

Using MIBIB_TE_X to Populate Open Archives*

Jean-Michel HUFFLEN

LIFC (EA CNRS 6942)

University of Franche-Comté

16, route de Gray

25030 BESANÇON CEDEX

FRANCE

jmhufflen@lifc.univ-fcomte.fr

<http://lifc.univ-fcomte.fr/home/~jmhufflen>

Abstract

After a brief recall of the notion of open archive, we recall that in France, an ‘official’ open archive’s repository is used by the Government to evaluate public laboratories and study how to budget them, so this repository should be updated as soon as laboratory members’ papers are accepted. We show how we use MIBIB_TE_X to update both bibliography files and this repository.

Keywords BIB_TE_X, MIBIB_TE_X, XML, XSLT, open archive, HAL.

Streszczenie

Po krótkim przypomnieniu znaczenia pojęcia „Open Archives“ zwrócimy uwagę na to, że we Francji rząd używa „oficjalnego“ repozytorium Open Archives do oceny publicznych instytucji naukowych i na tej podstawie przyznaje im budżety. Tak więc repozytorium powinno być aktualizowane bezpośrednio po zaakceptowaniu publikacji pracownika instytucji. Pokażemy, jak MIBIB_TE_X jest używany do aktualizowania zarówno plików bibliograficznych jak i repozytorium.

Słowa kluczowe BIB_TE_X, MIBIB_TE_X, XML, XSLT, archiwa otwarte, HAL.

0 Introduction

For a decade, Digital Edition and Open Archives have deeply changed the ways of publishing books and journals, especially about research works. Now, such a work can be immediately available as soon as it is validated. This point may be crucial, in particular within domains like medicine: a new treatment can be immediately known as soon as it has been proven efficient. Readers interested in the *modus operandi* of sites devoted to open access archives can consult the manual of the EPrints program [2], frequently used to manage such sites. Aspects related to law and organisation are described — in French — in [1].

As another ‘official’ point related to using open archives in France, the AERES¹, the governmental agency in charge of public research decided to grade French institutions for research by looking into the

HAL² site of open archive. HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research papers — whether they are published or not — and for PhD dissertations.

Since 1996, our institution, the LIFC³, has managed bibliography data bases⁴ as .bib files usable by the BIB_TE_X bibliography processor [8]. The challenge was to enter as many bibliographical entries as possible into HAL, from our .bib files, with particular attention to recent entries. In Section 1, we make precise the specification of this task, then Section 2 explains how we proceed. Reading this article only requires basic knowledge about BIB_TE_X, MIBIB_TE_X⁵ and XML⁶.

² *Hyper-Article en Ligne*, that is, ‘hyper-article on-line’. This site is located at <http://hal.archives-ouvertes.fr/>.

³ *Laboratoire d’Informatique de l’université de Franche-Comté*, that is, ‘Laboratory of Computer Science of the University of Franche-Comté’.

⁴ Located at our institution’s Web site: <http://lifc.univ-fcomte.fr/lifc/publications/>.

⁵ MultiLingual BIB_TE_X.

⁶ eXtensible Markup Language. Readers interested in an introductory book to this formalism can refer to [9].

* Title in Polish: *Zastosowanie MIBIB_TE_Xa do zasilania danymi Open Archives*.

¹ *Agence d’Évaluation de la Recherche et de l’Enseignement Supérieur*, that is, ‘agency evaluating research and university courses’.

```

<?xml version="1.0" encoding="ISO-8859-1"?>
<HAL xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:noNamespaceSchemaLocation="hal.xsd">
  <CONNEXION LOGIN="..." PASSWORD="..."/>
  <ARTICLE_RECENT>
    <META_ART>
      <REFERENCE_BIBLIO>
        <COMM_ACT>
          <TITOUV>Intelligent Tutoring System conference</TITOUV>
          <AUDIENCE>2</AUDIENCE>      <!-- International scope. -->
          <DATEPUB>2002</DATEPUB>
          <PAGE>31-40</PAGE>
          <EDCOM>Springer-Verlag</EDCOM>
          <EDSCI>S. A. Cerri and G. Guardères and F. Paraguaçu</EDSCI>
          <SERIE>LNCS</SERIE>
          <TITCONF>Intelligent Tutoring System conference</TITCONF>
          <DATECONF>2002</DATECONF>
          <VILLE>Paris</VILLE>
          <PAYS>FR</PAYS>
        </COMM_ACT>
      </REFERENCE_BIBLIO>
      <DOMAIN>...</DOMAIN>
      <TITLE>
        Social Network Analysis Used for Modelling Collaboration in Distance Learning Groups
      </TITLE>
      <ABSTRACT>
        We describe a situation of distance learning based on collaborative production occurring
        within groups over a significant time span. For such a situation, we suggest giving
        priority to monitoring and not to guiding systems. [...]
      </ABSTRACT>
      <LANGUE>en</LANGUE>
      <REF_INTERNE>rc02:ip</REF_INTERNE>
    </META_ART>
    <AUTLAB>
      <AUTEURS>
        <AUTEUR><LABIDS>3199</LABIDS><NOM>Reffay</NOM><PRENOM>Christophe</PRENOM></AUTEUR>
        <AUTEUR><LABIDS>3199</LABIDS><NOM>Chanier</NOM><PRENOM>Thierry</PRENOM></AUTEUR>
      </AUTEURS>
      <LABORATOIRES>      <!-- The institutions known by HAL are given a LABID identifier. -->
        <LABORATOIRE LABID="3199"/>
      </LABORATOIRES>
    </AUTLAB>
    <DEPOTS>
      <RIGHT>1</RIGHT>
      <FULLTEXT>
        <DEPOT NOM="http://lifc.univ-fcomte.fr/RECHERCHE/P7/pub/ITS02/reffayIts.pdf" FORMAT="PDF"/>
      </FULLTEXT>
    </DEPOTS>
  </ARTICLE_RECENT>
</HAL>

```

Figure 1: Metadata handled by HAL.

1 Tasks' specification

When an end user enters the information related to a new bibliographical item into HAL, the address of a file containing the corresponding work must

be provided — in which case this file is downloaded by HAL if entering this item succeeds — unless there is some public access — a URL⁷ — to this resource.

⁷ Uniform Resource Locator.

This file may be a PDF⁸ or PS⁹ file, other formats such as L_{ATEX} or RTF¹⁰ are accepted¹¹. Some metadata — e.g., the work’s title and its co-authors — must be provided. This operation may be performed by means of a graphical interface, or by sending an XML file grouping all these metadata. Such a file must be conformant to a *schema* — expressed using XML Schema [12] — given in [3]. Fig. 1 is an example of such an XML file suitable for HAL.

At first glance, most of these metadata could be deduced from information stored in the entries of .bib files. However, what is missing can be difficult to express in BIB_{TEX}. For example, an author must be affiliated to an institution. More precisely, HAL distinguishes a *recent* article — *all* the authors must be affiliated to an institution — and a *retro* article — the rule only holds on for one of the co-authors. In addition, authors may have written an article when they were members of our institution, and may have written another article when they were not, in which case this article is written in collaboration with other members of our institution.

Other information should be supplied for articles included in journals or presented at conferences: does the journal have an editorial board? does the conference provide proceedings¹²? Another information related to *scope* must be supplied: is a journal’s scope or a conference’s *national*, that is, limited to one country, or *international*¹³?

Last but not at least, the ADDRESS field of a BIB_{TEX} entry does not have a precise format. For a conference, it is mandatory for HAL to make precise the country, the town being optional. To sum up, most of information needed by HAL is included into .bib files, but converting .bib files into files suitable for HAL is not immediate.

2 Using MIBIB_{TEX} and xSLT

The .bib files used for the laboratory’s bibliography were already split into ‘national’ conferences, international conferences, ‘national’ journals, etc. First we asked authors to add either an URL field for works publicly available or a PDF field for other works, in

which case this PDF field gives the local address of the document. An abstract should be given, too, but the entries of our .bib files already included an ABSTRACT field¹⁴.

Let us recall that we developed MIBIB_{TEX} [4], a re-implementation of BIB_{TEX} with particular focus on multilingual features. When MIBIB_{TEX} parses a .bib file, it results in an XML tree. Besides the nbst¹⁵ language, used by the bibliography styles taking as much advantage as possible of MIBIB_{TEX}’s features is close to xSLT¹⁶, so it should be suitable for deriving XML texts using HAL format. Unfortunately, there was something missing in nbst: reading additional information from another XML file, as performed by the doc function in XQuery [6]. In addition, nbst is close to xSLT 1.0 [10], whereas some features of xSLT 2.0 [11] — using sequences of values [5] — were interesting for our task.

That is why we put a 2-step process into action: first, a .bib file is converted into an XML tree by means of MIBIB_{TEX}; second, this XML tree is processed by the Saxon program [7], Michael Kay’s processor of xSLT 2.0. By comparison with ‘original’ MIBIB_{TEX}, some additional checks are performed about the ADDRESS field. Allowed values are:

- $\langle \text{country name} \rangle$
- $\langle \text{town name} \rangle, \langle \text{country name} \rangle$
- $\langle \text{town name} \rangle, \dots, \langle \text{country name} \rangle$

where $\langle \text{country name} \rangle$ can be given in French or English. The second step uses three additional XML files:

- the first file groups all the members of the laboratory, making precise the corresponding period of time:

```
<person-db>
  <first>Jean-Michel</first>
  <last>Hufflen</last>
  <period>
    <from year="1994" month="sep"/>
    <!-- Membership's end is specified by a to
         element.
    -->
  </period>
</person-db>
```

- the second file is used for the affiliation of external people appearing within .bib files;
- the third file is used to specify the scope of journals and conferences:

⁸ Portable Document Format.

⁹ PostScript.

¹⁰ Rich Text Format.

¹¹ If no such file can be provided, the article can be only known by HAL as a *note*.

¹² For example, BACHO_{TEX} conferences provide proceedings. As a counter-example originating from T_{EX}’s conferences, the meetings of the German-speaking group DANTE (*Deutschsprachige ANwendervereinigung T_{EX}, e.V.*) do not.

¹³ In fact, this scope information may be left unspecified. But let us not forget that these informations will be used to evaluate researchers: it is preferable for them to make precise that the scope is international, whenever accurate.

¹⁴ Let us recall that when BIB_{TEX} builds a ‘References’ section, the fields not involved in this process, that is, unknown by the bibliography style used, are ignored.

¹⁵ New Bibliography Styles.

¹⁶ eXtensible Stylesheet Language Transformations.

```

<journal>
  <name>Biuletyn GUST</name>
  <scope level="3"/>
  <!-- ('3' is for 'national'.) -->
</journal>

```

This last part required huge work. In fact, .bib files were populated ‘manually’, and the conventions for journals’ and conferences’ names were not uniform. A first pass allowed us to collect all the names, then we decided about a ‘canonical’ name for each journal and conference in a second step. Finally, we had to ask people involved in these journals and conferences whenever we did not know the corresponding scope.

3 Conclusion

The final action is planned for the middle of May 2010. Presently, tests are promising. We decided to enter recent articles as many as possible, that is, if we can find all the affiliations. Otherwise, that is a ‘retro’ article, provided that it has been written prior to 2006. Whilst we are performing the last tests, we encourage people to put the PDF files of their articles on the laboratory’s central machine.

4 Acknowledgements

Some parts of this work have been done in collaboration with Fabrice Bouquet, Pierre-Alain Masson and Jean-Michel Caricand. Thanks to them since they trusted and helped me since the first steps of this work. Many thanks to Jerzy B. Ludwiczowski, who has translated the abstract and keywords in Polish.

References

- [1] Thierry CHANIER : *Archives ouvertes et publication scientifique. Comment mettre en place l'accès libre aux résultats de la recherche?* L'Harmattan. Décembre 2004.
- [2] *EPrints — Digital Repository Software*. March 2010. <http://www.eprints.org/software/>.
- [3] *Documentation Import XML-HAL*. Novembre 2007. <http://hal.archives-ouvertes.fr/>.
- [4] Jean-Michel HUFFLEN: “MIBIB_TE_X's Version 1.3”. *TUGboat*, Vol. 24, no. 2, pp. 249–262. July 2003.
- [5] Jean-Michel HUFFLEN: “XSLT 2.0 vs XSLT 1.0”. In: *Proc. BachoT_EX 2008 Conference*, pp. 67–77. April 2008.
- [6] Jean-Michel HUFFLEN: “Introduction to XQuery”. In: Tomasz PRZECHLEWSKI, Karl BERRY and Jerzy B. LUDWICHOWSKI, eds., *Proc. BachoT_EX 2009 Conference*, pp. 17–25. April 2009.
- [7] Michael H. KAY: *Saxon. The XSLT and XQuery Processor*. March 2009. <http://saxon.sourceforge.net>.
- [8] Oren PATASHNIK: *BIB_TE_Xing*. February 1988. Part of the BIB_TE_X distribution.
- [9] Erik T. RAY: *Learning XML*. O'Reilly & Associates, Inc. January 2001.
- [10] W3C: *XSL Transformations (XSLT). Version 1.0*. W3C Recommendation. Edited by James Clark. November 1999. <http://www.w3.org/TR/1999/REC-xslt-19991116>.
- [11] W3C: *XSL Transformations (XSLT). Version 2.0*. W3C Recommendation. Edited by Michael H. Kay. January 2007. <http://www.w3.org/TR/2007/WD-xslt20-20070123>.
- [12] W3C: *XML Schema*. December 2008. <http://www.w3.org/XML/Schema>.