



Electronic Records Archives (ERA)

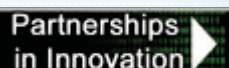


EXPLORE ERA

- [Welcome Message](#)
- [About ERA](#)
- [Program Management Information](#)
- [Research](#)
- [Advisory Committee on the Electronic Records Archives \(ACERA\)](#)

ERA RESOURCES

- [Papers and Presentation by the National Archives](#)
- [Other Helpful Sites](#)



[Go to the Electronic Records Archives \(ERA\) Main Page](#)


[Print](#)

[E-mail](#)

[Bookmark/Share](#)

Preservation and Migration of Electronic Records: The State of the Issue

Kenneth Thibodeau

The problem of preserving electronic records

The two-edged sword of continuing progress and rapid obsolescence of information technology is the most often cited, but perhaps not the most significant challenge archives face in the endeavor to preserve electronic records. Organizations rely more and more on digital technology to produce, process, store, communicate, and use information in their activities. Thus, the quantity of records being created in electronic form increases. In the experience of the National Archives and Records Administration of the United States, it increases exponentially. The technological challenge is compounded by the continuing extension of information technology in terms of the types of information objects it produces, and again in terms of its applicability to different spheres of activity and different types of actions within those spheres. The resultant records are increasingly diverse and complex. The impact is not only on individual records, but on the archival fonds as a structured whole.

Approaches to the problem of preserving electronic records

The field of information technology has, by and large, ignored the problems of long term preservation. If anything, one could say that the market has tended to exacerbate the problem of preserving electronic records. The pressures of competition have led the industry to obey Moore's law, replacing both hardware and software on a frequency of two years or less.

In one area, however, there has been some improvement in recent years: that of digital storage media. From the 1980s there was a trend towards storage media that were more fragile and less stable over time. In recent years, this trend, if not reversed, has been offset somewhat by the introduction of more stable and reliable media. Current research and development efforts offer the prospects of improved options for long term storage of digital information, notably in the areas of ion-milling and holographic media. But archival concern with digital media should not be limited to their durability. The ICA *Guide to Managing Electronic Records* sets out seven criteria for media used for preserving electronic records:

- open standards for digital recording on the medium,
- robust methods for preventing, detecting and reporting errors,
- sufficient market penetration,
- known longevity,
- known susceptibility to degradation or deterioration,
- a favorable cost/benefit ratio, and
- availability of methods for recovering from loss.

Whatever relief archives may find in the area of digital storage is more than offset by the increasing diversity, complexity and spread of electronic records. In recent years, increasing attention has been devoted to problems of digital preservation in a variety of spheres and professions. Several different approaches have been proposed. A few have been tried in test mode, fewer in actual practice. In practice, the experience of archives is largely limited to relatively simple technical formats, such as flat files. Some institutions have developed computer applications for preserving potentially complex databases. These include CONSTANCE at the National Archives of France, AERIC at NARA, ERICSON at the National Archives of Canada, and similar systems in Sweden, the United Kingdom and elsewhere. Significant preservation projects addressing the actual preservation of digital formats, at various stages of research or development, include the bundles proposal of the British Standards Institute, the CEDARs project at the University of Leeds, England, the Victoria Electronic Records System in Australia, the emulation experiment at the Royal Library in The Netherlands, the Universal Preservation Format sponsored by the WGBH Educational Foundation in Boston, and the Highly Integrated Information Processing and Storage technology being developed at Carnegie-Mellon University in the U.S. Current initiatives are pursuing quite a variety of approaches. The proposed solutions can be categorized into five broad categories:

- preserving the original technology used to create or store the records;
- emulating the original technology on new platforms;
- migrating the software necessary to retrieve, deliver, and use the records;
- migrating the records to up-to-date formats; and
- converting records to standard forms.

These approaches define a spectrum ranging, in broad terms, from no change in the records or the technological context in which they exist to one in which the original hardware and software have disappeared and the digital format of the records has changed. Each of these methods has pros and cons. None of them is entirely satisfactory. On the one hand, in general, one can say that the closer one stays to the original technology and original digital format of the records, the less the problem of authenticity; however, it is also obvious that the closer one stays to original technology, the more complex and more impractical the approach becomes over time. More complex because, as records continue to accumulate over time, there will be more and more varieties of technology that the archives would have to maintain. More impractical because, first, support for obsolete technologies will eventually disappear and, second, the distance and difference between the preserved technology or technical artifacts B including the records B and the best available technology for preserving, managing, retrieving and delivering the records will increase continuously. On the other hand, while moving ahead as technology progresses can eliminate such practical problems, it can entail loss or corruption of records.

The need for an archival approach to preserving electronic records

All of these approaches to preserving electronic records have in common the objective of solving technological problems related to the passage of time. None of them actually focus on the objective of preserving records. This technological orientation is misdirected because success in solving technological problems does not necessarily imply any success, or even relevance, in addressing archival requirements for the preservation of records. Logically, archival principles and objectives should dictate the requirements that technical solutions must satisfy. Archival requirements for preservation must be based on the conception of electronic records, not as the products of computer applications, but as the instruments and by-products of the practical activity of a records creator. The ultimate criterion for success in the preservation of electronic records is not whether they remain true to some given technological materialization, but whether they continue to provide authentic

evidence of the activities in which they were created.

An architecture for archival preservation

Clearly, the archival profession needs to determine specific requirements for the preservation of different types of records, and also to guarantee respect for provenance and the integrity of archival fonds over time. The InterPARES project, directed by Professor Duranti, brings together archivists from universities and archival institutions, along with computer and information scientists and engineers, from around the world in a concerted effort to delineate specific archival requirements for preserving authentic electronic records. InterPARES is working to define the archival requirements for authenticity on the basis of archival science and diplomatics.

Simultaneously, the InterPARES Preservation Task Force is examining technical issues related to digital preservation and developing a formal model of the preservation function as viewed from the perspective of the juridical or physical person responsible for preserving electronic records. While this work is still in progress, there are several ideas which have been proposed that are worth citing at this time. One key idea is that, strictly speaking, it is not possible to preserve electronic records; it is only possible to maintain the ability to reproduce electronic records. It is always necessary to retrieve from storage the binary digits that make up the record and process them through some software for delivery or presentation. Analogously, a musical score does not actually store music. It stores a symbolic notation which, when processed by a musician on a suitable instrument, can produce music. B Presuming the process is the right process and it is executed correctly, it is the output of such processing that is the record, not the stored bits that are subject to processing. This concept has important consequences. It shifts priority in preservation of electronic records from their storage over time, to the integral processes of putting the records into archival storage, getting them out of storage, and delivering them to future researchers. The recognition that electronic records must inevitably be reproduced accentuates the importance of being able to demonstrate the integrity and authenticity of the records. This entails extending the traditional concept of an unbroken chain of custody into one of an unbroken process of preservation. As defined in the ICA Guide, An electronic record is preserved if and only if it continues to exist in a form that allows it to be retrieved, and, once retrieved, provides reliable and authentic evidence of the activity which produced the record. Demonstrating the authenticity of electronic records depends on verifying that:

- the right data was put into storage properly;
- either nothing happened in storage to change this data or alternatively any changes in the data over time are insignificant;
- all the right data and only the right data was retrieved from storage;
- the retrieved data was subjected to an appropriate process, and
- the processing was executed correctly to output an authentic reproduction of the record.

Parallel to the InterPARES project, the National Archives and Records Administration is sponsoring research into the development of an information management architecture designed to address archival requirements for the preservation of electronic records. This architecture implements the proposed ISO standard for an Open Archival Information System (OAIS). The architecture extends that general reference model by articulating archival requirements. To address the basic problem of continuing change in technology over time, the architecture postulates that archival information systems should independent of the particular technology used to implement them at any time. That is, an archival information system should be built in such a way that it is possible to replace any component of hardware or software used in the system with minimal impact on the rest of the system and with no impact on the preserved collections of records.

Collection-based persistent object preservation

The information management architecture is being developed in the U.S. National Partnership for Advanced Computational Infrastructure. The Partnership is a collaboration of 46 institutions nationwide, and 6 foreign affiliates, with the San Diego Supercomputer Center serving as the leading edge technical resource. The research is addressing archival requirements for preservation of records, including respect for provenance. Rather than focus on technological problems, the method focuses on the objects that are to be preserved. In this case, the objects are records and also collections of records, as organized within archival fonds at all levels of hierarchy.

The method of collection-based persistent object preservation consists of identifying the properties of the objects to be preserved; expressing those properties in explicit, abstract models; and applying those models to transform the objects into an independent technological format suitable for long-term preservation. In the archival domain the development of this method started with the conception of the essential properties of records expressed in the ICA Guide on electronic records; that is, *A record is recorded information produced or received in ... an institutional or individual activity and that comprises content, context and structure sufficient to provide evidence of the activity regardless of the form or medium*. The essential *structure* of a record is its documentary form. This form may be expressed in the digital format in which the record is stored, but it is not necessarily identical to the digital format. Therefore, a transformation of the record which replaces one digital method with another one that is more suitable to long-term retention, preserves the record so long as it maintains the essential documentary form of the record. The immediate *context* of a record is its archival bond: the position of a record with respect to other records in the archival fonds. In our research, we have extended the list of essential properties of records beyond content, structure and context to include the *appearance* of the record. We are also addressing a special type of content that is unique to electronic records: hyperlinks.

Persistent Object Preservation expresses the structure of records using eXtensible Markup Language (XML) Document Type Definitions. The method encapsulates records using the metadata defined in these models, transforming records into a format that is independent of any specific technology. The research has demonstrated that this method can be applied to collections of records as well as to individual records. That is, one can construct a Document Type Definition to capture and preserve the structure of any archival collection, of arbitrary complexity, from individual files through series and classes to entire archival fonds.

The research is exploring different ways of preserving the appearance of records. One way is to use a technology known as Multi-Valent Documents to capture and retain a bitmapped image of the document. MVD enables the image to be retained not as a version of the document, but as a layer of the document object modeled as an acyclic directed tree. Another possible means of preserving appearance is through the eXtensible Style Sheet Language (XSSL) available in the XML standard. Using style sheets to capture the attributes of appearance is especially advantageous for types of applications, such as databases and geographic information systems, where stored data elements may participate in many different records. In such systems the records are likely to be expressed as views, forms, or reports which extract specific subsets of the data and present them in predefined formats. A different style sheet can be defined for each of these formats.

The method extends beyond the preservation of archival collections of records over time. It also addresses the key archival functions; notably, the accessioning of records into the archival repository, the establishment of intellectual control over the records, and the delivery or dissemination of the records to researchers. This extension of the persistent object approach is consistent with the basic premise of object oriented methodology which starts with the recognition that an object has behaviors or methods, as well as attributes. One of the essential behaviors of a record is that it occupies a specific position in relation to other records in the archival fonds. This behavior expresses the immediate context of the

record and is the basis for arriving at its significant context; that is, the activity of which the record provides evidence. The transformation of records into a persistent object format not only enables the records to be preserved indefinitely into the future, it also makes it possible to benefit from advanced technologies, which have not even been invented yet, to search, access and deliver the records in the future. This is made possible though the separation of context, structure and appearance in explicit schemas expressed in simple textual form. Over time, it will not be necessary to migrate the materials stored in persistent object form to new technologies, but only to interpret the schema metadata so that it can be used in future technologies.

Viability of the persistent object preservation method

The initiative to develop the collection-based persistent object method for preserving electronic records is still in the stage of research and development, and will remain in this stage for some time. Nonetheless, there are substantial reasons, in both the technical and the archival domains, to assume that it will be successful. In the domain of technology, two facts should be highlighted. First, the research is not developing any special technologies to suit archival needs. Rather, it is building archival solutions on the basis of technologies which are seen as essential to the next generation Internet and information infrastructure and as keys to electronic commerce and electronic government. Archives should benefit, therefore, from widespread market support for the enabling technologies. Second, while the research addresses archival requirements specifically, the method has broad application in other areas, such as digital libraries, museums and collections of scientific data. Thus, archival institutions can collaborate with organizations in these other domains to develop from the enabling technologies solutions for long-term preservation and access. In the archival domain, the promise of the persistent object preservation method has been demonstrated in several empirical tests, applying the method to a variety of collections across a broad quantitative scale. These demonstrations involved bringing the collections into the archival information system from external sources; examining the documents, databases, images, geographic information systems and other digital objects tested in order to generate XML models; transforming the records and capturing collection organization according to these models; storing the transformed collections and related meta-data; and retrieving and presenting the preserved records using technologies completely different than those which had originally been used to create and store the records.

Conclusion

The persistent object preservation method offers several advantages to archives. It provides a coherent and comprehensive framework that can be specifically tailored to archival requirements. Through abstraction of the context, structure and appearance of the contents of digital objects, it provides a single, but highly adaptable method that serves at once the need for preserving authentic electronic records over time, for adhering to archival principles, such as provenance, and for performing core archival functions. Moreover, the persistent object framework permits the simultaneous adoption of other techniques if the need arises. Clearly a substantial amount of research, analysis, testing and evaluation needs to be completed before this method reaches its full potential. Nonetheless, the positioning of this method in the center of major developments in computer science and information technology offers great potential for making of electronic records not so much a problem for preservation, but an opportunity for archives to achieve their objectives to a greater extent and at a higher level than has been possible before now.

The U.S. National Archives and Records Administration
8601 Adelphi Road, College Park, MD 20740-6001
Telephone: 1-86-NARA-NARA or 1-866-272-6272

