

Internet Archive (IA) Book Digitization and Quality Assurance Processes

Robert Miller
May 2007
Rev B

Confidential Document for use by
Internet Archive (IA), LPs (LP) and Digitization Sponsors (DS)

Overview- Technical and Operational details.

This section details in general, of the technical specs and the operation of how the equipment works.

1. Image capture-Color
2. Output formats
 - a. Color images in JPEG2000 format in pixels per inch listed below.
 - b. OCR in 2 XML formats: ABBYY and DJVU formats. ABBYY 6.0 is used, with its quality. As new versions and alternative vendors become available, a review will be coordinated between LP and DS before implementation. OCR XML character format is UTF-8.
 - c. XML for metadata from MARC.
 - d. XML for operational metadata collected during scanning.
 - e. Searchable PDF.
 - f. XML structural metadata for monographs include page numbers when apparent on the pages that is checked by the scanner operator.
These formats will be delivered from the Internet Archive servers the Internet via HTTP, FTP, RSYNC, or OAI.
3. DPI vs. Size
 - a. Example of DPI vs. size, chosen to optimally image a given size book.

DPI	Height (inch)	Width (inch)
300	14.6	9.7
400	10.9	12.5
500	8.7	6.58
600	7.3	4.9

Scanning Equipment For Bound and non-bound materials "Scribe"

Background-

The Internet Archive has tested and evaluated many commercially available scanning devices, but felt that due to the great variety of paper types, binding types and collections to be digitized, an in-house developed Scanning solution would provide the safest and, ultimately, a cost effective way of scanning books. The equipment shown below has been field tested and has successfully scanned millions of pages with virtually no damage caused by the equipment to materials being digitized.

Non-Destructive Scanning Station -

The Scribe workstation* - which is comprised of a frame that holds two cameras on a rail, to capture both the verso and recto pages of the book to be digitized, a cradle that the book sits in (the spring supported cradle is 'v' shaped so there is minimal stress put on the book), a glass platen that is raised and lowered by means of a foot pedal to allow for the pages to be turned, two banks of lights that illuminate the book and two small computers to run the cameras and pre-process the images. The captured images are, upon the completion of the book being scanned and the QA process being completed, are uploaded via RSYNC to processing computers located in California.



* Mechanical/Electrical parameters- 7 amps per Scribe, 800 watts of heat generated per Scribe, standard UK/USA/ voltage, appx 100 sq feet per Scribe of work area, Scribe footprint is 68" long x 37" wide x 79 high; (Dimensions for shipping with out crate- 60" long-rails removed, x 32" wide- remove monitor/arm x 79in high; Dimensions for shipping with crate- 68" long x 42" wide x 86" high, 609 pounds)- NOTE check all doors for access- the weight is appx 300 lbs. For a scanning center of 10 Scribes appx 1000 sq feet of space is required.

Workflow- Image capture

1. The Scribe scanner currently captures page images using a pair of digital single-lens-reflex (DSLR) cameras, either a 16.7 mega-pixel [Canon 1DS-Mark II](#) or a Canon EO 5D, 12.8 mega-pixel, with a Canon EF 100mm f 2.8 macro lens (<http://www.usa.canon.com/app/pdf/lens/EFLensChart.pdf>). IA is always evaluating new cameras and if a better solution comes along, after coordinating with LP and DS.
2. The lighting system for illuminating the target books consists of eight (8) 5000 Kelvin, 36 degree, 35 watt museum-grade Solux bulbs, and provides a smooth daylight spectrum with a high color-rendering index.
3. Lighting compensation program- To help make the lighting even across images being scanned.
4. Reference targets- A color target (ColorChecker 24) and a white card is shot with each book for reference, which can be used for ICC-based color management.
5. Image transfer- Page images are downloaded in real time to a scanner management and image processing computer which also run the camera management software that releases the camera shutters.
6. Equipment Calibration
 1. Cameras are calibrated per Manufacturer's spec. Cameras out of spec or standard performance will be sent back to Manufacturer for repair.
 2. Lights used in scanning process are replaced as is necessary. Light comp algorithms are run daily on each unit of scanning equipment.
 3. Scribe stations are calibrated and aligned before being used.

Overview- Process steps between Library and Internet Archive (IA)

This document outlines the general flow of information and sequence of steps between the IA and the library. At the end of this document is the QA section that details how the process is checked, maintained and updated to ensure compliance.

Process steps-

A. Meta Data- for all new collections that require attribution or at the beginning of a scanning center set up.

1. Meta data and set up form- includes Contributing Library, Digitization Sponsor, Collection Name, Contact details, etc. to be filled out by Library and sent to Jae@archive.org .
2. Meta data set up is then incorporated into IA scanning/loading screens. This is to ensure proper attribution and organization of the materials is being completed.
3. If the Z39.50 set up is not being used to locate Library catalog records, alternate confirmation of how records are to be located must be agreed to by LP and IA. A test pick list of at least 50 records, see below, and should be generated to test that IA can locate the proper MARC records.

Collection	BOOK ID	BIB ID	Title	Author	Vol.	Item Barcode	Call	Item Location	Item Status	Item Date
BANF 4.11.07	historicaldictio00wilebost	2368353	historical dictionary of Jainisr	Kristi L., 19		3.9999E+13	BL1303.W45 2004	BPL1ST	i	
BANF 4.11.07	padaindexreverse00yamabost	1523716	ada index to early Jain canormazaki, Mori			3.7777E+13	L1310.5 Y36 1996	BPL1ST	i	
BANF 4.11.07	studienzumsuyaga00bollbost	560193	Jainas u.d. anderen Weltansollee, Willem E			3.7777E+13	BL1312.3.S886 B6	BPL1ST	i	
BANF 4.11.07	thisisreality00potabost	1686912	This is reality	Potanjali		3.7777E+13	L132.Y6 P267 196	BPL1ST	i	
BANF 4.11.07	Assemblyoflisten00carrbost	1003270	sembly of listeners : Jains in			3.7777E+13	BL1320 .A87 1991	BPL1ST	i	

4. Any changes or requests by the Library that would impact the Meta data must be in writing; for example, a new collection, a new funding source or a new sub-collection.

B. Preservation and handling-

1. Library preservation personnel will meet with IA to establish and agree on how to handle materials to be scanned, how to deal with obvious rejects, how to flag and tag materials not able to be scanned, agreement on error codes and the like. Deviations from this process will

be in writing and where possible all steps are documented with a visual example for reference.

2. A review between IA/LP on what materials can/can't be scanned is conducted. Questionable materials may be tested or tried before being put into the scanning plan. All selection of materials will follow these guidelines, See section F.
3. Rare materials are defined as materials that would normally not be in circulation or should not be included in the general scanning population.

C. Retrieval, packing and shipment of Materials by LP-

1. Materials, meeting IA specifications, see F, will be delivered to the Internet Archive.
2. Materials will be packed on cart and wrapped for transport and shipping, unless otherwise agreed to. Any special procedures will be determined in advance.
3. A paper and a digital (excel) copy are to travel with each shipment of materials delivered to Scanning Center.

(Not applicable at University of Toronto)

D. Receipt of shipment by IA prior to scanning

1. IA will match the count on the pick list to the materials on the cart. If count matches, books will move to the loading area. If the book count does not match- IA will alert Library.

E. Inspection of materials to be scanned by IA prior to scanning

1. IA, as it loads each item to be scanned, will inspect each item for possible factors that would impact it's ability to be scanned. Rejected materials will be marked in a pre-agreed format and returned with books scanned. These items may be scanned at a later date as new processes become available or a different cost structure is put in place.

F. Digitization Criteria to be used by Library and IA- see addendum

1. Criteria to be used for determining if the materials may be digitized are listed below and will include, but not necessarily be limited to the following:
 1. Any special preservation standards that came on cart from the LP.
 2. Materials that have multiple titles/physical volume (eg, 63 vols, 20-40 pamphlets per vol, of "Forestry Pamphlets" without analytics) will be reviewed to ensure all proper meta data is understood.
 3. Materials will be screened for size- materials not fitting the requirements shown below will be returned unscanned.

4. Size-

- 9.7 wide by 14.5 high at the max. Books as small as 3 inches x 3 inches may possibly be scanned.
- less than 3 inches thick are ok, greater than 3 inches will need to be reviewed.
- Books should on average should be 200 pages or larger. If a collection is mostly under 100 pages a review should be undertaken between IA and DS to ensure the quoted price per scanned page can be met.

Comment [MC1]: Robert will determine Cost for special production channel in June 07

5. Book Style-

1. Side bound Monographs, no single sheets, no top bound books.
2. Rebound books need to be checked for how tight the gutter or binding is, or if the text runs outside the margin (it will show the cradle).
3. When there are more than one title within a bound book; each of these has to be clearly marked with a paper strip; each of these will be counted as a separate book. Deviations from this must be approved between IA/LP.
4. Fold outs are rejected now, but will be included as part of a special production channel.
5. Soft cover books are ok if they are bound.
6. Covers that are almost separated from book and appear too fragile may be rejected unless agreed to in advance by LP.
7. Material condition-Materials should be of similar condition or quality to what would be put into circulation. Materials deemed not robust to go into circulation should be reviewed with Scanning Center Coordinator before scanning.
8. Tight bindings that will not lay open for digitization per IA specification limits will be rejected.

Comment [MC2]: Robert will determine Cost for special production channel in June 07

6. Paper style-

1. Most paper styles can be scanned, except highly acidic paper that disintegrates to the touch. Note that if a hard-to-scan paper is to be digitized; a review of time to scan versus any 'damage' will be undertaken.
2. All pages should be pre-cut. Unless otherwise instructed, books with uncut pages will not be scanned.
3. Pages should be able to be lifted and turned with normal effort. Sticky pages or the like will not be scanned.

4. Pages should not be excessively dusty, have excessive mildew or be moldy.
 5. Microfilmed reproductions should be reviewed with Scanning Center before being scanned. If the paper/image looks like a film negative (Xerox), these won't be scanned. If the microfilm reproductions have more than one page on each leaf these will be rejected.
7. Gutters/Margins-
1. Any book where the text is less than a quarter inch off the gutter, on an approximate 75-degree angle will be unscannable.
 2. Text that runs to the edge of the page or margin can be scanned but the presentation will be poor, as the cradle will show. The LP must approve this.
8. Bibliographic data-
1. Multiple books or multi volume set-The usual problem with bibliographic info is that there is only one bib record for a set. The LP will make a decision as to how this should be handled.
 2. Books that are out of approved copyright range. These will be set aside, unless otherwise agreed to by LP and IA.
 3. A book will be set aside when a Marc record can't be located when a Call ID/Book ID or equivalent is inputted.
9. Rejection codes to be sent back to Library are:

Code	Definition
BI	Fragile or no binding (includes items in clam shells or phase boxes)
CAT	Cataloging error
DAM	Damaged
DAT	still in copyright
FO	Foldouts
LG	too large
MAR	margins too tight
MIS	Missing pages
MUL	multiple titles bound together
NA	not available
LAN	Outside language parameters
LIST	Picklist error
LINK	unsuccessful link to metadata
NOS	not on shelf – missing/lost
OUT	not on shelf – checked out
PAG	Pagination problems: section(s) bound out of order or upside down
PAP	Brittle paper, tissue paper

SKW	Skewed text – to point of being unreadable
SM	too small
UNC	Uncut pages (more than 5)
FOR	non-book format
VEL	Vellum
WD	Withdrawn
SPH	requires special handling
DUP	exact duplicate of another on list

G. IA loading process-

1. The book id is loaded into IA Screen to located appropriate MARC record- records not found will be cause for reject of book.
2. If IA receives a series that is cataloged under one bib record, without volume numbers, IA will either add the year as a volume number or add the volume number in IA's set up screen. This would be analogous to, in the physical world, where the series are shelved together, a person locating the book(s) from the call number and then scanning the items on the shelf for the volume they want. IA will not in a series, delete any information from the descriptions field or add to the descriptions field on the MARC record.
3. QA check is done to ensure book and MARC record match.
4. Unique identifier is created and MARC record is attached to that identifier.
5. Book is placed in queue for scanning.
6. As mentioned above, any books determined not to be scannable are set aside and a rejection form is attached.
7. Each book is also given a color-coded flag that shows the Book identifier.
8. Any special scanning instructions are included with book.

H. IA scanning process-

1. Materials to be scanned are placed in queue for scanners, typically on book carts.
2. The flag inside the material to be scanned is matched to digital file to ensure a proper match.
3. Images are scanned into appropriate digital file.
4. Images are QA'd following adjustments to ensure proper preservation or presentation.
5. Digital file is closed and uploaded to IA processing center.

I. IA processing

1. Uploaded images are processed to create storage files and access files.
2. File Formats, are covered in the Technical Spec and Equipment section above.
3. Files delivered for download to LP are listed above.
4. Detailed list of files
 - a. ID.pdf
 - b. ID_jp2.zip - will not be accessible, this is only for long term preservation
 - i. zipped folder of the book without bookplate and watermark is specific to the sponsor, contributor; e.g. this will vary for each sponsor/library.
 - ii. [ID]_nnnn.jp2 (where first image index number is the front cover and the last scan # is the back cover)
 - c. ID_lib_jp2.zip
 - i. zipped folder of the book with bookplate and watermark
 - ii. [ID]_lib_nnnn.jp2 (where image index the first number is front cover and the last scan # is the back cover)
 - d. ID_marc.xml
 - e. ID_meta.mrc
 - f. ID_meta.xml
 - g. ID_metasource.xml
 - h. ID_raw_jp2.zip , unprocessed storage format, no watermark/book plate
 - i. Scandata.zip
5. Metadata will reside in meta.xml file, and will include the following required fields for the library:
 - a. Identifier
 - b. ARK (begins ark:/13960/*) ; this is an experimental field for California Digital Library
 - c. Collection-Library (from pick list); this applies only for UC Libraries
 - d. Identifier-bib (unique identifier -- local catalog number; from pick list), this only applies to UC libraries
 - e. Contributor
 - f. Title
 - g. Volume
 - h. Creator (if in MARC record)
 - i. Publisher (if in MARC record)
 - j. Collection (possibly multiple collection fields)
 - k. Operator
 - l. Scanner
 - m. Scandate

n. Identifier-access (URL for accessing this book)

6. **Processing Background-**

The digitized image is captured initially as a camera raw file (CR2). This is run through a JPG 2000 compression to generate a raw JPG 2000 for storage. The raw JPG 2000 is then turned into a processed master which is used to generate the access formats.

- i. Storage format- raw JPG 2000 is a compressed, lossy, uncropped, non-rotated, non-deskewed, non-light comp'd JPG 2000 file; which is the storage file. Image sizes vary depending on the complexity of the page, but are typically in the 900 KB range, yielding an approximate compression ratio of 15:1 relative to the camera raw image (CR2 is appx 15MB/image.)
 - ii. Processed master- lossy, cropped, rotated, de-skewed, light comp'd JPG 2000. Image sizes may vary depending on complexity of the page, but are typically in the 800 KB range, yielding an approximate compression ratio of 15:1 relative to the camera raw image (CR2 is appx 15MB/image).
 - iii. Access format- the processed JPG 2000 masters are compressed in a JPG 2000 format which feeds into the OCR and book generation tools. Image sizes may vary depending on the complexity of the page, but are typically in the 760 KB range, yielding an approximate compression ratio of 20:1; relative to the camera raw image (CR2 is appx 15MB/image). PDF and DjVu; both of which are OCR'd.
- b. Quality settings will vary based on vendor tools used. For example a quality setting of 50 on a scale of 1-100 was used for the Luratech. This setting was determined based on user surveys.

K. Turnaround for processing by IA- typically 72 hours from arrival to return of book cart.

1. The goal is to derive and upload a book within 24 hours after scanning.
2. An internal IA QA step is performed inside the scanning center. Criteria for QA are outlined below in the Quality Section.
3. If the scanning lot is rejected, then the process outlined in the Quality Section K is undertaken-
4. A scanned item is then published online within 48 hours after scanning.

5. Materials scanned are then 'curated' by IA and are available for downloading after that by the LP.
6. Approved Materials having been scanned are then ready to be checked out and returned to the LP.

L. IA Scanning Center Check-Out Process

1. Scanning coordinator packs book into shipping cart/container per guidelines established between LP and IA.
2. Creates and attaches the report communicating books rejected for scanning and identify failure.
3. Books transferred to LP.

QA Plan-

Overview- There are four major phases to the QA process:

- 1) Before the materials are uploaded-At the book loading and scanning station; the scanners looks for; amongst other things, missing pages, crop/deskew problems, page marking (title page, front/back cover, tissue paper, first page of table of contents and notes any defects in the book (i.e. Missing/torn pages).
- 2) After the images are uploaded, derived and available via an URL- a statistical sampling and QA is conducted within the Scanning Center. Per ANSI z1.4 1993 Table 1, General Level 2. Details; see below.
- 3) Before the curation and bill is generated- An internal random audit is conducted outside the scanning center before the final curation approval and bill is generated.
- 4) After the materials are received by the LP and the DS. Errors brought to IA's attention will be dealt with in a timely (appx monthly) basis. A decision will be made by IA as to whether it is best to rescan the material or fix it post-derive. Rescanning is to be avoided as it requires generating a new URL and is usually the most expensive solution. The timeframe for the library to identify errors that will be fixed by the IA at no-charge shall be detailed in the digitization plan.

Scanning Center QA: IA uses ANSI z1.4 1993 Table 1, General Level 2 <http://www.proqc.com/dl/aql.pdf>

Each day the scanning center will review a set of books from the previous days scanning. The number of books to QA depends on the total number of books in the set.

Books in set	9-15	16-25	26-50	51-90	91-150	151-280
number to QA	3	5	8	13	20	32

The scanning center coordinator is responsible for choosing a representative set to reflect a mix of scanners/scribes and conforming to the statistical chart.

A. QA Process steps

Books are inspected for the criteria below (see Freeze codes B1-B 3) ; on-line, using the relevant files for each coded; found in the digital book record including pdfs to look for errors or defects. Errors or defects, if found, are noted and added to the IA meta manager form. An automatic scoring is then performed and a “pass/fail” grade is assigned to the lot.

bin	min	max	level	sample	major_accept	major_reject	minor_accept	minor_reject
1	9	15	2	3	0	1	0	1
2	16	25	2	5	0	1	0	1
3	26	30	2	8	0	1	1	2
4	31	90	2	13	1	2	1	2
5	91	150	2	20	1	2	2	3
6	151	200	2	32	2	3	2	3
7	201	1,000,000	2	99	7	10	7	10

Explanation- If 125 books were scanned in a period to be inspected, bin 5 would be selected. According to the truth table above, if there were 1 Major error or less and 2 Minor errors or less, the lot is passed. If there are 2 or more major defects or 3 or more Minor rejects, the lot fails. See below for what happens after this. The major/minor detail is show below in B5. Note: for major defect found during QA, they will be repaired on that book even if the lot passes.

- If the lot passes, the Scanning Center will approve all books (Curate).
- If a “fail” is generated, the Scanning Center Coordinator will review the errors/defects to ascertain if the errors were generated from outside the Scanning Center (for example a missing access file error would be sent to engineering for review) or from within the Scanning Center (for example a missing page).
 - If the error was generated from within the Scanning Center the Coordinator would follow a pre-determined set of process steps ultimately culminating in a recommendation to deviate or approve the lot or a portion of the lot with appropriate corrective actions identified. At this stage the Book’s Director or the Headquarters QA staff person is involved and must approve a deviation. A corrective action report will be generated for rejected lots. This will be reviewed with management for longer-term solutions or corrective action. This is done daily.

B. Codes used and shown on the QA report are:

B1. FREEZE CODES, part I. (100-113, and 130-138 will ultimately be rescanned and the original URL made dark (not publicly available), a new URL generated and this will be communicated to the Library by IA via email)
Books in general that can't be corrected post-derive, will have to be rescanned. If it eventually turns out that IA can't rescan or fix a book, it will not be billed.

-- Formats --

101 Test book

102 [DjVu](#) is missing or corrupt

103 PDF is missing or corrupt

104 Flip book is missing or corrupt

Resolution of errors found here- Material is rederived, if that doesn't correct the problem, the material is rescanned.

-- Uploading or piping problems --

110 Truncated file(s)

111 Book deleted from scribe before upload completed

112 Missing files(s)

113 Cr2.tar file is malformed

Resolution of errors found here- Material is rescanned.

-- Metadata --

120 Book is not in public domain

121 Date is 1923 or later

122 Date is unclear

123 Date is 1923 or later-

Resolution of errors found here- If material is in copyright, item is removed from Search engine. If material is in question, Library is consulted and appropriate action is taken.

-- Images --

130 Cropped text

131 Blurred page(s)

132 Missing page(s), goal is zero pages missing per book. IA will note if pages are missing in the books to be scanned.

133 Front cover missing

- 134 Back cover missing
- 135 Book was scanned twice; identified copy is darkened and removed from search engine.
- 136 Text is washed out or overly dark (bad light-comp)
- 137 Evidence of scanner (fingers/shadows/etc) visible on page
- 138 Glass not centered in gutter; text is distorted or cropped

Resolution of errors found here- For items 130-134, material is rescanned. For items 136-138 a review is made by IA to deviate, accept or rescan.

B2. FREEZE CODES, part II.

Use these codes for books that have fixable problems, but are not yet in billable condition.

- 140 Book and metadata do not match
- 141 n/a
- 142 Tissue pages marked incorrectly
- 143 Anomaly in image format is under investigation
- 144 Left/right pages are reversed

Resolution of errors found here- For items 140 and 142, post derive correction is attempted. For item 144, if post derive correction won't work, material is rescanned.

B3. INFORMATIONAL CODES

- 150 Bibliographic data missing:
- 151 De-commissioned
- 152 Copyright evidence reported incorrectly. Info corrected in QA.
- 153 Bibliographic record from library is truncated
- 154 Possible error in bibliographic record from library
- 155 Foreign language character encoding is incorrect
- 156 Incorrect of missing collection-library or bib-id

- 160 Light/dark pages (intermittent)
- 161 Light/dark pages (throughout),
- 162 Pages skewed
- 163 Color cards show in access formats
- 164 White cards show in access formats
- 165 Both white cards and color cards show in access formats
- 166 Image of cradle is visible at front or back
- 167 Different crop-box sizes in same spread

168 Bad crop at page edges
169 Duplicate page spreads scanned

170 Page types not marked or marked incorrectly
171 Title page not marked b/c book does not have title page
172 Scan factors not marked or noted

198 This would be a good display book
199 Approved with no problems noted

Resolution of errors found here- For items 150, 152, 155 and 156, post derive correction is possible. For errors, 153 and 154, errors must be reviewed with Library. For errors, 160, 161, 162, 166, 167, 168 a review is conducted to see if the material is acceptable for OCR'able texts. This is accomplished by using the OCR function on the word in question. If the material can't be OCR'd based on the IA software being used, then the book would be rejected For errors, (only IA internal requirements) 166, 170 and 171 a post derive correction is attempted, if unsuccessful materials must be rescanned.

Error resolution falls into three forms of correction; post-derive, rescanning or a consultation with the library. If materials can't be corrected with either of these methods, than the book is rejected.

The following errors may be attempted to be corrected by post-derive treatment:
110, 111, 112, 113, 140, 142, 144, 150, 152, 155, 156, 160, 161, 162, 166, 167, 168, 170, 171

The following errors may require rescanning:
130, 131, 132, 133, 134, 136, 137, 138

The following errors require consultation with the Library.
120, 121, 122, 123, 153, 154

B 4- Rescanning process-

1. For materials that are to be rescanned, a request to pull those books requiring rescanning is submitted to the Library; usually once a month. Materials are pulled, scanned and the original item is removed from the Internet Archive search engine and a new URL is assigned. This new URL is sent to the LP and the DS along with the old URL for reference. A bug report could be the means to track this process.

B5- Error Codes & Classes

class/id	Description	type	defects	Totals
Formats				
101	Test book	major	0	
102	DjVu is missing or corrupt	major	0	
103	PDF is missing or corrupt	major	0	
104	Flip book is missing or corrupt	major	0	
				major: 0
				minor: 0
				total: 0
				status: OK
Uploading or piping problems				
110	Truncated file(s)	major	0	
111	Book deleted from scribe before upload completed	major	0	
112	Missing files(s)	major	0	
113	Cr2.tar file is malformed	major	0	
				major: 0
				minor: 0
				total: 0
				status: OK
Metadata				
120	Book is not in public domain	major	0	
121	Date is 1923 or later	major	0	
122	Date is unclear	major	0	
123	Date is 1923 or later	major	0	
				major: 0
				minor: 0
				total: 0
				status: OK
Images				
130	Cropped text	major	0	
131	Blurred page(s)	major	0	
132	Missing page(s)	major	0	
133	Front cover missing	major	0	

134	Back cover missing	major	0
135	Book was scanned twice; this copy darkened.	minor	0
136	Washed-out text (bad light-comp)	minor	0
137	Evidence of scanner (fingers/shadows/etc) visible on page	minor	0
138	Glass not centered in gutter; text is distorted or cropped	minor	0
140	Book and metadata do not match	minor	0
141	n/a	major	0
142	Tissue pages marked incorrectly	major	0
143	Anomaly in image format is under investigation	major	0
144	Left/right pages are reversed	major	0
		major:	0
		minor:	0
		total:	0
		status:	OK

Bibliographic

150	Bibliographic data missing: MetaFetch was not run in scanning center (post MF done in QA)	minor	0
151	Bibliographic data missing: MetaFetch was run but did not merge (post MF done in QA)	minor	0
152	Copyright evidence with reported incorrectly. Info corrected in QA.	minor	0
153	Bibliographic record from library is truncated	minor	0
154	Possible error in bibliographic record from library	minor	0
155	Foreign language character encoding is incorrect	minor	0
156	Incorrect or missing collection-library or bib-id	major	0
		minor:	0
		total:	0
		status:	OK

Consistency

160	Light/dark pages (intermittent)	minor	0
161	Light/dark pages (throughout)	minor	0
162	Pages skewed	minor	0
163	Color cards show in access formats	minor	0

164	White cards show in access formats	minor	0
165	Both white cards and color cards show in access formats	minor	0
166	Image of cradle is visible at front or back	minor	0
167	Different crop-box sizes in same spread	minor	0
168	Bad crop at page edges	minor	0
169	Duplicate page spreads scanned	minor	0
170	Page types not marked or marked incorrectly	minor	0
171	Title page not marked b/c book does not have title page	minor	0
172	Scan factors not marked or noted	minor	0
		major:	0
		minor:	0
		total:	0
		status:	OK

2. To track errors found by LP or the DS a bug tracking system could be used. Response times for error resolution will be determined based on type of error and time of response. In general, IA will attempt to resolve errors brought to its attention within 30 days of an error being identified.

C. Meta Manager, the post scanning reporting tool

This is the reporting tool that the Library may use to search and review books that have been scanned, uploaded, QA'd and then curated. The curation stage is the last stage in the IA process where the books are made viewable to the Library. This may happen on a non-scheduled basis but is typically done several times a month.

Internet Archive **Meta-manager** | view metadata | selected columns | default columns | or all columns | select columns | help -- user: [agal](#) - [logout](#)

Welcome to the [Internet Archive's](#) Metadata Manager. This is a view of selected metadata for selected items in the archive. Please refer to the [help](#) page for more information.

Results (1-25) of 637 [[prev](#)] [[next](#)] - limit to [25, 50, 100] - show all | [un-filter](#) | link to [this query](#) | other [stats](#) | [qa report](#)

where: [hw_identifier="lond"](#) | [imagecount: 94,809](#) | [size: 237,631,886 KB](#) | [redrows: 0](#)

Operate on selections: [[QA](#)] [[task](#)] [[update](#)] [[export](#)] [[ftpdist](#)]

sort	e	date	identifier	title	mediatype	collection	imagecount	sponsor	uploader	postcopyright
1	<input type="checkbox"/>	1699	femalelegacy02textlond	The Female's Legacy, Containing 4...	texts	lendtest	116	MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
2	<input type="checkbox"/>	1817	pitclubofscotland02textlond	The Pitt Club of Scotland. A poem	texts	lendtest	60	MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
3	<input type="checkbox"/>	1806-08	verseshymns02textlond	Verses, hymns and poems, etc.	texts	lendtest	180	MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
4	<input type="checkbox"/>	1943 12*	selectofables02textlond	Select Fables	texts	lendtest	244	MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
5	<input type="checkbox"/>		fourhundred00nateclond		texts	americana	20	Internet Archive		
6	<input type="checkbox"/>	1817	pitclubofscotland00textlond	The Pitt Club of Scotland. A poem	texts	lendtest	66	MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
7	<input type="checkbox"/>	1806-08	verseshymns02textlond	Verses, hymns and poems, etc.	texts	lendtest	178	MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
8	<input type="checkbox"/>	1699	femalelegacy0004textlond	The Female's Legacy, Containing 4...	texts	lendtest	116	MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
9	<input type="checkbox"/>	1843 12*	selectofables00textlond	Select Fables	texts	lendtest	244	MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
10	<input type="checkbox"/>	1799	pleasuresofhope00complond	The pleasures of hope; with other...	texts	lendtest		MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
11	<input type="checkbox"/>	1778	poeticaltrifles00textlond	Poetical Trifles. By *** **** **...	texts	lendtest		MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
12	<input type="checkbox"/>	[1796]	cleaningsfromthe00textlond	Cleanings from the works of Laun...	texts	lendtest		MSN	Nathan@archive.org	NOT_IN_COPYRIGHT
13	<input type="checkbox"/>	1731	guelimussus0004textlond	Guelimius Susanna valedicens, [A...	texts	lendtest	24	MSN	robert@archive.org	NOT_IN_COPYRIGHT
14	<input type="checkbox"/>		brideoftheisles0004textlond	The bride of the Isles; a mask [...]	texts	lendtest	98	MSN	natetrain@msn.com	NOT_IN_COPYRIGHT
15	<input type="checkbox"/>		comediesandtragedies00textlond		texts	americana	64	Internet Archive		
16	<input type="checkbox"/>	182-7	comediesofprova003vanblond	The provoked husband ; a comedy, ...	texts	lendtest	88	Internet Archive		NOT_IN_COPYRIGHT
17	<input type="checkbox"/>		citizenscomplaint00textlond	The Citizens Complaint for want o...	texts	lendtest	32	MSN	natetrain@msn.com	NOT_IN_COPYRIGHT
18	<input type="checkbox"/>		grave_0004lond	The Grave	texts	lendtest	30	MSN	natetrain@msn.com	NOT_IN_COPYRIGHT
19	<input type="checkbox"/>		battleofhesham0007comlond	[The Battle of Hesham ... with re...	texts	lendtest	50	MSN	natetrain@msn.com	NOT_IN_COPYRIGHT
20	<input type="checkbox"/>		guineavoyage0004textlond	The Guinea Voyage, a poem.	texts	lendtest	118	MSN	natetrain@msn.com	NOT_IN_COPYRIGHT
21	<input type="checkbox"/>		apologyforthe0004textlond	An Apology for the believers in t...	texts	lendtest	656	MSN	scanner-stefan@archive.org	NOT_IN_COPYRIGHT
22	<input type="checkbox"/>		thecastleofindolence00textlond	[The castle of indolence, 1748. I...	texts	lendtest	66	MSN	natetrain@msn.com	NOT_IN_COPYRIGHT
23	<input type="checkbox"/>		thecorinthy0004textlond	[The Corinthy; or, the Gentlemen...	texts	lendtest	718	MSN	natetrain@msn.com	NOT_IN_COPYRIGHT
24	<input type="checkbox"/>		shakespearesgen0004textlond	Shakespeare's Genius justified; b...	texts	lendtest	502	MSN	scanner-stefan@archive.org	NOT_IN_COPYRIGHT
25	<input type="checkbox"/>		adoptedchild0004textlond	[The Adopted Child, a musical dra...	texts	lendtest	40	MSN	natetrain@msn.com	NOT_IN_COPYRIGHT

Results (1-25) of 637 [[prev](#)] [[next](#)] - limit to [25, 50, 100] - show all | [un-filter](#) | link to [this query](#) | other [stats](#) | [qa rep-ort](#)

steve@archive.org

Fields seen by LP and DS in the Meta-Manager view will also be inspected, plus several internal fields. These fields will include:

- identifier
- title
- creator
- collection
- image count
- contributor
- sponsor
- sponsor date
- scandate
- curatenote
- curate date

D. Library card, required to view the Meta Manager

An IA library card and an email are required to view the metamanager (see steps listed below). Here is the process to access the meta data page.

- 1 Go to www.archive.org
2. Go to Patron info
3. Click on "get a virtual library card"
4. Have Jae Jae@archive.org create the sponsor view and attach the library card info to that view
5. Books may then be viewed that have been curated

Changes to these Processes

Proposed process changes that would impact LP or DS will be communicated for review and discussion prior to being implemented.

Contact List- For IA Scanning Center

Contacts-IA- Note some of these will change based on specific IA scanning center. Contact Robert Miller as primary point of contact for initial program review.

Name	Role	Phone	Email	Fax
Jae Mauthe	Meta Data setup	415 810 5972	jae@archive.org	
See Digitization Plan	Pick List and shipment			
See Digitization Plan	Shipment			
Robert Miller	IT	415 640 1092	Robert@archive.org	
See Digitization Plan	Scanning center manager			
Robert Miller	Project manager	415 640 1092	robert@archive.org	
Marcus	Metamanager	415 561 6767	marcus@archive.org	
Marcus Lucero	Quality	415 561 6767	marcus@archive.org	
Robert Miller	Download	415 640 1092	robert@archive.org	
Marcus Lucero	Metadata	415 561 6767	marcus@archive.org	
Robert Miller	Eng questions	415 640 1092	robert@archive.org	

Addendum- Ramping up a Scanning Center

Each of the segments below is to confirm process, performance and expectations as a scanning center is ramped up. Note: This primarily pertains to new scanning centers versus new collections.

IA Engineering pilot- 5 books scanned by supervisor. Confirm pick list, packing list, z39.50, MARC record, IA deriving/processing/posting with attribution. URLs are QA'd by Jae. URLs are then sent to LP and DS for review.

IA Production pilot- 50 books scanned by scanner (5 books from each Scribe). Confirm shipping method to IA, material handling, pick list, packing list, z39.50, MARC record, deriving/processing/posting with attribution, return shipping procedure, metamanager and OAI. URL's are QA'd by supervisor. URLs are sent to LP and DS for review.

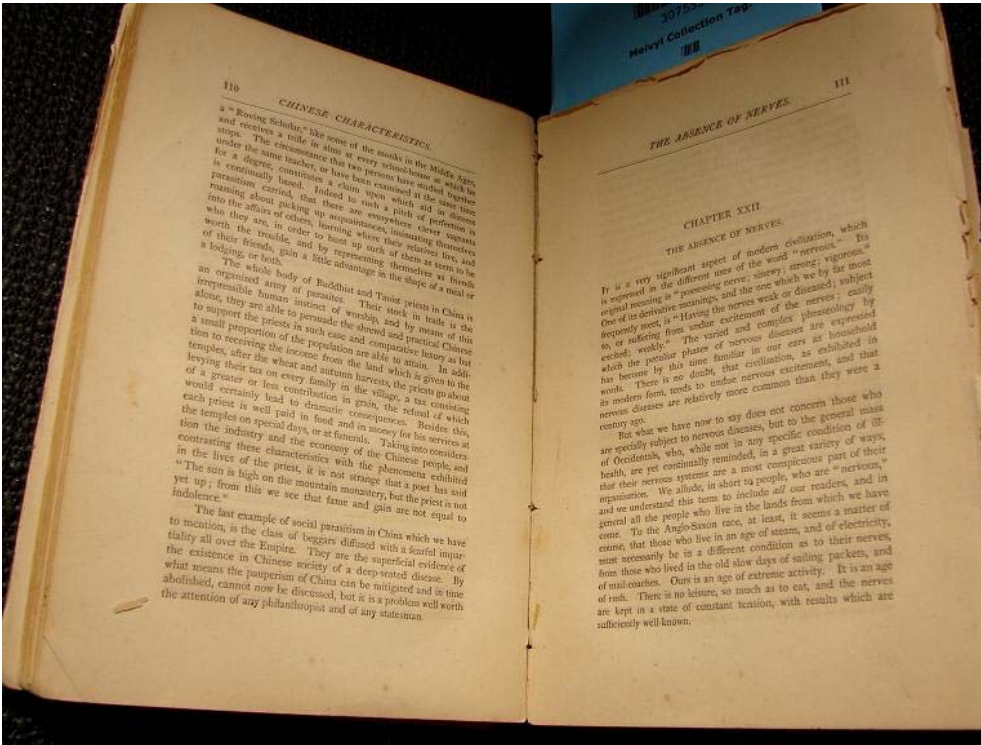
IA Production- establish agreement on quantity and expected turnaround.

Addendum- Examples of rejected books. Books in question maybe tested for scannability

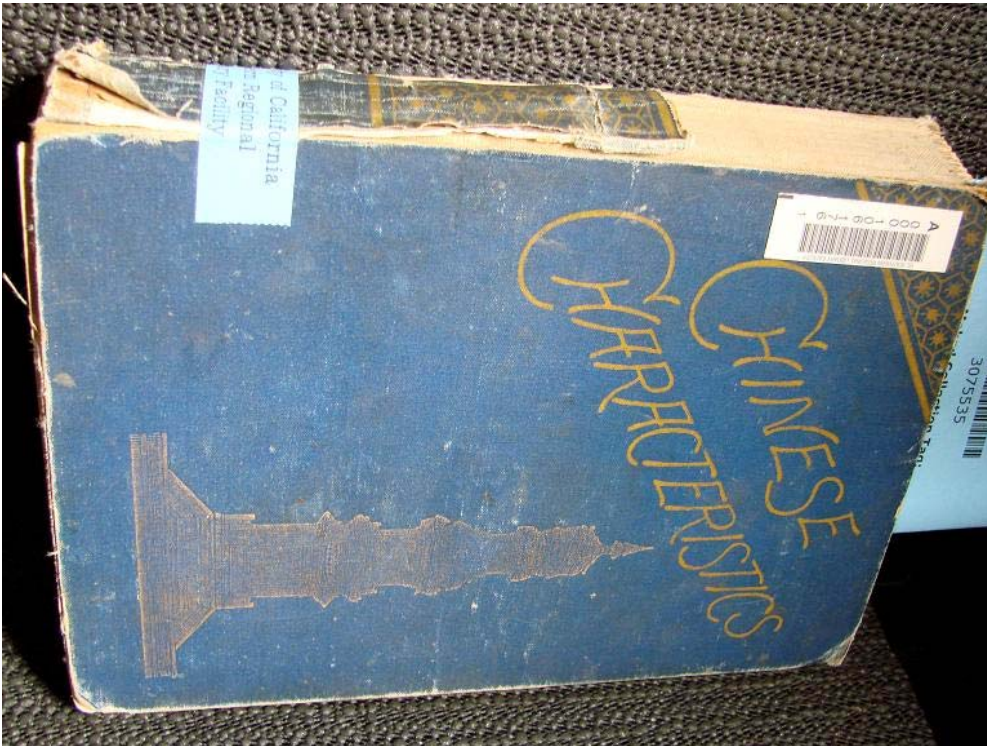
Code	Definition
BI	Fragile or no binding (includes items in clam shells or phase boxes)
CAT	Cataloging error
DAM	Damaged
DAT	still in copyright
FO	Foldouts
LG	too large
MAR	margins too tight
MIS	Missing pages
MUL	multiple titles bound together
NA	not available
LAN	Outside language parameters
LIST	Picklist error
LINK	unsuccessful link to metadata
NOS	not on shelf – missing/lost
OUT	not on shelf – checked out
PAG	Pagination problems: section(s) bound out of order or upside down
PAP	Brittle paper, tissue paper
SKW	Skewed text – to point of being unreadable
SM	too small
UNC	Uncut pages (more than 5)
FOR	non-book format

VEL	Vellum
WD	Withdrawn
SPH	requires special handling
DUP	exact duplicate of another on list

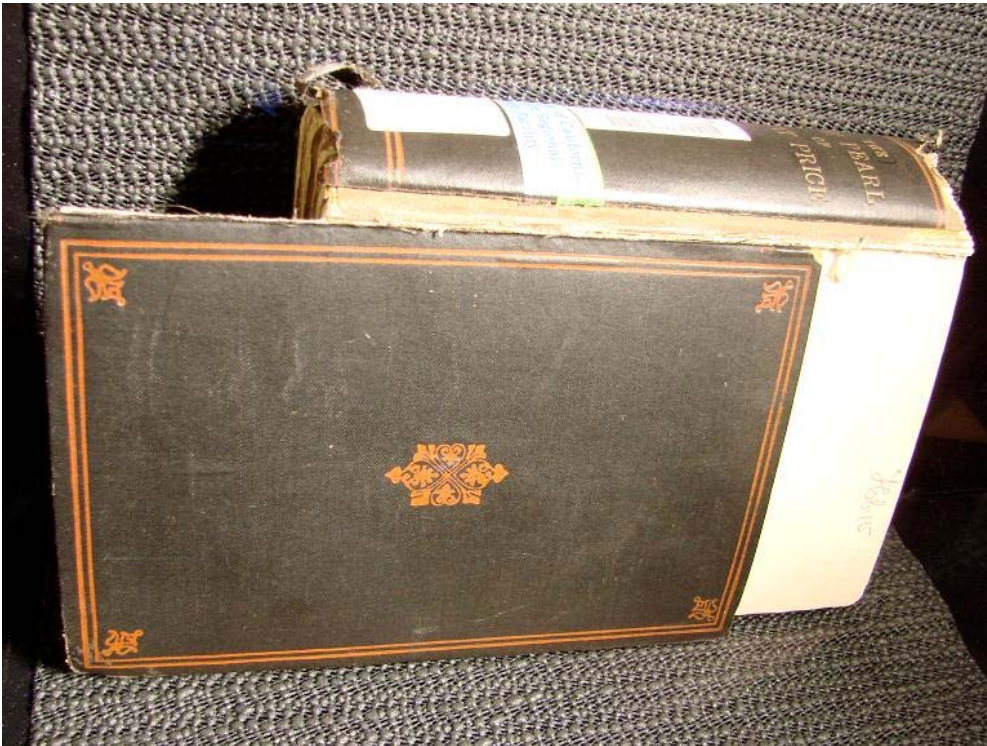
Brittle books- BI, DAM, MIS, PAG, UNC, SPH



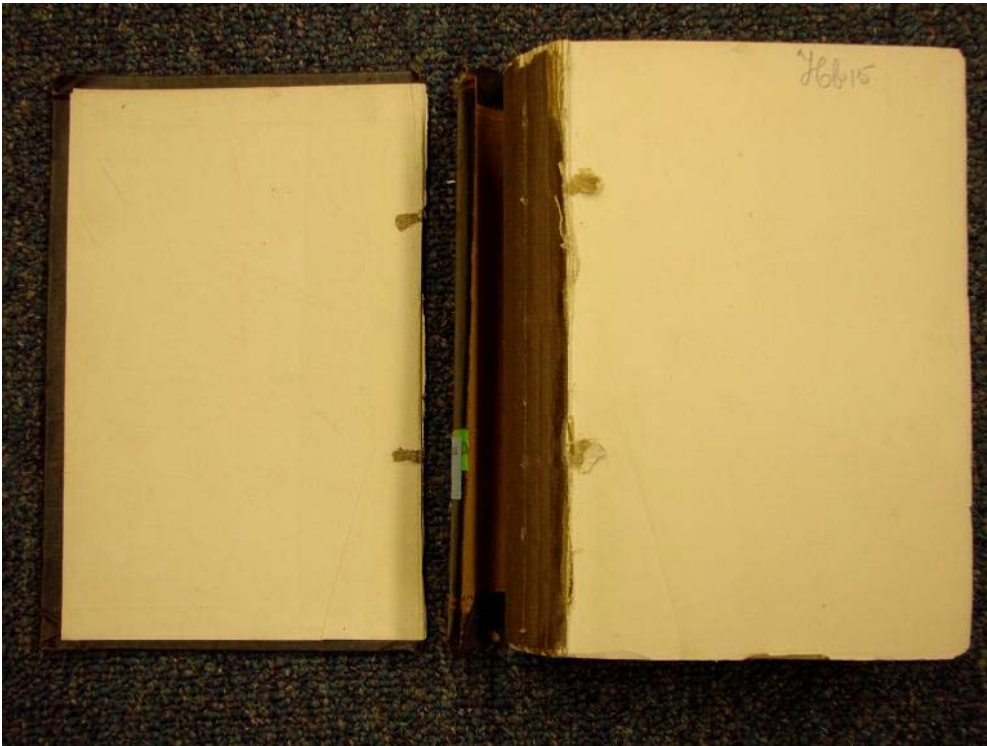
Brittle book- BI, DAM, SPH



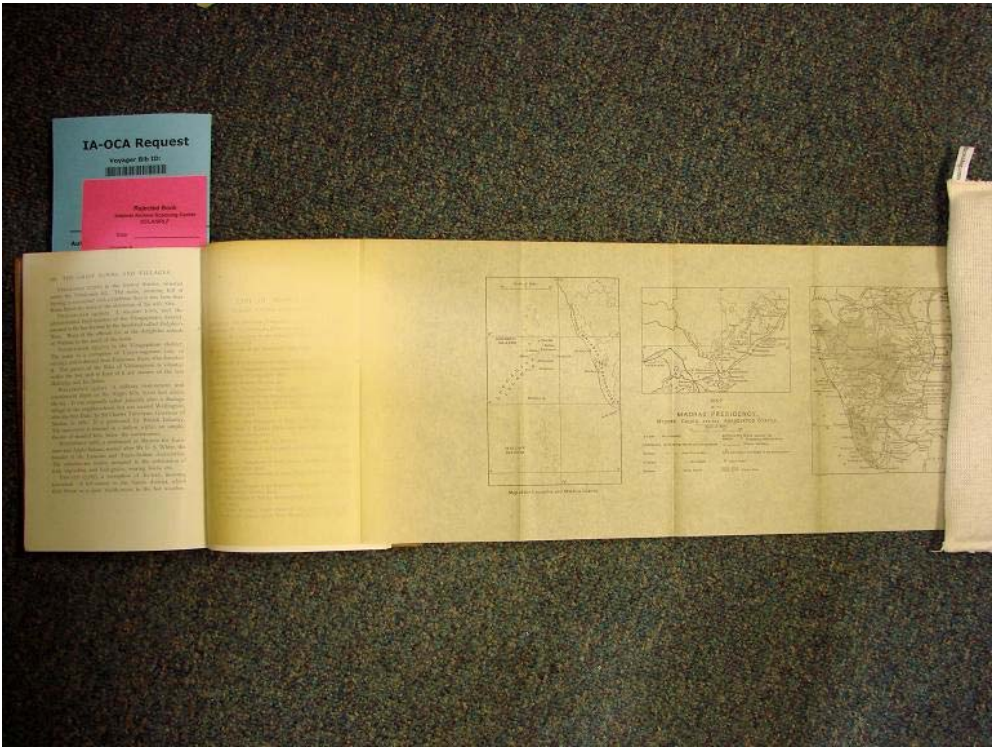
Brittle book- BI, DAM, SPH



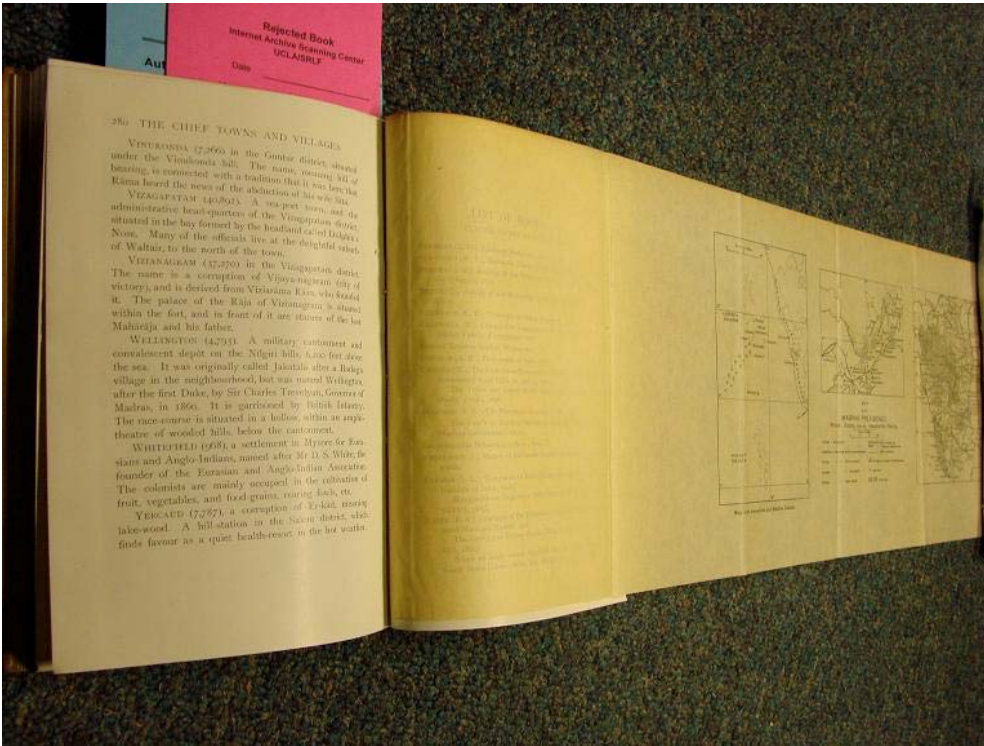
Fold out-FO



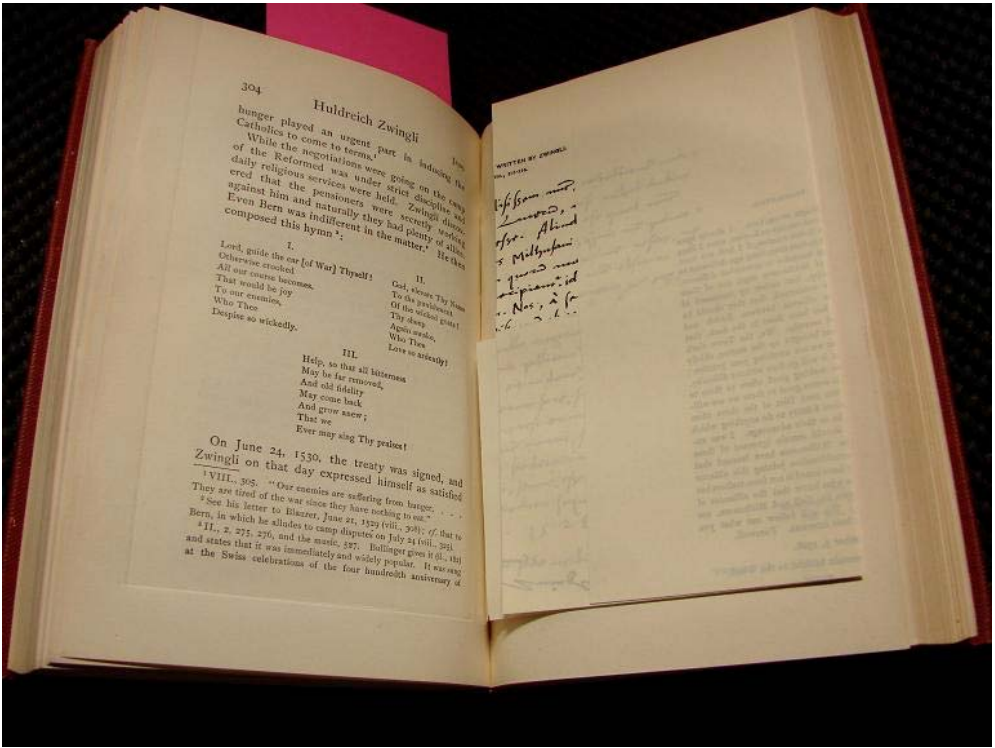
Fold out-FO



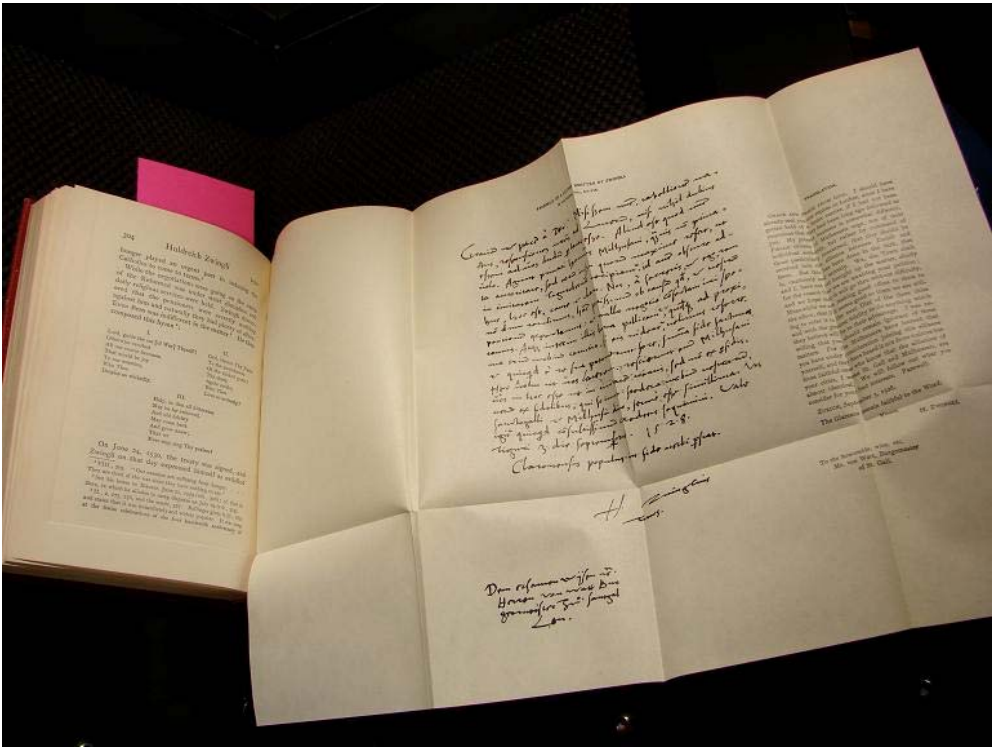
Fold out-FO



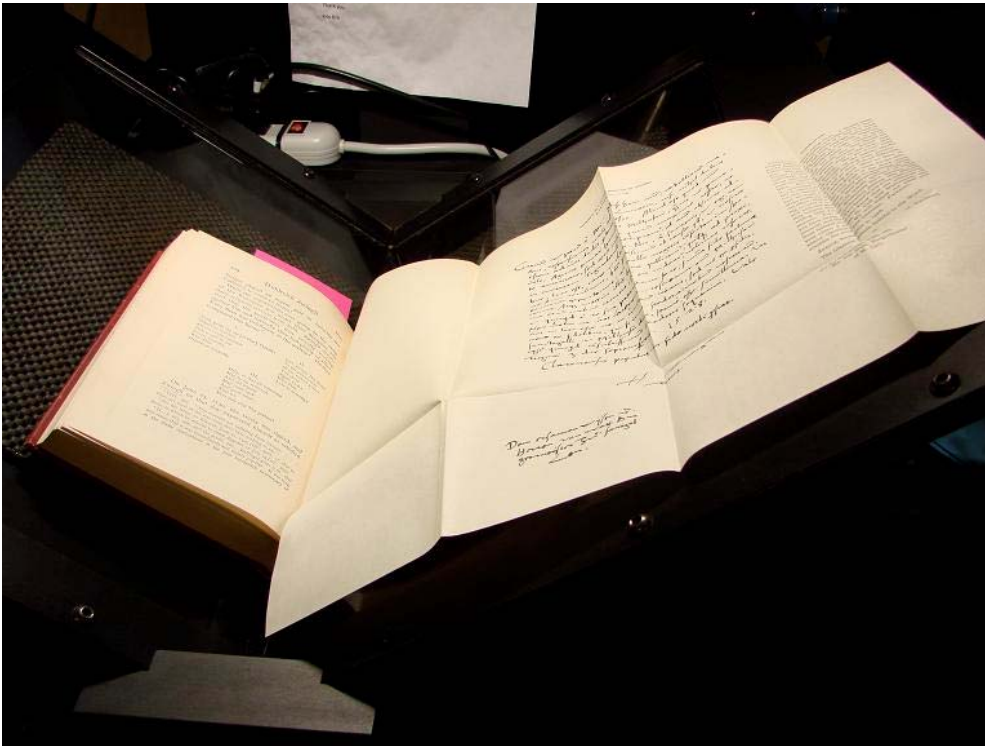
Fold out-FO



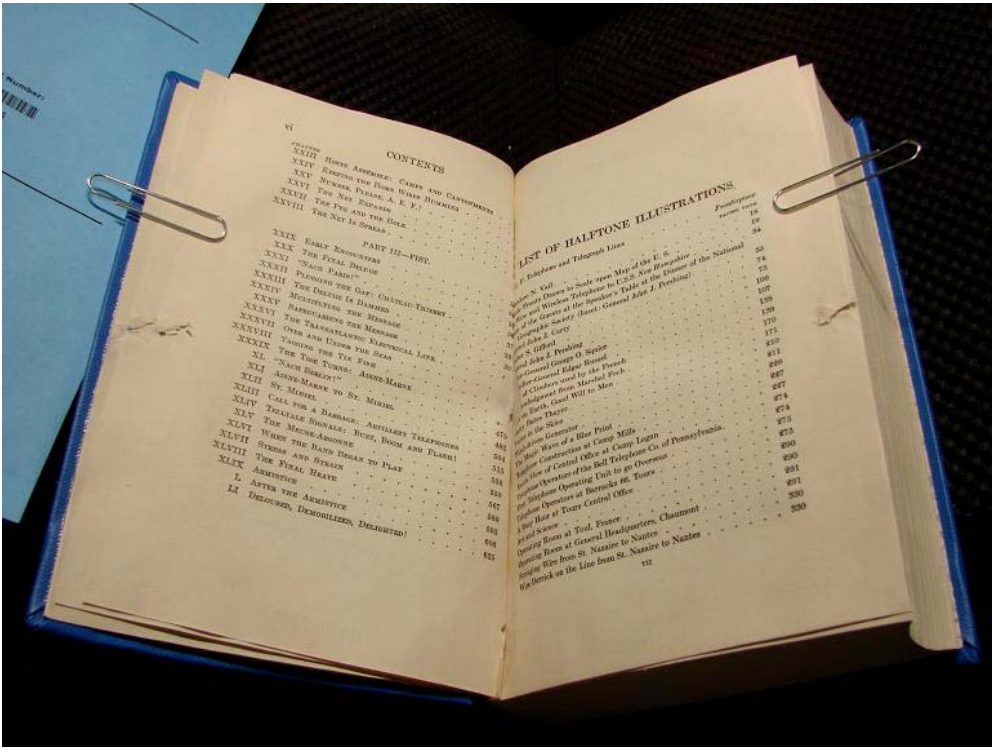
Fold out- FO



Fold out-FO

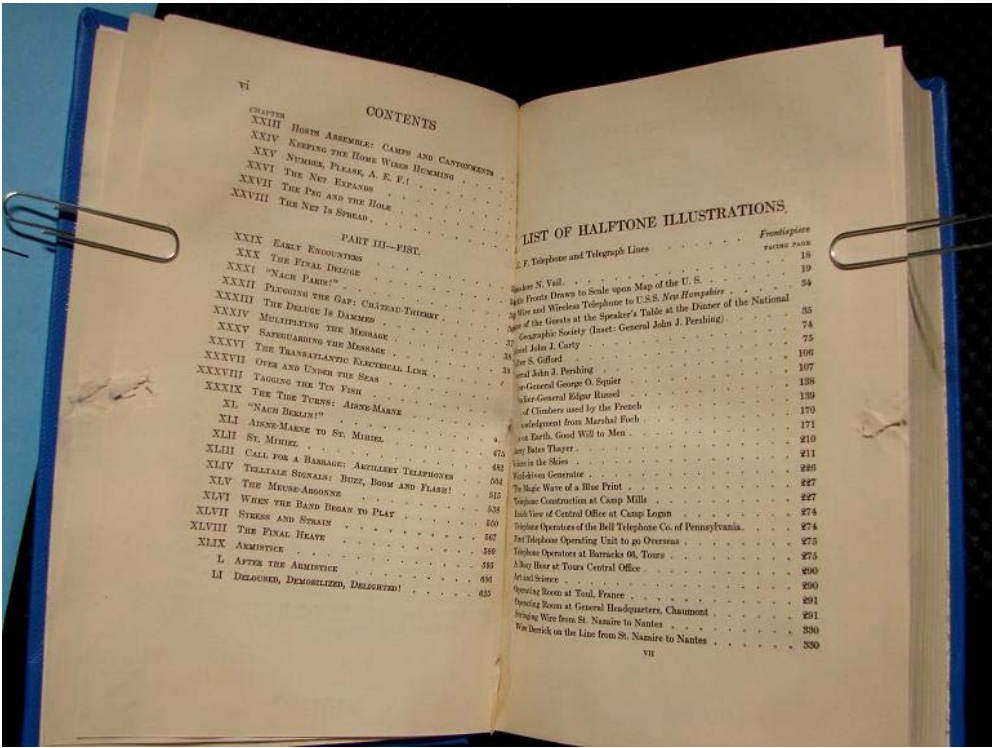


Glued pages- DAM, MIS, PA, PAG, SPH- replace photo with a new shot not using paperclips; paperclips were for photographic purposes only

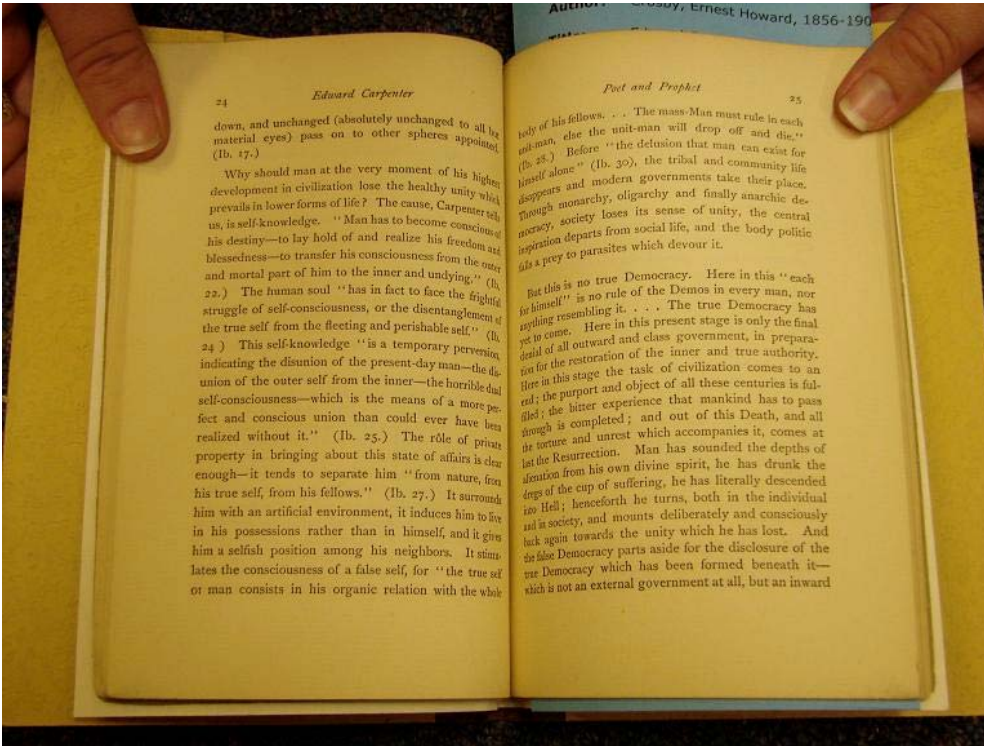


Glued pages- DAM, MIS, PA, PAG, SPH

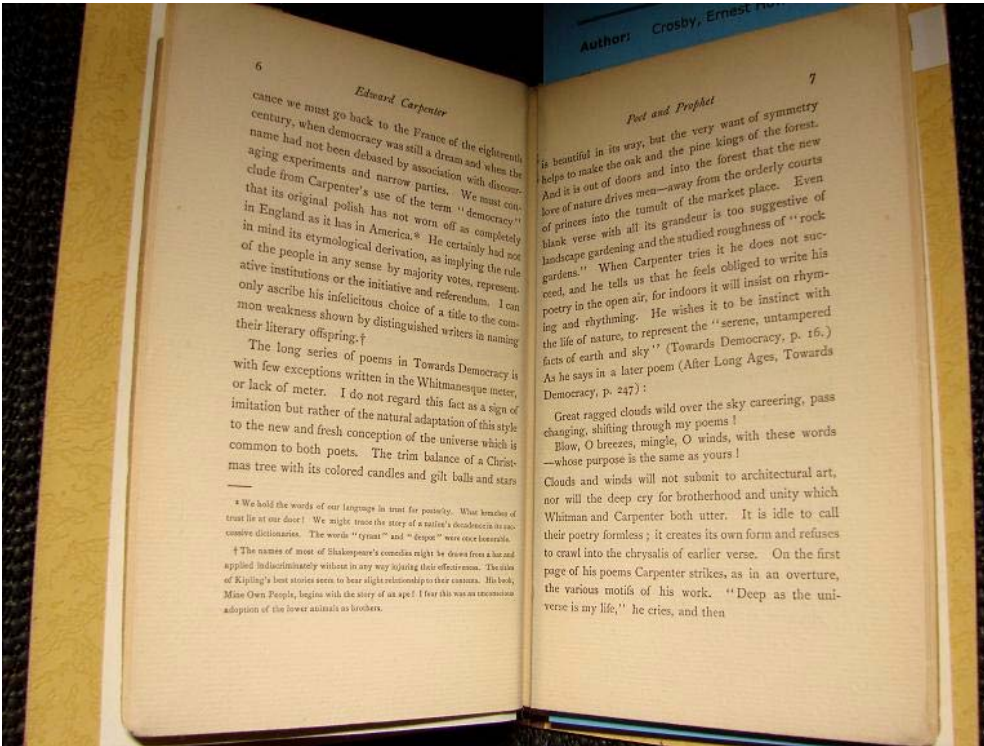
CONTENTS		
XIII	Hours Attendance, Calls and Correspondence	423
XIV	Records of the House of Representatives	424
XV	The New Zealand, A. E. F.	425
XVI	The New Zealand, A. E. F.	426
XVII	The New Zealand, A. E. F.	427
XVIII	The New Zealand, A. E. F.	428
XIX	The New Zealand, A. E. F.	429
PART III—FIRST.		
X	The Final Report	430
XI	The Final Report	431
XII	The Final Report	432
XIII	The Final Report	433
XIV	The Final Report	434
XV	The Final Report	435
XVI	The Final Report	436
XVII	The Final Report	437
XVIII	The Final Report	438
XIX	The Final Report	439
LIST OF HALF-TONE ILLUSTRATIONS		
1	Map of the U. S.	54
2	Map of the U. S.	55
3	Map of the U. S.	56
4	Map of the U. S.	57
5	Map of the U. S.	58
6	Map of the U. S.	59
7	Map of the U. S.	60
8	Map of the U. S.	61
9	Map of the U. S.	62
10	Map of the U. S.	63
11	Map of the U. S.	64
12	Map of the U. S.	65
13	Map of the U. S.	66
14	Map of the U. S.	67
15	Map of the U. S.	68
16	Map of the U. S.	69
17	Map of the U. S.	70
18	Map of the U. S.	71
19	Map of the U. S.	72
20	Map of the U. S.	73
21	Map of the U. S.	74
22	Map of the U. S.	75
23	Map of the U. S.	76
24	Map of the U. S.	77
25	Map of the U. S.	78
26	Map of the U. S.	79
27	Map of the U. S.	80
28	Map of the U. S.	81
29	Map of the U. S.	82
30	Map of the U. S.	83
31	Map of the U. S.	84
32	Map of the U. S.	85
33	Map of the U. S.	86
34	Map of the U. S.	87
35	Map of the U. S.	88
36	Map of the U. S.	89
37	Map of the U. S.	90
38	Map of the U. S.	91
39	Map of the U. S.	92
40	Map of the U. S.	93
41	Map of the U. S.	94
42	Map of the U. S.	95
43	Map of the U. S.	96
44	Map of the U. S.	97
45	Map of the U. S.	98
46	Map of the U. S.	99
47	Map of the U. S.	100
48	Map of the U. S.	101
49	Map of the U. S.	102
50	Map of the U. S.	103
51	Map of the U. S.	104
52	Map of the U. S.	105
53	Map of the U. S.	106
54	Map of the U. S.	107
55	Map of the U. S.	108
56	Map of the U. S.	109
57	Map of the U. S.	110
58	Map of the U. S.	111
59	Map of the U. S.	112
60	Map of the U. S.	113
61	Map of the U. S.	114
62	Map of the U. S.	115
63	Map of the U. S.	116
64	Map of the U. S.	117
65	Map of the U. S.	118
66	Map of the U. S.	119
67	Map of the U. S.	120
68	Map of the U. S.	121
69	Map of the U. S.	122
70	Map of the U. S.	123
71	Map of the U. S.	124
72	Map of the U. S.	125
73	Map of the U. S.	126
74	Map of the U. S.	127
75	Map of the U. S.	128
76	Map of the U. S.	129
77	Map of the U. S.	130
78	Map of the U. S.	131
79	Map of the U. S.	132
80	Map of the U. S.	133
81	Map of the U. S.	134
82	Map of the U. S.	135
83	Map of the U. S.	136
84	Map of the U. S.	137
85	Map of the U. S.	138
86	Map of the U. S.	139
87	Map of the U. S.	140
88	Map of the U. S.	141
89	Map of the U. S.	142
90	Map of the U. S.	143
91	Map of the U. S.	144
92	Map of the U. S.	145
93	Map of the U. S.	146
94	Map of the U. S.	147
95	Map of the U. S.	148
96	Map of the U. S.	149
97	Map of the U. S.	150
98	Map of the U. S.	151
99	Map of the U. S.	152
100	Map of the U. S.	153
101	Map of the U. S.	154
102	Map of the U. S.	155
103	Map of the U. S.	156
104	Map of the U. S.	157
105	Map of the U. S.	158
106	Map of the U. S.	159
107	Map of the U. S.	160
108	Map of the U. S.	161
109	Map of the U. S.	162
110	Map of the U. S.	163
111	Map of the U. S.	164
112	Map of the U. S.	165
113	Map of the U. S.	166
114	Map of the U. S.	167
115	Map of the U. S.	168
116	Map of the U. S.	169
117	Map of the U. S.	170
118	Map of the U. S.	171
119	Map of the U. S.	172
120	Map of the U. S.	173
121	Map of the U. S.	174
122	Map of the U. S.	175
123	Map of the U. S.	176
124	Map of the U. S.	177
125	Map of the U. S.	178
126	Map of the U. S.	179
127	Map of the U. S.	180
128	Map of the U. S.	181
129	Map of the U. S.	182
130	Map of the U. S.	183
131	Map of the U. S.	184
132	Map of the U. S.	185
133	Map of the U. S.	186
134	Map of the U. S.	187
135	Map of the U. S.	188
136	Map of the U. S.	189
137	Map of the U. S.	190
138	Map of the U. S.	191
139	Map of the U. S.	192
140	Map of the U. S.	193
141	Map of the U. S.	194
142	Map of the U. S.	195
143	Map of the U. S.	196
144	Map of the U. S.	197
145	Map of the U. S.	198
146	Map of the U. S.	199
147	Map of the U. S.	200
148	Map of the U. S.	201
149	Map of the U. S.	202
150	Map of the U. S.	203
151	Map of the U. S.	204
152	Map of the U. S.	205
153	Map of the U. S.	206
154	Map of the U. S.	207
155	Map of the U. S.	208
156	Map of the U. S.	209
157	Map of the U. S.	210
158	Map of the U. S.	211
159	Map of the U. S.	212
160	Map of the U. S.	213
161	Map of the U. S.	214
162	Map of the U. S.	215
163	Map of the U. S.	216
164	Map of the U. S.	217
165	Map of the U. S.	218
166	Map of the U. S.	219
167	Map of the U. S.	220
168	Map of the U. S.	221
169	Map of the U. S.	222
170	Map of the U. S.	223
171	Map of the U. S.	224
172	Map of the U. S.	225
173	Map of the U. S.	226
174	Map of the U. S.	227
175	Map of the U. S.	228
176	Map of the U. S.	229
177	Map of the U. S.	230
178	Map of the U. S.	231
179	Map of the U. S.	232
180	Map of the U. S.	233
181	Map of the U. S.	234
182	Map of the U. S.	235
183	Map of the U. S.	236
184	Map of the U. S.	237
185	Map of the U. S.	238
186	Map of the U. S.	239
187	Map of the U. S.	240
188	Map of the U. S.	241
189	Map of the U. S.	242
190	Map of the U. S.	243
191	Map of the U. S.	244
192	Map of the U. S.	245
193	Map of the U. S.	246
194	Map of the U. S.	247
195	Map of the U. S.	248
196	Map of the U. S.	249
197	Map of the U. S.	250
198	Map of the U. S.	251
199	Map of the U. S.	252
200	Map of the U. S.	253
201	Map of the U. S.	254
202	Map of the U. S.	255
203	Map of the U. S.	256
204	Map of the U. S.	257
205	Map of the U. S.	258
206	Map of the U. S.	259
207	Map of the U. S.	260
208	Map of the U. S.	261
209	Map of the U. S.	262
210	Map of the U. S.	263
211	Map of the U. S.	264
212	Map of the U. S.	265
213	Map of the U. S.	266
214	Map of the U. S.	267
215	Map of the U. S.	268
216	Map of the U. S.	269
217	Map of the U. S.	270
218	Map of the U. S.	271
219	Map of the U. S.	272
220	Map of the U. S.	273
221	Map of the U. S.	274
222	Map of the U. S.	275
223	Map of the U. S.	276
224	Map of the U. S.	277
225	Map of the U. S.	278
226	Map of the U. S.	279
227	Map of the U. S.	280
228	Map of the U. S.	281
229	Map of the U. S.	282
230	Map of the U. S.	283
231	Map of the U. S.	284
232	Map of the U. S.	285
233	Map of the U. S.	286
234	Map of the U. S.	287
235	Map of the U. S.	288
236	Map of the U. S.	289
237	Map of the U. S.	290
238	Map of the U. S.	291
239	Map of the U. S.	292
240	Map of the U. S.	293
241	Map of the U. S.	294
242	Map of the U. S.	295
243	Map of the U. S.	296
244	Map of the U. S.	297
245	Map of the U. S.	298
246	Map of the U. S.	299
247	Map of the U. S.	300
248	Map of the U. S.	301
249	Map of the U. S.	302
250	Map of the U. S.	303
251	Map of the U. S.	304
252	Map of the U. S.	305
253	Map of the U. S.	306
254	Map of the U. S.	307
255	Map of the U. S.	308
256	Map of the U. S.	309
257	Map of the U. S.	310
258	Map of the U. S.	311
259	Map of the U. S.	312
260	Map of the U. S.	313
261	Map of the U. S.	314
262	Map of the U. S.	315
263	Map of the U. S.	316
264	Map of the U. S.	317
265	Map of the U. S.	318
266	Map of the U. S.	319
267	Map of the U. S.	320
268	Map of the U. S.	321
269	Map of the U. S.	322
270	Map of the U. S.	323
271	Map of the U. S.	324
272	Map of the U. S.	325
273	Map of the U. S.	326
274	Map of the U. S.	327
275	Map of the U. S.	328
276	Map of the U. S.	329
277	Map of the U. S.	330
278	Map of the U. S.	331
279	Map of the U. S.	332
280	Map of the U. S.	333
281	Map of the U. S.	334
282	Map of the U. S.	335
283	Map of the U. S.	336
284	Map of the U. S.	337
285	Map of the U. S.	338
286	Map of the U. S.	339
287	Map of the U. S.	340
288	Map of the U. S.	341
289	Map of the U. S.	342
290	Map of the U. S.	343
291	Map of the U. S.	344
292	Map of the U. S.	345
293	Map of the U. S.	346
294	Map of the U. S.	347
295	Map of the U. S.	348
296	Map of the U. S.	349
297	Map of the U. S.	350
298	Map of the U. S.	351
299	Map of the U. S.	352
300	Map of the U. S.	353
301	Map of the U. S.	354
302	Map of the U. S.	355
303	Map of the U. S.	356
304	Map of the U. S.	357
305	Map of the U. S.	358
306	Map of the U. S.	359
307	Map of the U. S.	360
308	Map of the U. S.	361
309	Map of the U. S.	362
310	Map of the U. S.	363
311	Map of the U. S.	364
312	Map of the U. S.	365
313	Map of the U. S.	366
314	Map of the U. S.	367
315	Map of the U. S.	368
316	Map of the U. S.	369
317	Map of the U. S.	370
318	Map of the U. S.	371
319	Map of the U. S.	372
320	Map of the U. S.	373
321</		



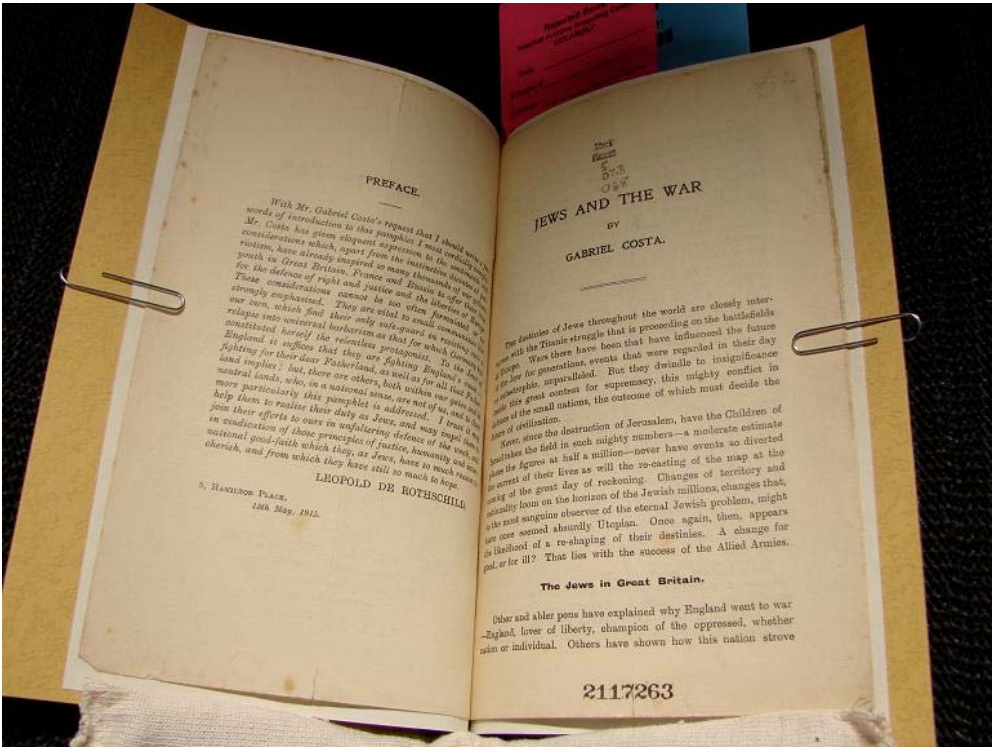
Tight Binding- MAR, PA, PAG, SPH



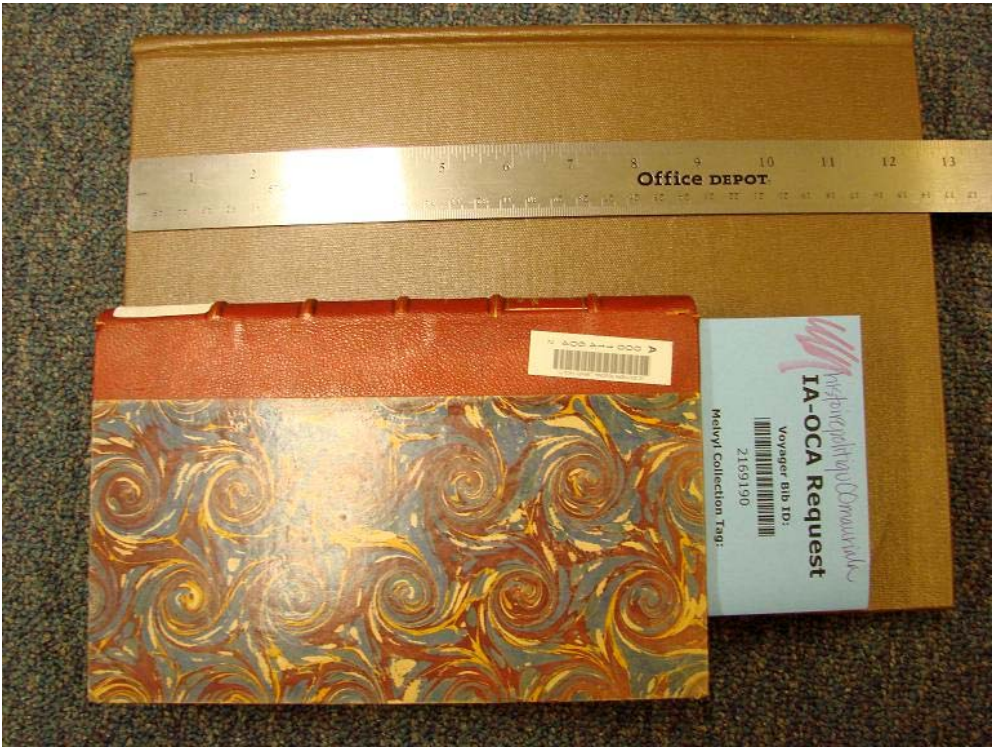
Tight Binding- MAR, MIS, PA, PAG, SPH



Tight Binding- MAR, MIS, PA, PAG, SPH- replace picture without paperclips (which are for photographics purposes only)



Size too big- LG, SPH



Size too big-LG, SPH



Book to thick- LGH, SPH



Book too thick-LGH, SPH

