

Humanities data

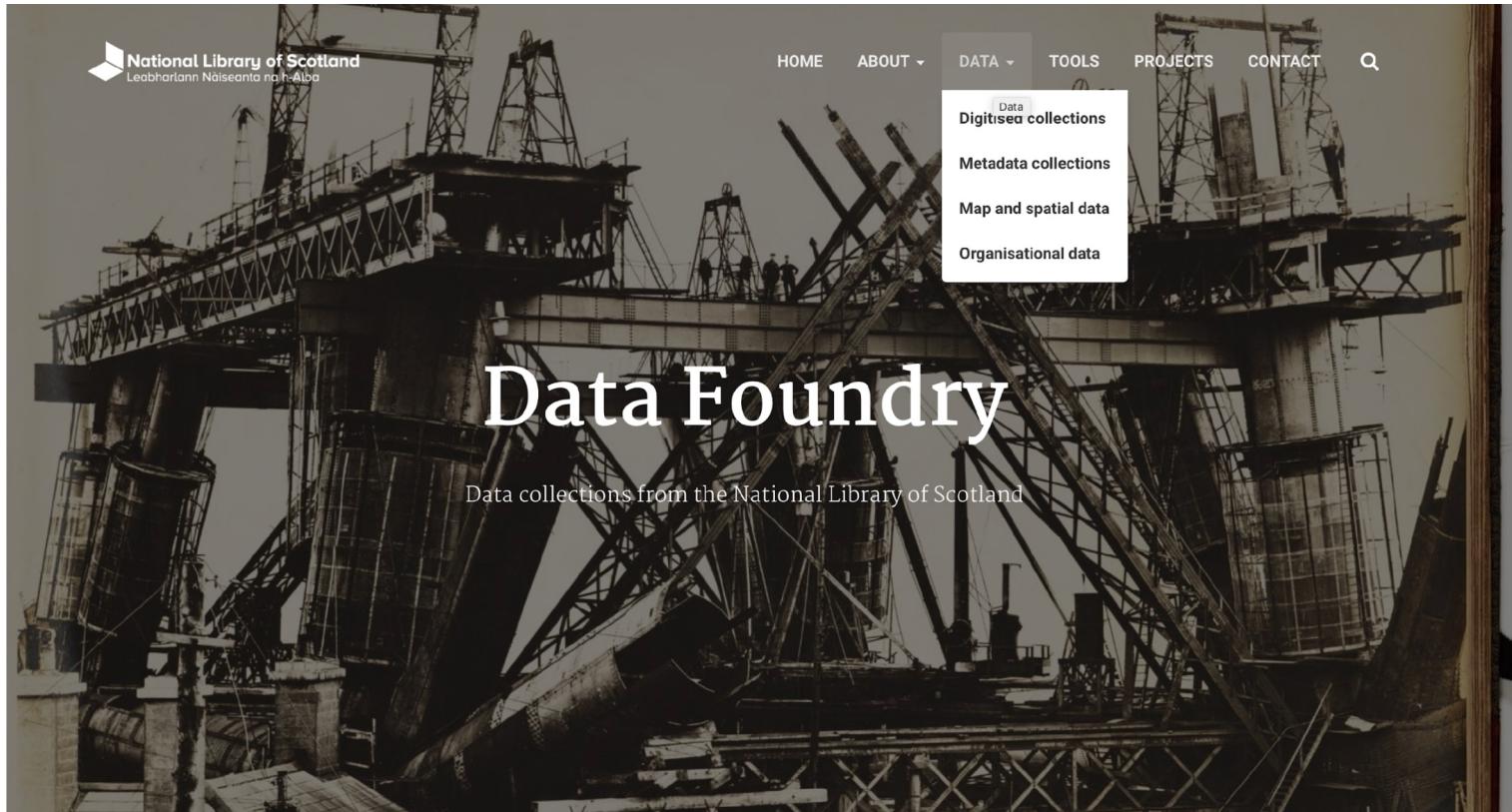
Humanities scholars will sometimes describe elaborate visualizations to me, involving charts and graphs and change over time.

'Great,' I respond. 'Let's see your data.'

'Data?' they say. 'Oh, I don't have any data'.

- Miriam Posner, Humanities Data: A Necessary Contradiction.
<http://mirriamposner.com/blog/humanities-data-a-necessary-contradiction/> (15 July 2015)

Curated cultural heritage data

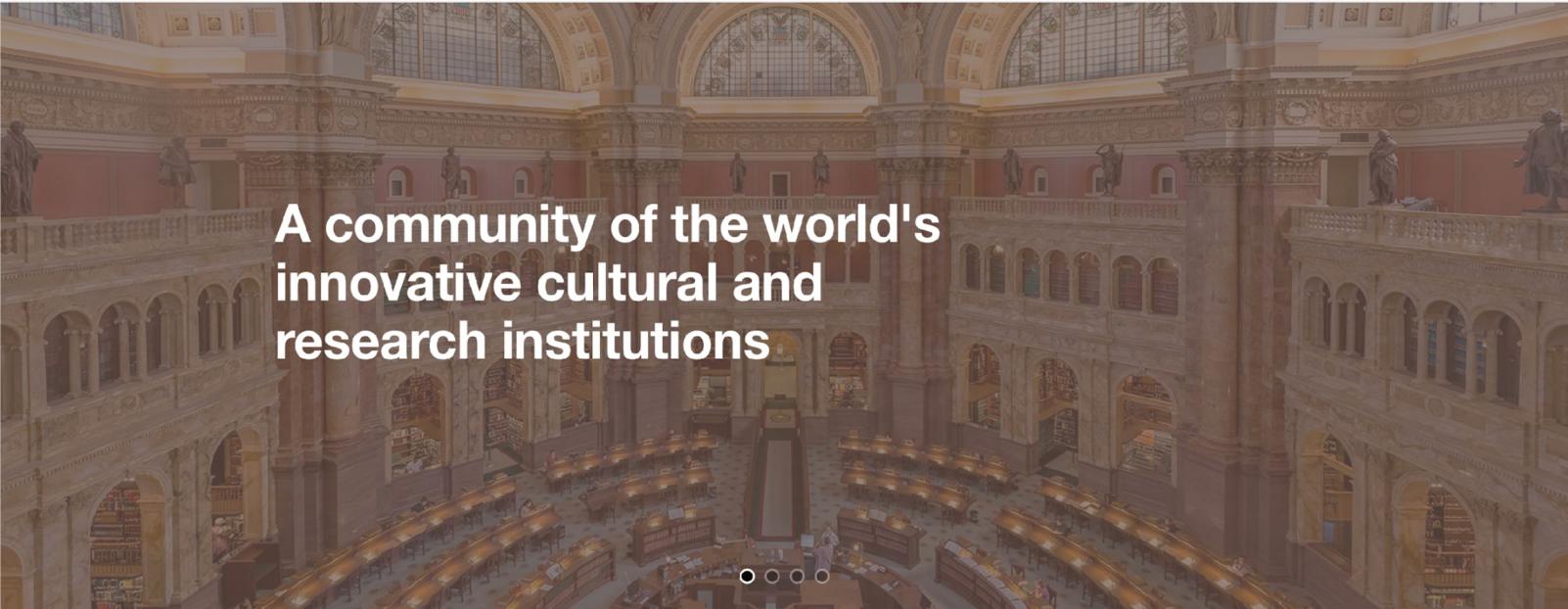


<https://data.nls.uk>

More curated data

**International
GLAM Labs
Community**

[Home](#) [Member Map](#) [Publications](#)



A community of the world's innovative cultural and research institutions

<https://glamlabs.io>

Collections as data

Williams

[Library](#) » [Research Guides](#) » [Collections as Data](#) » [Home](#)

Collections as Data

A portal to Williams Libraries collections as datasets. Explore pre-packaged data or request bespoke datasets for your research.

HOME

[WHAT IS COLLECTIONS AS DATA?](#)

[COLLECTIONS AS DATA VALUES](#)

[WORKING WITH COLLECTIONS AS
DATA](#)

[DATA ABOUT AND FROM WILLIAMS
LIBRARIES](#)

[HOW TO CITE OUR DATA](#)

[ASK A QUESTION OR REQUEST DATA](#)

What is Collections as Data?

Collections as Data is the concept of using digital collections, digital objects, and their descriptive metadata, made available as datasets, to perform computational analysis.

Examples:

- A text file or corpus comprised of all Williams unrestricted theses for a particular academic department
- All digital issues of the Williams Record
- Metadata describing translated books in the catalog along with which languages they were translated from/to and how many times they were checked out

Collections as Data might consist of the content of the library, archival resources, unrestricted records of the College themselves or might consist of information about those resources. From this data, a researcher may be able to derive quantitative measures, new datasets, or even predictive models.

Often these datasets need to be extracted from our systems and cleaned or reformatted to make them usable for different types of analysis. If you are interested in data that is not already available on-demand, please see the

Catalogues as data

National bibliographies have been identified as a crucial resource for historical research on the publishing landscape, but using them requires addressing challenges of data quality, completeness, and interpretation. We call this approach bibliographic data science.

Lahti, Leo, Jani Marjanen, Hege Roivainen, and Mikko Tolonen. "Bibliographic Data Science and the History of the Book (c. 1500–1800)." *Cataloging & Classification Quarterly* 57, no. 1 (2019): 5-23.

Catalogues as data

The use of bibliographic metadata as a research object has, however, proven to be challenging [...] [The] Primary motivation for cataloging has been to preserve as much information of the original document and its physical creation as possible. This includes potential errors caused by the printer. If, for instance, a place name is wrongly spelled on the title page, for cataloging purposes it is relevant also to preserve that misspelling. For anyone desiring to work on quantitative approach to bibliographic metadata, this is a crucial point to understand and respect.

.

Lahti, Leo, Jani Marjanen, Hege Roivainen, and Mikko Tolonen. "Bibliographic Data Science and the History of the Book (c. 1500–1800)." *Cataloging & Classification Quarterly* 57, no. 1 (2019): 5–23.

Against data as such

... we re-examine the intellectual foundations of digital humanities... first and foremost that we reconceive all data as *capta*..

Capta is 'taken' actively while data is assumed to be a given, able to be recorded and observed. From this distinction, a world of differences arises. Humanistic inquiry acknowledges the situated, partial, and constitutive character of knowledge production, the recognition that knowledge is constructed, taken, not simply given as a natural representation of pre-existing fact.

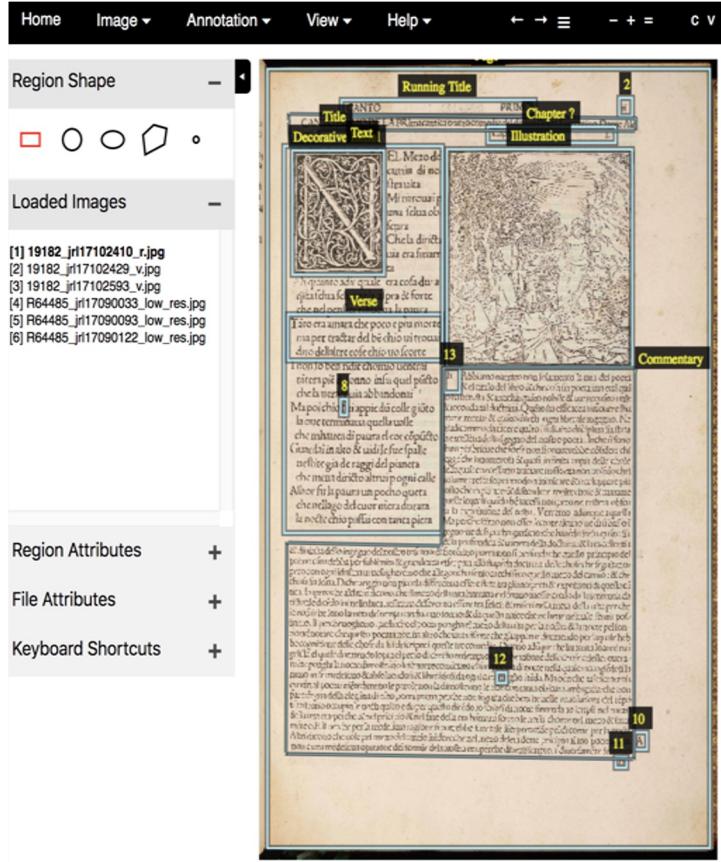
- Johanna Drucker, "Humanities approaches to graphical display." *Digital Humanities Quarterly* 5, no. 1 (2011): 1-21

Texts as data vs. the material text

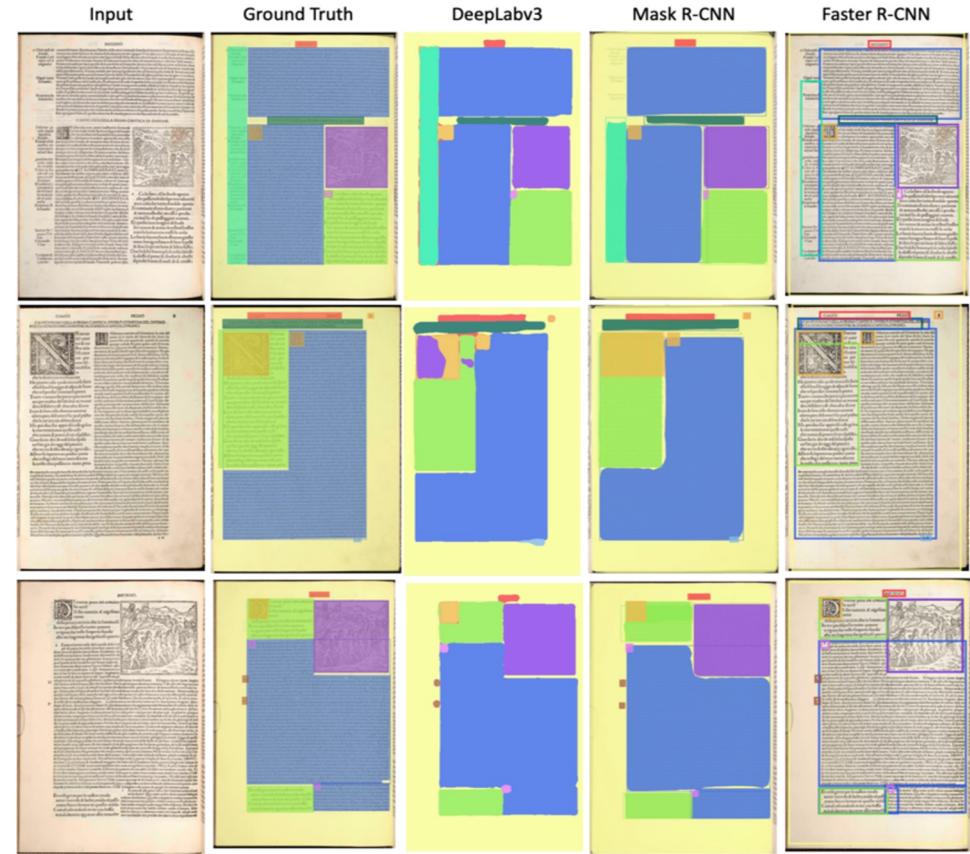
Digital reproductions are central to our work as scholars and teachers [but] many scholars have cautioned against treating digital images as surrogates for the original, citing the loss of texture, size, smell, colour and context. With the proliferation of more and higher-resolution reproductions online, however, subtle but interesting variations between copies at different institutions have become more visible, ironically helping foment copy-specific work in book history and thus spurring on the digitization of even more material.

- Zachary Lesser and Whitney Trettien. "Material / digital." In *Shakespeare / Text: Contemporary Readings in Textual Studies, Editing and Performance*, edited by Claire M. L. Bourne , 402–423. The Arden Shakespeare, 2021

AHRC-funded Envisioning Dante project



Creating training data with VGG VIA – credit: Charlotte Alton



Testing various machine learning models – Suny Shtedritski, with Andrea Micaldi

Digitisation and its outputs

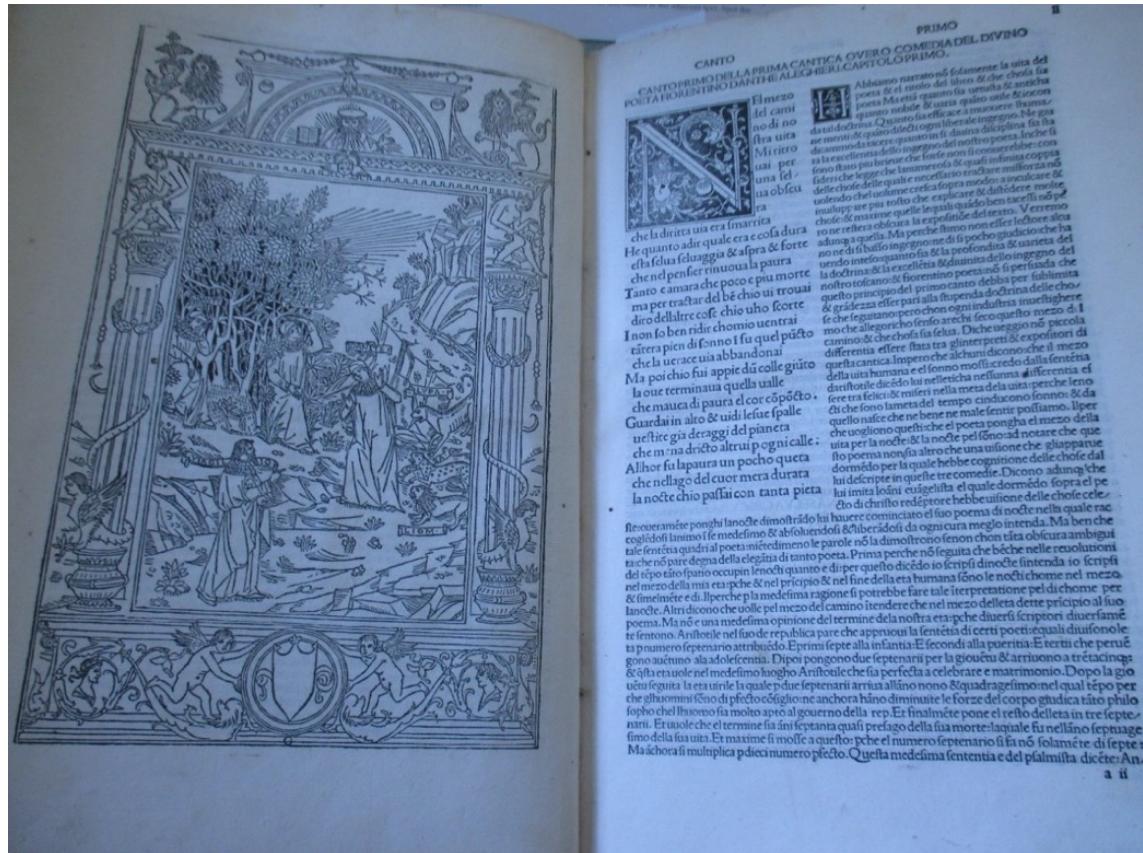


Copy to Dropbox Download

ENV DANTE / TIFF / 17280
from James Robinson (The University of Manchester)

| | | | | |
|-----------------------------|-----------------------------|-----------------------------|-----------------------------|---------------------------|
| JRL231751283.tif 89.9 MB | JRL231751284.tif 82 MB | JRL231751285.tif 85.1 MB | JRL231751286.tif 75.4 MB | JRL231751287.tif 77 MB |
| JRL231751288.tif 77.3 MB | JRL231751289.tif 78.2 MB | JRL231751290.... 78.1 MB | JRL231751291.tif 78.8 MB | JRL231751292.tif 78 MB |
| [Download icon] | [Download icon] | [Download icon] | [Download icon] | [Download icon] |

Digitising the page: a caution



Oxford Taylor Institute Library 4 Dante Arch. Fol. It. 1491

We only occasionally see one page of a book at a time; the two pages making an opening are really the unit of the book, and this was thoroughly understood by the old book producers.

- William Morris, 'The Ideal Book', *The Library*, 1893, 179-186 (183)

3000+ digitized chapbooks



Original OCR: no
clean-up



47,329 ALTO XML
files at page level



47,329 image files



METS metadata
files at item level



1,432,928 lines
and 10,818,763
words

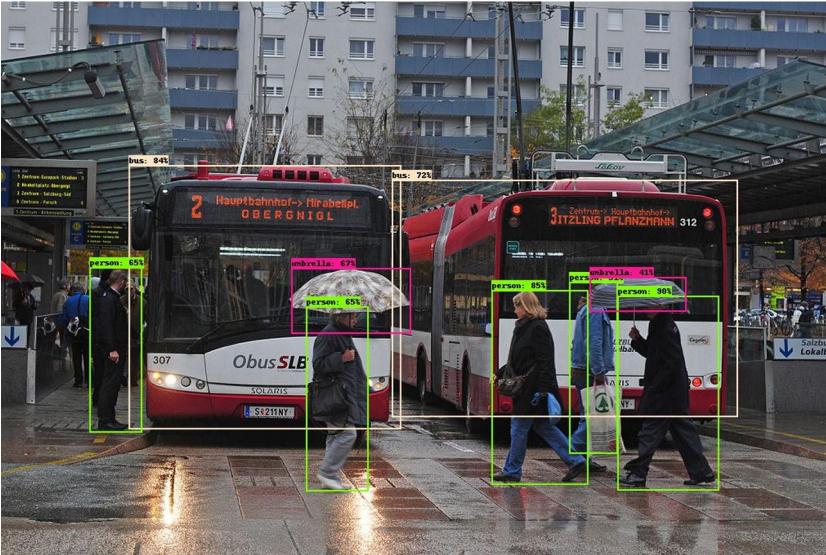


Covers period
c1700-c1899

data.nls.uk/data/digitised-collections/chapbooks-printed-in-scotland

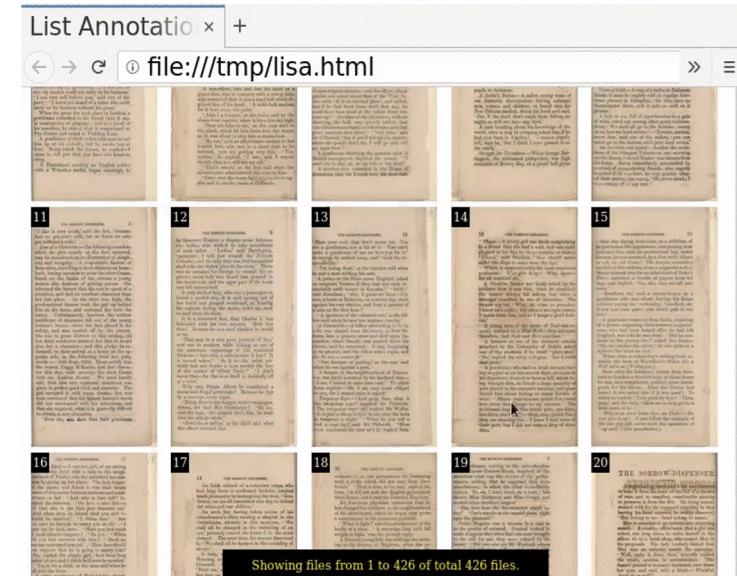
Object Detection

1. We select an object-detection model

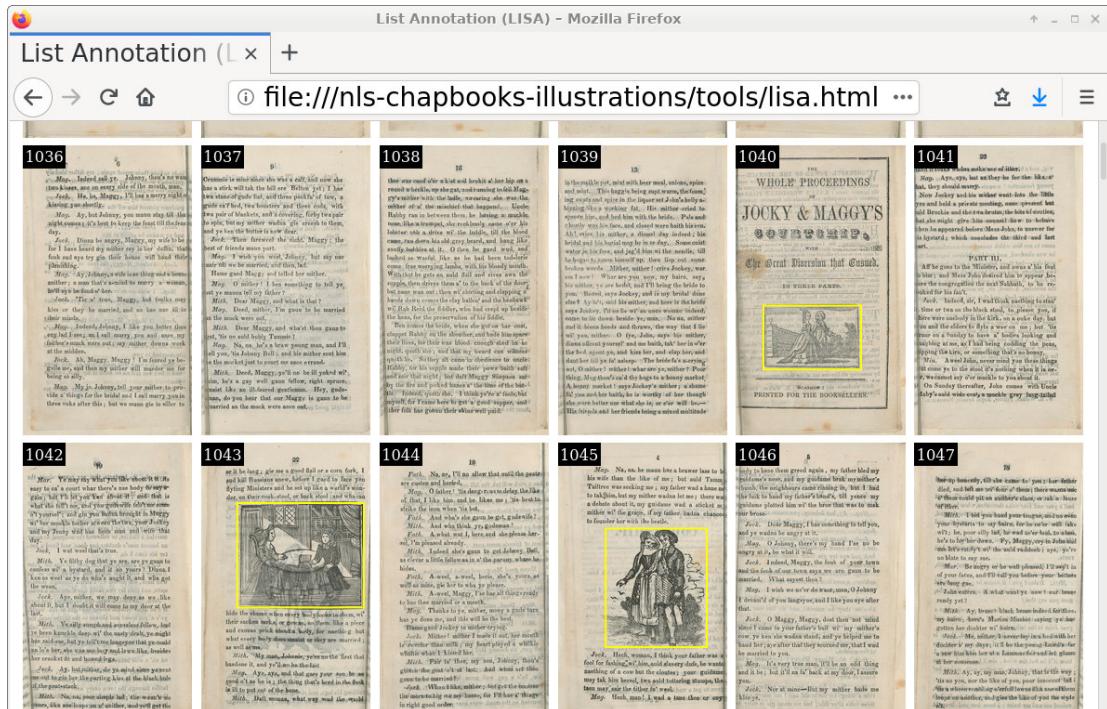


Detecting people, buses and umbrellas with the EfficientDet object detector trained on the COCO dataset

2. We retrain with VGG LISA on book illustrations



3. We successfully predict 3600 illustrations



4. We test and release the retrained model for reuse

The screenshot shows the 'nls-chapbooks-illustrations' repository on GitLab. The command used to run the illustration detector is:

```
python model_inspect.py --runmode=saved_model_infer \
--model_name=efficientdet-d0 \
--saved_model_dir=$MYSTORE/nls-chapbooks-illustrations/models/full-model-training/saved_model \
--input_image=$MYSTORE/nls-chapbooks-illustrations/data/images/test_images/*.jpg \
--hparams=$MYSTORE/models/full-model-training/hparams.yaml \
--output_image_dir=$MYSTORE/detections/
```

The following image shows the result of applying illustration detector on a new image (i.e. image not seen during training):

Cross Validation to Estimate Performance of Illustration Detector

Illustration Detector in a Python notebook

The screenshot shows a Google Colab notebook interface. The title bar reads "Illustration Detection for Chapbooks Printed in Scotland.ipynb". The left sidebar shows a file tree with a single folder named "sample_data". The main area is divided into two tabs: "+ Code" and "+ Text". The "+ Text" tab is active and contains the following content:

1. Download and Install the Required Tools

The illustration detector developed in the [VGG Chapbooks Project](#) is based on the [EfficientDet](#) object detector. The VGG Chapbooks Project code repository contains all the data, pre-trained object detector and tools required in this tutorial. Therefore, we download the [code repository](#) and setup the environment in this colab document. This setup is essential for all the remaining sections of this tutorial and therefore must be executed before running commands from any other section.

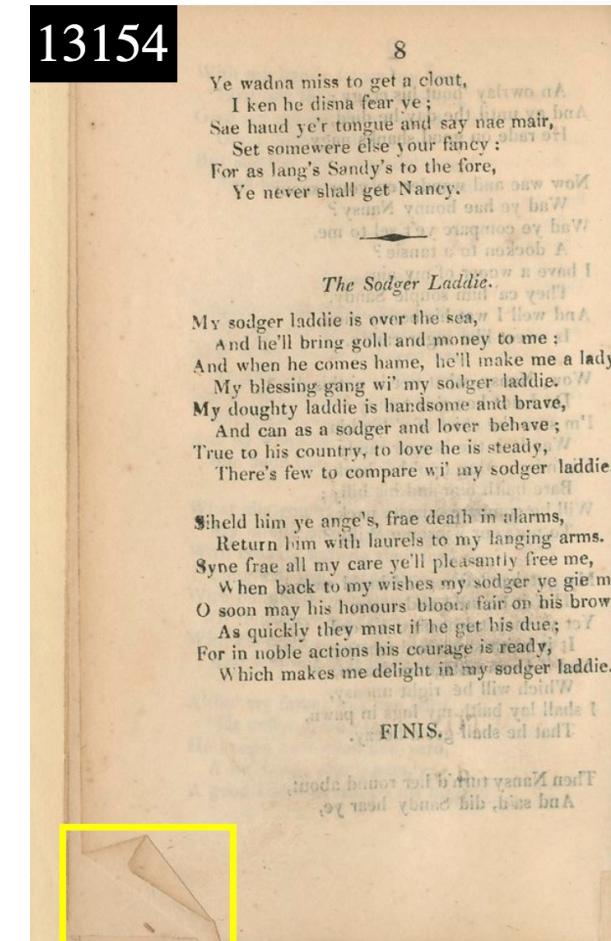
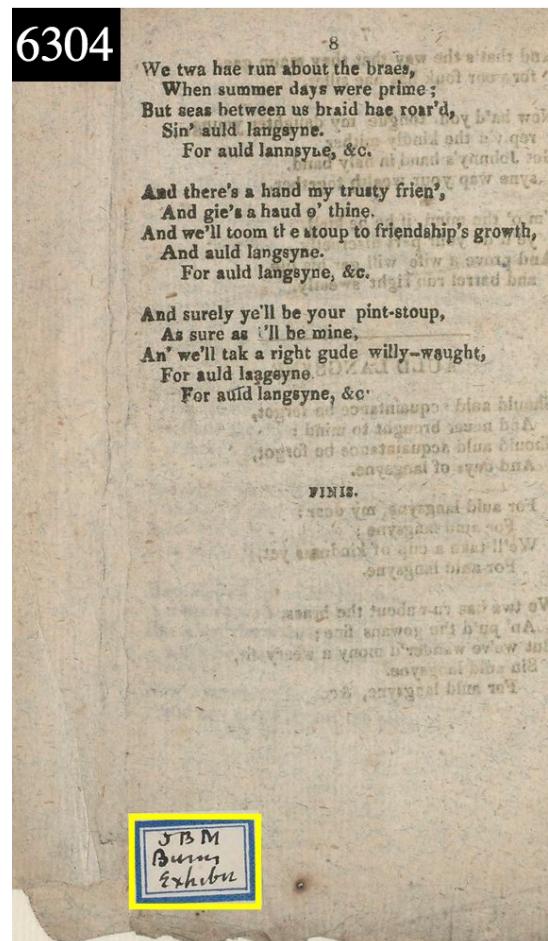
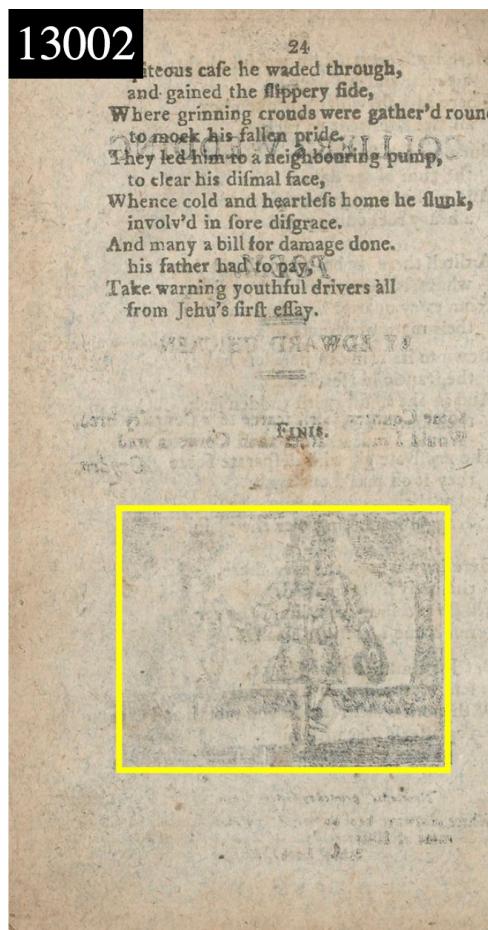
```
[ ] ## Download VGG Chapbooks project code repository and setup environment
import os
import sys
import tensorflow.compat.v1 as tf
import cv2
from google.colab.patches import cv2_imshow
import datetime
import json

if 'nls-chapbooks-illustrations' not in os.getcwd():
    !git clone --recurse-submodules https://gitlab.com/vgg/nls-chapbooks-illustrations.git
    os.chdir('nls-chapbooks-illustrations/autonl/efficientdet')
    !git pull origin master # update EfficientDet code to the latest version
    !pip install -r requirements.txt
    !pip install -U 'git+https://github.com/cocodataset/cocoapi.git#subdirectory=PythonAPI'

## We define some utility functions that will be used throughout this tutorial
def show_thumbnail(img_fn, tsize=500):
    ...
    Show a thumbnail sized version of an image in Colab
    ...
    img = cv2.imread(img_fn)
    w, h, c = img.shape
    if w > tsize or h > tsize:
        if w > h:
            new_width = tsize
            new_height = int( (w/h) * new_width )
```

The bottom status bar shows "Disk" with a progress bar at 40.86 GB available.

False positives – a printed offset, a provenance stamp and a dog-ear





Visual Analysis of Chapbooks Printed in Scotland

Abhishek Dutta, Giles Bergel and Andrew Zisserman

Overview

Chapbooks were short, cheap printed booklets produced in large quantities in Scotland, England, Ireland, North America and much of Europe between roughly the seventeenth and nineteenth centuries. A form of popular literature containing songs, stories, poems, games, riddles, religious writings and other content designed to appeal to a wide readership, they were frequently illustrated, particularly on their title-pages. This paper describes the visual analysis of such chapbook illustrations. We automatically extract all the illustrations contained in the National Library of Scotland Chapbooks Printed in Scotland dataset, and create a visual search engine to search this dataset using full or part-illustrations as queries. We also cluster these illustrations based on their visual content, and provide keyword-based search of the metadata associated with each publication. The visual search, clustering of illustrations based on visual content, and metadata search features enable researchers to forensically analyse the chapbooks dataset and to discover unnoticed relationships between its elements. We release all annotations and software tools described in this paper to enable reproduction of the results presented and to allow extension of the methodology described to datasets of a similar nature.



<https://www.robots.ox.ac.uk/~vgg/research/chapbooks/>

The problem of sustaining DH data

Project name

Arts and Humanities Data Service: Enabling Digital Resources for the Art and Humanities

Project principal investigator(s)

Sheila Anderson

Decommission Date

September 2017

Archive URL(s)

<http://www.ahds.ac.uk/>

Additional links

[Internet Archive](#)

Overview

The Arts and Humanities Data Service was a UK national service funded by the JISC and AHRC to collect, preserve and promote the electronic resources which result from research and teaching in the arts and humanities.

<https://ahds.ac.uk>

The problem of sustaining humanities data



Holding page: British Book Trade Index

The Bodleian Libraries have taken The British Book Trade index resource offline, owing to a change in guidance within the University relating to cyber security.

The British Book Trade Index resource has therefore been taken offline as a precaution, in the light of the new guidance, while we develop new approaches to being able to support and deliver it.

We acknowledge how disruptive this has been to the many scholars and communities who use this resource. Alternative ways to access The British Book Trade Index resource are below, while we determine routes and funding to take the resource forward.

More on the background here: [Bodleian Service Updates](#)

An export of the BBTI database and John Feather's Bibliography is available via the [Oxford University Research Archive](#).

The above dataset has been redeveloped [as a stand-alone site](#) by Prof. Janelle Jansted and Martin Holmes at University of Victoria.

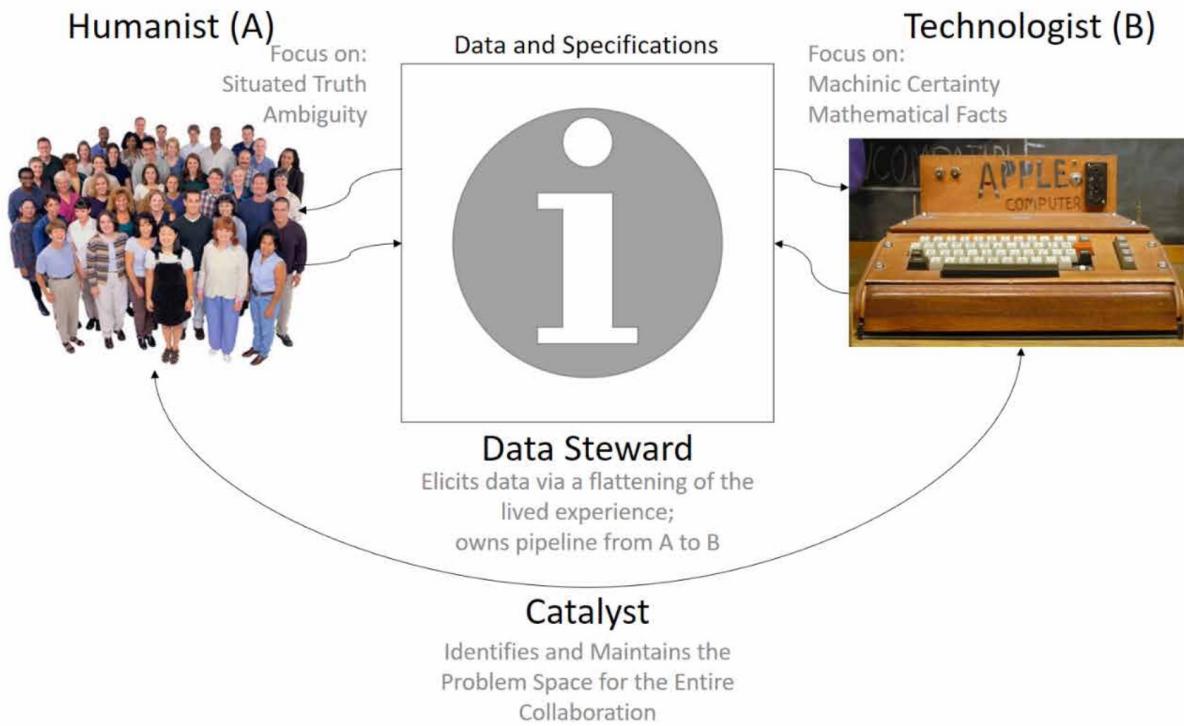
The dataset has also been redeveloped [as a stand-alone site](#) by Prof. Allison Muri of the University of Saskatchewan.

Users may also wish to consult the [Dicey & Marshall catalogue](#).

Please contact digitalsupport@bodleian.ox.ac.uk with any questions.

<https://bbti.bodleian.ox.ac.uk>

Data stewardship



We propose that effective collaborations in Digital Art History require more than just a humanist and a technologist to succeed. Indeed, we find that there are four different roles that need to be filled: Humanist, Technologist, Data Steward, and Catalyst.

It is important to note here that we are talking about roles, not people. Two or more of these roles can be performed by the same person.

Berg-Fulton, Tracey, Alison Langmead, Thomas Lombardi, David Newbury, and Christopher Nygren. "A Role-Based Model for Successful Collaboration in Digital Art History." *International Journal for Digital Art History: Issue 3, 2018: Digital Space and Architecture 3* (2019): 153

The role of the data steward

The Data Steward is responsible for ensuring that the essential character of the historical evidence is not lost throughout the process of converting primary source material into computable information, and also that this data will be suitable for the technological processes that it will undergo.

By designating this role as the party responsible for surfacing both technological and historical assumptions throughout the process of preparing the data for computation, we make explicit the requirement to observe and discuss the constraints and limitations of this entire process.

Berg-Fulton, Tracey, Alison Langmead, Thomas Lombardi, David Newbury, and Christopher Nygren. "A Role-Based Model for Successful Collaboration in Digital Art History." *International Journal for Digital Art History: Issue 3, 2018: Digital Space and Architecture* 3 (2019): 153

Further reading

Alan Galey, 'The Enkindling Reciter: E-Books in the Bibliographical Imagination,' *Book History* 15 (2012): 210-47

Martin Paul Eve, 'You have to keep track of your changes': The Version Variants and Publishing History of David Mitchell's Cloud Atlas, *Open Library of Humanities* 2, no. 2 (2016): 1-34

Bonnie Mak, 'Archaeology of a Digitization,' *Journal of the Association for Information Science and Technology* 65, no. 8 (2014): 1515–26

Zachary Lesser and Whitney Trettien, 'Material / digital,' In *Shakespeare / Text: Contemporary Readings in Textual Studies, Editing and Performance*, edited by Claire M. L. Bourne, (London: The Arden Shakespeare, 2021), 402–423

Whitney Trettien, 'The History of the Early Modern Book in the Digital Age,' in Adam Smyth (ed.), *The Oxford Handbook of the History of the Book in Early Modern England*, Oxford Handbooks (2023; online edn, Oxford Academic, 18 Sept. 2023)

Ryan Cordell, "Q i-tjb the Raven": Taking Dirty OCR Seriously,' *Book History* 20 (2017): 188–225

Matthew Kirschenbaum and Sarah Werner. "Digital scholarship and digital studies: The state of the discipline." *Book History* 17, no. 1 (2014): 406-458