# Watch out for the potholes: The bumpy road of digitizing a historical collection

Tony Chu, PhD, MLIS '14

*Center of Alcohol Studies*
*Rutgers, The State University of New Jersey*

## The collection for digitization

The major task of the digitization project the author involves at the Rutgers Center of Alcohol Studies library is for digitizing the historical collection of Ralph G. Connor Alcohol Research Reference Files (CARRF). The Ralph G. Connor Alcohol Research Reference Files is an archive of instruments (questionnaires, interview schedules, surveys, etc) that have been used in studies of substance abuse. The primary focus of this collection is on various aspects of drinking behavior and alcoholism.

The society for Study Problems founded this collection, which is now named in memory of Dr. Ralph G. Connor, in 1960. The Files were developed and maintained by Dr. Connor at Eastern Washington College, and following his death they were transferred to the instruments to the Center of Alcohol Studies at Rutgers University. The Center has continued to add new instruments to the collection, and periodically issues new updates to the descriptive inventory providing brief information on the materials.

## The tasks of digitization and metadata construction

When building a digital library, the digitization of a historical collection can be done by (1) digitizing and indexing each of the print documents one-by-one, (2) digitizing all the print documents, importing them into the new system, and then indexing each of them, or (3) digitizing all the print documents, importing the existing metadata from the old system, and matching the two parts in the new system. The third option is often considered because the metadata is typically available as an online public access catalog (OPAC). For efficiency, this option is also prefered by reducing the operational time.

Issues could arise with the third option. Using the digitization of our library's historical collection as a case study, we documented two such issues. First, digitizing the print documents, our staff had included the archive ID and title of each document as part of the filename, which helped with general identification, but did so with a lack of delimiters necessary to parse the information with data mining and statistical software such as Statistical Package for the Social Sciences (SPSS). For matching other information with the file ID number, this caused problems when attempting to separate the ID number from the document title.

The second problem was about matching the digitized content with the metadata. The metadata of the online public access catalog were extracted to be included in an Excel file. Each record/row was identified by a unique linking ID number, namely the CARRF id number which refers to the print document archived. On the other hand, the file names of

the digitized documents were saved as another Excel file. In the latter, the linking ID's were manually read, normalized, separated from the file names, and recorded in new columns. For this part, because a physical archive folder may contain several documents which all link to the same metadata record in the old system, an ID could relate to various types of digitized documents. In other words, the existing metadata are not granular enough to specify and describe all the documents of the historical collection. The records of these two files did not appear to simply match each other on one-to-one basis.

The data manipulating software application SPSS was used for reading and matching these two Excel files by the ID's. Before the match/file merger was applied, the duplicates of the ID's in the file that included the file names of the digitized documents were removed. The result of matching showed that the metadata from our existing database system had not been updated to cover the entire collection.

## Possible resolutions

The preliminary steps of linking the file names of the digitized documents and the metadata in the existing system have been taken. For a simple one-to-one merger, a decision was made to only link the file names of the digitized surveys with the metadata. This merged table is being loaded into the new content management system and eventually will be linked with the digital objects of the surveys.

For those archived documents which were not indexed and covered by the existing metadata, workflows need to be created to manually index the digitized documents.