

Experiments with text extraction



(Sprint 6)

The task

6 Done

!

Spike an OCR demonstration using Tesseract library

#40 opened by piriypordes

!

Discuss approach for prototyping

#28 opened by piriypordes

!

Initial discussion of IIIF with Sys Dev

#39 opened by piriypordes

!

Write up findings from the workshop

#27 opened by piriypordes

!

Talk to CEE about image and document types

#30 opened by piriypordes

!

Work out how to work with Findability team

#31 opened by piriypordes

So, what's a 'spike' exactly?

*“We carry out a Spike when we're about to embark on a new piece of work and are finding it hard to **make decisions about direction in business terms, technical terms, or both**. Similar in concept to a hackathon, a Spike is an investigatory piece of work which should:*

- *Answer **a single question***
- *Either be technical or customer-focused*
- *Reduce uncertainty and create a way forward*

*The key thing to note is that **a Spike doesn't directly contribute to an increment in the working software**. As a result, we don't estimate spikes, instead we timebox them.”*

*Flewelling, P., 2018. **The Agile Developer's Handbook**. 1st ed. Pakt Publishing.*

Our 'single question'

Can we quickly evidence
the ability for OCR
technologies to extract text
from images for the
purpose of **enhancing
metadata, particularly for
search?**



WIKIPEDIA
The Free Encyclopedia

[Main page](#)
[Contents](#)
[Featured content](#)
[Current events](#)
[Random article](#)
[Donate to Wikipedia](#)
[Wikipedia store](#)

[Interaction](#)

[Help](#)
[About Wikipedia](#)
[Community portal](#)
[Recent changes](#)
[Contact page](#)

Not logged in [Talk](#) [Contributions](#) [Create account](#) [Log in](#)

Article [Talk](#)

Read

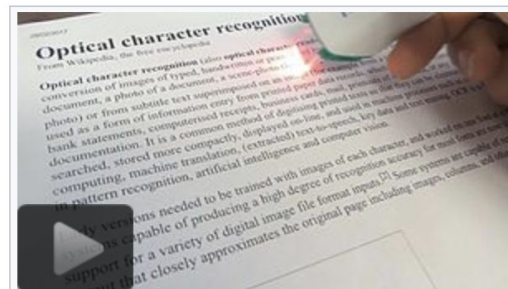
[Edit](#) [View history](#)

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]

Widely used as a form of information entry from printed paper data records – whether



Video of the process of scanning and real-time optical character recognition (OCR) with a portable scanner.

OCR - with a focus on typed text

**This was our
first exposure to
OCR
technologies
(*not our area of
expertise)**



Google Cloud

Vision API



Google Open Source Tesseract OCR

*Tesseract engine was originally developed as proprietary software by HP Labs (starting in 1985). Released as Open Source in 2005. Sponsored by Google since 2006.

Google Cloud Vision API

Good OCR, and quite a bit more

1. Good OCR
2. 'Web entities' classification
3. "Douglas Haig, 1st Earl Haig"
4. Pages with matched images

1

No. Napa 4015/07 APPENDIX 1. GENERAL HEAD QUARTER, Батаи Але я
Франски 15th September, 1918. I am informed that the officers in your
Department, and especially the officer in charge of Research in the Engineer -
in - Chief's branch, have been so good as to give great assistance to an
officer of my head - quarters recently sent Home in connection with Sound
Ranging. I would be glad if you would accept yourself, and convey to the
officers concerned, my thanks for the courtesy shown and the assistance
rendered, which has been of material advantage towards furthering military
operations. I have the honour to be, ottany 0 General,
comanding - in - Chief, British Armies in France. 1 The Postmaster - General
, General Post Office, LONDON

#D4CFC9, RGB(212, 207, 201)

Dominant colours

Aspect ratio (0.8)

Aspect ratio (1.2)

Aspect ratio (1)

Crop hints

Adult (Very Unlikely)
Spoof (Very Unlikely)
Medical (Very Unlikely)
Violence (Very Unlikely)
Racy (Very Unlikely)

Safe search

Image-4.jpg

Labels

Text (89%)

Font (70%)

Paper (67%)

Document (66%)

2

Web entities

World War I (0.58)

Letter (0.53)

Primary source (0.48)

Document (0.469)

War (0.24)

Image (0.18)

Letterhead (0.2)

BT archives (0.18)

Baseball (0.18)

Sports (0.18)

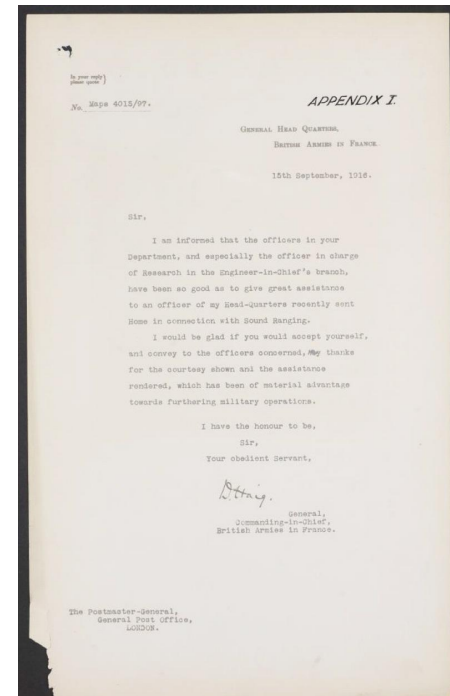
Douglas Haig, 1st
Earl Haig (0.03)

4

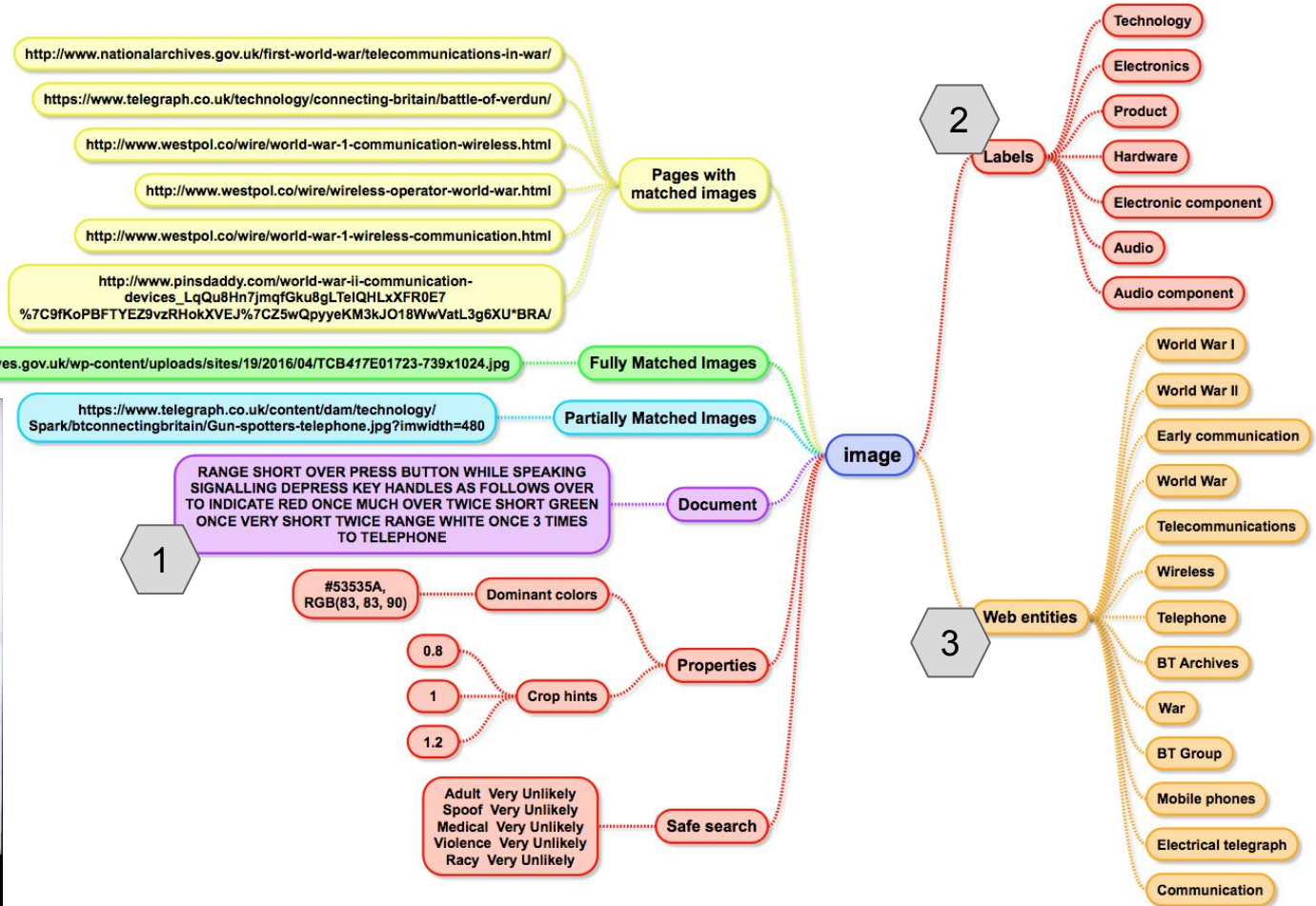
Pages with
matched images

staging.bates.org.au/...
<http://staging.bates.org.au/images/59/59-114T.jpg>

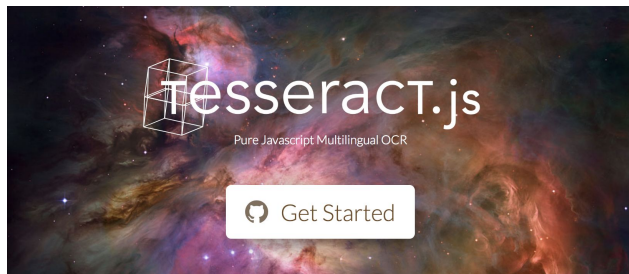
w.nationalarchives.gov.uk/...
<http://www.nationalarchives.gov.uk/wp-content/uploads/sites/19/2016/04/Image-4-653x1024.jpg>



An object



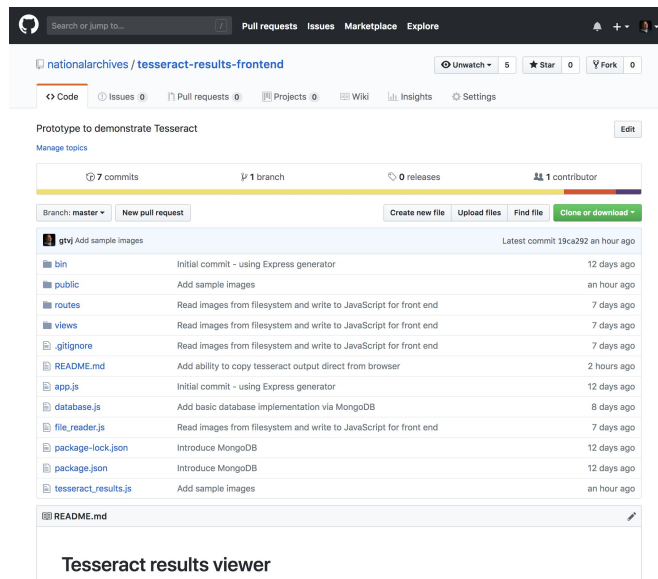
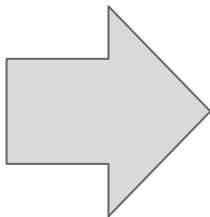
Tesseract OCR



Tesseract.js is a pure Javascript port of the popular [Tesseract OCR engine](#).

This library supports **over 60 languages**, automatic text **orientation and script detection**, a simple interface for reading paragraph, word, and character **bounding boxes**. Tesseract.js can run either in a **browser** and on a server with **NodeJS**.

Check out the [Example code and API docs on GitHub](#).



<https://github.com/nationalarchives/tesseract-results-frontend>

**Tesseract
results ranged
from very
good...**

PARLIAMENTARY
AFFAIRS

1. The Cabinet were informed of the business to be taken in the House of Commons during the following week.

It was noted that further impetus was urgently needed for the public presentation of the Government's European policies and that it was intended that material for this purpose should emerge from the speeches to be made by the Government spokesmen during the debate on European Community affairs which would be taking place on 27 November.

The Cabinet were informed that there was considerable concern, which was shared by many Government supporters, on the question of teachers' pensions, which was the subject of an Opposition Motion due to be discussed on 28 November. The Secretary of State for Education and Science was considering what steps could be taken and announced in the debate to allay this concern; and it would be helpful if in advance of the debate she could have informal discussions with Government supporters in the House of Commons.

PARLIAMENTARY 1c The Cabinet were informed of the business to be taken in the AFFAIRS Heuse of Commons during the following week. It was noted that further impetus was urgently needed for the public presentation of the Government's European policies and that it was intended that material for this purpose shoold emerge from the speeches to be made by the Government spokesmen during the debate on Eur0pean Community affairs which would be taking place on 27 November. . The Cabinet were informed that there was considerable concern, which was shared by many Government supporters, on the question of teachers' pensions, which was the subject of an OppositiOn Motion due to be discussed on 2.8 November. The Secretary of State for Education and Science was considering what steps could be taken and announced in the debate to allay this concern: and it w0uld be helpful if in advance of the debate she could have informal - discussions with Government supporters in the House of Commons.

...to good

(b) The National Insurance Fund will be in deficit next year by about £40 millions and thereafter increasingly each year. We shall need to raise the stamp as high as we can to cover these deficits, which threaten to undermine the insurance basis of the scheme. We must raise this; with all its difficulties it is the only way to meet the major problem which lies ahead. But there is a limit to the amount which we could put on the stamp without repercussions on wages.

[h] The Natisnal Insurance Fund will be in deficit next year by sheet £41] millicns and thereafter increasingly each year. We shall nee-:1 te raise the stamp as high as we can to ccver these deficits. which threaten tc undermine the insurance basis cf the scheme. We must raise this; with all. its difficulties it is the rent}: way te meet the major prehlem which lies ahead. But there is a limit tn the amcnnt which we cculd Put an the stamp without repercussions en wages

...to not so good

When in drink he becomes talkative. He has told me of journeys he made to the Continent (when connected with the 'communicating' branch) with the late Earls Curzon and Grey, and, on one occasion, mentioned one ^{1/2 937} JESSER-DAVIS, a King's messenger, with whom he apparently was frequently working.

He boasts of his friendship with Mr Harry Preston of Brighton, at whose hotel he frequently stayed with his wife. He often refers to his house in Kensington and explains his non-residence there to its being in the hands of the

'H'hen in drink he heeeneea talkative. I-le haa told me at" journeys he made to the continent [when connected with the 'cummunieating' branch} 1with the late Earle Cur-an: and Grey, and, on cue aceeeian, mentianed mefifififirilfiafiggfl :1 King's messenger, with when he apparently wee fragment-1y working. He haaete of hie frienaehin with Mr Henry" Preaten of Brighton. at. 'lhDEIB hltfil 11E frEQHEnt-ly starred with hie wife. He aft-in refer-a to hie lieuae in Kenaingttn and explains his nan-residence there te ita being in the handle at the .-

Conclusions?

- Some **very promising results in good conditions**:
 - Explore pre-optimization of images (<https://scikit-image.org>)
 - Explore using ML for confidence, training or even clean up
 - **The graph** might be interesting
 - Is there potential benefit in using OCR to expose relationships between records?
 - We have only scratched the surface.
-