# Overview

**Day 1 (9 Oct.)**
- Introduction to OCR
- Challenges of OCR in practice
- Activity: Ready-made OCR tools

**Day 2 (16 Oct.)**
- Review of programming basics
- Introduction to OCR packages
- Working with OCR in Python & R

Course Material

# What Is OCR?

- **O**ptical **C**haracter **R**ecognition
- OCR is the technique to process images of text, such as written or printed documents, and produce machine-readable documents.
- Machine-readable documents are encoded in formats that computers can process, allowing the text to be searched, edited and analysed.
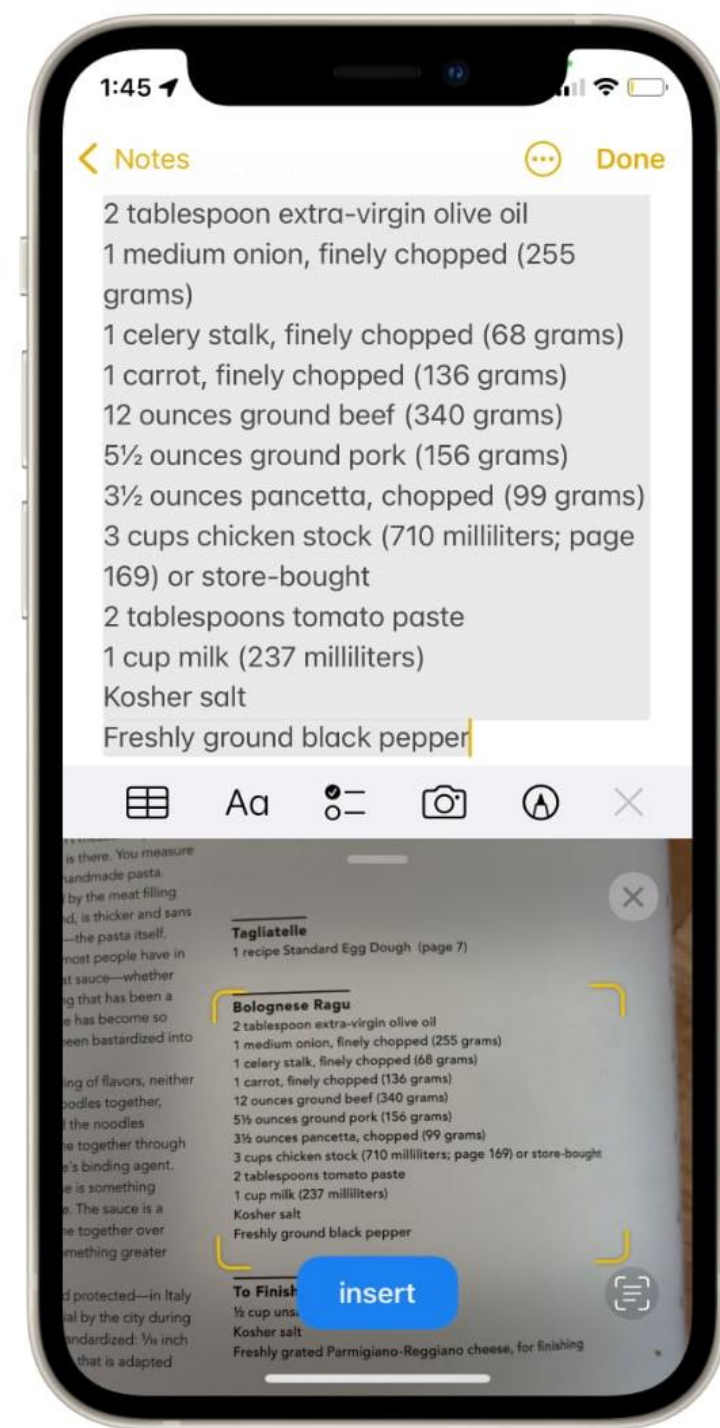
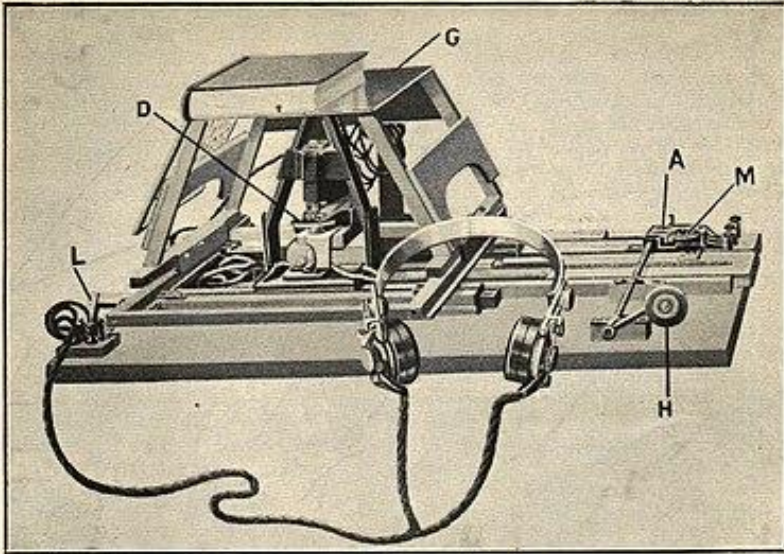# OCR In the Wild
# Real-World Examples

- Scanning your passport at the airport.

- Generate machine-readable text for text-to-speech technology.

- Making digitalised physical archives searchable.

- Creating a dataset of for text mining or text analysis.



www.cdcs.ed.ac.uk

THE UNIVERSITY of EDINBURGH
Centre for Data, Culture & Society

# OCR Timeline: A Brief History

**Early Foundations**
1914: Emanuel Goldberg invents a machine that converts text into telegraph code; Edmund Fournier d'Albe invents the Optophone, which converts text to audio tones matched with letters

1931: Goldberg invents the "Statistical Machine," which searches archives through optical code recognition



Kurzweil Reading Machine (circa 1978)

**Mid-Century Machinery**
1974: Ray Kurzweil develops the "Reading Machine," an early form of OCR technology compatible with text-to-speech software
1976: Kurzweil sells to Xerox

**The Internet Era**
1990s-00s: OCR goes mainstream; improved accuracy; mass digitization projects
00s-present: increased accessibility and accuracy

# OCR workflow

1. Select the images to be scanned
2. Scan the images with OCR software
3. Inspect and 'clean' the processed files
4. Save the results for further use

28/02/2017

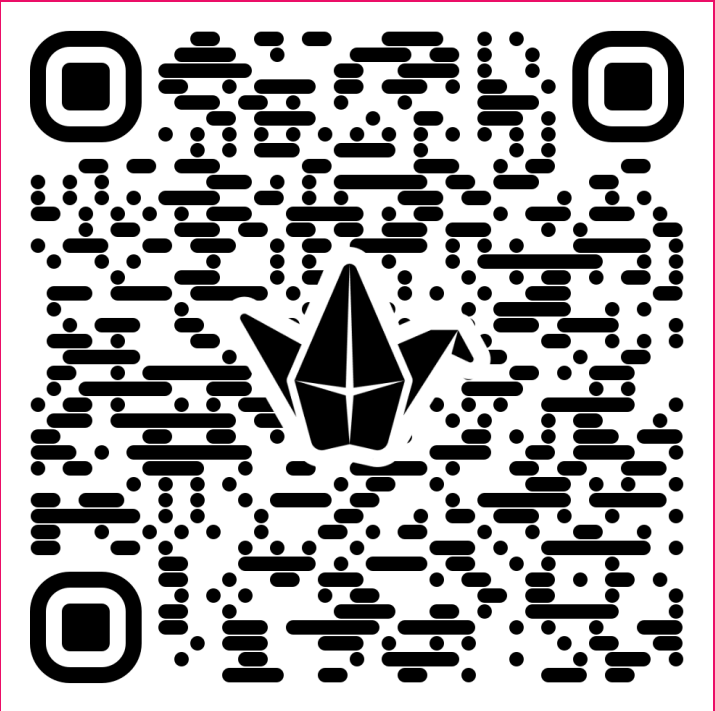# Optical character recognition

From Wikipedia, the free encyclopedia

**Optical character recognition** (also **optical character reader, OCR**) is the mechanical or electronic conversion of images of typed, handwritten or printed text into machine-encoded text, whether from a scanned document, a photo of a document, a scene-photo (for example the text on signs and billboards in a landscape photo) or from subtitle text superimposed on an image (for example from a television broadcast).[1] It is widely used as a form of information entry from printed paper data records, whether passport documents, invoices, bank statements, computerised receipts, business cards, mail, printouts of static-data, or any suitable documentation. It is a common method of digitising printed texts so that they can be electronically edited, searched, stored more compactly, displayed on-line, and used in machine processes such as cognitive computing, machine translation, (extracted) text-to-speech, key data and text mining. OCR is a field in pattern recognition, artificial intelligence and computer vision.

Early versions needed to be trained with images of each character, and worked on one font at a time. Systems capable of producing a high degree of recognition accuracy for most fonts are now common. Some systems are capable of reproducing formatted output that closely approximates the original page including images, columns, and other non-textual components.

[2] Support for a variety of digital image file format inputs.

# Activity: OCR Workflow



1. Identify a dataset (images of text) that you might use in your research

2. Write the steps to obtained encoded text from your dataset. (e.g. software, OCR module)

3. Share your dataset, workflow, and plan with your small group.

4. Identify potential issue in each step. And discuss how you might address them.

# Challenges of OCR

Accuracy depends on **dataset quality, visual complexity** and **software capability**. Some common sources of error include:

- Human errors and typos

- Age and damage (stained or blurry)

- Mixed text and images, or multiple languages

- Cursive handwriting

# Challenges of OCR

Possible solutions:

- Select good quality dataset to begin with.

- Pre-process your dataset to improve its quality

- Correct errors in OCR-produced files, if they are predictable.

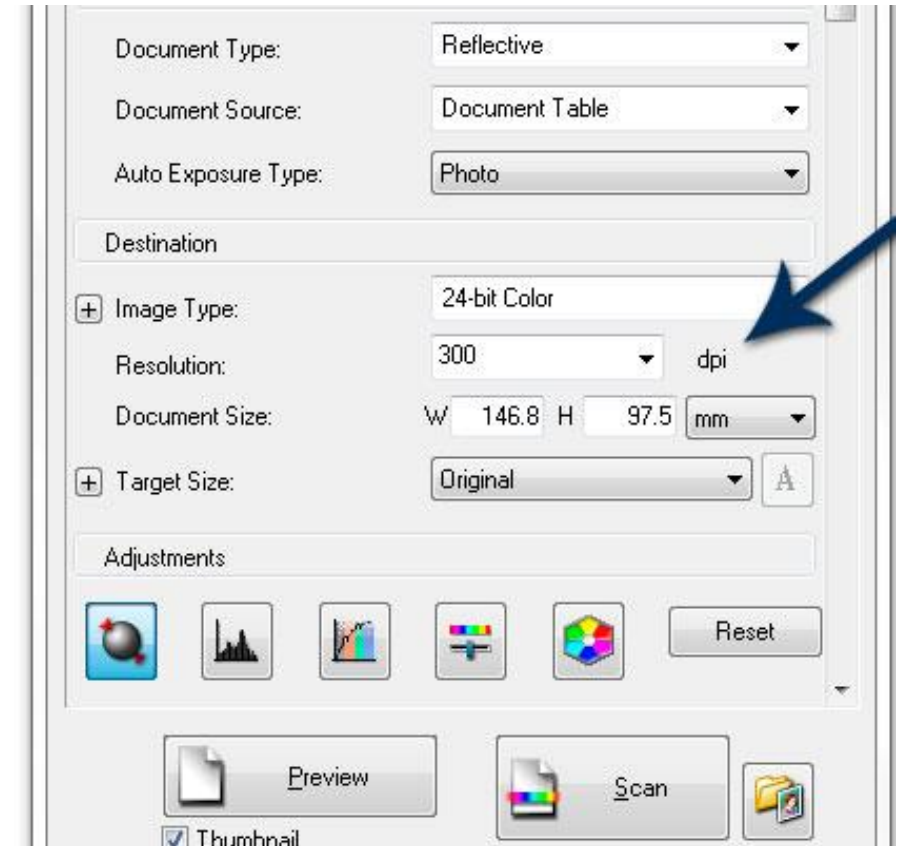- Improve OCR engine capability

# Data Selection

## Image resolution

300 DPI is often used as a benchmark for good quality printing reproducibility for photographs, but this may vary
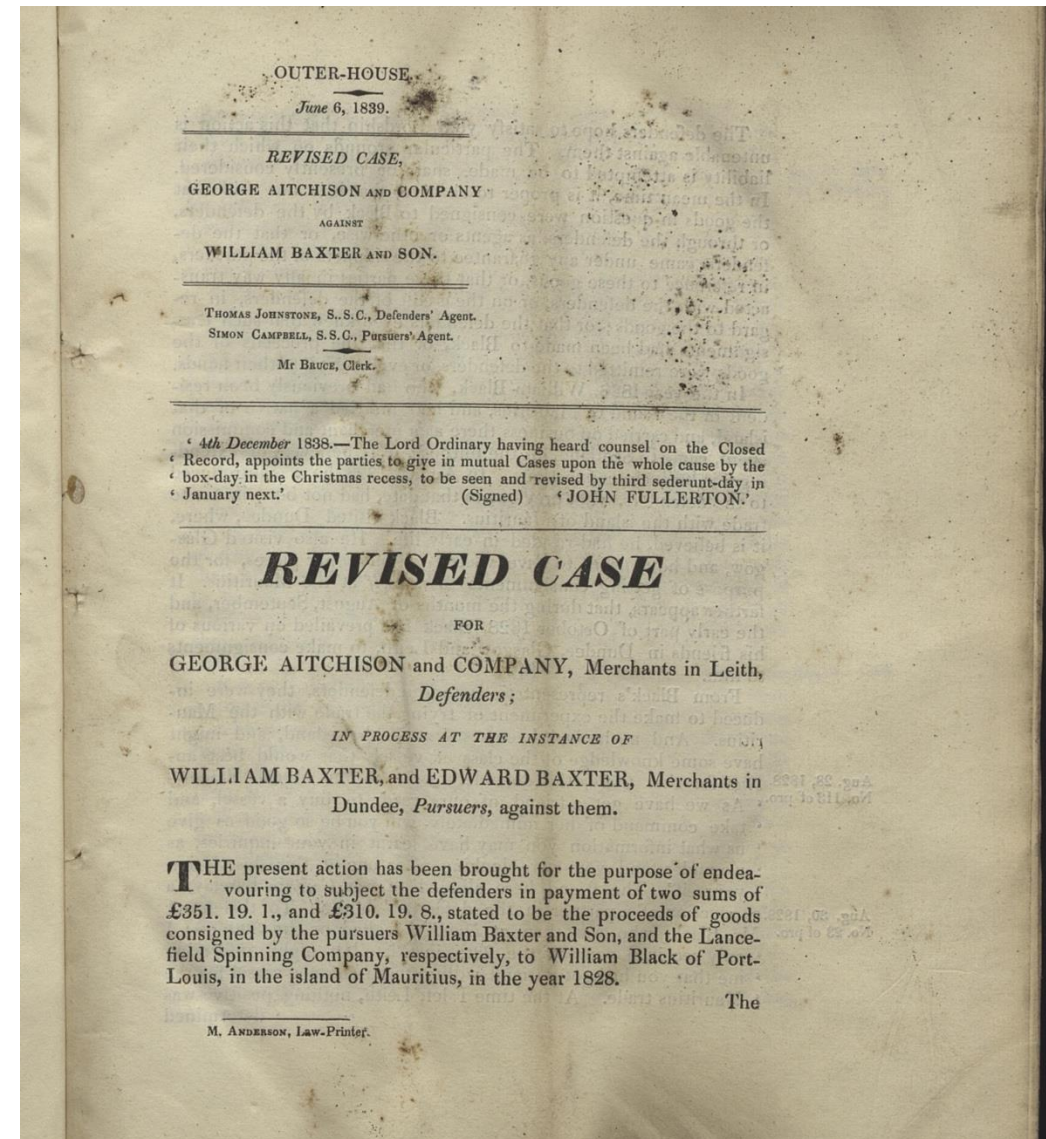
Ex. National archive

# Data Selection

## Manual inspection

# Activity:
# Inspect the Images

- Identify issues you might encounter when processing the following images.

- Are there any steps you could take to preprocessing the document that might improve the output accuracy?

# Pre-processing

- **Manually**: Adjust the colour, contrast, or formats of the files.

- **OCR tools**: Some OCR software have built-in correction functions such as patches.

# Scanning

Pay attention to the limitations of the software

- **File Size**

- **File Format**

- **Text orientation**

- **Languages**

# Activity: Scanning

- https://tools.pdf24.org/en/ocr-pdf

- https://www.onlineocr.net/

- https://www.sodapdf.com/ocr-pdf/

- https://www.sejda.com/ocr-pdf

- https://ocr.space/

- https://avepdf.com/pdf-ocr

1. Identify a picture or scanned pdf

2. Scan the documents with OCR software

3. Compare the results and discuss what are the advantages and limitations of these options?

(if you can't find one, please try the pictures in the previous slides, or you can also try the Edinburgh archive database:

https://archives.collections.ed.ac.uk/)

# **Cleaning**

- Manually remove errors

  Predictable errors can be fixed with codes , such as Regex.

# Regex

- **Concept** used in many **different programming environments** for **pattern matching**.

- Powerful tool **to find, manage, and transform data and files.**

- Use a sequence of characters to define a **search to match strings**
  - **Match on types of characters** (e.g. 'upper case letters', 'digits', 'spaces', etc.).
  - **Match patterns** that repeat any number of times.

www.cdcs.ed.ac.uk

# Regex

| Symbol | Description | Symbol | Description |
|--------|-------------|--------|-------------|
| ^ | Start of line + | ? | 0 or 1 + |
| \A | Start of string + | {3} | Exactly 3 + |
| $ | End of line + | {3,} | 3 or more + |
| \Z | End of string + | {3,5} | 3, 4 or 5 + |
| \b | Word boundary + | \ | Escape Character + |
| \B | Not word boundary + | \n | New line + |
| \< | Start of word | \r | Carriage return + |
| \> | End of word | \t | Tab + |
| \s | White space | . | Any character except new line (\n) + |
| \S | Not white space | (a\|b) | a or b + |
| \d | Digit | [abc] | Range (a or b or c) + |
| \D | Not digit | [^abc] | Not a or b or c + |
| \w | Word | [0-7] | Digit between 0 and 7 + |
| \W | Not word | [a-q] | Letter between a and q + |
| * | 0 or more + | [A-Q] | Upper case letter + between A and Q + |
| + | 1 or more + | | |

- https://programminghistorian.org/en/lessons/cleaning-ocrd-text-with-regular-expressions

- https://programminghistorian.org/en/lessons/understanding-regular-expressions

- https://librarycarpentry.org/lc-data-intro/01-regular-expressions/

# Printed Text Recognition using Python pytesseract

Tesseract is an OCR engine developed for various operating systems

- It has become open-source since 2005
- It is of the most accurate open-source OCR engines available
- Originally only support English, but more languages have been added
- Can run from command line interface or embedded in main coding languages (R & Python)

DATA CULTURE SOCIETY

# Handwritten Text Recognition with Python trOCR

- Select good quality dataset to begin with
- Pre-process your dataset to improve its quality
- Correct errors in OCR-produced files, if they are predictable.
- Improve OCR engine capability

THE UNIVERSITY of EDINBURGH
Centre for Data, Culture & Society

www.cdcs.ed.ac.uk

[Digital Scholar Lab](https://...)