# Sound, to Print, to Cloud: saving the news with the Milo Ryan Phonoarchive Finding Aid

Tigh Bradley, MLIS
University of Washington

## In short

This is a project to transform a print catalog/finding aid into a digital one, with accuracy beyond the power of Optical Character Recognition (OCR) techniques alone.

## The gap

The University of Washington Media Arcade/Reformatting and Digitization Lab is creating digital copies of approximately 1500 hours of audio from the Milo Ryan Phonoarchive. This sound collection includes the largest extent corpus of daily news broadcasts from World War II in the US. There are notable contributions by Edward R. Murrow and "Murrow's Boys", who transmitted news reports live from the theater of war in Europe.

The UW is preserving the sound recordings in a digital format. Rich descriptions of the contents of the tapes are needed for the collection to be discovered and used by researchers and members of the public.

Fortunately, the required information exists in print. However, simply scanning the pages and applying OCR does not provide sufficient quality reproduction of the text. This project utilizes algorithmic processing of the scanned text so that accurate, descriptive digital metadata can be attached to the sound recordings, greatly enhancing the approachability and utility of the historic radio broadcasts.



**Figure 1.** Original print finding aid.



**Figure 2.** New Digital finding aid.

## Methods and Materials

The original print catalog, *History in Sound* by Professor Milo Ryan, comprises some 580 pages of printed text, which has been scanned and processed with OCR. This text contains numerous errors not appearing in the original book. These include extra breaks between lines and omitted or incorrect characters.

The OCR processed text was studied for systemic errors, and Python scripts were implemented to eliminate them. The Python scripts were also utilized to extract metadata features such as program titles and transmission dates from the body text.

## Results

Through a combination of programming and manual inspection of the catalog entries, errors in the text were reduced to fewer than 1 out of 100 catalog fields based on randomized sample tests. The final product is a database of 4,763 digital catalog records of the tapes which faithfully reproduce the text from the print catalog.

| Error Type Corrected | Count |
| --- | --- |
| Extra line breaks | ~7000 |
| Punctuation interrupting a word | ~2000 |
| Lowercase s for uppercase S | ~500 |
| Numeral 0 for letter O | ~500 |
| Wrong character(s) instead of letter | ~250 |
| Bracket { instead of parenthesis ( | ~200 |
| Whitespace instead of letter | ~200 |
| Other | ~200 |

**Table 1.** Summary of errors in OCR text repaired during this project.



**Figure 3.** The original tapes are being digitized via reel to reel machines.

## Discussion

Several steps were necessary for this project:
1) Evaluate the original OCR text by eye.
2) Use of regular expressions to find-and-replace pervasive errors in the text.
3) Programming of scripts to fix more intricate pervasive errors and extract metadata from program descriptions.
4) Human review of catalog entries to assure matching to original text. While algorithmic methods improved the OCR text, about 20% of the entries had to be hand corrected after processing.

## Impact and Future Work

The completed finding aid is being used by the UW Media Arcade/Reformatting and Digitization Lab for already, to catalog the archival holdings of the Milo Ryan Phonoarcive online at ArchiveSpace.

The methods developed in this project may be used to create digital versions of other print catalogs, finding aids, dictionaries, and encyclopedias. They will be especially useful for encoding the second print finding aid for this collection, *History in Sound Part 2*.
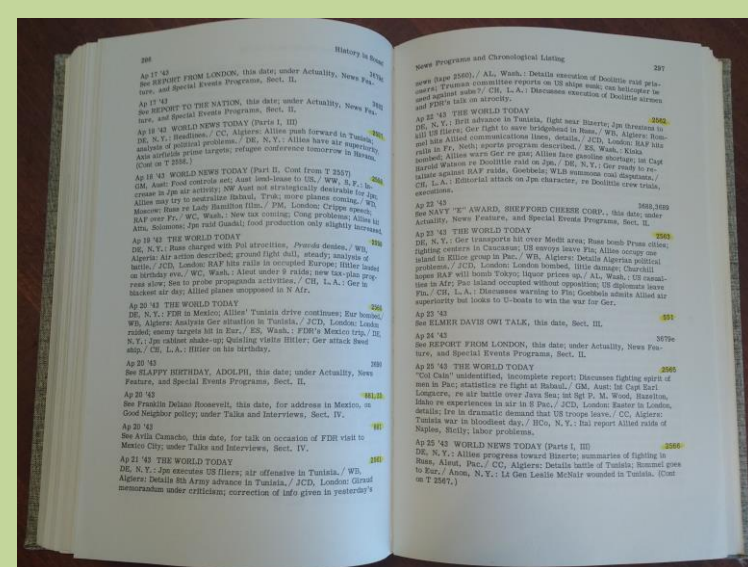
**Contact**

Tigh Bradley
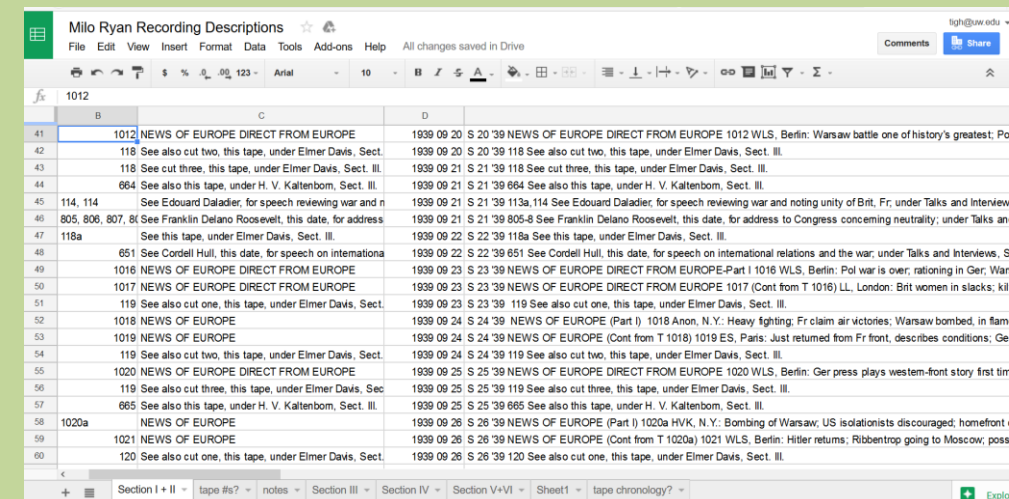University of Washington Information School
Email: tigh@uw.edu