

AcceSciTech: A Global Approach to Make Scientific and Technical Literature Accessible

Alex Bernier¹ and Dominique Burger²

¹ Association BrailleNet
9, quai Saint-Bernard
75252 Paris Cedex 5, France
alex.bernier@upmc.fr

² INSERM UMRS_968, Université Pierre et Marie Curie
9, quai Saint-Bernard
75252 Paris Cedex 5, France
dominique.burger@upmc.fr

Abstract. In this paper, we introduce AcceSciTech, a research and development project which addresses the challenges faced in providing access to scientific and technical literature for the visually impaired and, more broadly, for those who are not able to read conventional print. Based on XML formats, this project aims to develop a coherent set of tools to produce, edit, deliver and to render complex documents accessible to impaired people.

Keywords: accessible publishing, Braille, DAISY, ebooks, EPUB, PEF, print-disabled persons, scientific documents, workflow, XML.

1 Introduction

In this paper, we introduce AcceSciTech and its first results, a three years research and development project funded by the French national research agency started in September 2012 involving three partners : Université Pierre et Marie Curie (Paris), BrailleNet and Eurobraille.

The AcceSciTech project addresses the challenges faced in providing access to scientific and technical literature for the visually impaired and, more broadly, for those who are not able to read conventional print. The project aims to develop coherent tools to create, edit, deliver, to render and to interact with complex documents. In this context, scientific and technical literature refers to all written text (articles, handouts, slides, books, etc.) that include complex elements such as mathematical formulae, graphs, charts, diagrams, notes, etc. These documents can cover both "exact sciences" and social sciences, and be of university level or aimed at the general public.

Visually impaired readers access text in a number of ways: by touch through the use of Braille (either as static text embossed on paper or dynamic text generated by a Braille display device); by listening to text read by a narrator or a speech synthesiser; or by sight through large print. These different means of

accessing information are not exclusive: a student, for example, may find it useful to learn to read a sentence in Braille while listening to the same sentence. XML-based formats (mainly DAISY [1] and EPUB [2]) are used in this project because they allow to separate the content and the presentation of a document. Also these formats provide a general framework to render documents using various means (touch, listening or sight). EPUB and DAISY also provide mechanisms to enable, for example, the user to navigate from section to section or from chapter to chapter with ease, or to opt whether or not to read certain elements of a book such as footnotes or page numbers.

In order to introduce the problems associated with accessing scientific and technical literature, the following questions may be asked:

- How can one efficiently produce highly structured documents that originate from various sources (scans, PDF or LaTeX files provided by a publisher, etc.)?
- For documents containing many images, in which publishing environment should experts be invited to provide descriptions as and when the user needs them?
- How does one translate a mathematical formula into spoken words and/or Braille?
- How can a young visually impaired pupil with limited computer skills access a mathematics textbook with ease?
- How can one read and navigate with ease in a complex document using a Braille terminal or new touchscreen-based devices?

In many countries, the law grants approved organisations the right to ask publishers to provide source files used in the production of printed works subject to copyright, so that they may be adapted and distributed to people with disabilities. In order to take advantage of this favourable legal environment, an optimal technical environment that will facilitate and increase the production and distribution of accessible texts must be created. The AcceSciTech project proposes to do this, in the context of scientific and technical literature.

In the first sections sections, we will focus on production steps required to create accessible documents: after introducing the formats on which are operating the production chain in section 2, we will describe pre-processing steps required to detect and enhance the structure of a document in section 3, the various possibilities offered to edit accessible documents in section 4, and, in section 5, the softwares used to produce final documents, ready to be delivered to the end-user. In section 6, we will introduce the Helene platform, the core of an online library service for visually impaired people. Two new modules added to this platform will be described, allowing to dynamically generate and to distribute documents via Web services. Last, in section 7, we will highlight the features required in a rendering system based on touch to ease the reading of complex documents.

2 Formats

The first thing we did to find efficient solutions to answer the question “How can one efficiently produce highly structured documents that originate from various

sources?” was to find common formats able to take the variety of representable documents into account. The different tools developed in the AcceSciTech project mainly rely on two formats: ZedAI for the production and EPUB 3 for the distribution. We extensively described these two formats and why they can be candidates to be part of high-quality document production chain in a previous paper ([3]). Let us introduce them briefly:

ZedAI (Z39.98 Authoring and Interchange [4]) is designed by the DAISY consortium to replace the XML DTBook format defined in the previous DAISY specifications (Z39.86-2005 [1]). On the DAISY consortium website, ZedAI is presented as ”a specification that defines an XML-based framework with which content producers can represent various types of information in an extensible, standards-compliant way, suitable for the transformation into multiple output formats” [5]. The important point here is that ZedAI is not a format defined by a grammar like XML DTBook, but ZedAI is a framework to define grammars.

EPUB 3 is a standard distribution format (i.e. for the end-user of the content) for digital publications and documents. An EPUB 3 document is basically a single compressed file (the container) embedding HTML5 files, CSS, images, metadata and other resources. The accessibility related features (synchronisation mechanisms, skippable and escapable contents) available in DAISY 2005 have been fully integrated into this new version of the EPUB format.

Both ZedAI and EPUB 3 embed vocabulary association mechanisms based on RDF. This means that it is possible to transfer semantic inflections of elements from a ZedAI document into EPUB 3 using common concepts shared by ZedAI and EPUB 3 or linking ZedAI-specific vocabularies to the EPUB 3 document when concepts defined in these vocabularies are not already available in EPUB.

ZedAI and EPUB 3 are also good candidates to handle complex documents because their specifications allow to use MathML for representing formula and SVG for the graphical contents.

3 Pre-processing Steps

Because it is rich and extensible, we have selected ZedAI to be the core format of the production tool chain we design in the AcceSciTech project. The main requirement for this production chain is to be able to manage a wide range of source documents formats. The pre-processing steps are those done to convert a document provided at the input of the chain into a ZedAI file. They are as follow:

Document acquisition: the objective of this step is to obtain an electronic version of the document by digitising the printed version or by getting a source file from the publisher/author. If a source file is provided, it can be in a wide range of formats: LaTeX, PDF, Microsoft Word, XML (using various DTD), etc. This step is implemented thanks to a high speed document scanner for

book digitisation and using the PLATON platform provided by the French national library which allows accredited organisations to request source files corresponding to printed books. BrailleNet, as one of the accredited organisations, has noted two major difficulties using PLATON in a real production context: it hard to predict when the files will be uploaded by the publisher and once uploaded, in which format. Legally, publishers have to upload source files corresponding to a printed book no later than two months after the document has been requested. But in practice, some publishers do not respect the delay, and when they do, their responsiveness is quite variable. Regarding the formats, approximately 70% of the uploaded files are in PDF, 25% in XML, 4% in Microsoft Word and the last 1% in various other formats (LaTeX, InDesign, etc).

Conversion into readable format: if the document has been digitised or provided using an image-based format, it should be processed by an Optical Characters Recognition (OCR) software to be converted into readable text format. This step is implemented using the ABBYY FineReader SDK [7] for "simple" documents and the InftyReader [8] for documents with math.

Structure detection: to enable users to easily navigate inside a document, its structure should be known (chapter, section, footnotes, etc). Nowadays, some OCR tools (like FineReader) are able to detect some structural elements of documents like headings, footnotes, tables, etc. Unfortunately, this kind of detection is not efficient enough with complex layout or with complex document structure. So, an additional tool is currently being developed, to improve the results of structure detection. This tool is using the ALTO format as input and produce an annotated version as output, based on matching with document description models.

Conversion into ZedAI: finally, the resulting documents of previous processes should be converted into the core format. We are developing various converters to manage the input formats (annotated-ALTO, XML from publishers), integrated in one global tool. Some experimentation is also done to test existing converters from LaTeX: TRALICS [9] and TeX4ht [10].

4 Edition Tools

Because creating fully accessible documents will still require manual work (for example to describe images), AcceSciTech also addresses issues related to publishing, by providing various possibilities to edit document. We are working on three solutions:

A collaborative environment with an user interface dedicated to adapting books.

This environment is based on a Wiki platform allowing users to work simultaneously on the same book divided into pages. The Wiki supports XML DTBook and ZedAI is added. Production management tools are also currently developed to improve the global efficiency of the platform, making possible to easily follow what kind of work has been done on a book (proof-reading, structure improvement, image description, etc). Support for math

formula is included using the TeX language. The Wiki is designed to adapt educational material which can be partially adapted by several teachers.

Oxygen with a ZedAI module: Oxygen [11] is an XML editor. Natively, it provides WYSIWYG features to edit DocBook, XHTML and TEI files. It can be extended to support new formats. For example, MathFlow is a plug-in which allows users to edit equations in the MathML format. As part of the AcceSciTech project, a module to handle ZedAI in Oxygen is under development. This solution is targeting users with some knowledge of XML and is well designed to be integrated in a production workflow, to adapt whole documents.

ODT2DAISY [12] is a plug-in for OpenOffice.org/LibreOffice.org to convert ODF-text files into DAISY. Currently, ODT2DAISY supports the XML DT-Book formats. We plan to modify it to handle natively the ZedAI format. Math equations can be created with an equation editor and are internally stored using the MathML format. This edition solution will be helpful for the users who are not familiar with XML and who would like to use a "standard" word-processor to create specific documents.

5 Production of Documents in Distribution Formats

The previous steps aimed to create a structured and rich ZedAI document. Because ZedAI is an authoring format, such a document should be converted in a format more appropriate for distribution purposes like EPUB 3 or Portable Embossable Format (PEF) [20]. This task is done using the DAISY Pipeline [6] 2, an "XML-centric open source cross-platform framework for the automated processing of various digital formats" [13] developed by the DAISY consortium. It uses XProc [14], an XML Pipeline language designed by W3C to create conversion workflows which can be based on the ZedAI format. Two working group have been set up: one dedicated to the production of Braille document and the other to the audio generated by Text-To-Speech (TTS) documents.

The AcceSciTech project will contribute to the developments of the DAISY Pipeline 2. Our efforts will focus on the two following aspects: production of complex Braille documents and efficient production using TTS.

Our main goal regarding the Braille production module is to support the French Braille code specification (both contracted and non-contracted codes) and to be able to produce Braille math from MathML equations. To include MathML support in the Braille module, three open-source softwares are evaluated: UMCL [15], LibLouis [16] + LibLouisXML [17] and NAT Braille [18]. Our evaluation covers efficiency (processing speed), input format and output code coverage. Depending of the results, one of these tools will be selected to be integrated into the DAISY Pipeline Braille module, or it will be decided to improve one of them or develop a new one.

The possibility to produce audio documents from XML DTBook using TTS is already provided by the Narrator module available in the first version of the DAISY Pipeline. We are participating in in the design of a new version of this

module to be integrated in the DAISY Pipeline 2. Our main objectives are the following:

- Providing support for the ZedAI Format
- Optimising the production process to take advantage of multi-CPU computers
- Handling math by contributing a library to convert MathML to SSML (a language designed to be processed by a TTS software [19])

6 Delivering Documents

After the adaptation and production processes, accessible documents are ready to be delivered to end-users. We use an online library (based on the Helene platform provided by BrailleNet), dedicated to visually impaired people.

The online library is based on Koha, an open-source library management system we have extended to support features required to manage digital documents [21]. More complex is the delivered document, more useful is it for the end-user to be able to customise it corresponding to his/her specific needs. For example, it can be helpful to get a manual with the main text pronounced by a voice, footnotes by a different one, and math equations by a third. The user should be able to select voices according to his/her preferences.

It is already possible with the Helene platform to customise books on-demand, to deliver PDF version (for large-print), DAISY text-only and XHTML. Our objective is now to add this feature for audio books produced by TTS. The main issue here is to properly manage resources allowed to produce the requested books: scheduling, monitoring and caching mechanisms are implemented to make the system efficient and robust to deliver customised books to a large number of users.

Currently, the online library is accessible via Internet using a Web browser: the user search a book in the catalogue, download it and should copy it manually on its reading system. These operations require to be familiar with computers. It can be helpful, especially for young children, to make this process easier: this is the purpose of DAISY Online Delivery Protocol (DODP) [22] which is being implemented on the Helene platform. DODP specifies Web services to allow direct communication via Internet between a connected reading system and an online library management system. Two contents elections methods are proposed: "Out-of-band" and "Browse Content". The first is well-suited to deliver periodicals contents or books selected by a librarian: basically, the user has only to start his/her reading system and the list of contents available will be presented automatically. The second method allow the user to search the catalogue (using the keypad of the reading system) and to browse it by categories thanks to a dynamic menu mechanism.

7 Rendering Complex Documents and New Interactions Possibilities

Last, a multi-modal reading software is developed to render complex documents in Braille, audio and on a regular screen. It should have the following features:

Support DAISY (2.02 and 3.0) and EPUB 3 formats
 Fine synchronisation of audio, Braille and visual outputs
 MathML rendering
 DAISY Online client

This software is designed to be multi-platform with constraints of embedded systems in mind: our objective is to integrate it in a portable Braille notebook (Esytime by Eurobraille) and into iOS and Android platforms with the goal to create new prospects for improving the reading efficiency of visually impaired users. Indeed, Esytime embeds a video-card to be plugged on a screen, a text-to-speech software and a new hardware mechanism able to detect fingers movement of the user on the Braille display thanks to optical sensors fixed on every Braille cell of the device. The sensors allow to recognise sequences of movement made by the user and to process these sequences to execute appropriate commands. Almost all the iOS and Android phones/tablets use touchscreen as the primary input device. To take advantage of these new possibilities, we are designing gesture libraries (in one and two dimensions) to experiment new interactions methods aiming to improve reading speed and to ease the navigation inside complex documents.

8 Conclusion

One of the main concerns of AcceSciTech is to reduce the gap between the number of accessible documents for print-disabled people and the number of documents available into the mainstream publishing market. This gap is significantly greater for scientific and technical literature because of the high costs involved to produce this kind of accessible material. Based on standard and extensible formats, AcceSciTech aims to create tools and to participate to the development of existing open-source projects, to contribute to the global effort required to make complex documents more accessible.

Acknowledgement. This work is supported by the Agence Nationale de la Recherche (AcceSciTech project: 2012 CORD 008 02). We are grateful to the Alcatel-Lucent company for their support and cooperation.

References

1. National Information Standards Organization (NISO). ANSI/NISO Z39.86-2005. Revision of ANSI/NISO Z39.86-2002. Specifications for the Digital Talking Book, <http://www.daisy.org/z3986/2005/Z3986-2005.html>
2. International Digital Publishing Forum. EPUB 3: <http://idpf.org/epub/30>
3. Bernier, A., Burger, D.: XML-Based Formats and Tools to Produce Braille Documents. In: Miesenberger, K., Karshmer, A., Penaz, P., Zagler, W. (eds.) ICCHP 2012, Part I. LNCS, vol. 7382, pp. 500–506. Springer, Heidelberg (2012)

4. National Information Standards Organization (NISO). ANSI/NISO Z39.98-2012. Authoring and Interchange Framework for Adaptive XML Publishing Specification, <http://www.daisy.org/z3998/2012/z3998-2012.html>
5. DAISY Consortium: ZedAI Introduction, <http://www.daisy.org/zw/ZedAIIntroduction>
6. DAISY Pipeline, A framework for document-related pipelined transformations, for the DAISY Consortium community: <http://code.google.com/p/daisy-pipeline/>
7. ABBYY FineReader SDK: <http://www.abbyeu.com/sdk/>
8. Fukuda, R., Ohtake, N., et al.: Optical recognition and Braille transcription of mathematical documents. In: ICCHP, pp. 711–718 (2000)
9. TRALICS: a LaTeX to XML translator: <http://www-sop.inria.fr/marelle/tralics/>
10. TeX4ht, <http://tug.org/tex4ht>
11. Oxygen XML Editor, <http://www.oxygenxml.com/>
12. Strobbe, C., Engelen, J., Spiewak, V.: Generating DAISY Books from OpenOffice.org. In: Miesenberger, K., Klaus, J., Zagler, W., Karshmer, A. (eds.) ICCHP 2010, Part 1. LNCS, vol. 6179, pp. 5–11. Springer, Heidelberg (2010)
13. Deltour, R.: XProc at the heart of an ebook production framework. XLL Prague (2013)
14. XProc: An XML Pipeline Language, <http://www.w3.org/TR/xproc/>
15. Archambault, D., Fitzpatrick, D., Gupta, G., Karshmer, A.I., Miesenberger, K., Pontelli, E.: Towards a Universal Maths Conversion Library. In: Miesenberger, K., Klaus, J., Zagler, W.L., Burger, D. (eds.) ICCHP 2004. LNCS, vol. 3118, pp. 664–669. Springer, Heidelberg (2004)
16. Liblouis, A Braille translation and back-translation library: <http://code.google.com/p/liblouis/>
17. LiblouisXML - A Braille transcription software for XML documents: <http://code.google.com/p/liblouisxml/>
18. NAT Braille: A free universal Braille translator, <http://natbraille.free.fr/>
19. Speech Synthesis Markup Language (SSML) Version 1.0., <http://www.w3.org/TR/speech-synthesis/>
20. Håakansson, J.: PEF 1.0 - Portable Embosser Format Public draft, <http://files.pef-format.org/specifications/pef-2008-1/pefspecification.html>
21. Bernier, A., Burger, D.: Helene, an Open-Source and DAISY Compliant Digital Library Management System. In: CSUN Conference, March 14-19 (2011)
22. DAISY Online Delivery Protocol: <http://www.daisy.org/projects/daisy-online-delivery/drafts/20100402/do-spec-20100402.html>