# Metadata Creation, Transformation and Discovery for Social Science Data Management: The DAMES Project Infrastructure

**Abstract**
This paper discusses the use of metadata, underpinned by DDI (Data Documentation Initiative), to support social science data management. Social science data management refers broadly to the discovery, preparation, and manipulation of social science data for the purposes of research and analysis. Typical tasks include recoding variables within a dataset, and linking data from different sources. A description is given of the DAMES project (Data Management through e-Social Science), a UK project which is building resources and services to support quantitative social science data management activities. DAMES provides generic facilities for performing (and recording) operations on data. Specific resources include support for analysis through micro-simulation, and support for access to specialist data on occupations, educational qualifications, measures of ethnicity and immigration, social care, and mental health.

The DAMES project tools and services can generate, use, transform, and search metadata that describe social science datasets (including microdata from social survey datasets and aggregate-level macrodata). On DAMES, these metadata are described by various standards including DDI Version 2, DDI Version 3, JSDL (Job Submission Definition Language), and the purpose-designed JFDL (Job Flow Definition Language). The paper describes how DAMES uses metadata with a range of resources that are integrated with a job execution infrastructure, a Web portal, and a tool for data fusion.

*Jesse M. Blum, Guy C. Warner, Simon B. Jones, Paul S. Lambert, Alison S. F. Dawson, Koon Leai Larry Tan, Kenneth J. Turner[1]*

## 1. Introduction
This paper discusses the use of metadata within the infrastructural provisions of the UK-based DAMES project (Data Management through e-Social Science, www.dames.org.uk). The remit of this project is described, along with the requirements for data services and the role of metadata in this work.

A review [1] has been conducted of existing social science metadata standards of relevance to the DAMES project. This concludes that DDI (Data Documentation Initiative [12]) is the obvious standard with which to coordinate the use of social science metadata. The authors also observed that the transition from DDI Version 2 to DDI Version 3 supported many additional metadata structures which in principle would assist in the documentation of social science datasets relevant to DAMES. This paper reports on progress in implementing metadata tools within DAMES, outlining prototype services and the means by which they exploit DDI and other metadata formats.

### 1.1 The DAMES Project Approach to Data Management in Social Science
Data management is taken to mean the discovery, preparation, and manipulation of social science data for the purposes of research and analysis. More particularly, this refers to social survey research. Here, data management is typically undertaken by applied researchers themselves (e.g., a secondary survey data analyst), and/or by a data distributor (e.g., a survey collection agency). Typical data management tasks include recoding variables within a dataset, harmonising and standardising variables, and linking data from different sources in order to enhance analysis. Data management tasks are highly significant in social science research. They account for a very substantial proportion of the time spent undertaking research. They have the potential to fundamentally change the conclusions from an analysis. They are also critical to the replicability (or otherwise) of the research workflow (e.g., [6]).

There is nevertheless relatively little methodological literature on the topic of data management (at least, compared to more voluminous methodological material on collection and analysis of social science data). Existing resources include instruction materials on good practices for data management in the context of specific relevant software packages (e.g., [9, 10]). Some Internet sites and publications offer advice on dealing with specific data resources, variables, and measures (e.g., the UK Survey Resources Network, surveynet.ac.uk). Capacity-building activities also offer training and advice (e.g., [11]).

Our observation is that most such materials are ad hoc, in the sense that they deal with specific problems and requirements, and are difficult for non-specialist researchers to engage with. The DAMES philosophy is that progress would be made in social science research if

data management aspects of the research process could be followed more systematically by applied researchers. DAMES is working to make this engagement much easier than is presently common.

DAMES is therefore creating tools and services for management of social survey data. The facilities should be available for general use by applied researchers, and they should support and promote good practice in data management operations. Good practice includes providing clear documentation for data management tasks. For example, reusable traces should be defined for the data preparation commands used in a mainstream software package (e.g., [10]). Good practice also includes making use of previous research efforts in the field. For example, there should be consideration of the ways in which outcome measures have been operationalised in previous studies, e.g., recoding or numeric standardisation, and implementing those existing operationalisations within new studies.

The work on DAMES is contributing online infrastructural services that are accessible to most social science researchers, and that address the challenges of documentation and engagement with data management in previous research. These goals are being pursued through three strategies:

1. Developing Web sites and portal systems (using LifeRay [13]) that offer easy-to-use points of entry into DAMES services

2. Developing online services for researchers to deposit and search heterogeneous Data Management Information Resources (DMIRs)

3. Developing online services that allow researchers to undertake relatively challenging data management tasks (such as matching complex data files) and to document these

Structured metadata records are central to delivering these strategies for two reasons. First, they are essential to providing adequate, searchable records for the complex, heterogeneous data management information resources relevant to (2). DMIRs can take many forms such as quantitative data tables (e.g., aggregate statistics on occupational titles), command files for specific statistical packages, and unformatted notes and documentation. Second, structured metadata records can provide the point of connection between the DMIRs in (2) and the services needed for (1) and (3). The Web site and portal framework can use structured metadata to selectively display suitable information about data resources. The documentation of a data management task itself requires metadata about the component datasets used.

The challenge of generating effective online services in DAMES is increased by the very large volume of data management information resources that are relevant and interesting in social science research. As examples, DAMES is supporting specialist research domains including the study of occupations, educational qualifications, ethnicity and immigration, social care, and mental health inequalities (*www.dames.org.uk/themes. html*).

The DAMES infrastructure thus provides access to a selected range of data management information resources. The principal challenges are:

- To provide (and document) repeatable data management activities

- To provide accessible guidance for using complex, heterogeneous data resources

- To provide accessible guidance for otherwise neglected data management techniques (e.g., standardisation of variables, matching data files)

- To facilitate and encourage data management activities of a higher standard than is currently common

A data curation tool has been developed to collect and maintain metadata about information resources. A data fusion tool has been developed to support data management tasks such as linking data files by exploiting relevant metadata. These are discussed later, focusing on the use of structured metadata.

**1.2 The GEODE Project Approach to Metadata**
The GEODE project (Grid-Enabled Occupational Data Environment, www.geode.stir.ac.uk) was a precursor to DAMES, focusing exclusively on helping social science researchers obtain better access to occupational information resources. The DAMES project is continuing this service, whilst generalising it to support other forms of data management information resources.

GEODE collected a diverse range of occupational information resources, and created metadata descriptions for them. This was done by first asking a data depositor to complete a small online form describing the data. This information was stored in a simple XML metadata format. The information could then be supplemented by a much more substantial set of metadata, in DDI 2.1 format, which was added manually to the pool of occupational metadata. A subset of the DDI 2.1 tags was used to document these records, though more tags could be accommodated. The tags are defined by the GEODE-M schema (www.geode. stir.ac.uk/geode_m_curation.html). The principal structural DDI 2.1 tags are all used (<codebook>, <docDscr>, <stdyDscr>, <fileDscr>, <dataDscr> and <otherMat>), as well as most of their immediate sub-elements (such as

<citation>, <docSrc>, <stdyInfo>, <varGrp> and so on). However, many of the more detailed or application-specific sub-elements were not needed in GEODE -- <topcClas>, <westBL> and <qstn> to name a few. GEODE-M was used to perform occupational matching on micro-datasets. This links values of semantically equivalent variables with the associated occupational information resource resulting in new classification variable values (e.g., for another occupational scheme), thereby helping researchers to prepare their analyses. References [7] and [8] describe the process of adding metadata to occupational resources and the use of metadata within GEODE services.

DAMES is extending the work of GEODE by generating services for several new specialist information resources [14]. DAMES is also developing a number of new resources that feature significant data management requirements. Reference [1] observed that the GEODE metadata model is likely to be too restrictive to apply readily to other domains. A more generalised approach is needed to support the multiple and highly varied domains in DAMES that require a wider range of data management tasks. Automated metadata construction is desirable given the large volume of datasets. It is also desirable to support the curation of more complex data formats, given the wide range of data management information resources.

## 2. Infrastructure Overview

As motivated in the preceding section, the DAMES infrastructure centres upon supporting access to a selected range of data management information and data manipulation resources. Each of these resources is described using searchable metadata for improving resource discovery by social scientists. This section outlines the architecture and operation of the DAMES infrastructure.

To facilitate accessibility to social science researchers, DAMES provides a Grid-based infrastructure accessed through an online portal, with a familiar and comfortable mode of working for the researcher. Through the portal, the researcher is able to use portlets customised for discovering and accessing specialised information resources [14], and for carrying out data management activities such as data curation and data fusion. Figure 1 shows an overview of the DAMES infrastructure.

The infrastructure provides the user-facing portlets with service-oriented access to resources stored at (potentially) physically distributed locations:

- Datasets may be held in *external repositories* or uploaded to the DAMES *file store*, with curation information entered into the metadata store.

- The *metadata store* is provided by an eXist XML database [15]. This is populated with curation

information concerning explicitly uploaded datasets and remote datasets. More importantly, in the context of this paper, the metadata store holds details of datasets resulting from data management activities (identities of input datasets, processing activities carried out, the rationale for the activity). This use of metadata is discussed in more detail in the following sections.

- The *file store* is supported by iRODS [16]. This provides transparent access to distributed file store through an advanced network and server. The file store is populated with uploaded datasets, and also the results of data management activities such as data fusion and standardisation of variables.

- The *search services* allow discovery of curated datasets in both the file store and external locations, through access to the metadata store. This aspect of the DAMES infrastructure is not yet fully developed; it is planned to offer advanced searches based on Semantic Web concepts.

- The *compute resources* would typically be in a Condor pool [17] to which tasks are submitted for carrying out requested data management activities. These requests are expressed in JFDL (Job Flow Definition Language [5]), an extended version of JSDL (Job Submission Definition Language [3]) described in the following sections. The job flows themselves become part of the metadata documenting the resultant datasets.

- *Enact fusion* is one example of the data manipulation services in the DAMES infrastructure. These services receive task descriptions created by the researcher through the portal for searching and task customisation. This results in fetching relevant datasets from local and/or external resources through the *file access* service, submitting jobs to *compute resources* for execution, returning result datasets to the file store, and curating information in the metadata store as appropriate.

The following scenario illustrates common usages of the DAMES infrastructure. A social science researcher wishes to fuse Scottish Household Survey data with privately collected study data on Internet usage. Having registered with the UK Data Archive and obtained the appropriate End User License from them, the researcher uses the data curation and data fusion tools available via the DAMES portal to upload the data and generate a derived dataset. The metadata about this derived dataset is then made public through the portal. Another researcher now searches the portal for data related to the Scottish Household Survey,
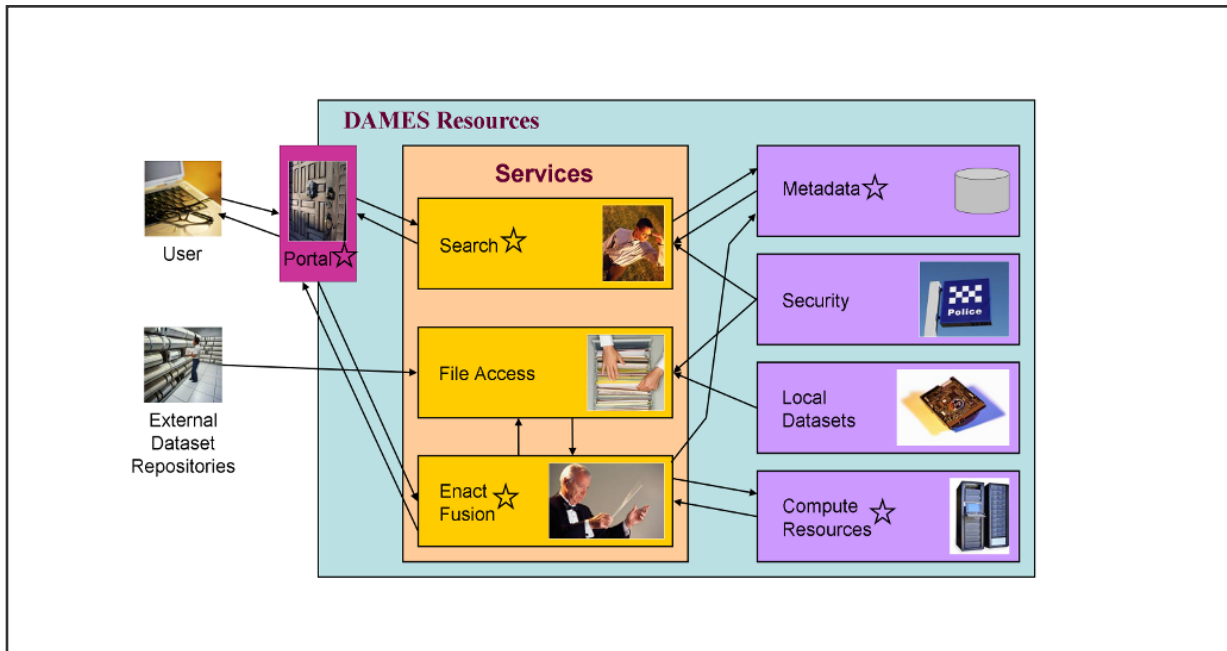
**Figure 1**. The DAMES Infrastructure – key use is made of DDI in the components marked ☆

and finds the metadata file related to this derived dataset. Realising that this would be suitable, the researcher obtains a licence and is then able to access the derived dataset. The derived dataset might then be re-purposed, using the DAMES tools to fuse it with data from a third study.

It is important to note that the DAMES infrastructure provides only tools and support to the researcher, and that it is the responsibility of the researcher to establish the scientific merit of the data fusion.

### 3. Metadata Curation and Generation

The DAMES infrastructure is still in development at the time of writing. An infrastructure has been implemented to support DDI 3. This is also supported by a Grid computing standard called JSDL (Job Submission Description Language [3]).

A number of metadata standards were reviewed in [1], with the conclusion that DDI Version 3 was the most appropriate one to use for curating datasets and

resources, and for describing datasets generated or modified by the DAMES tools. Many resources that are currently available for use in DAMES (and in particular resources available from GEODE [2]) have DDI Version 2 metadata



**Figure 2**. The first stage of data curation using the DAMES infrastructure

records. It is therefore necessary to support this version as well. However, the following focuses on metadata curated and generated using Version 3.

### 3.1 Data Entry and Curation
The DAMES services distinguish private datasets, public datasets, and other resources (such as classification schemes and processing instructions). Users may populate their own private storage space with data and resources and not have them exposed to other users. The DAMES services can process private data, public data, or a combination of the two. All datasets resulting from the use of services are considered private. They are therefore available only to the user that generated the data. Users may propose making data and resources public (uploaded or generated through tools). However, since uploading data and resources to the DAMES infrastructure does not require or generate any metadata, the users must first perform curation.

Curation is handled through a Web portal wizard. It takes users through various stages to generate DDI Version 3 metadata that describe the items being curated. Users can upload existing metadata, which are used to automate much of the curation process. Curated data and resources can be proposed for publication to the DAMES management team. The management team then inspects the items and the metadata, and can make the metadata public or suggest improvements to the users. Figure 2 illustrates the first stage of curating a new data resource.

Metadata records are made public, but not the actual data and resources. This is done to avoid making generated datasets public that are derived from datasets which the user does not have authority to publish. As explained below, the metadata for derived datasets refer to the original datasets and their metadata. The DAMES services access the originals on behalf of the users and are therefore limited to the access permissions of the users. If the users do not have access to the originals, they are limited to inspecting the metadata of derived datasets in order to determine how to obtain appropriate permissions. The possibility of automatically inferring access rights, particularly for derived datasets, is being investigated. In principle it would be possible to export the DDI metadata records, although this facility is currently not available through the DAMES services.

### 3.2 Tool Inputs and Outputs
Whilst DDI metadata are appropriate for describing resources, an instruction language is also needed to describe composable tasks that services can use in a Grid environment. JSDL was chosen owing to its widespread use in existing Grid systems, and the need for DAMES services to run on large-scale existing Grids. However, using JSDL raised two key problems: how to handle job flows, and how to transform JSDL records into meaningful social science metadata.

The first problem with JSDL is that does not support relationships between jobs. This is important because DAMES services require interactions between numerous jobs to complete their processing. For example, a service for fusing two datasets accessed from different external databases might require jobs for staging in each dataset, imputing variables, mapping variables, and fusing the data. Some of these jobs could run in parallel, while others might depend on the results of earlier jobs. The mechanisms offered by Condor [4] were considered, but these are hampered by lack of support for data flows. The result is that data cannot be staged in effectively from various external data sources on behalf of users. It is also hard to specify how to effectively stage out the resulting datasets and metadata.

The solution adopted was to describe service job flows by combining the purpose-designed JFDL (Job Flow Definition Language [5]) with JSDL. JSDL is used to describe service provision issues and data staging. In addition, JFDL is used to describe the relationships between datasets and the transformational tasks needed to realise DAMES services. The JSDL/JFDL fragment in Listing 1 shows job flow and data flow for an example service (omitting the obvious closing tags here).

```
<jsdl:JobDefinition xmlns:jsdl="…" xmlns:jfdl="…">
   <jsdl:JobDescription>
      <jsdl:JobIdentification>
      …
      <jsdl:Application> <jsdl:ApplicationName>DAMES::Fusion …
         <jfdl:JobFlow>
            <jfdl:Dataset file="Gamma1.csv" id="acsv"/>
         …
            <jfdl:Method submitFile="submit.imputeA"
            id="imputeA"/>
            <jfdl:Job id="j1">
               <jfdl:DataSequence>
               …
               <jfdl:Method ref="imputeA"/>
            <jfdl:Job id="j3">
               <jfdl:Parent_job ref="j1"/>
               <jfdl:Parent_job ref="j2"/>
            …
               <jfdl:Method ref="fusion"/>
      <jsdl:DataStaging>
         <jsdl:FileName>Gamma1.csv
         <jsdl:CreationFlag>overwrite
         <jsdl:Source/>
         <jsdl:DataStaging> <jsdl:FileName>output.csv …
         <jsdl:Target>
            <jsdl:URI>irods:/HOME/MyFusion/ad7a835325cd2bd
            d9f549c0271274796/F.csv …
```

***Listing 1.*** *An Example JSDL/JFDL Fragment*

The second problem with using JSDL (and JFDL) to describe DAMES tasks is how to transform these records into meaningful social science metadata. When DAMES services succeed in outputting datasets created from JSDL/JFDL inputs, the JSDL/JFDL instructions are metadata describing how the output datasets were created. However, these records lack descriptive qualities such as the purpose of the investigation, the investigator, and a logical view of the input and output datasets. The records are also not in a standard format for sharing by social scientists.

To solve these problems the DAMES team has developed a solution that uses XSLT (Extensible Stylesheet Language Transformations) to transform JSDL/JFDL instances, along with DDI 3 instances describing user and job profiles, into DDI 3 dataset records. Following translation, the DDI metadata are stored in the metadata database for future searches by researchers. In addition, the JSDL/JFDL file is stored so that the job can be re-run if necessary. The JSDL/JFDL fragment in Listing 1 results in the DDI 3 fragment shown in Listing 2 (again, omitting closing tags here).

The DAMES team has developed a data fusion service to test the infrastructure. Social science researchers can use this service to specify two datasets to be grouped by common variables (imputing data if necessary). A wizard was developed as a Web portlet to generate the JSDL/

JFDL file and submit it to the service. The service stages in the datasets and processing algorithms on behalf of the user, runs the jobs, generates the metadata, and makes the output dataset and metadata available to the researcher. Additionally, researchers can inspect the status of failed jobs to determine the causes of failure.

Figure 3 illustrates the metadata cycle implied by the data fusion service: Searching the metadata discovers datasets for processing. The outcome of processing is curated



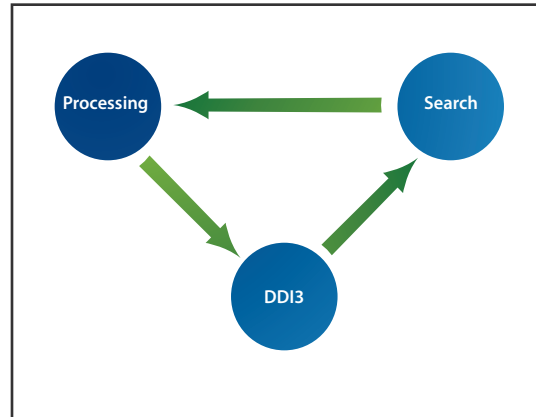**Figure 3.** The DAMES Metadata Cycle

```
<ns1:DDIInstance xmlns:ns1="ddi:instance:3_0_CR3" ...>
  ...
  <g:Group time="T0" instrument="I0" panel="P0" geography="G0"
  dataSet="D3" language="L0">
    <r:Citation>
     ...
    <g:Purpose>...
    ...
    <a:Archive>
        <r:name>DAMES Data Archive ...
    <g:Concepts>
      <!-- A conceptual component describing the data that are
      being fused is given for each mapped and imputed variable -->
      ...
    <g:DataCollection>
      <d:DataCollection>
        <d:ProcessingEvent isIdentifiable="true" id="[Fusion
        Method Id]" isDerived="false">
          <d:Coding isIdentifiable="true" id="[Fusion Method
          Submit File]">
            <d:GenerationInstruction>
              <d:SourceVariable ...>
                <r:URN...
                <d:Mnemonic>Donor Dataset
              <d:SourceVariable ...>
                <r:URN>...
                <d:Mnemonic>Recipient Dataset ...
            <r:Command>
              <r:CommandFileformalLanguage="[Language
```

```
                                      such as SPSS or STATA]">
                ...
                <r:URI>[URI for the command file]
                ...
<g:LogicalProduct>
    …
    <l:CategorySchemeReference ... URI="[URI to category
    scheme for dataset]">
    ...
<g:StudyUnit>
    <!-- Donor dataset -->
    ...
<g:StudyUnit>
    <!-- Recipient dataset -->
    ...
<!-- A comparison is added for each mapped donor & recipient
variable -->
<cm:Comparison>
    ...
    <cm:VariableMap>
      <cm:SourceSchemeReference
      isExternal="true" URI="[URI
      to DonorDatasetVariable]" />
      <cm:TargetSchemeReference
      isExternal="true" URI="[URI to
      RecipientDatasetVariable]" />
    ...
```

*Listing 2. An Example DDI 3 Fragment*

along with its own metadata, ready for re-discovery in subsequent searches.

The data fusion tool outputs datasets described by metadata records that make use of the grouping and processing event description capabilities of DDI 3 (as shown in Listing 2). The metadata files group the corresponding input datasets as referenced study units, along with references to command files containing the processing instructions. They also contain citation, grouping purpose and information about the fusion tool, along with concepts and logical products composed from references to the input datasets' metadata. The principal structural DDI 3 tags currently used by the DAMES infrastructure are <DDIInstance>, <Citation> and <Group>, with detailed metadata held within <Group> elements using <Purpose>, <Concepts>, <DataCollection> <LogicalProduct> and <StudyUnit> tags.

## 4. Discussion and Conclusion

This paper has presented the DAMES infrastructure which is being developed to support the provision of data access and processing services for social science researchers. We have implemented prototype tools based on the infrastructure for resource curation and data fusion. The curation tool generates standardised metadata about social science data and resources to operationalise and fuse the data. The data fusion tool uses existing metadata to perform data linkage tasks and generates metadata for its output datasets. The DAMES infrastructure integrates these tools in a framework that manages depositing, accessing, and processing data. Prototype and early versions of the services are accessible through the DAMES Web site at www.dames.org.uk/resources.html. This paper has shown that DDI 3 features are useful to the curation and fusion services. The paper has also shown that DDI 3 can be integrated with approaches making use of additional metadata schemas purposed for service execution such as the JFDL for job flows description and JSDL for job submission description.

### Notes
1 Jesse M. Blum, Research Assistant, Computing Science and Mathematics, email jmb@cs.stir.ac.uk. Guy C. Warner, Research Assistant, Computing Science and Mathematics, email gcw@cs.stir.ac.uk. Simon B. Jones, Lecturer, Computing Science and Mathematics, email sbj@cs.stir.ac.uk. Paul S. Lambert, Lecturer, Applied Social Science, email paul.lambert@stir.ac.uk. Alison S. F. Dawson, Research Fellow, Applied Social Science, email a.s.f.dawson@stir.ac.uk. Koon Leai Larry Tan, Research Assistant, Computing Science and Mathematics, email klt@cs.stir.ac.uk. Kenneth J. Turner, Professor, Computing Science and Mathematics, email kjt@cs.stir.ac.uk. University of Stirling, Stirling, FK9 4LA, UK

### References
1. J. M. Blum and K. J. Turner. The DAMES Metadata Approach, Technical Report CSM-177, Department of Computing Science and Mathematics, University of Stirling, Dec. 2008, ISSN 1460-9673.

2. P.S. Lambert. An illustrative guide: Using GEODE to link data from SOC-2000 to NS-SEC and other occupation-based social classifications, edition 1.1, GEODE project Technical Paper Number 2, University of Stirling, 2007.

3. A. Anjomshoaa et al. Job Submission Description Language (JSDL) Specification Version 1.0, Global Grid Forum, Nov. 2005

4. D. Thain, T. Tannenbaum and M. Livny. Distributed Computing in Practice: The Condor Experience, *Concurrency and Computation: Practice and Experience*, 17(2–4):323–356, Feb.–Apr. 2005.

5. J. Blum and G. Warner. Job Flow Definition Language (JFDL 1.0), available from *www.cs.stir.ac.uk/~jmb/dames/schemas*, accessed Oct. 2009.

6. A. Dale. Quality Issues with Survey Research, *Int. J. of Social Research Methodology*, 9(2):143–158, 2006.

7. P. S. Lambert and K. L. L. Tan. Instructions for using the GEODE Portal, Edition 1.1, GEODE Project Technical Paper No. 1, University of Stirling, available from *www.geode.stir.ac.uk*, 2007.

8. P. S. Lambert, K. L. L. Tan, K. J. Turner, V. Gayle, K. Prandy and R. O. Sinnott. Data Curation Standards and Social Science Occupational Information Resources, *Int. J. of Digital Curation*, 2(1), 73-91, 2007.

9. R. Levesque and SPSS Inc. Programming and Data Management for SPSS Statistics 17.0, SPSS Inc., Chicago, 2008.

10. J. S. Long. The Workflow of Data Analysis Using Stata, CRC Press, Boca Raton, 2009.

11. V. van den Eynden, L. Corti, M. Woollard and L. Bishop. Managing and sharing data: A best practice guide for researchers, UK Data Archive, Colchester, 2009.

12. M. Vardigan, P. Heus and W. Thomas. Data Documentation Initiative: Towards a standard for the social sciences, *Int. J. of Digital Curation*, 3(1):107–113, 2008.

13. Liferay. Open source enterprise portal with Web content management system, collaboration and social networking – Liferay, available from *www.liferay.com/web/guest/ products/portal*, accessed Oct. 2009.

14. P. S. Lambert, V. Gayle, K. L. L. Tan, J. M. Blum, A. Bowes, S. Jones, K. J. Turner, G. Warner, R. O. Sinnott and E. Bihagen. Grid enabled specialist data environments: Forward planning for the GE*DE services for specialist data on occupations, educational qualifications, and ethnicity, Technical Paper 2008-1, University of Stirling, available from *www.dames.org.uk*, 2008.

15. eXist. An open source database management system built using XML technology, available from *www.exist-db. org*, accessed Oct. 2009.

16. iRODS. Data grids, digital libraries, persistent archives, and real-time data systems, available from *www.irods.org*, accessed Oct. 2009.

17. Condor. High throughput computing, University of Wisconsin, available from *www.cs.wisc.edu/condor*, accessed Oct. 2009.