# Archival Quality in Digital Preservation Repositories

## Constructing a New Approach to Metrics

*"I aspired to authenticity, but I never got beyond verisimilitude."*

# Outline of Presentation

▸ **Concepts**

▸ **Research Design**

▸ **Implications**

Archival Quality in Digital Preservation    10 August 2010

# Archival Quality – A Value Proposition

▸ **Archival nature**

  ▸ 1939 on : distinguishing characteristics of archives

▸ **Preservation media and procedures**

  ▸ 1961 on : technical characteristics of longevity

  ▸ 1985 on : protection against loss

▸ **Reliability [InterPARES]**

  ▸ 1995 on : completeness and process control

▸ **Significant properties**

  ▸ 2001 on : migration of essential elements

# Information Quality – DL Evaluation

- **IQ research establishes framework of attributes and clusters**
  - Wang & Strong (1996); Lin (2006) – MIS
  - Bovee (2003) – Accounting
  - Stvilia (2007); Rieh (2002) – Information Science
  - Knight (2008) – IQ/DQ community
- **Digital library evaluation establishes [mostly weak] end-user evaluation models and methods**
  - Saracevic (2005) – retrieval effectiveness
  - Saracevic (2007) – weaknesses in relevance research
  - Harley (2004); Pisciotta (2005) – image based user studies

Archival Quality in Digital Preservation    10 August 2010

# Research Environment

- **From vertical integration to distributed management**
  - "take what we can get"

- **HathiTrust**
  - 27 partners
  - 6 million+ volumes
  - Infrastructure, business model, TRAC certification



- **Google hysteria**
  - Data-poor reaction to a variety of socio-political-technical phenomena

# Les Archives de La France [Laborde, 1867]

# Google Book Search: Image and Text

**ARCHIVES DE LA FRANCE**

LEURS VICISSITUDES

PENDANT LA RÉVOLUTION

LEUR RÉGÉNÉRATION

SOUS L'EMPIRE

Le changement radical qu'ont subi les ar-
chives de la France pendant la Révolution est
tellement lié avec le cours des événements
politiques, que je suis amené, bien malgré
moi, en dehors de mes goûts & de mes habi-
tudes, à exprimer mon opinion sur le fait im-
mense qui s'appelle 89. Je ne l'aborderai
qu'autant qu'il se rattachera intimement au
sort des archives en servant à expliquer les
mesures fatales prises contre elles, & encore
je ne veux pas entrer dans cette voie sans faire
mes réserves. Je suis de ceux qui croient
qu'une nouvelle société pouvait se former
pour ainsi dire d'elle-même & sans martyriser

LES
ARCHIVES DE LA FRANCE
I. EUKS \iCISSITUDKW
PENDANT LA RÉVOLUTION
1 KUHI!É', V. KKRAIiON
sous L'EMPIRE
Le changement radical qu'ont subi les ar-
chives de la France pendant la Révolution est
tellement lié avec le cours des événements
politiques, que je suis amené, bien malgré
moi, en dehors de mes goûts & de mes habi-
tudes, à exprimer mon opinion sur le tait im-
mense qui s'appelle 89. Je ne l'aborderai
qu'autant qu'il se rattachera intimement au
sort des archives en servant à expliquer les
mesures fatales prises contre elles, & encore
je ne veux pas entrer dans cette voie sans faire
mes réserves. Je suis de ceux qui croient
qu'une nouvelle société pouvait se former
pour ainsi dire d'elle-même & sans martyriser

# Archival Quality / Large-Scale Digitization

**At Present:**
**Quality Standards**

- Material centered
- One size fits all
- Vertical integration
- <span style="color:red">Process control</span>
- Compromise is failure

**Trends Forward:**
**Acceptable Loss**

- User centered
- Fitness for use
- Third party creators
- <span style="color:red">Acceptance testing</span>
- Good enough is a value

# Use Cases

- ## Reading online
  - ### Digital page images
  - ### Text legibility; illustration interpretability; graphic accuracy

- ## Reading volumes printed on demand
  - ### Whole or substantial parts of volumes
  - ### Accuracy, completeness, consistency

- ## Processing full-text data
  - ### Underlying text content
  - ### Accuracy thresholds, readiness for analysis; "non-consumptive"

- ## Managing print collection
  - ### Surrogacy of the whole
  - ### Low cumulative error; non-critical errors; completeness; redundancy

# Two Views of Validation

- **Objective measurement of phenomena**
  - Definition of metrics
  - Testing of metrics
  - Statistical verification and confidence

- **Logical consistency from user's perspective**
  - Generalized error models
  - Few, but fatal, errors
  - Personalization of error perception

# Outline of Presentation

▸ Concepts

▸ **Research Design**

▸ Implications

Archival Quality in Digital Preservation    10 August 2010

# Research Question 1

What is "intrinsic quality" within the context of digitized books and serials? [or anything bound]

- Hierarchy of information errors based on prior research (IQ/DQ + UM, Google)
- Define and test measures of attribute error
  - Frequency and severity on ordinal scales
- Define and measure correlation effects across measures (co-occurrence)
- Build and test IQ indexes (accuracy, consistency, completeness, redundancy)
  - Cluster  and factor analysis

- Outcome: valid quality metrics + indices

Archival Quality in Digital Preservation    10 August 2010

# Incidence of Critical Error in HathiTrust

*University of Michigan Quality Review, 2006-10*

| Critical Error Type | Cause | May 2006-April 2007 | | May 2007-April 2008 | | May 2008-April 2009 | | May 2009-April 2010 | | TOTAL |
|---|---|---|---|---|---|---|---|---|---|---|
| Thick text | scanning | 189 | 0.57% | 70 | 0.19% | 19 | 0.06% | 144 | 0.81% | 422 |
| Broken text | scannng | 518 | 1.57% | 121 | 0.33% | 76 | 0.26% | 64 | 0.36% | 779 |
| Blurred text | scanning | 252 | 0.76% | 40 | 0.11% | 10 | 0.03% | 54 | 0.30% | 356 |
| Obscured text | source | 57 | 0.17% | 35 | 0.09% | 21 | 0.07% | 8 | 0.04% | 121 |
| Warped page | post-scan | 47 | 0.14% | 37 | 0.10% | 14 | 0.05% | 22 | 0.12% | 120 |
| Cropped text block | post-scan | 424 | 1.28% | 246 | 0.67% | 100 | 0.34% | 67 | 0.38% | 837 |
| Cleaning | post-scan | 208 | 0.63% | 214 | 0.58% | 1256 | 4.23% | 439 | 2.46% | 2117 |
| Colorization | post-scan | 3250 | 9.83% | 272 | 0.74% | 35 | 0.12% | 19 | 0.11% | 3576 |
| | | | | | | | | | | |
| Volumes ingested | | 288,044 | | 460,620 | | 2,523,049 | | 1,665,167 | | 4,936,880 |
| Volumes reviewed (20 pages/vol.) | | 33,047 | | 36,981 | | 29,677 | | 17,850 | | 117,555 |
| Ingested/Received | | 11.47% | | 8.03% | | 1.18% | | 1.07% | | 2.38% |

# Two Examples ["… flattening & thickening of meaning…"]

*Heather MacNeil*

## Warped Page



Critical Warp

## Thick Text



Critically Thick

# Errors in Source or Scanning

## Source Crop



The text is printed incorrectly and runs off the page. This is clearly a publisher crop and not a human error.
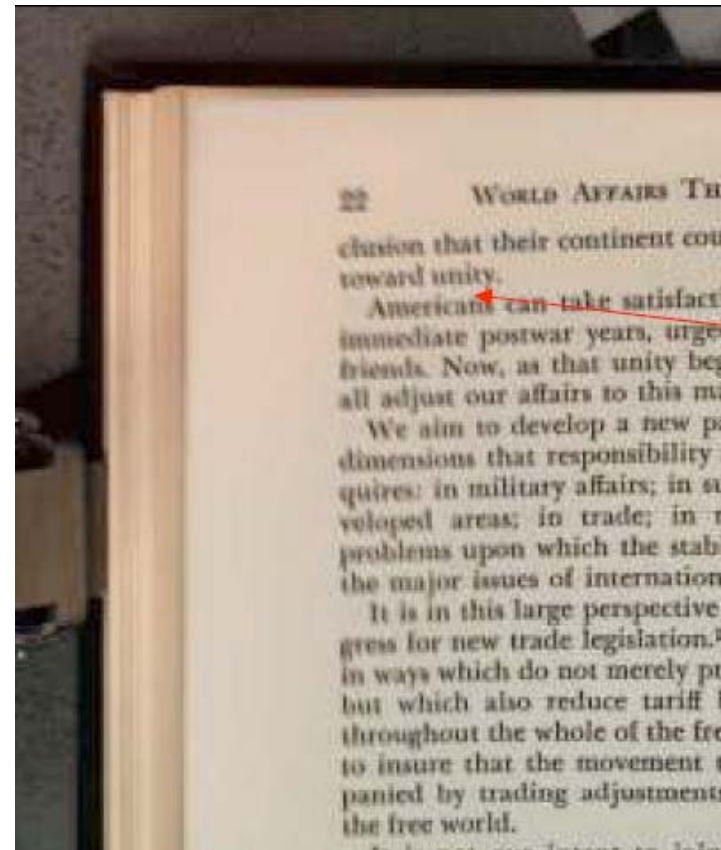
## Scan Crop

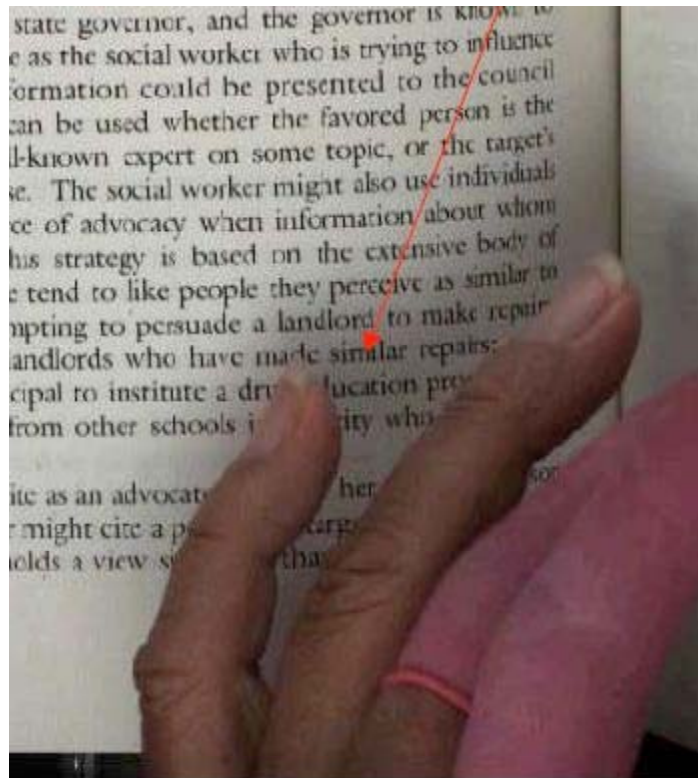# Errors in Source or Scanning

## Source Blur



## Scanning Blur



Archival Quality in Digital Preservation    10 August 2010

# Fingers in Manual Scanning

**Traces of human error**

**Traces  digitally cleaned**

Archival Quality in Digital Preservation    10 August 2010

# Error Model

*LEVEL 1: DATA/INFORMATION*
**1.1   Image: thick [character fill, excessive bolding, indistinguishable characters]**
**1.2   Image: broken [character breakup, unresolved fonts]**
1.3   Full-text: OCR errors per page-image
1.4   Illustration: scanner effects [moiré patterns, halftone gridding, lines]
1.5   Illustration: tone, brightness, contrast
1.6   Illustration: color imbalance, gradient shifts
*LEVEL 2: ENTIRE PAGE*
**2.1   Blur [movement]**
**2.2   Warp [text alignment, skew]**
**2.3   Crop [gutter, text block]**
**2.4   Obscured/cleaned [portions not visible]**
**2.5   Colorization [text bleed, low text to carrier contrast]**
2.6   Full-text: patterns of errors at the page level (e.g., indicative of cropping errors in digitization processing)
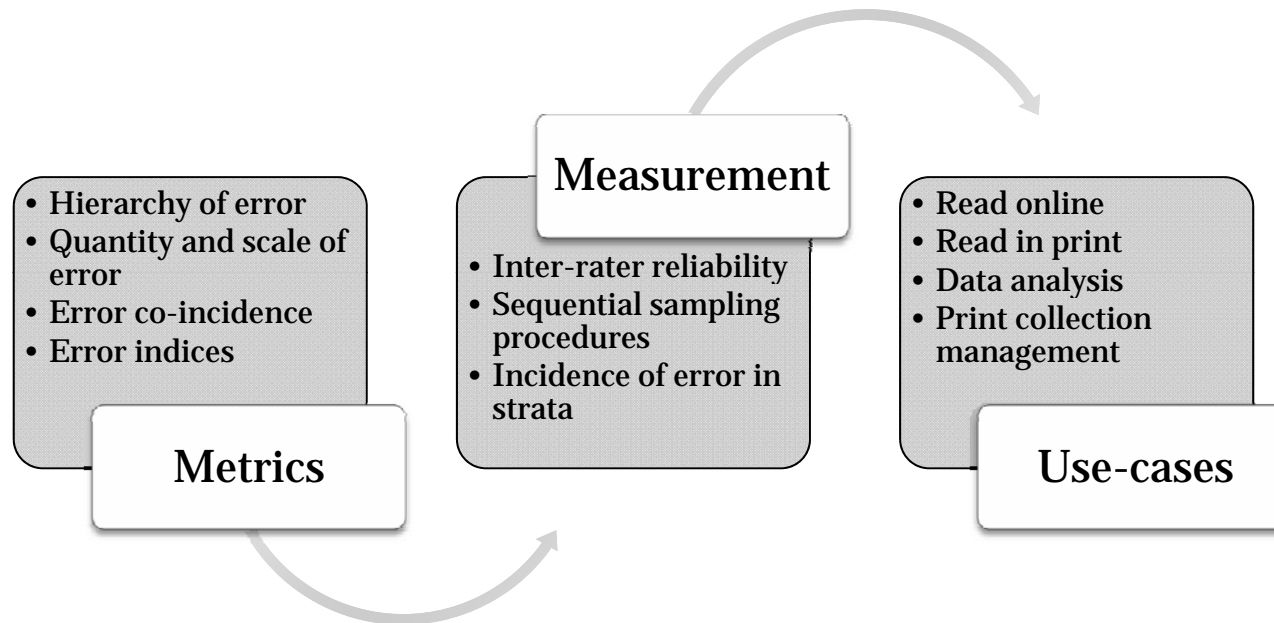*LEVEL 3: WHOLE VOLUME*
3.1   Order of pages [original source or scanning]
3.2   Missing pages [original source or scanning]
3.3   Duplicate pages [original source or scanning]
3.4   False pages [images not contained in source]
3.6   Full-text: patterns of errors at the volume level (e.g., indicative of OCR failure with non-Roman alphabets)

# Research Question 2

What is the estimated error-incidence in various clusters of HathiTrust content?

- Apply measures and indices (Q1) within selected strata
  - E.g., pub date; illustrations; source of digitization
- Extensive manual review of many random samples (some including original digitized books)
  - Examine differences between examining entire volume and samples from digital volumes
  - Compare digitized book with original book
- Assess and manage inter-coder inconsistencies in a distributed review model

▶ Outcome: costs and limits of manual review

▶ Outcome: identify potential for automated processing of quality review

▶ Outcome: mechanisms for branding quality using PREMIS metadata framework

# Research Workflow

Metrics

- Hierarchy of error
- Quantity and scale of error
- Error co-incidence
- Error indices

Measurement

- Inter-rater reliability
- Sequential sampling procedures
- Incidence of error in strata

Use-cases

- Read online
- Read in print
- Data analysis
- Print collection management

# Outline of Presentation

▸ Concepts

▸ Research Design

▸ **Implications**

# Implications for Preservation/DL Practice

▶ **Tools and techniques for measuring quality**

▶ **Expose content quality as part of certification process**

▶ **Limitations of use case scenarios**

    ▶ Fruitless pursuit of complete user satisfaction

▶ **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 content within.**

# Implications for Archival Theory

- Digital "archiving" through preservation is theoretically defensible
- Establish the archival nature of digitized surrogates
- Establish preservation value of digital surrogates
- Reaffirm relationship of provenance and reliability
- Archival quality defined through use

- Question: To what extent can or should a fundamental archival principle be measured?

# Acknowledgements

▶ **Planning support from Andrew W. Mellon Foundation**

▶ **John Wilkin, Exec. Director, HathiTrust**

▶ **Planning team**
  ▶ Jeremy York (HathiTrust)
  ▶ Emily Campbell (Mlibrary)
  ▶ Nikki Calderone (School of Information)
  ▶ Devan Donaldson (School of Information)
  ▶ Sarah Shreeves (University of Illinois)
  ▶ Robin Dale (Lyrasis)

# References (1)

- Bovee, M, Srivastava, R and Mak, B (2003). A Conceptual Framework and Belief-Function Approach to Assessing Overall Information Quality. *International Journal of Intelligent Systems.* 18 (1): 51-74.

- Cockburn, A (2000). *Writing Effective Use Cases.* Boston: Addison-Wesley.

- Harley, D et al. (2006). Use and Users of Digital Resources: A Focus on Undergraduate Education in the Humanities and Social Sciences. Berkeley: Center for Studies in Higher Education.

- Knight, S (2008). *User Perceptions of Information Quality in World Wide Web Information Retrieval Behaviour.* (PhD Dissertation) Perth, Australia: Edith Cowan University.

- Lin, X (2006). Quality Assurance in High Volume Document Digitization: A Survey. *Proceedings of the Second International Conference on Document Image Analysis for Libraries* (DIAL'06), 27-28 April, Lyon, France, pp. 319-326.

- Pisciotta, H et al. (2005). Penn State's Visual Image User Study," portal: Libraries and the Academy 5 (January): 33-58.

# References (2)

- Rieh, S (2002). Judgment of Information Quality and Cognitive Authority in the Web. J*ournal of the American Society for Information Science and Technology* 53 (2): 145-161.
- Saracevic, T. (2000). "Digital Library Evaluation: Toward an Evolution of Concepts." Library Trends 49 (2): 349-369.
- Saracevic, T. (2004). *How Were Digital Libraries Evaluated?* Paper first presented at the DELOS WP7 Workshop on the Evaluation of Digital Libraries, p. 9.
- Stvilia, B, et al. (2007). A Framework for Information Quality Assessment. *Journal of the American Society for Information Science and Technology* 58 (12): 1720-1733.
- Tanner, S, Munoz, T., and Ros, P. (2009). Measuring Mass Text Digitization Quality and Usefulness. *D-Lib Magazine* 15 (July/August): 209.
- Wang, R and Strong, D (1996). Beyond Accuracy: What Data Quality Means to Data Consumers. *Journal of Management Information Systems* 12 (4): 5-34.
- York, JJ (2009). This Library Never Forgets: Preservation, Cooperation, and the Making of HathiTrust Digital Library. *Proc. IS&T Archiving* 2009, Arlington, VA, pp. 5-10.

# Thank you for your attention!

Paul Conway, Associate Professor

School of Information, University of Michigan          pconway@umich.edu