

Opening up research data at Essex: experiments with EPrints

Thomas Ensom & Alexis Wolton

University of Essex



Today we're going to be reporting on our work developing a pilot EPrints data repository at the University of Essex

Background

JISC

- JISC-funded Managing Research Data programme project
- Developing MRD infrastructure and policy at the University of Essex
- Utilising UK Data Archive expertise in data management
- University of Essex has an EPrints IR, we have piloted a separate instance for data
- How best to adapt it for data?



Research Data @Essex is a JISC-funded project aiming to develop a sustainable research data management and sharing infrastructure, built on best practise guidance from the research data management community and UK Data Archive expertise.

The University has an EPrints institutional repository, and an important part of the project is setting up a data instance building on the same implementation. Today we will be talking about our approach to adapting it to better suit collections of data.

Accommodating diverse data

- We spoke to researchers from four pilot departments:
 - **Language and Linguistics**
 - **Biological Sciences**
 - **Computing and Electronic Systems**
 - **Business School**
- Interviews, inventory and sample data gathered for testing
- On-going contact and consultation with researchers throughout development

We have been working with four pilot departments, covering a broad range of disciplines Essex.

This has involved interviewing researchers and asking for sample data collections to trial ingest into the test-bed repository. We have continued to work with these researchers as the project has continued.

Design ethos

- Minimising barriers for researchers to deposit
- ...while satisfying requirements for re-use (i.e. sufficient metadata and documentation)

- Yes, we want deposit to be as easy as possible
 - But we want the data to be more than just a ticked box
 - But filling repository with rubbish is pointless

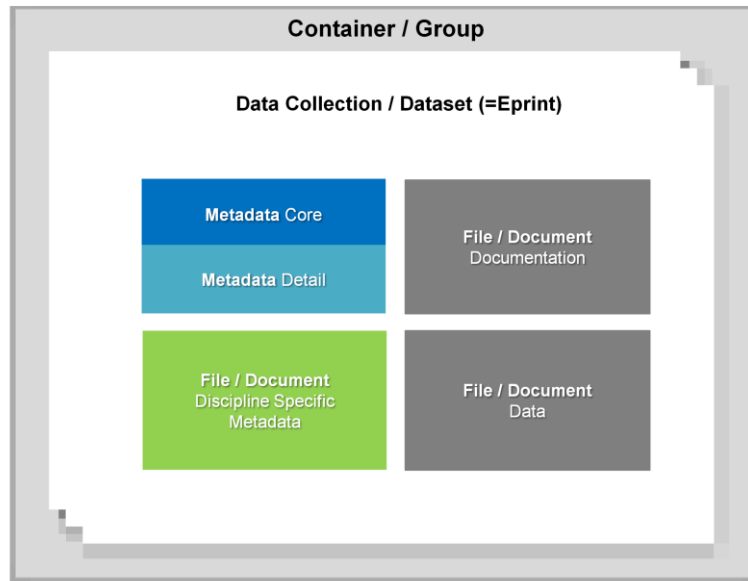
How should we compromise?

Our design maxim has been to minimise barriers while enabling re-use.

Yes, we do want the deposit process to be as straightforward as possible for the user, but we also want publishing data to be more than just a tick in a box – we need rich metadata and as much documentation as possible.

We realise though, that asking for too much you might end up with a load of junk. Can we find a compromise?

Groups, data collections & files



Research data differs greatly from research publications, in level of complexity. An article is typically a single file, while a dataset or data collection could (conceivably) be hundreds of files with multiple relationships between them. So let's define our terms in the EPrints context.

A data collection is our 'eprint', the key unit. This could be anything from a set of audio interviews with transcripts, to a single spreadsheet. Within each collection there is a set of descriptive metadata, and a series of files. These files can be of the types: data, documentation and metadata. Data collections can be grouped inside larger containers. For example, a series of datasets produced as part of an umbrella project. We are trying to decide whether these higher level groupings should be formal or user instigated.

Metadata

- Extended the default EPrints metadata profile to better suit research data
- Based on existing schema to enable interoperability
- Minimal mandatory elements based on DataCite metadata, to enable DOI minting further down the road
- Rich metadata based on
 - DataShare, for Edinburgh digital repository
 - INSPIRE, and EU standard for data with geospatial content
 - DDI (Data Documentation Initiative), from the social science community but now being by others e.g. to describe biomedical data



We have developed a metadata profile built on the DataCite schema - we intend to mint DataCite DOIs (we intend to mint DataCite DOIs further down the road). To improve descriptive richness, we also examined several other schema including:

DataShare – work done at Edinburgh University for sharing research datasets

INSPIRE – for geospatial data, but also providing a neat generic description of research data

DDI - a metadata schema originally from the social science community, but now finding applications in biomedical research due to it's depth and power

This will be published within the next few months

Implementation

- The next phase of the project looked at rendering a Data Collection in EPrints
- We knew how we wanted to describe our data, but we faced a number of challenges turning our metadata into a useable EPrints screen:
 - how to display a metadata schema that had been extended by nearly 50%
 - How to clearly present the 3 tiers of our data collection
 - how best to group together and display the different files that make up the collection.
- We looked around at what others were doing / had done
 - Ecrystals – orders files according to type
 - Kulture – uses ‘Containers’ that inherit metadata

We’ve talked about the requirements gathering exercise and the metadata profiles that were been generated as a consequence.

The next phase was to render a Data Collection in Eprints

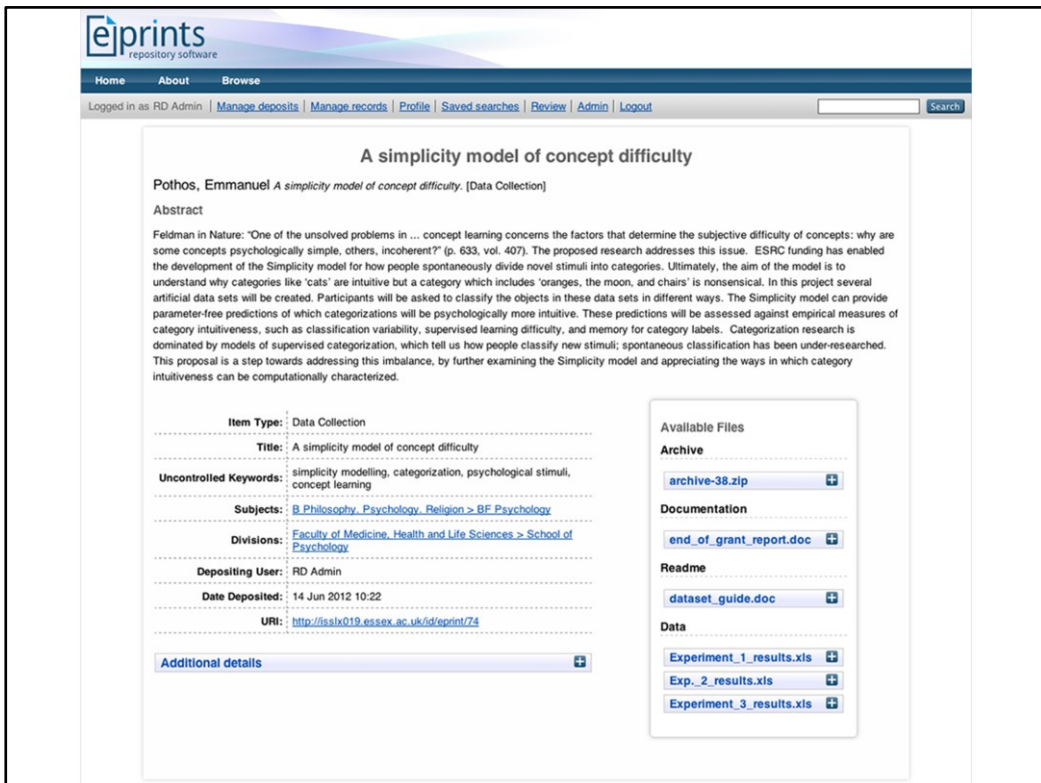
We were faced with a number of challenges

We knew how we wanted to describe our data, but we faced a number of challenges turning our metadata into a useable Eprints screen

- how to display a metadata schema that had been extended by nearly 50%
- How to clearly present the 3 tiers of our data collection
- how could we group together and display the different files that make up the collection.

We looked around at what others were doing / had done in the past.

We were particularly impressed by the way Ecrystals ordered files according to file type,
And were interested in the way Kulture was using ‘Containers’ that inherited metadata



This is the current RD@Essex citation screen on our test server – it's a work in progress.

Looks a lot like base Eprints, but has key differences

We wanted an organised, tidy screen, but without sacrificing any detail.

We've added two extra components to do the work, hooking into default Eprints javascript to control the amount of metadata onscreen at any one time

We wanted to work with, building on top of the solid base that already exists, but adding Research Data specific elements to it.

The different elements of the screen:

The screenshot shows the EPrints repository interface. At the top, the 'eprints' logo is visible. Below it, there are navigation links: Home, About, Browse. A secondary navigation bar includes: Logged in as RD Admin | Manage deposits | Manage records | Profile | Saved searches | Review | Admin | Logout. The main content area features the title 'A simplicity model of concept difficulty' by Pothos, Emmanuel. Below the title is an abstract. A central box highlights the item's metadata:

Item Type:	Data Collection
Title:	A simplicity model of concept difficulty
Uncontrolled Keywords:	simplicity modelling, categorization, psychological stimuli, concept learning
Subjects:	B Philosophy, Psychology, Religion > BF Psychology
Divisions:	Faculty of Medicine, Health and Life Sciences > School of Psychology
Depositing User:	RD Admin
Date Deposited:	14 Jun 2012 10:22
URI:	http://issth019.essex.ac.uk/id/eprint/74

Below the metadata box is a link for 'Additional details'. To the right, an 'Available Files' section lists:

- Archive:** archive-38.zip
- Documentation:** end_of_grant_report.doc
- Readme:** dataset_guide.doc
- Data:** Experiment_1_results.xls, Exp_2_results.xls, Experiment_3_results.xls

2. Core metadata

Remains mostly the same as with a base EPrints install
 Visible here is the new Data Collection item type we're using:

Data Collection

The screenshot shows the eprints repository interface. At the top, there is a navigation bar with links for Home, About, and Browse. Below this, a user is logged in as 'RD Admin'. The main content area features the title 'A simplicity model of concept difficulty' by Pothos, Emmanuel. An abstract follows, discussing concept learning and the development of a simplicity model. The metadata section includes fields for Item Type (Data Collection), Title, Uncontrolled Keywords, Subjects, Divisions, Depositing User, Date Deposited, and URI. A red box highlights a collapsed 'Additional details' field. To the right, an 'Available Files' section lists several files: archive-38.zip, end_of_grant_report.doc, dataset_guide.doc, and three experiment result files (Experiment_1_results.xls, Exp_2_results.xls, Experiment_3_results.xls).

3. Metadata detail

Rendered as a collapsed box by default

Unrolled forms the complete metadata record

Additional details

Alternative title: [blank]

Creators	Email
Pothos, Emmanuel	e.m.pothos@swansea.ac.uk

Corporate Creators: Emmanuel Pothos

Contributors	Contribution	Name	Email
Research team head		Pothos, Emmanuel	e.m.pothos@swansea.ac.uk

Funders: Economic and Social Research Council

Grant Number: RES-000-23-1541

Geographic coverage: [blank]

East Longitude: 4.446

North Latitude: 51.843

South Latitude: 51.476

West Longitude: -3.73

Collection Methodology: Laboratory-based data collection with non-clinical participants (mostly members of the local university student community). Experiment 1 involved presenting stimuli as cards (N=169), Experiments 2 (N=180) and 3 (N=195) were computer-based (participants sa

Lineage: Results files were anonymised. All experiments had been approved by the Department of Psychology, Swansea University ethics committee. No ethical issues were raised during ethics monitoring or the actual project.

Additional Information: [blank]

Projects: [blank]

Status: Published

Use constraints: [blank]

Publisher: [auto]

Contact Email Address: e.m.pothos@swansea.ac.uk

Comments and Suggestions: [blank]

Experiment_1_results.xls

Exp_2_results.xls

Experiment_3_results.xls

Metadata full unrolled

Shows the extent of metadata we've added

The screenshot shows the Eprints repository interface. At the top, the 'eprints' logo is visible. Below it, there are navigation links: Home, About, Browse. A user is logged in as 'RD Admin', with links for Manage deposits, Manage records, Profile, Saved searches, Review, Admin, and Logout. The main title of the record is 'A simplicity model of concept difficulty' by Pothos, Emmanuel. The abstract discusses concept learning and the development of a simplicity model. On the right side, there is a section titled 'Available Files' which lists files categorized by type: Archive (archive-38.zip), Documentation (end_of_grant_report.doc), Readme (dataset_guide.doc), and Data (Experiment_1_results.xls, Exp_2_results.xls, Experiment_3_results.xls).

Item Type: Data Collection

Title: A simplicity model of concept difficulty

Uncontrolled Keywords: simplicity modelling, categorization, psychological stimuli, concept learning

Subjects: [8 Philosophy, Psychology, Religion > BF Psychology](#)

Divisions: [Faculty of Medicine, Health and Life Sciences > School of Psychology](#)

Depositing User: RD Admin

Date Deposited: 14 Jun 2012 10:22

URI: <http://eprints019.essex.ac.uk/eprint/74>

Additional details

Available Files

Archive

- [archive-38.zip](#)

Documentation

- [end_of_grant_report.doc](#)

Readme

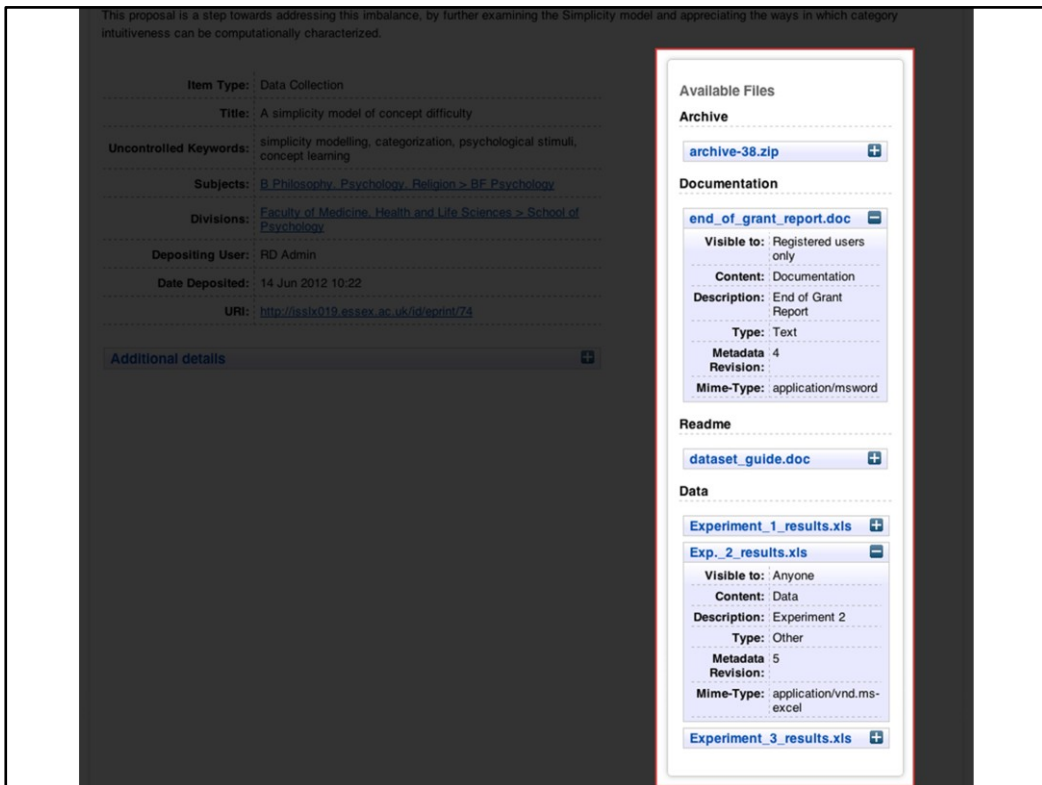
- [dataset_guide.doc](#)

Data

- [Experiment_1_results.xls](#)
- [Exp_2_results.xls](#)
- [Experiment_3_results.xls](#)

4. Documents associated with each Eprint/ Data Collection

We wanted to sort uploaded files according to a type: Archive, Documentation, Readme and Data



4 Documents associated with each Eprint/ Data Collection extended

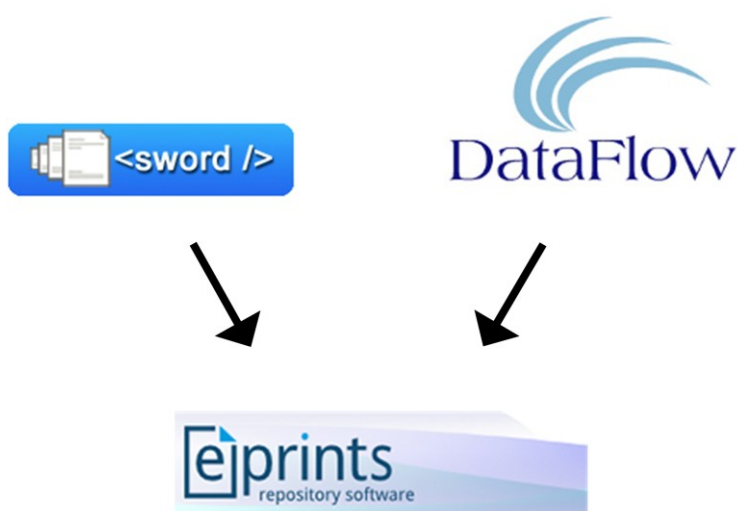
We wanted file metadata to be viewable, but not immediately so again we've used collapsible boxes to keep the screen tidy

Quite a bit of debate as to the best way to sort the files – inspired by the ecrystals layout we tried initially to order by file extension, quickly realised this wasn't going to work as

Different content types could have multiple mime types ie .doc, .pdf .xls etc

We decided it's better to sort using a metafield - content

Data sharing and deposit



While metadata schema was being finalised, we were also asking:
What are the practicalities of technically managing University data?

We were interested in four key areas:

1. Pre-deposit - collaborative local storage environments?
2. Moving data from local storage and depositing data into EPrints
3. Displaying the data files and metadata to users
4. Persistent identification

Sword2:

Interested because Sword could potentially facilitate an easier deposit and therefore encourage researchers to add data.

Dataflow/Datastage

Interested because it's an interesting idea putting together a collaborative environment and Sword based deposit.

Do they help us at this stage of their development? Not really. We need to upload multiple files with very complex metadata. At the moment, it makes much more sense to continue to use the GUI.

We'll be keeping an eye on developments.

What next?

- Continue deposit & ingest testing with real data collections
- Implement DataCite DOIs
 - Using UK Data Archive methodology
- As alternative to discipline specific metadata, test the use of assignable fields for addition of non-standard metadata
 - Gives freedom
 - Recommendation from Southampton's IDMB project
- Should access be at the file or eprint level?
- What data licence options can/should we provide?



Work will continue on testing what we have so far, using real data collections.

We will implement a system for minting data DOIs, adapting a methodology developed at the UK Data Archive

An idea proposed by the IDMB project, Southampton, is to allow depositors to create their own fields using blank fields. This is something we'd like to explore. OR an alternative approach – discipline specific metadata could be included with data and documentation as additional files e.g. XML

Considering access control options required to cover every scenario, including use of embargos and other item level restrictions.

Related is how to licence data – another chance to draw on UK Data Archive expertise

Challenges for the EPrints community

- Dealing with **complex collections**
 - Large file sizes
 - Large number of files
 - Multiple versions of the same files
 - Inter-dependent files e.g. GIS database
- Adding metadata to files during deposit – how to mass apply
- Looking forward to SWORD2 for data (see <http://swordapp.org/2012/07/data-deposit-scenarios>)
- Can we visualise data eprints? Many different file types/formats.

Questions?

researchdataessex.posterous.com

