

Mass digitization and OCR

Aly Conteh

Collections Digitisation & Strategy Programme
Manager
British Library

Polish Digital Libraries 2010

About the **British Library**

THE BRITISH LIBRARY

Explore the world's knowledge



We hold 14 million books, 920,000 journal and newspaper titles, 58 million patents, 3 million sound recordings, and so much more. Start exploring here.

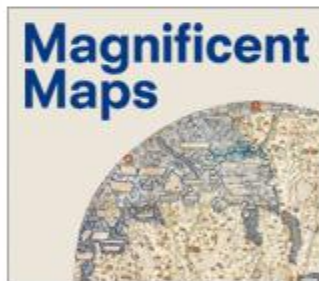
SEARCH

GO

Search tips and advanced searching

- British Library**
10,000 pages on our main website
- Online Gallery**
30,000 treasures from our collection
- Catalogue records**
14 million items in our collections
- Journal articles**
9 million articles from 20,000 journals

Quick links



▶ Exhibition closed

- ▶ Opening times, maps
- ▶ Reader Registration
- ▶ Reading Rooms
- ▶ Help for researchers
- ▶ Online catalogues
- ▶ Information in foreign languages
- ▶ For higher education
- ▶ For entrepreneurs
- ▶ For librarians
- ▶ For publishers: legal deposit etc.
- ▶ Collection Care
- ▶ Press Room
- ▶ Contact us

What's on

Site highlights

Your library

News

- 17 Sep 2010
2020 Vision launched
- 16 Sep 2010
New Learning Centre
- 15 Sep 2010
UK SoundMap
- 15 Sep 2010
Olwyn Hughes Archive
- 3 Sep 2010
Innovation Season at the British Library



▶ Business & IP Centre



▶ Online Gallery



▶ Learning



▶ Support us

British Library websites

Please choose...





150 million items





3.5 billion pages

A close-up photograph of a stack of old, worn books. The spines of the books are visible, showing signs of age and use. A blue callout box with a black border is overlaid on the books, containing the text "825 million pages" in a bold, blue, sans-serif font.

825 million pages

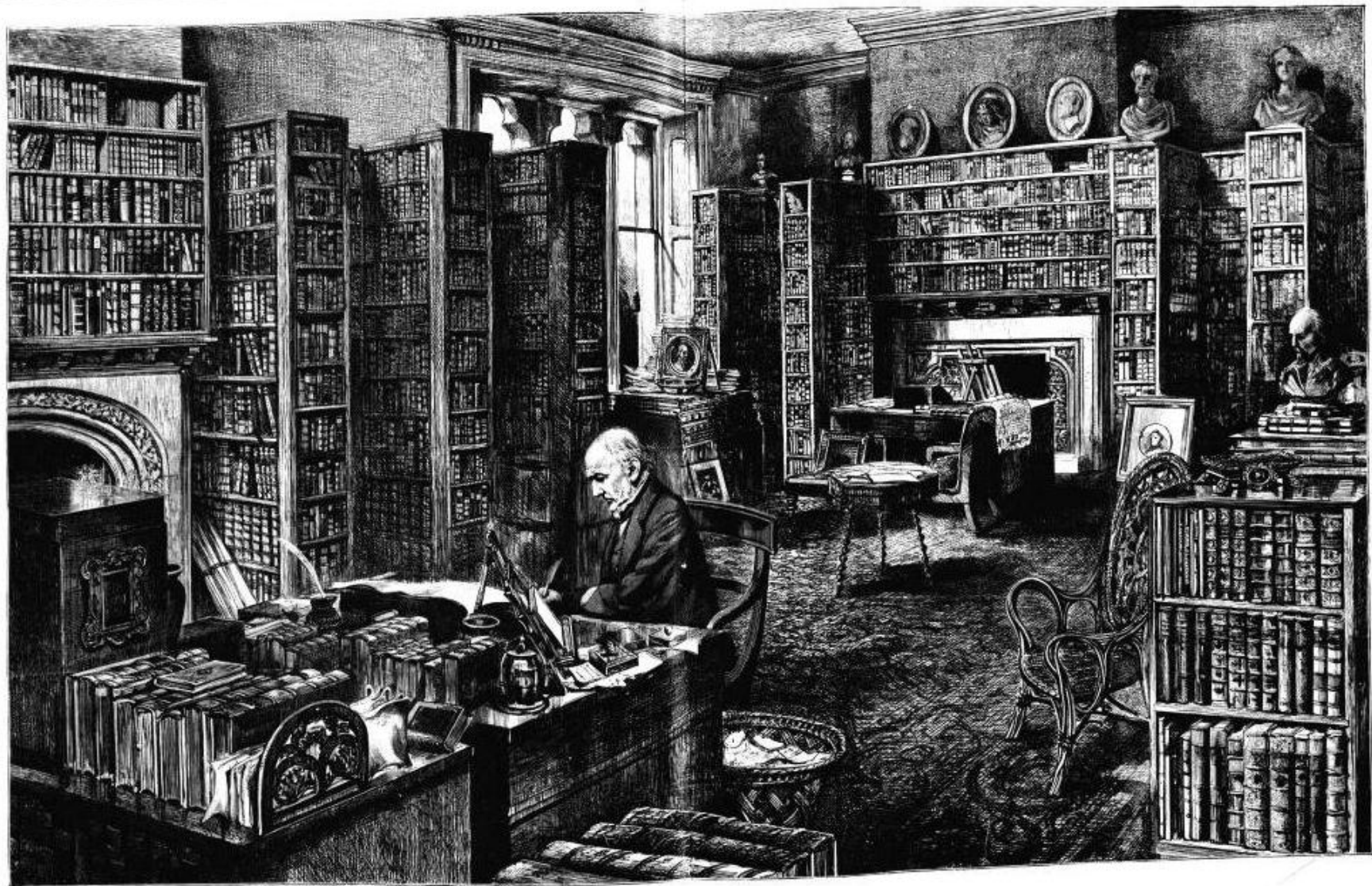
What have we **done** so far

þa gyt he him a setton segen
ðanne heah of eorðe heafod lætan holan bode
geafon on gær se g h m þa g seomor se g
munnende mod men ne cunnan se g an a
sode sele we ðanne healed under heopene
þa þam hlafte on fong

Ða wæs on byrigum beowulf se ylðingas be
leod cuning longe þrage folcum se fwe
se fæder ellor hwearf aldor of eorðe
of þam eft on poc heah healf dene heold
þendor lifde samol 7 sud weow glæde se ylð
dingas ðan fæder bearn forð se g med w
porold pocum weoroda weofa heoro gær
hyrd gær 7 halga til hyrde ic þe ðan eorð
heald se ylðingas healf se bedda þa wæs hyrd
gær e hie sped gær þa weofa weofa weofa
ham he gær magas seorne hyrdon oððe
se g se g se g se g se g se g se g se g se g
on mod bearn þe healf weofa healf weofa

Beowulf

1994



THE TEMPLE OF PEACE: MR. GLADSTONE IN HIS STUDY AT HAWARDEN
FROM A PHOTOGRAPH BY MACGARDY AND GOGAN, WREXHAM

1	2	3	4
1			

4 million pages pre-1900 Newspapers

ARCHIVAL SOUND RECORDINGS ACCENTS AND DIALECTS



8,000 hours of Sound Recordings

What **challenges** do we face with our mass digitisation programme

...the amount of **data** we
generate

...25 million pages of 19th
century books

...old standard of uncompressed
TIFF files

...creates over **2PB** worth of data

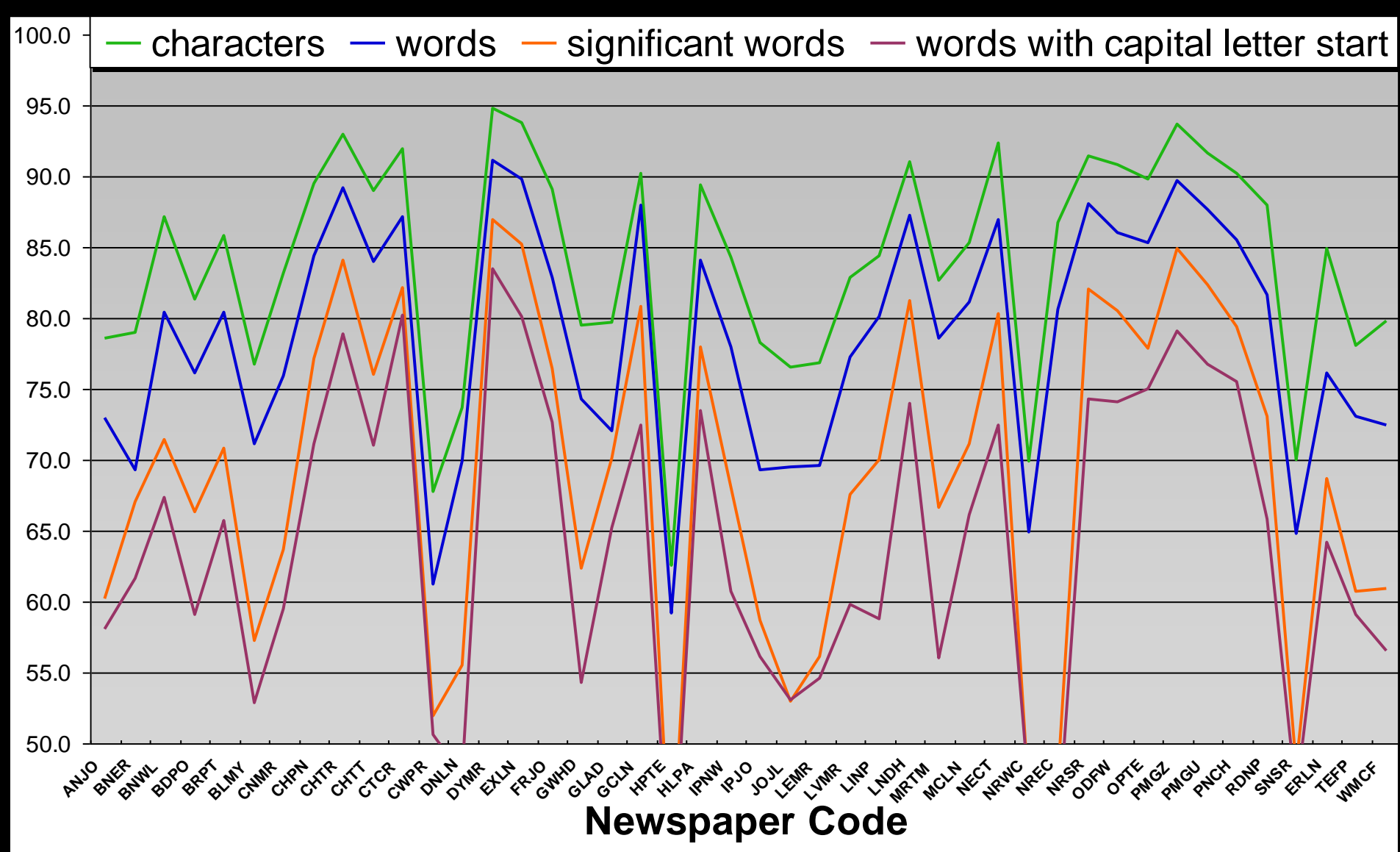
...the solution for us was to use
the **JPEG2000** format instead

...using **lossy!** compression,
under 40TB

...is **visually** lossless
compression the answer?



...Need for better **text** extraction

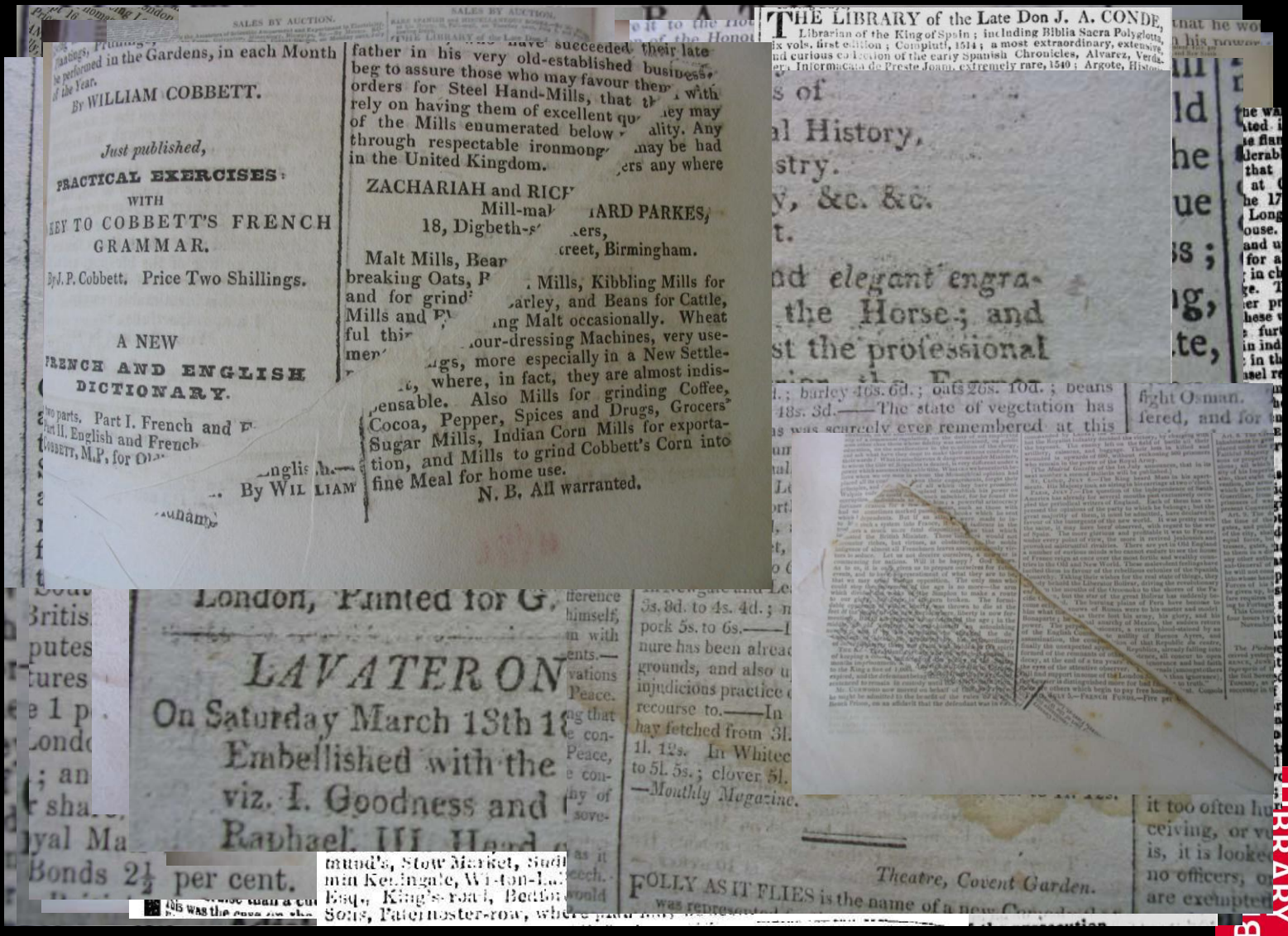


Measuring **OCR** accuracy

**...four issues that affect
performance of OCR?**

PHYSICAL CHARACTERISTICS OF SOURCE MATERIAL

- Bleed through
- Stains
- Tight binding
- Holes/tears
- Printer errors
- Stamps
- Creases
- Dirt
- Animals
- Inconsistent inking
- Repairs
- Paper quality
- Lamination



SINGULAR AND SERIOUS ACCIDENT.—On Wednesday noon Mr. Charles Wyber, of the Borough-road, went into the Fleet Prison to visit a friend, and joined a party in the Coffee-room, who entered into the foolish amusement of tossing up a penny-piece to the top of the room, and catching the same in the mouth upon its descent to the floor. Mr. Wyber was considered a perfect adept at this game; but the penny-piece at last found its way into the throat, where it remained for upwards of half an hour. A Surgeon tried to force it upwards, but being unable to do so, he contrived to move it downwards into the stomach. Mr. Wyber was comparatively much relieved by the penny-piece being removed out of the throat, and was enabled, in the evening, to be carried to his house in a hackney-coach.

A SINGULAR AND SERIOUS ACCIDENT.- -On the 17th of the month of
noon Mr. Charles Wyber, of the Borough-road, near
Fleet Prison to visit a friend, and joined a party in a
room, who entered into the foolish and unskillful game of
penny-piece to the top of the room, and catching the
the moth upon its descent to the floor. Mr. Wyber was con-
sidered a perfect adept at this game but time flew
last found its way into the throat, where it was retained
wards of half an hour. A Surgeon tried to force
but being unable to do so, he contrived to push it
into the stomach. Mr. Wyber was comparatively
relieved by the penny-piece being removed out of the
was enabled, in the evening, to be carried to
hackney-coach.

The numbers

There are **152** words

61 words are incorrectly identified

Giving us 60% **word** accuracy

But all **words** are not equal

SINGULAR AND SERIOUS ACCIDENT.—On Wednesday noon Mr. Charles Wyber, of the Borough-road, went into the Fleet Prison to visit a friend, and joined a party in the Coffee-room, who entered into the foolish amusement of tossing up a penny-piece to the top of the room, and catching the same in the mouth upon its descent to the floor. Mr. Wyber was considered a perfect adept at this game; but the penny-piece at last found its way into the throat, where it remained for upwards of half an hour. A Surgeon tried to force it upwards, but being unable to do so, he contrived to move it downwards into the stomach. Mr. Wyber was comparatively much relieved by the penny-piece being removed out of the throat, and was enabled, in the evening, to be carried to his house in a hackney-coach.

They had the internet in 1816 !

to the discussion on
new loan from the
an arrangement on
vernment.
of LAUDERDALE
0,000, or 2,000,000
ount at present.
MILITARY.
that the subject of
had given notice
consequence of the
appeared from what
the law had been
all the efforts of the
power above the
pon themselves the
ation with the civil
at there was an im-
y-authority and the
statement rendered
aid of the Secretary
w of the land. It
count, who united
ad of the Police of
tue of which office
country, should,
without the means
channel, it was more dangerous than at any other period to re-
sume cash-payments. The experience of every war proved the
embarrassment consequent on a return from war to peace, and
the experience of the most successful wars proved this in the
greatest degree. There were persons now living, who recol-
lected that on the peace of 1763 the greatest apprehensions were
entertained, whether the finances of the country would be equal
to the interest of the National Debt. When any one considered
the duration of the late war, the extensive establishments which
had been kept up, the character of a war, which was directly
aimed at the disinterese of the country, they would not be sur-
prised that the change should be accompanied by considerable
embarrassment: But looking to the evils as merely of a tem-
porary nature, their Lordships were bound to proceed with
great caution on a subject so immediately connected with them
as that now under their consideration. Considerable difference
of opinion existed, whether the measure of 1797 had on the
whole been productive of good or bad effects;—on that subject
he would not now give any opinion;—but there was one
disadvantage which he had always considered as result-
ing from the measure of 1797, namely, the difficulty of
again returning to the old system. However, he had not
the least difficulty in saying, that the Bank ought to return
to cash payments as soon as possible. He took the 5th of
July, 1818, in preference to a shorter period, because he was
convinced that it would be more advantageous to the credit and
interest of the country, and the interest of the Bank, to consi-
der at once what would be the first period when they could re-
sume with safety their payments, than to enact a measure for a
never
person
puffe
if Pa
in fav
Lo
retur
was r
Thou
to re
time
and t
woul
woul
occa
Th
hims
the l
arity
How
Ti
affid
of Ju
M
cult
red t
Lo
Bill.
Ti

and DVD in 1803!

A Coroner's Inquest was held on Wednesday on the body of Elizabeth Colebird, late wife of George Colebird, whose death was occasioned by an affray at her house in St. Giles's. It appeared that a man named Burk, with several other persons, rushed into the room in quest of her husband, with whom they had previously quarrelled, and that she being far advanced in pregnancy, was prematurely delivered of a child which was born alive, but, together with the mother, died soon after. There was no evidence to prove that the blows she had received were the cause of her death, and two surgeons, who had examined the body, gave it as their opinion that she died by fright. The jury, after an investigation of three hours, returned a verdict—*Died by the Visitation of God.*

MEETING OF CREDITORS AT GUILDHALL.

The Morning Chronicle (London, England), Friday, June 10, 1803; Issue 10625

...Summary of issues

Geometric **distortions** lead to text being missed or incorrectly identified

Quality of **source** material has a notable impact on accuracy levels

The need to focus on **significant** words

False positives can be introduced by using modern lexicons

And there are others.....

How are we **addressing** the OCR
issues?

IMPACT

Innovating **OCR** software and language technology

Sharing **expertise** and building capacity across Europe

Ensuring that **tools** and services will be sustained after the end of the project

Improving Access to Text

IMPACT



IMPACT is supported by the European Community under the FP7 ICT Work Programme. The project is coordinated by the National Library of the Netherlands.

26 partners: Libraries, Research Institutes, Industry Partners



Facts and figures

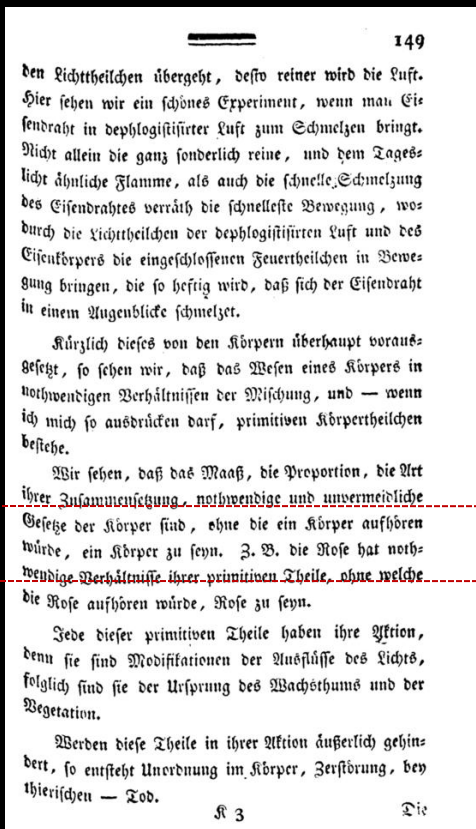
- Project supported by the European Community under the FP7 ICT Work Programme.
- coordinated by the National Library of the Netherlands (KB)
- EU funding: € 11 500 000
- Start date: 1 January 2008
- Duration: 48 months
- From 2012: sustainable Centre of Competence
- Web site: www.impact-project.eu

...what will **IMPACT deliver**

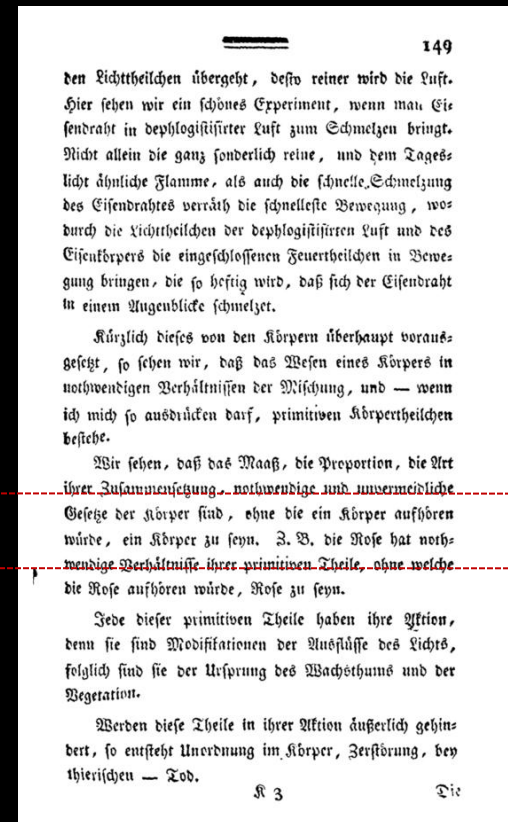
...advancement in the **state of the art** in image enhancement, segmentation and binarisation

Straightening the lines

Original image



Dewarping v.1.1



...improvements with **OCR**
products

...development of **language** tools

Building Named Entities resources (Source Material 34,500 pages)

Dictionary of National Biography

Alumni Oxonienses: the members of the University of Oxford 1500-1714

Alumni Oxonienses: the members of the University of Oxford, 1715-1886

Wilson's Mercantile Directory of Great Britain and Ireland

Cassell's Gazetteer of Great Britain and Ireland

...tools to **build** capability

...and a research **dataset**

Thank You



www.bl.uk

aly.conteh@bl.uk

twitter: @aconteh



<http://www.impact-project.eu/>

Twitter: @impactocr

<http://impactocr.wordpress.com/>