work in this area.

The final section of the report details the outreach work that the project team carried out during phase 3 through presentations, publications and blog posts.

## Introduction

In order to manage research data effectively for the long term we need to consider how we incorporate digital preservation functionality into our Research Data Management (RDM) workflows. The idea behind Filling the Digital Preservation Gap was to investigate Archivemata and explore how it might be used to provide digital preservation functionality within a wider infrastructure for Research Data Management.

Phase 1 of the project investigated the need for digital preservation as part of a wider infrastructure for research data management and looked specifically at how the open source digital preservation system Archivemata could fulfil this function. Archivemata was installed and tested locally and the project team assessed how it would handle research data of various types. Areas for improvement were highlighted and a plan put in place for enhancing Archivemata to make it more suitable for incorporating into an infrastructure for research data management. The details of this work have been fully documented in a report that was produced at the end of phase 1:

*Filling the Digital Preservation Gap. A Jisc Research Data Spring project. Phase One report* - July 2015. Jenny Mitcham, Chris Awre, Julie Allinson, Richard Green, Simon Wilson<sup>1</sup>

Phase 2 of the project built on the phase 1 work, using the findings of the feasibility study to take practical steps to enhance Archivemata for use for research data preservation. We worked with Artefactual Systems on five discrete developments, and recognising that improving documentation for Archivemata lifts one of the barriers to uptake, we also funded a small piece of documentation work. The development work was designed to address some of the features of research data (trying to reduce the bottlenecks around creating SHA256 checksums for large datasets and implementing a way to highlight unidentified files) and the system integrations and workflows required for the research data infrastructure (allowing better integration with repository systems, automating the generation of a Dissemination Information Package (DIP) as part of the re-ingest process and enabling third party tools to access information for reporting purposes). Additionally in phase 2 the project team worked on their own implementation plans, mapping out how Archivemata would be implemented at each institution. These, along with full details of the development and dissemination work are available in our phase 2 report:

*Filling the Digital Preservation Gap. A Jisc Research Data Spring project. Phase Two report* - February 2016. Jenny Mitcham, Chris Awre, Julie Allinson, Richard Green, Simon Wilson<sup>2</sup>

As each phase of this project builds upon the previous phases it is suggested that readers familiarise themselves with the phase 1 and 2 reports in order to fully understand the context of the project. This report, for example, references the implementation plans in the phase 2 report which provide details of the workflows we were trying to establish.

---

<sup>1</sup> <http://dx.doi.org/10.6084/m9.figshare.1481170>

<sup>2</sup> <http://dx.doi.org/10.6084/m9.figshare.2073220>

This report describes the work that has been carried out during phase 3 of Filling the Digital Preservation Gap. Phase 3 ran for a period of six months from 14th March to the 14th September 2016. Work in phase 3 had the following aims:

- To establish proof of concept implementations of Archivematica at the Universities of Hull and York, integrated with other research data systems at each institution
- To investigate the problem of unidentified research data file formats and consider practical steps for increasing the representation of research data formats in PRONOM<sup>3</sup>
- To continue to disseminate the outcomes of the project both nationally and internationally and to a variety of different audiences

Our work in these areas will be discussed in detail in this report.

## 1. Implementation

Given the available budget and timescales at play, the purpose of our implementation work at Hull and York Universities was to establish a proof of concept rather than a production installation of Archivematica for research data. Our priority was to demonstrate that the workflows and integrations described in our phase 2 report were possible. Although Hull and York have some common infrastructure (namely Fedora and Hydra repositories), there are also differences in systems and facilities (for example CRIS and data storage) so our Archivematica implementations are not identical but instead fit with priorities and workflows at an institutional level. Whilst being able to benefit from discussions across the project team on how best to implement Archivematica we also feel it is useful to demonstrate how Archivematica can be established in two different institutional contexts.

To aid comparison between the two implementations, their key features are summarised in the table below.

<b>Infrastructure element or function</b>	<b>University of Hull implementation</b>	<b>University of York implementation</b>
<b>Metadata deposit</b>	Via Box, in a file associated with the data file(s)	Via PURE
<b>Data deposit</b>	Via Box, though a shared folder	Bespoke web form developed in Ruby/Rails
<b>Transfer via ...</b>	Box Watcher (a script that watches the shared Box folder and automatically initiates the transfer to Archivematica)	Watched directory (when something is placed in this directory it initiates a transfer to Archivematica)
<b>SIP packaging</b>	BagIt	SIP organised into a folder structure understandable by Archivematica. No

<sup>3</sup> <http://www.nationalarchives.gov.uk/PRONOM/>

		packaging (e.g., BagIt) applied.
<b>Ingest into ...</b>	Archivematica	Archivematica
<b>Archivematica workflow via ...</b>	Automation Tools (fully automated)	Automation Tools (fully automated)
<b>DIP created?</b>	Always	On request
<b>DIP sent to ...</b>	Hydra repository (via DIP processor to unpack DIP and generate Hydra objects)	Hydra application backed with Fedora 4 repository
<b>DIP discovered via ...</b>	Hydra repository	Data catalogue (to be defined)
<b>AIP sent to ....</b>	University of Hull Research Storage Service	University of York filestore

Table 1: A comparison of the proof of concept implementations at Hull and York, highlighting the similarities and differences

The following section of the report describes each implementation in turn.

## Implementation at Hull

As noted in our phase 2 report, “Hull needs functionality that is capable of providing “preservation on request” for other types of digital content in addition to research data and so the proof-of-concept implementation for phase 3 “... needs to be a pathway through the workflows which address this bigger picture.” With this in mind, the Hull implementation uses the Box collaborative cloud storage system<sup>4</sup> as its point of ingest. Box is available to all Hull staff and students through an institutional subscription, and is integrated with the campus single sign on system, so can be used to store files from many different workflows. Thus, Box takes the place of the Ingest folders from the original architecture described. We have also focused on how we use this ingest route, as it is more relevant to our use case over direct ingesting to Hydra, which we will re-visit at a later date.

---

<sup>4</sup> <https://www.box.com>



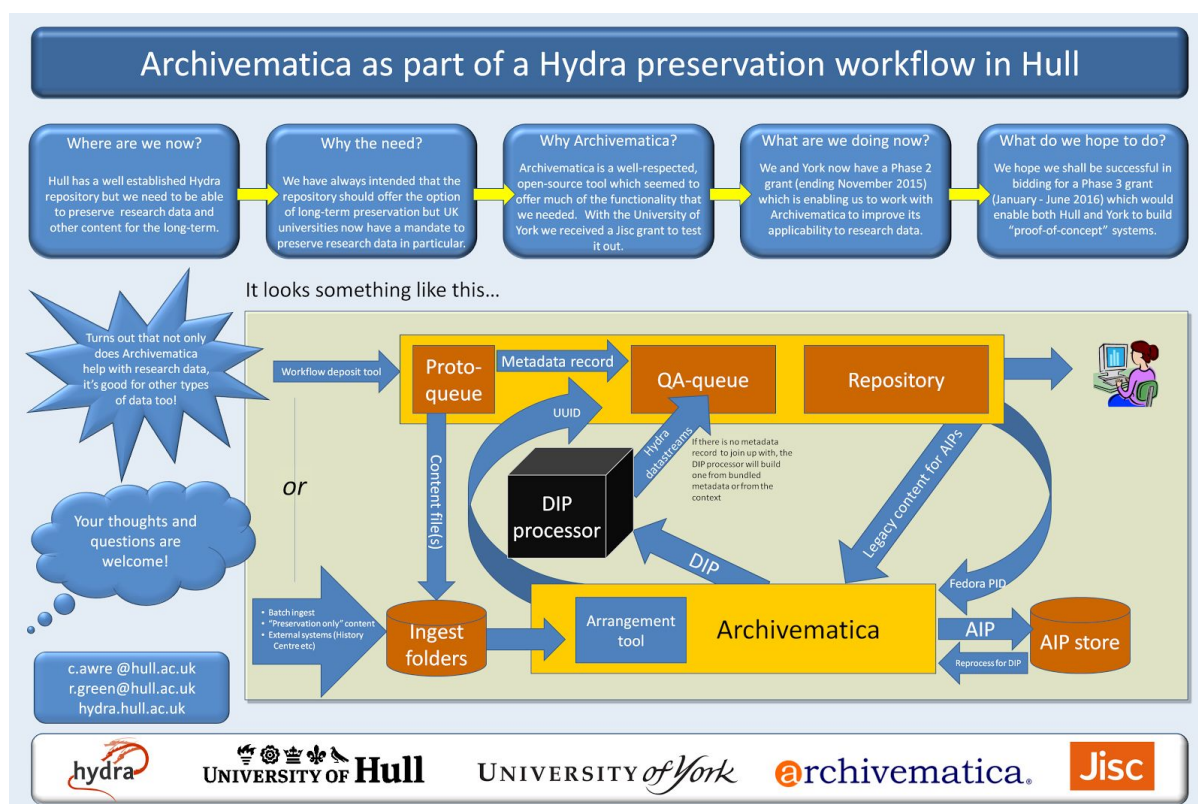


Figure 1: Hull RDM preservation architecture

## How it works

Depositors assemble the material that they wish to deposit in a Box folder within their own Box space. This material may be a single file or multiple files, it may also be a folder structure containing multiple files in an organised hierarchy; the resulting repository content may be a single object or multiple objects. These objects describe the file(s) that have been placed in the University's Research Storage Service and may optionally contain copies of the data for download. Users will be provided with instructions for preparing up to four simple files (depending on their needs) for inclusion in their folder and these will determine the way in which their data will be handled (see full instructions in Appendix 1). In practice we anticipate that many users will find they need only to generate one of these files, that which provides descriptive metadata to go with the content (the additional files deal with more complex deposit requirements). This descriptive file is a simple text file with a number of defined fields (including all the Dublin Core fields) of which only "title" and "author(s)" are mandatory.

When the user has assembled their data they use the facility in Box to share the relevant folder with a user called "Archivematica". The proof of concept system keeps a watch for new shares of this type and will process any that appear. Our "Box watcher" service will check the contents of the shared folder and, if it is valid, will create a "bag" (as in the BagIt standard<sup>5</sup>) which is passed to Archivematica. Archivematica then processes this bag of information to create an Archival Information Package (AIP), which is passed to preservation storage in the Research Storage Service, and a Dissemination Information Package (DIP) which is passed on for further processing. The final stage of Hull's proof of concept system sees the contents of this DIP transformed into one or more Hydra objects, which are then

<sup>5</sup> <https://en.wikipedia.org/wiki/BagIt>

passed into the quality assurance queue of the Hydra repository for checking prior to provision of access as appropriate.

The additional coding for this proof of concept system was carried out by Cottage Labs<sup>6</sup> and is available from Hull's github repository<sup>7</sup>.

## Future work

The next steps for the University of Hull are twofold:

- To test the proof of concept with a variety of real life examples of data, giving specific attention to the different use cases, from single files through to a directory with hierarchy. The output from these tests will be disseminated separately, and inform the definition of a production instance of the proof of concept that we can provide to researchers wishing or required to archive their data.
- To use the same system as the basis for a digital archive for the City of Culture 2017, encompassing collections of business materials, materials from artists, and materials generated by those attending events. This will use the proof of concept as a kernel for a broader system, adding additional functionality to assist with the processing of the materials for archival purposes.

The project has also made us aware of the value of getting individual tools to do what they do best and combining these as required to achieve a solution that is greater than the sum of its parts. We will also be continuing to liaise with York, as well as other members of the Hydra and Archivemata UK communities on ongoing development of the tools to facilitate their future development.

## Implementation at York

York's proof of concept implementation is very close to what was envisaged during our phase 2 project: "Intended to be an integrated solution, the RDMonitor tool will use information from the PURE Web Services about datasets described in PURE; from Archivemata about the storage of the datasets themselves and from our Fedora repository about the access copies of data". The term 'RDMonitor' has been dropped in favour of 'Research Data York' to mirror the name of our local RDM service. At the time of writing the phase 2 report we had not investigated using Automation Tools<sup>8</sup>, nor were we aware of the Puree gem<sup>9</sup> from Lancaster University Library, both of these have saved much local development effort and demonstrate our desire to use community solutions over local ones wherever practical.

---

<sup>6</sup> <http://cottagelabs.com>

<sup>7</sup> <https://github.com/organizations/uohull>

<sup>8</sup> <https://github.com/artefactual/automation-tools>

<sup>9</sup> <https://github.com/lulibrary/puree>



The second URL generated by R DYork (mentioned above) is a data access URL. This URL will be made available alongside metadata about the dataset within our data catalogue. When a user reaches this page, one of two things will happen:

- If this is the first time the data has been requested, the user will be asked to supply a valid email address and will then be alerted when the data is available.
- If a prior request has been made, the data will be automatically available for download at the URL.

Research data staff are alerted by email to the initial request for access, so they can make sure there are no restrictions on the data. They then approve the request (if appropriate) through the R DYork application and another automatic process is initiated to request the creation of a Dissemination Information Package (DIP) by Archivematica. Files from the DIP are ingested into the Fedora repository sitting behind the R DYork application.

The code for York's proof of concept system is available from the Library & Archives Technical Team github repository<sup>10</sup>. A screencast showing the working application is also available<sup>11</sup>.

### Future work

Almost in parallel with the Research Data Spring projects Jisc were planning a Research Data Shared Service<sup>12</sup>. The programme to build this service has now started and the resulting system will be managed and hosted by Jisc, and will offer three core modules : repository, preservation and reporting. The Phase 1 and 2 reports from Hull and York have been influential for scoping the preservation module and demonstrating how repository and preservation systems could be integrated. After a tendering exercise there are commercial and open source offerings for each module, including Archivematica (for preservation) and Hydra (for the repository). Over 20 pilot institutions have been recruited (including York) and all have identified preservation as a priority.

York plan to continue the implementation work we have started but ensure it continues to align with the Jisc Shared Service. With longer timescales at play for the establishment of the Shared Service we are also keen to move the prototype we have developed into production in the meantime so that our research data team have a system through which they can manage deposits to R DYork.

There are several areas we would like to address in order to refine and polish our implementation work and move it into production:

- Complete work on the fully automated DIP re-ingest process.
- Carry out further testing and implement improvements as required.
- Finalise the data model (see appendix 2) and align this with broader Hydra activity at York.
- Deploy code to a production architecture.

---

<sup>10</sup> <https://github.com/digital-york/researchdatayork>

<sup>11</sup> [https://www.youtube.com/watch?v=3cl5W\\_7gYvM](https://www.youtube.com/watch?v=3cl5W_7gYvM)

<sup>12</sup> <https://www.jisc.ac.uk/rd/projects/research-data-shared-service>

- Investigate the University of Lancaster's 'preservation' gem<sup>13</sup> as a standard means of interacting with Artefactual's Automation Tools.
- Consider extracting some functionality from the RDYork application into separate gems to better able reuse of the code-base. The code to extract the file structure from the Archivematica METS file has already elicited external interest, for example
- Investigate Hull's 'Box Watcher' approach with Google Drive.
- Extend functionality to support data deposits that aren't made via PURE, in particular those from PhD students (who don't have PURE accounts at York).
- Consider how to marry the automated workflows defined here with the need to solve the file format identification problem (described in a later section of this report).

In addition, we will be testing out Lancaster's DMAOnline<sup>14</sup> reporting tool with data from PURE and Archivematica.

## Lessons Learned

At the start of phase 1 we set out to answer the question: 'is Archivematica suitable for preserving research data?', and at the end of phase 3 our conclusion remains that it is. Developments in the software during the course of phases 2 and 3 (both those that the project funded and others sponsored by other institutions), make Archivematica better able to deal with larger datasets, easier to automate, and introduce options for querying Archivematica to assess RDM readiness. The proof of concept implementations at both Hull and York successfully demonstrate fully automated ingest workflows and both make use of Artefactual Systems' Automation Tools in order to do this.

Some lessons learnt from the project include:

- It can be hard to balance tight project deadlines (ours) with product release cycles (Archivematica) where an externally developed system is being used.
- Translating the phase 2 implementation plans into development plans benefitted from conscious and detailed discussion on how ideas translated into required code.
- Establishing the server infrastructure on which we needed to implement the proof of concept required negotiation and regular communication with local IT staff, particularly around firewall and security settings. Making the clear distinction between a test proof of concept and a production system at an early stage will help define the level of work required.
- When you don't have much time, work with good developers. If this resource is not available (with free development cycles to match the project timeframe) in house, have external resource available to tap into.
- Archivematica is reasonably easy to install in a standard configuration but there are documentation gotchas, particularly if you want to diverge or delve further into the flexibility of the system. Improving developer documentation remains a challenge for Archivematica.
- Developing out in the open, seeking peer review and discussion work is definitely the right approach.

---

<sup>13</sup> <https://github.com/lulibrary/preservation>

<sup>14</sup> <http://www.dmao.info/>

## 2. Research data file formats

In our phase 1 report we discussed the nature of research data and highlighted the many and varied software applications and file formats that are in use for research. In the report it was noted that “It is clear that as we start to ingest research data into a digital archive we will encounter many files that will not be automatically identified”<sup>15</sup> but this hypothesis was not formally tested.

Our work during phase 2 partially addressed this problem by funding a development within Archivemata which would enable users of the system to more easily locate those files with a format that can not be identified automatically. Highlighting these files would enable the digital curator to carry out follow up actions. Further development work to enhance and improve this functionality was also proposed and discussed.

In phase 3 we continued this theme alongside our implementation work. Preserving digital data isn't solely reliant on the implementation of a digital preservation system, it is also necessary to think about related challenges that will be encountered and how they may be addressed.

### Profiling research data

Since May 2015 the University of York has been accepting datasets from researchers into its institutional repository. Many of these deposits came as a result of the EPSRC Expectations<sup>16</sup> thus the majority of the data was scientific in nature, much of it originating from the Department of Chemistry. We used this collection of research data to test the hypothesis from our phase 1 report regarding the variety of formats and difficulties that would be encountered when trying to automatically identify formats.

We ran DROID<sup>17</sup> (a file identification tool from The National Archives) over the files in the collection to establish which file formats within the collection could be automatically identified and the results of this exercise were published as a blog post<sup>18</sup>.

Our initial sample size wasn't large, and actual analysis was carried out on a sample smaller than anticipated. When the analysis was underway it became clear that DROID does not look inside all of the zip files that we hold<sup>19</sup>. Datasets submitted to Research Data York typically arrive zipped up and we do not extract the files as a matter of course. The sample included 8 .rar files which contained another 1291 digital objects which were not included in

---

<sup>15</sup> A fuller discussion on the nature of research data and the common research data software applications in use at the University of York can be found in our phase 1 report

<sup>16</sup> <https://www.epsrc.ac.uk/about/standards/researchdata/expectations/>

<sup>17</sup>

<http://www.nationalarchives.gov.uk/information-management/manage-information/preserving-digital-reCORDS/droid/>

<sup>18</sup> <http://digital-archiving.blogspot.co.uk/2016/05/research-data-what-does-it-really-look.html>

<sup>19</sup> A response received from The National Archives on 19th May 2016 states that “DROID can only probe within the following archival container formats: zip, gzip, tar, plus the web archival formats .arc and .warc. It is on our development roadmap for DROID to add the ability to probe further archival containers, including .rar, 7z, .bz, .iso, which would provide the requested functionality, but we don't currently have a target release date for this”

the statistics (though a quick check showed that the percentage of identified files in this additional sample was strikingly similar to that reported).

The headlines were quite startling. Of the 3752 files analysed, only 1382 (37%) were assigned a file format identification by DROID. This success rate appeared to be very low but we were keen to compare this figure with similar profiles from research data and other born digital collections, so, within the blog post, a question was posed:

*Is identification of 37% of files a particularly bad result or is it similar to what others have experienced?*

As part of this project, a comparable test was carried out on the born digital holdings of Hull University Archives and the results were very different with an impressive 98% of files identified<sup>20</sup>. Bentley Historical Library, Norfolk Record Office also engaged in this exercise and shared their results<sup>21</sup>. This was a great example of how the digital preservation community can work together and learn from each other.

These profiling exercises certainly seemed to demonstrate that 37% of files identified at York was a particularly low result, but considering the small sample size for the York research data it was decided that it would be beneficial to look at DROID profiles of other research datasets. The University of Hull carried out a DROID profiling exercise on the research data held on their research storage area network<sup>22</sup>, which currently holds work in progress as well as final outputs, and the results were quite different to those for York research data with 89% of files assigned an identification by DROID. The sample size for this study was very large (leading to issues when trying to analyse the results) and a look at the top ten identified formats showed that 70% of the files on this storage area were TIFF images (from research groups carrying out medical imagery for the most part). This is perhaps not typical of research datasets and this will certainly have contributed to the high format identification rate.

Lancaster University also carried out a similar study, looking specifically at research datasets that had been deposited in their institutional repository<sup>23</sup>. A file format identification rate of 46% was not too far off initial findings from York, though as can be seen in the table below, the method of identification differed.

---

<sup>20</sup> The results of this exercise were published as a blog post:

<http://digital-archiving.blogspot.co.uk/2016/08/research-data-is-different.html>

<sup>21</sup> Results from Bentley Historical Library:

<http://archival-integration.blogspot.co.uk/2016/06/born-digital-data-what-does-it-really.html>; results from Norfolk Record Office:

<http://digital-archiving.blogspot.co.uk/2016/09/file-format-identification-at-norfolk.html>

<sup>22</sup> The results of this exercise are not published elsewhere

<sup>23</sup> These findings from Lancaster University were shared directly with project team and are not yet published

Institution and test data	No of files in sample	% of files identified	% identified by signature	% identified by container	% identified by extension
<b>RESEARCH DATA</b>					
University of York Research Data	3,752	37%	48%	5%	47%
University of Hull Research Data	10,174,380	89%	data not collected		
University of Lancaster Research Data	24,069	46%	90%	1%	9%
<b>OTHER DIGITAL ARCHIVES</b>					
Hull University Archives Born Digital Holdings	270,867	98%	data not collected		
Bentley Historical Library Born Digital Holdings	731,949	90%	88%	10%	2%
Norwich Record Office Born Digital Holdings	49,117	96%	83%	15%	2%

Table 2: A summary of DROID analysis of results for several institutions. Note that for ease of comparison, figures quoted are rounded up or down to nearest whole number.

As well as the low file format identification rate for York research data, another finding reported in the blog post related to the method of identification for those research data files that were identified. Only 53% of the identified files in York’s research dataset were identified by signature or container (methods which suggest more reliable and accurate identification) and this contrasts to 98% at both Bentley Historical Library and Norwich Record Office and 91% for research data at Lancaster University. The remaining identifications were carried out by file extension. File extension is not the most reliable identification method, given that unrelated files coming from different software applications can share the same extension (for example .dat files as discussed below). In the York results 47% of identified files were identified in this way which does suggest that further human intervention may be required to validate those identifications.

Another interesting finding from the work to identify data deposited with Research Data York was the large number files with no file extension. Of the files that were not automatically identified, files with no extension made up 26% of the total. Given the reliance on identification by extension within this sample, this does rather reduce our chances of identifying these files. The second largest group of unidentified files (12%) were files with a .dat extension. This is a fairly common file extension and represents a whole range of files



produced by a variety of different software applications, hardware or operating systems<sup>24</sup> so this isn't a problem that could be solved by the creation of a single file format signature.

Of course there may be other things at play here, for example it would be interesting to explore whether the date of the sample data is a factor in the results. The York research data sample was all deposited within the last 18 months and been modified relatively recently. As stated in the blog post "The data is mostly fairly recent, as suggested by the last modified dates on these files, which range from 2006 to 2016 with the vast majority having been modified in the last five years"<sup>25</sup>. As illustrated in the graph below there is only a small quantity of data from 2006 and the next oldest dates from 2009. Not all of the other studies quoted included details of file dates, but the Bentley example demonstrates that the data was collected over a longer period of time and the range of dates present is wider: "The vast majority of the data was last modified in the past 15 years, and our peaks are in in 2006 and 2008."<sup>26</sup>

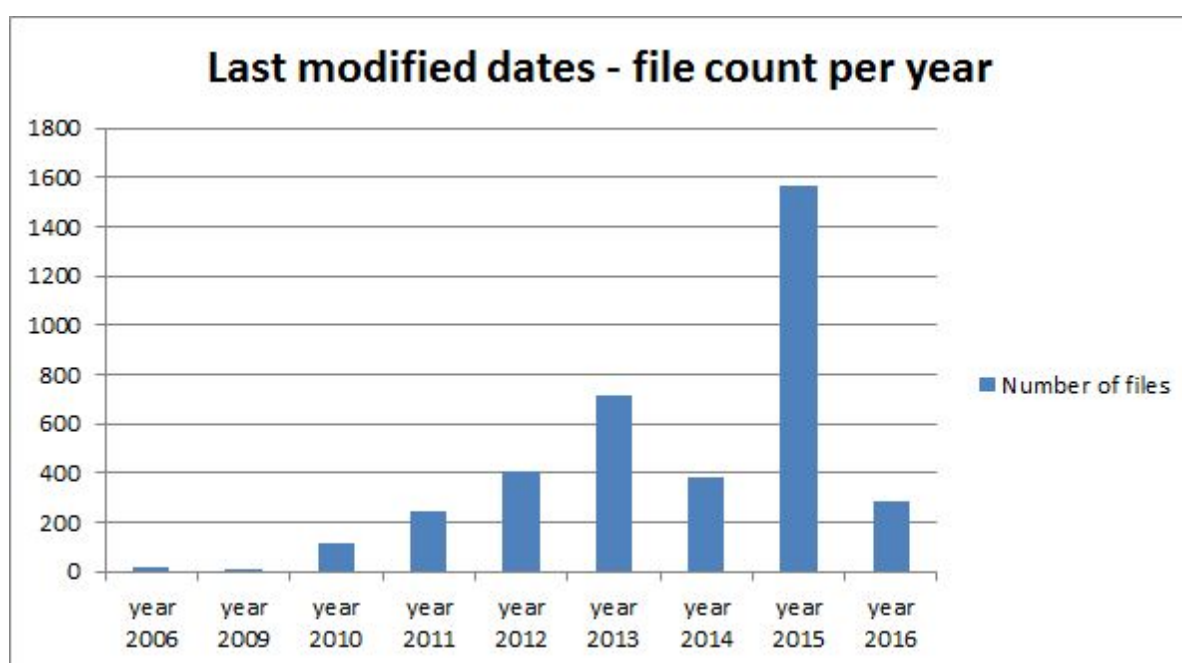


Figure 3: Last modified dates in York's research data sample

We could speculate that given the pace of technological change, it is possible that data that has been worked on more recently may be more likely to have been saved in newer file formats and those newer formats may not yet have made their way into the PRONOM registry.

Our initial blog post was a fairly quick and informal snapshot of research data at York but as others have contributed to the picture it becomes clear that the landscape is a varied one. It would be useful to carry out a more formal piece of work in this area with a larger sample of contributing institutions and guidelines in place regarding the PRONOM signature file version and internal settings in use within DROID. This would allow for easier comparison of results across the institutions and more solid conclusions to be reached.

<sup>24</sup> See for example the list at Just Solve the File Format Problem (by no means exhaustive): [http://fileformats.archiveteam.org/wiki/Category:File\\_formats\\_with\\_extension\\_.dat](http://fileformats.archiveteam.org/wiki/Category:File_formats_with_extension_.dat)

<sup>25</sup> <http://digital-archiving.blogspot.co.uk/2016/05/research-data-what-does-it-really-look.html>

<sup>26</sup> <http://archival-integration.blogspot.co.uk/2016/06/born-digital-data-what-does-it-really.html>

However, it is clear from work carried out in this area so far that there is scope for increasing the number of research data formats within the PRONOM database. This issue is addressed in the next section of the report.

## Identifying Research Data File Formats

As mentioned in our phase 1 report “As we move towards a proof of concept for archiving research data we should continue to engage with the team who maintain PRONOM and promote discussions within the wider digital preservation community about a sustainable or more automated way to address this problem.”

One of the goals of this phase of the project was therefore to try and increase the representation of research data file formats within the PRONOM database whilst encouraging wider engagement in the digital preservation and research data community.

### Signature creation at The National Archives

The majority of file format signature development work for PRONOM is undertaken at The National Archives (TNA) in the UK. They are understandably driven by their own priorities - so that they can manage the formats that they hold - but they also carry out work for the community on request. Anyone may submit information and sample files to the PRONOM team so that signature development work can be carried out but the timescales for this unfunded community work will vary greatly (and will also depend on the quality and accuracy of the information submitted)<sup>27</sup>.

Due to the short timescales at play for this project, we directly funded some signature development work to ensure that it would be completed within the project time frame.

At York, the research data profiling work described above and the top 20 research data applications at York as mentioned in our phase 1 report helped inform our selection of file formats for submission to PRONOM. We chose Gaussian input files<sup>28</sup> and JEOL NMR Spectroscopy files<sup>29</sup>. A longer discussion about these formats and the submission to PRONOM is available as a blog post<sup>30</sup>. The two new signatures were released on 29th June 2016 as part of DROID signature file version 85<sup>31</sup>.

Priorities at Hull were informed through a survey of research data file formats generated, and subsequent discussions with researchers self-identifying as having more unusual formats. This led to the creation of six new signatures. Firstly the addition of an ESRI ArcMap Document<sup>32</sup>. We had noted in our phase 1 report that ESRI files are well represented in PRONOM but that the .mxd file was not included so this work was designed to fill that gap.

---

<sup>27</sup> A blog post from Paul Young provides a good summary of the community engagement around PRONOM: <http://blog.nationalarchives.gov.uk/blog/identifying-digital-file-formats-collaborative-effort/>

<sup>28</sup> <http://www.nationalarchives.gov.uk/PRONOM/fmt/894>

<sup>29</sup> <http://www.nationalarchives.gov.uk/PRONOM/fmt/895>

<sup>30</sup> <http://digital-archiving.blogspot.co.uk/2016/07/new-research-data-file-formats-now.html>

<sup>31</sup> <http://www.nationalarchives.gov.uk/aboutapps/pronom/release-notes.xml>

<sup>32</sup> <http://www.nationalarchives.gov.uk/PRONOM/fmt/916>

Five further signatures were also developed for AmiraMesh file formats<sup>33</sup>. These new signatures were released on 27th July 2016 in version 86.

### Creating our own PRONOM signatures

It was interesting working directly with TNA on signature development work and discussions with the PRONOM team led to a deeper understanding of how new entries in PRONOM are created. Following on from this, the project team decided to try to create their own file format signatures for PRONOM. The file format identification problem is a large and complex one and it was recognised that it could be solved more quickly with wider and more active community engagement. As noted in Paul Young's blog post<sup>34</sup>, the PRONOM team commit to releasing 100 new PRONOM records per year. One way we could increase this figure and prioritise research data formats would be through more direct engagement with this problem. Clearly we are not the first to engage with this task, but we wanted to establish whether it was practical for 'an average digital archivist' to attempt file signature creation.

York's experiences at signature development are discussed in full in a blog post<sup>35</sup>. The resulting signature for an OMNIC Spectral Data File<sup>36</sup> was made available in DROID signature file version 88 on the 27th September 2016.

Hull were also able to create their own PRONOM signature and this work by Transforming Archives trainee Dave Heelas has again been documented in a blog post<sup>37</sup>. This work led to the creation of a signature for Final Draft Document 5-7<sup>38</sup> again in signature version 88.

Inspired by the work of this project, Andrea Byrne from Archives New Zealand also took on the challenge and again blogged about her progress<sup>39</sup>. Her detailed post describes some of the detective work necessary to understand the files she chose to investigate as well as some of the benefits of working in an open and community-focused way. The resulting file signatures are for AppleSingle 1 and 2<sup>40</sup> and are now available in signature version 88 along with a revised entry for AppleDouble Resource Fork 1<sup>41</sup>.

The key points to note from this recent community signature creation work are as follows:

- It is possible for data curators and digital archivists to actively contribute to signature development (and it is not a task only suited to those who are more technically minded).

---

<sup>33</sup> <http://www.nationalarchives.gov.uk/PRONOM/fmt/917>,  
<http://www.nationalarchives.gov.uk/PRONOM/fmt/918>,  
<http://www.nationalarchives.gov.uk/PRONOM/fmt/919>,  
<http://www.nationalarchives.gov.uk/PRONOM/fmt/920>,  
<http://www.nationalarchives.gov.uk/PRONOM/fmt/921>

<sup>34</sup> <http://blog.nationalarchives.gov.uk/blog/identifying-digital-file-formats-collaborative-effort/>

<sup>35</sup> <http://digital-archiving.blogspot.co.uk/2016/08/my-first-file-format-signature.html>

<sup>36</sup> <http://www.nationalarchives.gov.uk/PRONOM/fmt/963>

<sup>37</sup> <http://hullhistorycentre.blogspot.co.uk/2016/08/from-plans-to-digital-content-daves.html>

<sup>38</sup> <http://www.nationalarchives.gov.uk/PRONOM/fmt/964>

<sup>39</sup>

<http://openpreservation.org/blog/2016/09/08/making-the-switch-from-user-to-user-and-contributor-my-first-file-format-signature/>

<sup>40</sup> <http://www.nationalarchives.gov.uk/PRONOM/fmt/967>,

<http://www.nationalarchives.gov.uk/PRONOM/fmt/968>

<sup>41</sup> <http://www.nationalarchives.gov.uk/PRONOM/fmt/966>

- Start with the low-hanging fruit - binary files are much easier to develop signatures for than ASCII. Starting with files represented in the institution's own collections is also a more meaningful starting point and it is often possible to locate additional sample files to test the identification work.
- Signature development is not always straightforward - input from the PRONOM team at TNA is required to check, refine and test the signatures before incorporating them into a PRONOM signature release<sup>42</sup>.
- There are benefits to be had from sharing experiences around signature development with the wider digital preservation community (for example via the PRONOM google group<sup>43</sup>).
- Enhanced and improved documentation<sup>44</sup> around signature development would help the community engage more as would moving the signature development utility<sup>45</sup> from prototype to production.

As mentioned above, this flurry of work specifically targeted binary formats, however it should be noted that many of the research data formats in the York sample were ASCII files. File format signatures can be developed for ASCII files (if the files conform to a set structure which can be defined) but these are typically more complex to create than signatures for binary files, using regular expressions to describe their identifying characteristics<sup>46</sup>. It is also helpful to think beyond file identification of single digital objects and consider recognising groups of files that might together make up a more complex digital object, for example a website or software application.

## Discussion

### Why PRONOM?

The file format problem is not a new one for the digital preservation community. A recent attempt to address this issue as a community, Just Solve the File Format Problem<sup>47</sup>, has produced a useful resource but does not currently include many of the research data file formats highlighted in our profiling work. Furthermore, the real need as highlighted by our project is for a resource or registry which provides a means of automatically identifying file formats and allocating a unique identifier. Just Solve the File Format Problem was a project to bring together disparate sources of information about file formats, not specifically to identify them.

PRONOM, first released in 2002, is designed both to store information to aid automatic identification of files and to provide a means of uniquely identifying the format (via the Persistent Unique Identifier or PUID) and has demonstrated a longevity that other format

---

<sup>42</sup> In a blog post entitled 'A week of file format research' David Clipsham describes the work he undertook in the course of a week to help refine submitted signatures before incorporating them into PRONOM: <http://openpreservation.org/blog/2016/08/31/a-week-of-file-format-research/>

<sup>43</sup> <https://groups.google.com/forum/#!forum/pronom>

<sup>44</sup> Helpful documentation does exist but it does not describe the full process of signature development (including use of the signature development utility and testing using DROID):

<http://www.nationalarchives.gov.uk/documents/information-management/pronom-file-signature-research.pdf>

<sup>45</sup> <http://www.nationalarchives.gov.uk/pronom/sigdev/index.htm>

<sup>46</sup> Andy Jackson discusses further identification methods for text-based formats in a blog post:

<http://anjackson.net/2016/06/08/frontiers-in-format-identification>

<sup>47</sup> <http://fileformats.archiveteam.org/>

registries have not been able to match<sup>48</sup>. PRONOM is designed to work alongside DROID, a file identification tool also developed by TNA. However, there are many other tools that can be used to identify files - see for example those listed in the tool registry COPTR<sup>49</sup>.

Archivematica gives the operator the choice of the PRONOM based tools FIDO<sup>50</sup> or Siegfried<sup>51</sup> to carry out format identifications, or there is the option of identification by file extension<sup>52</sup>.

One of the benefits of using a PRONOM based tool for identification is that a PUID will be assigned to the file to match the identifier of the format in PRONOM. In theory this should allow reporting and search and retrieval based on file format<sup>53</sup> and would help facilitate the sharing of information about file formats with others, for example with other digital repositories about strategies for the preservation of a particular collection of files. Within Archivematica the PUID is stored within the PREMIS metadata for the digital object and is used by the Format Policy Registry (FPR)<sup>54</sup> which defines which preservation or dissemination actions should occur to a particular type of file. If a file is not identified, using the tools provided within Archivematica the system will not store an identifier for the format and further automated preservation actions through the FPR will be limited and non-specific.

### What to accept?

One possible solution to the file format problem as described would be to limit the types of files that would be accepted within the digital repository. This is a tried and tested approach for certain disciplines and data archives<sup>55</sup> and is one that is also consistently supported by the digital preservation literature. For example The National Digital Stewardship Alliance's Levels of Digital Preservation<sup>56</sup> recommends at level one that an organisation carrying out digital preservation activities should "... encourage use of a limited set of known open formats ...". This is problematic for institutions facing the task of preserving research data. Researchers within an institution will use such a wide range of specialist hardware and software and it will be hard for the repository and research support staff to provide appropriate advice on suitable formats. For much of the data there will be no obvious preservation format for that data.

We also wish to encourage the innovative and cutting-edge research that is going on within our institutions so limiting the range of formats deposited would be counter-productive. At the University of York we encourage researchers (through both our RDM training sessions and RDM webpages<sup>57</sup>) to consider file formats throughout their project and think about the longevity and accessibility of the formats they select, but this is only guidance and ultimately we leave it up to the researcher to decide what formats to deposit their data in. We accept

---

<sup>48</sup> See for example the Unified Digital Format Registry which announced it would no longer continue in April 2016: <http://udfr.org/>

<sup>49</sup> [http://coptr.digipres.org/Category:File\\_Format\\_Identification](http://coptr.digipres.org/Category:File_Format_Identification)

<sup>50</sup> [http://coptr.digipres.org/FIDO\\_\(Format\\_Identification\\_for\\_Digital\\_Objects\)](http://coptr.digipres.org/FIDO_(Format_Identification_for_Digital_Objects))

<sup>51</sup> <http://coptr.digipres.org/Siegfried>

<sup>52</sup> Identification by file extension is carried out with a basic python script in Archivematica.

<sup>53</sup> This functionality is not yet available in Archivematica

<sup>54</sup>

<https://www.archivematica.org/en/docs/archivematica-1.5/user-manual/preservation/preservation-planning/#fpr>

<sup>55</sup> See for example depositor guidelines for the UK Data Archive:

<https://www.ukdataservice.ac.uk/manage-data/format/recommended-formats> and Archaeology Data Service: <http://archaeologydataservice.ac.uk/advice/FileFormatTable>

<sup>56</sup> <http://ndsa.org/activities/levels-of-digital-preservation/>

<sup>57</sup> <http://www.york.ac.uk/library/info-for/researchers/data/organising/#tab-1>

these formats and will preserve them on a best efforts basis. Understanding the file format moves us one step closer to preservation and reuse over the longer term.

### Automation v. human intervention

It is perhaps worth highlighting the slight paradox in our project approach. Our proof of concept implementation work set out to establish a pragmatic and parsimonious<sup>58</sup> approach to digital archiving involving freely available tools and a high degree of automation in order to limit staff intervention in the process. However, if the file format identification problem is to be solved, this will of course require a substantial amount of staff time to highlight areas where work is required, prioritise and research the formats and submit them for signature development. This problem could be addressed at a higher level with targeted project funding to improve identification rates for research data, but given the rate of technological change it will always be an ongoing issue that institutions will need to engage with periodically.

### Future work

Our project has highlighted the challenges around file format identification for research data and this information will feed into the ongoing Jisc Research Data Shared Service project<sup>59</sup> through which more resource is available to address the problem. Through the Shared Service project a piece of work will be carried out by the Open Preservation Foundation<sup>60</sup> to scope this issue and assess solutions to facilitate the preservation of research data. This piece of work will include an investigation of identification methods that are not based on PRONOM (for example Apache Tika<sup>61</sup> and the Unix File utility<sup>62</sup>) and consider how information about formats that are identified by these other methods can be submitted to PRONOM in order that we can uniquely identify them and share information across preservation systems. Some prior work to compare and aggregate the contents of different format registries already exists<sup>63</sup> and it is anticipated that there will be benefits to revisiting this problem.

## Recommendations

As discussed above, we have further investigated the extent of the file format problem for research data and carried out some work to increase the representation of research data file formats within PRONOM.

We recognise that the research data formats that this project has added to PRONOM are not a solution in themselves but should be seen as just the start of a bigger piece of community work in this area. We are keen to see this work continued by the digital preservation and

---

<sup>58</sup> The term Parsimonious Preservation was coined by Tim Gollins and refers to the simple and affordable steps you can take to start preserving digital material using freely available tools:

<http://www.nationalarchives.gov.uk/documents/information-management/parsimonious-preservation.pdf>

<sup>59</sup> <https://www.jisc.ac.uk/rd/projects/research-data-shared-service>

<sup>60</sup> <http://openpreservation.org/>

<sup>61</sup> [http://coptr.digipres.org/Apache\\_Tika](http://coptr.digipres.org/Apache_Tika)

<sup>62</sup> [http://coptr.digipres.org/Fine\\_Free\\_File\\_Command](http://coptr.digipres.org/Fine_Free_File_Command)

<sup>63</sup> See for example Andy Jackson's format registry aggregator: <http://www.digipres.org/formats/> and recent work by Tim Allison:

<http://openpreservation.org/blog/2016/10/04/apache-tikas-regression-corpus-tika-1302/>

research data community in order to enable better research data identification nationally and internationally. The following recommendations aimed at a number of stakeholders suggest a way forward to help us achieve this as a community:

#### For data curators

- Actively engage with issues around file format identification. This could include:
  - Sharing file format profiles
  - Noting and sharing information about files that are wrongly identified by the available tools<sup>64</sup>
  - Noting and sharing information about files that are not identified by the available tools
  - Engaging in signature development work
  - Making sample files available to test corpuses (if appropriate permissions can be sought)
- Greater engagement with researchers on the value and necessity of recognising and recording the file formats they will use/generate to inform effective data curation.

#### For TNA

- Provide enhanced documentation describing how to engage in the process of signature creation.
- Greater active engagement with the community about what is currently being worked on and what the priorities are (transparency of process).
- Greater active engagement with the community to seek example file formats for testing and development purposes.
- Facilitate better integration between PRONOM and third party digital preservation tools and systems - for example automatic notifications which allow the tools to pick up new signatures as they are released<sup>65</sup>.
- Facilitation of crowdsourcing efforts to suit need.

#### For digital preservation tool providers

- Digital preservation tools should incorporate better methods for engaging with the file format identification problem. For example:
  - Highlighting unidentified files as an integrated element of the ingest process and to increase awareness for data curators of the implications of this for future workflow and preservation actions
  - Re-running file identification tools
  - Sharing file identification profiles and reports
  - Efficient methods of incorporating up-to-date PRONOM sigs

#### For educators

- Those working in digital preservation should have a basic understanding of how file identification works and how they can contribute to the available registries and tools. The mechanics of file identification and how to contribute to the tools and registries should be taught in digital preservation training courses.

---

<sup>64</sup> A good example of this level of sharing was the winning poster at iPRES 2016 "To Act or not To Act - Handling File Format Identification Issues in Practice" from ETH Zurich:  
<https://twitter.com/ipres2016/status/783217695835783168>

<sup>65</sup> This notification currently works for DROID but not other PRONOM-based tools

### For funders

- Time and money are barriers to engagement with this problem - particularly given the fact that it is too big for one institution to solve alone. There is often an interest and willingness to engage, but people struggle to make the case for this community work against local priorities. Funding for this work would reduce the barriers to engagement.

### For digital preservation membership organisations

- Encourage community effort in this area - for example by organising training days, workshops and hackathons<sup>66</sup> around signature development.
- Facilitate the sharing of DROID profiles to establish what the priorities should be for signature development work.

### For researchers

- Supply adequate metadata about submitted datasets. Clear and accurate metadata about file formats and hardware/software dependencies will aid file format identification and future preservation work.
- Be open to sharing sample files for test data corpuses and to aid signature development where appropriate.

---

<sup>66</sup> See for example the recent and successful JHOVE hack day organised by the OPF:  
<http://openpreservation.org/blog/2016/10/19/jhove-online-hack-day-report/>



### 3. Outreach

During the six months of phase 3 of this project, substantial efforts were made to ensure that we kept people informed about what we were doing, as well as promoting the existence of our phase 1 and 2 project reports. We attended a wide range of different events, engaging with audiences from the UK and beyond. As in previous phases of this project, we have been encouraged by the level of interest the project has generated and have received positive feedback on the work we are doing.

Outreach channels consisted largely of presentations and posters at organised events and blog posts published on the University of York's Digital Archiving blog<sup>67</sup>; our project was also featured in other publications. Outreach work continued beyond the active project period to include dissemination of our phase 3 work.

#### Events

##### International Digital Curation Conference (IDCC16) - Amsterdam (22-24 February 2016)

*Jisc Research Data Management Shared Service Workshop: An institutional perspective* - Jenny Mitcham (presented in the Jisc Research Data Management Shared Service Pilot workshop)

*"Filling the digital preservation gap" for Research Data* - Jenny Mitcham, Julie Allinson, Chris Awre, Richard Green, Simon Wilson<sup>68</sup>

##### 'Digital Preservation: Strategic Issues' - National Library of Wales (25 February 2016)

*A collaborative approach to "filling the digital preservation gap" for Research Data Management* - Julie Allinson

##### UK Archives Discovery Forum - Kew (17 March 2016)

Poster: *Filling the Digital Preservation Gap* - Jenny Mitcham and Simon Wilson

---

<sup>67</sup> <http://digital-archiving.blogspot.co.uk/>

<sup>68</sup>

<http://www.dcc.ac.uk/sites/default/files/documents/IDCC16/Parallel%20B/Session%203/Jen%20Mitcham.pdf>

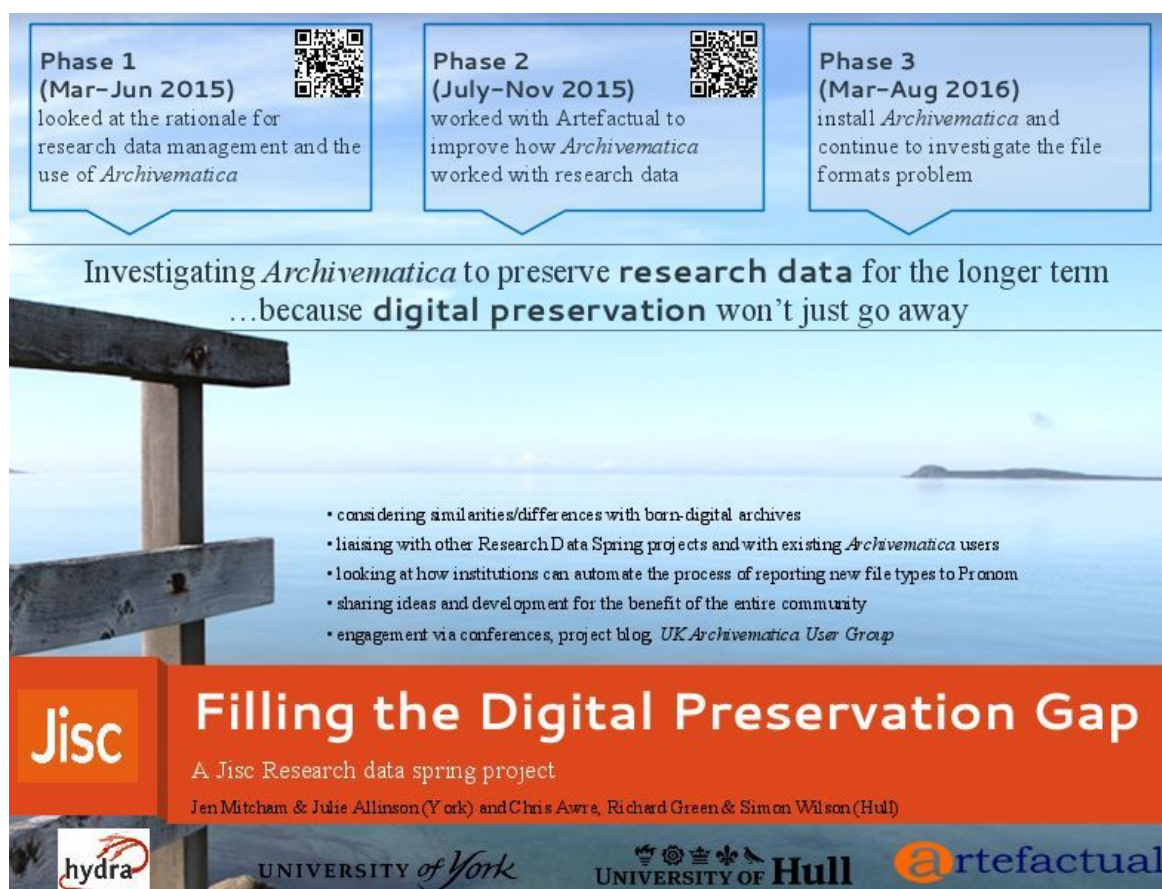


Figure 4: UK Archives Discovery poster

### UK Archivematica group meeting - York (22 March 2016)

*"Filling the digital preservation gap" for research data: Results from phase 2 and plans for phase 3* - Jenny Mitcham<sup>69</sup>

### Research Data, Records and Archives: Breaking the Boundaries - Edinburgh (18 April 2016)

*A collaborative approach to "filling the digital preservation gap" for Research Data Management* - Chris Awre<sup>70</sup>

### Open Repositories (OR16) - Dublin (13-16 June 2016)

*Prototyping a digital preservation pipeline with Archivematica, Fedora 4 and Hydra* - Julie Allinson and Justin Simpson (Artefactual Systems)

### Jisc and CNI conference - Oxford (6 July 2016)

*Addressing the preservation gap at the University of York* - Jenny Mitcham

<sup>69</sup>

<https://docs.google.com/presentation/d/1NbzpxB27yNWRQP2oUr-8qeApFq5lSUc2jiAMqD8sWOO/edit?usp=sharing>

<sup>70</sup>

[https://media.ed.ac.uk/media/Chris+Awre+%28Head+of+Information+Services%2C+University+of+Hull%29/1\\_ip395t3t/41567111](https://media.ed.ac.uk/media/Chris+Awre+%28Head+of+Information+Services%2C+University+of+Hull%29/1_ip395t3t/41567111)

**Hydra Virtual Connect (7 July 2016)**

*Hydra For Research Data* - Julie Allinson and Matthew Phillips (University of Durham)

**TNA Digital Transformation Day - Kew (25 July 2016)**

*Going digital: a case study from the Borthwick Institute for Archives* - Jenny Mitcham

**Jisc Research Data Network meeting - Cambridge (6 September 2016)**

*Implementing Archivemata for Research Data Preservation at York and Hull* - Jenny Mitcham

**UK Archivemata group meeting - Lancaster (14 September 2016)**

*Filling the Digital Preservation Gap: update on phase 3 work* - Julie Allinson and Jenny Mitcham<sup>71</sup>

**iPRES conference - Bern (3-6 October 2016)**

*Preserving Research Data: Linking Repositories and Archivemata* - Jenny Mitcham, Matthew Addis (Arkivum), Julie Allinson, Chris Awre, Richard Green, Simon Wilson<sup>72</sup>

**Hydra Connect - Boston (3-6 October 2016)**

*Hydra, research data and Archivemata* - Julie Allinson, Richard Green

Poster: *Filling the Digital Preservation Gap: integrating Archivemata and Hydra for Research Data Management*

**Research Data Spring Showcase - Birmingham (20 October 2016)**

Lightning talk and demo: *Filling the Digital Preservation Gap* - Julie Allinson, Chris Awre, Jenny Mitcham

---

<sup>71</sup>

<https://docs.google.com/presentation/d/1b-yF36iD3llqB1WxVbOyiwgSqmvUWk0TE3k97ejFuuM/edit?usp=sharing>

<sup>72</sup>

[http://www.ipres2016.ch/frontend/organizers/media/iPRES2016/\\_PDF/IPR16.Proceedings\\_4\\_Web\\_Broschuere\\_Link.pdf](http://www.ipres2016.ch/frontend/organizers/media/iPRES2016/_PDF/IPR16.Proceedings_4_Web_Broschuere_Link.pdf)

# Filling the Digital Preservation Gap : integrating Archivemata and Hydra for Research Data Management

**Where are we now?**  
Hull and York are Fedora users. Hull has an established Hydra repository. Both of us need to preserve research data and other content long-term.

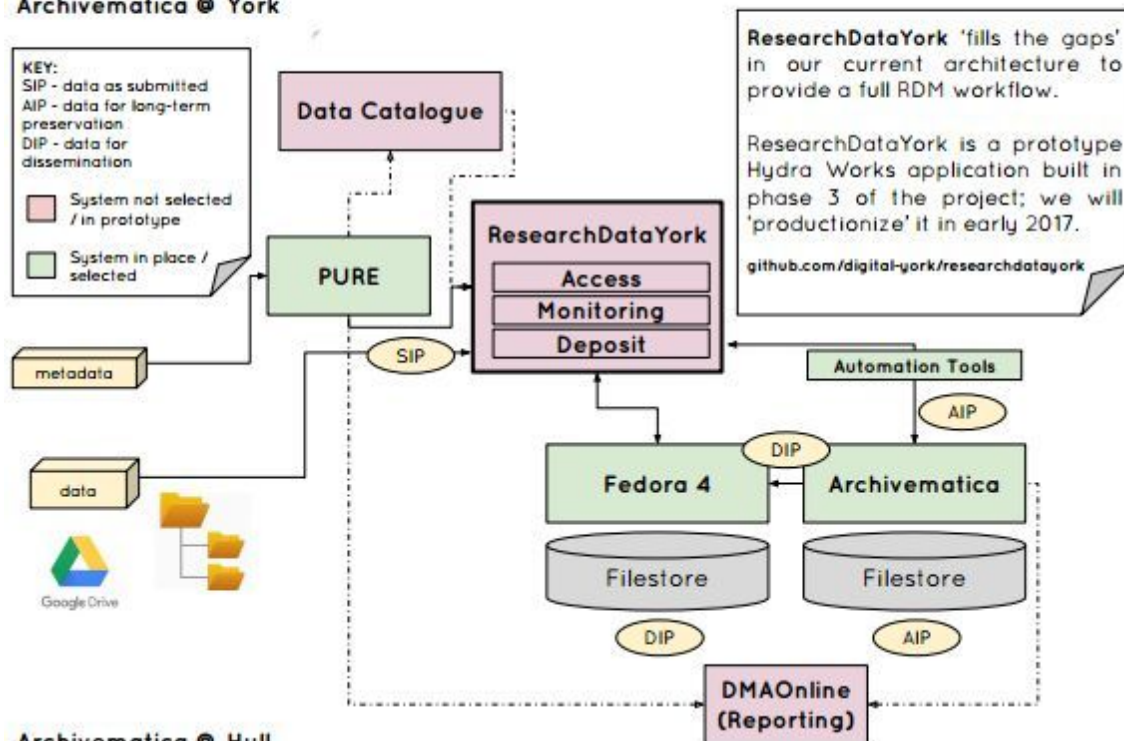
**Why the need?**  
We both are committed to long-term preservation in the repository but UK universities have a mandate to preserve research data in particular.

**Why Archivemata?**  
Archivemata is a well-respected, open-source tool which seemed to offer much of the functionality we needed.

**What we have done?**  
We have had 3 phases of Jisc project funding, including a feasibility study and work with Archivemata to improve its applicability to research data.

**What are we doing now?**  
We are finalising "proof-of-concept" systems to demonstrate Archivemata and Hydra being used for preserving research data.

### Archivemata @ York



### Archivemata @ Hull

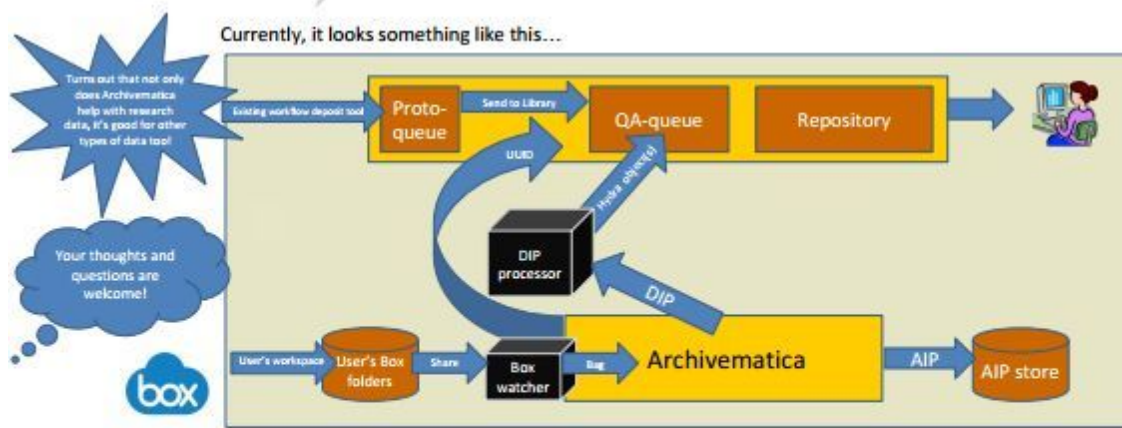


Figure 5: Poster at Hydra Connect

## Other publications

### The National Archives

Details of our project were featured as a case study within *Consultation on a new Strategic Vision for the Archives Sector*

<http://www.nationalarchives.gov.uk/documents/archives/consultation-strategic-vision-for-archives.pdf>

### Nestor

An edited version of one of our project blog posts “My first file format signature” was translated into German for inclusion in a new nestor publication series. Nestor is the German competence network for digital preservation<sup>73</sup>:

[http://files.dnb.de/nestor/kurzartikel/thema\\_03-Meine\\_erste\\_Dateiformatsignatur.pdf](http://files.dnb.de/nestor/kurzartikel/thema_03-Meine_erste_Dateiformatsignatur.pdf)

## Blogs

The project team have been blogging about the project on the University of York’s Digital Archiving blog: <http://digital-archiving.blogspot.co.uk/>

Blog posts relating to the project (either describing or informing our phase 3 work) are listed below:

Title of blog post	Date of release (2016)	No of views <sup>74</sup>
<b>New "Filling the Digital Preservation Gap" report released</b> <a href="http://digital-archiving.blogspot.co.uk/2016/02/new-filling-digital-preservation-gap.html">http://digital-archiving.blogspot.co.uk/2016/02/new-filling-digital-preservation-gap.html</a> -	5th February 2016	513
<b>Kicking off phase 3 of "Filling the Digital Preservation Gap"</b> <a href="http://digital-archiving.blogspot.co.uk/2016/04/kicking-off-phase-3-of-filling-digital.html">http://digital-archiving.blogspot.co.uk/2016/04/kicking-off-phase-3-of-filling-digital.html</a>	1st April 2016	338
<b>Research data - what does it *really* look like?</b> <a href="http://digital-archiving.blogspot.co.uk/2016/05/research-data-what-does-it-really-look.html">http://digital-archiving.blogspot.co.uk/2016/05/research-data-what-does-it-really-look.html</a>	31st May 2016	1302
<b>Modelling Research Data with PCDM</b> <a href="http://digital-archiving.blogspot.co.uk/2016/07/modelling-research-data-with-pcdm.html">http://digital-archiving.blogspot.co.uk/2016/07/modelling-research-data-with-pcdm.html</a>	4th July 2016	971

<sup>73</sup> [http://www.langzeitarchivierung.de/Subsites/nestor/EN/Home/home\\_node.html](http://www.langzeitarchivierung.de/Subsites/nestor/EN/Home/home_node.html)

<sup>74</sup> as of 17th October 2016

<b>New research data file formats now available in PRONOM</b> <a href="http://digital-archiving.blogspot.co.uk/2016/07/new-research-data-file-formats-now.html">http://digital-archiving.blogspot.co.uk/2016/07/new-research-data-file-formats-now.html</a>	4th July 2016	597
<b>Research data is different</b> <a href="http://digital-archiving.blogspot.co.uk/2016/08/research-data-is-different.html">http://digital-archiving.blogspot.co.uk/2016/08/research-data-is-different.html</a>	5th August 2016	510
<b>My first file format signature</b> <a href="http://digital-archiving.blogspot.co.uk/2016/08/my-first-file-format-signature.html">http://digital-archiving.blogspot.co.uk/2016/08/my-first-file-format-signature.html</a>	19th August 2016	614
<b>Filling the Digital Preservation Gap - a brief update</b> <a href="http://digital-archiving.blogspot.co.uk/2016/08/filling-digital-preservation-gap-brief.html">http://digital-archiving.blogspot.co.uk/2016/08/filling-digital-preservation-gap-brief.html</a>	30th August 2016	309
<b>UK Archivemata group at Lancaster</b> <a href="http://digital-archiving.blogspot.co.uk/2016/09/uk-archivemata-group-at-lancaster.html">http://digital-archiving.blogspot.co.uk/2016/09/uk-archivemata-group-at-lancaster.html</a>	16th September 2016	265
<b>File format identification at Norfolk Record Office (a guest post)</b> <a href="http://digital-archiving.blogspot.co.uk/2016/09/file-format-identification-at-norfolk.html">http://digital-archiving.blogspot.co.uk/2016/09/file-format-identification-at-norfolk.html</a>	21st September 2016	363
<b>Some highlights from iPRES 2016</b> <a href="http://digital-archiving.blogspot.co.uk/2016/10/some-highlights-from-ipres-2016.html">http://digital-archiving.blogspot.co.uk/2016/10/some-highlights-from-ipres-2016.html</a>	11th October 2016	307

Table 3: Project blog posts released during Phase 3 and number of page views recorded

## Project website

The project website is hosted at the Borthwick Institute for Archives at the University of York. This is available at <http://www.york.ac.uk/borthwick/projects/archivemata/>.

In the period from the 9th December 2015 (when these stats were last reported in the Phase 2 report) to the 16th October 2016, there were 559 pageviews representing 480 unique visits to the page. Average time spent on the page was 4 minutes and 11 seconds.

The screenshot shows the Borthwick Institute for Archives website. The header includes the University of York logo, the name 'Borthwick Institute for Archives', a search bar, and navigation links for 'Borthwick' and 'University'. A breadcrumb trail reads 'Home > Borthwick Institute for Archives > Projects > Archivemata'. A left-hand navigation menu lists various site sections, with 'Projects' expanded to show 'Planning your funded research project'. The main content area features a banner titled 'Filling the Digital Preservation Gap' with the subtitle 'Report on Archivemata for research data now available.' and the Jisc logo. Below the banner is a paragraph of text and a diagram illustrating the digital preservation workflow.

In March 2015 we started work on phase 1 of a project to explore the potential of the digital preservation solution [Archivemata](#) to help manage research data that academics within the University produce. This project is being carried out with funding from [Jisc](#) as part of their [Research Data Spring](#) program and is being carried out in collaboration with the [University of Hull](#).

The diagram illustrates the digital preservation workflow. It starts with a 'SIP' (Source Information Package) box. An arrow points from the SIP to a 'Transfer' box, which then leads to an 'Ingest' box. From 'Ingest', the flow goes to 'Archival Storage', which then leads to 'Access'. Below this main flow, there are three storage components: 'Transfer Store', 'AIP Store', and 'DIP Store'. Arrows indicate that data from 'Transfer' goes to 'Transfer Store', from 'Ingest' to 'AIP Store', and from 'Access' to 'DIP Store'. The 'AIP Store' and 'DIP Store' boxes contain examples of storage systems: 'eg. arkivum, local filestore, local preservation repository ...' for AIP Store, and 'eg. local access repository, archives management system, local filestore ...' for DIP Store. A feedback loop arrow points from the 'Transfer Store' back to the 'SIP' box.

Figure 6: The project website

## Project reports

In mid July at the second sandpit workshop our phase 1 project report was made available via Figshare (<http://dx.doi.org/10.6084/m9.figshare.1481170>). This report has been viewed 3113 times in the intervening period and downloaded 620 times<sup>75</sup>.

In February 2016 our phase 2 project report was made available via Figshare (<https://dx.doi.org/10.6084/m9.figshare.2073220>). This report has been viewed 3265 times in the intervening period and downloaded 465 times<sup>76</sup>.

These reports are also available from the University of Hull repository.

<sup>75</sup> Statistics collected on 17th October 2016

<sup>76</sup> Statistics collected on 17th October 2016

## Glossary

**AIP:** Archival Information Package - processed information sent to the archival store for preservation

**API:** Application Programming Interface - protocol that allows integration between software for example to allow third-party developers to create additional functionality for a piece of software

**Automation Tools:** A set of python scripts, that are designed to automate the processing of transfers in an Archivematica pipeline<sup>77</sup>

**Box:** A service used at the University of Hull for secure file sharing, storage and collaboration<sup>78</sup>

**DC:** (in the context of this report) Dublin Core metadata

**DIP:** Dissemination Information Package - information created from the material being archived intended for sending to a user

**DMAOnline:** Data Management Administration Online - a Data Spring Project based at Lancaster University it seeks to provide a single dashboard view of its RDM activities<sup>79</sup>

**DROID:** (Digital Record Object Identification) An automatic file format identification tool from The National Archives

**EPSRC:** Engineering and Physical Sciences Research Council

**Fedora:** (in the context of this report) An open-source digital repository platform<sup>80</sup>

**Figshare:** A repository where researchers, institutions and publishers can share research outputs<sup>81</sup>

**Hydra:** A repository solution based on a number of “best-of-breed” open-source components, including Fedora<sup>82</sup>

**JHOVE:** JSTOR/Harvard Object Validation Environment is an extensible software framework for performing format identification, validation, and characterization of digital objects<sup>83</sup>

**json:** JavaScript Object Notation - lightweight data-interchange format that is easily read by humans and parsed by machines and is supported by all modern browsers

---

<sup>77</sup> <https://github.com/artefactual/automation-tools>

<sup>78</sup> <https://www.box.com/en-gb/home>

<sup>79</sup> <http://www.dmao.info>

<sup>80</sup> <http://www.fedora-commons.org/>

<sup>81</sup> <http://www.figshare.com>

<sup>82</sup> <http://projecthydra.org/>

<sup>83</sup> <http://openpreservation.org/technology/products/jhove/>



**METS:** The METS metadata schema is a widely adopted standard for encoding descriptive, administrative, and structural metadata

**PCDM:** Portland Common Data Model is a flexible, extensible domain model that is intended to underlie a wide array of repository and DAMS applications. Currently being implemented by the Hydra and Islandora communities.

**PREMIS:** The PREMIS Data Dictionary for Preservation Metadata is the international standard for metadata to support the preservation of digital objects and ensure their long-term usability. Version 3 of the standard has just been released.<sup>84</sup>

**PRONOM:** A resource provided by the National Archives in the UK providing definitive information about file formats, software products and other technical components required to support long-term access to electronic records and other digital objects of cultural, historical or business value.<sup>85</sup>

**PUID:** PRONOM Unique Identifier<sup>86</sup> - an extensible scheme for providing persistent unique identifiers for records in the PRONOM registry

**PURE:** Research information system from Elsevier used at York.

**Rails:** Ruby on Rails - web application framework that provides structures for a database, web service or web pages. It uses json or XML for data transfer and html for display.

**RDM:** Research Data Management

**SHA256:** (and SHA-512, md5) hash algorithms that create the unique digital signature or checksum that can be used to prove a file has not changed over time. A single change to a file would produce a different hash value using the same algorithm.

**SIP:** Submission Information Package - information sent from its producer for archiving

**UUID:** a universally unique identifier

---

<sup>84</sup> <http://www.loc.gov/standards/premis/>

<sup>85</sup> <https://www.nationalarchives.gov.uk/PRONOM/Default.aspx>

<sup>86</sup> <http://www.nationalarchives.gov.uk/aboutapps/pronom/puid.htm>

## Appendix 1: Draft instructions for use of Hull implementation

These draft instructions are aimed at researchers at the University of Hull and are intended to give the reader a sense of how a Hull production system would be used. Whilst most of the functionality referred to has been implemented in our proof of concept system, some has not: in particular, only a subset of the Dublin Core metadata terms are currently supported and only dataset, journal article, book, book chapter and photograph are functional from the list of possible content models.

---

### Depositing digital content for preservation and discovery using Box folders

#### Overview

The process described here allows a user to deposit materials for long term digital preservation and to have a record of them placed in the University's digital repository, Hydra. Optionally, copies of the files can be made available through Hydra for download.

#### Deposit options

Material to be deposited may be a single file or multiple files, it may also be a folder structure containing multiple files in an organised hierarchy.

The resulting repository content may be a single object or multiple objects. These objects describe the file(s) that have been placed in the University's preservation store and may optionally contain copies of the data for download.

As a depositor, in addition to your files for deposit, you must provide one or more simple text files, as described below, that indicate how you wish your content to be dealt with.

#### One or more files to be represented in a single object

This first option allows you to deposit one or more files as a group. The files will be kept together and there will be a single object in the repository that describes them all. Optionally, any or all of the files can be copied to the repository for download.

1. Create a folder in your Box space in which to assemble your content.
2. Copy into this folder the files that you wish to deposit.

3. You must now add to the Box folder one or more simple text files that will determine how your content is processed through our automated system. Any simple text editor can be used to create the files, for instance Notepad in Windows or TextEdit on an Apple.

Create a file called "description.txt"

- It **must** contain a title that describes the materials deposited
- It **must** contain the name(s) of the creator(s) of the material in the form "Lastname, firstname" or "Lastname, Firstname Initial(s)." Multiple authors should be separated by semicolons.
- It **must** contain a list of files ("visibleFiles") that are to be made available for download in Hydra. The list can be blank (no visible files) or may contain the word "all" (all files to be visible). Multiple filenames must be separated by semicolons.
- It may contain a "contentModel" field that describes the primary type of material being deposited. If this field is omitted, the default is "dataset". 1) (See list in Table
- It may contain additional Dublin Core metadata fields (see list in Table 2), for instance a description of the materials deposited or a specific citation that you would like used in references. **Not all DC fields are currently implemented.**

The following are valid examples of "description.txt". Note that the field labels begin with a lower case letter and are followed by a colon.

```
title: Data package with all files visible in Hydra

creator: Mouse, Mickey; Mouse, Minnie; Duck, Donald A.

description: Specimen data package with a small number of files.

contentModel: dataset

visibleFiles: all
```

```
title: Data package with all files visible in Hydra and with additional
metadata

creator: Mouse, Mickey; Mouse, Minnie; Duck, Donald A.

description: Specimen data package with a number of files, all to be
downloadable through Hydra.

visibleFiles: all

citation: Mouse, Mickey; Mouse, Minnie; Duck, Donald A. (2016) "Research
data from the Disney Project" University of Hull

subject: Cartoons; Disney, Walt; Animation
```

```
title: Data package with selected files visible in Hydra

creator: Mouse, Mickey; Mouse, Minnie; Duck, Donald A.

description: Specimen data package with a small number of files. Some
```

```
files for preservation only - not to be made available in Hydra, three
copy files available for download.

visibleFiles: file1.jpg; file2.pdf; file3.log
```

```
title: Data package with no files visible in Hydra

creator: Mouse, Mickey; Mouse, Minnie; Duck, Donald A.

description: Specimen data package with a small number of files. Files
for preservation only - copies not to be made available in Hydra, just a
metadata record.

contentModel: dataset

visibleFiles:
```

artwork	drawing	letter	policyOrProcedure
book	event	licence	presentation
bookChapter	genericContent	map	regulation
conferencePaper	guidance	meetingPapers	report
conferenceAbstract	handbook	movingImage	software
conferencePoster	internetPublication	musicScore	sound
dataset	journalArticle	newsletterArticle	
diagram	learningMaterials	photograph	

Table 1: Acceptable content model types

abstract	coverage	hasFormat	isVersionOf	requires
accessRights	created	hasPart	language	rights
accrualMethod	creator	hasVersion	license	rightsHolder
accrualPeriodicity	date	identifier	mediator	source
accrualPolicy	dateAccepted	instructionalMethod	medium	spatial
alternative	dateCopyrighted	isFormatOf	modified	subject
audience	dateSubmitted	isPartOf	provenance	tableOfContents
available	description	isReferencedBy	publisher	temporal
bibliographicCitation	educationLevel	isReplacedBy	references	title
conformsTo	extent	isRequiredBy	relation	type
contributor	format	issued	replaces	valid

Table 2: Dublin Core metadata acceptable in description files

4. If you have citations of a paper or papers related to data that you are depositing, for instance in Bibtex format, you may add this/these to your folder in a file called "citations.bib". **Other formats are not currently supported.**
5. If there is information you wish to pass to the repository staff in the University Library, create another simple text file called "readme.txt" in your Box folder. The text in this folder will be visible to them when they come to check the repository object that is created and they can then contact you if necessary.
6. If you are satisfied that everything is as you want it to be, use the Box "Share" option for your folder to share it with the email address [archivematica@hull.ac.uk](mailto:archivematica@hull.ac.uk).
7. Essentially, that's it! If you check back from time to time you will see that your folder name has the status of your material appended. So, for instance, a folder called "Mydata" will go through a number of stages starting with "Mydata – processing" and ending with "Mydata – all processing complete". Once you reach this last stage you can, should you wish, delete the folder; do not delete it until then!

### More than one file, each to be represented in its own object

This second option allows you to create multiple objects from the files that you deposit. Optionally, any or all of the files in these objects can be copied to the repository for download.

1. Create a folder in your Box space in which to assemble your content.
2. Copy into this folder the files that you wish to deposit.
3. You must now add to the Box folder a comma-separated-values (csv) file that will determine how your content is processed through our automated system. A csv file is most easily created in a spreadsheet application

Create a set of headings in row 1 corresponding to the metadata fields that you intend passing into the system and then a row of metadata for each file to be deposited.

- You **must** first pass the filename that each row of metadata applies to
- You **must** pass a title that describes the materials deposited in that object
- You **must** pass the name(s) of the creator(s) of the material in each object in the form “Lastname, firstname” or “Lastname, Firstname Initial(s).” Multiple authors for the same material should be separated by semicolons.
- Each file row **must** contain a “visibleFiles” entry if the file is to be available for download in Hydra. The visibleFiles entry should be set to “all” if a copy of the file is to be made available for download. If it is left blank, no copy file will be available.
- You may pass a “contentModel” field that describes the type of material being deposited. If this field is omitted, the default is “dataset”. (See list in Table1)
- It may contain additional Dublin Core metadata fields (see list in Table 2), for instance a description of the materials deposited or a specific citation that you would like used in references. **Not all DC fields are currently implemented.**

The following is an example of “description.csv” as laid out in a spreadsheet app. Note that the field labels begin with a lower case letter.

filename	title	creator	description	subject	contentModel	visibleFiles
picture.jpg	My picture	Mouse, Mickey	Darling Minnie	Minnie Mouse	photograph	
story.ppt	My presentation	Duck, Donald A.	My research work	Nautical history	presentation	all
data.xls	My spreadsheet	Mouse, Minnie; Mouse, Mickey	My data			all

Assuming that the named files are present in your folder, processing this description.csv would result in three objects in Hydra with corresponding data in the preservation store. The first object would contain the metadata for a preserved image file, but the image would not be available for download. The second object would describe the presentation and make a copy available for download. The final object would describe a spreadsheet of data and make a copy available for download; this would be recognised as a dataset within the repository even though the contentModel was not specified.

4. If there is information you wish to pass to the repository staff in the University Library, create another simple text file called “readme.txt” in your Box folder. The text in this folder will be visible to them when they come to check the repository object that is created and they can then contact you if necessary.
5. If you are satisfied that everything is as you want it to be, use the Box “Share” option for your folder to share it with the email address [archivematica@hull.ac.uk](mailto:archivematica@hull.ac.uk).
6. Essentially, that’s it! If you check back from time to time you will see that your folder name has the status of your material appended. So, for instance, a folder called “Mydata” will go through a number of stages starting with “Mydata – processing” and ending with “Mydata – all processing complete”. Once you reach this last stage you can, should you wish, delete the folder; do not delete it until then!

### A hierarchy of files to be represented in a single object

This final option allows you to deposit a hierarchy (folder structure) of files. The files will be kept together in their folder structure which is turned into a zip file. There will be a single object in the repository that describes them and, optionally, the zip file can be copied to the repository for download.

1. Create a folder in your Box space in which to assemble your content.
2. Copy into this folder the folder that contains the hierarchy so that you get a folder within the Box folder
3. You must now add to the Box folder one or more simple text files that will determine how your content is processed through our automated system. Any simple text editor can be used to create the files, for instance Notepad in Windows or TextEdit on an Apple.

Create a file called “description.txt”

- It **must** contain a title that describes the materials deposited
- It **must** contain the name(s) of the creator(s) of the material in the form “Lastname, firstname” or “Lastname, Firstname Initial(s).” Multiple authors should be separated by semicolons.
- It **must** contain a “visibleFiles” entry to show whether a copy of the zip file should be made available for download in Hydra. The entry can be blank (file not available) or may contain the word “all” (file to be downloadable).
- It may contain a “contentModel” field that describes the primary type of material being deposited. If this field is omitted, the default is “dataset”. (See list in Table 1)
- It may contain additional Dublin Core metadata fields (see list in Table 2), for instance a description of the materials deposited or a specific citation that you would like used in references. **Not all DC fields are currently implemented.**

The following are valid examples of “description.txt”. Note that the field labels begin with a lower case letter and are followed by a colon.

```
title: Hierarchical data package with zip file visible in Hydra and with
additional metadata

creator: Mouse, Mickey; Mouse, Minnie; Duck, Donald A.

description: Specimen data package with a tree of files. Files for preservation
only - not to be made available in Hydra.

visibleFiles: all

citation: Mouse, Mickey; Mouse, Minnie; Duck, Donald A. (2016) "Research data
from the Disney Project" University of Hull

subject: Cartoons; Disney, Walt; Animation

contentModel: dataset
```

```
title: Hierarchical data package with additional metadata - zip file for
preservation only - no copy downloadable through Hydra

creator: Mouse, Mickey; Mouse, Minnie; Duck, Donald A.

description: Specimen data package with a tree of files. Files for preservation
only - not to be made available in Hydra.

visibleFiles:

citation: Mouse, Mickey; Mouse, Minnie; Duck, Donald A. (2016) "Research data
from the Disney Project" University of Hull

subject: Cartoons; Disney, Walt; Animation
```

4. If you have citations of a paper or papers related to data that you are depositing, for instance in Bibtex format, you may add this/these to your folder in a file called "citations.bib".
5. If there is information you wish to pass to the repository staff in the University Library, create another simple text file called "readme.txt" in your Box folder. The text in this folder will be visible to them when they come to check the repository object that is created and they can then contact you if necessary.
6. If you are satisfied that everything is as you want it to be, use the Box "Share" option for your folder to share the Box folder (not the folder contained within it) with the email address [archivematica@hull.ac.uk](mailto:archivematica@hull.ac.uk).
7. Essentially, that's it! If you check back from time to time you will see that your folder name has the status of your material appended. So, for instance, a folder called "Mydata" will go through a number of stages starting with "Mydata – processing" and ending with "Mydata – all processing complete". Once you reach this last stage you can, should you wish, delete the folder; do not delete it until then!



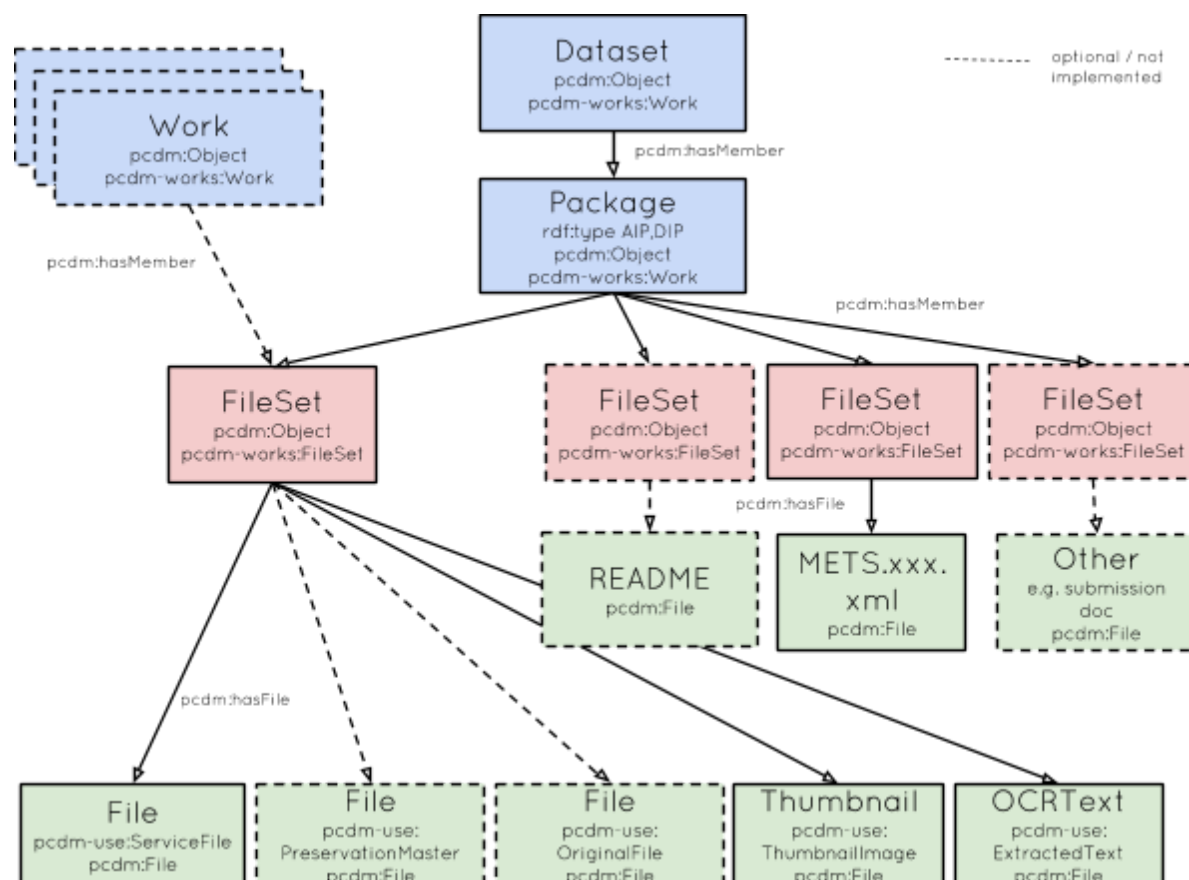
## Appendix 2: A Draft PCDM-based Data Model for Datasets

The Portland Common Data Model (PCDM) is “a flexible, extensible domain model that is intended to underlie a wide array of repository and DAMS applications”<sup>87</sup>. It is informing the latest developments of Hydra and Islandora in their implementations on top of Fedora 4. Aligning our model with PCDM means that we are designing with interoperability in mind, and adding to the growing list of example PCDM models.

For the ‘live’ document, please see:

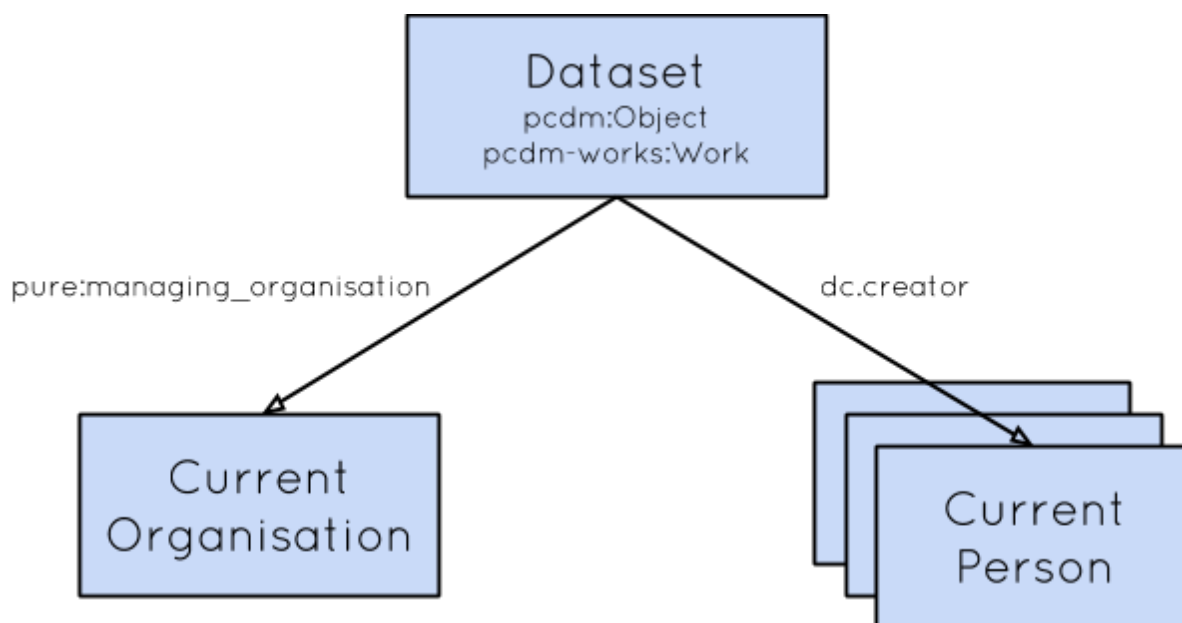
<https://docs.google.com/document/d/1QPw9kLqRFzI5aStRr3nlBqqj5BzkI3eFZY4zNjOmo8w/edit#heading=h.1jggb1ac1v2v>

### PCDM Model



<sup>87</sup> <https://github.com/duraspace/pcdm>

## People and Organisations Model



## Namespaces

am: <<http://dlib.york.ac.uk/ontologies/oais-archivematica#>>

dlib: <<http://dlib.york.ac.uk/ontologies/generic#>>

dc: <<http://purl.org/dc/terms/>>

identifiers: <<http://id.loc.gov/vocabulary/identifiers/>>

sorg: <<http://schema.org/>>

foaf: <<http://xmlns.com/foaf/0.1/>>

## Models

### Dataset

A Dataset is a discrete set of related data files. The files can be of any type. In PURE a single Dataset metadata record has a one-to-one relationship to a Dataset as described here. A Dataset does not, however, *require* an accompanying PURE record. Local descriptive metadata can be added with the exception of pureUuid which *must* be derived from PURE.

Code:

<https://github.com/digital-york/dlibhydra/blob/master/lib/dlibhydra/models/works/dataset.rb>

Class	Property	Expected Object Type	Usage
<b><u>Dataset</u></b>	rdf:type	URI	dcat:dataset
	dc:title	Literal (String)	1 (datacite mandatory) From PURE

	dc.creator	Current Person	1..n (datacite mandatory) From PURE
	identifiers:doi	URI	0..1 (datacite mandatory) From PURE
	dc.publisher	Literal (String)	1 (datacite mandatory) From PURE
	dc.available	Literal (String)	1 (datacite mandatory) From PURE
	dlib:pureUuid	Literal (String)	1 From PURE
	pure:managingUnit	Current Organisation	1 From PURE
	pure:pureLink	URI	0..n From PURE
	dc:accessRights	Literal (String)	1 From PURE
	dlib:embargoRelease Date	Literal (String)	1
	dlib:lastAccess	Literal (String)	0..1
	dlib:retentionPolicy	URI/Object	0..1
	dlib:for_indexing	Literal (String)	1 From PURE
	dlib:restriction_note s	Literal (String)	0..1

	pcdm:hasMember	Package	0..n
--	----------------	---------	------

Data available from PURE but not currently added:

- Publications and projects
- People with roles other than 'Creator'
- Description
- Geographical
- Temporal
- Date of production
- Workflow / Visibility

Data not available from PURE Web Services:

- Embargo
- Contact person
- Relations to other datasets
- Restrictions narrative

Class	Property	Expected Object Type	Usage
<b><u>CurrentPerson</u></b>	rdf:type	URI	sorg:Person pure:PurePerson
	skos:prefLabel	Literal (String)	1
	foaf:familyName	Literal (String)	1
	foaf:givenName	Literal (String)	1
	pure:pureUuid	Literal (String)	1
	pure:puretype	Literal (String)	1

Class	Property	Expected Object Type	Usage
<b><u>CurrentOrganisation</u></b>	rdf:type	URI	sorg:Organisation pure:PureOrganisation
	skos:prefLabel	Literal (String)	1

	foaf:name	Literal (String)	1
	pure:pureUuid	Literal (String)	1
	pure:puretype	Literal (String)	1

### **Package**

A Package represents a deposit of data files that form part of a larger Dataset. A Package with the `rdf:type am:ArchivalInformationPackage (AIP)` will contain sufficient metadata to identify the location of the stored AIP (aipUuid as a minimum). A Package with the `rdf:type am:DisseminationInformationPackage (DIP)` will contain references to the files comprising the DIP and the dipUuid as a minimum). A single Package can be both AIP and DIP. Folder structure for the original deposit shall be retained in such a way that it can be re-created to the end user, for example in an Archivematica METS file. Other files may be included in the Package, for example submission documentation or a readme file. These should be distinguished from the data files themselves.

Code:

<https://github.com/digital-york/dlibhydra/blob/master/lib/dlibhydra/models/works/package.rb>

Class	Property	Expected Object Type	Usage
<b>Package</b>	<code>rdf:type</code>	URI	dlib:Package am:ArchivalInformationPackage am:DisseminationInformationPackage
	<code>skos:prefLabel</code>	Literal (String)	1
	<code>am:aipUuid</code>	Literal (String)	0..1 from Archivematica
	<code>am:aipStatus</code>	Literal (String)	1 from Archivematica
	<code>am:aipSize</code>	Literal (String)	0..1 from Archivematica

	am:aipCurrentPath	Literal (String)	0..1 from Archivemata
	am:aipCurrentLocation	Literal (String)	0..1 from Archivemata
	am:dipUuid	Literal (String)	0..1 from Archivemata
	am:dipStatus	Literal (String)	0..1 from Archivemata
	am:dipSize	Literal (String)	0..1 from Archivemata
	am:dipCurrentPath	Literal (String)	0..1 from Archivemata
	am:originPipeline	Literal (String)	0..1 from Archivemata (dip)
	am:dipCurrentLocation	Literal (String)	0..1 from Archivemata
	am:aipResourceUri	Literal (String)	0..1 from Archivemata
	am:dipResourceUri	Literal (String)	0..1 from Archivemata
	dlib:requestorEmail	Literal (String)	0..n