# The Center for Research Libraries

# Text-Mining Opportunities and Challenges

## November 13, 2012

The *Center for* Research Libraries
GLOBAL RESOURCES NETWORK

# CRL Global Resources Forum

Analysis

Intelligence

Library investment

# Today's agenda

- **Recent Trends in Text-Mining and Library Services: Research Library Perspective**

  David Magier, Princeton University Library

- **Elsevier Perspective on Mining the Scientific Literature**

  Judson Dunham, Elsevier

- **Gale's Support of Text- and Data-Mining Requests**

  Ray Abruzzi, Cengage Learning

- **Text-Mining in the Cloud**

  Ann Okerson, Center for Research Libraries

*The* **Center** *for* **Research Libraries**
GLOBAL RESOURCES NETWORK

# Presenters

- **David Magier**

  Associate University Librarian
  for Collection Development

  Princeton University Library

# Recent Trends in Text Mining and Library Services:

# Research Library Perspective

David Magier
Assoc. Univ. Librarian for Collection Development
Princeton University Library

# What we do with text

- Use it to enable more granular **discovery**
  - to find and read specific chunks of content
  - "consumptive" use

- Use lots of it to analyze large corpora of text
  - subject text to statistical analyses to discover **meaningful correlations and trends**
  - "non-consumptive" use

# Examples

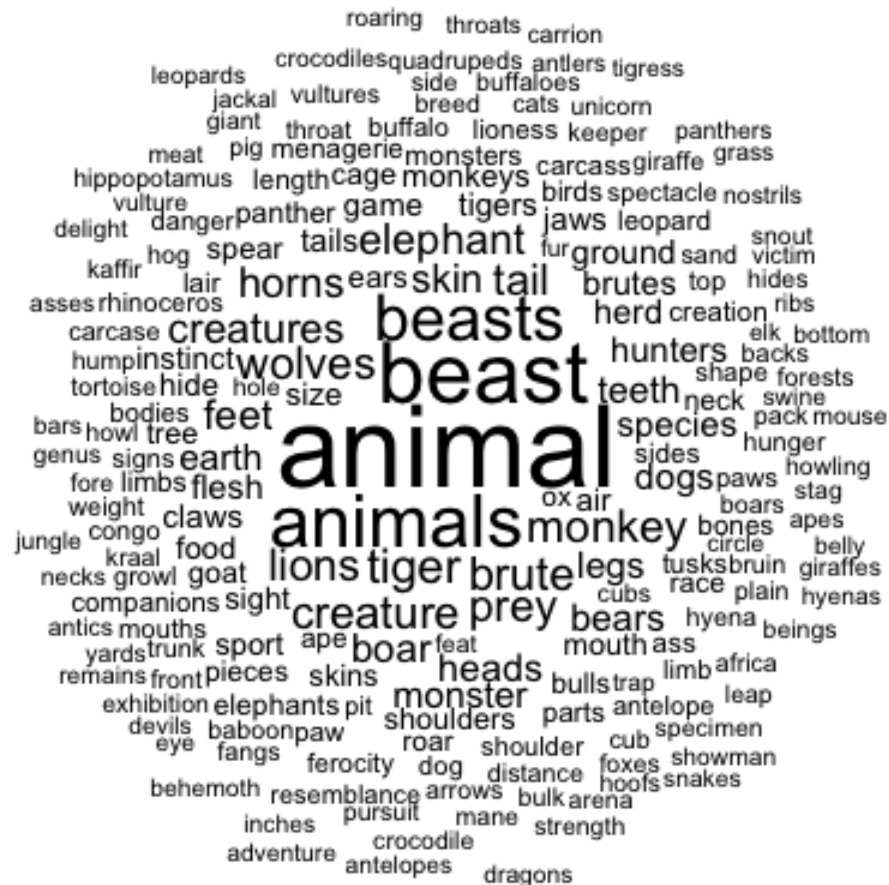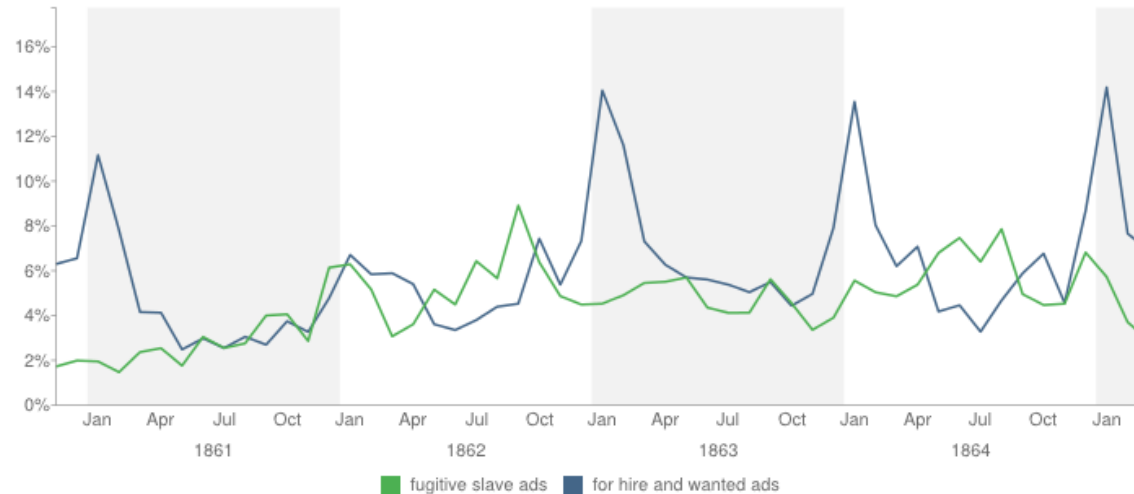| | |
|---|---|
| **Author/Artist:** | Martin, Meredith, 1976- |
| **Title:** | The rise and fall of meter : poetry and English national culture, 1860-1930 / Meredith Martin. |
| **Edition:** | 1st ed. |
| **Published/Created:** | Princeton, N.J. : Princeton University Press, c2012. |
| **Physical description:** | 274 p. : ill. ; 24 cm. |
| **Contents:** | Introduction: the failure of meter -- Modern instability -- Metrical communities -- Meter as culture -- A note on historical prosody -- The history of meter -- A metrical history of England -- A grammatical history of England -- Grammatical instability -- Metrical instability -- The stigma of meter -- Metrical irrelevance -- The British empire of letters -- Marking instress -- Acute stress in "The wreck of the Deutschland" -- Mistrusting the ear -- The institution of meter -- Metrical mastery -- Inventing the "Britannic" -- Dynamic reading -- Mastery for the masses -- The English ear -- A prosodic entity -- The discipline of meter -- Patriotic pedagogy -- Matthew Arnold's metrical intimacy -- Henry Newbolt's cultural metrics -- Private meters, public rhythms -- The sound of the drum -- The trauma of meter -- Wartime, poetics -- Sad death for a poet! -- Therapeutic measures -- Bent-double -- The kindred points of heaven and home -- The before- and afterlife of meter -- Metrical modernism -- Make it old : Robert Bridges and obsolescence -- Alice Meynell's "English metres" -- Toward a critical prosody. |
| **ISBN:** | 9780691152738 (hc. : alk. paper)<br>069115273X (hc. : alk. paper)<br>9780691155074 (pbk. : alk. paper<br>0691155070 (pbk. : alk. paper |

# Matthew Jockers:
## 500 Themes from a corpus of 19th-Century Fiction

roaring throats carrion
crocodilesquadrupeds antlers tigress
leopards side buffaloes
jackal vultures breed cats unicorn
giant throat buffalo lioness keeper panthers
meat pig menagerie monsters carcassgiraffe grass
hippopotamus length cage monkeys birds spectacle nostrils
vulture danger panther game tigers jaws leopard
delight snout
hog spear tails elephant fur ground sand victim
kaffir lair horns ears skin tail brutes top hides
asses rhinoceros herd creation ribs
carcase creatures beasts elk bottom
humpinstinct wolves hunters backs
tortoise hide hole size beast shape forests
bodies feet teeth neck swine
bars howl tree species pack mouse
genus signs earth animal sides hunger
fore limbs flesh dogs paws howling
weight ox air boars stag
jungle congo claws animals monkey bones apes
kraal food circle belly
necks growl goat lions tiger brute legs tusksbruin giraffes
companions sight cubs race plain hyenas
antics mouths creature prey bears hyena beings
yards trunk sport ape boar feat mouth ass
remains front pieces skins heads bullstrap limb africa
exhibition elephants pit monster parts antelope leap
devils baboonpaw shoulders cub specimen
eye fangs roar shoulder foxes showman
ferocity dog distance hoofs snakes
behemoth resemblance arrows bulk arena
inches pursuit mane strength
adventure crocodile
antelopes dragons

*Macroanalysis: Digital Methods and Literary History* (UIUC Press, 2013),
http://www.matthewjockers.net/macroanalysisbook/macro-themes/

# Robert Nelson: Mining the Dispatch



**Mining the *Dispatch***

INTRODUCTION | TOPICS | ABOUT

Legend: ■ fugitive slave ads  ■ for hire and wanted ads

"Mining the Dispatch," seeks to explore—and encourage exploration of—the dramatic and often traumatic changes as well as the sometimes surprising continuities in the social and political life of Civil War Richmond. It uses as its evidence nearly the full run of the Richmond *Daily Dispatch* from the eve of Lincoln's election in November 1860 to the evacuation of the city in April 1865. It uses as its principle methodology topic modeling, a computational, probabilistic technique to uncover categories and discover patterns in and among texts. On this site you'll be able to view and generate graphs and charts that reveal some of the changing patterns in the topics that dominated the news during the Civil War in the capital of the Confederacy's newspaper of record.

Robert K. Nelson, rnelson2@richmond.edu, Digital Scholarship Lab, the University of Richmond

# The Corpus:
## GoogleBooks Ngram Viewer

**Google** books **Ngram Viewer**

Graph these **case-sensitive** comma-separated phrases: | Albert Einstein,Sherlock Holmes,Frankenstein |

between 1800 and 2000 from the corpus English ⇕ with smoothing of 3 ⇕ .

Search lots of books

■ Albert Einstein  ■ Sherlock Holmes  ■ Frankenstein

http://books.google.com/ngrams

# Paper Machines:
## a text-mining visualization plug-in for Zotero from Harvard's MetaLab



**Four topics shown for all documents mentioning India (top) and/or Ireland (bottom).**

# The Corpus:
## The Library as Data-Set: Text Mining at Million-Book Scale



The Library as Dataset: Text Mining at Million-Book Scale

YaleUniversity · 1,065 videos

189 views

2 likes, 0 dislikes

Subscribe 28,231

David Mimno:
http://www.youtube.com/watch?v=o_jgjinLRlQ

# The Corpus:
## HathiTrust Research Center

**HATHI TRUST** Digital Library

Home | About | Collections | My Collections

Our Partnership | Our Digital Library | **Our Research Center** | News and Publications | Help

About › Our Research Center

## Our Research Center

### Overview

The HathiTrust Research Center (HTRC) enables computational access for nonprofit and educational users to published works in the public domain and, in the future, on limited terms to works in-copyright from the HathiTrust.

The HTRC is a collaborative research center launched jointly by Indiana University and the University of Illinois, along with the HathiTrust Digital Library, to help meet the technical challenges of dealing with massive amounts of digital text that researchers face by developing cutting-edge software tools and cyberinfrastructure to enable advanced computational access to the growing digital record of human knowledge.

Leveraging data storage and computational infrastructure at Indiana University and the University of Illinois at Urbana-Champaign, the HTRC will provision a secure computational and data environment for scholars to perform research using the HathiTrust Digital Library. The center will break new ground in the areas of text mining and non-consumptive research, allowing scholars to fully utilize content of the HathiTrust Library while preventing intellectual property misuse within the confines of current U.S. copyright law.

http://www.hathitrust.org/htrc

# Where does the library fit?

- Who creates tools?

- How does content become "a corpus"?

- **Who can connect the scholars to the content AND the tools they need?**

# Presenters

- **Judson Dunham**

  Senior Product Manager, Elsevier

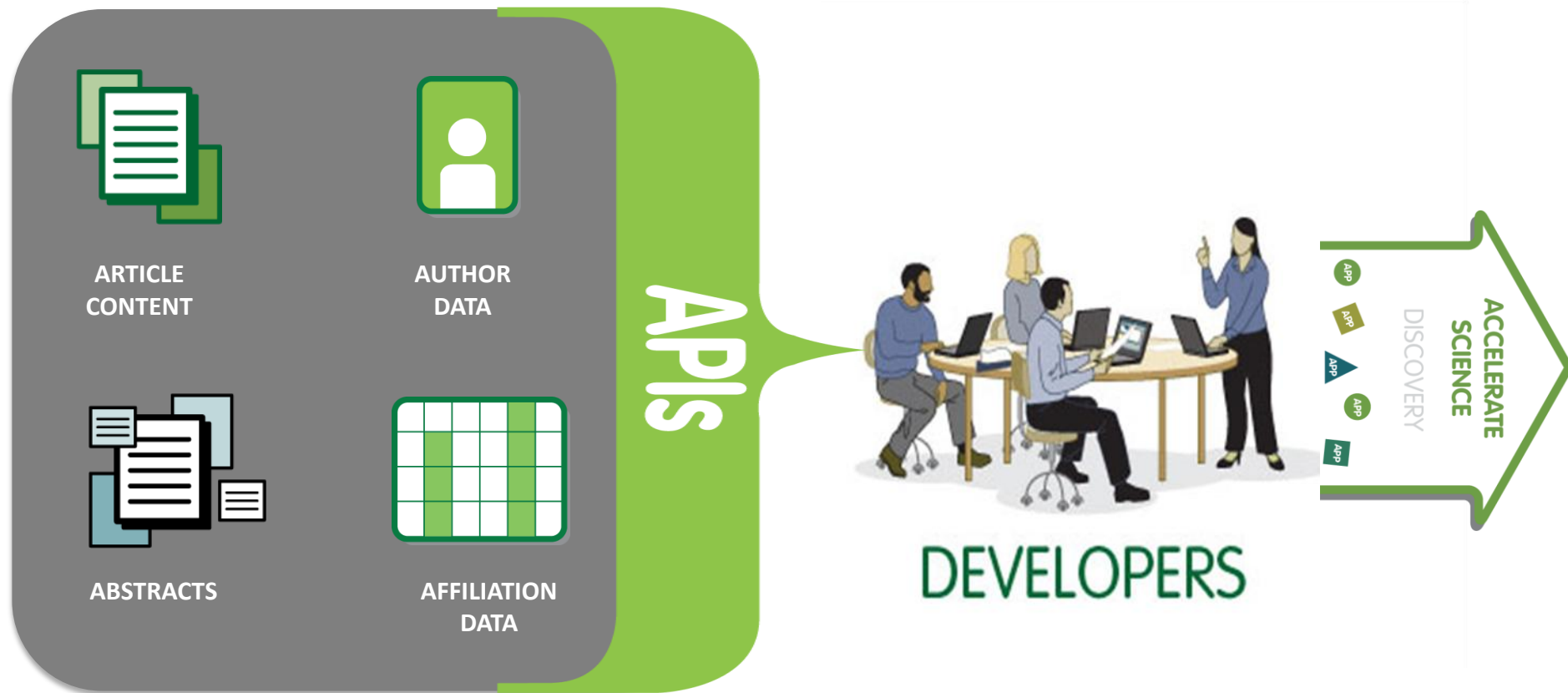# Elsevier Perspective on Mining the Scientific Literature

Judson Dunham
Sr. Product Manager, Elsevier

# 2010

ARTICLE CONTENT

AUTHOR DATA

ABSTRACTS

AFFILIATION DATA

APIs

DEVELOPERS

APP
APP
APP
APP
APP

DISCOVERY

ACCELERATE SCIENCE

## Collaboration at a different level

# 2012

# What researchers want from content mining

- Assurance of right to mine subscribed or freely available content across publishers
- Clear understanding of what they can do with the outputs
- Standardized, mineable formatting of data
- Stable, reliable systems for retrieving and storing content

# How Elsevier is working to address this

## Simple process to support text mining

- Researchers/Librarians can request text mining access for specific projects on case by case basis
- This is a pilot period – there is still much to learn about researcher needs

## Provide efficient methods for content retrieval

- **Access to APIs for individual researchers** via click through agreements on the web
- **Slimmed down content formats** optimized for text mining

# How Elsevier is working to address this

## Clear guidelines on reuse and redistribution

- Simple **CC-BY-NC** license on output when redistributing results of text mining in a research support or other non-commercial tool

- **DOI link back to mined article** whenever feasible when displaying extracted content

- Clear guidelines and permissions on **snippets around extracted entities** to allow context when presenting results

# USE CASE: Studying patterns of data sharing and reuse in the published literature

## Text mining small subset of content

- Using literature as a data source to explore a research issue (in this case, how is data in repositories shared, referenced and reused in the literature)

- Requires 1,000's of documents (not 1,000,000's)

- Will result in papers published in journals, as well as a publicly available database providing information about reuse of specific datasets (mined from the literature)

**Solution:** Access to subscribed content via API's, permission to reuse results in non-commercial tool

*Heather Piwowar - Postdoc at Duke University*

# USE CASE: Extraction of DNA sequences from millions of papers to facilitate biocuration and discovery

## Text mining large volume of content

Identify all words in full-text scientific articles resembling DNA sequences, which are extracted and then mapped to public genome sequences. They can then be displayed on genome browser websites and used in data-mining applications.



*Max Haeussler – Post Doc, UCSC*

# USE CASE: Extraction of DNA sequences from millions of papers to facilitate biocuration and discovery



Link to paper with Metadata and abstract

Mined Sequence (bold) with flanking sentence

*Max Haeussler – Post Doc, UCSC*

# USE CASE: Extraction of DNA sequences from millions of papers to facilitate biocuration and discovery



*Max Haeussler – Post Doc, UCSC*

# Still some unsolved problems

## Decentralized content aggregation:

- Retrieving and storing data locally = shipping loads of content around the web/world

- Still tons of formatting and preparation work for *every single text miner*

- Significant computational infrastructure required to do large scale text mining

**A cloud-based service for running large-scale text mining without having to build infrastructure or worry about the "plumbing" would be valuable.**

# Still some unsolved problems

## Preservation of knowledge:

- Storing and redistributing (valuable) text mining results requires long term infrastructure, hosting and other services

- Many academic research efforts end "when the postdoc moves on" – where does the mined knowledge go?

- **A service offering a reliable, persistent, managed, open venue for the hosting of content mining results would be valuable**

# END

- j.dunham@elsevier.com/@judso

**Thanks:**

- Max Haeussler (UCSC)
- Heather Piwowar (UBC/Duke)
- Bradley Allen (Elsevier Labs)

# Presenters



- **Ray Abruzzi**

  Director, Strategic Planning
  Learning and Research Solutions

  Cengage Learning

# TRANSFORMING LEARNING

## TRANSFORMING LIVES

CENGAGE Learning™

Ray Abruzzi
Director, Strategic Planning

## Gale Scholarly Collaborations

- **Gale has been involved with or is currently supporting over 30 projects to-date involving data-mining and/or textual analysis**

- **Efforts range from individuals to institutional to project-based groups**

- **Gale supports these efforts through a variety of means, most often providing:**

  - **Raw data consisting of entire collections, subsets of collections, or subsets of multiple collections, in multiple, specific formats**

  - **Technical advice and support**

  - **Connecting scholarly efforts through our involvement in multiple projects**

TRANSFORMING
LEARNING
TRANSFORMING
LIVES

CENGAGE
Learning

## Examples: IMPACT (Improving Access to Text)

**Objectives: Significantly improve access to historical text**

- Innovate OCR technology
  - By exploring the challenges using different approaches
  - By developing cutting-edge approaches such as collaborative correction

- Provide innovative language technologies to remove the historical language barrier

- Ensure the interoperability of the results
  - By defining an overall technical architecture and monitoring technical integration across all parts of the project

- Take away the barriers that stand in the way of the mass digitization of the European cultural heritage

- Provide Best Practice guidance about the operational context for digitization

- **Deliver a coherent program of dissemination, training and demonstration aimed at capacity-building in and beyond participating institutions**

## Examples: IMPACT (Improving Access to Text)

**Gale's participation:** Working with the British Library, IMPACT was provided with sample images and OCR'd text from the *19th Century British Library Newspapers*.

**What they accomplished:** This content formed part of the overall project dataset which was used by researchers to benchmark, develop/optimize and test language and OCR tools for historical material.

**What the benefit was:** This research led to increased OCR success for nine European languages/historical material and generated new areas of research interest.

CENGAGE
Learning™

## Examples: 18th Connect (and NINES)

**Objectives:** NINES and 18thConnect are communities of scholars in 19th-century and 18th-century studies, respectively, organized around web sites that aggregate metadata and provide links to digital scholarly materials in their fields.

- In addition, for 18th-century studies, Gale has contributed page images (which are shown in snippets) and the OCR-generated texts from the ECCO collection to 18thConnect.

- These image snippets and the corresponding OCR text are made visible in 18thConnect's TypeWrite tool which allows 18thConnect's users to correct the OCR.

- Gale has given 18thConnect the opportunity to offer as a reward the full, typed texts (without page images) to anyone who corrects it. **Corrected texts may also be loaded back into ECCO to improve searching.**

- **18thConnect has also run new OCR engines on the page images from ECCO with the intent on creating improved OCR XML which could then be loaded back into the ECCO database, again resulting in a better search experience.**

TRANSFORMING
LEARNING
TRANSFORMING
LIVES

CENGAGE
Learning

## Examples: 18th Connect and NINES

**Gale's participation:** Provided OCR text and metadata for *Eighteenth Century Collections Online* and a range of 19th century content (some 19C material is metadata only).

**What they accomplished:** The metadata for ECCO is fully loaded and allows all scholars to search through it in conjunction with metadata from other scholarly databases. **The TypeWrite program is fully functional though no corrected texts have yet been returned to Gale.**

**What the benefit was**: At the current status of the projects, the 18thConnect and NINES websites provide a federated search-like approach to research across a range of metadata from key scholarly databases, allowing users to do one search instead of doing the same search in various standalone databases.

- If a user of these websites finds a work of interest within ECCO and their institution is a subscriber/owner of ECCO, they can then pass directly into the native application and see the entire work within ECCO.

TRANSFORMING
LEARNING
TRANSFORMING
LIVES

CENGAGE
Learning

## A typical request

- "My name is XXX and I have a gig at XXX working for their Cultural Evolution of Religion project. I'm also the customer who is requesting access to the text files in the metadata of ECCO documents. See attached three images. Two are screenshots of my Advanced Search window, one of which I have highlighted each of the fields I have toggled. The only search term is 'earthquake', as you'll see. The difference between the two is that one isolates only documents in Religion and Philosophy whereas the other aims to capture all documents (in the same set of years) that contain 'earthquake'. The third image is a screenshot of my Gale Learning/ECCO user account page, which represents the search results in the these two searches. The narrow search issued 498 docs whereas the broad search issued 1453. Thanks for your help with this. If you can download the XML files from what I've given you now, I'll be very happy as it takes some time to click through on each page and copy and paste the Gale Doc ID numbers. **I need this by the end of the week to meet my project deadlines. Thank you.**"

TRANSFORMING
LEARNING
TRANSFORMING
LIVES

CENGAGE
Learning

## Benefits to the academic community and to Gale

- **Gale is pleased to support academic efforts and projects of this nature**

  - **Of course, Gale is keen to generate goodwill in the academic community and to further scholarship**

- **The *expected outcomes* and *anticipated results* of these projects expand Gale's knowledge of both the changing needs of researchers (our customers) and broaden our understanding in terms of the kinds of tools Gale should be building into our resources**

  - **NCCO incorporates several new tools which enable a degree of textual analysis, and Gale is expanding these tools to our entire Digital Collections program over the next two years, starting with ECCO in 2013**

- **Gale may be able to improve or modify existing collections with improved data or tools based on the outcomes of various projects**

TRANSFORMING
LEARNING
TRANSFORMING
LIVES

CENGAGE
Learning

# Challenges—Rights and Licenses

- **Rights: Gale licenses content from over 250 institutions to create our Digital Collections program**

  - **In many cases, this license permits very specific/limited uses of content and the data associated with that content**

  - **Gale often needs to seek permission to use this data in ANY project that does not fall under the provisions of the contract—institutions are reluctant to sign up for "any and all" uses for their content**

  - **Seeking permission or an amendment to the contract can take days, weeks, or months**

  - **Gale consistently seeks the broadest rights for our content licenses, with varying degrees of success**

  - **The process to acquire these rights for non-commercial projects-- which may or may not yield tangible, actionable benefits to Gale-- drains time and resources from other efforts**

  - **A license or contract with the person, institution or group is also needed, to ensure that the terms of use are not violated**

TRANSFORMING
LEARNING
TRANSFORMING
LIVES

CENGAGE
Learning

## Challenges: Actionable outcomes, tangible benefits

- **These are projects, often based on theories and experimentation—results may or may not yield tangible benefits**

- **Marginal improvements to searchability or OCR processes are not cost-justified. In other words, Gale cannot re-process millions of pages of data to make them "slightly" better**

- **Even when successful, results are often specific to the outcomes being sought, and are not transferable to any product of program offered from Gale**

- **The variety of data formats and tools and the processes utilized in these projects may not synch up with Gale—there may not be a ready means of ingestion for the outcomes, and the cost of creating one may far exceed the benefits of doing so**

TRANSFORMING
LEARNING
TRANSFORMING
LIVES

CENGAGE
Learning

## Summing up

- **Gale continues to work with the digital humanities communities and support projects as we are able.**

- **The number and complexity of requests for data are growing, and Gale cannot donate resources to meet all of these requests**

- **Gale is actively seeking rights to use/utilize content (data) more broadly from source institutions to expedite our support of these projects**

- **Gale has built some initial tools for textual analysis and data-mining into our interfaces, and continues to expand on those tools to meet the needs of users at the point of research**

**Thank you—Questions?**

Ray Abruzzi, Director, Strategic Planning.
Ray.abruzzi@cengage.com

# Presenters

**Ann Okerson**

Senior Advisor on Electronic Strategies,

Center for Research Libraries

# Text Mining in the Cloud

*Ann Okerson*

*Center for Research Libraries*

**Webinar - November 2012**

access property quickly every lot scholarly thousands before tens modern big sense new look said Intellectual Trust Dr people knows scale one because more suggest available just published between all publishers open links want tool journals use publicly find gene create other drugs bought Any technique world allow scientific allowed ways report university science scientists bn through genetics last similar Murray-Rust DNA sequence papers already Kiley Professor Bergman Elsevier Piwowar able Text2genome diseases never year individual computers powerful genes funded help know sequences led things patterns Though healthcare articles data mining literature information potential once mining make content text research researchers

# What we've been talking about

- "Automated processing of large amounts of digital data/content for purposes of information retrieval, extraction, interpretation, analysis" (Bernie Reilly)
- "Automated tools, techniques or technology to process large volumes of digital content that is often not well structured – to identify and select relevant information; to extract information from the content, to identify relationships within/between/across documents and incidents or events for meta-analysis" (Efke Smit)
- Access to data at a degree of granularity greater than an individual reader can make; and ability to process that granular data
- Mining terminology is a little vexing – text mining, data mining, content mining

# Nothing new under the sun

- Don Swanson (1924-), American information scientist, known for his work in literature-based discovery in the biomedical domain
  - For some years, Dean of the U of Chicago Library School
- His method is called "Swanson linking:" connecting two pieces of knowledge previously thought to be unrelated
- Part of "Arrowsmith System project," which sought to determine meaningful links between Medline articles
  - It may be known that illness A is caused by chemical B, and that drug C is known to reduce the amount of chemical B in the body
  - However, because the respective articles were published separately from one another (called "disjoint data"), the relationship between illness A and drug C may be unknown
- Don called this "undiscovered public knowledge"

# Then, the next generation

- Data are extracted and processed behind the scenes:

- The Perseus Project (started in the 80s) automatically analyses grammar and vocabulary of ancient Greek texts and makes predictions about the function of a specific word in a specific passage

- So if the reader is not sure about a specific word (*nomoisi*) and clicks on it, a parallel window pops up to show several different dictionary entries all spelled *nomos*, gives statistical information, meanings, and links to further information.

- See the text and then the interpretations:

# Herodotus, *The Histories*

Your current position in the text is marked in blue. Click anywhere in the line to jump to another position.

book:

chapter:

section:

Hdt. 1.94.1

Click on a word to bring up parses, dictionary entries, and frequency statistics

## This text is part of:

Greek and Roman
Materials

Greek Prose

Greek Texts

Herodotus

Λυδοὶ δὲ νόμοισι μὲν παραπλησίοισι χρέωνται καὶ Ἕλληνές, χωρὶς ἢ ὅτι τὰ θήλεα τέκνα καταπορνεύουσι, πρῶτοι δὲ ἀνθρώπων τῶν ἡμεῖς ἴδμεν νόμισμα χρυσοῦ καὶ ἀργύρου κοψάμενοι ἐχρήσαντο, πρῶτοι δὲ καὶ κάπηλοι ἐγένοντο.

Herodotus, with an English translation by A. D. Godley. Cambridge. Harvard

## View text chunked by:

# Greek Word Study Tool

**νόμος**                                    anything assigned, a usage, custom, law, ordinance
(Show lexicon entry in Middle Liddell Slater) (search)

| νόμοισι †   noun pl masc dat epic ionic aeolic | | no user votes | 92.5% | [vote] |

† This form has been selected using statistical methods as the most likely one in this context. It may or may not be the *correct* form. (More info)

Word Frequency Statistics (more statistics)

| Words in Corpus | Max | Max/10k | Min | Min/10k | | Corpus Name |
|---|---|---|---|---|---|---|
| 184,947 | 706 | 38.173 | 0 | 0 | | Herodotus, The Histories |

---

**νόμος**                                    that which is in habitual practice, use
(Show lexicon entry in LSJ) (search)

| νόμοισι   noun pl masc dat epic ionic aeolic | | no user votes | 2.5% | [vote] |

Word Frequency Statistics (more statistics)

| Words in Corpus | Max | Max/10k | Min | Min/10k | | Corpus Name |
|---|---|---|---|---|---|---|
| 184,947 | 706 | 38.173 | 0 | 0 | | Herodotus, The Histories |

---

**νομός**                                    place of pasturage,
(Show lexicon entry in LSJ Middle Liddell Autenrieth) (search)

| νόμοισι   noun pl masc dat epic ionic aeolic | | no user votes | 2.5% | [vote] |

Word Frequency Statistics (more statistics)

| Words in Corpus | Max | Max/10k | Min | Min/10k | | Corpus Name |
|---|---|---|---|---|---|---|
| 184,947 | 706 | 38.173 | 0 | 0 | | Herodotus, The Histories |

---

**νομός**                                    herbage
(Show lexicon entry in Slater) (search)

# Potential for research

- A critical mass of content now provides huge opportunities
  - Archives so large that there's no way to manually analyze them
  - A lot of this work starting in the past decade
  - Research needs rising, more demand in disciplines
- Efke Smit (May 2011):  majority of scholarly publishers receive requests for mining less than 10 x per year; 5 for research purposes
- Publishers generally liberal in granting permissions; wish to understand this phenomenon and often go case by case
  - Reluctant to compete with their own products and services
- English is dominant language of scientific publication
- Don Swanson's dream being realized

# Happening now

- Newer, better software tools are being created to handle large corpora of content
- Biomed most advanced, along with corporate pharmaceutical and chemistry
- Growing number of users in government intelligence, business, social media, news, financial markets, legal arena
- Policy discussions at the national level, in the UK in particular, with the Hargreaves and JISC reports
- Lots of negative focus; given low, though rising demand, how did we so quickly get into a wrangle with publishers about TM?
- Automatic assumption that publishers have too many economic interests and won't cooperate

# Licenses

- Copyright/fair use status of text mining not clear
- Licenses can clarify what copyright does not
  - Licenses should permit automated queries for research and education – however it is done
  - Explicit language "text mining" or implicit language "for educational and research purposes"
- CDL have been seeking rights since 2008 and have secured some agreements using their standard language
- Overall, more and more library licenses have TM provision, but users not aware that library can help
- We will move from case by case to global rights and permissions – not so different after all from ILL

# Next Steps

- Libraries can become more aware of campus needs and offer support/expertise
- Libraries can encourage publishers in licensing and researcher support.
- Topics span not just one publisher but cross publishers
  - There are numerous different publisher platforms – we need standards and standard content formats
  - How to do this?  Cross-Rev, Cross-Mark, Cross-Check, Cross-Mine?
- Facilitating researchers' needs requires infrastructure development in various ways
- Collaboration is required across publishing, libraries, the research community
- Librarians can participate in facilitating such activities

# Contact Information

- **Bernie Reilly, CRL**
  breilly@crl.edu
- **Virginia Kerr, CRL**
  vkerr@crl.edu
- **Ann Okerson, CRL**
  aokerson@crl.edu
- **David Magier, Princeton**
  dmagier@princeton.edu
- **Judson Dunham, Elsevier**
  j.dunham@elsevier.com
- **Ray Abruzzi, Cengage Learning**
  ray.abruzzi@cengage.com

*The* **Center** *for* **Research Libraries**
GLOBAL RESOURCES NETWORK

# Upcoming Events

**CRL Webinar on News Collections in a Digital Age:**
January 14, 2013


**Annual Meeting of the CRL Council of Voting Members:**
April 19, 2013


**CRL Workshop at ALA Annual Meeting:**
**Global Resources Roundtable: "Beyond the Fold – Scholarly Access to News in the Digital Age"**
June 27-28, 2013, Newberry Library, Chicago

*The* **Center** *for* **Research Libraries**
GLOBAL RESOURCES NETWORK

# For More Information

- Fill out our follow-up survey at
  http://www.surveymonkey.com/s/CRL_Text-Mining_Webinar_Nov_2012

- Check out other CRL presentations on our YouTube channel
  www.youtube.com/crldotedu

- Visit CRL website www.crl.edu

- Sign up for *CRL Connect:* www.crl.edu/connect

- Find CRL on Facebook and Twitter

*The* **Center** *for* **Research Libraries**
GLOBAL RESOURCES NETWORK