

NAME

ltx2crossrefxml.pl – create XML files for submitting to crossref.org

SYNOPSIS

```
ltx2crossrefxml [--debug] [-c config_file] [-o output_file] [-input-is-xml]
                latex_file1 latex_file2 ...
```

OPTIONS

-c *config_file*

Configuration file. If this file is absent, defaults are used. See below for its format.

-o *output_file*

Output file. If this option is not used, the XML is output to stdout.

-rpi-is-xml

Do not transform author and title input strings, assume they are valid XML.

--debug

Output some progress reports.

The usual `--help` and `--version` options are also supported. Options can begin with either `-` or `--`, and ordered arbitrarily.

DESCRIPTION

For each given *latex_file*, this script reads `.rpi` and (if they exist) `.bbl` and `.aux` files and outputs corresponding XML that can be uploaded to Crossref (<https://crossref.org>). Any extension of *latex_file* is ignored, and *latex_file* itself is not read (and need not even exist).

Each `.rpi` file specifies the metadata for a single article to be uploaded to Crossref (a `journal_article` element in their schema); an example is below. These files are output by the `resphilosophica` package (<https://ctan.org/pkg/resphilosophica>), `aomart` package (<https://ctan.org/pkg/aomart>), the TUGboat publication procedure (<https://tug.org/TUGboat/repository.html>) and other packages. They can also be created by hand or by whatever other method you implement.

Any `.bbl`, `.aux`, and `.bib` files are used for the citation information in the output XML. See the CITATIONS section below.

Unless `--rpi-is-xml` is specified, for all text (authors, title, citations), standard TeX control sequences are replaced with plain text or UTF-8 or eliminated, as appropriate. The `LaTeX::ToUnicode::convert` routine is used for this (<https://ctan.org/pkg/bibtexperl>). Tricky TeX control sequences will almost surely not be handled correctly.

If `--rpi-is-xml` is given, the author and title strings from the `rpi` files are output as-is, assuming they are valid XML; no checking is done.

Citation text from `.bbl` files is always converted from LaTeX to plain text.

This script just writes an XML file. It's up to you to do the uploading to Crossref; for example, you can use their Java tool `crossref-upload-tool.jar` (<https://www.crossref.org/education/member-setup/direct-deposit-xml/https-post>).

For the definition of the Crossref schema currently output by this script, see https://data.crossref.org/reports/help/schema_doc/5.4.0/index.html with additional links and information at <https://www.crossref.org/documentation/schema-library/metadata-deposit-schema-5-4-0/>.

CONFIGURATION FILE FORMAT

The configuration file is read as Perl code. Thus, comment lines starting with # and blank lines are ignored. The other lines are typically assignments in the form (spaces are optional):

```
$variable = value ;
```

Usually the value is a "string" enclosed in ASCII double-quote or single-quote characters, per Perl syntax. The idea is to specify the user-specific and journal-specific values needed for the Crossref upload. The variables which are used are these:

```
$depositorName = "Depositor Name";
$depositorEmail = 'depositor@example.org';
$registrant = 'Registrant'; # required, organization name
$fullTitle = "FULL TITLE"; # required, journal name
$issn = "1234-5678"; # required, ISSN
$abbrevTitle = "ABBR. TTL."; # optional, abbreviated journal name
$coden = "CODEN"; # optional
```

For a given run, all .rpi data read is assumed to belong to the journal that is specified in the configuration file. More precisely, the configuration data is written as a `journal_metadata` element, with given `full_title`, `issn`, etc., and then each .rpi is written as `journal_issue` plus `journal_article` elements.

The configuration file can also define a Perl function `LaTeX_ToUnicode_convert_hook`. If it is defined, it is called at the beginning of the procedure that converts LaTeX text to Unicode, which is done with the `LaTeX::ToUnicode` module, from the `bibtexperl` package (<https://ctan.org/pkg/bibtexperl>). The function must accept one string (the LaTeX text), and return one string (presumably the transformed string). The standard conversions are then applied to the returned string, so the configured function need only handle special cases, such as control sequences particular to the journal at hand. (See TUGboat's `ltx2crossrefxml-tugboat.cfg` for an example.)

The configuration file can also define a hash `BibentryToCrossref` that maps Crossref entry types to BibTeX entry types used in the bibliography processing (see CITATIONS), for example

```
%BibentryToCrossref = ('WEBPAGE' => 'other',
                        'MISC' => 'other');
```

The keys in this hash must be in the upper case, while the entries must be in the lower case.

RPI FILE FORMAT

Here's the (relevant part of the) .rpi file corresponding to the `rpsample.tex` example in the `resphilosophica` package (<https://ctan.org/pkg/resphilosophica>):

```
%authors=Boris Veytsman\and A. U. Th{\o }r\and C. O. R\"espondent
%title=A Sample Paper:\\ \emph {A Template}
%year=2012
%volume=90
%issue=1--2
%startpage=1
%endpage=1
%doi=10.11612/resphil.A31245
%paperUrl=http://borisv.lk.net/paper12
%publicationType=full_text
```

Other lines, some not beginning with %, are ignored (and not shown). For more details on processing, see the code.

The %paperUrl value is what will be associated with the given %doi (output as the resource element). Crossref strongly recommends that the url be for a so-called landing page, and not directly for a pdf (<https://www.crossref.org/education/member-setup/creating-a-landing-page/>). Special case: if the url is not specified, and the journal is *Res Philosophica*, a special-purpose search url using pdcnet.org is returned. Any other journal must always specify this.

The %authors field is split at \and (ignoring whitespace before and after), and output as the contributors element, using sequence="first" for the first listed, sequence="additional" for the remainder. The authors are parsed using BibTeX::Parser::Author (<https://ctan.org/pkg/bibtexperllibs>).

If the %publicationType is not specified, it defaults to full_text, since that has historically been the case; full_text can also be given explicitly. The other values allowed by the Crossref schema are abstract_only and bibliographic_record. Finally, if the value is omit, the publication_type attribute is omitted entirely from the given journal_article element.

Each .rpi must contain information for only one article, but multiple files can be read in a single run. It would not be difficult to support multiple articles in a single .rpi file, but it makes debugging and error correction easier to keep the input to one article per file.

MORE ABOUT AUTHOR NAMES

The three formats for names recognized are (not coincidentally) the same as BibTeX:

```
First von Last
von Last, First
von Last, Jr., First
```

The forms can be freely intermixed within a single %authors line, separated with \and (including the backslash). Commas as name separators are not supported, unlike BibTeX.

In short, you may almost always use the first form; you shouldn't if either there's a Jr part, or the Last part has multiple tokens but there's no von part. See the btxdoc ("BibTeXing" by Oren Patashnik) document for details. The authors are parsed using BibTeX::Parser::Author (<https://ctan.org/pkg/bibtexperllibs>).

In the %authors line of a .rpi file, some secondary directives are recognized, indicated by | characters. Easiest to explain with an example:

```
%authors=|organization|\LaTeX\ Project Team \and Alex Brown|orcid=123
```

Thus: 1) if `|organization|` is specified, the author name will be output as an `organization` contributor, instead of the usual `person_name`, as the Crossref schema requires.

2) If `|orcid=value|` is specified, the *value* is output as an ORCID element for that `person_name`.

These two directives, `|organization|` and `|orcid|` are mutually exclusive, because that's how the Crossref schema defines them. The `=` sign after `orcid` is required, while all spaces after the `orcid` keyword are ignored. Other than that, the ORCID value is output literally. (E.g., the ORCID value of 123 above is clearly invalid, but it would be output anyway, with no warning.)

Extra `|` characters, at the beginning or end of the entire `%authors` string, or doubled in the middle, are accepted and ignored. Whitespace is ignored around all `|` characters.

CITATIONS

Each `.bbl` file corresponding to an input `.rpi` file is read and used to output a `citation_list` element for that `journal_article` in the output XML. If no `.bbl` file exists for a given `.rpi`, no `citation_list` is output for that article.

The `.bbl` files are processed to create the `unstructured_citation` references defined by Crossref, that is, the contents of the citation (each paragraph in the `.bbl`) as a single flat string without markup of any kind, including font changes.

Bibliography text is unconditionally converted from TeX to XML, via the method described above. It is not unusual for the conversion to be incomplete or incorrect. It is up to you to check for this; e.g., if any backslashes or pairs of dollar signs remain in the output, it is most likely an error.

Furthermore, it is assumed that the `.bbl` file contains a sequence of references, each starting with `\bibitem{KEY}` (which itself must be at the beginning of a line, preceded only by whitespace), and the whole bibliography ending with `\end{thebibliography}` (similarly at the beginning of a line). A `.bbl` file not following this format will not produce useful results. The `.bbl` file can be created by hand, or with BibTeX, or any other method, as long as it has this format.

The key attribute for the `citation` element is taken as the *KEY* argument to the `\bibitem` command. The sequential number of the citation (1, 2, ...). The argument to `\bibitem` can be empty (`\bibitem{}`), and the sequence number will be used on its own. Although TeX will not handle empty `\bibitem` keys, it can be convenient when creating a `.bbl` purely for Crossref.

The `.rpi` file is also checked for the bibliography information, in this same format.

Crossref's structured citations are added as follows:

1. If an `.aux` file is present, it is checked for any `\bibdata` commands. The `bib` files in these commands are read, and the information there is used to generate XML entries. The script uses `kpsewhich` to look for the `bib` files, so the usual BibTeX conventions for the search paths are followed.
2. For any citation the corresponding entry in the `bib` file is processed.
3. The Crossref entry type is determined according to the algorithm describe below ("CITATION ENTRY TYPES").

4. The entry fields are used to populate structured citation.

CITATION ENTRY TYPES

The current Crossref schema

<https://data.crossref.org/reports/help/schema_doc/5.4.0/schema_5_4_0.html> defines `type` attribute for a citation. Unfortunately the list of possible types does not fully coincide with the list of BibTeX entry types. Therefore the script uses the following algorithm to determine the Crossref entry type for a citation:

1. If the entry has the field `crossrefentrytype`, it is used.
2. Otherwise if BibTeX entry type appears in the hash `BibentryToCrossref` in the configuration file ("CONFIGURATION FILE FORMAT"), its value is used.
3. Otherwise the default mapping is used. The script knows many BibTeX entry types, and should do a good job in most cases.

EXAMPLES

```
ltx2crossrefxml.pl ../paper1/paper1.tex ../paper2/paper2.tex \  
-o result.xml
```

```
ltx2crossrefxml.pl -c myconfig.cfg paper.tex -o paper.xml
```

AUTHOR

Boris Veytsman <<https://github.com/borisveytsman/crossrefware>>

COPYRIGHT AND LICENSE

Copyright (C) 2012–2026 Boris Veytsman

This is free software. You may redistribute copies of it under the terms of the GNU General Public License (any version) <<https://www.gnu.org/licenses/gpl.html>>. There is NO WARRANTY, to the extent permitted by law.