



DIGITIZATION AND PRESERVATION AS AN EDITOR

Marc Cormier
Director, Humanities Publishing
Gale
Farmington Hills, MI

A test case for preservation and digitization: The Daily Mail Historical Archive

Sourcing archival material

Image Capture

OCR

Metadata

XML

QA



The source of
Gale's Daily Mail
Historical Archive

LIBRARY
HSTLRB



Which Edition?

- There is no single edition!
- Content is usually similar, but “local” advertising, story selection will differentiate each edition.
- Digitizing all editions, while valuable is expensive
- Bringing all editions to a research environment can be messy
- Alternate editions can often be used to fill gaps in coverage of your chosen edition



Balancing priorities of preservation and access

- Unique version of the newspaper published on board the Atlantic liners that carried passengers between the UK and New York.
- The only known set of these issues is held in the Daily Mail offices, and it had long been forgotten and subject to water and rodent damage
- We decided to include these in the project for preservation purposes – otherwise they would have been genuinely lost to history.

Working with the source material: challenges and opportunities

SCANNING FROM ORIGINAL



SCANNING FROM MICROFILM



Working with the source material: challenges and opportunities

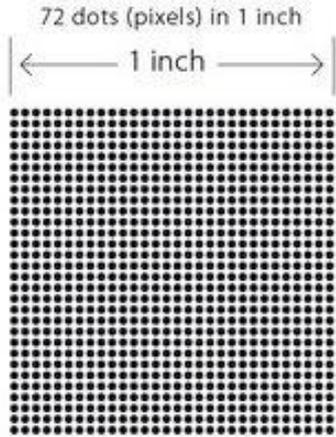
SCANNING FROM MICROFILM

- Many newspapers have been microfilmed, as part of preservation efforts from the 1960s onwards.
- Digitising from microfilm is cheaper and logistically easier than scanning physical copies
- Although it does mean that you will only get black & white images.

SCANNING FROM ORIGINAL

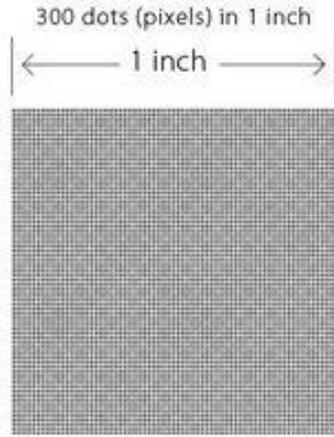
- Allow you to create high-quality colour images using the latest scanning technology
- it requires a VERY understanding source library to allow you to handle the bound volumes.
- For optimal results, you need to be able to dis-bind the volumes to avoid getting curvature on the image.

Image Capture



72 dpi

72 dots per-inch



300 dpi

300 dots per-inch

Image courtesy of www.zenfolio.com

- Higher DPI = higher clarity, but it also means large file sizes and slow load times.
- For Daily Mail, we used the LOC guidelines of 400dpi, a balance between readability and file size.
- For more intricate documents like mediaeval or scientific manuscripts, we may scan as high as 1,200dpi in order to permit close viewing of marginalia, ink and script details, illustrations, etc.

Image options: focusing on the reading experience

BITONAL



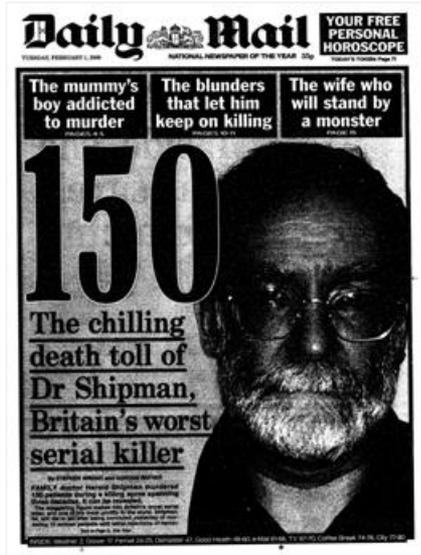
- Bitonal: B/W scan provides high contrast but can hinder images and illustrations
- Grayscale: permits subtlety of tone and shades but can pick up grain from paper
- Bitonal is regularly used for historic newspapers but the Daily Mail project forced us to reconsider, based on when photographs first appear.

GREYSCALE



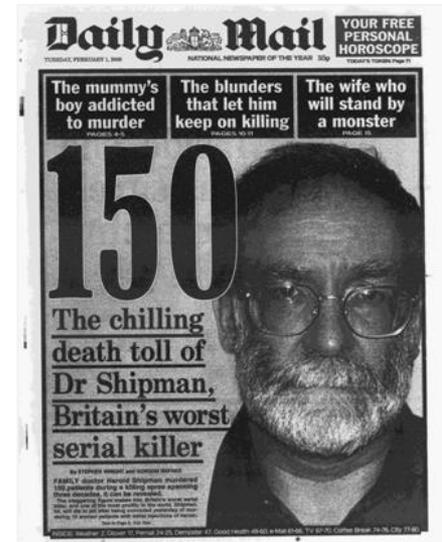
Image options: focusing on the reading experience

BITONAL



- The Daily Mail becomes increasingly more illustrated throughout the 20th C.
- Bitonal is used up to 1962, when color and true halftone photographs began appearing more regularly.
- Users will notice a contrast in the way the images appear either side of this date.

GREYSCALE





Optical Character Recognition

- The quality of the OCR text usually says more about the condition of the original materials than it does about the performance of the OCR software.
- Typically, older newspapers produce much less satisfactory results than modern ones.
- Wartime newspapers noticeably produce poorer results than adjacent periods of history; especially 'bleed-through' of text from the other side of the page.

How accurate is *your* OCR?

WHAT DOES “ACCURACY” REALLY MEAN?

- The word ‘accuracy’ is misleading; OCR software works from a *confidence* rating, and NOT true accuracy
- The software calculates a confidence level from 0-9 for each character it detects, but does not know whether a character has been converted correctly or not.
- The *Daily Mail Historical Archive* had a team of over 400 people creating and reviewing the data for the archive, but with over 1.2 million pages to digitise and convert, it was not physically possible to clean up the OCR for every article.
- Typically, only small-scale or crowd-sourced digitisation projects have a realistic opportunity of producing 100% perfect OCR.

METADATA

High quality metadata makes it much easier to find the specific types of information a researcher is looking for, and to place effective parameters and filters around a search query, such as date ranges limited to specific article types.

Daily Mail THURSDAY, FEBRUARY 24, 1991 28p **FEMAIL Magazine** SEE PAGES 23-24

5AM: BUSH STOPS WAR

Iraqis 'accept all the UN terms for peace'

From **GEORGE GORDON** in Washington
PRESIDENT Bush declared an end to fighting in the Gulf War early today. He announced a ceasefire to begin at midnight (GMT British time). And within hours a number of reports from the United Nations said that Iraq's foreign minister **TARIQ ALI** had said he agreed to the UN terms — which stated: "The government of Iraq agrees to fully comply with the 12 UN resolutions."

It is already reported, on page 1, that the UN Security Council has decided to allow the UN to inspect Iraq's oil fields for oil. The suspension of all offensive military operations was announced as Iraq's compliance with the terms of the ceasefire.

A cease-fire was announced at midnight today and the UN Security Council has decided to allow the UN to inspect Iraq's oil fields for oil. The suspension of all offensive military operations was announced as Iraq's compliance with the terms of the ceasefire.

Violated

The ceasefire was made to start immediately with Iraq. The UN Security Council has decided to allow the UN to inspect Iraq's oil fields for oil. The suspension of all offensive military operations was announced as Iraq's compliance with the terms of the ceasefire.

President officials said last night the fighting was about to end and Iraq's foreign minister, **TARIQ ALI**, said he had agreed to the UN terms. The UN Security Council has decided to allow the UN to inspect Iraq's oil fields for oil. The suspension of all offensive military operations was announced as Iraq's compliance with the terms of the ceasefire.

Turn to Page 2, Col. 1

VICTIMS OF A TERRIBLE BATTLE ERROR

Nine Desert Rats killed by American 'friendly fire'

NINE young Desert Rats have been killed by American 'friendly fire', it was learnt yesterday.

An A-1H aircraft's mission and shells hit the targets in their own Desert fighting cars. And nine of the tank-buster teams, which is used to give the point guard resistance.

One of those who died was **Facilier Conrad Cole**, aged 17.

Strafed

The A-1H aircraft's mission and shells hit the targets in their own Desert fighting cars. And nine of the tank-buster teams, which is used to give the point guard resistance.

Turn to Page 4, Col. 2

INSIDE: Weather 7, Night Bomber 21, TV and Radio 25-28, Books 60-61, Entertainment 62, Letters 84, City 86-87, Coffee Break 94, Sport 93-96

Who?
(George Gordon)

What?
(Article Title: 5AM: Bush Stops War)
(Article Type: News)
(679 words)

When?
(28 February 1991)

Where?
(Daily Mail, London, England)
(Page 1)
(Issue 29, 453)

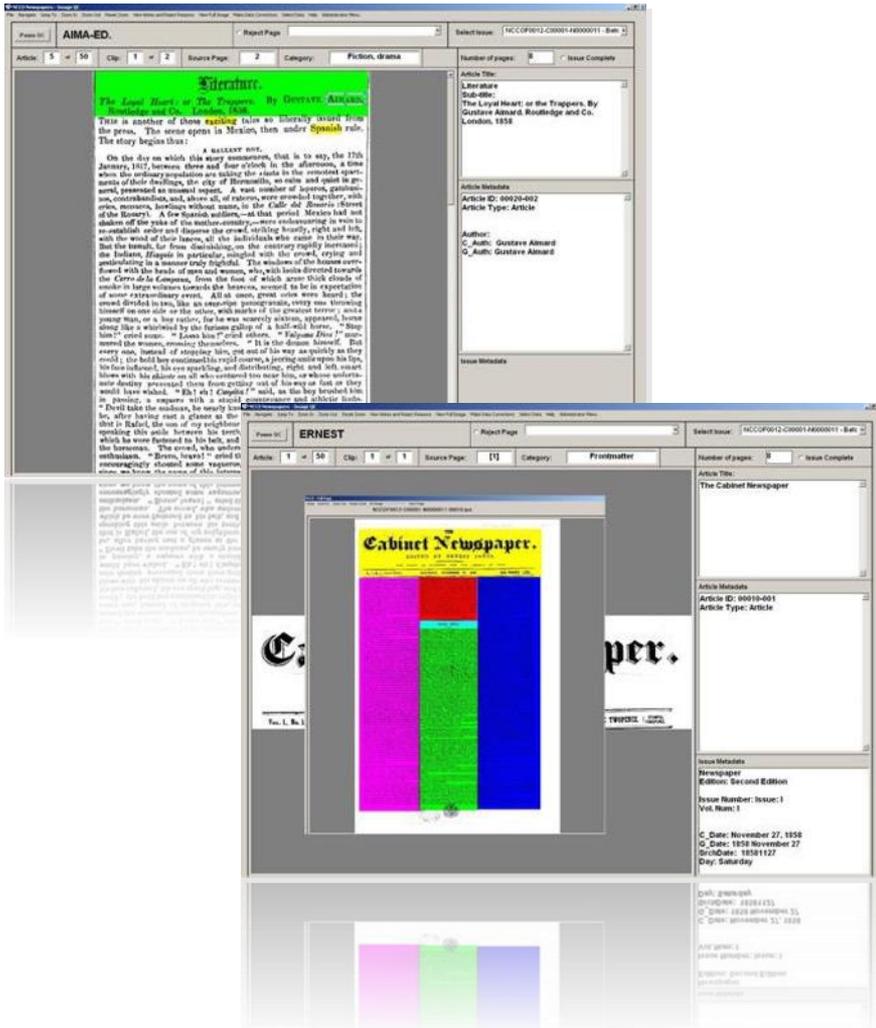
XML

- Creating XML is by far the costliest and most labour-intensive part of any digitisation project.
- Word coordinates are essential to hit-term highlighting and zoning of article copy.
- DTD defines the data structure for the archive, outlining a list of permissible legal elements and attributes.
- All data that is captured for a project must fit the DTD, or it does not pass verification – there are no exceptions.

```
<wd pos="547,1273,622,1298">From</wd>
<wd pos="638,1271,784,1297">WALTER</wd>
<wd pos="799,1272,890,1299">FARR</wd>
<wd pos="906,1273,989,1303">(Daily</wd>
<wd pos="1004,1273,1065,1299">Mail</wd>
<wd pos="1080,1272,1182,1303">Special</wd>
<wd pos="1199,1273,1412,1304">Correspondent)</wd>
<wd pos="1446,1286,1649,1318">Washington,</wd>
<wd pos="1668,1288,1766,1319">Friday.</wd>
<wd pos="574,1325,802,1367">BRITAIN,</wd>
<wd pos="831,1328,894,1362">the</wd>
<wd pos="925,1326,1058,1362">United</wd>
<wd pos="1087,1328,1214,1369">States,</wd>
<wd pos="1239,1328,1366,1372">China,</wd>
<wd pos="1394,1331,1529,1373">Russia,</wd>
<wd pos="1553,1333,1624,1368">and</wd>
```

Quality Assurance

- Article zoning and segmentation defines the borders of each article, page property, section property, section property, etc.
- “Clip” view allows users to see article segmentation through highlighting
- Manual QA is used to ensure segmentation is accurate and that article split points are correct.



TEXT DATA MINING – driving new insights from traditional content

Articles

Public Health, the Journal of t...
 Subscribers' Copies of Public...
 Notices (p.1)

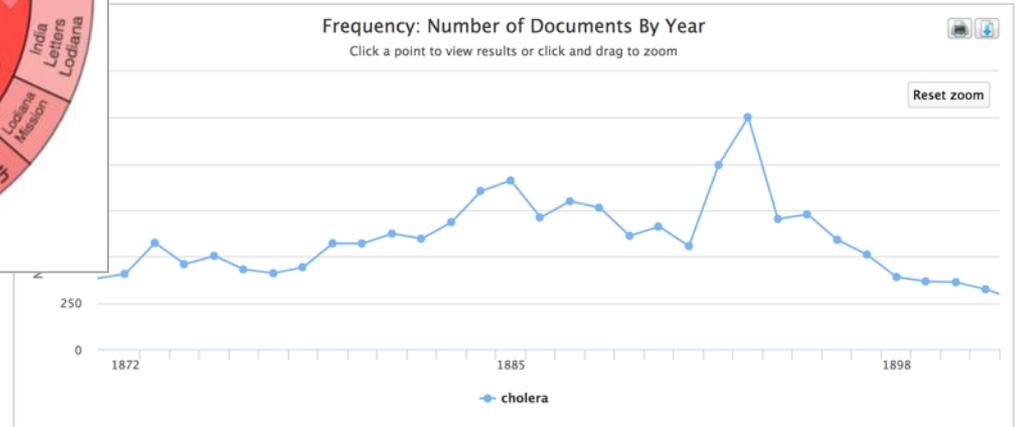
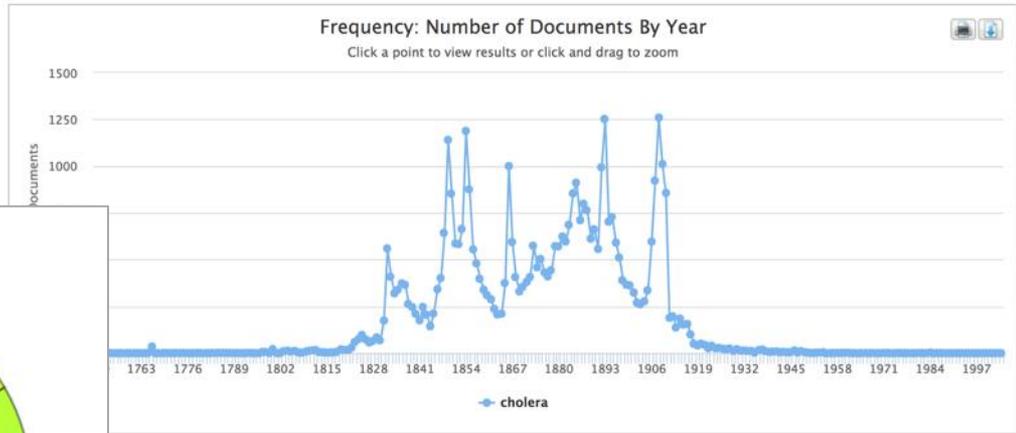
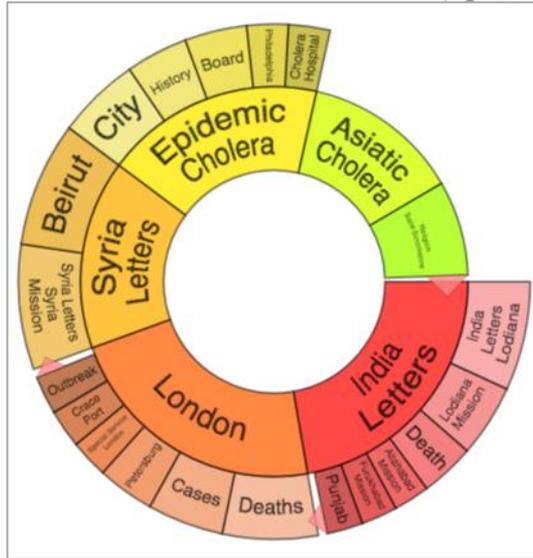
The Incorporated Society of ...
 Annual Dinner (p.1)

The Organisation of Meat Ins...
 The Incorporated Society of ...
 International Congress of Hy...
 Bacteriological Investigations...

The Statistics of Diphtheria in...
 On the Present State of Kno...
 On the Serum Therapeutics o...
 The Bearing of Experimental ...
 Purification by Fire (p.20)

Report on the Epidemic of En...
 Preservation of Fresh Milk (p...
 The Corrected Death Rates o...
 Diagnosis of Influenza (p.32)
 Public Health (p.32)

The Budapest Congress (p...
 The Corrected Death Rate of ...
 Epidemic of Enteric Fever, Tr...
 On an Outbreak of Anomalou...
 A Case of Hyperpyrexia (p.38)



```

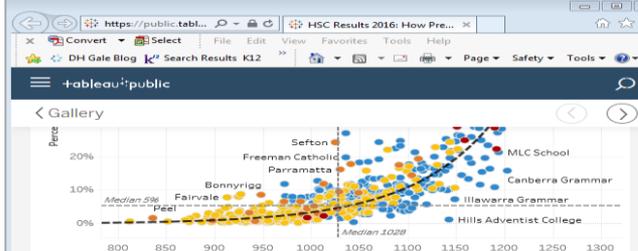
Administration: C:\windows\System32\WindowsPowerShell\v1.0\Powershell.exe
line 13 char 5
1 <<<
2 <<<
3 <<<
4 <<<
5 <<<
6 <<<
7 <<<
8 <<<
9 <<<
10 <<<
11 <<<
12 <<<
13 <<<
14 <<<
15 <<<
16 <<<
17 <<<
18 <<<
19 <<<
20 <<<
21 <<<
22 <<<
23 <<<
24 <<<
25 <<<
26 <<<
27 <<<
28 <<<
29 <<<
30 <<<
31 <<<
32 <<<
33 <<<
34 <<<
35 <<<
36 <<<
37 <<<
38 <<<
39 <<<
40 <<<
41 <<<
42 <<<
43 <<<
44 <<<
45 <<<
46 <<<
47 <<<
48 <<<
49 <<<
50 <<<
51 <<<
52 <<<
53 <<<
54 <<<
55 <<<
56 <<<
57 <<<
58 <<<
59 <<<
60 <<<
61 <<<
62 <<<
63 <<<
64 <<<
65 <<<
66 <<<
67 <<<
68 <<<
69 <<<
70 <<<
71 <<<
72 <<<
73 <<<
74 <<<
75 <<<
76 <<<
77 <<<
78 <<<
79 <<<
80 <<<
81 <<<
82 <<<
83 <<<
84 <<<
85 <<<
86 <<<
87 <<<
88 <<<
89 <<<
90 <<<
91 <<<
92 <<<
93 <<<
94 <<<
95 <<<
96 <<<
97 <<<
98 <<<
99 <<<
100 <<<
101 <<<
102 <<<
103 <<<
104 <<<
105 <<<
106 <<<
107 <<<
108 <<<
109 <<<
110 <<<
111 <<<
112 <<<
113 <<<
114 <<<
115 <<<
116 <<<
117 <<<
118 <<<
119 <<<
120 <<<
121 <<<
122 <<<
123 <<<
124 <<<
125 <<<
126 <<<
127 <<<
128 <<<
129 <<<
130 <<<
131 <<<
132 <<<
133 <<<
134 <<<
135 <<<
136 <<<
137 <<<
138 <<<
139 <<<
140 <<<
141 <<<
142 <<<
143 <<<
144 <<<
145 <<<
146 <<<
147 <<<
148 <<<
149 <<<
150 <<<
151 <<<
152 <<<
153 <<<
154 <<<
155 <<<
156 <<<
157 <<<
158 <<<
159 <<<
160 <<<
161 <<<
162 <<<
163 <<<
164 <<<
165 <<<
166 <<<
167 <<<
168 <<<
169 <<<
170 <<<
171 <<<
172 <<<
173 <<<
174 <<<
175 <<<
176 <<<
177 <<<
178 <<<
179 <<<
180 <<<
181 <<<
182 <<<
183 <<<
184 <<<
185 <<<
186 <<<
187 <<<
188 <<<
189 <<<
190 <<<
191 <<<
192 <<<
193 <<<
194 <<<
195 <<<
196 <<<
197 <<<
198 <<<
199 <<<
200 <<<
201 <<<
202 <<<
203 <<<
204 <<<
205 <<<
206 <<<
207 <<<
208 <<<
209 <<<
210 <<<
211 <<<
212 <<<
213 <<<
214 <<<
215 <<<
216 <<<
217 <<<
218 <<<
219 <<<
220 <<<
221 <<<
222 <<<
223 <<<
224 <<<
225 <<<
226 <<<
227 <<<
228 <<<
229 <<<
230 <<<
231 <<<
232 <<<
233 <<<
234 <<<
235 <<<
236 <<<
237 <<<
238 <<<
239 <<<
240 <<<
241 <<<
242 <<<
243 <<<
244 <<<
245 <<<
246 <<<
247 <<<
248 <<<
249 <<<
250 <<<
251 <<<
252 <<<
253 <<<
254 <<<
255 <<<
256 <<<
257 <<<
258 <<<
259 <<<
260 <<<
261 <<<
262 <<<
263 <<<
264 <<<
265 <<<
266 <<<
267 <<<
268 <<<
269 <<<
270 <<<
271 <<<
272 <<<
273 <<<
274 <<<
275 <<<
276 <<<
277 <<<
278 <<<
279 <<<
280 <<<
281 <<<
282 <<<
283 <<<
284 <<<
285 <<<
286 <<<
287 <<<
288 <<<
289 <<<
290 <<<
291 <<<
292 <<<
293 <<<
294 <<<
295 <<<
296 <<<
297 <<<
298 <<<
299 <<<
300 <<<
301 <<<
302 <<<
303 <<<
304 <<<
305 <<<
306 <<<
307 <<<
308 <<<
309 <<<
310 <<<
311 <<<
312 <<<
313 <<<
314 <<<
315 <<<
316 <<<
317 <<<
318 <<<
319 <<<
320 <<<
321 <<<
322 <<<
323 <<<
324 <<<
325 <<<
326 <<<
327 <<<
328 <<<
329 <<<
330 <<<
331 <<<
332 <<<
333 <<<
334 <<<
335 <<<
336 <<<
337 <<<
338 <<<
339 <<<
340 <<<
341 <<<
342 <<<
343 <<<
344 <<<
345 <<<
346 <<<
347 <<<
348 <<<
349 <<<
350 <<<
351 <<<
352 <<<
353 <<<
354 <<<
355 <<<
356 <<<
357 <<<
358 <<<
359 <<<
360 <<<
361 <<<
362 <<<
363 <<<
364 <<<
365 <<<
366 <<<
367 <<<
368 <<<
369 <<<
370 <<<
371 <<<
372 <<<
373 <<<
374 <<<
375 <<<
376 <<<
377 <<<
378 <<<
379 <<<
380 <<<
381 <<<
382 <<<
383 <<<
384 <<<
385 <<<
386 <<<
387 <<<
388 <<<
389 <<<
390 <<<
391 <<<
392 <<<
393 <<<
394 <<<
395 <<<
396 <<<
397 <<<
398 <<<
399 <<<
400 <<<
401 <<<
402 <<<
403 <<<
404 <<<
405 <<<
406 <<<
407 <<<
408 <<<
409 <<<
410 <<<
411 <<<
412 <<<
413 <<<
414 <<<
415 <<<
416 <<<
417 <<<
418 <<<
419 <<<
420 <<<
421 <<<
422 <<<
423 <<<
424 <<<
425 <<<
426 <<<
427 <<<
428 <<<
429 <<<
430 <<<
431 <<<
432 <<<
433 <<<
434 <<<
435 <<<
436 <<<
437 <<<
438 <<<
439 <<<
440 <<<
441 <<<
442 <<<
443 <<<
444 <<<
445 <<<
446 <<<
447 <<<
448 <<<
449 <<<
450 <<<
451 <<<
452 <<<
453 <<<
454 <<<
455 <<<
456 <<<
457 <<<
458 <<<
459 <<<
460 <<<
461 <<<
462 <<<
463 <<<
464 <<<
465 <<<
466 <<<
467 <<<
468 <<<
469 <<<
470 <<<
471 <<<
472 <<<
473 <<<
474 <<<
475 <<<
476 <<<
477 <<<
478 <<<
479 <<<
480 <<<
481 <<<
482 <<<
483 <<<
484 <<<
485 <<<
486 <<<
487 <<<
488 <<<
489 <<<
490 <<<
491 <<<
492 <<<
493 <<<
494 <<<
495 <<<
496 <<<
497 <<<
498 <<<
499 <<<
500 <<<
501 <<<
502 <<<
503 <<<
504 <<<
505 <<<
506 <<<
507 <<<
508 <<<
509 <<<
510 <<<
511 <<<
512 <<<
513 <<<
514 <<<
515 <<<
516 <<<
517 <<<
518 <<<
519 <<<
520 <<<
521 <<<
522 <<<
523 <<<
524 <<<
525 <<<
526 <<<
527 <<<
528 <<<
529 <<<
530 <<<
531 <<<
532 <<<
533 <<<
534 <<<
535 <<<
536 <<<
537 <<<
538 <<<
539 <<<
540 <<<
541 <<<
542 <<<
543 <<<
544 <<<
545 <<<
546 <<<
547 <<<
548 <<<
549 <<<
550 <<<
551 <<<
552 <<<
553 <<<
554 <<<
555 <<<
556 <<<
557 <<<
558 <<<
559 <<<
560 <<<
561 <<<
562 <<<
563 <<<
564 <<<
565 <<<
566 <<<
567 <<<
568 <<<
569 <<<
570 <<<
571 <<<
572 <<<
573 <<<
574 <<<
575 <<<
576 <<<
577 <<<
578 <<<
579 <<<
580 <<<
581 <<<
582 <<<
583 <<<
584 <<<
585 <<<
586 <<<
587 <<<
588 <<<
589 <<<
590 <<<
591 <<<
592 <<<
593 <<<
594 <<<
595 <<<
596 <<<
597 <<<
598 <<<
599 <<<
600 <<<
601 <<<
602 <<<
603 <<<
604 <<<
605 <<<
606 <<<
607 <<<
608 <<<
609 <<<
610 <<<
611 <<<
612 <<<
613 <<<
614 <<<
615 <<<
616 <<<
617 <<<
618 <<<
619 <<<
620 <<<
621 <<<
622 <<<
623 <<<
624 <<<
625 <<<
626 <<<
627 <<<
628 <<<
629 <<<
630 <<<
631 <<<
632 <<<
633 <<<
634 <<<
635 <<<
636 <<<
637 <<<
638 <<<
639 <<<
640 <<<
641 <<<
642 <<<
643 <<<
644 <<<
645 <<<
646 <<<
647 <<<
648 <<<
649 <<<
650 <<<
651 <<<
652 <<<
653 <<<
654 <<<
655 <<<
656 <<<
657 <<<
658 <<<
659 <<<
660 <<<
661 <<<
662 <<<
663 <<<
664 <<<
665 <<<
666 <<<
667 <<<
668 <<<
669 <<<
670 <<<
671 <<<
672 <<<
673 <<<
674 <<<
675 <<<
676 <<<
677 <<<
678 <<<
679 <<<
680 <<<
681 <<<
682 <<<
683 <<<
684 <<<
685 <<<
686 <<<
687 <<<
688 <<<
689 <<<
690 <<<
691 <<<
692 <<<
693 <<<
694 <<<
695 <<<
696 <<<
697 <<<
698 <<<
699 <<<
700 <<<
701 <<<
702 <<<
703 <<<
704 <<<
705 <<<
706 <<<
707 <<<
708 <<<
709 <<<
710 <<<
711 <<<
712 <<<
713 <<<
714 <<<
715 <<<
716 <<<
717 <<<
718 <<<
719 <<<
720 <<<
721 <<<
722 <<<
723 <<<
724 <<<
725 <<<
726 <<<
727 <<<
728 <<<
729 <<<
730 <<<
731 <<<
732 <<<
733 <<<
734 <<<
735 <<<
736 <<<
737 <<<
738 <<<
739 <<<
740 <<<
741 <<<
742 <<<
743 <<<
744 <<<
745 <<<
746 <<<
747 <<<
748 <<<
749 <<<
750 <<<
751 <<<
752 <<<
753 <<<
754 <<<
755 <<<
756 <<<
757 <<<
758 <<<
759 <<<
760 <<<
761 <<<
762 <<<
763 <<<
764 <<<
765 <<<
766 <<<
767 <<<
768 <<<
769 <<<
770 <<<
771 <<<
772 <<<
773 <<<
774 <<<
775 <<<
776 <<<
777 <<<
778 <<<
779 <<<
780 <<<
781 <<<
782 <<<
783 <<<
784 <<<
785 <<<
786 <<<
787 <<<
788 <<<
789 <<<
790 <<<
791 <<<
792 <<<
793 <<<
794 <<<
795 <<<
796 <<<
797 <<<
798 <<<
799 <<<
800 <<<
801 <<<
802 <<<
803 <<<
804 <<<
805 <<<
806 <<<
807 <<<
808 <<<
809 <<<
810 <<<
811 <<<
812 <<<
813 <<<
814 <<<
815 <<<
816 <<<
817 <<<
818 <<<
819 <<<
820 <<<
821 <<<
822 <<<
823 <<<
824 <<<
825 <<<
826 <<<
827 <<<
828 <<<
829 <<<
830 <<<
831 <<<
832 <<<
833 <<<
834 <<<
835 <<<
836 <<<
837 <<<
838 <<<
839 <<<
840 <<<
841 <<<
842 <<<
843 <<<
844 <<<
845 <<<
846 <<<
847 <<<
848 <<<
849 <<<
850 <<<
851 <<<
852 <<<
853 <<<
854 <<<
855 <<<
856 <<<
857 <<<
858 <<<
859 <<<
860 <<<
861 <<<
862 <<<
863 <<<
864 <<<
865 <<<
866 <<<
867 <<<
868 <<<
869 <<<
870 <<<
871 <<<
872 <<<
873 <<<
874 <<<
875 <<<
876 <<<
877 <<<
878 <<<
879 <<<
880 <<<
881 <<<
882 <<<
883 <<<
884 <<<
885 <<<
886 <<<
887 <<<
888 <<<
889 <<<
890 <<<
891 <<<
892 <<<
893 <<<
894 <<<
895 <<<
896 <<<
897 <<<
898 <<<
899 <<<
900 <<<
901 <<<
902 <<<
903 <<<
904 <<<
905 <<<
906 <<<
907 <<<
908 <<<
909 <<<
910 <<<
911 <<<
912 <<<
913 <<<
914 <<<
915 <<<
916 <<<
917 <<<
918 <<<
919 <<<
920 <<<
921 <<<
922 <<<
923 <<<
924 <<<
925 <<<
926 <<<
927 <<<
928 <<<
929 <<<
930 <<<
931 <<<
932 <<<
933 <<<
934 <<<
935 <<<
936 <<<
937 <<<
938 <<<
939 <<<
940 <<<
941 <<<
942 <<<
943 <<<
944 <<<
945 <<<
946 <<<
947 <<<
948 <<<
949 <<<
950 <<<
951 <<<
952 <<<
953 <<<
954 <<<
955 <<<
956 <<<
957 <<<
958 <<<
959 <<<
960 <<<
961 <<<
962 <<<
963 <<<
964 <<<
965 <<<
966 <<<
967 <<<
968 <<<
969 <<<
970 <<<
971 <<<
972 <<<
973 <<<
974 <<<
975 <<<
976 <<<
977 <<<
978 <<<
979 <<<
980 <<<
981 <<<
982 <<<
983 <<<
984 <<<
985 <<<
986 <<<
987 <<<
988 <<<
989 <<<
990 <<<
991 <<<
992 <<<
993 <<<
994 <<<
995 <<<
996 <<<
997 <<<
998 <<<
999 <<<
1000 <<<

```

```

0000454_18591230_Issue.xml [C:\Users\bcosta\AppData\Local\Temp\7z01B50.tmp\0000454_18591230_Issue.xml] - <Xygen> XML Editor
File Edit Find Project Options Tools Document Window Help
KPath 2.0 - [Execute XPath on 'Current File']
Project
sample.xml
The Master File support is disabled
Outline
Document name filter
Issue
page 1
page 3504709198
page 3
page 4
page 5
page 6
page 7
page 8
r:\metaData.xml | o:\openSearch.xml | a:\adft_nutrition.xml | 0000454_18591230_Issue.xml | ArtemisManuscript.dtd
1 <?xml version="1.0" encoding="UTF-8"?>
2 <DOCTYPE issue SYSTEM "TMIIssue.dtd">
3 <issue>
4 <metadataInfo>
5 <newspaperID>BSNP0454</newspaperID>
6 <mode>YGT</mode>
7 <assetID>348500486</assetID>
8 <divCollectionID>BNCNC0005</divCollectionID>
9 <cor>87.58</cor>
10 <language code="English" primary="Y">English</language>
11 <sourceLibrary>
12 <libraryName>British Library</libraryName>
13 <libraryLocation>London, United Kingdom</libraryLocation>
14 <copyrightStatement>Copyright © The British Library Board</copyrightStatement>
15 </sourceLibrary>
16 <PSMID>0000454_18591230</PSMID>
17 <id>1307</id>
18 <issue>
19 <id>XXVI</volNum>
20 <id>A</partNum>
21 <year>1859</year>
22 <month>December</month>
23 <day>30</day>
24 <composed>December 30, 1859</composed>
25 <searchableDateStart>18591230</searchableDateStart>
26 </id>
27 <id>Friday</div>
28 <id>B</ip>

```



```

File Edit Packages Windows Help
dataComplaints <- dataComplaints[rep(".*2014$", dataComplaints$DateReceived), ]
#Load State Polygons
states <- map_data("state")
#Extract complaints for each state
mapData <- as.data.frame(table(dataComplaints[3]))
mapData$Var1 <- tolower(state.name[match(mapData$Var1, state.abb)])
mapData <- mapData[complete.cases(mapData), ]
#Rename columns for merging
colnames(mapData)[2] <- "val"
colnames(mapData)[1] <- "region"
#Merge polygon and values dataframes
map_data <- merge(states, mapData, by="region", all=T)
map_data <- map_data[order(map_data$order), ]
#Plot Map
ccMap <- ggplot(map_data, aes(x=long, y=lat))
geom_polygon(aes(x=group, fill=val), col="NA", lwd=0) +
scale_fill_gradient(name="# Complaints", low="#FFDDDDDD", high="#FF0000") +
borders("state", colour="#666666", alpha=1)
map_data <- map_data[order(map_data$order), ]
plot(ccMap)
# Code for question 2
# CREATE PIE CHART OF COMPLAINT TYPES
countProd <- table(dataComplaints$Product)
countProd

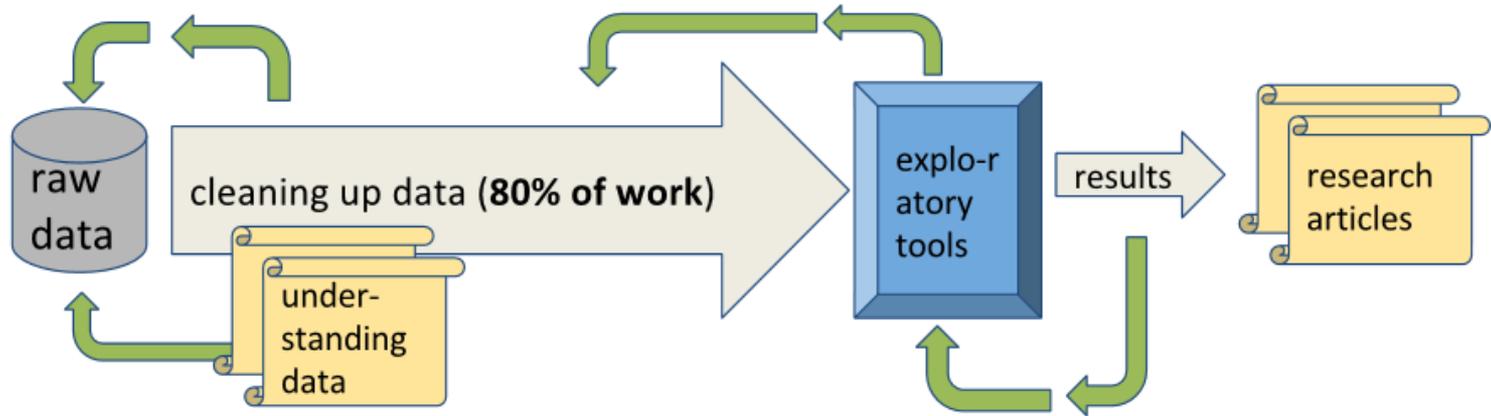
```

```

Notepad++ v6.9.2 new features and bug
1. Add most wanted feature: Log Morr
2. Add new feature: Find in Folder
3. Fix status bar display bug in his
4. Fix open in explorer problem whil
5. Fix smart highlighter issue after
Included plugins:
1. NppExport v0.2.8
2. Plugin Manager 1.3.5
3. Converter 3.0
4. Nima Tool 1.9

```

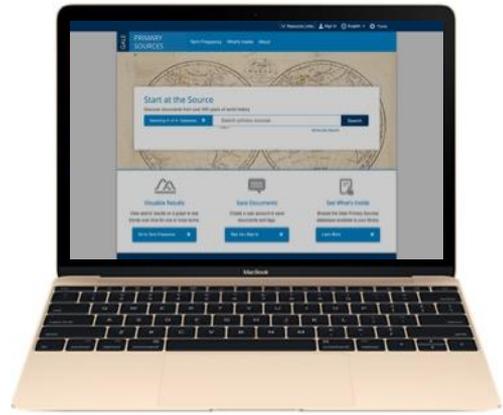
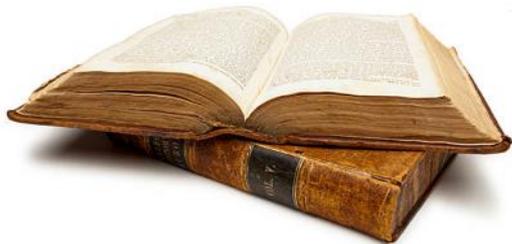
Text Data Mining (TDM) format: start with clean content



➔ 80% of your time for data cleanup, another 80% for algorithms, ...

Text and data mining Eighteenth Century based on ESTC & ECCO
COMHIS Collective
BSECS Conference 2017, Oxford
<http://j.mp/comhis-bsecs>

Digital Humanities ask different things from a text

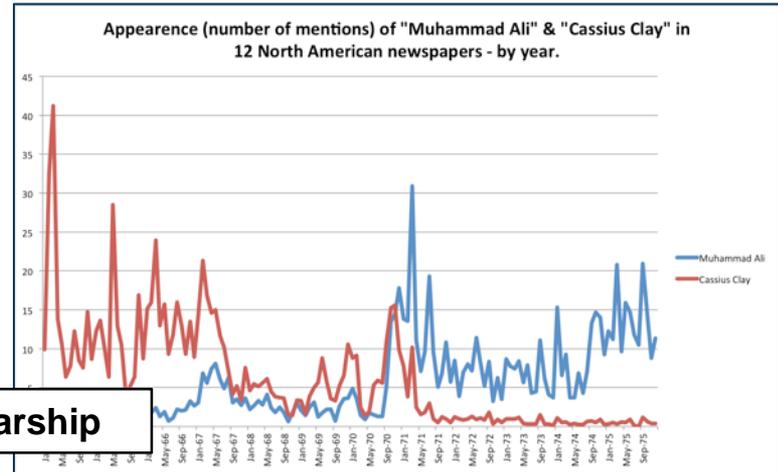


Close Reading vs. Distant Reading



What Close Reading Actually Means

Close Reading – Traditional Scholarship



Distant Reading – Digital Scholarship

Get started



Create a new project



Quick analysis



Create a new content set



Create a new analysis



Create a new result set

CONTINUE YOUR RESEARCH

In Search of Vida Marie
Created 4/14/2013, last accessed 8/15/2017The Language of the Blitz - Daily Mail, 1939-1940
Created 4/16/2013, last accessed 8/15/2017
Find the frequency of terms used to report UK current events of the following year.Tracking D-Day: Analytical Insights of the Longest
Created 2/16/2013, last accessed 8/15/2017OSCAR WILDE ON TRIAL
Created 2/20/2013, last accessed 8/15/2017
THIS PROJECT FOCUSES ON MATERIALS ABOUT WILDE'S TRIAL AND HISOscar Wilde on Trial, 1900-1900
Created 2/16/2013, last accessed 2/16/2017Oscar Wilde in the news
Created 2/16/2013, last accessed 2/16/2017
A project containing news articles about the illustrious Oscar WOscar Wilde on Trial, 1900-1900
Created 2/16/2013, last accessed 2/16/2017
A project containing news articles about the illustrious Oscar WOscar Wilde on Trial, 1900-1900
Created 2/16/2013, last accessed 2/16/2017
A project containing news articles about the illustrious Oscar W

The Language of the Blitz - Daily Mail, 1939-1940

Summary Items Items Data Item Trends Item Trend/Date Documents Documents Data

Select parameters for the left to customize the visualization.

