

Scanning Techniques to Create Accessible PDF Documents

Many instructors provide their classes with materials in the form of scanned copies of journal or newspaper articles or book chapters, usually as PDF files. While this practice provides information that may be much more recent than textbook chapters, it does bring to the forefront the issue of how to scan these items to create accessible PDF files.

Accessible PDFs have been rendered into searchable text through the Optical Character Recognition (OCR) process available through Adobe Acrobat and some scanner software. Once the scanned image of text is recognized, it can be read aloud by screen reader technology, such as VoiceOver, NVDA, JAWS, etc. These screen reader applications are essential tools for people with various disabilities, including visual, learning, cognitive, or mobility disabilities.

Guidelines for Good Scans:

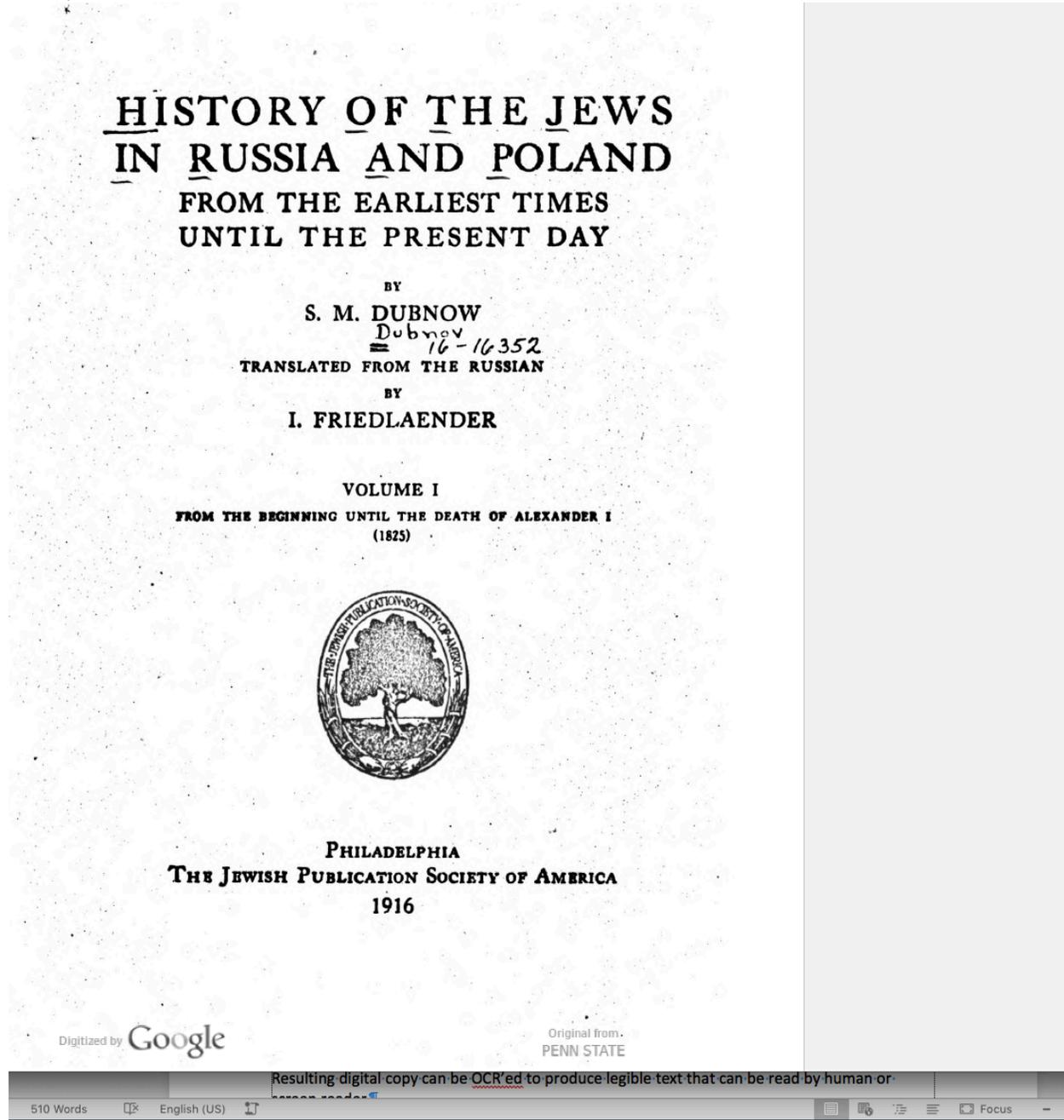
- Use “clean,” high quality source material with good contrast and without blurriness. Determine if an accessible, digital version of the material already exists. Using an accessible, electronic version of the text is preferable to scanning, which introduces a certain level of replication error.
- Choose source material that is an original printing of the text, not a scan that must be scanned another time. Scanning a scan creates another level of replication error, decreasing the crispness of the dots that create text and images and increasing the likelihood and intensity of printer or paper artifacts. To make text that can be read by screen readers, programs use Optical Character Recognition (OCR) to render a scanned image into text characters that can be recognized as words by the screen reader. Printer or paper artifacts can confuse the OCR process, either stopping it if too many errors are present or creating misreads.
- Review the source material to verify that the text is complete, and clean, and not cut off at gutter or edge. Scanned documents will not improve on the quality of the original. Skewed or crooked source material, highlighting, underlining, marginalia, or stains complicate the OCR process, leading to errors or inability to produce text that can be read by screen readers.
- The OCR process can also be hampered by highly textured or aged paper, which creates artifacts in the scanned copy.
- When scanning multiple pages from a bound source, press the book firmly against the scanner bed to prevent visible gutters and scan such that only one page of source material is copied per digital page. This technique can reduce or eliminate problematic artifacts introduced from the gutter formed at the spine of bound documents and the page’s curvature that can interfere with optical character recognition. Avoid shadows from subsequent or previous pages of bound source material.

- Some scanners allow users to adjust the scanning resolution. The optimal scanning resolution is usually 300-400 DPI.
- Scan in color (24-bit) ONLY if color is essential for the purpose of the document.
- If possible, save the scan as “Searchable PDF.” Name the file clearly and understandably, preferably with more than one identifier in the title, such as the course identifier, semester, author, etc.

Analyzing Scans

Problematic Scans

Example 1:



Friedlaender, I. (1916). *History of the Jews in Russia and Poland: From the Earliest Times Until the Present Day*. Philadelphia, PA: The Jewish Publication Society of America.

Scan Analysis: The scan above is from an old book published in 1916. The paper age or quality has caused the speckling artifacts in the scan. The page also has handwritten notes. The modern notes indicating that the page was digitized by Google and Originated from Penn State

may cause problems for the OCR because of the font color used, particularly if the scanner reads "Google" as colorized and the technician has not indicated that the scan should be processed in black and white.

Verdict: When this scan undergoes OCR, these artifacts will generate misreads, areas in which the OCR process is unable to correctly render the image into accessible text.

Correction: To create an accessible version of this text, the page should be rescanned from a cleaner copy of the book or retyped.

Example 2:

CONTENTS

11

CHAPTER	PAGE
IX. THE BEGINNINGS OF THE RUSSIAN RÉGIME	
1. <u>The Jewish Policy of Catherine II (1772-1796)</u> . . .	308
2. Jewish Legislative Schemes during the Reign of Paul I.	321
3. Dyerzhavin's "Opinion" on the Jewish Problem.	328
X. THE "ENLIGHTENED ABSOLUTISM" OF ALEXANDER I.	
1. "The Committee for the Amelioration of the Jews"	335
2. The "Jewish Constitution" of 1804.	342
3. The Projected Expulsion from the Villages.	345
4. The Patriotic Attitude of Russian Jewry during the War of 1812.	355
5. Economic and Agricultural Experiments.	359
XI. THE INNER LIFE OF RUSSIAN JEWRY DURING THE PERIOD OF "ENLIGHTENED ABSOLUTISM"	
1. Kahal Autonomy and City Government.	366
2. The Hasidic Schism and the Intervention of the Government	371
3. Rabbinism, Hasidism, and Enlightened "Berlin-erdom"	379
XII. THE LAST YEARS OF ALEXANDER I.	
1. "The Deputation of the Jewish People"	390
2. Christianizing Endeavors	396
3. "Judaizing" Sects in Russia.	401
4. Recrudescence of Anti-Jewish Legislation.	408
5. The Russian Revolutionaries and the Jews.	409

Friedlaender, I. (1916). *History of the Jews in Russia and Poland: From the Earliest Times Until the Present Day*. Philadelphia, PA: The Jewish Publication Society of America.

Scan Analysis: This scanned page also exhibits speckling caused by the age and texture of the source material, and it also has highlighting and underlining. An edge of the book's gutter is also evident in the vertical line at the lower left. Notice the blurriness of the text as well.

Verdict: The artifacts, highlighting, and blurriness of the text will cause OCR misreads.

Correction: To generate an accessible document, this text should be rescanned from a cleaner copy of the source material or retyped if one is not available.

Example 3:

Generated for Meadzer, Darcy C. (University of Oklahoma) on 2015-02-05 14:33 GMT / http://hdl.handle.net/2027/pst.000017378338
Public Domain, Google-digitized / http://www.hathitrust.org/access_use#pd-google

allan 3 v History 3/17/63

TRANSLATOR'S PREFACE

It is not my intention to expatiate in these prefatory remarks on the present work and its author. A history of the Jews in Russia and Poland from the pen of S. M. Dubnow needs neither justification nor recommendation. The want of a work of this kind has long been keenly felt by those interested in Jewish life or Jewish letters, never more keenly than to-day when the flare of the world conflagration has thrown into ghastly relief the tragic plight of the largest Jewry of the Diaspora. As for the author, his power of grasping and presenting the broad aspects of general Jewish history and his lifelong, painstaking labors in the particular field of Russian-Jewish history fit him in singular measure to cope with the task to which this work is dedicated.

In what follows I merely wish to render account of the English translation and of the form of the original which it has endeavored to reproduce.

The translation is based upon a work in Russian which was especially prepared by Mr. Dubnow for THE JEWISH PUBLICATION SOCIETY OF AMERICA. Those acquainted with modern Jewish literature in the Russian language know that the author of our book has treated the same subject in his general history of the Jewish people, in three volumes, and in a number of special studies published by him in the periodical *Yevreyskaya Starina* ("Jewish Antiquity"). Upon this

Digitized by Google

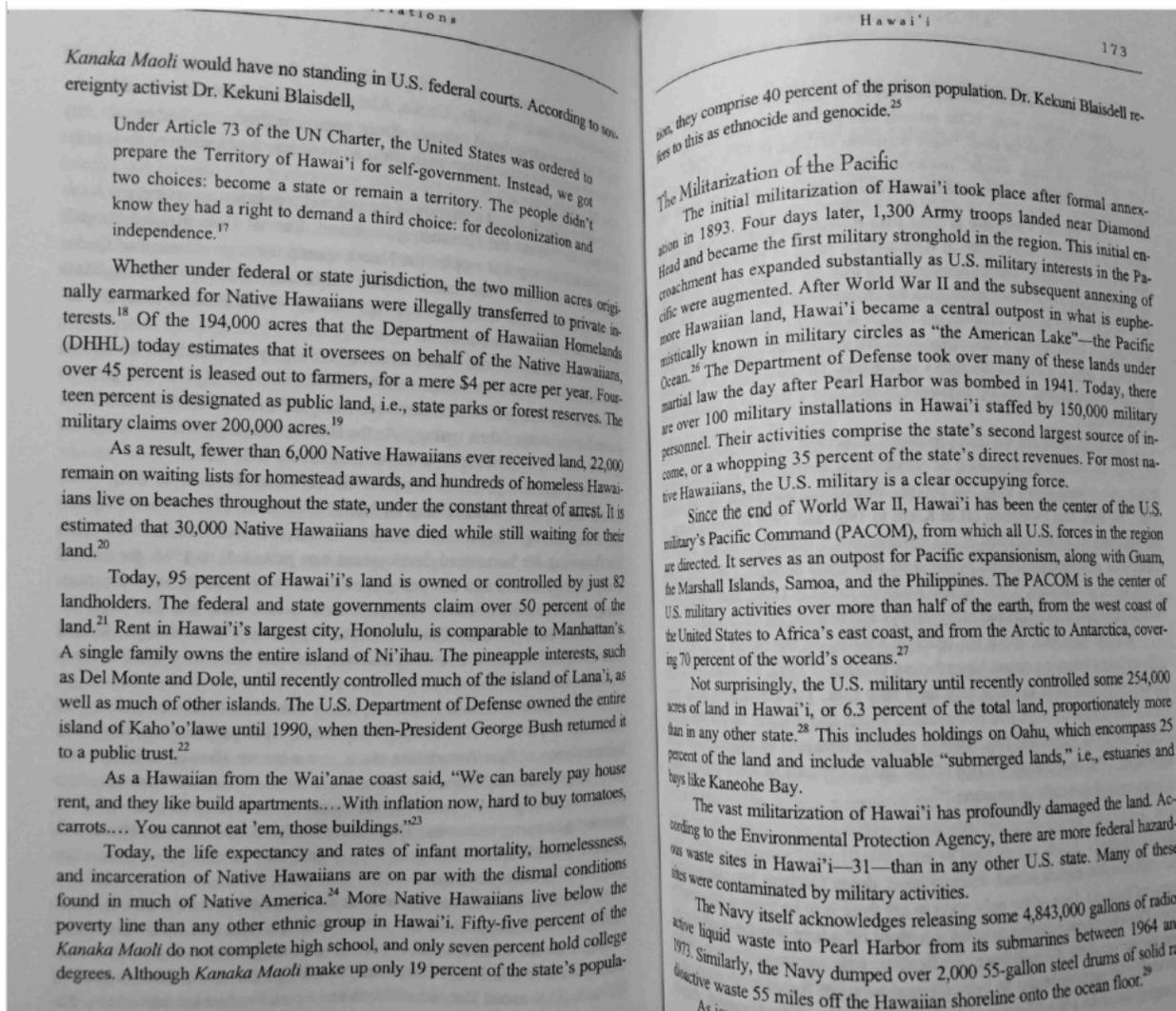
Original from
PENN STATE

Scan Analysis: While this page exhibits many of the same flaws as the previous example scans, it also has vertical text in a very light font color and vertical, blurry handwriting. Near the “Original from Penn State” notation is an area of strong black dots that the OCR will attempt to resolve into readable text.

Verdict: The OCR process on this page will result in misreads in multiple, disjointed areas that will need considerable manual rework.

Correction: Correction: To generate an accessible document, this text should be rescanned from a cleaner copy of the source material or retyped if one is not available.

Example 5:



LaDuke, Winona. (2017). *All Our Relations: Native Struggles for Land and Life*. Chicago, IL: Haymarket Books.

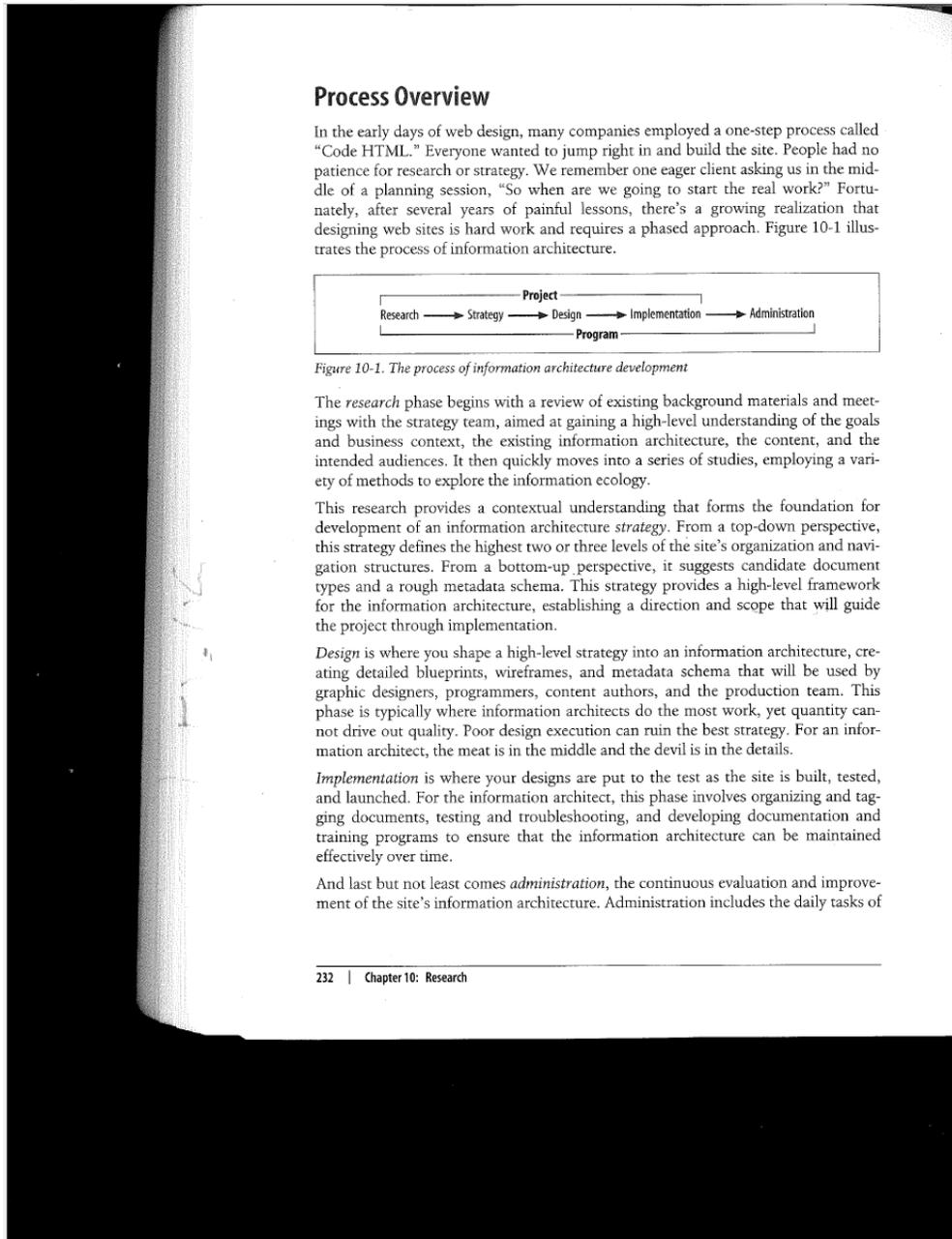
Scan Analysis: This scan was generated by using a scanning application from a mobile device. Two pages of the source material is visible on the single scan. The curvature of the book is very evident, skewing the text. Some of the text has been cut off on the lower right side, resulting in

partial words. The page is pigmented, contributing to the darkness of the scan. Shadows of text from the previous and subsequent pages show through.

Verdict: The OCR process on this page will result in misreads. Incomplete words will require manual editing.

Correction: Rescan the source material on a scanner, pressing down firmly to reduce the curvature of the book. Scan to create a single page of source material per page of scanned image. Increase the scanner's brightness setting to decrease the pigmentation of the page and shadows of previous and subsequent page text.

Example 1:



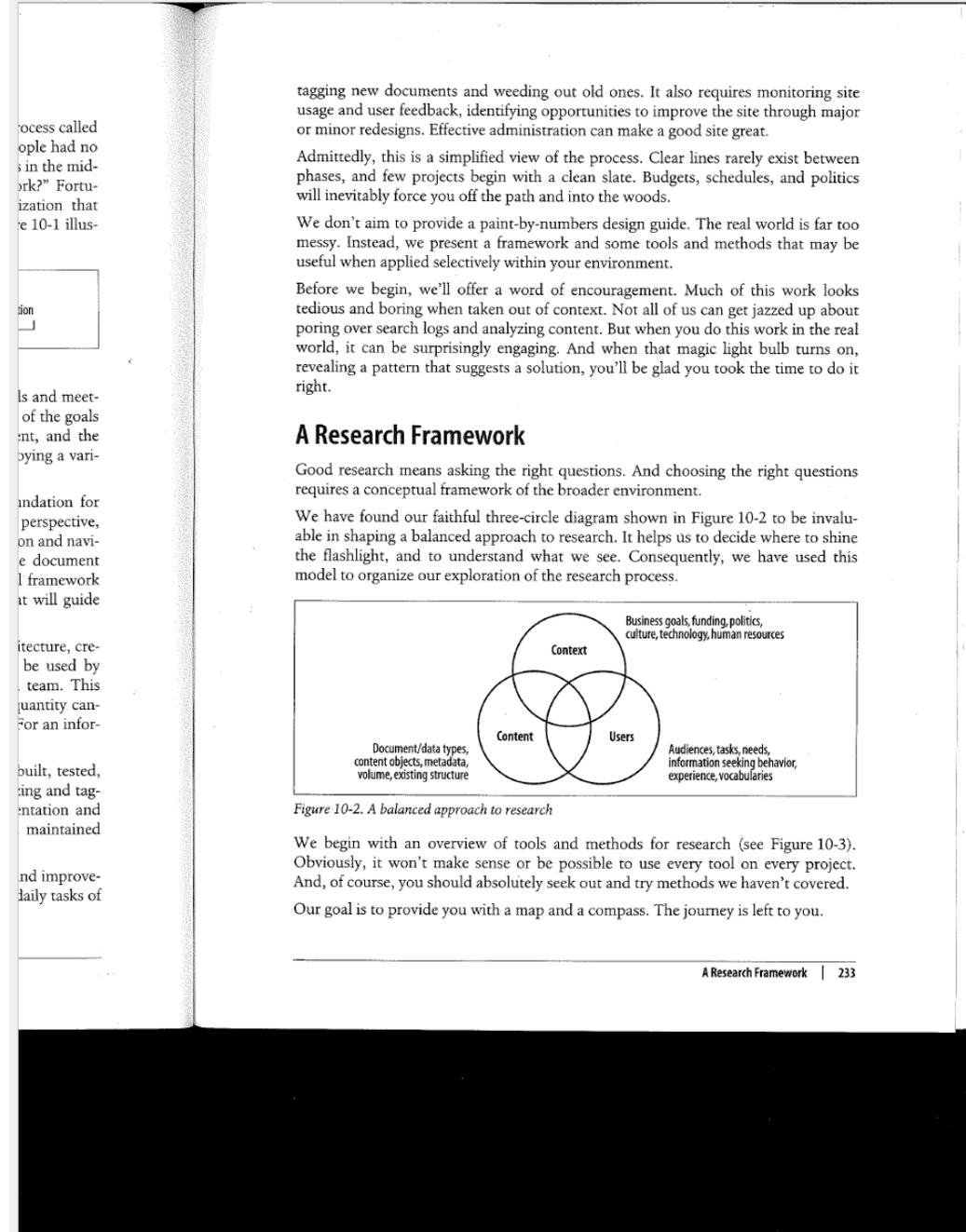
Morville, Peter, et. al., (2006). *Information Architecture for the World Wide Web: Designing Large-Scale Web Sites, 3rd Edition*. Sebastopol, CA: O’Reilly Media.

Scan Analysis: The scan has black banding from the difference between the book’s size and the scanner’s scanning area. Shadows of previous pages and some artifacts are visible along the left side. A shadow of the gutter between the left and right pages of the book can be seen along the right margin. The text itself is crisp and clean.

Verdict: Although the page is not perfectly scanned, it can be edited by Adobe Acrobat before the OCR process to remove the problematic aspects.

Correction: Adobe Acrobat's Crop Tool can extract the text from the banding and shadowing, resulting in a page that can be resolved through OCR into accessible text.

Example 2:



Scan Analysis: A partial page of the source material is visible along the page's left. Some shadows of subsequent pages can be seen on the right margin. A black band lies along the bottom of the page from a size mismatch between the book's size and the scanning area. The text, however is clean and free from artifacts.

Verdict: This scan can be corrected using Adobe Acrobat's tools.

Correction: Adobe Acrobat's Crop Tool can extract the text from the banding and shadowing, resulting in a page that can be resolved through OCR into accessible text.