



Peer Reviewed

Title:

Curating Scientific Research Data for the Long Term: A Preservation Analysis Method in Context

Author:

[Conway, Esther](#), Science & Technology Facilities Council, Rutherford Appleton Laboratory
[Dunckley, Matthew](#), Science & Technology Facilities Council, Rutherford Appleton Laboratory
[Giaretta, David](#), Science & Technology Facilities Council, Rutherford Appleton Laboratory

Publication Date:

10-05-2009

Series:

[iPRES 2009: the Sixth International Conference on Preservation of Digital Objects](#)

Publication Info:

iPRES 2009: the Sixth International Conference on Preservation of Digital Objects, California Digital Library, UC Office of the President

Permalink:

<http://www.escholarship.org/uc/item/14h35961>

Multimedia URL:

<http://www.cdlib.org/services/uc3/iPres/video.html?file=ipres/Conway2&title=Esther%20Conway%3A%20Curating%20Scientific%20Research%20Data%20for%20the%20Long%20Term%3A%20A%20Preservation%20Analysis%20Method%20in%20Context>

Abstract:

The challenge of digital preservation of scientific data lies in the need to preserve not only the dataset itself but also the ability it has to deliver knowledge to a future user community. A true scientific research asset allows future users to reanalyze the data within new contexts. Thus, in order to carry out meaningful preservation we need to ensure that future users are equipped with the necessary information to re-use the data. This paper presents an overview of a preservation analysis methodology which was developed in response to that need on the CASPAR and Digital Curation Centre SCARP projects. We intend to place it in relation to other digital preservation practices discussing how they can interact to provide archives caring for scientific data sets with the full arsenal of tools and techniques necessary to rise to this challenge.

Supporting material:

Presentation

Copyright Information:

iPRES 2009

THE SIXTH INTERNATIONAL CONFERENCE ON THE PRESERVATION OF DIGITAL OBJECTS

Proceedings

October 5-6, 2009
Mission Bay Conference Center
San Francisco, California



California Digital Library

Curating Scientific Research Data for the Long Term: A Preservation Analysis Method in Context

Esther Conway, Matthew Dunckley and David Giaretta

Science and Technology Facilities Council

Rutherford Appleton Laboratory, Chilton, Didcot, Oxfordshire OX11 0QX esther.conway@stfc.ac.uk

Abstract

The challenge of digital preservation of scientific data lies in the need to preserve not only the dataset itself but also the ability it has to deliver knowledge to a future user community. A true scientific research asset allows future users to reanalyze the data within new contexts. Thus, in order to carry out meaningful preservation we need to ensure that future users are equipped with the necessary information to re-use the data. This paper presents an overview of a preservation analysis methodology which was developed in response to that need on the CASPAR and Digital Curation Centre SCARP projects. We intend to place it in relation to other digital preservation practices discussing how they can interact to provide archives caring for scientific data sets with the full arsenal of tools and techniques necessary to rise to this challenge.

Introduction

This paper presents a brief overview of the preservation analysis methodology which was developed on the CASPAR [1] and Digital Curation Centre SCARP projects [2]. After describing the main stages and purpose of the method we intend to place it in relation to other digital preservation and existing archival practices discussing how they can interact to provide archives caring for scientific data sets with the full arsenal of tools and techniques necessary to rise to this challenge. We then intend to illustrate the benefits of preservation analysis with examples from the SCARP case studies and solutions we have implemented on the CASPAR project with a focus on the following key areas.

Maximizing return on investment

Good preservation analysis is essential in order to design a truly reusable asset. This methodology capitalizes on a community's expertise and knowledge by appreciating the nature of data use, evolution and organizational environment. When scientific data is used by a community it develops a history. During its lifetime, the custody of a data set may pass through several bodies generating rich documentation which explains the scientific purpose of the dataset and how it has evolved over time. Organizations also develop around branches of science publishing or producing grey materials over time

which prove to be important for the interpretation or analysis of data.

This method seeks to design the optimal research asset by capturing key information which facilitates reuse. We aim to demonstrate with comparative examples from the DCC SCARP project case studies how judicious analysis permits the design of Archival Information Packages (AIP) which deliver a greater return of investment by both improving the probability of the data being reused and potential outcome of that reuse.

Provision of measureable and testable solutions

The methodology incorporates a number of analysis stages into an overall process capable of producing an actionable preservation plan for scientific data, which satisfies a well defined preservation objective. In this paper we will discuss how the creation of an archival information package using preservation actions selected on this basis ensures a measurable and testable solution. Using an example from the SCARP Atmospheric Sciences case study [3] we will show how we developed such an objective and how this facilitates a testable solution giving an archive the necessary assurance in any preservation action taken.

Management of research assets through the modeling of preservation networks

The analysis method facilitates modeling of information networks based on the archival information package solution. Using illustrative examples from the CASPAR testbeds we intend to show how these network models are a representation of the digital objects, operations and relationships which allow a preservation objective to be met for a future designated community. The models provide a sharable, stable and organized structure for digital objects and their associated requirements. They expose the risks, dependencies and tolerances within an archival information package. This allows for the automation of event driven or the periodic review of archival holdings by knowledge management technologies.

The clear definition of relationships also facilitates the identification of reusable solutions which can be deposited within registry repositories of representation

information, thus sharing preservation efforts within and across communities.

Informing preservation activities in the wider institutional environment

Preservation analysis clearly exposes the requirements and issues associated with an individual data set thereby supplying audit, certification and repository planning activities with essential information. Archival Information Packages are likely to be held by or transferred between institutional repositories. These repositories need to be designed planned and managed, competing for resources within complex organizational structures. In this paper we intend to conclude by touching upon how preservation analysis can inform audit and certification processes such as DRAMBORA [4] and TRAC [5] or planning activities for new repositories such as PLATTER [6]. By doing this we allow preservation analysis at the data set level to be placed within the context of institutional planning and operations.

Overview of the Preservation Analysis Method

The challenge of digital preserving scientific data lies in the need to preserve not only the dataset itself but also the ability it has to deliver knowledge to a future user community. This entails allowing future users to reanalyze the data within new contexts. Thus, in order to carry out meaningful preservation we need to ensure that future users are equipped with the necessary information to re-use the data.

The Digital Curation Centre SCARP and CASPAR projects have a strong focus on the preservation and curation requirements for scientific data sets. These projects engaged with a number of archives based at the STFC [7] Rutherford Appleton Laboratory. In particular we carried out extensive analysis work to consider the preservation requirements of the British Atmospheric Data Centre [8], the World Data Centre [9] and the European Incoherent Scatter Scientific Association (EISCAT) [10]. During these studies it became clear that there was a need for a consistent preservation analysis methodology.

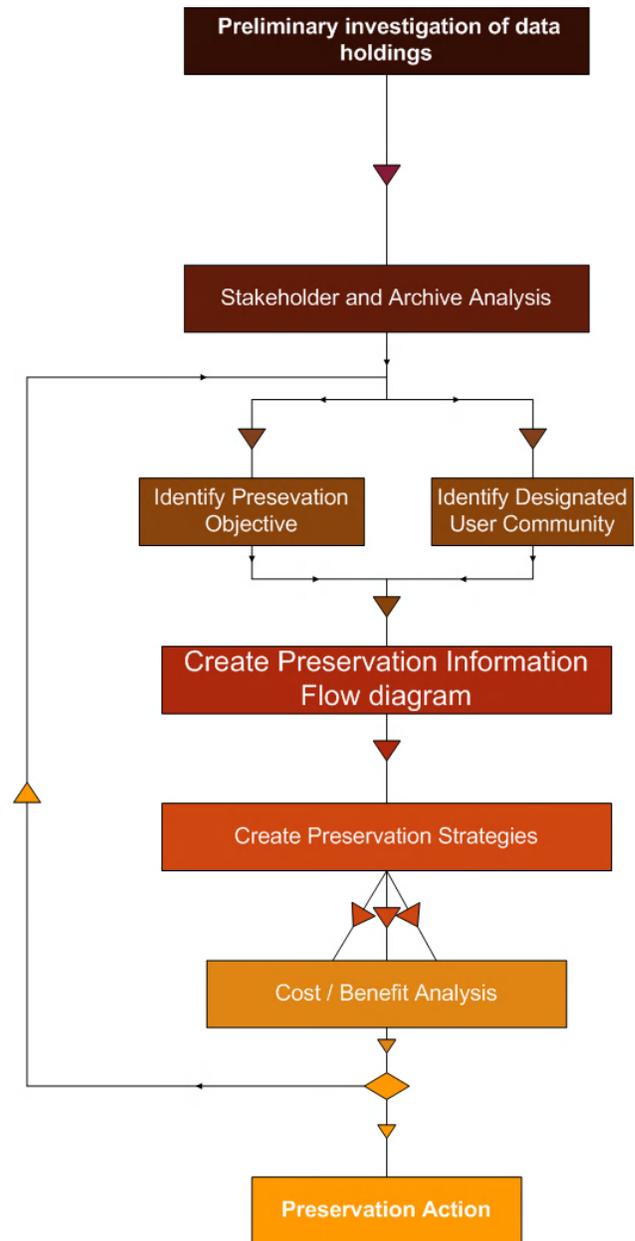


Figure 1. Preservation Analysis Workflow

In the resulting methodology we sought to incorporate a number of analysis techniques, tools and methods into an overall process capable of producing an actionable preservation plan for scientific data archives. Figure 1 illustrates the stages of this methodology. In the rest of this paper we shall discuss the stages in detail, illustrated with examples of work with the scientific archives.

Preliminary investigation of data holdings

The initial step is to undertake a preliminary investigation of the data holdings of the target archive.

The CASPAR project developed a questionnaire [11] containing key questions which allowed the preservation analyst to initiate discussion with the archive. It critically allows the analyst to.

- Understand the information extracted by users from data
- Identify Preservation Description and Representation information
- Develop a clearer understanding of the data and what is necessary for its effective re-use
- Understand relationships between data files and what constitutes a digital object within the archive

While it is appreciated that this questionnaire is not an exhaustive list of questions which one may need to ask about a preservation target, it still provides sufficient information to commence the analysis process. The full questionnaire and results from the Ionosonde WDC holdings [12] can be obtained from the CASPAR website.

Stakeholder and Archive Analysis

After carrying out the questionnaire process for each an archive it is necessary to carry out a stakeholder analysis. This is because

- Stakeholders may hold different views of the knowledge a data set was capable of providing an end user
- Stakeholders can identify different end users whose skill sets and knowledge base vary
- Stakeholders may have produced or be custodians of information vital for re-use of data

The stakeholder analysis classifies stakeholders into a number of categories each with their own concerns. In addition to identifying the stakeholders from the different categories it is also beneficial to understand how an archive has evolved and been managed. This process can be used to illuminate the different uses of data over time and highlight the existence of associated Representation Information.

Defining a preservation objective

The analysis carried out before this point may present one with a natural easily defined preservation objective or alternatively there may be a greater number of options which overlap and are more difficult to define. It is important to note that this type of analysis cannot advise you as to which preservation option to choose but merely clarifies the options available to you. Preservation objectives should be

- Specific well defined and clear to anyone with a basic knowledge of the domain
- Actionable the objective should be currently achievable.
- Measurable it is critical to be able to know when the objective has been attained in order to assess if any preservation strategy developed is adequate.
- Realistic based on findings from the previous stages of analysis

Defining a designated user community

The Designated Community is defined in OAIS [13] as “An identified group of potential Consumers who should be able to understand a particular set of information. The Designated Community may be composed of multiple user communities”

An archive defines the Designated Community for which it is guaranteeing to preserve some digitally encoded information and must therefore create AIPs with appropriate Representation Information.

The designated community will possess a skills and knowledge base which allow them to successfully interact with a set of information stored within an AIP in order to extract required knowledge or recreate the required performance or behavior. In common with the preservation objective the analysis up to this point may present one with a range of community groups which the archive may chose serve.

The definition of the skill set is vital as it limits to the amount of information which must necessarily be contained within an AIP in order to satisfy a preservation objective. In order to do this the definition of the designated community must be

- Clear with sufficient detail to permit meaningful decisions to made regarding information requirements for effective re-use of the data.
- Realistic and stable in so far as there is reasonable confidence in the persistence of the knowledge base and skill set.

Preservation information flows

Once the objective and community have been identified and described an analyst should be in position to determine the information required to achieve an objective for this community. An analyst proceeds by identifying risks which are to be addressed by preservation action. We advocate the creation of an OAIS preservation information flow diagram at this juncture.

An OAIS preservation information flow diagram is graphical representation and analysis tool which is a hybrid of an information flow diagram and the OAIS information model. It provides a convenient format to

facilitate group discussions over preservation plans and strategies. A preservation information flow diagram we created for the MST data is shown below:

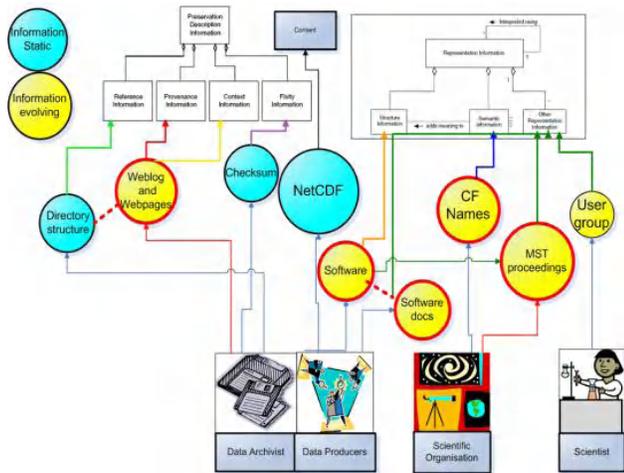


Figure 2. OAIS information flow diagram for the MST data set

The OAIS reference model specifies that within an archival system, a data item has a number of different information items associated with it, each performing a different role in the preservation process. The preservation objective for a designated community should be satisfied when each item of the OAIS information model has been adequately populated with sufficient information. The information model provides a checklist which ensures that the preservation objective can be met. All information objects must be mapped to at least one of the element of the OAIS information model.

In addition to information objects and the standard OAIS information model the diagram contains a number of other components.

As multiple strategies can be developed a number of competing preservation plans are available. A preservation plan should consist of a unique

- Set of information objects
- Set of supply relationships
- Set of preservation strategies

Each plan will allow an archive to carry out a series of clear preservation actions in order to create an AIP. The archive should now be in a position to take a number of plans to the cost/benefit/risk analysis stage where they can be evaluated and a preferred option chosen.

Cost/Benefit/Risk Analysis

The final stage of the workflow is where plan options can then be assessed according to

- Costs to the archive directly as well as the resources knowledge and time of archive staff
- Benefits to future users which ease and facilitate re-use of data
- Risks inherent to the preservation strategies and accepted impact to the archive.

Once this analysis is complete the optimal plan can be selected and progressed to preservation action. If no plans are deemed suitable then the process must begin again with an adjustment to the preservation objective and/or the designated community to be served.

Maximizing return on investment

When we examined different archives we discovered there was the potential to create different research assets for data depending and how one might choose to support reuse of a data set. We observed the following sorts of factors influencing the use and re-use of data over time:

- Birth and development of a science
- Events which influence data use such as the second world war or global warming
- Development of technologies and the emergence of global networks
- Publication of journals technical manuals, interpretative handbooks, conference proceeding, minutes of user group meetings, software etc.
- Emergence of branches of science and associated organisations
- Stewardship of data and the influence of different custodians

This is not an exhaustive list as many factors influencing data re-use are domain specific as is the categorization of the stakeholders. Naturally most of these can only be expected to be dealt with in the most cursory way in any practical study nevertheless even this can be extremely important in understanding the situation. As after this evaluation you should be in position to scope what types of reuse may be realistically achieved.

By comparing two data sets we can demonstrate how an appreciation of a data sets previous use and stakeholder relationships can inform the design of an AIP allowing an archivist to select the optimal set of information to realize return through improving the probability of the data being reused and the outcome of that use.

The Mesosphere Troposphere Stratosphere (MST) [14] data set is extremely well documented and tightly managed. Access to the data is restricted, with end users required to report back on how they have used the data. The Archivist is the key manager of these data for a number of reasons

- He is the project scientist involved in production of the data
- He is a field expert and practising scientist in close contact with relevant scientific organisations
- He provides support, runs and keeps records of user group meetings.

When we consider these factors we can see that it is reasonable to try to capture information from current users which facilitate the re-use of data by future scientists. The information has been captured in user group minutes, conference proceedings and scientific papers resulting from the study of MST data. The capture and determination of value of these resources is possible as a result of the archivist's domain knowledge and close connection to users.

We now contrast the scenario above with that of the ionosonde data holdings of the World Data Centre. Whilst being a skilled individual with some domain knowledge he does not have the same strong connection with users. The data currently comes from 252 geographically diverse locations and current users are simply required to provide an e-mail address to gain access. As a result it would be completely impractical to capture user generated information even if it might facilitate re-use. We make the judgment that the archive would realize a return on the large investment required to source capture and assess information from these user groups.

The added value of information which contributes to reuse must be assessed against the likely cost of its capture and maintenance. If creation of asset reliant on such information is deemed viable and cost effective an archive may then proceed by developing a preservation objectives and which accommodates the information in scope.

Measurable and Testable Solutions

We shall now examine how a preservation objective developed for the MST data set in accordance with the methodology produces a testable solution.

The importance of the MST dataset lies in the fact that it is an irreplaceable earth observational record: once lost, these data cannot be replaced by the repetition of the experiment. The dataset is valuable because it

- Contains data from the UK's most powerful and versatile wind-profiling instrument
- Provides information about atmospheric stability, turbulence, humidity fields, precipitation and a variety of atmospheric phenomena
- Contains measurements of winds up to many kilometres from the ground

- Contains a record of winds sampled continuously, with a cycle time of a few minutes over a long period of time
- Provides a record not only of the horizontal but also the vertical air velocity which additionally has high temporal and spatial resolution.

Give the challenge of digital preservation lies in the need to preserve not only the dataset itself but also the ability it has to deliver knowledge to a future user community. In order to carry out meaningful preservation we need to ensure that future users are equipped with the necessary Information to re-use the data

Use of atmospheric data may involve format conversion to make it interoperable with other data sets. Specialist visualisation of time related data may be needed to study some kinds of atmospheric behaviour. Interpolation, subsetting or different forms of statistical analysis may also be employed. Identification of atmospheric phenomena and behaviours is achieved through the creation or application of data to established models of dynamic atmospheric systems. The same data may be used to study different aspects of atmospheric behaviour; listed below are some which we identified in peer reviewed literature resulting from studies which employed the MST data.

- **Precipitation.** Clouds contain moisture; when the droplets in clouds coalesce they become sufficiently large to cause precipitation. A recent evolution of knowledge surrounding the MST radar data allows the data to be used to study precipitation
- **Convection.** Convection is the transfer of heat by movement within a substance. The MST radar data permits you to study the convective circulation of air within the atmosphere.
- **Gravity Waves.** Gravity waves are generated in the troposphere by frontal systems or by airflow over mountains. The geographic position of the MST radar site is ideal for studying this phenomenon.
- **Rossby Waves.** Rossby waves are a subset of inertial waves. These atmospheric waves are large scale motions with wavelengths of up to 6000 km. The continual monitoring of a discrete region of the atmosphere over a long period allows for the analysis of such waves.
- **Mesoscale and Microscale Structures.** The frequency of observation and the resolution of the MST radar also permit analysis to be carried on mesoscale (~50km) and microscale (atmospheric lasting a matter of minutes) structures.

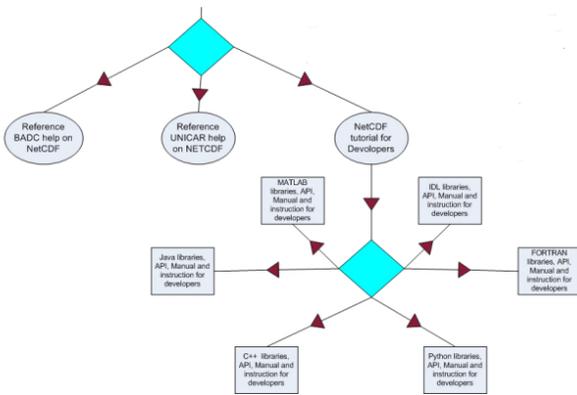


Figure 4 Preservation network model of a NetCDF reusable solution

If we look at the network section above, which performs the function of allowing a user to extract the desired parameters from NetCDF formatted files. There are eight different strategies a user can employ all of which must fail before there is a critical failure of the solution. As this section of the network has a specific well defined function which is to allow a user to extract parameters form NetCDF formatted files, the solution can be deposited within repository of representation information such as RRORI [15]. It can then be reused as part of wider solution for different atmospheric data sets which utilize the format.

Informing preservation activities in the wider institutional environment

Preservation analysis clearly exposes the requirements and issues associated with an individual data set thereby supplying audit, certification and repository planning activities with vital information. We intend to conclude by touching upon how preservation analysis can inform audit and certification processes such as DRAMBORA [4] and TRAC [5] or planning activities for new repositories such as PLATTER [6]. By doing this we allow preservation analysis at the data set level to be placed within the context of institutional planning and operations.

DRAMBORA

DRAMBORA describes digital curation as being characterized as a risk-management activity; where the job of digital curator is to rationalize the uncertainties and threats that inhibit efforts to maintain digital object authenticity and understandability, transforming them into manageable risks. Six stages are implicit within the process. Initial stages require auditors to develop an organizational profile, describing and documenting the repository's mandate, objectives, activities and assets.

We view preservation analysis as being a complimentary activity as it exposes the risks and objectives of a research assets created via the preservation analysis process. Drambora can then proceed to assess these exposed risks in terms of their likelihood and potential impact taking advantage of the clarity provided by preservation network modeling. Auditors are encouraged to conceive of appropriate risk management responses to the identified risk. The DRAMBORA process can then enable effective resource allocation, enabling repository administrators to identify and categorize the areas where shortcomings are most evident or have the greatest potential for disruption.

TRAC

The TRAC checklist is an auditing tool used to assess the reliability, commitment and readiness of institutions to assume long-term preservation responsibilities. The repository checklist is under the care of the Center for Research Library's [16] who are utilizing it in several independent projects. Preservation analysis provides evidence that a repository meets audit and certification criteria. The checklist is divided in the three sections

- A. Organizational infrastructure
- B. Digital object management
- C. Technologies, technical infrastructure and security

Below are some example of how preservation analysis provide this evidence:

- By establishing duties the repository needs to perform
- By identifying skills that staff require proving an adequate professional development program is in place
- Providing a definition of a repositories designated community
- Through preservation network modeling supporting periodic review
- By preservation analysis providing transparency.
- Repository will have identified properties it needs to preserve for digital objects.
- The repository will have clearly specified the information that needs to have associated with digital material at the time of deposit.
- The repository will have documented preservation strategies it has employed.

PLATTER

The Planning Tool for Trusted Repositories (PLATTER) provides a basis for a digital repository to plan the development of its goals, objectives and

performance targets over the course of its lifetime. PLATTER is designed to be complimentary to existing audit and certification tools.

The process is centered on a group of strategic objective plans. Preservation analysis can again inform this process by providing information on data complexity, data specialization, acquisition through clearly identifying that which needs to be covered in a deposit agreement, specialization, the technical by having highlighted the technical risks to the digital object within an AIP and of course the preservation plan

Conclusions

This paper presented an overview of a preservation analysis methodology which was developed on the CASPAR and DCC SCARP projects. Placing it in relation to other digital preservation practices discussing how they can interact to provide archives caring for scientific data sets with the full arsenal of tools and techniques necessary to rise to this challenge of preserving a long term research asset. We have shown how the use of preservation analysis can provide greater return on investment, measurable solutions, assist the management of research assets and support audit/certification and repository planning activities.

Wider application, trialing and further development of the preservation analysis methodology outlined here would be desirable to test its validity in a broader range of disciplines and organisational settings. In addition the production of training materials and support for archivists who wish to adopt our approach for data preservation would be of benefit.

Archives can find it difficult to articulate and specify reasons for the preservation of data. We additionally recommend that the organisations such as DCC develop further guidance on setting preservation objectives and establishing valid business cases for the preservation of scientific data.

Acknowledgements

Work partially supported by European Community under the Information Society Technologies (IST) program of the 6th FP for RTD - project CASPAR and the Joint Information Systems Committee (JISC) for the Digital Curation Centre SCARP project.

We would also like to thank our colleagues at STFC David Hooper, Sam Pepler, Matthew Wild, Steve Crothers, Chris Davis, Rita Blake, Ruth Bamford, Simon Lambert, Stephen Rankin and Brian McIlwrath.

References

- [1] CASPAR Project <http://www.casparpreserves.eu/>
- [2] Digital Curation Centre SCARP project <http://www.dcc.ac.uk/scarp/>
- [3] Conway, E.; Curating Atmospheric Data for long term use: Infrastructure and Preservation Issues for the Atmospheric Sciences community, 2 June 2009 http://www.dcc.ac.uk/docs/publications/case-studies/SCARP_B4832_Atmospheric.pdf
- [4] Drambora <http://www.repositoryaudit.eu/>
- [5] Trustworthy Repositories Audit & Certification (TRAC): Criteria and Checklist; February 2007 <http://www.crl.edu/PDF/trac.pdf>
- [6] PLATTER <http://www.digitalpreservationeurope.eu/platter/>
- [7] Science and Technology Facilities Council <http://www.stfc.ac.uk/>
- [8] British Atmospheric Data Centre <http://badc.nerc.ac.uk/home/index.html>
- [9] World Data Centre for Solar Terrestrial Physics www.ukssdc.ac.uk/wdcc1/wdc_menu.html
- [10] European Incoherent Scatter Radar <http://www.eiscat.rl.ac.uk>
- [11] CASPAR Questionnaire <http://www.casparpreserves.eu/Members/cclrc/Reference/Documents/caspar-test-case-questionnaire/>
- [12] CASPAR Ionosonde case study; March 31 2007 <http://www.casparpreserves.eu/other-caspar-products/other-caspar-products/ionosonde-case-study.pdf>
- [13] ISO 2002. Reference Model for an Open Archival Information System (OAIS). *Recommendation for Space Data Systems Standard, CCSDS Blue Book.* <http://public.ccsds.org/publications/archive/650x0b1.pdf>
- [14] The Natural Environment Research Council (NERC) Mesosphere-Stratosphere-Troposphere (MST) Radar at Aberystwyth <http://mst.nerc.ac.uk/>
- [15] CASPAR/DCC Representation Information Registry <http://registry.dcc.ac.uk:8080/RegistryWeb/Registry/>
- [16] The Center for Research Libraries <http://www.crl.edu/>