

OpenITI, ver. 2021.2.5

—Release Notes—

Corpus Metadata

The current release metadata is available in the `OpenITI_metadata_2021-2-5.csv` and `OpenITI_metadata_2021-2-5_merged.csv` (merged¹ version) files.

Folder Structure

- **data:** main data folder with *Author* > *Book* > *Versions* structure;
- **metadata**
 - `OpenITI_metadata_2021-2-5.csv`: metadata file, with a row for each text in the corpus, including a row for each part of the multi-part books in the corpus;
 - `OpenITI_metadata_2021-2-5_merged.csv`: metadata file, with one row for each multi-part book in the corpus;
- **release_notes**
 - `OpenITI_release-notes_2021-2-5.pdf`: these release notes;
 - `release-notes_files_2021-2-5.zip`: csv files including the provided list of new changes in the current release (see the description of each file in this release notes).

¹ The corpus contains a number of texts that are too big for GitHub and had to be split into multiple files (currently only two versions of the book *Bihār al-anwār* (1111Majlisi.BiharAnwar)). The metadata file contains statistics on each part of this split text. We also provide a separate metadata file in which the statistics for the separate parts of those books that we have merged. Since the merged metadata for these split files does not refer to an existing file, the `local_path` field for these virtual texts will be “NA”.

Corpus Statistics

Category	Stats
Number of unique titles	6,337
Number of authors	2,636
Number of books (all versions/editions)	10,344

Length of texts (all books)

	Number of words	Number of pages (300 w/p)
Total	2,050,063,811	6,833,547
Min.	47	1
1st Qu.	7,834	27
Median	35,727	120
Mean	198,189	661
3rd Qu.	144,569	482
Max.	11,912,693	39,709

Length of texts (primary books)

	Number of words	Number of pages (300 w/p)
Total	1,009,950,381	3,366,502
Min.	48	1
1st Qu.	6,967	24
Median	29,719	100
Mean	159,399	532
3rd Qu.	119,113	398
Max.	11,912,693	39,709

Annotation statistics

<i>Number of texts with extension .mARkdown</i>	353
<i>Number of texts with extension .completed</i>	718
<i>Number of texts with extension .inProgress</i>	9

Book Ids

The list of the new book ids in this version is available in the **ids.csv** file. It includes the newly added book ids and modified ids. The URI includes the information of the new book (i.e., date, author, and book title).

Modified URIs

List of modified URIs in the current release is available in **modified_uris.csv**. Changes typically affect such fields as year, author, and title. These changes are applied to the entire metadata (book IDs remain unchanged).

Annotation Update

The list of texts that have been structurally annotated or the annotation has changed (can be tracked by the file extensions) since our previous release (version [2021.1.4](#)) is provided in **annotation_update.csv**. This file shows URIs of texts together with their current extension, which is a part of the **local_path** in the metadata file.

For more information on the OpenITI mARkdown and the extensions please see [here](#).

Credits

Current contributors (*alphabetically*):

- Sohail Merchant (*metadata app*)
- Lorenz Nigst (*corpus management; structural annotation*)
- Maxim Romanov (*OpenITI co-PI; conceptual development; mARkdown*)
- Sarah Bowen Savant (*OpenITI co-PI; KITAB Project PI*)

- Masoumeh Seydi (*technical development*)
- Peter Verkinderen (*technical development; preparing new texts for the corpus*)

- Mathew Barber (*structural annotation*)
- Gowaart Van Den Bossche (*structural annotation*)
- Hamidreza Hakimi (*structural annotation*)
- Aslisho Qurboniev (*structural annotation*)
- Simon Loynes (*structural annotation*)

Past contributors:

- Maroussia Bednarkiewicz (*structural annotation*)
- Christoph Gümmer (*structural annotation*)
- Jonas Köpsel (*structural annotation*)
- Cornelis [Eric] Van Lit (*structural annotation*)
- Cornelia Neubauer (*structural annotation*)
- Leonie Nückell (*structural annotation*)
- Fatemeh Shams (*structural annotation*)