# Automated Enhancement of Controlled Vocabularies: Upgrading Legacy Metadata in CONTENTdm
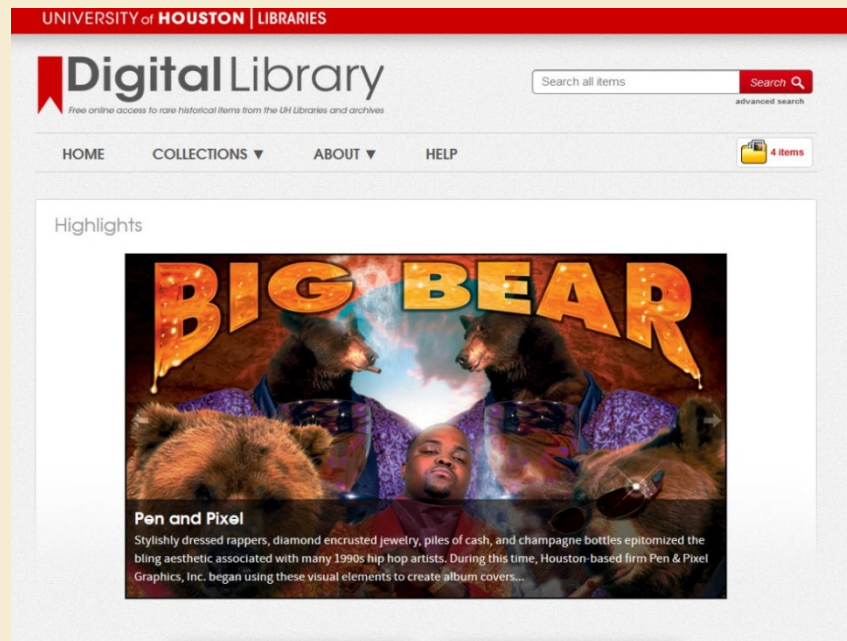
Andrew Weidner, Metadata Librarian

Annie Wu, Head of Metadata and Digitization Services

Santi Thompson, Head of Digital Repository Services

UNIVERSITY of **HOUSTON** | **LIBRARIES**
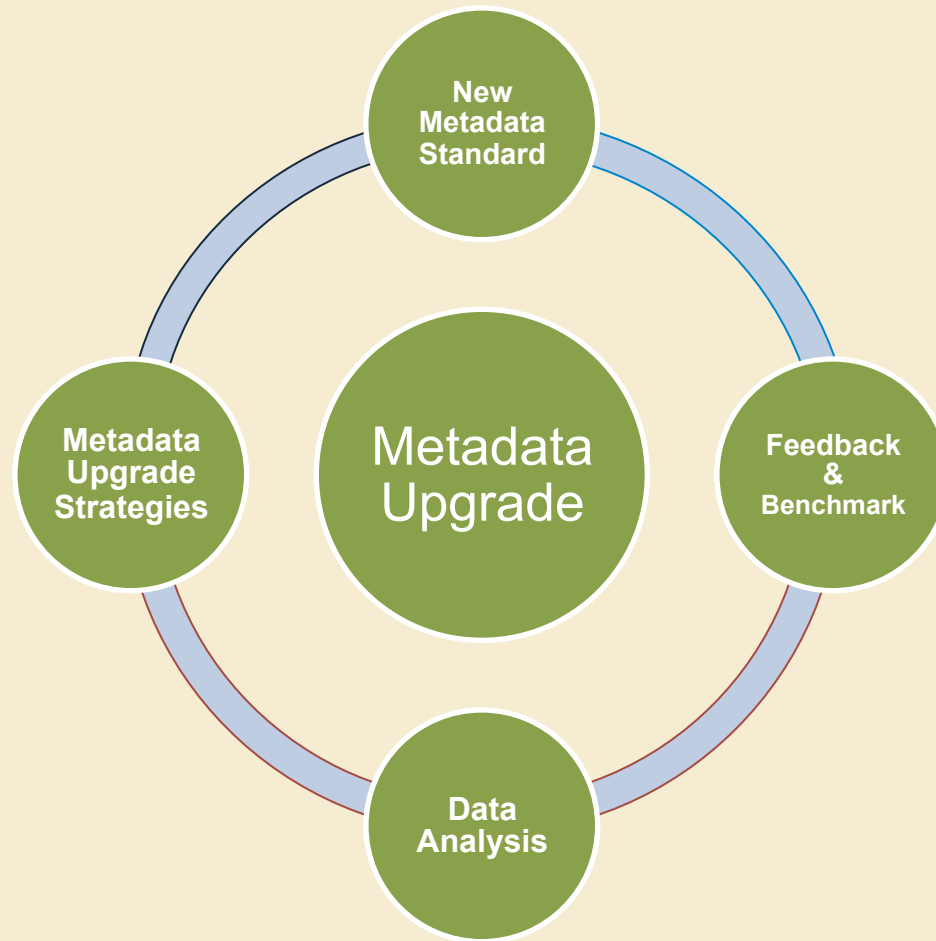
# UH Digital Library



- Launched in 2009
- 60 collections
- 55,000+ digital objects
- Images, Text, Audio/ Video

# Metadata Upgrade: Rationale

- Provide discoverability of digital objects via our discovery platform

- Improve retrieval and display of digital metadata through our newly designed UHDL interface

- Ensure metadata interoperability with other digital collections

- Be ready for linked data environment

More Robust, Reliable and Retrievable Metadata!

# Metadata Upgrade: Overview



- New Metadata Standard
- Feedback & Benchmark
- Metadata Upgrade
- Metadata Upgrade Strategies
- Data Analysis

# Metadata Upgrade: Strategy

## *Phase One*

- Standardize field labels

- Standardize mapping to DC

- Add new fields to UHDL

# Metadata Upgrade: Strategy

## Phase One

- Standardize field labels

- Standardize mapping to DC

- Add new fields to UHDL

## Phase Two

- Standardize collection names

- Add or edit collection level fields

- Cross linking

# Metadata Upgrade: Strategy

## Phase One

- Standardize field labels

- Standardize mapping to DC

- Add new fields to UHDL

## Phase Two

- Standardize collection names

- Add or edit collection level fields

- Cross linking

## Phase Three

- Templates to titles

- Apply ISO date

- Add DCMI type terms

- Headings control

# **Phase Three: Challenges**

- Limitations of existing interfaces
  - No easy way to automate or batch edit legacy metadata in Project Client or Web Admin interfaces
  - Must edit item-by-item, which is time consuming and resource intensive

# Phase Three: Challenges

- Limitations of existing interfaces
  - No easy way to automate or batch edit legacy metadata in Project Client or Web Admin interfaces
  - Must edit item-by-item, which is time consuming and resource intensive

- Limitations of tab-delimited process
  - Exporting metadata, editing tab delimited files, and re-uploading content can strain server and software
  - OpenRefine can tell us where the problems are, but the work still has to be done in the Project Client or Web Admin

# Phase Three: Challenges

- Limitations of existing interfaces
    - No easy way to automate or batch edit legacy metadata in Project Client or Web Admin interfaces
    - Must edit item-by-item, which is time consuming and resource intensive

- Limitations of tab-delimited process
    - Exporting metadata, editing tab delimited files, and re-uploading content can strain server and software
    - OpenRefine can tell us where the problems are, but the work still has to be done in the Project Client or Web Admin

- Limitations on staff time

# Phase Three: Solutions

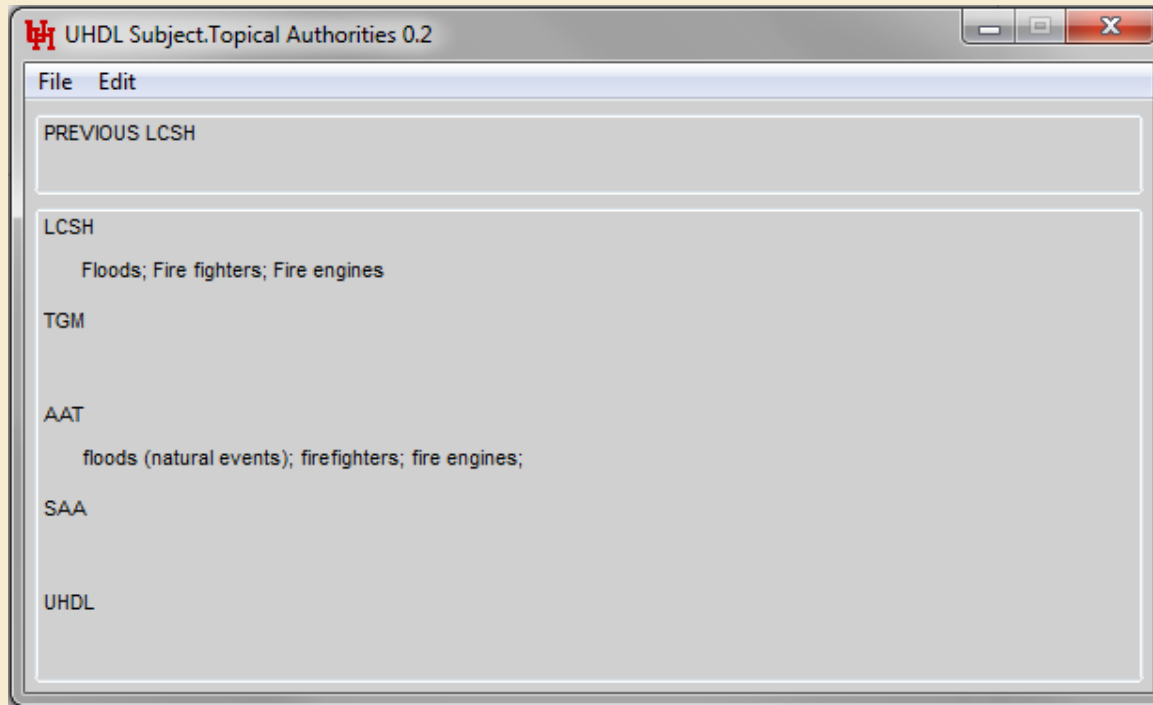- Work within Project Client

# Phase Three: Solutions

- Work within Project Client

- Develop tools layered on top of the Project Client to automate key tasks

# Phase Three: Solutions

- Work within Project Client

- Develop tools layered on top of the Project Client to automate key tasks
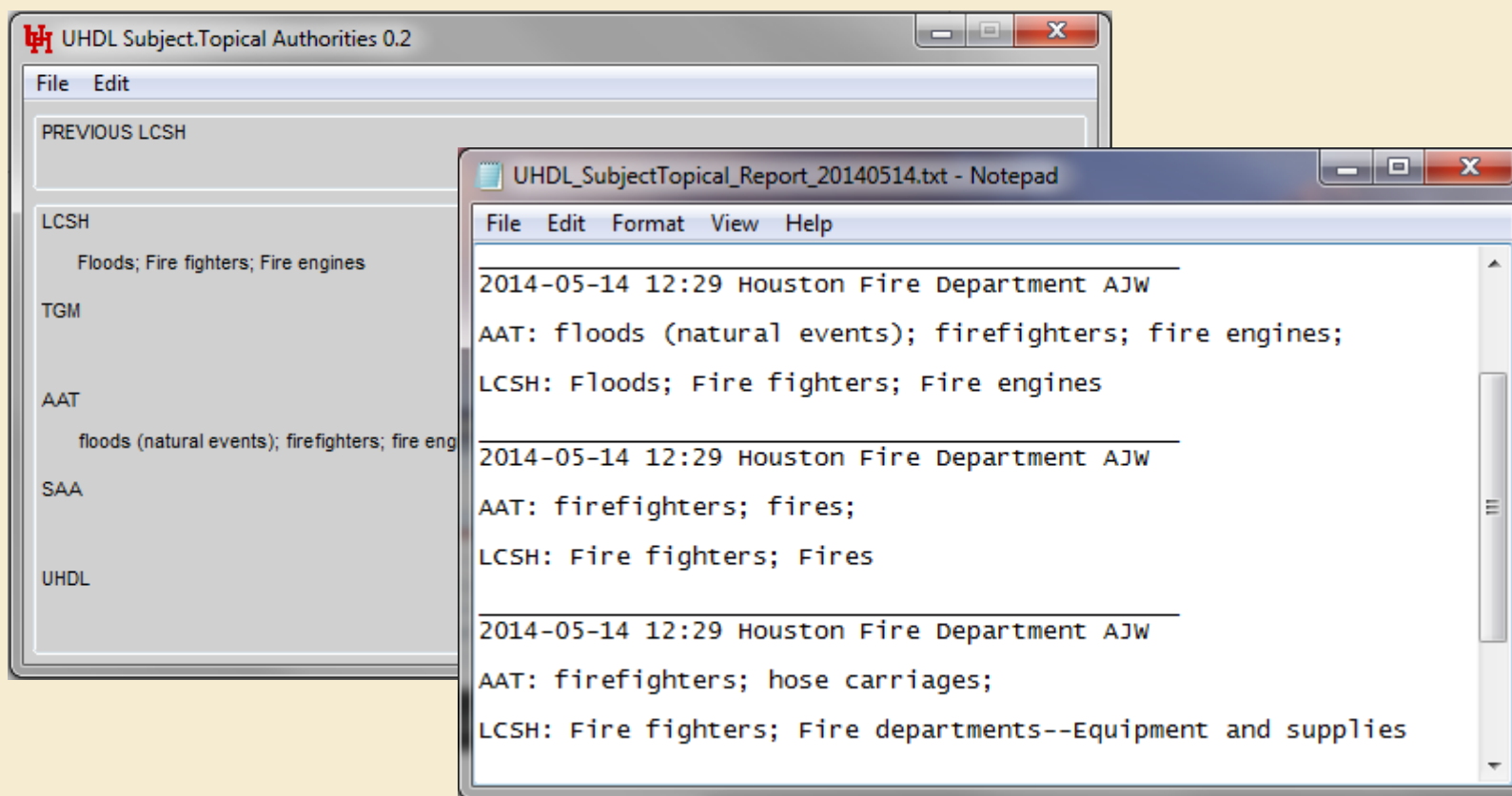
# Upgrade Tools: Subject App



- Automates mapping from various vocabularies to LCSH

- Facilitates harvesting of URIs from LC linked data service (id.loc.gov)

# Upgrade Tools: Subject App

File   Edit   Search   View   Encoding   Language   Settings   Macro   Run   Plugins   Window   ?                                                                X

UHDL_SubjectTopical.txt

```
241    TGM Schools  Schools http://id.loc.gov/authorities/subjects/sh85118451    CPH 20140513
242    TGM Cotton   Cotton  http://id.loc.gov/authorities/subjects/sh85033302    CPH 20140513
243    TGM Government facilities   Government buildings     http://id.loc.gov/authorities/subjects/sh85108620    CPH 20140513
244    AAT buildings    Buildings   http://id.loc.gov/authorities/subjects/sh85017769    CPH 20140513
245    AAT farming Agricultural systems     http://id.loc.gov/authorities/subjects/sh85002407   CPH 20140513
246    AAT hospitals    Hospitals   http://id.loc.gov/authorities/subjects/sh85062285    CPH 20140513
247    TGM Jails    Jails    http://id.loc.gov/authorities/subjects/sh85069253    CPH 20140513
248    AAT banks    Banks and banking, Central  http://id.loc.gov/authorities/subjects/sh85011647    CPH 20140513
249    AAT city blocks City blocks http://id.loc.gov/authorities/subjects/sh85014924    CPH 20140513
250    TGM Waterworks   Waterworks  http://id.loc.gov/authorities/subjects/sh85145757    CPH 20140513
251    TGM Buildings    Buildings   http://id.loc.gov/authorities/subjects/sh85017769    CPH 20140513
252    TGM Portrait photographs    Portraits   http://id.loc.gov/authorities/subjects/sh85105182    CPH 20140513
253    TGM Shipping     Shipping    http://id.loc.gov/authorities/subjects/sh85121579    CPH 20140513
254    TGM Advertising Advertising http://id.loc.gov/authorities/subjects/sh85001086    CPH 20140513
255    TGM Business & finance   Business    http://id.loc.gov/authorities/subjects/sh85018260    CPH 20140513
256    AAT buildings (structures)   Buildings   http://id.loc.gov/authorities/subjects/sh85017769    CPH 20140513
257    AAT skyscrapers Skyscrapers http://id.loc.gov/authorities/subjects/sh85123270    CPH 20140513
258    TGM Skyscrapers Skyscrapers http://id.loc.gov/authorities/subjects/sh85123270    CPH 20140513
259    TGM Banks    Banks and banking, Central  http://id.loc.gov/authorities/subjects/sh85011647    CPH 20140513
260    TGM Railroad stations    Railroad stations    http://id.loc.gov/authorities/subjects/sh85111042    CPH 20140513
261    TGM Industry     Industry & trade summary     http://id.loc.gov/authorities/names/n92090429    CPH 20140513
262    AAT factories    Factories   http://id.loc.gov/authorities/subjects/sh85046823    CPH 20140513
263    TGM Advertisements  Advertising http://id.loc.gov/authorities/subjects/sh85001086    CPH 20140513
264    TGM Construction     Construction    http://id.loc.gov/authorities/subjects/sh99005337    CPH 20140513
265    AAT architectural drawings  Designs and plans   http://id.loc.gov/authorities/subjects/sh87005580    CPH 20140513
266    AAT banks (buildings)    Bank buildings  http://id.loc.gov/authorities/subjects/sh85011556    CPH 20140513
267    AAT stretcher    Litters http://id.loc.gov/authorities/subjects/sh85077674    AJW 20140513
268    AAT floods (natural events) Floods  http://id.loc.gov/authorities/subjects/sh85049168    AJW 20140514
```

Normal text file                            length : 25183    lines : 269            Ln : 1   Col : 1   Sel : 0 | 0                    Dos\Windows        UTF-8            INS

File   Edit   Search   View   Encoding   Language   Settings   Macro   Run   Plugins   Window   ?
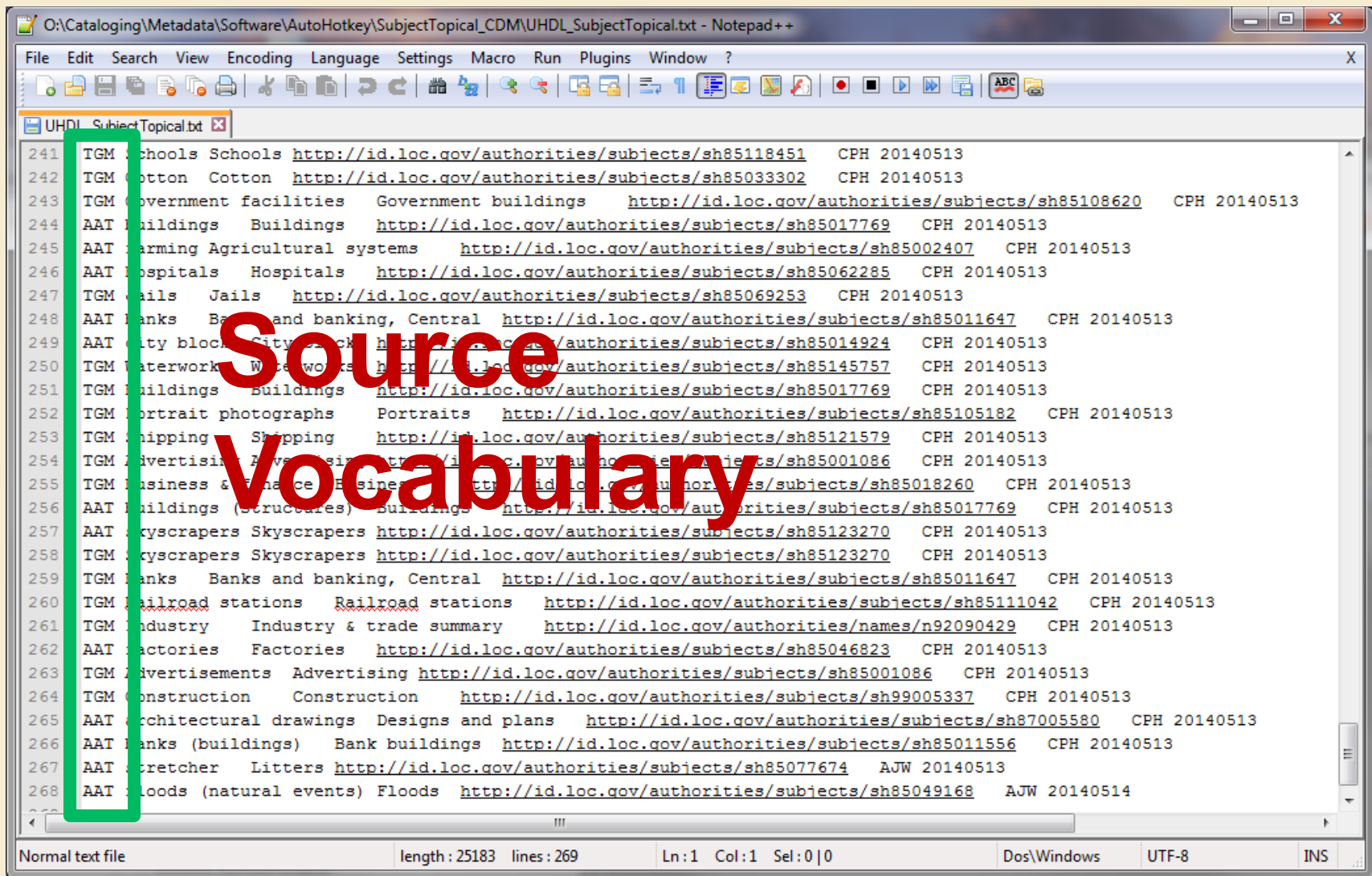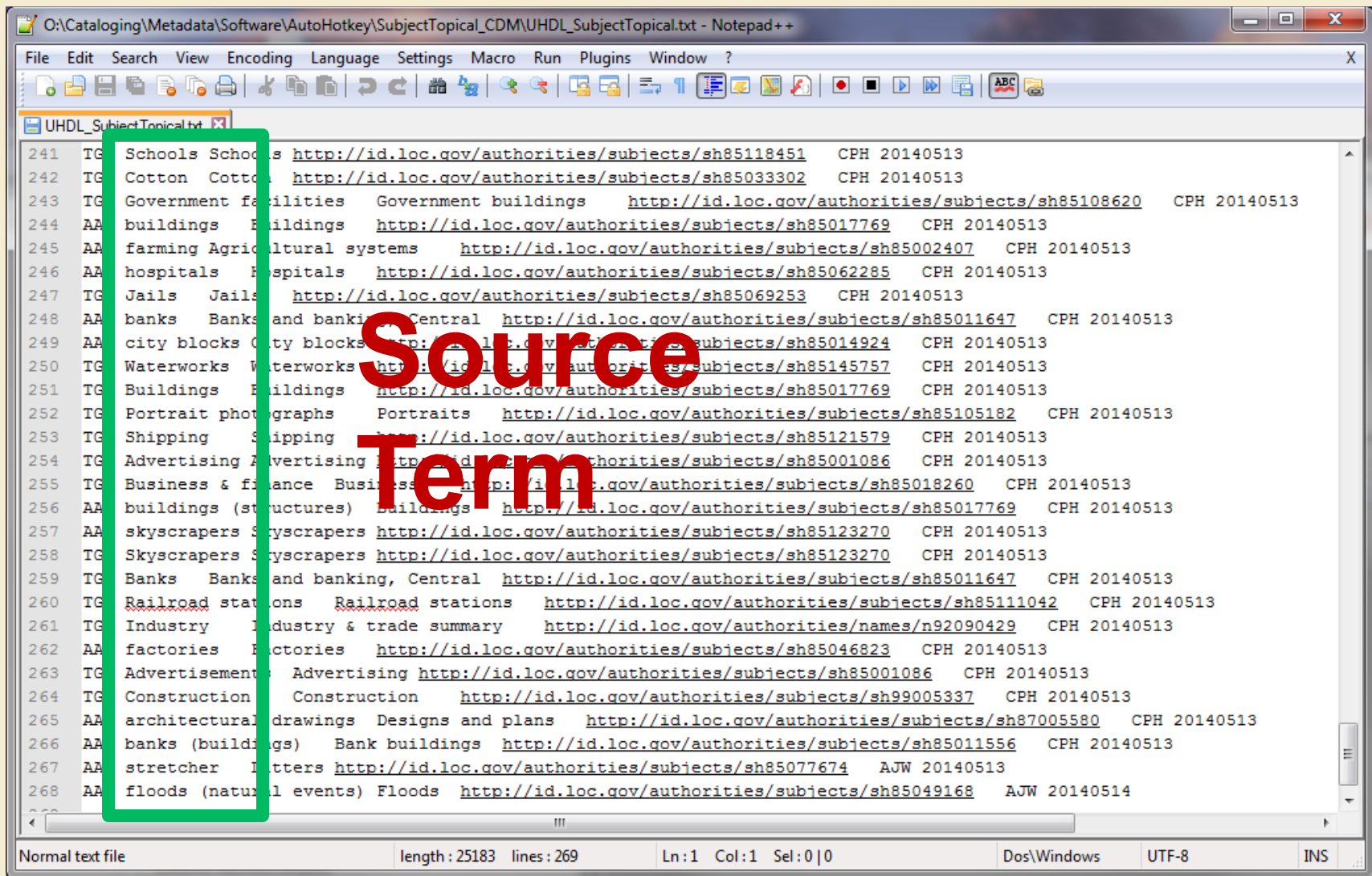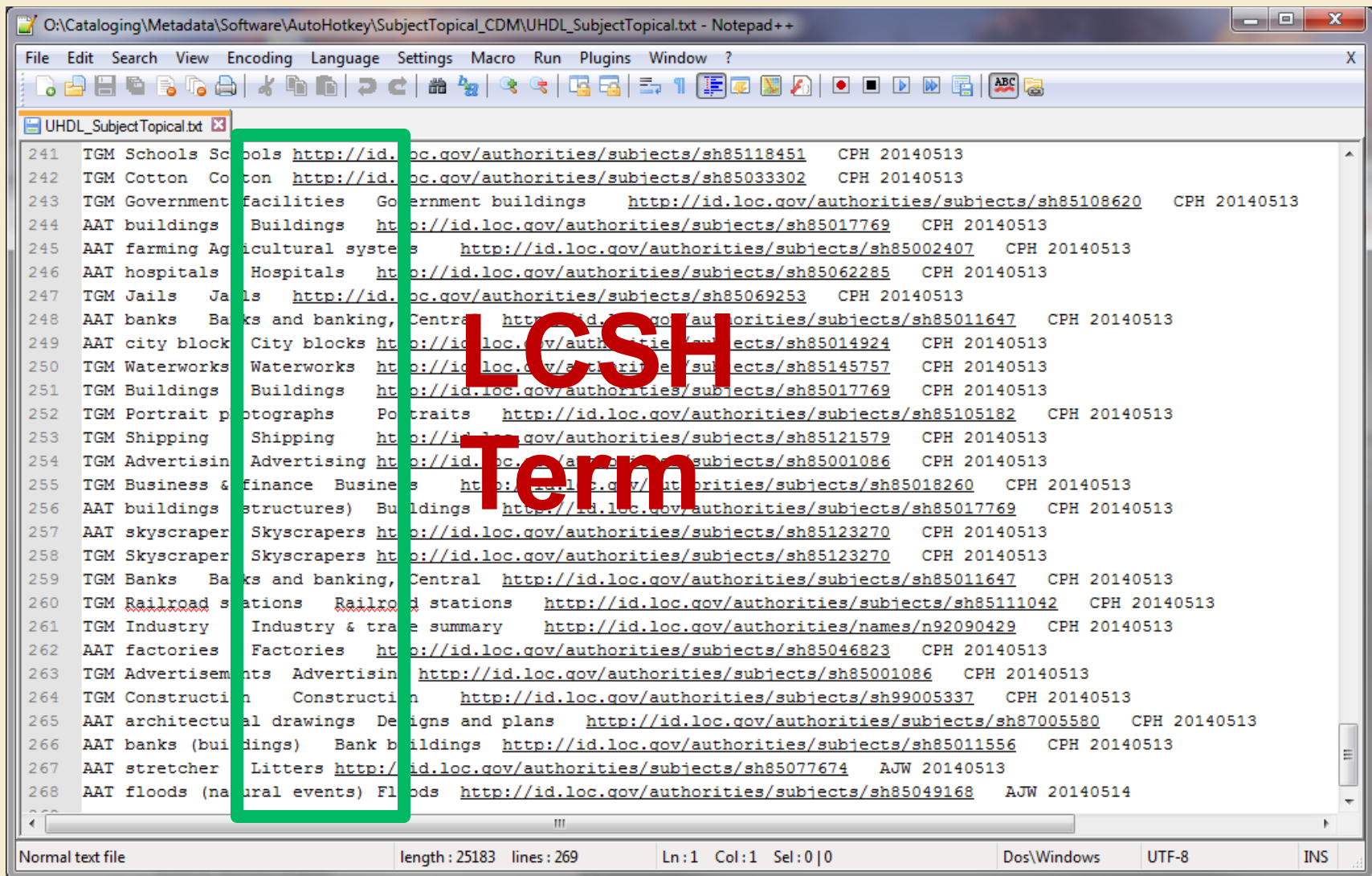
UHDL_SubjectTopical.txt

```
241  TGM Schools  Schools  http://id.loc.gov/authorities/subjects/sh85118451   CPH 20140513
242  TGM Cotton   Cotton   http://id.loc.gov/authorities/subjects/sh85033302   CPH 20140513
243  TGM Government facilities  Government buildings   http://id.loc.gov/authorities/subjects/sh85108620   CPH 20140513
244  AAT Buildings    Buildings   http://id.loc.gov/authorities/subjects/sh85017769   CPH 20140513
245  AAT Farming Agricultural systems    http://id.loc.gov/authorities/subjects/sh85002407   CPH 20140513
246  AAT Hospitals    Hospitals   http://id.loc.gov/authorities/subjects/sh85062285   CPH 20140513
247  TGM Jails    Jails   http://id.loc.gov/authorities/subjects/sh85069253   CPH 20140513
248  AAT Banks    Banks and banking, Central   http://id.loc.gov/authorities/subjects/sh85011647   CPH 20140513
249  AAT City blocks  City blocks  http://id.loc.gov/authorities/subjects/sh85014924   CPH 20140513
250  TGM Waterworks   Waterworks   http://id.loc.gov/authorities/subjects/sh85145757   CPH 20140513
251  TGM Buildings   Buildings   http://id.loc.gov/authorities/subjects/sh85017769   CPH 20140513
252  TGM Portrait photographs   Portraits   http://id.loc.gov/authorities/subjects/sh85105182   CPH 20140513
253  TGM Shipping    Shipping   http://id.loc.gov/authorities/subjects/sh85121579   CPH 20140513
254  TGM Advertising Advertising  http://id.loc.gov/authorities/subjects/sh85001086   CPH 20140513
255  TGM Business & Finance Business   http://id.loc.gov/authorities/subjects/sh85018260   CPH 20140513
256  AAT Buildings (structures)  Buildings   http://id.loc.gov/authorities/subjects/sh85017769   CPH 20140513
257  AAT Skyscrapers Skyscrapers  http://id.loc.gov/authorities/subjects/sh85123270   CPH 20140513
258  TGM Skyscrapers Skyscrapers  http://id.loc.gov/authorities/subjects/sh85123270   CPH 20140513
259  TGM Banks    Banks and banking, Central   http://id.loc.gov/authorities/subjects/sh85011647   CPH 20140513
260  TGM Railroad stations    Railroad stations    http://id.loc.gov/authorities/subjects/sh85111042   CPH 20140513
261  TGM Industry    Industry & trade summary    http://id.loc.gov/authorities/names/n92090429   CPH 20140513
262  AAT Factories    Factories   http://id.loc.gov/authorities/subjects/sh85046823   CPH 20140513
263  TGM Advertisements  Advertising  http://id.loc.gov/authorities/subjects/sh85001086   CPH 20140513
264  TGM Construction    Construction    http://id.loc.gov/authorities/subjects/sh99005337   CPH 20140513
265  AAT Architectural drawings  Designs and plans   http://id.loc.gov/authorities/subjects/sh87005580   CPH 20140513
266  AAT Banks (buildings)   Bank buildings  http://id.loc.gov/authorities/subjects/sh85011556   CPH 20140513
267  AAT Stretcher   Litters  http://id.loc.gov/authorities/subjects/sh85077674   AJW 20140513
268  AAT Floods (natural events) Floods   http://id.loc.gov/authorities/subjects/sh85049168   AJW 20140514
```
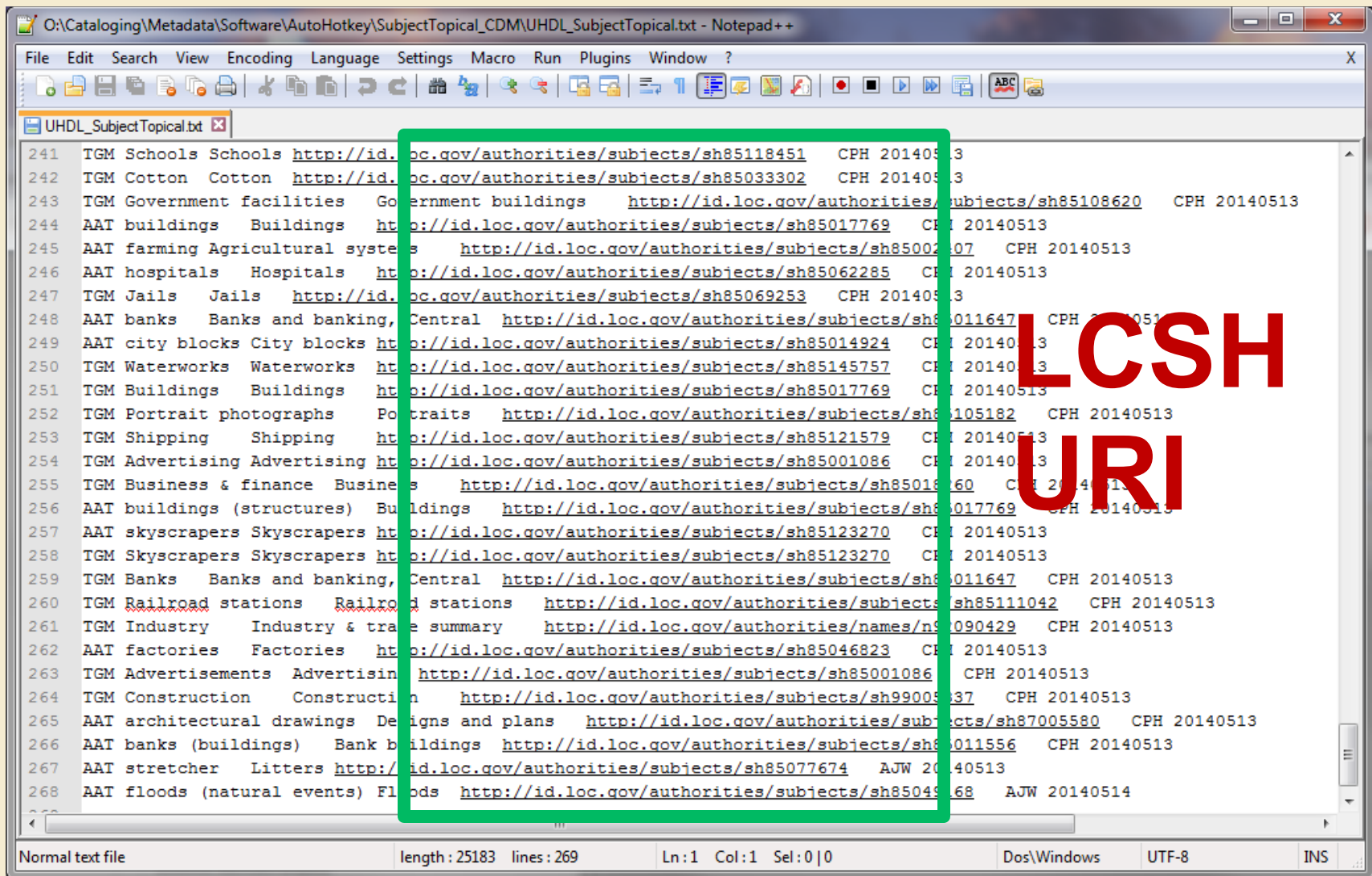
**Source Vocabulary**

Normal text file     length : 25183   lines : 269     Ln : 1   Col : 1   Sel : 0 | 0     Dos\Windows     UTF-8     INS

File   Edit   Search   View   Encoding   Language   Settings   Macro   Run   Plugins   Window   ?                                                                X
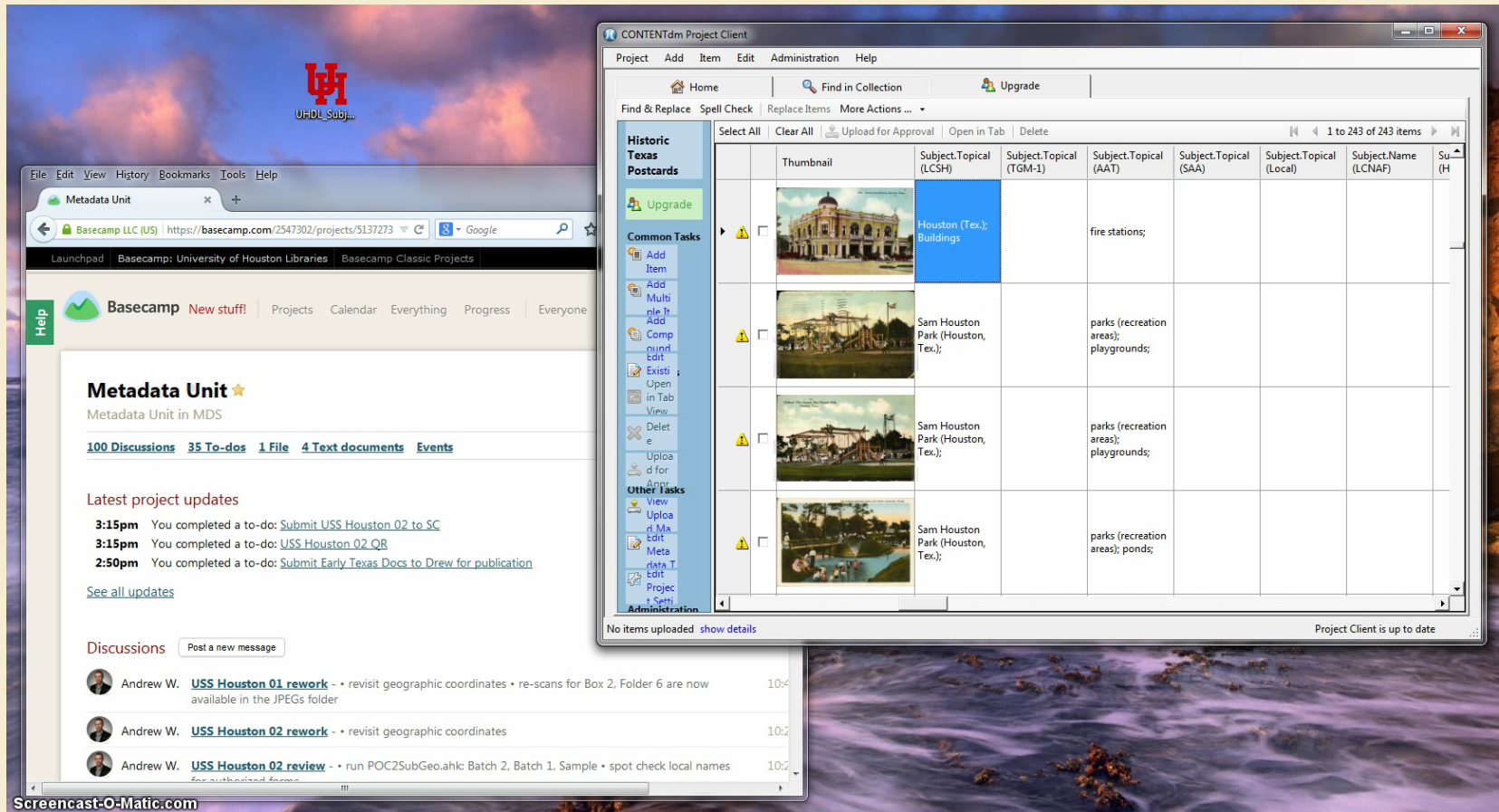
UHDL_SubjectTopical.txt

```
241   TGM Schools  Schools  http://id.loc.gov/authorities/subjects/sh85118451    CPH 20140513
242   TGM Cotton   Cotton   http://id.loc.gov/authorities/subjects/sh85033302    CPH 20140513
243   TGM Government facilities   Government buildings    http://id.loc.gov/authorities/subjects/sh85108620    CPH 20140513
244   AAT buildings    Buildings    http://id.loc.gov/authorities/subjects/sh85017769    CPH 20140513
245   AAT farming Agricultural systems    http://id.loc.gov/authorities/subjects/sh85002307    CPH 20140513
246   AAT hospitals    Hospitals    http://id.loc.gov/authorities/subjects/sh85062285    CPH 20140513
247   TGM Jails    Jails    http://id.loc.gov/authorities/subjects/sh85069253    CPH 20140513
248   AAT banks    Banks and banking, Central    http://id.loc.gov/authorities/subjects/sh85011647    CPH 20140513
249   AAT city blocks City blocks http://id.loc.gov/authorities/subjects/sh85014924    CPH 20140513
250   TGM Waterworks   Waterworks   http://id.loc.gov/authorities/subjects/sh85145757    CPH 20140513
251   TGM Buildings    Buildings    http://id.loc.gov/authorities/subjects/sh85017769    CPH 20140513
252   TGM Portrait photographs    Portraits    http://id.loc.gov/authorities/subjects/sh85105182    CPH 20140513
253   TGM Shipping    Shipping    http://id.loc.gov/authorities/subjects/sh85121579    CPH 20140513
254   TGM Advertising Advertising http://id.loc.gov/authorities/subjects/sh85001086    CPH 20140513
255   TGM Business & finance  Business    http://id.loc.gov/authorities/subjects/sh85018860    CPH 20140513
256   AAT buildings (structures)  Buildings    http://id.loc.gov/authorities/subjects/sh85017769    CPH 20140513
257   AAT skyscrapers Skyscrapers http://id.loc.gov/authorities/subjects/sh85123270    CPH 20140513
258   TGM Skyscrapers Skyscrapers http://id.loc.gov/authorities/subjects/sh85123270    CPH 20140513
259   TGM Banks    Banks and banking, Central    http://id.loc.gov/authorities/subjects/sh85011647    CPH 20140513
260   TGM Railroad stations    Railroad stations    http://id.loc.gov/authorities/subjects/sh85111042    CPH 20140513
261   TGM Industry    Industry & trade summary    http://id.loc.gov/authorities/names/n91090429    CPH 20140513
262   AAT factories    Factories    http://id.loc.gov/authorities/subjects/sh85046823    CPH 20140513
263   TGM Advertisements  Advertising http://id.loc.gov/authorities/subjects/sh85001086    CPH 20140513
264   TGM Construction    Construction    http://id.loc.gov/authorities/subjects/sh99005337    CPH 20140513
265   AAT architectural drawings  Designs and plans    http://id.loc.gov/authorities/subjects/sh87005580    CPH 20140513
266   AAT banks (buildings)    Bank buildings http://id.loc.gov/authorities/subjects/sh85011556    CPH 20140513
267   AAT stretcher    Litters http://id.loc.gov/authorities/subjects/sh85077674    AJW 20140513
268   AAT floods (natural events) Floods    http://id.loc.gov/authorities/subjects/sh85049268    AJW 20140514
```

**LCSH URI**

UNIVERSITY of HOUSTON | LIBRARIES

File   Edit   Search   View   Encoding   Language   Settings   Macro   Run   Plugins   Window   ?                                                                                X

UHDL_SubjectTopical.txt

```
241   TGM Schools  Schools http://id.loc.gov/authorities/subjects/sh85118451    CPH 20 40513
242   TGM Cotton   Cotton  http://id.loc.gov/authorities/subjects/sh85033302    CPH 20 40513
243   TGM Government facilities    Government buildings    http://id.loc.gov/authorit es/subjects/sh85108620    C H 20140513
244   AAT buildings     Buildings   http://id.loc.gov/authorities/subjects/sh85017769    CPH 20140513
245   AAT farming Agricultural systems     http://id.loc.gov/authorities/subjects/sh8 002407    CPH 20140513
246   AAT hospitals     Hospitals   http://id.loc.gov/authorities/subjects/sh85062285    CPH 20140513
247   TGM Jails    Jails   http://id.loc.gov/authorities/subjects/sh85069253    CPH 20 40513
248   AAT banks    Banks and banking, Central   http://id. oc. v/authorities/subjects sh85011647    CPH 20140513
249   AAT city blocks City blocks http://id.loc.gov/auth rit s  ie /s 014924    CPH 20140513
250   TGM Waterworks   Waterworks  http://id.loc.gov/auth  it es/se rr 85145757    CPH 20140513
251   TGM Buildings    Buildings   http://id.loc.gov/authorities/subjects/sh85017769    CPH 20140513
252   TGM Portrait photographs     Portraits   http://id.loc.gov/authorities/subjects sh85105182    CPH 20140513
253   TGM Shipping     Shipping    http://id.loc.gov/auth rities/subjects/sh85121579    CPH 20140513
254   TGM Advertising Advertising http://id.loc.gov/auth rit   ie /sh85001086    CPH 20140513
255   TGM Business & finance  Business    http://id.loc. v au or t es ubjects/sh8 018260    CPH 20140513
256   AAT buildings (structures)   Buildings   http://id.loc.gov/authorities/subjects sh85017769    CPH 20140513
257   AAT skyscrapers Skyscrapers http://id.loc.gov/authorities/subjects/sh85123270    CPH 20140513
258   TGM Skyscrapers Skyscrapers http://id.loc.gov/authorities/subjects/sh85123270    CPH 20140513
259   TGM Banks    Banks and banking, Central   http://id.loc.gov/authorities/subjects sh85011647    CPH 20140513
260   TGM Railroad stations    Railroad stations    http://id.loc.gov/authorities/subj cts/sh85111042    CPH 20140 13
261   TGM Industry    Industry & trade summary    http://id.loc.gov/authorities/name /n92090429    CPH 20140513
262   AAT factories   Factories   http://id.loc.gov/authorities/subjects/sh85046823    CPH 20140513
263   TGM Advertisements  Advertising http://id.loc.gov/authorities/subjects/sh8500 86    CPH 20140513
264   TGM Construction    Construction    http://id.loc.gov/authorities/subjects/sh9 005337    CPH 20140513
265   AAT architectural drawings  Designs and plans   http://id.loc.gov/authorities/ ubjects/sh87005580    CPH 2 140513
266   AAT banks (buildings)   Bank buildings  http://id.loc.gov/authorities/subjects sh85011556    CPH 20140513
267   AAT stretcher   Litters http://id.loc.gov/authorities/subjects/sh85077674    AJ  20140513
268   AAT floods (natural events) Floods  http://id.loc.gov/authorities/subjects/sh8 049168    AJW 20140514
```

**User Info**

UNIVERSITY of HOUSTON | LIBRARIES

# Upgrade Tools: Subject App



https://www.dropbox.com/s/wg9nwn73k5d4hzt/01_SubjectApp.avi?dl=0

# Upgrade Tools: Name App

Logan, William
M., 1802-1839;
Elliot, Peter

Sublett, Philip
Allen, 1802-1850;
Burditt, Jesse;
Augustus, G. W.;
Hotchkiss,
Archibald,
1794-1882;
Thomas. John:

Bell, Thomas B.;
Potter, Robert,
1799-1842;
Harden, Samuel
H.

- Large collections with legacy metadata

- All names entered in LCNAF field and formatted to look like an LCNAF name

- App moves names to correct authority

# Benefits of Metadata Upgrade

- Best Practices
- Data interoperable
- Enhance retrieval of digital data
- New and better standard for future metadata creation and integration



Image source: http://123rf.com

# Thank You!



Andrew Weidner:
[ajweidner@uh.edu](mailto:ajweidner@uh.edu)

Annie Wu:
[awu@uh.edu](mailto:awu@uh.edu)

Santi Thompson:
[sathompson@uh.edu](mailto:sathompson@uh.edu)

# Questions?