# Mind the Gap!

Moving from
Aspiration to Experience
in UK Data Management

Leslie Carr,
EPrints, Web & Internet Science Research Group,
University of Southampton

Every day, two or three people fall into the foot-wide gap on platform 2 of Mumbra Railway Station, forty kilometres from Mumbai.

**Nilesh Nikade**

Posted On Monday, November 23, 2009

# Openness Agendas

- Open Access to Research Outputs
- Open Educational Resources
- Open Research Data

- All underpinned by active preservation and curation policies and workflows
- **Repositories are tremendously important**

# Disclaimer

- The following slides refer to implementation work done on the EPrints repository platform (version 3.3) however

  - Other repository platforms do exist…

  - The concept of **repository** is important, not the brand names!

# Repository Scope

## Research Activity

Data Collection

Data Analysis

### Research Business

Research Management
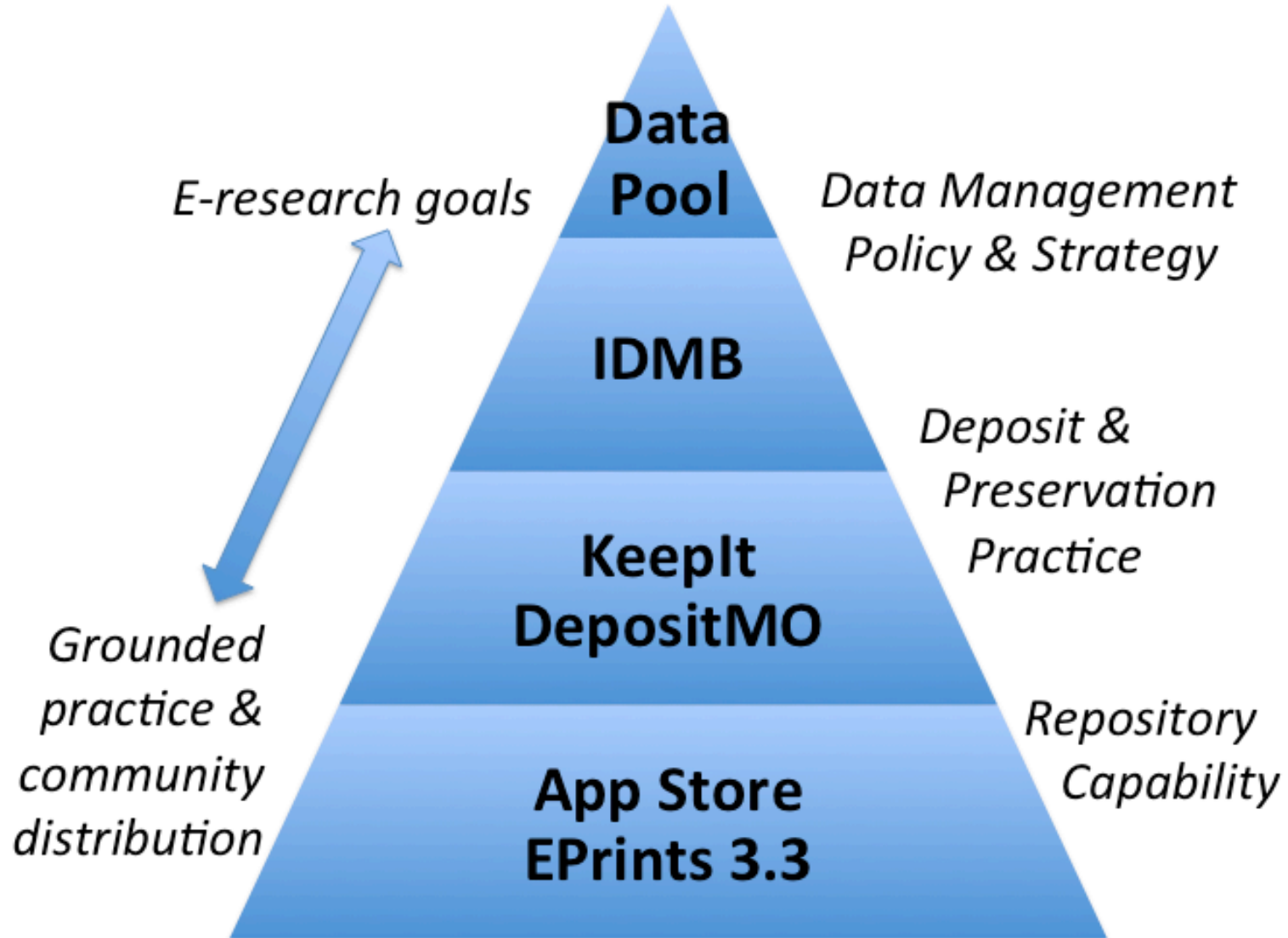
Research Processes

Research Impact

#### Research Content

Research Publications

Educational Resources

Scientific Data

Research Data Management Practice

# Ontological Exactitude *vs...*

# *...vs* Pragmatic Simplicity

- Tempting to get very explicit about subject data classification & control
  - this is a responsibility of data gathering/analysis, not archiving
  - an IR subject taxonomy cannot be augmented for data taxa
    - this is impractical as there are hundreds just for archaeology!

- Keywords attribution should mainly be done by data capturing tools
  - it's too late by the time the repository is involved.
  - may be extracted from files (see the **DepositMO** project)

- Specific structured metadata can be provided
  - As an explicit document whose content is of type *metadata*
  - As a set of linked data relations (URIs)

# IR Pragmatic Policies

- Metadata fields are mainly HUMAN-oriented and advisory

- Explanations of equipment used and processes undertaken should be provided in the README documents

    – Although they could be represented as complex metadata or workflow documents, they are unlikely to be successful in attracting any popular use in the short term
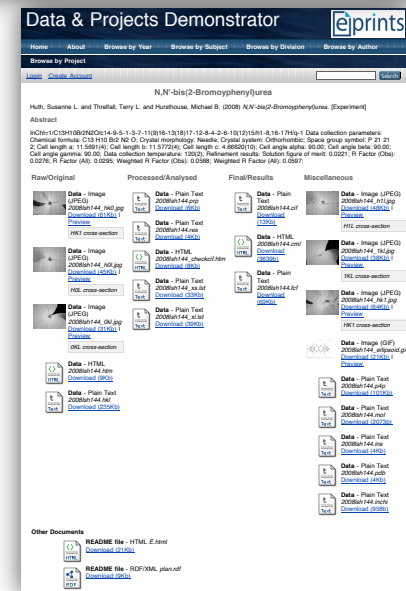
# Data-Oriented EPrint Records

- Dataset
  - A one-off, independent collection of data.
  - Contains one (or possibly more) data documents and probably README information to explain how to interpret the data document.
  - Metadata reflects the collection and processing of the data
    (*e.g.* a range of contributors instead of a creator)

- Experiment
  - A structured process that may result in many datasets, intermediate analyses and final results.
  - Contains various data documents and README information. The experimental methodology/trial protocol/observational process can be described in a normal descriptive text.
  - Metadata reflects the experimental activities and context.

# Data-Carrying EPrint Records

- Publication
  - Any publication-style eprint – journal article, conference paper, report *etc*
  - Publications may contain a data document with associated README information.
  - Metadata reflects the publication, not the data gathering
    - list of authors rather than research technicians



Carbon Nanotubes in a Photonic Metamaterial

Nikolaenko, Andrey E. and De Angelis, Francesco and Boden, Stuart A. and Papasimakis, Nikitas and Ashburn, Peter and Di Fabrizio, Enzo and Zheludev, Nikolay I. (2010) *Carbon Nanotubes in a Photonic Metamaterial.* Physical Review Letters, 104 (15). ISSN 0031-9007

**Abstract**

Hybridization of single-walled carbon nanotubes with plasmonic metamaterials leads to photonic media with an exceptionally strong ultrafast nonlinearity. This behavior is underpinned by strong coupling of the nanotube excitonic response to the weakly radiating Fano-type resonant plasmonic modes that can be tailored by metamaterial design.

**Published Version** - PDF *2010_Nikolay_Carbon_nanotubes_in_a_photonic_metamaterial.pdf*
Download (1597Kb) I Preview

**Data** - Image (PNG) *image10a.png*
Download (1836Kb) I Preview

*Gold metamaterial. Samples from N Zheludev*

**Data** - Image (PNG) *image10b.png*
Download (1886Kb) I Preview

*Carbon nanotube bundles on gold metamaterial. Samples from N Zheludev*

**Additional Metadata** - Microsoft Excel *metadata10.xlsx*
Download (34Kb)

*Standard ZEISS Orion metadata as imprinted on image: first column attribute name, second column attribute value, third column units. The second and third columns are repeated (in the fourth and fifth columns) to refer to the second data image. First row column headings. Attributes are: Field of View / Working Dist / Blanker Current / Image Size / Dwell Time / Mag (4x5 Polaroid) / Date / Time*

# Kinds of Publication Content

- different versions or lifecycle stages
  - draft / submitted / accepted / published / updated / redacted
- associated material that augments the publication
  - supplemental material
- slides or video of a talk about the publication
  - presentation
- an image that visually represents the publication
  - coverimage

# Kinds of Data Content

- ## Data
  - Primary information, experimental results, equipment readings, survey responses

- ## Additional Metadata
  - Subject-specific secondary information about the conditions under which the data was obtained (*e.g.* focal length, beam energy etc)

- ## README
  - Experiment-specific explanation of methodology, process or rules of interpretation of the result data
  - If a separate README document is not justified, a README field is provided for convenience.

**Data** - Archive (ZIP) *transcripts.zip*
Download (94Kb)

*Anonymised transcripts from the self-report experiments*

**Data** - Plain Text *2008lsh144.cif*
Download (13Kb)

**Data** - Image (JPEG) *2008lsh144_h1l.jpg*
Download (48Kb) | Preview

*H1L cross-section*

**Additional Metadata** - Microsoft Excel *metadata01.xlsx*
Download (32Kb)

*Standard ZEISS Orion metadata as imprinted on image: first column attribute name, second column value, third column units. First row column headings. Attributes are: Field of View / Working Dist / Blanker Current / Image Size / Dwell Time / Mag (4x5 Polaroid) / Date / Time*
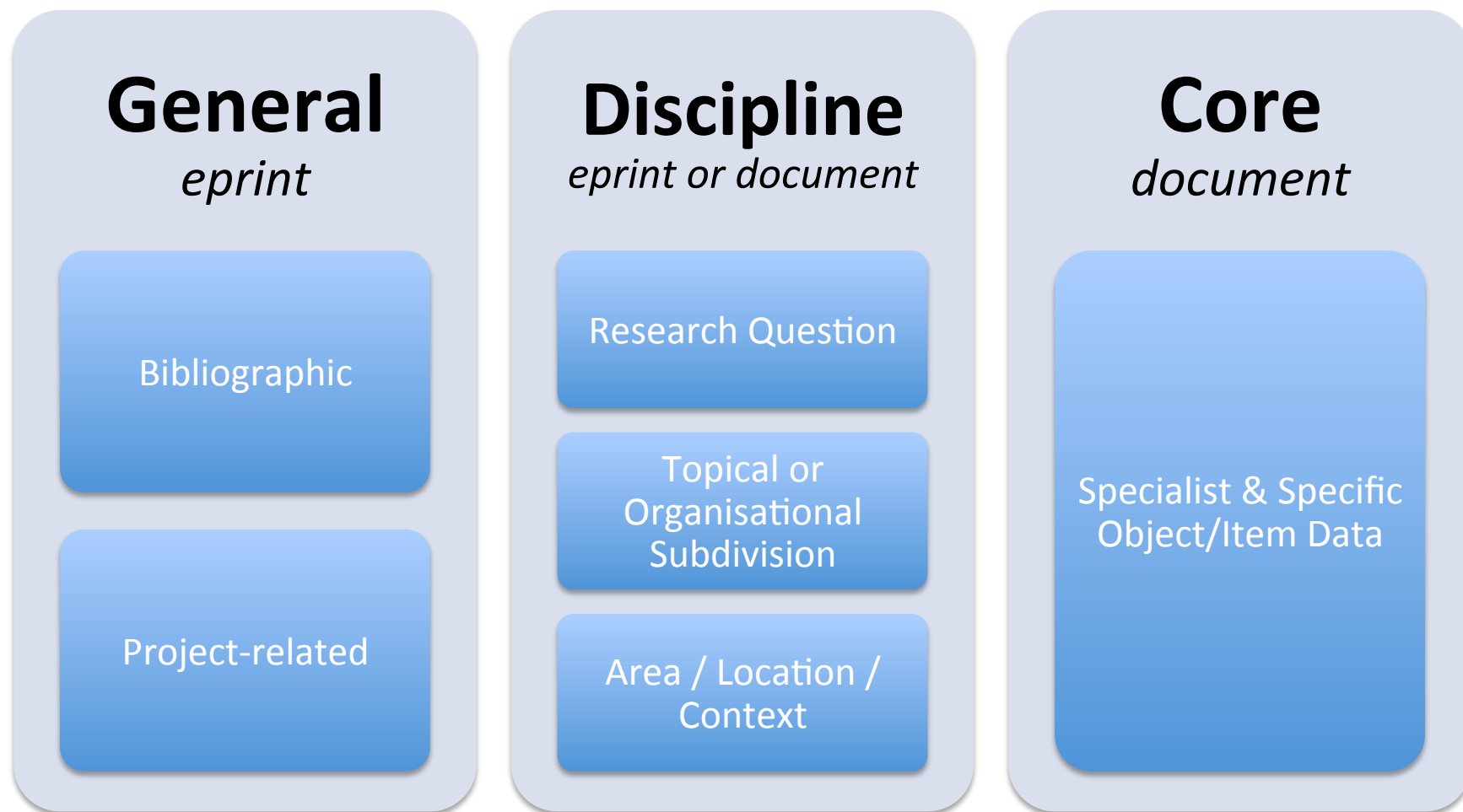
**README file** - HTML *E.html*
Download (21Kb)

**README file** - RDF/XML *plan.rdf*
Download (9Kb)

**README file** - Image (PNG)
Download (43Kb) | Preview

# Institutional Data Management Blueprint 3-layer **metadata** guidance
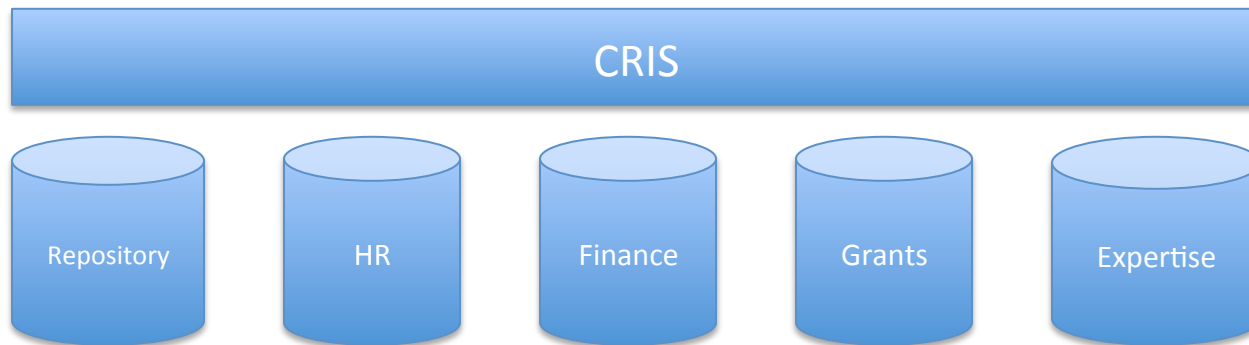
## General
*eprint*

Bibliographic

Project-related

## Discipline
*eprint or document*

Research Question

Topical or Organisational Subdivision

Area / Location / Context

## Core
*document*

Specialist & Specific Object/Item Data

*Research-specific info should be delegated to the document level, or the data files themselves*

# Research Management

- A CRIS (Current Research Information System) pulls together information from all of an institution's research-relevant databases

  - Used for research reporting

- Repositories should support the CERIF standard to co-operate as components of a CRIS environment

| CRIS | | | | |
|------|------|------|------|------|
| Repository | HR | Finance | Grants | Expertise |

# EPrints + CRIS

- EPrints internally accommodates CERIF data
  - Not just *publications* but also *projects* and *organisations*
  - Allows data interchange with external CRIS systems
  - Allows repository to act as a simple CRIS

- CERIFed repositories have many separate datasets, all linked together via explicit relationships



**Formate assay in body fluids: application in methanol poisoning.**

Makar, A B and McMartin, K E and Palese, M and Tephly, T R (1975) *Formate assay in body fluids: application in methanol poisoning.* Biochemical medicine, 13 (2). pp. 117-26. ISSN 0006-2944

PDF - Published Version
Download (1211Kb) | Preview

**Abstract**

A sensitive and specific assay for formic acid in body fluids has been developed. The assay is based on the reaction of formate with bacterial formate dehydrogenase coupled to a diaphorase-catalyzed reduction of the nonfluorescent dye resazurin to the fluorescent substance resorufin. Formate concentrations of 0.5 µg/ml of reaction mixture can be accurately measured. Small volumes of body fluids can be used for the analysis of both methanol and formate. The procedure described is simple and allows for the economical and rapid determination of formate. It can be used in studies concerned with the disposition of formate, as it relates to methanol metabolism. Also, it may be useful in studies where formate might exist as a metabolic intermediate of certain drugs or chemicals.

**Projects**

[118] Performance of Nonlinear Controllers
[428] High Performance and Robust Systems

**Item Type:** Article

---

**Performance of Nonlinear Controllers**

We are concerned with controlling uncertain nonlinear systems via adaptive techniques. We are particularly interested in evaluating the performance of adaptive controllers, and comparing them against eg. robust designs. This has involved developing techniques which allow lower and upper bound estimates to be made of eg. LQ performance. Uniquely in adaptive control theory, we are accounting for the control effort in the cost. Our original focus of attention is in controlling systems containing significant static functional uncertainties (as opposed to the more standard set-up where the uncertainties considered are parametric). The approach considered involves the introduction of function approximators for on-line modelling of the static uncertainties. We have developed a framework for describing the classes of uncertainties for which such controls are valid -- contrasting to the robust theory, uncertainties are measured by spatial L2 weighted norms contrasting to usual static uncertainty models which are formed by pointwise bounds. The interest in performance arose as we tried to quantify which function approximator structures are `best'. This wonderfully ill-posed question is very rich. Currently we have been able to exhibit some structures whose associated LQ performance scales badly as the resolution of the approximator is increased, and also to construct controllers and approximator structures which scale well. Unfortunately, the class of approximator based controllers scale poorly includes some of the standard designs. Our focus of attention is now on using the framework developed for addressing the above question to compare the performances of more classical designs.

| Contributors | Type | Name | ID |
|---|---|---|---|
| | Principal Investigator | French, Mark | maf@ecs.soton.ac.uk |
| | Co-Investigator | Harris, Chris | ch@ecs.soton.ac.uk |
| | Co-Investigator | Rogers, E | ecr@ecs.soton.ac.uk |

**Grant Reference** GR/R27594/01
**Funders** [21] Engineering and Physical Sciences Research Council
**Commencement Date** 01 April 2001
**Completion Date** 31 May 2004
**URI** http://www.isis.ecs.soton.ac.uk/control/projects/adaptive/adaptive.htm
**Id** 118

article record → links to → project record

# Research Data Collection

- Repository as data collection service
- Social scientists view 'the Web' as their next big data challenge
  - and methodological challenge!
- Support researchers for long-term data collection (months / years)
  - forum contributions
  - Google search results
  - Youtube videos
  - Facebook entries
  - Tweets
  - Blog comments

# The challenges of managing social media research data: a researcher's perspective

- *Dr Anne Alexander is Co-ordinator of the Cambridge Digital Humanities Network, a network of researchers at Cambridge who are interested in how the use of digital tools is transforming scholarship in the humanities and social sciences.*

- "But as a researcher, interacting with the material on Facebook brings a huge number of challenges. A particular problem I am struggling with at the moment is how to 'freeze' a dynamic digital environment such as a Facebook wall in order to capture some of the data I am interested in studying. Beyond taking screen shots and saving as a pdf, this is difficult to do, and it is risky to assume that the material you want to look at will still be available in six months' time, let alone years hence."

# Twitter



- Two repository datasets
  - individual tweets
  - timeline
- With rendering and rudimentary analysis
  - to support long-term data collection
- Package downloadable from EPrints Bazaar

500,000 tweets about "Dr Who"
collected over 6 months

# EPrints Bazaar

- 1-click installation

# App *as* Repository? *for* Repository?

- Publication workflows and technologies are similar for all disciplines

- Scientific workflows and technologies are different for every discipline

- Every lab or research group should have its own app that integrates seamlessly with its own practices and deposits automatically into the repository
  - SWORD
  - JISC DepositMO

# Final Thought
# What is a Repository ANYWAY?

- Do all these agendas (open access, research reporting, e-research) just confuse the IR mission statement?

- A repository is not just a piece of information management software

- It is a socially embedded technological phenomenon that promotes new relationship to research information
    - International programs of 'advocacy'
    - Institutionally embedded, with teams of librarians trained to use, and to train researchers to use, repositories
    - Personal engagement with end-users