# Rescuing Legacy Transportation Data @ NTL

Lisa Curtin, USDOT Intern & MSIS Candidate 2020
University of Tennessee School of Information Sciences
lisa.curtin2@gmail.com

Mary Moulton, Digital Librarian
National Transportation Library, USDOT
mary.moulton@dot.gov,

U.S. Department of Transportation

## WHY RESCUE OLD DATA?

- Expand temporal and spatial coverage of existing datasets
- Meet open government and open data priorities
- Make data reusable and reduce duplication of effort
- Make data full-text searchable & findable
- Improve accessibility for screen-reader users

## HOW TO RESCUE DATA

- Use commercial Optical Character Recognition (OCR) software such as ABBYY or OmniPage
- Use open-source or design DIY solutions to correct for OCR errors
- Pursue cutting-edge machine learning solutions for data rescue and error correction
- Crowdsourcing to correct OCR errors

## ABSTRACT

Legacy data—data collected or compiled in the past and stored in obsolete formats (e.g., paper and floppy disks)—can be made accessible via the process of data rescue, which may be as simple as scanning documents or, to achieve machine readability, as complex as developing artificial intelligence algorithms. Machine readability not only facilitates the discoverability of data and their reuse for new purposes but also expands the coverage of data sets and supports access for screen reader users. At the National Transportation Library, legacy data exist primarily as scanned-to-PDF documents that, albeit digitized, are not machine-readable due to the quality of scanning. Using Adobe Acrobat Pro's text-recognition function, the Library's current ability to rescue machine-readable data from PDF images has an error rate exceeding 25%. Consequently, in an average table of data, at least 25% of cells require manual correction, which though feasible for occasionally rescuing small datasets is infeasible for large-scale data rescue. As an informal survey of other federal information agencies and a literature review revealed, however, lower error rates are possible with alternative methods of data rescue, including using commercial optical character recognition software and contracting external data rescue service providers.

**LEGACY DATA ARE DATA COLLECTED OR COMPILED IN THE PAST, STORED IN AN OBSOLETE FORMAT.**

**DATA RESCUE IS THE PROCESS OF MAKING LEGACY DATA ACCESSIBLE.**

## WHAT DO OTHER AGENCIES DO?

- The National Science Foundation funds legacy data rescue projects
- The Library of Congress and the Smithsonian utilise public crowdsourcing platforms
- The U.S. Geological Survey uses OmniPage to OCR documents and dabbles in crowdsourcing projects
- The National Library of Medicine's PubMed Central converts documents to HTML
- The National Institute of Standards and Technology contracts with PubMed Central to convert documents to HTML
- The National Library of Education contracts with the Internet Archive for digitization

**NTL's Current Data Rescue Capability Using Adobe Acrobat Pro**



TABLE 27 - LANE MILE TABLE
SAN DIEGO FREEWAY

Original scans are relatively low resolution and often skewed, contain image artifacts and faded text, and in this case lack lines separating columns and rows.



By exporting OCR'd data to MS Word, most data is copied, but many values are incorrect due to artifacting and the lack of column and row separators.

Pasting this data directly into Excel and manually comparing it to the original table reveals a 43% error rate.



With some pre-formatting in MS Word including applying "normal" text formatting, correcting spacing errors, and removing artifacts, the error rate is reduced to 26%.

While correcting a quarter of all values on a single table is not too difficult, this is not a scalable solution for rescuing data on a large scale.