Agenda:

- The Spencer project

- Re-examining the workflow

- Implementing OpenRefine

- Other uses for OpenRefine in digital collections

# Spencer American Popular Sheet Music Project

# Spencer American Popular Sheet Music Project

Project started in 1999

TexTreasures Grant to digitize 1,000 pieces

Descriptive metadata loaded into ILS

Static HTML was programmatically generated and placed on server

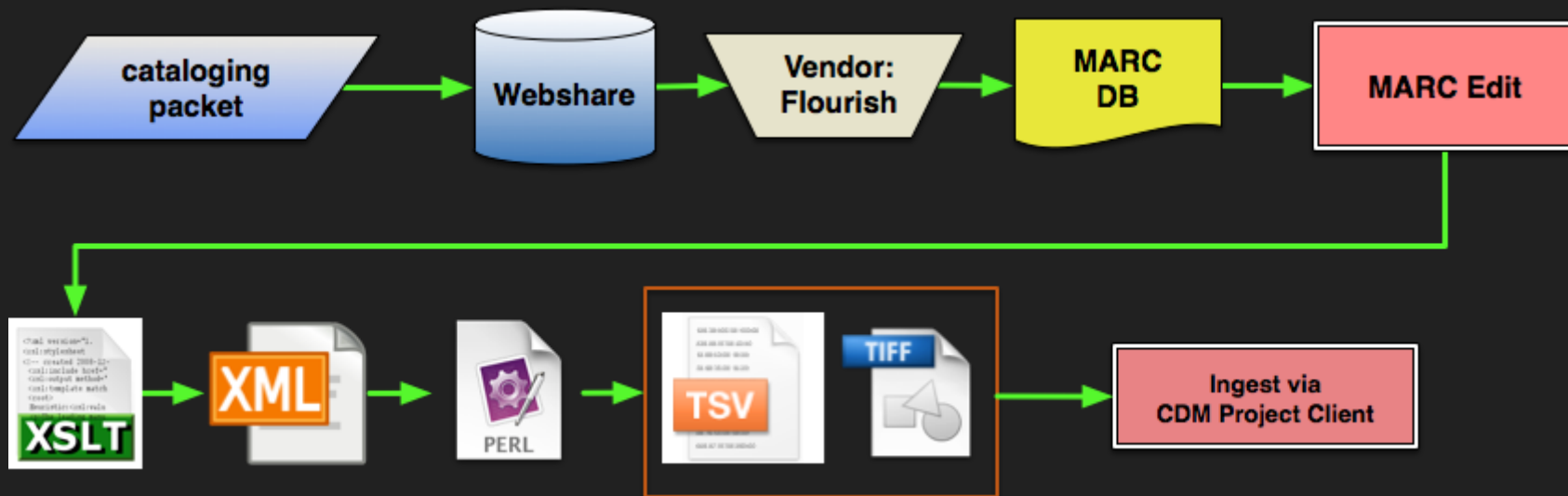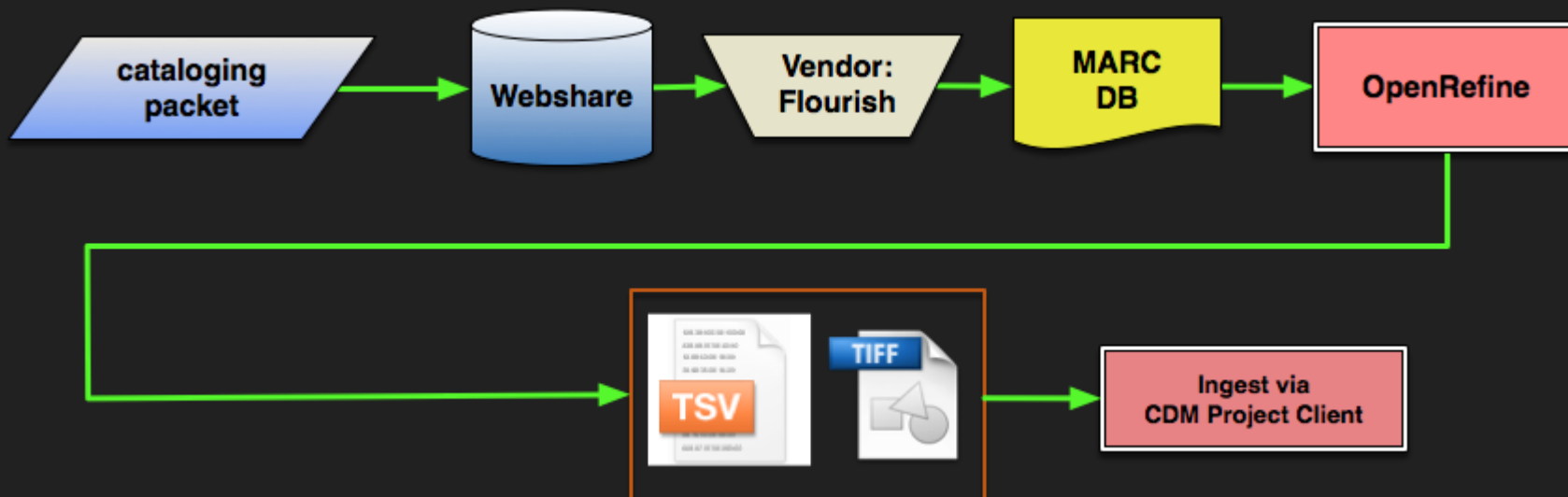# Spencer American Popular Sheet Music Project



CONTENTdm®



IN-HOUSE
OUTSOURCE ✓



Flourish
Music·Contract·Cataloging

# Spencer American Popular Sheet Music Project
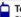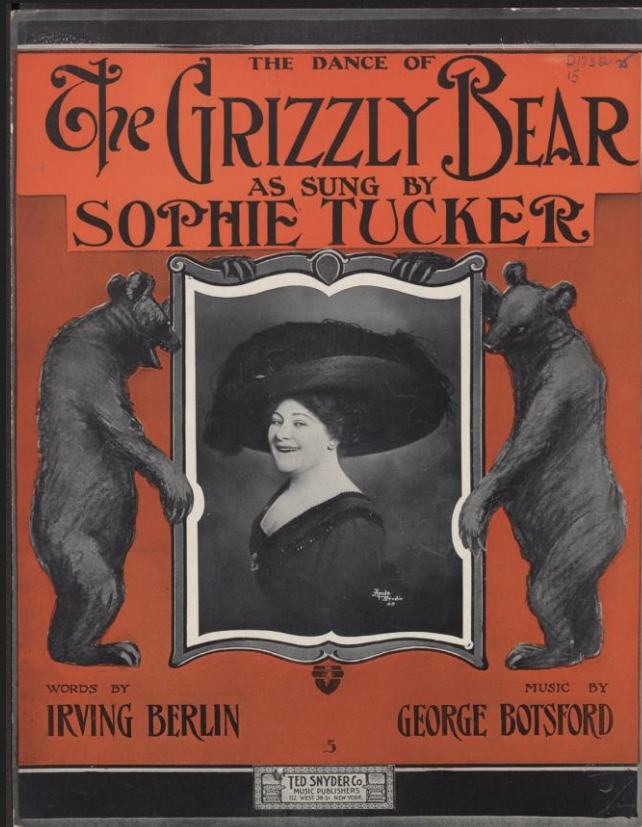
# Spencer American Popular Sheet Music Project

#18853426

D1732  Botsford, George, 1874 - 1949.
15-2      Grizzly bear. Words by Irving Berlin.
        As sung by Tim McMahon. N. Y., Ted Snyder
SPENCER  Co. (Inc.), 1910.
        no pl. no.

        Lyrics with piano
        Fc: orange, bears, photo Tim McMahon.
        Bc: excerpts 2 songs, illus. covers.

Author  Botsford, George, 1874-1949.
Title  The dance of the grizzly bear / words by Irving Berlin ; music by George Botsford.
Publication Info  New York : Ted Snyder Co., c1910.

Connect to:
    Connect to this title in the Baylor University Libraries Digital Collections
    Connect to this title in the Baylor University Libraries Digital Collections

| LOCATION | CALL # | STATUS |
| --- | --- | --- |
| Crouch SCR Spencer | Spencer D1732 .15 | BY APPT ONLY |
| Crouch SCR Spencer | Spencer D1732 .15-2 | BY APPT ONLY |

📱 Text the
    call number

Description  1 score (5, [1] p.) ; 35 cm.
Note  For voice and piano.
        D1732 .15: illustrated t. p. in orange, black and white with a drawing of 2 bears / Frew; photo. of Sophie Tucker.
        D1732 .15-2: Illustrated t. p. in orange, black and white with a drawing of 2 bears / Frew; photo. of Tim McMahon.
        D1732.15: "As sung by Sophie Tucker."
        D1732.15-2: "As sung by Tim McMahon."
        D1732.15: Advertisement for Alexander's ragtime band march and twostep / Berlin on p. [6].
        D1732.15-2: Advertisement for Draggy rag and Call me up some rainy afternoon / Berlin on p. [6].
Local Note  Frances G. Spencer Collection of American Sheet Music.
        Spencer subject: Dance - other.
        Spencer subject: Famous people - Sophie Tucker.
Subject  Tucker, Sophie, 1884-1966 -- Portraits.
        McMahon, Tim -- Portraits.
        Songs with piano.
        Popular music -- United States -- 1901-1910.
        Grizzly bear -- Songs and music.
        Ragtime music.
Local Subject  Dance.
        Famous people - Sophie Tucker.
Alt Author  Berlin, Irving, 1888-1989. Lyricist.
        Frew, 1875-1955. Illustrator.
        Tucker, Sophie, 1884-1966. Performer.
        McMahon, Tim. Performer.
Alt Title  Grizzly bear
Note  First line of text: Out in San Francisco where the weather's fair
        First line of chorus: Hug up close to your baby
Alt Title  Frances G. Spencer Collection of American Sheet Music.

THE DANCE OF
The GRIZZLY BEAR
AS SUNG BY
SOPHIE TUCKER

D1732 .15

WORDS BY
IRVING BERLIN

MUSIC BY
GEORGE BOTSFORD

.5

TED SNYDER Co.
MUSIC PUBLISHERS
112 WEST 38 St. NEW YORK

# MARC

```
LEADER 00000ncm  2200000Ia 4500
001    426135815
003    OCoLC
005    20100607104730.0
008    090717s1910    nyurga        n    zxx d
035    (OCoLC)426135815
040    SST|cSST|dIYU
048    vn01|aka01
049    IYUU
099    Spencer D1732 .15
099    Spencer D1732 .15-2
100 1  Botsford, George,|d1874-1949.
245 14 The dance of the grizzly bear /|cwords by Irving Berlin ;
       music by George Botsford.
246 3  Grizzly bear
246 1  |iFirst line of text:|aOut in San Francisco where the
       weather's fair
246 1  |iFirst line of chorus:|aHug up close to your baby
260    New York :|bTed Snyder Co.,|cc1910.
300    1 score (5, [1] p.) ;|c35 cm.
500    For voice and piano.
562    |cD1732 .15: illustrated t. p. in orange, black and white
       with a drawing of 2 bears / Frew; photo. of Sophie Tucker.
562    |cD1732 .15-2: Illustrated t. p. in orange, black and
       white with a drawing of 2 bears / Frew; photo. of Tim
       McMahon.
562    |cD1732.15: "As sung by Sophie Tucker."
562    |cD1732.15-2: "As sung by Tim McMahon."
562    |cD1732.15: Advertisement for Alexander's ragtime band
       march and twostep / Berlin on p. [6].
```

```
562    |cD1732.15-2: Advertisement for Draggy rag and Call me up
       some rainy afternoon / Berlin on p. [6].
590    Frances G. Spencer Collection of American Sheet Music.
590    Spencer subject: Dance - other.
590    Spencer subject: Famous people - Sophie Tucker.
600 10 Tucker, Sophie,|d1884-1966|vPortraits.
600 10 McMahon, Tim|vPortraits.
650  0 Songs with piano.
650  0 Popular music|zUnited States|y1901-1910.
650  0 Grizzly bear|vSongs and music.
650  0 Ragtime music.
690    Dance.
690    Famous people - Sophie Tucker.
700 1  Berlin, Irving,|d1888-1989.|4lyr
700 1  Frew, John.|4ill
700 1  Tucker, Sophie,|d1884-1966.|4prf
700 1  McMahon, Tim.|4prf
793 0  Frances G. Spencer Collection of American Sheet Music.
856 41 |uhttp://digitalcollections.baylor.edu/u?/fa-spnc,32641
       |zConnect to this title in the Baylor University Libraries
       Digital Collections
856 41 |uhttp://digitalcollections.baylor.edu/u?/fa-spnc,30638
       |zConnect to this title in the Baylor University Libraries
       Digital Collections
915    Flourish
```

Gentlemen, you can't fight in here.
This is the War Room!

just-a-happy-camper

# Creating a new project

- Rename and reorder MARC fields

- Join values

- Split values

- Re-format dates

- Remove unnecessary punctuation, delimiters, etc.

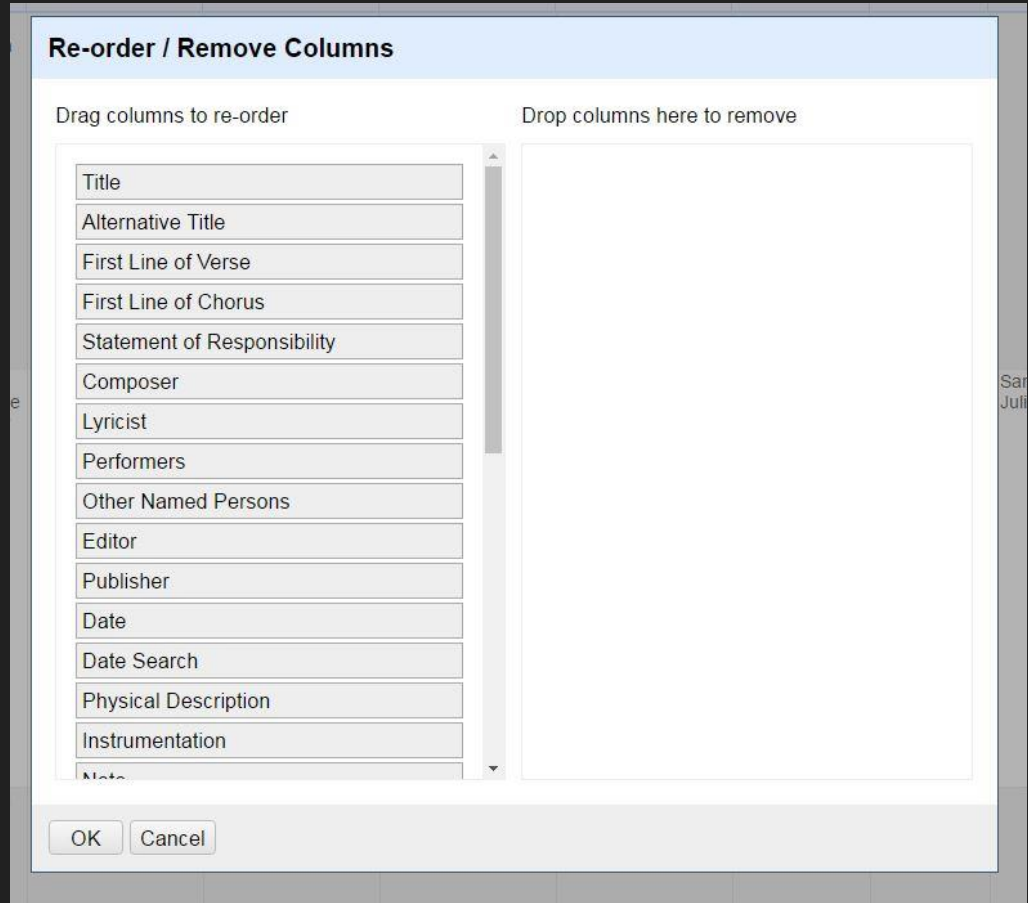- Add new fields for the digital collection

# Renaming columns

Columns are the primary units of interaction. The drop down menu of functions at the column level allows us to rename, reorder, or transform columns.

Column names must exactly match our CDM field names in order for upload the metadata.

MARC 100 →  Composer

# Re-ordering columns

Columns must also exactly match the order that the corresponding fields appear in our CONTENTdm collection. Once all the fields have been re-named, they can be re-ordered under the All columns menu.

# Joining Values

Transform data with Google Refine Expression Language (GREL)

Joining the 245$a and 245$b to create the Title field

Splitting values

The 246 must be split into two or three fields:

- Alternative Title

- First line of verse

- First line of chorus

Splitting values

Know your data!

Extract and save operation history

Apply this to new data sets that need the same kind of clean up

# Identifying clean up in existing CONTENTdm collections

- Text faceting

- Custom text facets

- Identifying duplicates

# Invaluable Resources

http://openrefine.org/

http://freeyourmetadata.org/

https://github.com/OpenRefine/OpenRefine/wiki/GREL-Functions

Verborgh, Ruben, and Max De Wilde. *Using OpenRefine*. Birmingham: PACKT Publishing, 2013.

Van Hooland, Seth, and Ruben Verborgh. *Linked Data for Libraries, Archives, and Museums: How to Clean, Link, and Publish Your Metadata*. Chicago: Neal-Schuman, 2014.