



**HAL**  
open science

## New features in ALTO XML

Jean-Philippe Moreux, Frederick Zarnedt

► **To cite this version:**

Jean-Philippe Moreux, Frederick Zarnedt. New features in ALTO XML. IFLA World Library and Information Congress, 2014, Lyon, France. hal-04251649

**HAL Id: hal-04251649**

**<https://bnf.hal.science/hal-04251649>**

Submitted on 30 Oct 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

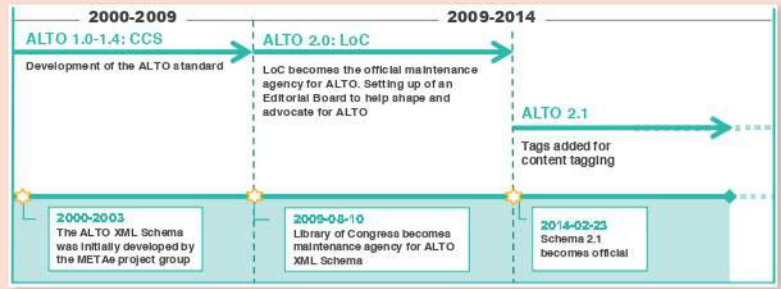
L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - ShareAlike 4.0 International License

# New features in ALTO XML: Tags for layout, structure, roles, named entities...

ALTO (Analyzed Layout and Text Object) is a widely used XML based open standard for describing the layout and content of text images. ALTO is most often used with METS XML and TIFF or JPEG2000 images to create digital surrogates for books, journals, magazines and newspapers. ALTO XML is implemented in text digitisation projects around the world.



ALTO is an open standard, first administered and maintained by Content Conversion Specialists (CCS), one of the MetaE project members. In 2009, CCS transferred its administration to the ALTO Editorial Board whose members are from libraries and private organizations around the world. Library of Congress is the maintenance agency for ALTO Schema.



## 3 Use Cases for Tags in the Digital Library Community (Institutions, Service Providers, Research Projects)

### Logical Labelling of Structural Elements

The ALTO format captures the layout and the full text of a page. Although this OCRing provides full text retrieval, that retrieval can benefit from marking additional structural elements such as:

- running titles,
- page numbers
- captions
- marginalia, footnote

Unlike entities which record the physical and logical structure of an entire document (in a format such as METS), these elements are specific to the page.



Gallica.bnf.fr / Bibliothèque nationale de France

### Layout/Content Tagging

Before the actual character recognition starts, most OCR software analyses the page structure to identify regions. Modern OCR software supports numerous region types that ALTO can record: tables, graphics, music scores, maps, etc. These regions can be labelled either by the software or manually by operators. Like labelling elements, labelling regions helps improve retrieval.



Gallica.bnf.fr / Bibliothèque nationale de France

### Named Entities & Roles

Named-entity recognition (NER) is a subtask of information extraction that seeks to locate and classify discrete textual elements into predefined categories (names of persons, organizations, locations, expressions of times, quantities, etc.). Role tagging describes how entities were involved in the content creation: author, editor, publisher, distributor, etc.



Chronicling America: Historic American Newspapers, Lib. of Congress

## Implementers of ALTO 2.1

### European Newspapers project

Named Entities Recognition: EN project has implemented an ALTO 2.1 output for its NER task. Persons, locations and organizations are recognized from hundred thousands of European newspapers.

### BnF Mass Digitization program

Markup of important words: The next BnF large-scale digitization program makes use of NER tagging to identify meaningful words in OCR results and to concentrate the post-correction task on these specific words.

