Validating the Quality of Digitized Books for Long-Term Preservation

[Some] Findings on Digitization Error and Archival Quality



- Archival quality in an archival context
- Selective findings (more with poster and on website)
- Implications for practice and theory

Project Website: <u>http://hathitrust-quality.projects.si.umich.edu/</u>



45 People Are Involved

- Planning: Andrew W. Mellon Foundation [\$49,000]
- Research: Institute of Museum and Library Services [\$674,722]

Research Team (Michigan)

- Paul Conway, School of Information, PI
- Ed Rothman, LSA Statistics, co-PI
- John Wilkin, HathiTrust (5%)
- Leighann Ayers, oversight and space (5%)
- Jeremy York, liaison (10%)
- Jackie Bronicki, MLibrary, coordination (100%)
- Ken Guire, CSCAR, statistician (10%)
- Ryan Rotter, MLibrary, system design (75%)
- Stacy Maat, SI, physical review coordinator (25%)
- Melissa Chalmers, SI, user research (50%)
- Eugene Malin & Nina Elias, SI, website (10%)
- Sarah Jones, qualitative coder
- Jennifer Wright, qualitative coder
- Cherie Edmonds, qualitative coder
- Sarah Helm, qualitative coder
- > Jenny Vainberg, qualitative coder
- Nadia Sion, qualitative coder
- Naomi Scheinemann, qualitative coder
- Siqi Wei, qualitative coder

Research Team (Minnesota)

- ▶ John Butler, leadership
- ▶ Jason Roy, liaison
- > Ahnna Mahoney, coordination
- Kelly Frosch, qualitative coder
- Megan Scherer, qualitative coder
- Rui Zhao, qualitative coder
- Christina Graber, qualitative coder

Advisory Board

- ✓ Roy Tennant, OCLC
- ✓ Oya Rieger, Cornell University
- ✓ Ed Van Gemert, Wisconsin
- ✓ Robin Dale, LYRASIS
- ✓ Besiki Stvilia, Florida State Univ.

13 School of Information Volunteers

Julia Corrin, Natalie Bond, Graham Hukill, Jacqueline DiOrio, Janice Wong, Laura Brubacher, Brendan Coates, Laura Andrews, Holly Little, Monique Lowe, Molly Des Jardin, Amelia Lowry, David Fulmer



Archival Quality - A Value Proposition

- Archival nature
 - 1939 on : distinguishing characteristics of archives
 - > 2000 on: significant properties of digital objects
- Preservation
 - > 1961 on : media longevity [e.g., microfilm and acid-free paper]
 - 1985 on : processes to protect against loss [archival processes]
 - 1990 on: digitization image quality [archival master]
- Reliability [InterPARES]
 - ▶ 1995 on : completeness and process control

"... degree of completeness and degree of control of the procedure of creation are the **only** two factors that determine reliability of records." [Duranti 1995, p. 6]



Archival Quality in Archival Theory

Seamus Ross : digital libraries are digital archives

"A library ... is first of all an archive or repository in which society We easily observe that they may be fibraries by name, can find what it has already learned." [Kaplan 1964] but they are archives by nature." [Ross 2007, p. 8].

Terry Cook: context of creation

"This new paradigm for [has] a renewed focus on the context, purpose, intent, interrelationships, functionality, and accountability of the record, its creator, and its creation processes, wherever these occur." [Cook 1997, p. 48]

Geoffrey Yeo: boundary objects & persistent representation

"Records are "**persistent representations** of activities... or other **occurrents**... created by participants or observers of those occurrents or by their authorized **proxies**..." [Yeo 2008, p. 136]

INSTITUTE of

Archival Quality is Archival Science [Thomassen 2001]

- Object of Archival Science:
 - defining the nature of "process-bound information," which is "information itself and the processes that have generated and structured that information."
- Aims of Archival Science:
 - * "the establishment and maintenance of archival quality, that is to say: of the optimal visibility and durability of the records, the generating work processes, and their mutual bond."
- Methodology of Archival Science:
 - "maintaining the formal quality of process-bound information, by ensuring its availability, readability, completeness, relevance, representativeness, topicality, authenticity and reliability."

Thomassen, Theo. (2001). "A First Introduction to Archival Science," Archival Science 1: 373-385.



Research Environment - Digitization

- New preservation reality: from vertical integration to distributed management
 - Preservation programs used to exercise end-to-end control of reformatting
 - Now: preserving digitized content: "take what we can get"
 - End-user trust turns on validating "fitness-for-use"
- Two Research Questions for our project
 - What is quality? [definition, measurement, distribution]
 - What difference does lack of quality make for users? [barriers, acceptance testing]

Testbed: HathiTrust Digital Library

[http://www.hathitrust.org/]

Andrew W. Mellon Foundation [planning]

IMLS NLG [research, reporting]





SAA Research Forum -- August 2012

Rethinking **Quality** for Preservation and Access



SAA Research Forum -- August 2012

lbrary.

SERVICES

Research Workflow [2011-13]



SERVICES

9

Phase 1 [2011] - Metrics of Digitization Error

Level of Abstraction

LEVEL 1: DATA/INFORMATION

- 1.1 Text: thick text [fill, excessive]
- **1.2 Text: broken text [character breakup]**
- 1.3 Illustration: scanner effects [moiré, gridding]
- 1.4 Illustration: tone, brightness, contrast
- **1.5 Illustration: color imbalance, gradient shifts** LEVEL 2: ENTIRE PAGE
- 2.1 Blur [distortion]
- 2.2 Warp [text alignment]
- 2.3 Skew [page alignment]
- 2.4 Crop [gutter, text block]
- 2.5 Obscured/cleaned [portions not visible]
- 2.6 Colorization [text bleed, low contrast]
- LEVEL 3: WHOLE VOLUME
- 3.1 Fully obscured [foldouts or objects]
- 3.2 Missing pages [one or more]
- 3.3 Duplicate pages [one or more]
- 3.4 Order of pages
- 3.5 False pages [not part of Original Content]

Possible Cause of Error

Source or post-processing Source or post-processing Scanning or post-processing Scanning, post-processing or source Scanning, post-processing or source

Scanning or source

Post-processing Scanning, source or post-processing Source or post-processing Scanning or post-processing Source or post-processing

Scanning Source or scanning Source or scanning Source or scanning Scanning or post-processing



Phase 1 [2011] - Error Severity Scale

o - Error is **undetectable** on the page.

- 1 Error has negligible affect on Original Content.
- 2 Error alters appearance of Original Content.
- 3 Error has affects readability of Original Content.
- 4 Error requires significant inference to read Original Content.
- 5 Error renders Original Content undecipherable.



Error Type - Thick

Avour Pristan 4100 " afan'd idin , that a main main reason aroant offranti ananar "Historia way beliew of the manif " ungamous sweet survey and "at the moment of Westroying " start crimitant structure une sa-th 4. Administ Man and and Must " as the well maintainer eid That farme Now if this Coroner didesay "very rationally rinliges a there this, for which we have no more " every melanchely fit does not than an ewspaper authority mind, " deprive a man of the capacity ! I say that he delivered a doctrine "of discerning tight from wrong to completely: at variance with the " and therefore, if a real futurie . Law of the Land, and that he "kill himself in a facid interval; was 2 guilty of a breach of his "he is a self-murderer as much duty." The law adopts no such " as another man." azion Braczerowe, in his Let the public judge, then, of Fourth Book, and 14th Chap- the manner in which this Coroner ter, after calling suicide pretend- performed his duty upon this oced hereism; but real cowardice, casion. You see; even if a nopreceds to say, that the Law of torious lanatie, a man who has Bighand has ranked this amongst been a lunatic for years, kill him. the highest crimes, making it a self in a lucid interval, the law peculiar species of felony. Then sends his body to be buried in the he gees on thus : "The party highway with a stake driven "must be in his senses, else it through it, and makes his goods "the no crime. But this excuse and chattels forfest to the King ; "ought not to be strained to that " boping," says Bilkckstone, " that "length? to which our Coroners " his care for either his own re-"Furies are upt to carry it ; " putation, or the welfare of his . "manhely, that the very act of "family, would be some mative "dicide in an evidence of vin- "to restrain him from so desperate "statity; an if every man who "and wicked an act" But what "able contrary to reason had no is there to restrain any man, if

Original from

Error Type - Broken

the of the war me Government has and - the the terrie are to keep at even with in and the court net must get some better ter ber bei ber ber ber ber Englishmen is here here ere cover to Shamese schools, and, of in som men der bereichte beite Betre subjects is a le source de seu d'a mung Englishman is the second relation while be as much as HAS NO REPORTED AND WE KNOW IT THE COUNTY. so we see to use traing. where, and arithmin is a largelagest and while encourage them the state and an instant topad that the so any ter see to part and became that keen in the light and cartary fine assess? Mr. ing of the contract we begin in the schools of Sign. and which is a structure match which he once got is in massion of other than and the made on of Hindas. and a The reacts has hered granter in their and a second second the substance while managementation is the management of the second second

Error Type - Warp

the only 1st class passenger on board. She ought to January 19th. Still blowing had. It being Sun-tow the common of the control of the being Sunday, the crews of the foreign men-of-war have low on shore, and the amount of drunken men, and row in the streets, could not have been surpased by In the streets, could not have been surpassed by Portsmouth itself in the good old days, when it seemed to be Jack's first duty to let the inhabitants realize that he was on shore January 20th, Coaled and watered. There is no Place that I have yet visited where coaling and vatering requires to the state of t Vatering requires more attention. The prices must De fixed beforehand, of course. Then, if after all you Set your full weight of coal or full quantity of water, you will be harden. The harden give Some good advice about this on the consul-General can give this of the consul-General can give the consul-General can ome good advice about this. The one thing to void, is to enter into an one with any of The rascals who come on board to rout for custom The rascals who come on board to tout for custom areas as him washing to the second scheme of through a ere. The washing you send ashore, if through a roker, is 25. 6d. a dogen b is cont to a man whose , etc. Ine washing you send ashore, if througu a roker, is 3s. 6d. a dozen; if sent to a man whose a ≥une I have left at the Consent Canaral's office, it is Anne I have left at the Consul General's office, it is a dozen. The same undertakes the A state 1 have left at the Consul General's office, it is a dozen. The same man undertakes the asanng in both cases. Thursday, 21st, I called Masher (late Grand Viziar) on His Highness Kiamil Ura io a verv interesting Asher (late Grand Vizier). He is a very interesting most kind and communicative, and An. He was most kind and communicative, and duent Frontiat On family 2 and an, He was most kind and communicative, and eaks the most fluent English. On January 22nd, returned the call at the Consul e next day, he returned the call at the Consul eneral's house. His view of the at the Consul The next day, he returned the call at the Consum coneral's house. His view of the Armenian difficulty arm not at liberty to publick heat I may say that it Am not at liberty to publish, but I may say that it m not at liberty to Publish, but I may say that it majority of Europeane Cimmetances which have fers to to carlo from that which is commonly held by majority of Europeans. Circumstances which have t12 Original from Googi PRINCETON UNIVERSITY

Error Type – Crop







Phase 2 [2011-12] - Error Detection, Coding, Analysis

- Random Samples of Digital Volumes
 - Populations: pre-1923 N= 1.3 million post-1923 N=6.5 million
 - Samples: 1,000 volumes per study
 - 1. Google (pre-1923, English)
 - 2. Google (post-1923, English, no serials)
 - 3. Internet Archive (pre-1923, English)
 - 4. non-Roman scripts (250 volumes in four alphabets)
 - Systematic sampling strategy within each volume
 - Up to 100 pages per volume, evenly distributed front to back
 - Up to 25% of a volume, evenly distributed
- Coding: page-level, whole volume, physical inspection
 - Coding of 456,217 page-images [for 11 errors]
 - Double coding of 10% of each sample [ca. 45,000 page-images]
 - Coding of 2,000 whole volumes [for volume level errors]
 - Coding of 1,490 physical volumes for book/bib. characteristics



Findings: Comparison of Most Frequent Errors

Total coding of 178,297 page-images digitized by Google

	Sever	ity = 0	Severity = 1		Severi	ty = 4	Severity = 5		
	<<<< 19	23 >>>>	<<<< 192	23 >>>>	<<<< 192	23 >>>>	<<<< 192	23 >>>>	
Text									
Thick	62.04%	67.52%	25.66%	21.00%	0.19%	0.40%	0.11%	0.42%	
Broken	61.00%	73.37%	29.96%	19.18%	0.19%	0.41%	0.25%	0.36%	
Page									
Crop	99.37%	98.85%	0.27%	7.05%	0.02%	0.04%	0.15%	0.25%	
Warp	29.22%	45.78%	60.18%	48.93%	0.04%	0.04%	0.05%	0.06%	
Obscure	16.88%	56.83%	78.05%	41.69%	0.08%	0.02%	0.46%	0.16%	
			182,205		490		972		
Portion of Total Error (pre-1923)		96.9%		82.5%		87.9%			
				113,682		795		1,077	
Portion of Tota	l Error (post	-1923)		90.5%		87.9%		86.5%	



SAA Research Forum -- August 2012

Findings: Comparison of Most Frequent Errors

Total coding of 85,535 page-images digitized by Internet Archive

	Severity = 0	Severity = 1	Severity = 4	Severity = 5
	<<<< 1923	<<<< 1923	<<<< 1923	<<<< 1923
Text				
Thick	93.23%	4.12%	0.01%	0.00%
Broken	81.71%	11.77%	0.19%	0.10%
Illustration				
Tone	69.00%	24.65%	0.17%	0.01%
Page				
Crop	99.61%	0.22%	0.02%	0.07%
Warp	41.13%	57.00%	0.00%	-
Obscure	56.93%	39.61%	0.02%	0.07%
Colorization	47.56%	45.22%	0.00%	0.03%
Skew	90.24%	9.33%	-	-
Blur	94.22%	4.28%	0.04%	0.03%
		95,293	199	206
Portion of "Big Five" Errors		57.40%	52.00%	77.40%



Findings: Distribution of Severe Error

- Proportion of volumes with severe error
- Level 4 or 5 severity in any error type on any page-image

Pages w/ Severe Error		Nu	mber of Volu	mes		Cumulative Percent				
		Go	Google		et /	Archive	Google			
<<<< 1923 >>>>		<<<< 19 23 >>>>		<<<< 1923		<<<< 1923	<<<< 1923 >>>>			
0	0	555	637	876		93.19%	59.55%	69.16%		
1	1	167	131	43		97.99%	77.47%	83.39%		
2	2	76	50	7		98.51%	85.62%	88.82%		
3	3	39	29	2		98.72%	89.81%	91.97%		
4	4	24	11	0		98.72%	92.38%	93.16%		
5	5	12	8	3		99.04%	93.67%	94.03%		
6	6	12	6	3		99.36%	94.96%	94.68%		
7	7	8	3	0		99.36%	95.82%	95.01%		
8	8	6	3	2		99.57%	96.46%	95.33%		
9	9	1	2	0		99.57%	96.57%	95.55%		
10	10	4	3	0		99.57%	97.00%	95.87%		
11 to 21	11 to 28	20	28	1		99.68%	99.10%	98.97%		
22 to 68	38 to 168	8	10	3		100.00%	100.00%	100.00%		
		932	921	940						



Phase 1 [2011] - Metrics of Digitization Error

Level of Abstraction

LEVEL 1: DATA/INFORMATION

- 1.1 Text: thick text [fill, excessive]
- 1.2 Text: broken text [character breakup]
- 1.3 Illustration: scanner effects [moiré, gridding]
- 1.4 Illustration: tone, brightness, contrast
- **1.5 Illustration: color imbalance, gradient shifts** LEVEL 2: ENTIRE PAGE
- 2.1 Blur [distortion]
- 2.2 Warp [text alignment]
- 2.3 Skew [page alignment]
- 2.4 Crop [gutter, text block]
- 2.5 Obscured/cleaned [portions not visible]
- 2.6 Colorization [text bleed, low contrast]
- LEVEL 3: WHOLE VOLUME
- 3.1 Fully obscured [foldouts or objects]
- 3.2 Missing pages [one or more]
- 3.3 Duplicate pages [one or more]
- 3.4 Order of pages
- 3.5 False pages [not part of Original Content]

Possible Cause of Error

Source or post-processing Source or post-processing Scanning or post-processing Scanning, post-processing or source Scanning, post-processing or source

Scanning or source

Post-processing Scanning, source or post-processing Source or post-processing Scanning or post-processing Source or post-processing

Scanning Source or scanning Source or scanning Source or scanning Scanning or post-processing



Whole Book Errors - Preliminary Findings [one sample]

- Average number of pages per volume with whole book error.
- Certainty of loss of Original Content on some part of a page.

	Variable	n	Mean	Std. Dev.	Minimum	Maximum
	Pages per Volume	997	397.49	272.75	8	1628
Mean = 4.93	Whole Book Error					
pages/volume	Obscured Content	997	3.36	20.82	0	366
	Missing Page(s)	997	0.67	6.529	0	155
	Duplicate Page(s)	997	0.62	4.48	0	92
	Out of Order	997	0.24	2.185	0	43
	False Page (s)	997	0.04	0.343	0	8
Mean = 5.56 pages/volume	Page-level Quality Error	s in a V	olume			
	Sure Loss 1/3 page	996	2.44	11.69	0	177
	Sure Loss 2/3 page	996	0.72	4.22	0	68
	Sure Loss all page	996	0.58	0.046	0	59
	Unsure Loss 1/3 page	996	1.51	7.94	0	156
	Unsure Loss 2/3 page	996	0.11	0.603	0	9
	Unsure Loss all page	996	0.206	1.49	0	28



20

Impact of Physical Characteristics on Error

Number of Pages per Volume with at least 1 severe error

Characteristic	Characteristic of Volume			≥ 7 Pages	Chi-Square (p value)	
	:	< 1860	70	19	<0.0001	
Publication Year	:	1866-1899	70	11		
	:	≥ 1900	110	7		
Cuttor Margin	:	More than 1.0 cm	184	20	0.002	
Gutter Margin	:	Less than 1.0 cm	66	17		
Craphic Contant	:	No	153	20	0.004	
Graphic Content	:	Yes	97	17		
Dorts of Dagos Missing	:	No	241	37	0.05	
Parts of Pages Missing	:	Yes	9	0	0.06	
	:	Fully Intact	190	27		
Binding	:	Loose	41	4	0.05	
	:	Not Intact	12	6	0.06	
		Missing All/Part	3	0		

Next step in analysis: map physical characteristics to specific error types.



SAA Research Forum -- August 2012

More Data Gathering - 2012

- Special focus on graphics and illustrations
 - Compile examples of digitization error
 - Get diagnosis from panel of imaging scientists
 - Explore options for correcting error

- Digitization processes documentation and analysis
 - Google digitization/post-processing techniques
 - Internet Archive post-processing techniques
 - HathiTrust ingest and re-ingest routines
 - Costs and limitations of manual review of digitization error



Study Phase 3 [2012-13] - Two Major User Studies

- Reading online Error threshold for user rejection
 - Concepts: Text legibility; illustration interpretability
 - **Scholarship:** IQ (intrinsic); Relevance clues (object); Readability
 - Population: Digital humanities scholars who use books as primary sources
 - Goal: Identify thresholds of acceptability (limbo bar)
- Managing library print collections
 - **Concepts:** Low cumulative error; completeness; redundancy
 - Scholarship: "Last copy" criteria and policy
 - Population: HathiTrust members: collection development and preservation librarians
 - Goal: Certify individual volumes as "fit for use" (high bar)



Study Phase 3 [2012] - Process Study

Predicting error from text to/from page-image

- Spatial mapping of error text landscapes [training set] with hOCR text file [Google, IA, JSTOR]
- > Partnership with SI Professor Qiaozhu Mei's research team

the Beck case, the evidence on which he was convicted has become discredited to a point at which no jury would maintain its verdict of guilty. The reluctance is not to confess that an innocent man is being punished, but to proclaim that a guilty man has escaped. For if escape is possible deterrence shrinks almost to nothing. There is not better established rule of criminology than that it is not the severity of punishment that deters, but its certainty. And the flaw in the case of Terrorism is that it is impossible to obtain enough certainty to deter. The police are compelled to confess every year, when they publish their statistics, that against the list of crimes reported to them they can set only a percentage of detections and convictions. And the list of reported crimes can form only a percentage, how large or small it is impossible to say, but probably small, of the crimes actually committed; for it is the greatest mistake to suppose that everyone who is robbed runs to the police: on the contrary, only foolish and ignorant or very angry people do so without very serious consideration and great rely In most cases it costs nothing and a good deal to prosec in Heartbreak House, w



INSTITUTE of

SERVICES

SAA Research Forum -- August 2012

Deliverables [2013]

- Report findings on website
 - Tables, analysis, links to data
 - http://hathitrust-quality.projects.si.umich.edu/
- Publish peer-reviewed articles & proceedings
 - > American Archivist, Archival Science, JASIST, IJDL, IJIQ
 - JCDL, iPres, IS&T Archiving
 - College & Research Libraries, First Monday
- Distribute Quality Review web-application(s)
 - Three tools that use sampling strategies
 - One tool for volume-by-volume certification



Summary

• What is quality?

- Absence of page-image error relative to expected uses.
- > Presence of intrinsic character sufficient to inspire trust.
- "Fit for purpose" exploring the limits of "one size fits all"
- How bad is it?
 - Very low incidence of very severe error?
 - Likely findable with automated processing
 - High incidence of low-severity text error (Google)
 - Very low incidence of whole volume error
 - Unlikely findable with machine processing algorithms
 - Very high incidence (likely) of scanner effects on book illustrations
- Why does error occur?
 - Physical book characteristics have little or no impact
 - Faith of digitizers in post-scan image processing at scale



Implications for Practice

- Lowering the bar on image quality is not necessarily an ethical or professional compromise
- New tools and techniques for measuring quality will emerge from this study
- Communicating error to users is important
- Need for automated quality validation routines
 - Error models as first steps toward machine processing
 - Distinguishing errors that matter from those that don't
- Proposition: Certification of trustworthy repositories must encompass the qualities of the content within.



Implications for Archival Theory

- An archival principal [archival quality] can be described empirically.
 - Scoping the "intrinsic value" of copies [Boon 2010]
- Reaffirm value of digital surrogates as preservable products
 - Preservation trumps access as a compelling archival rationale
- Establish the archival nature of digitized surrogates
 - "Archivalness" derives from creation processes [reliability]
 - Provenance derives in part from digital curation
 - Appraisal of value through assessment of use



References

- Boon, Marcus. (2010). *In Praise of Copying*. Cambridge: Harvard University Press.
- Conway, Paul. (2011). "Archival Quality and Long-term Preservation: A Research Framework for Validating the Usefulness of Digital Surrogates." *Archival Science* 11 (3).
- Cook, Terry. (1995). "Electronic Records, Paper Minds: The revolution in information management and archives in the post-custodial and post-modernist era." *Archives and Manuscripts* 22 (2): 300-328.
- Cook, Terry. (1997). "What is Past is Prologue: A History of Archival Ideas Since 1898, and the Future Paradigm Shift." *Archivaria* 43 (Spring): 17-63.
- Duranti, Luciana. (1995). "Reliability and Authenticity: The Concepts and Their Implications," *Archivaria* 39 (Spring): 5-10.
- Henry, Charles & Smith, Kathlyn. (2010). "Ghostlier demarcations: Large-scale text digitization projects and their utility for contemporary humanities scholarship." In *The idea of order: Transforming research collections for 21st century scholarship*. Washington: CLIR, pp. 106-115.



References

- Kaplan, Abraham. (1964). "The Age of Symbol: A Philosophy of Library Education," *Library Quarterly* 34 (1964): 297.
- Rieger, Oya. (2008). *Preservation in the age of large-scale digitization: A white paper*. Washington: CLIR.
- Ross, Seamus. (2007). *Digital Preservation, Archival Science and Methodological Foundations for Digital Libraries*, Keynote Address at the 11th European Conference on Digital Libraries (ECDL), Budapest (September 17).
- Taylor, Hugh. (1987-88). "Transformation in the Archives: Technological Adjustment or paradigm Shift?."*Archivaria* 25 (Winter): 12-28.
- Thomassen, Theo. (2001). "A First Introduction to Archival Science," *Archival Science* 1: 373-385.
- Yeo, Geoffrey. (2008). "Concepts of Record (2): Prototypes and Boundary Objects." *American Archivist* 71 (Spring-Summer): 118-143.
- York, Jeremy. (2010). "Building a future by preserving our past: The preservation infrastructure of HathiTrust digital library." *Proceedings of 76th IFLA General Congress and Assembly*, 10-15 August, Gothenberg, Sweden.



SAA Research Forum -- August 2012

VALIDATING QUALITY IN LARGE-SCALE DIGITIZATION

use cases

results

reports

news

metrics

about

home

Project Summary

IN LESS THAN A DECADE the large-scale digitization of books has begun transforming the way we read and learn and changing how research libraries manage and preserve their collections.

DIGITIZED BOOKS made by third-party vendors are being preserved in online repositories. In this new preservation environment, the quality of what is preserved becomes an important factor in inspiring trust that digitized books are fit for the purposes envisioned for them.

INNOVATIVE RESEARCH presented at this website is developing and testing methods for measuring the severity of detectable errors in digitized books and validating the impact of error on the end-user. Here you will find information on the project, selected findings, and links to the project's reports, presentations, publications, and products.

HATHITRUST DIGITAL LIBRARY serves as a testbed of digitized books and serials for the project, which has three overlapping phases.

- Phase 1 (2011) Define a model of digitization error and a severity scale for recording observed error consistently and accurately.
- Phase 2 (2011-12) Apply the research methodology to representative samples of digitized volumes.
- Phase 3 (2012-13) Validate the results of the error analysis for specific use-case scenarios.

THE SCHOOL OF INFORMATION at the University of Michigan is leading a multi-year collaboration with the University of Michigan District Library and the present Libraries. The

Copyright © 2012 Regents of the University of Michigan

Institute of Museum and Library Services | School of Information | MLibrary | HathiTrust | University of Minnesota Libraries

Questions?

Thank you for your attention!

Project Website: <u>http://hathitrust-quality.projects.si.umich.edu/</u>

Paul Conway, Associate Professor University of Michigan School of Information

pconway@umich.edu

