# Repurposing Archival Metadata with the Python CSV Writer
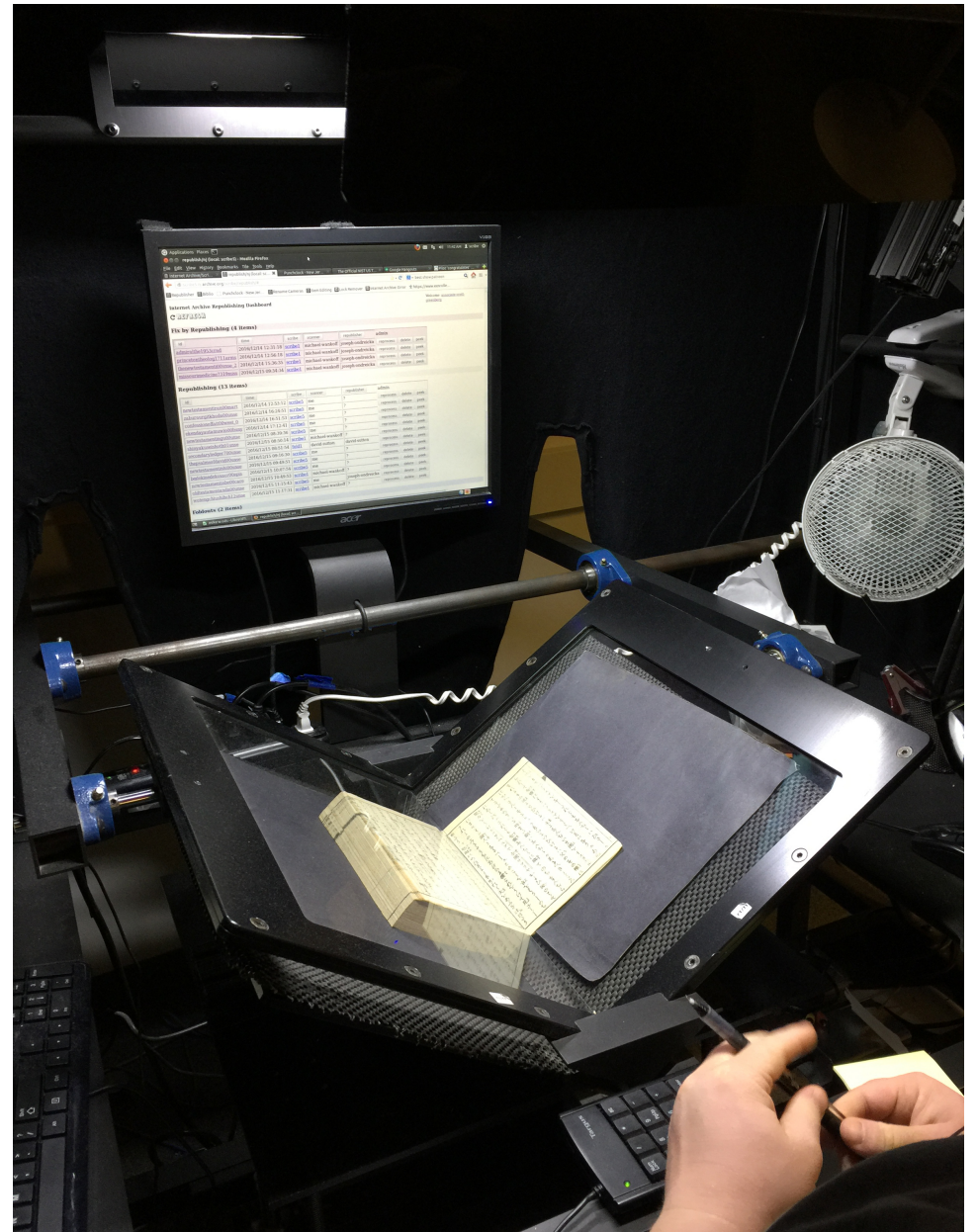
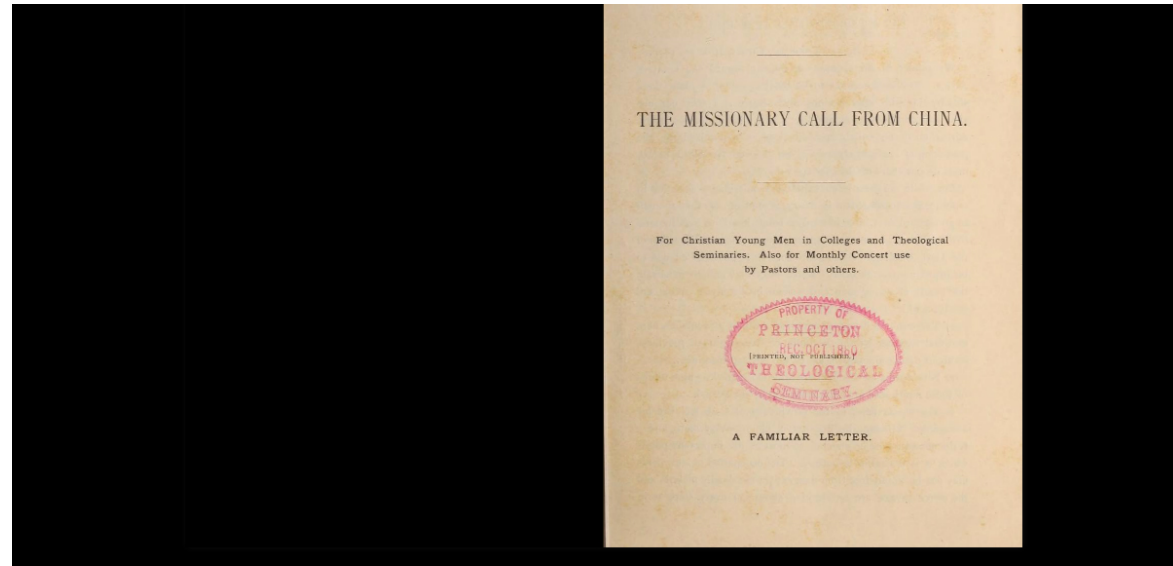Jackie Rider

Programming for Cultural Heritage 2016

December 15, 2016

For my project, I wrote a Python script that will extract selected XML metadata from an EAD encoded finding aid and export it as a CSV file.

Princeton Theological Seminary Library contracts with the Internet Archive to digitize some of the Seminary's special collections. Part of the digitization process includes creating metadata for each digital object and supplying that metadata in a spreadsheet to technicians at the Internet Archive.

The Internet Archive then displays the metadata online with the scanned digital image of the collection item.



THE MISSIONARY CALL FROM CHINA.

For Christian Young Men in Colleges and Theological Seminaries. Also for Monthly Concert use by Pastors and others.

PROPERTY OF PRINCETON THEOLOGICAL SEMINARY.

A FAMILIAR LETTER.

📖 The missionary call from China : a familiar letter
by Chapin, L. D; American Board of Commissioners for Foreign Missions

Published 1880
Topics Missions

SHOW MORE

"For Christian young men in colleges and theological seminaries, also for monthly concert use by pastors and others."

"Printed, not published."

Publisher [Boston : A.B.C.F.M.]
Pages 6
Language English
Call number SCP #45,100
Digitizing sponsor Princeton Theological Seminary Library
Book contributor Princeton Theological Seminary Library
Collection majorityworldcollection; Princeton; americana

Full catalog record MARCXML

This book has an editable web page on Open Library.

DOWNLOAD OPTIONS

| ABBYY GZ | 1 file |
| DAISY | 1 file |
| EPUB | 1 file |
| FULL TEXT | 1 file |
| KINDLE | 1 file |
| PDF | 1 file |
| SINGLE PAGE ORIGINAL JP2 TAR | 1 file |
| SINGLE PAGE PROCESSED JP2 ZIP | 1 file |
| TORRENT | 1 file |

SHOW ALL

15 Files
9 Original

Some collections have been processed and have EAD finding aids; others do not. Student workers are paid to create the metadata for each archival object to be scanned. When saved and exported as a CSV file, metadata created for scanning can be repurposed when creating collection finding aids.

Thus, the Seminary saves time and money by paying one student to create metadata for each object only once.

This metadata is pulled from Library of Congress Subject Headings and Name Authority files and complies with the archival descriptive standard Describing Archives: A Content Standard. Normalized across collections, it can link to other collections internally or content from other institutions with digital collections in the Internet Archive.

# First, I grabbed the XML EAD finding aid file:



```
NeumannFrederick.xml                                    NeumannFrederick.xml

 1  <?xml version="1.0" encoding="utf-8"?>
 2  <ead xmlns="urn:isbn:1-931666-22-9" xmlns:xlink="http://www.w3.
    org/1999/xlink"
 3    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
 4    xsi:schemaLocation="urn:isbn:1-931666-22-9 http://www.loc.gov/ead/ead.
    xsd">
 5    <eadheader countryencoding="iso3166-1" dateencoding="iso8601"
    langencoding="iso639-2b"
 6      repositoryencoding="iso15511">
 7      <eadid countrycode="US" mainagencycode="US-NjPT"/>
 8      <filedesc>
 9        <titlestmt>
10          <titleproper>Frederick Neumann papers<num>2003.3.30</num></
            titleproper>
11        </titlestmt>
12        <publicationstmt>
13          <publisher>Princeton Theological Seminary. Library. Special
            Collections</publisher>
14          <p id="logostmt"><extref xlink:actuate="onLoad"
15              xlink:href="http://www.ptsem.edu/library/images/PTSLOGO-
                Library-large.jpg"
16              xlink:show="embed" xlink:type="simple"/></p>
17          <address><addressline>Princeton Theological Seminary</addressline>
            <addressline>Library</addressline><addressline>PO Box 821</
```

Line 1, Column 1                                    Spaces: 2          XML

Then, I drilled down through the hierarchy of nested EAD elements to the metadata elements I want to extract:



```xml
                Neumann
297             lectures</unittitle>
298             <unitdate normal="1953/1999" type="inclusive">1953–1999</
                unitdate>
299             <container id="aspace_c0127391d856da02ce279c10b8b03ab5" label=
                "Text" type="carton"
300                 >17</container>
301           </did>
302         </c02>
303         <c02 id="aspace_d498008816c17344a37b24d160ff2499" level="file">
304           <did>
305             <unittitle>Printers' galleys</unittitle>
306             <unitdate normal="1965/1979" type="inclusive">1965–1979</
                unitdate>
307             <container id="aspace_eeceef9ac6fb50fcf1c64cd1e629272b" label=
                "Text" type="carton"
308                 >18</container>
309           </did>
310         </c02>
311       </c01>
312     </dsc>
313   </archdesc>
314 </ead>
315
```

Next, I created a Python file and imported the ElementTree Class from the xml.etree module and the csv module, and asked the xml module to load the xml file and parse it:

```python
#Import ElementTree Class from xml.etree module
import xml.etree.ElementTree as etree

#import csv module
import csv

#ask xml module to load xml file and parse it
tree = etree.parse('/Users/jrider/GitHub/PFCH16_Final/NeumannFrederick.xml')

#return the root xml element and store it in root variable
root = tree.getroot()
```

Then, I looped through the XML file to get to my nested metadata: unittitle, unitdate, and container.



```python
        if 'c01' in still_another_element.tag:

            for yet_another_element in still_another_element:

                if 'c02' in yet_another_element.tag:

                    for one_more_element in yet_another_element:

                        if 'did' in one_more_element.tag:

                            for even_one_more_element in one_more_element:

                                if 'unittitle' in even_one_more_element.tag:

                                    print(even_one_more_element.text)

                                if 'unitdate' in even_one_more_element.tag:

                                    print(even_one_more_element.text)

                                if 'container' in even_one_more_element.tag:

                                    print(even_one_more_element.attrib['type'])
```

# I exported that metadata to a CSV file, which can accompany items being sent to the Internet Archive for scanning:

```python
if 'did' in one_more_element.tag:
    csv_row = {'unittitle':'','unitdate':'','container':''}

    for even_one_more_element in one_more_element:

        if 'unittitle' in even_one_more_element.tag:

            print(even_one_more_element.text)
            csv_row['unittitle']=even_one_more_element.text
        if 'unitdate' in even_one_more_element.tag:

            print(even_one_more_element.text)
            csv_row['unitdate']=even_one_more_element.text

        if 'container' in even_one_more_element.tag:

            print(even_one_more_element.attrib['type'])
            csv_row['container']=even_one_more_element.attrib['type']




metadatawriter.writerow([csv_row['unittitle'],csv_row['unitdate'],csv_row['container']
```

| | | |
|---|---|---|
| A-C | 1930-1960 | carton |
| C-E | 1930-1980 | carton |
| F-J | 1929-1968 | carton |
| J-O | 1933-1980 | carton |
| P-Z | 1931-1990 | carton |
| Daily planners and diaries | 1950-1967 | carton |
| Sermon notecards | 1940-1960 | box |
| Published articles, works by others, biographical materials | 1875-1997 | carton |
| Faith and Knowledge, God's Fifth Columnist, the Jewish Question | 1937-1994 | carton |
| Der Römerbrief | 1943-1996 | carton |
| Sermons, Sermon on the Mount | 1943-1990 | carton |
| Publisher correspondence, publicity materials, PTS Neumann lectures | 1953-1999 | carton |
| Printers' galleys | 1965-1979 | carton |

Result is normalized and repurposed metadata. Last step: use the Python command module to convert all of this script to one command line a student worker can easily run without extensive coding:

The command/sys module presents an argument that pulls the scripts together into one command. Working in the terminal, a student worker simply changes the xml file being parsed.

```
11
12  tree = etree.parse(sys.argv[1])
13
14  #return the root xml element and store it in root variable
15  root = tree.getroot()
16
17  with open('metadata.csv', 'w') as csvfile:
18      metadatawriter = csv.writer(csvfile, delimiter=',')
19      #use for loop to loop through root element
20      for a_element in root:
21
```

# Success!!

```
● ● ●                  📁 PFCH16_Final — -bash — 80×24

Sermon notecards
1940-1960
box
Published articles, works by others, biographical materials
1875-1997
carton
Faith and Knowledge, God's Fifth Columnist, the Jewish Question
1937-1994
carton
Der Römerbrief
1943-1996
carton
Sermons, Sermon on the Mount
1943-1990
carton
Publisher correspondence, publicity materials, PTS Neumann
               lectures
1953-1999
carton
Printers' galleys
1965-1979
carton
jackierdersmbp3:PFCH16_Final jrider$ python3.5 read_Neumann.py '/Users/jrider/Gi
tHub/PFCH16_Final/NeumannFrederick.xml'▯
```