

A Generic XML Language for Characterising Objects to Support Digital Preservation

Christoph Becker, Andreas Rauber
Vienna University of Technology, Austria
{becker,rauber}@ifs.tuwien.ac.at

Volker Heydegger, Jan Schnasse,
Manfred Thaller
University of Cologne, Germany
{herrmanv,jan.schnasse,manfred.thaller}@uni-koeln.de

ABSTRACT

The dominance of digital objects in today's information landscape has changed the way humankind creates and exchanges information. However, it has also brought an entirely new problem: the longevity of digital objects. Due to the fast changes in technologies, digital documents have a short lifespan before they become obsolete. Digital preservation, i.e. actions to ensure longevity of digital information, thus has become a pressing challenge. Different strategies such as migration and emulation have been proposed; however, the decision between available tools for format migration is very complex. Preservation planning supports decision makers in reaching accountable decisions by evaluating potential strategies against well-defined requirements. Especially the evaluation of different migration tools for digital preservation has to rely on validating the converted objects and thus on an analysis of the logical structure and the content of documents. This paper presents the eXtensible Characterisation Languages (XCL) that support the automatic validation of document conversions and the evaluation of migration quality by hierarchically decomposing a document and representing documents from different sources in an abstract XML language. We present the context of the development of these languages and tools and describe the overall concept and features of the languages and how they can be applied to the evaluation of digital preservation solutions.

Categories and Subject Descriptors

H.3 [Information Storage and Retrieval]: H.3.7 Digital Libraries; I.7 [Document and Text Processing]

Keywords

Digital Preservation, Preservation Planning, Content Characterisation, Migration, XML languages

1. INTRODUCTION

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SAC'08 March 16-20, 2008, Fortaleza, Ceará, Brazil
Copyright 2008 ACM 978-1-59593-753-7/08/0003 ...\$5.00.

Today's information landscape has almost completely gone digital. Digital objects are the primary way we create, exchange, and discover information. However, the rapid changes in technologies, file formats, and information systems make the longevity of digital information a challenging problem. Digital objects are inherently ephemeral. They become obsolete for a variety of reasons, be it simple media failure, obsolescence of file formats and tools, or the inability to find metadata that is necessary to understand the content. Yet, in many instances, there is no analog counterpart to which one could refer if digital information is lost.

Digital preservation denotes the efforts to preserve digital objects for a given purpose over long periods of time. The urgency of digital preservation has recently been reemphasised by the results of a survey among archiving professionals[17]. In the last years, numerous research initiatives have started around the world that aim at mastering this challenge.

The two strategies generally considered to prevail today are migration and emulation. While migration operates on the objects and transforms them to more stable or more widely adopted representations, emulation operates on the environment of an object, trying to simulate the original environment that the object needs.

Consider a large collection of documents written in an old version of Word several years ago. This application does not run on modern operating systems. One could try to emulate the original operating system; but the emulation software still depends on current hardware, and emulating the hardware might be even harder. On the other hand, one could also migrate the documents to a current version of Word, or to PDF. While this would lose the original look-and-feel of the application, it would probably preserve the content and layout. With PDF, one would exchange the tool support provided by text processing software, as well as metadata such as an edit history, for a file format that is generally considered to be quite stable.

An important part of ongoing efforts in many large international projects is the outreach to vendors for advocating document engineering technologies for sustainable documents. The effects can be seen in standards such as PDF/A [9] or the Open Document Format (ODF) [11]. However, many objects exist and many more are created every day that face the threats of obsolescence. Hence, ex-post actions for preserving access to content are necessary. Performing actions on objects always risks damaging the content; but preserving authentic records also means being able to prove authenticity[8, 15].

Various migration tools are available for standard file formats such as office documents; the picture is less positive for more exotic and complex compound objects. However, even within migration tools for office documents, variation regarding the quality of conversion is very high. Some tools fail to preserve the proper layout of tables contained in a document; others miss footnotes or hyperlinks. Finding out which information has been lost during a conversion, and if this loss threatens the value of the object for a given purpose, is a very time-consuming task. Some losses might be acceptable, while others threaten the authenticity of documents. For example, if migrating the collection of Word documents mentioned above results in a loss of page breaks, this might be irrelevant if the textual content is the only thing of interest. However, if there are page references in the text, this loss might be unacceptable.

A variety of tools performing preservation actions such as migration or emulation exist today; most often, there is no optimal solution. The complex situations and requirements that need to be considered when deciding which solution is best suited for a given collection of objects mean that this decision is a complex task. Preservation planning aids in the decision making process by evaluating available solutions against clearly defined and measurable criteria. This evaluation needs verification and comparison of documents and objects before and after migration to be able to judge migration quality in terms of defined requirements. It thus has to rely on an analysis of the logical structure of documents that is able to decompose documents and describe their content in an abstract form, independent of the file format. Especially considering migration actions working on large numbers of objects, it is essential to validate the authenticity of transformed objects automatically. When migrating a million documents from ODF to PDF/A, validation of these objects can not be done manually.

This paper presents the eXtensible Characterisation Languages (XCL) that support the automatic validation of document conversions and the evaluation of conversion quality by hierarchically decomposing documents from different sources and representing them in an abstract XML language. We outline the basic concepts underlying the languages. We then describe the main architecture and features of XCL and discuss its application in the context of digital preservation.

The remainder of this paper is structured as follows. The next section outlines related work in the area of document engineering, digital preservation and the usage of XML in document conversion and extraction. Section 3 presents the Extensible Characterisation Languages and their usage within the context of Preservation Planning. Section 4 draws conclusions and points out directions for future work.

2. RELATED WORK

Digital preservation has become a highly active area of research in the last decade, as many memory institutions realised that the information they are responsible for will cease to exist within years[19].

A large part of the discourse has focused on discussing the dominant strategies, migration and emulation. Lawrence et. al.[12] analyzed risks of migration actions. While emulation is in principle widely applicable, the complexities and costs associated with it form a major obstacle to its wider adoption. Bearman strongly argues against the usage of emulation in [1]. Rothenberg as one of the main proponents

of emulation in digital preservation calls for encapsulation techniques to support emulation[14]. Encapsulation as a complementary strategy packages the object to be preserved together with instructions on how it can be interpreted. Often the encapsulation layer is expressed in XML, which has a stronghold in digital preservation[6]. It is used not only for encapsulation, but also as a target file format for migration[13, 3] or as a metadata container such as in PREMIS.¹

In the discussion of file formats for long-term preservation it has recently been increasingly understood that the simple decision to use ‘PDF’ or ‘TIFF’ is highly problematic, as in both cases it is possible to create either files with very high as well as very low preservation value. In the case of extremely rich formats like PDF, this has led to the definition of subsets of the rules comprising the format which are considered safe for preservation, and furthermore to a tendency to identify informal ‘subformats’ of file formats[18, 9].

At the heart of a preservation endeavour lies preservation planning. It is a core entity in the Reference Model for an Open Archival Information System, the OAIS model, which is a widely used model for archives[4]. It is also a core part of the EU project ‘Preservation and Long-Term Access via Networked Services’ (PLANETS)² which is creating a distributed service oriented architecture for digital preservation. Strodl et. al.[16] present the PLANETS preservation planning methodology that aids in reaching well-founded decisions. The procedure defines measurable requirements for preservation strategies in a hierarchical form and evaluates them in a standardised testbed setting. Similarly, Ferreira [7] presents a Service Oriented Architecture for recommending and performing format migrations based on pre-specified requirements. However, so far most of the evaluation of tools against these requirements has to be done manually. For example, to evaluate if the layout of a document has been preserved during migration, a human has to look at the files and compare them with each other. This is not feasible for large collections; automated services have to be integrated that characterise content to support this evaluation.

A number of tools and services have been developed that perform content characterisation specifically for digital preservation. The National Library of New Zealand Metadata Extraction Tool³ extracts preservation metadata for various input file formats. Harvard University Library’s tool JHove⁴ enables the identification and characterisation of digital objects. However, both tools only extract metadata such as the presence of specific file format features in a document; they do not describe the content of a document.

Some solutions exist for transforming, matching, and comparing structured documents. Díaz describes the usage of XML for handling the conversion problems that arise in the exchange of business documents between organisations[5]. In the area of grid computing, the Global Grid Forum Data Format Description Language Working Group has been working on a language called DFDL[2] which extends XML Schema. The aim is to describe *the structure of binary and character encoded (ASCII/Unicode) files and data streams so that their format, structure, and metadata can be exposed*⁵.

¹<http://www.oclc.org/research/projects/pmwg/>

²<http://www.planets-project.eu>

³<http://meta-extractor.sourceforge.net/>

⁴<http://hul.harvard.edu/jhove>

⁵<http://forge.ggf.org/sf/projects/dfdl-gw>

3. THE EXTENSIBLE CHARACTERISATION LANGUAGES

3.1 Introduction

As outlined above, converting any number of documents from one format into another inevitably raises the issue of preserving authenticity. Moreover, to confidently choose between alternative target formats and tools, one has to evaluate their suitability in a given context. This leads to the following underlying questions.

1. Which information contained in the old format is also contained in the new format?
2. Which information relevant to the usage of the content of the old format is contained in the new format?
3. Is the conversion process $a(old,new)$ better than $b(old,new)$, i.e. does it preserve more of the relevant information contained within the object?

Comparing information in two different file formats implies the following requisites.

1. An abstract way of expressing the information in a format-neutral model. This is henceforth called an *extensible characterisation definition language (XCDL)*.
2. A way of extracting information in specific file formats and describing it using XCDL. While it would be theoretically possible to create an extraction tool for every given file format, in practice this is not feasible. A better solution is to define a comprehensive *extensible characterisation extraction language (XCEL)* and implement an extractor component that is able to interpret this language.
3. Algorithms for comparing two XCDL descriptions for degrees of equality.

Such a language should be defined so generic that it supports the description of arbitrary file formats and thus the extraction of characteristics from *any* given file.

An XCDL document describes the content of a specific file with format type X, tagged in XML according to the XCDL language specifications, and is processible through an XCDL interpreter. An XCEL document describes what information can be extracted from any given file of format X, enabling an XCEL processor to extract this information and express it in XCDL. XCEL thus creates a mapping between the declarative description of the information in a physical file and its abstract interpretation outside of a format specification. Both XCDL and XCEL are meta-languages defined using XML Schema. In contrast to other applications of XML in digital preservation, XCL does not migrate digital objects as a whole to XML nor store exclusively preservation metadata; it transforms the entire content of an object into an abstract unified form. A key application we focus on is the comparison of different representations of the same object in order to validate migration within the preservation planning procedure.

The next sections will describe the basic architecture of both characterisation languages. We will then outline an example of how they can be applied in practice.

3.2 The extraction language XCEL

An XCEL document comprises the following components.

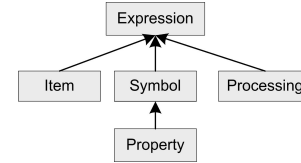


Figure 1: The structuring elements of XCEL

```

<!-- The complete IDAT chunk is expressed as one item that
prescribes all its children to appear in the given order -->
<item xsi:type="structuringItem" order="sequence"
  identifier="IDIO2" multiple="true">
  <symbol identifier="IDATLength" interpretation="uint32"
    length="4"> </symbol>
  <property identifier="IDATChunkType" interpretation="ASCII"
    length="4">
    <value>IDAT</value>
  </property>
  <!-- Set the length of the expression "IDATChunkData"
to the value given by the expression "IDATLength"-->
  <processing type="pushXCEL" xcelRef="IDATChunkData">
    <processingMethod name="setLength">
      <param valueRef="IDATLength"/>
    </processingMethod>
  </processing>
  <symbol identifier="IDATChunkData" interpretation="uint8">
    <name>normData</name>
  </symbol>
  <symbol identifier="IDATCRC" interpretation="uint8" length="4">
    <name>checksum</name>
  </symbol>
</item>
  
```

Figure 2: XCEL description of a PNG chunk

1. **Preprocessing** information includes configuration tasks affecting the behaviour of the XCEL interpreter.
2. The **format description** is the core part defining the structure of an object.
3. **Templates** describe recurring structures such as number formats in ASCII based file formats.
4. **Postprocessing** instructions define actions to be performed on the result of the format processing.

Figure 1 shows the abstract relations of XCEL expressions. This structure is based upon the assumption that any file format can be expressed as (a) a set of hierarchies of blocks of content, all of which can (b) be addressed from within but also outside of these hierarchies.

An XCEL format description starts with an *Item*, a container element that can have different content models, similar to the XML-Schema content models. A *Symbol* is an expression that adds a name or ID to a specific part of the byte stream. A *Property* is a *Symbol* with a predefined value. The *Processing* element models an expression that is used to call built-in methods for the extraction processor. This structure describes file formats in a tree where each branch describes one possible variation. It is the job of the XCEL processor to find out which branch matches to a given file.

Figure 2 shows the XCEL description of the IDAT chunk of a PNG[10] image. Every chunk in PNG consists of the consecutive parts **length**, **chunk type**, **chunk data** and **CRC**, where the **length** is a four byte unsigned integer that contains the length of the **chunk data** field, **chunk type** is a four byte ASCII keyword, **chunk data** is a field that can contain any data structure, and **CRC** contains a checksum.

The XCEL processing software ('Extractor') processes the binary files of given formats using the specific XCEL doc-

```

<normData id="n6">An important word</normData>
<property id="p8" source="raw">
  <name>Fontname</name>
  <valueSet id="v1">
    <labVal>
      <val>Times-Roman</val>
      <type>XCLLabel</type>
    </labVal>
    <dataRef ind="normSpecific">
      <ref id="n6" start="0" end="1"/>
      <ref id="n6" start="13" end="16"/>
    </dataRef>
  </valueSet>
  <valueSet id="v2">
    <labVal>
      <val>Times-Bold</val>
      <type>XCLLabel</type>
    </labVal>
    <dataRef ind="normSpecific">
      <ref id="n6" start="3" end="11"/>
    </dataRef>
  </valueSet>
</property>

```

Figure 3: XCDL representation of primary information and corresponding properties

uments created for these formats. Currently we have prepared XCEL documents for various file formats, focusing on the image, text and audio data domain (e.g. TIFF, PNG, GIF, WAV, and PDF). The Extractor is conceived in such a way as to be able to process any additionally created XCEL document without modifications of its core implementation. Thus, to enlarge the spectrum of supported file formats one only has to write an XCEL document for that format.

3.3 The definition language XCDL

Figure 3 provides a part of an XCDL description of a text document containing the phrase ‘An **important** word’. A common characteristic of content models is a separation between primary information and properties of that information. Within the textual domain this separation consists e.g. in the difference between the string ‘an important word’ and the means by which we indicate that the single words in that string are expressed in specific fonts. The corresponding XCDL representation is provided below. The *normData* tag wraps the primary information in a context-free manner, removing or transforming all format-specific information as well as its specific representation. We call this kind of representation normalised data. Text encoding is translated into UTF8 by default. The fonts are described within the *property* tags. In this case we have a property describing the fonts used. For each different font a value set is created. The font name appears as a labelled value, typed as *XCLLabel*. XCL labels refer to exactly defined terms in the XCL properties ontology. The *dataRef* tags define positions within the normalised data, indicating where the specific fonts are applied. This basic structure is common for all underlying content models: In the case of images, e.g., the primary stream of bytes describing the pixels can have properties, which are applicable to an image as a whole (e.g. a gamma correction) or only to parts thereof, as e.g., an embedded explanatory text in the image.

For preservation purposes, the properties extracted from a file may include a category of properties which are not needed to model the content of the file. Consider, e.g., a file format which compresses a part of the data it contains. A proper XCL extractor, which extracts the content of the

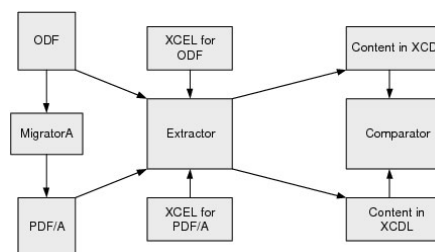


Figure 4: Using XCL to compare migrated documents

file and expresses it in XCDL, has to be able to apply that algorithm in order to decompress the content. Once this is done, the algorithm applied to the original file becomes irrelevant, as the content is simply the result of its application. For preservation purposes - basically tracking the history of a file and its authenticity - properties like ‘originally compressed by algorithm X’ can be expressed.

3.4 Comparing digital objects

Figure 4 shows a scenario for applying XCL in the context of format migration. After converting a document from ODF to PDF/A, the XCDL documents of the original and the transformed object can be compared using an interpretation software. A comparison tool (‘Comparator’) for XCDL documents is currently under development. Key objectives are the property-specific definition of metrics and their implementations as algorithms in order to identify degrees of equality between two XCDL documents. In its core functionality the comparator loads two XCDL documents, extracts the property sequences and compares them according to comparison metrics which are defined with respect to the types of the values in the value sets. In the example of Figure 3, the comparator looks up the defined metrics for property ‘Fontname’ and executes the comparison according to the metrics definition. This can be a simple binary comparison that checks the XCL ontology for the entries ‘Times-Roman’ and ‘Times-Bold’. For other properties such as as possible deviation of font size, absolute or relative difference measures can be used.

To verify the approach we migrated a benchmark corpus of PNG files⁶ to TIFF and compared the resulting XCDL documents. In contrast to other tools such as JHove or tiffInfo⁷, XCL was able to extract file properties as well as the normalised content from all files. Comparing the *normData* with a tool revealed that conversion of images with specific characteristics was not working properly.

For evaluating preservation strategies, preservation planning activities define requirements that a solution has to meet. Often, a complete and extensive comparison is not needed. The comparator offers the possibility to select and weigh only a subset of properties, thus enabling users to regulate the relevance of different properties with respect to their specific needs. By mapping the content structures in XCDL as well as the results from the *Comparator* to the requirements, performance comparisons across different preservation strategies can be defined and recommendations for a solution can be given in an automated way.

⁶<http://www.schaik.com/pngsuite>

⁷<http://remotesensing.org/libtiff/man/tiffinfo.1.html>

4. DISCUSSION AND OUTLOOK

While a range of tools exist today for migration between different document formats, the evaluation of suitability of these tools is highly complex. Digital preservation is in need for automated characterisation services that support preservation planning in evaluating potential strategies. These services need an abstract means of describing a document's content in an interoperable, format-independent way.

When comparing the content of two files stored in two different formats, we have to distinguish between the abstract content and the way in which it is wrapped technically. On a very abstract level, this will for a long time be impossible: Whether an image of a hand-written note contains the same 'information' as a transcription of that note in UTF-8 is philosophically interesting, but scarcely decidable on an engineering level. In a more restricted way, a solution is possible if we express the content stored in different file formats in terms of an abstract model of that type of content.

The extensible characterisation extraction and definition languages presented in this paper are an important step towards this goal. The extraction language XCEL allows the extractor component to extract the content of any document provided in a format for which an XCEL specification exists. The content is described in the description language XCDL and can thus be compared to other documents in a straightforward way. This differentiates the XCL approach from the approach used by JHove and similar projects. The XCL does not attempt to extract a set of characteristics from a file, but it proposes to express the complete informational content of a file in a format independent model.

The DFDL language, on the other hand, concentrates on exact typing of data formats for interoperability on the grid. Consider a binary file with the content '0000000000100000'. Using DFDL, it is possible to express that the file contains an unsigned 16 bit number in big endian encoding, i.e. 32. XCL is able to express that the file contains a 16 bit number in big endian encoding meaning `imageWidth=32`. Thus XCL also intends to describe the semantics of a file.

This paper outlined the basic architecture of the two characterisation languages, provided examples of how they can be applied in practice and discussed the potentials of the proposed approach in the context of digital preservation and preservation planning. Future work will be geared towards implementing automated verification and evaluation of different tools and integrating comparison services into the decision support software for preservation planning.

Acknowledgements

Part of this work was supported by the European Union in the 6th Framework Program, IST, through the PLANETS project, contract 033789.

5. REFERENCES

- [1] BEARMAN, D. Reality and chimeras in the preservation of electronic records. *D-Lib Magazine* 5, 4 (April 1999).
- [2] BECKERLE, M., AND WESTHEAD, M. GGF DFDL Primer. Tech. rep., Global Grid Forum Data Format Description Language Working Group, 2004.
- [3] BRANDL, S., AND KELLER-MARXER, P. Long-term archiving of relational databases with Chronos. In *First International Workshop on Database Preservation (PresDB'07)* (Edinburgh, March 2007).
- [4] CONSULTATIVE COMMITTEE FOR SPACE DATA SYSTEMS. *Reference Model for an Open Archival Information System (OAIS)*. CCSDS 650.0-B-1, 2002.
- [5] DÍAZ, L. M., WÜSTNER, E., AND BUXMANN, P. Inter-organizational document exchange: Facing the conversion problem with XML. In *SAC '02: Proceedings of the 2002 ACM Symposium on Applied Computing* (NY, 2002), ACM, pp. 1043–1047.
- [6] DIGITAL PRESERVATION TESTBED PROJECT. XML and digital preservation. Tech. rep., 2002.
- [7] FERREIRA, M., BAPTISTA, A. A., AND RAMALHO, J. C. An intelligent decision support system for digital preservation. *International Journal on Digital Libraries* 6, 4 (July 2007), 295–304.
- [8] GILLILAND-SWETLAND, A., AND EPPARD, P. Preserving the authenticity of contingent digital objects: The InterPARES project. *D-Lib Magazine* 6, 7/8 (July–August 2000).
- [9] ISO. *Document management - Electronic document file format for long-term preservation - Part 1: Use of PDF 1.4 (PDF/A) ISO/CD 19005-1*, 2004.
- [10] ISO. *Information technology - Computer graphics and image processing - Portable Network Graphics (PNG): Functional specification (ISO/IEC 15948:2004)*, 2004.
- [11] ISO. *Information technology - Open Document Format for Office Applications (OpenDocument) v1.0 ISO/IEC 26300:2006*, 2006.
- [12] LAWRENCE, G. W., KEHOE, W. R., RIEGER, O. Y., WALTERS, W. H., AND KENNEY, A. R. Risk management of digital information: A file format investigation. CLIR Report 93, Council on Library and Information Resources, June 2000.
- [13] RAMALHO, J. C., FERREIRA, M., FARIA, L., AND CASTRO, R. Relational database preservation through XML modelling. In *Proceedings of Extreme Markup Languages 2007* (Montréal, Québec, August 2007).
- [14] ROTHENBERG, J. *Avoiding Technological Quicksand: Finding a Viable Technical Foundation for Digital Preservation*. Council on Library and Information Resources, January 1999.
- [15] ROTHENBERG, J., AND BIKSON, T. Carrying authentic, understandable and usable digital records through time. Tech. rep., Report to the Dutch National Archives and Ministry of the Interior, 1999.
- [16] STRODL, S., BECKER, C., NEUMAYER, R., AND RAUBER, A. How to choose a digital preservation strategy: Evaluating a preservation planning procedure. In *Proceedings of the 7th ACM IEEE Joint Conference on Digital Libraries (JCDL'07)* (June 2007), pp. 29–38.
- [17] THE 100 YEAR ARCHIVE TASK FORCE. The 100 year archive requirements survey. <http://www.snia-dmf.org/100year/>, 2007.
- [18] THE LIBRARY OF CONGRESS. Preferences in summary for textual content. Website, accessed August 2007. http://www.digitalpreservation.gov/formats/content/text_preferences.sht%ml.
- [19] UNESCO. UNESCO charter on the preservation of digital heritage. Adopted at the 32nd session of the General Conference of UNESCO., October 17, 2003.