



Contents lists available at [SciVerse ScienceDirect](#)

Fusion Engineering and Design

journal homepage: www.elsevier.com/locate/fusengdes



Long term preservation of scientific data: Lessons from jet and other domains

R. Layne^{a,b,*}, A. Capel^a, N. Cook^{a,b}, M. Wheatley^a, JET EFDA Contributors¹
JET-EFDA, Culham Science Centre, OX14 3DB, Abingdon, UK

^a EURATOM-CCFE Fusion Association, Culham Science Centre, Abingdon, Oxon OX14 3DB, UK

^b Tessella plc, 26 The Quadrant, Abingdon Science Park, Abingdon, Oxon OX14 3YS, UK

H I G H L I G H T S

- ▶ Data from a fusion experiment must be preserved for decades, and will be required beyond the lifetime of the experiment.
- ▶ General principles of long-term data preservation have been identified in domains outside fusion.
- ▶ Key challenges are media obsolescence, data organisation, and file format obsolescence.
- ▶ Data from JET has to date been preserved for almost three decades: the steps taken to achieve this are discussed.
- ▶ Recommendations for new devices in planning for the long-term preservation of their data are made.

A R T I C L E I N F O

Article history:
Available online xxx

Keywords:
Data
Preservation
Curation

A B S T R A C T

All fusion devices generate large amounts of data, which must remain accessible throughout and beyond the lifetime of the device. This need for long-term data preservation creates major software and hardware challenges. Issues surrounding long-term data preservation will be discussed, and the main challenges in data curation will be outlined. The paper will then review the approaches taken at JET to ensure accessibility of all of JET's data over almost three decades. Using the JET experience and the wider context, the paper will then suggest how new experiments should plan for the long-term preservation of their data.

© 2012 EURATOM Published by Elsevier B.V. All rights reserved.

1. Introduction

All fusion devices generate large amounts of data: this includes design data, scientific data, engineering data, meta-data describing the data, and information about the process of data production. These devices all have a long but finite lifespan.

The lifecycle of a fusion device starts with design and construction of the device. The device will operate for years – potentially many decades. After this, the device will be decommissioned, while a period of post-operation analysis of the data begins. Data may be required for further analysis for follow-on devices.

The relationship between JET and ITER is an example of this cycle: JET is a mature device, which will reach the end of its

mission at some point in the coming decades, while ITER is still under construction. At some point, then, each fusion experiment will end: but the data generated by the experiment will still be needed by future researchers.

A key question which should concern those responsible for data management for fusion devices is: *how do we ensure our data is preserved for the decades the experiment is in operation, and potentially for decades beyond that?*

This paper will address this issue by examining the wider context and how this applies to the experience of JET.

Firstly, some general principles of long-term data curation will be outlined. The problem does not just apply to the fusion domain, and there are some common issues to be addressed by any organisation concerned with data preservation.

Secondly, JET's experience of data preservation will be outlined, with a particular focus on scientific data. The first JET pulse was in 1983, and all data back to that date is still available. The steps taken at JET to achieve this, and the plans for continued work in this area, will be outlined.

Finally, the general principles and JET's experience will be used to make recommendations for future devices in planning for the long-term preservation of their data, and selecting an appropriate

* Corresponding author at: EURATOM-CCFE Fusion Association, Culham Science Centre, Abingdon, Oxon OX14 3DB, UK. Tel.: +44 01235 465021; fax: +44 01235 464404.

E-mail address: richard.layne@ccfe.ac.uk (R. Layne).

¹ See the Appendix of F. Romanelli et al., Proceedings of the 23rd IAEA Fusion Energy Conference 2010, Daejeon, Republic of Korea.

file format for scientific data. These recommendations are deliberately technology neutral: it is not the aim of this paper to propose a single technology for scientific data storage.

2. General principles of data preservation

Long-term preservation of digital data is not just a fusion problem. Over the past decades there has been growing awareness of the issue [1], and research and development has been ongoing. Organisations in many research fields, and “memory institutions” such as national archives and libraries, are developing standards and technologies in this field.

There are many examples of problem caused by a lack of attention to data preservation which led to this work. One key example was NASA’s Viking Lander data [2]. Two landers were sent to Mars in 1975: datasets were compiled by scientists based on the collected data. The resulting data was stored on magnetic tape, in climate controlled conditions, but despite this, the physical tapes deteriorated. In addition, by the late 1990s, scientists were unable to decode the data format. This data was ultimately recovered by re-entering it from microfilms and printouts.

This mission was highly expensive and difficult or impossible to repeat: the same can be said of fusion experiments. The timescales are similar too: 1975 was eight years before the first JET pulse.

Experts in data preservation have identified three key challenges in the field: media obsolescence, distributed or disjointed data organisation, and file format obsolescence.

2.1. Media obsolescence

All storage media have a finite lifespan, and can be expected to need replacement during the lifetime of a data archive. For instance, the stated lifetime of CD and DVD storage is 20 years.

It is also possible that the hardware to access storage media may become obsolete: for example, 3.5 in. floppy disk drives are becoming rare on new PCs. This was a particular issue with early data storage systems using custom hardware (for instance, the 1986 BBC Domesday Book project [2]): there are also much more recent examples (for example, HD DVD).

2.2. Distributed and disjointed data organisation

It is common for data to be stored in different parts of an organisation. This is perhaps less of an issue for the main chain of scientific data from a large experiment, provided a strong policy on central data storage is in place and enforced.

There are however other types of data which should be considered: for instance, information about how processed data was arrived at, and meta-data about signals in a raw or processed data store. Design documents for the device and its components, and documentation for data formats and software packages, also needs to be considered as part of the organisation’s data archive.

2.3. File format obsolescence

An organisation will use many file formats: this could include third party formats (commercial or open source), and specialist formats which are developed in-house.

Over a decades-long timescale, commercial file formats will in many cases become obsolete – for example, files stored in many early word processing formats are now difficult to interpret. Specialist formats may also be problematic: if the format is documented inadequately, it will not be supportable in the absence of the original developers.

3. JET experience

JET has been operating since 1983, and all JET data back to the first pulse is still accessible. This section will review how the experience at JET has matched the general principles of data preservation at key points in its history. JET’s original data archive and the replacement systems developed in 2000 were described in a 2001 conference paper [3].

3.1. Original 1983 data archive

JET’s data archiving strategy was determined from the start by a strong directive by Hans-Otto Wüster, JET’s first director, who said in 1982:

“All diagnostics will collect data on every pulse and all of the data will be made available to every member of the JET staff”. [4]

This seems today an obvious statement to make, but in the early 1980s it was by no means a common approach: in many large experiments, data would be collected and stored by the scientists responsible for individual diagnostics.

The original scientific data archive for JET was split between the JPF (JET Pulse File), for raw data storage, and the PPF (Processed Pulse File) for processed data storage. The original systems were both custom-developed at JET, running on Norsk Data machines (for the JPF system), and an IBM mainframe (for the JPF and PPF systems).

The design of both systems was constrained by the available 1980s computing capacity and optimised for the available hardware. The use of custom developed systems gave us complete control of the data format and associated source code, and both systems were well documented.

3.2. Data archive replacement in 2000

In 2000, it was decided to decommission the JET IBM mainframe. As a result JET’s data archive was migrated to a UNIX platform [3]. At this point, 17 years after the first JET pulse, how did the experience at JET compare to the general principles of data preservation?

3.2.1. Media obsolescence

The IBM data storage hardware was well-managed with a hardware maintenance framework. This meant that the data storage media, and the hardware to manage the media, were regularly replaced during the lifetime of this system. During the process of transferring data to the new system, no data tapes were found to be damaged, and all JET data could still be read. As a result, the move to a new hardware platform for data storage was straightforward.

3.2.2. Distributed and disjointed data organisation

The management policy on centralised data storage was strongly enforced, and all data derived from JET data was stored in the PPF system to be accessible to all JET users. On the IBM mainframe, it was possible to scan user space for other types of data, recognise the format, and plan for the porting of this data. Data organisation is likely to be more problematic today, as it is much more common for information to be stored on users’ desktop computers.

3.2.3. File format obsolescence

This was the biggest issue in moving the JET data archive. The custom-developed JPF and PPF software was tied to the IBM platform, and required either a port of the low-level code, or replacement of the system, in order to be supportable on the new platform.

For the raw JPF data, it would have been very difficult to change the format of the data. The raw data format was closely coupled to the output from JET's plant systems: changing the format would have required changes to all main collection systems, data access and storage systems. As a result, the JPF file format was retained and the low-level data access code was ported from the IBM mainframe to UNIX. This was possible because all the code was in house, and the header files and comments and code itself have enough information to handle the complex structures involved. In addition, elements of the JPF software had already been ported between the original Norsk Data machines, Solaris and the IBM mainframe. While this was by no means a simple process, porting the full JPF client and server software was achieved within the required timescales and all data remains accessible.

A different approach was taken to the processed data in the PPF system. To ensure long-term accessibility, the decision was taken to move from the custom-developed format to widely used third party software, while retaining the original PPF data access interface to minimise disruption to users. For data storage, NetCDF was adopted as the new format: this is a widely used open source, flexible format with a large development team. A commercial RDBMS was used to store metadata about the files for fast indexing.

3.3. Data archive status in 2011

During the move to the new data archive in 2000, all of JET's data was found to be readable. The changes made at that point to refresh technology aimed to ensure that the data would remain readable into the future.

The 2000 data archive has now been in place for over 10 years, with incremental evolution during this time. We would expect further evolution for the rest of the lifetime of JET. Beyond the end of JET, the data will still be required for future experiments, so we would expect the archive in some form to outlive JET.

At the time of writing, there are over 16 million items in JET's data archive. JET's data is stored on an Oracle M8000 data server. New data is available in a rapid access cache in memory, and a RAID data cache has the capacity to store around one year of JET's operational data.

All data is archived to tape, with about one data tape filled per day. Three tape copies are stored: one in JET's tape silo for access by the data warehouse. The two other tape copies are stored in offsite backups: one at the other side of the Culham site, and the second at the Rutherford Appleton Laboratory 10 km from JET, with whom we have a reciprocal arrangement for offsite data storage. The three tape copies are written at different points in time, mitigating us against the risk of a single point of failure in writing our backup tapes.

To ensure data integrity, we have periodically read and checked every bit of data on moving from various backup and storage systems. All data (both raw and processed) is write once, eliminating the risk of data corruption on modification.

The data archive hardware and storage media has been regularly renewed since 2000. The current data server has been sized based on JET's current and predicted data curation rates, and the expected life of the hardware. If, at some point, JET becomes a purely archival site, it is possible the data server could be downsized.

3.4. Future JPF and PPF system evolution

The 2000 versions of the PPF and JPF systems are still operational and have coped well with the increase in JET data production rates. As with the hardware, there have been some incremental changes to both systems to cope with larger data items and new types of data.

The architecture adopted for the PPF system is still a good choice: there is no reason to switch from NetCDF for our underlying data storage in the short term. As the main consumer of JET data beyond the end of JET is likely to be ITER, it may make sense to adopt ITER's chosen solution for data storage in the long term, to ensure ease of ongoing maintenance.

The JPF system is still maintainable in its current form. In the long-term, it may become more difficult to change the underlying server code as the technical skills required to develop the complex C and FORTRAN 77 code involved become rarer in the job market. At some point moving the JPF data to a new system will have to be considered. This will be easier when JET is no longer producing new raw data, and the coupling of the JPF data to JET's plant systems is no longer an issue.

3.5. Other types of data

This paper focusses on JET's experience with scientific data, but it should be noted that the management of other types of data has been equally challenging over this time period. For instance, maintenance of design data is challenging due to the vendor-specific tools used to store and access this data. Databases must also be managed over long periods: a common approach to this has been to periodically refresh the database technology used to ensure continued access to the data [5].

4. Lessons for future devices

Fusion devices which are currently being constructed, or which are at an early stage in their lifecycle, can benefit from the experience of JET and the general principles of data preservation in designing or improving their data archives. It is important to note that it is much easier to address preservation issues early – ideally at the point of designing the archive – rather than try to fix deficiencies later.

Data preservation is an organisation wide issue – this is not just a CODAS/CODAC problem, or an IT problem, or a Physics problem. The danger of this is that the issue is not the responsibility of one department, and may easily be overlooked as people focus on their own responsibilities.

Long-term data preservation does not appear to be an area of wide discussion in the fusion community, although the nature of our experiments means that it must be considered. Much research is ongoing in other communities, which the fusion community could learn from.

The Space Science community originated an ISO standard for digital archives – the OAIS (Open Archival Information Systems) model [6]. In high energy physics, several workshops on Data Preservation and Long-Term Analysis have been held, with many papers available online [7]. Some of the most advanced research in this area has been in “memory institutions” including national archives and libraries. Many of these have already implemented long-term archives for digital information. In addition, useful tools have been made available by these institutions such as PRONOM [8], an online registry of information about file formats and other components required for long-term access to data.

Any data producing organisation should consider how to address the key challenges of data preservation. This applies to all types of data produced by the organisation, not just the scientific data discussed here.

4.1. Media obsolescence

Over a decades-long timescale, physical storage media and the hardware required to access the media will inevitably degrade or become obsolete. Organisations should plan a strategy to

deal with this, expecting to periodically refresh their data storage technology. Keeping abreast of technical developments so that the technology is replaced before it becomes obsolete is crucial.

4.2. Data organisation

Strong policies on central storage of data should be defined by a fusion lab and enforced. This should apply to all scientific data. It should also apply to other types of data such as design data, reports, meeting minutes, etc.: ITER's strategy of storing all such data centrally in ICP [9] is a good example of this approach.

A key concern in the scientific archiving community which is not currently being addressed in the memory institution domain is that of data provenance – ensuring the traceability and justification of published results. This should be built into the design of a scientific data archive from the start – this was not the case with the JET archive.

4.3. File format obsolescence

An organisation should store centrally information about all file formats in use, and what they are used for. This information should be regularly updated. Centralised repositories of file format information such as PRONOM [8] can help identify potential problems with supportability of these formats.

When selecting a format for data storage, the longevity and future potential of the format should be considered: benchmarking current performance of the format is not enough. Questions to be asked about any format include:

- How widely is it used
- How large is the development team and support community?
- Is it open source? If so, how well written is the software?
- Could your organization support the software itself if it had to?
- How easy would it be to migrate from this format to another, if necessary?

4.4. Active preservation

It is not sufficient simply to save bytes of data: it is vital to ensure data is still readable into the future. There are two main approaches to keeping data readable when the original format is no longer supportable.

Emulation involves emulating the hardware and software environment that your data is readable on. In this approach, the “data experience” is retained: i.e. the data appears identical to a user but the actual data values may not have been retained.

Migration involves moving data to a new format which is still supported. In this approach, the actual data values are retained. This is the approach used by JET with the PPF system in 2000. Migration is the most appropriate approach in the fusion domain.

5. Conclusions

The experience at JET, and the more general principles discussed here, show that long-term preservation is not something that can be achieved without planning and continuous attention. It is a complex problem, but it can be done, and has been achieved at JET over almost three decades. Data created from the first JET experiment is still readable today: photographs taken on the first day of JET operation show that a plot of the first plasma current produced in 1983 is identical to a plot produced today.

Those who are responsible for data archives in fusion should consider how their archive will be maintained in decades to come, when the technological landscape is likely to be very different and the experts currently responsible for the system are likely to be no longer in the workplace.

Acknowledgments

This work, supported by the European Communities under the contract of Association between EURATOM and CCFE, was carried out within the framework of the European Fusion Development Agreement. The views and opinions expressed herein do not necessarily reflect those of the European Commission. This work was also part-funded by the RCUK Energy Programme under grant EP/I501045.

References

- [1] P. Sinclair, Are you ready? Assessing whether organisations are prepared for digital preservation, in: *iPRES 2009: The Sixth International Conference on Preservation of Digital Objects*, California Digital Library, UC Office of the President, 2009.
- [2] *Mind the Gap: Assessing Digital Preservation Needs in the UK*, 2006, <http://www.dpconline.org/docs/reports/uknamindthegap.pdf>.
- [3] R. Layne, M. Wheatley, New data storage and retrieval systems for JET data, *Fusion Engineering and Design* 60 (2002).
- [4] J.P. Christiansen, Integrated analysis of data from JET, *Journal of Computational Physics* 73 (1987) 85.
- [5] R. Layne, et al., The central physics file, a high level fusion database, in: *6th Workshop on Fusion Data Processing, Validation and Analysis*, Madrid, 2010.
- [6] *Space Data and Information Transfer Systems – Open Archival Information System – Reference Model*, <http://www.iso.org/iso/catalogue-detail.htm?csnumber=24683>.
- [7] R. Mount, Data Preservation and Long-term Analysis in High-energy Physics, *Ariadne Issue* 58, 2008.
- [8] PRONOM: <http://www.nationalarchives.gov.uk/PRONOM/Default.aspx>.
- [9] C. Capuano, F. Carayon, V. Patel, ICP (ITER Collaborative Platform), 7th IAEA TM on Control, Data Acquisition, and Remote Participation, Aix-en-Provence, 2009.