HATHI TRUST
RESEARCH CENTER

# Text Mining HathiTrust as Big Data Exemplar in PRAGMA

**Beth Plale – @bplale**
**Professor, School of Informatics and Computing**
**Director, Data To Insight Center**
**Indiana University**



INDIANA UNIVERSITY

RESEARCH CENTER

ILLINOIS

Tweet us - @HathiTrust  #HTRC

# HathiTrust Digital Library

- HathiTrust is a partnership of academic & research institutions, offering a collection of millions of titles digitized from libraries around the world.

  – Founding members of HathiTrust along with University of Michigan are Indiana University, University of California, and University of Virginia

http://www.hathitrust.org

→ Distinguished from

RESEARCH CENTER

http://www.hathitrust.org/htrc

# HATHI TRUST Digital Library

Home | About | Collections | My Collections

## Currently Digitized

- 10,796,403 total volumes
- 5,658,745 book titles
- 281,890 serial titles
- 3,778,741,050 pages
- 484 terabytes
- 128 miles
- 8,772 tons
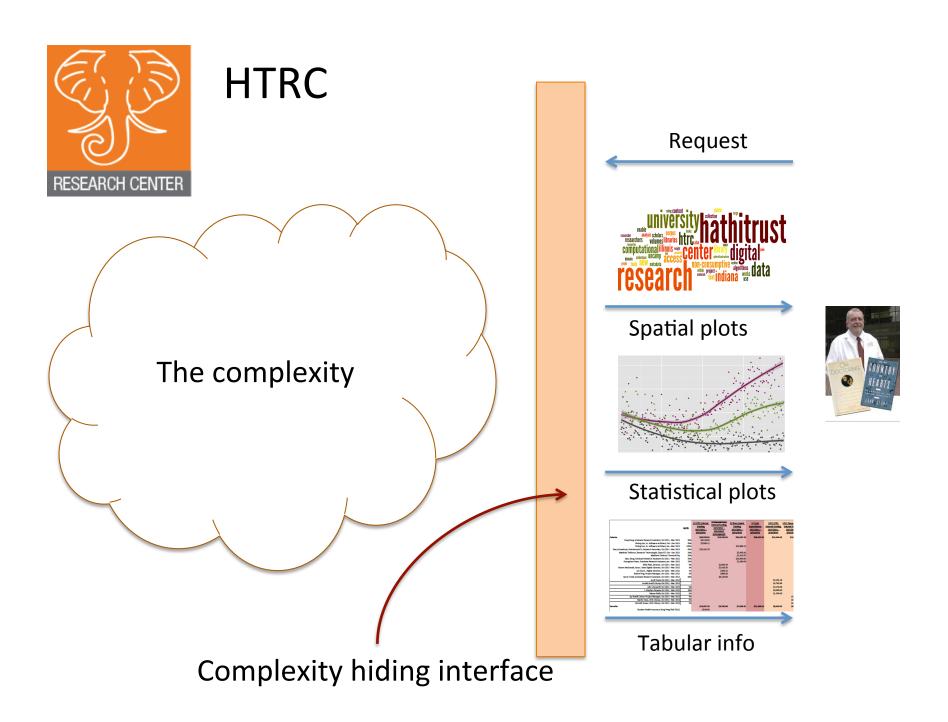- 3,450,939 volumes(~32% of total) in the public domain

View visualizations of HathiTrust call numbers, languages, and dates
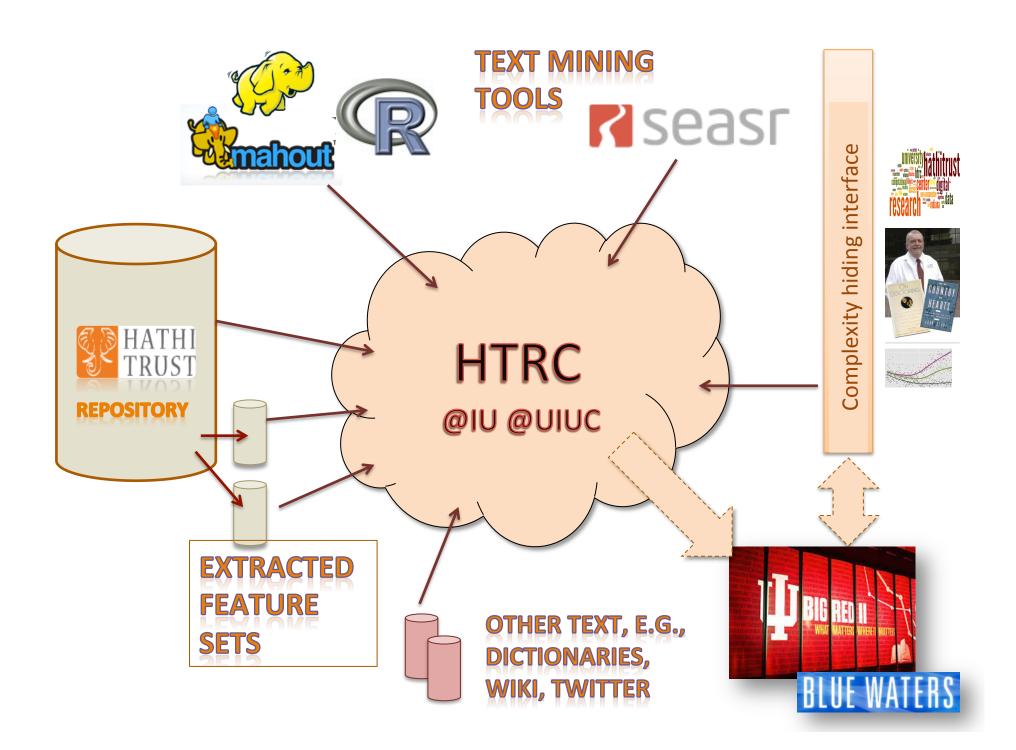statistics information ››

→ HathiTrust repository is a latent goldmine for text mining analysis, analysis of large-scale corpi through computational tools, and time-based analysis

→Restricted nature of HT content suggests need for new forms of access that preserve intimate nature of research investigation while honoring restrictions

→ Paradigm: computation takes place close to the data

# Mission of HT Research Center

- Research arm of HathiTrust
- Goal:  enable researchers world-wide to carry out computational investigation of HT repository through
  - Develop model for access: the 'workset'
  - Develop tools that facilitate research by digital humanities and informatics communities
  - Develop secure cyberinfrastructure that allows computational investigation of entire copyrighted and public domain HathiTrust repository
- Established:  July, 2011
- Collaborative effort of Indiana University and University of Illinois

HTRC

The complexity

Complexity hiding interface

Request

Spatial plots

Statistical plots

Tabular info

TEXT MINING TOOLS

seasr

HTRC
@IU @UIUC

Complexity hiding interface

HATHI TRUST
REPOSITORY

EXTRACTED FEATURE SETS

OTHER TEXT, E.G., DICTIONARIES, WIKI, TWITTER

BIG RED II
WHAT MATTERS WHERE MATTERS

BLUE WATERS

# HTRC architecture

- Philosophy: computation moves to data
- Web services (REST) architecture and protocols
- WS02 Registry for worksets and results
- Solr Indexes: full text, MARC, and new metadata
- noSQL (Cassandra) store as volume store
- Authentication using WS02 Identity Server
- Portal front-end, programmatic access
- Mining tools: currently SEASR

# Big Data Research Challenges in HTRC



Hooks to external collections / topical knowledge bases

e.g., Biodiversity Heritage Library

Integration: Categorizational and ontological sources

Text analytics at scale: Support for novel mining and NLP uses and mining

RESEARCH CENTER

Lifecycle issues in collection augmenting of data assets and services

Compute-to-data "Non-consumptive" text mining

# Big Data and PRAGMA

- Use PRAGMA member collections as anchor collections around which PRAGMA testbed could exercise experiments in data discovery and interoperability

- PRAGMA uses its testbed to exercise relevant activity of Research Data Alliance.  RDA WG on Type Registries is about to deliver; assess usefulness in PRAGMA use cases

- Data on program for PRAGMA 27:  devote time to address/expose policy and cultural issues in data sharing; are there tech solutions to policy issues that PRAGMA would be interested in?