

Capture metadata from the PDF file (Grobid Integration)

Through the integration with [GROBID Library](#), DSpace-CRIS since 7th March 2018 ([e65e0fa](#)) leverages machine learning technologies for extracting information directly from PDF publications. Metadata such as title, authors, abstract, keywords, identifiers, source, etc. are extracted through parsing of uploaded full-text.

The feature is implemented as a BTE dataloader and accessible in the import for a file section of the New submission

New submission

The screenshot shows a web interface for a 'New submission'. At the top, there are tabs for 'Search Form' and 'Results'. Below are sections for 'Default mode Submission', 'Search for identifier', and 'Upload a file'. The 'Upload a file' section contains instructions: 'Select a file to upload and its type from the drop-down menu. If "Preview Mode" is enabled, the list of the publications in the file will be shown to you to select the one for submission. If it is disabled, all publications will be imported in your MyDSpace page as "Unfinished Submissions" while the first one will go through the submission process.' There are labels for 'Select data type:', 'File:', and 'Collection:'. A dropdown menu is open, showing a list of file formats: 'Please, specify the file format', 'PubMed XML', 'CrossRef XML', 'arXiv XML', 'CINii XML', 'BibTeX', 'Research Information Systems (RIS)', 'EndNote', 'Comma Separated Values (CSV)', 'Tab Separated Values (TSV)', and 'PDF' (which is highlighted in blue). To the right of the dropdown are 'Process' and 'Exit' buttons.

i The PDF file is also automatically attached to the new item in the ORIGINAL bundle if the import is done directly **skipping the preview mode** of the BTE framework.

To enable the feature you need

1. to install grobid as a local accessible webservice. It is not needed to expose it over internet, in fact, it is recommended to keep it visible only to the dspace server to reduce security risks. Please refer to the [Grobid documentation for the installation](#). The integration has been tested against Grobid 0.5.0 but it should work with more recent versions that doesn't break the grobid REST contract
2. to enable the grobid data loader in the bte.xml file, see https://github.com/4Science/DSpace/blob/dspace-5_x_x-cris/dspace/config/spring/api/bte.xml#L136 and https://github.com/4Science/DSpace/blob/dspace-5_x_x-cris/dspace/config/spring/api/bte.xml#L21 (if you want to use it also from the command line import)
3. to configure the grobid server connection in the grobid.cfg file, see https://github.com/4Science/DSpace/blob/dspace-5_x_x-cris/build.properties#L342 (use `http://localhost:8070` for a default grobid installation)

The mapping between the extracted metadata and the dspace metadata is managed in the usual way of the BTE framework, where the mapping between the Grobid metadata and the BTE model is defined in the [Grobid Data loader field map](#) and the final mapping in the [output map](#)