*Gale Digital Scholar Lab*

Colloquium On Text & Data Mining in Libraries 2023

University of Toronto

Jess Ludwig, Director, Product Management

Gale, here for **everyone.**

GALE
A Cengage Company

# Digital Humanities at Gale

# Expand the Possibilities of Research

## GALE PRIMARY SOURCES



## GALE DIGITAL SCHOLAR LAB



Gale, here for **everyone.**

# Global Partnerships

**Gale Fellowships**

# Layers of Discovery

**Document Discovery & Use**

**Cross-Search**

**Simple Analysis**

**Text & Data Mining**

Gale, here for **everyone.**

GALE
A Cengage Company

Boolean operators and multiple search fields provide users any number of search criteria combinations.

Users will find many refining options to help limit the number of results to those that are most relevant to the definition of each corpus they create.

Gale, here for **everyone.**



GALE DIGITAL SCHOLAR LAB

My Content Sets  Build  Clean  Analyze

Basic Search  Search by keyword  Advanced Search

SEARCH OPTIONS
Advanced Search

# Advanced Search

**Search Terms**

| Terms | Field | Finds results that... |
|---|---|---|
| Search for | in Keyword | contain these terms in key fields; does not search entire document |
| And | in Keyword | contain these terms in key fields; does not search entire document |
| And | in Keyword | contain these terms in key fields; does not search entire document |

Search  Add a Row

**Search Tips**

*Operators*  *Special Characters*
AND, OR, NOT  Proximity  Nesting  Quotation Marks  Wildcards  Ignored

**Selected Databases to Search (50/50)**

All
Select All  Deselect All

☑ Amateur Newspapers from the American Antiquarian Society
☑ American Fiction, 1774-1920
☑ American Historical Periodicals from the American Antiquarian Society
☑ Archives of Sexuality and Gender
☑ Archives Unbound
☑ Associated Press Collections Online
☑ Brazilian and Portuguese History and Culture
☑ British Library Newspapers
☑ China and the Modern World
☑ Crime, Punishment, and Popular Culture, 1790-1920
☑ Daily Mail Historical Archive
☑ Declassified Documents Online: Twentieth-Century British Intelligence
☑ Eighteenth Century Collections Online
☑ Indigenous Peoples of North America
☑ International Herald Tribune Historical Archive, 1887-2013
☑ Liberty Magazine Historical Archive, 1924-1950
☑ Mirror Historical Archive, 1903-2000
☑ Nineteenth Century Collections Online
☑ Nineteenth Century U.S. Newspapers
☑ Nineteenth Century UK Periodicals
☑ Picture Post Historical Archive, 1938-1957
☑ Political Extremism and Radicalism
☑ Public Health Archives: Public Health in Modern America, 1890-1970
☑ Punch Historical Archive, 1841-1992
☑ Refugees, Relief, and Resettlement: Forced Migration and World War II
☑ Religions of America
☑ Sabin Americana: History of the Americas, 1500-1926
☑ Seventeenth and Eighteenth Century Burney Newspapers Collection
☑ Seventeenth and Eighteenth Century Nichols Newspapers Collection
☑ Slavery and Anti-Slavery: A Transnational

**Search Limiters**

by publication section:

by publication country:

by publication state/province:

by publication city:

by publication year(s):
⦿ All  ○ Before  ○ Within  ○ After  ○ Between
☐ Include documents with no known publication date.

by content type:
Select Content Type(s)

by document type:

by publication title:

by language:
Select Language(s)

Source Library
Select Source Library

Illustrated Works
Select Illustrated Work(s)

OCR Confidence Range
Min ____  —  Max ____
Documents can range from 0 - 100 in OCR Confidence Score

Search  Clear

**Build** Your Content Sets
- Use Gale Primary Sources archives
- Upload and use your own text files
- Download content sets for use elsewhere

Learn more about the Build step »

**Clean** Texts for Computational Analysis
- Apply stop words to your analysis
- Use flexible options to target specific character removal
- Reuse configurations across Content Sets

Learn more about the Clean step »

**Analyze** Content in Powerful New Ways
- Visualize data from up to 10,000 documents at a time
- Explore individual documents overlaid with analysis data
- Download the raw data and your visualizations

Learn more about the Analyze step »

GALE
A Cengage Company

Use a default stop word configuration, customize it, or create your own and then test with a sample from your content set.

Further customize what you'd like to remove or find and replace in the content set.

GALE DIGITAL SCHOLAR LAB

My Content Sets    Build    Clean    Analyze

New Configuration    Save As    Test Configuration

## Clean

**Select Cleaning Configuration to Edit**

Default Cleaning Configuration

Video: Cleaning a Content Set
Additional topics:
- Creating a Clean Configuration
- Applying During Analysis
- View All Clean Help »

01:09

### Stop words

Set the words you want the Analysis Tools to ignore.
Choose a Starter List

a
about
above
across
after
afterwards
again
against
all
almost
alone
along
already
also
although
always
am
among
amongst
amoungst
amount
an
and
another
any
anyhow
anyone
anything
anyway
anywhere
are
around
as
at
back
be
became

☑ Ignore stop words case

Clear List

### Text Correction

Options for automatic text correction that will be applied before each Analysis

☐ Turn on all options

▾ **Text Modification**
  ☐ All lower case

▾ **Characters**
  ☐ Remove all extended ASCII characters
  ☐ Remove all number characters

▾ **Special Characters**
  ☐ Remove all special characters
  Set specific characters ›

▾ **Punctuation**
  ☐ Remove all punctuation
  Set specific punctuation ›

▾ **Spacing**
  ☑ Remove all tabs
  ☑ Remove all line breaks
  ☑ Reduce multiple spaces to one space (ex: "hello there" becomes "hello there")

▾ **Document Sections**
  ☐ Remove body text
  ☑ Remove all non-body content
  Set specific sections ›

▾ **Replacements**
  Replace this...    With this...

Add a Row

### Configuration Notes

Space to make notes or describe the purpose of this configuration.

Configuration description/notes

Gale, here for **everyone.**

GALE
A Cengage Company

# Primary Sources Analysis

## EXPANDING RESEARCH POSSIBILITIES

- Sentiment Analysis

- Parts of Speech

- Named Entity Recognition

- Ngrams

- Topic Modeling

- Document Clustering

Gale, here for **everyone.**

## Legend

**View**

● Top 100 Ngrams

○ Ngram Search

[ Enter terms to match ngrams    🔍 ]

**Color Coding** ❓

● Frequency    ○ Rank

Low Frequency          High Frequency

Legend
Tool Setup
Run History



federal government
fourteenth amendment constitution          chief justice
railway labor act  school year      working class     plaintiffs exhibit  state alabama
amendment constitution
communist party  racial discrimination  constitutional right
interstate commerce  junior high  grand jury  attorney general  fifth amendment
labor act  state courts       rights act  case bar       findings fact
ques tion  labor relations  state law  equal protection
national labor relations  railway labor         york city  state board  public school
matter fact  equal protection clause  civil rights
elementary school  race color  cert denied  took place  constitutional rights
general assembly  voting rights  defendants exhibit  public schools  protection clause
member communist  civil rights act  high schools  board education  answer question
process clause  duly sworn  years ago  instant case
north carolina       act use  fourteenth amendment
white students  title vii  county board  writ certiorari  school board
black students  amend ment       fifth circuit  process law  national committee
federal courts  south carolina       present case
petition writ  national labor  school districts  trial judge  clause fourteenth
high school  state action  pet app  protection laws  school boards
county school  objection sustained  force violence  period time  time time  cause action
right vote  equal protection laws  govern ment  petition writ certiorari
legislative history  governments exhibit  state york  common law  civil action
elementary schools  racial segregation
state louisiana  questions presented
reasonable doubt
subject matter

Gale, here for **everyone.**

GALE
A Cengage Company

**Parts of Speech**
1650-1800 Coffee - Manners & Customs

Download · Help

Back

## Legend

### Parts of Speech

Click to toggle categories on and off

| NOUN | INTJ | PART |
| PROPN | CCONJ | PUNCT |
| PRON | SCONJ | NUM |
| VERB | DET | SYM |
| ADV | ADP | SPACE |
| ADJ | AUX | OTHER |

### Group Data By

Author

10 of 192 selected
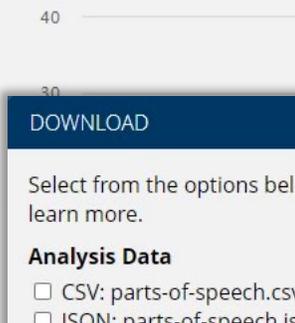Change selection

Graphs Shown: ❓ 10 ⌄

Grid Columns: 4 ⌄

Legend · Tool Setup · Run History

### Adam Petrie
1 documents | Tagged
26,243 Parts of Speech

### Alexander Pope
5 documents | Tagged
305,224 Parts of Speech

### Althea Fanshawe
1 documents | Tagged
34,760 Parts of Speech

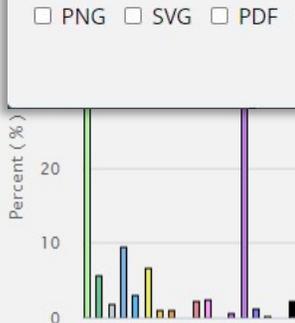### American in England
1 documents | Tagged
25,357 Parts of Speech

### Anne Fisher
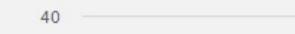1 documents | Tagged
111,368 Parts of Speech

### Argens, marquis d'
11 documents | Tagged
814,165 Parts of Speech

### Arthur Murphy
4 documents | Tagged
277,736 Parts of Speech

### Arthur William Costigan
2 documents | Tagged
87,586 Parts of Speech

---

**DOWNLOAD** ✕

Select from the options below. Visit the Learning Center to learn more.

**Analysis Data**
☐ CSV: parts-of-speech.csv
☐ JSON: parts-of-speech.json

**Document Metadata**
☐ JSON: document-metadata.json

**Parts of Speech Visualization**
☐ PNG ☐ SVG ☐ PDF ☐ JPEG

Cancel · Download
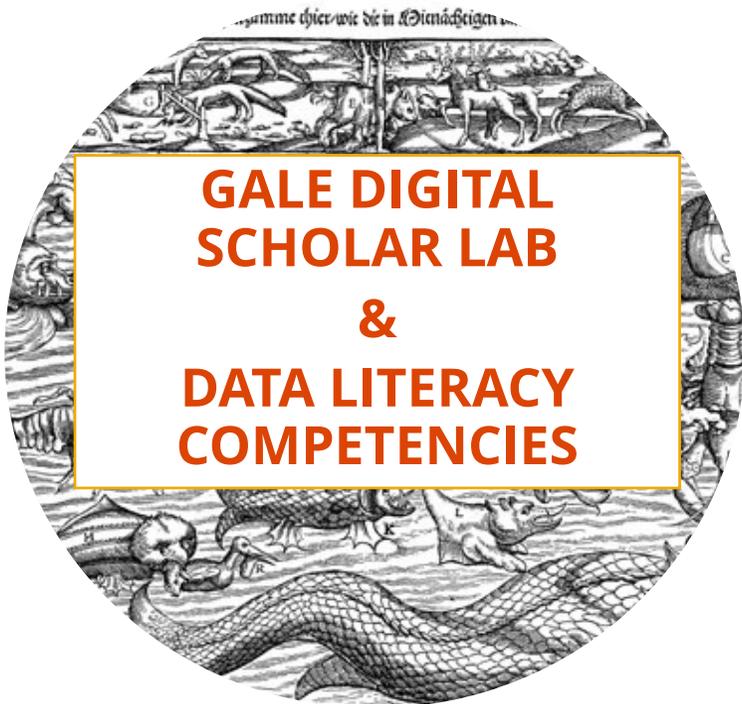
---

## ＜ Inspect ✕

### Document

The works of Alexander Pope, Esq; vol. IV. Part II. Containing the second part of his letters

**Open Document**

### Parts of Speech

| Category | Count | Percentage |
| --- | --- | --- |
| NOUN | 14,972 | 27.78% |
| PROPN | 4,901 | 9.09% |
| PRON | 1,339 | 2.48% |
| VERB | 6,044 | 11.21% |
| ADV | 1,789 | 3.32% |
| ADJ | 3,962 | 7.35% |
| INTJ | 399 | 0.74% |
| CCONJ | 349 | 0.65% |
| SCONJ | 0 | 0% |
| DET | 919 | 1.71% |
| ADP | 1,226 | 2.27% |
| AUX | 50 | 0.09% |
| PART | 368 | 0.68% |
| PUNCT | 16,038 | 29.76% |
| NUM | 652 | 1.21% |

**GALE DIGITAL SCHOLAR LAB**

Get Link

Digital Scholar Lab
Learning Center

Getting Started Gale Digital Scholar Lab
What is the Digital Scholar Lab?

01:15

Translate Article

## Build Your Content Sets

Collect your research materials in a Content Set by finding documents from the Gale Primary Sources databases or uploading your own.

Build Overview
Search Strategies
Understanding Search Results
Building a Content Set
Uploading Documents

Video Overview: Building a Content Set

03:21

## Clean Texts and Prep for Analysis

Use options like stop word lists and text correction for the removal or replacement of specified characters and terms to prepare for analysis.

Clean Overview
Creating a Clean Configuration
Testing Your Configurations
Applying During Analysis

Video Overview: Cleaning Your Texts for Analysis

00:56

## Analyze Content Sets with Powerful Tools

Use digital tools to analyze your Content Sets in ways that would have been too time consuming without the power of computational algorithms.

Analyze Overview
Selecting the Right Tool
Setting Up and Running
Tool: Document Clustering
Tool: Named Entity Recognition
Tool: Ngrams
Tool: Parts of Speech
Tool: Sentiment Analysis
Tool: Topic Modeling

Video Overview: Analyzing Content Sets

01:22

## Organize with My Content Sets

Use the My Content Sets pages to organize and manage your research with helpful tools to examine the makeup and ensure it has what you need.

My Content Sets Overview
Managing a Content Set
Downloading Content Sets

Video Overview: Managing Content Sets

02:02

---

Black America and the Law in the Mid-20th Century

The Rise of Electricity in the Late 19th and Early 20th Centuries

Food and Civility 1650-1800

---

## Sample Syllabi

These complete syllabi serve as course models to how digital humanities concepts, ideas, and activities can be incorporated into traditional Digital Humanities courses or Special Topics Humanities courses. We provide examples of course scope and sequences that integrate key phases of the Gale Digital Scholar Lab workflow that may culminate in a final project using the Lab. You can access our three sample syllabi here:

Digital Humanities
Exploring Topics, Online Course
Exploring Topics, 10-week Course

**GALE** A Cengage Company

**GALE DIGITAL SCHOLAR LAB & DATA LITERACY COMPETENCIES**

**Build**
- Interpret and critically evaluate data and their sources
- Read/understand data types and formats
- Find, select, access, or create datasets in order to test a hypothesis or answer a research question
- Create metadata to meet data publication requirements

**Clean**
- Clean / process / convert data

**Analyze**
- Analyze data
- Communicate data effectively to different audiences, in part by using visualizations
- Ethically collect / use / cite data
- Integrate and synthesize data into different contexts with other sources and prior knowledge

# Gale Digital Scholar Lab Case Studies

# McGill University: Riddle Culture in the 18th & 19th Centuries

Enigmatic Bills of Fare led to explorations of riddling culture in 18th and 19th centuries.



**Figure 1.** Printed Enigmatic Bill of Fare. Bodleian, John Johnson Collection



**Figure 2.** Manuscript Enigmatic Bill of Fare. Cadbury Research Library: Special Collections

Gale, here for **everyone.**

# Newcastle University: Writing New Worlds, 1688-1789



Ngram Word Cloud depicting the most frequently occurring words and phrases contained in the Equiano and Wheatley texts, with cleaning configuration applied.

# Loyola University: Anthropology 100: Globalization & Local Cultures



The availability of digital humanities data and the tools to analyze it has allowed me to truly engage ALL my students in research. – *Catherine Nichols, Advanced Lecturer in Cultural Anthropology and Museum Studies, Loyola University*

Gale, here for **everyone.**

# University of Washington: Applied Digital Humanities



Gale, here for **everyone.**

# Gale Digital Scholar Lab Roadmap

# User Engagement Group



**Catherine Nichols**
Loyola University

**Kathrina Perry**
University of Northampton

**Jose Intriago Suarez**
Marquette University

**Qiuyang Chen**
University of Warwick

**Oihane Etayo**
University of Warwick

**Caitlin Bagley**
Gonzaga University

# User Engagement Program

**PURPOSE**

- General Platform Feedback

- Validation & Usability feedback for in progress projects

- Discovery feedback for upcoming projects

Gale, here for **everyone.**

# New Data Sets

- Watergate in the News

- Watergate Declassified Docs

- Suffragettes

- Stonewall Riots

- The Murder of Olof Palme

- The Boxer Uprising

- Roberto Calvi Trial

- Sojourner Truth

Dataset 1: Watergate in the News

**Get a copy of the dataset**

Get a Copy

The Watergate in the News Dataset was compiled using Gale Primary Sources and Gale Digital Scholar Lab.
**Archives Used:** The International Herald Tribune Digital Archive; The Daily Mail; The Telegraph; The Sunday Times; The Times Digital Archive.
**Scope:** 1972 - 1975. 4741 documents.
**Source Libraries:** The New York Times Company (2194); Telegraph Media Group (1091); Times Newspapers Limited (1049); Associated Newspapers Limited (406).
**Document Types:** Article (4257); Editorial (251); Letter to the editor (144); Back matter (60); Front matter (29).

Gale, here for **everyone.**

GALE
A Cengage Company

# Group Collaboration

**WHAT**

A collaborative space in which students can work together to create content sets, analyse their materials and build their project as a group.
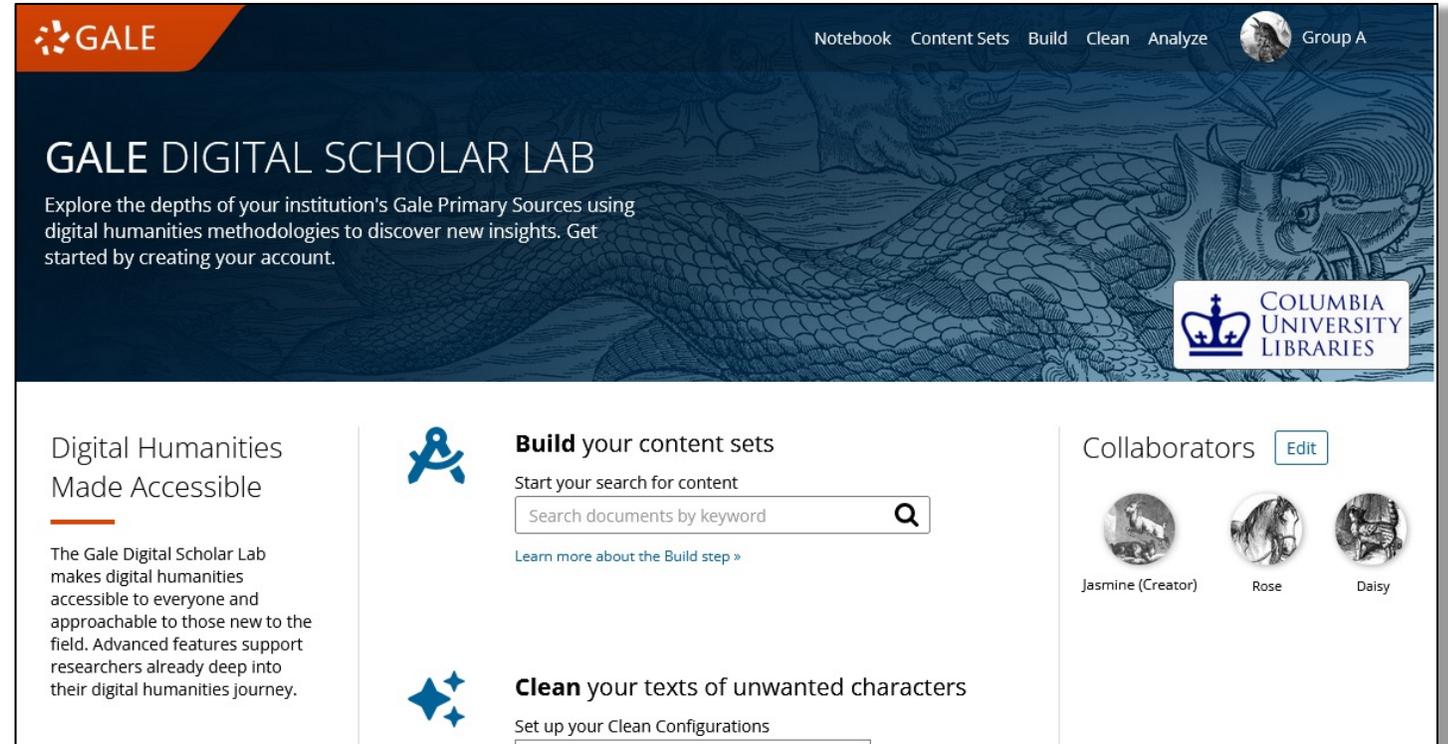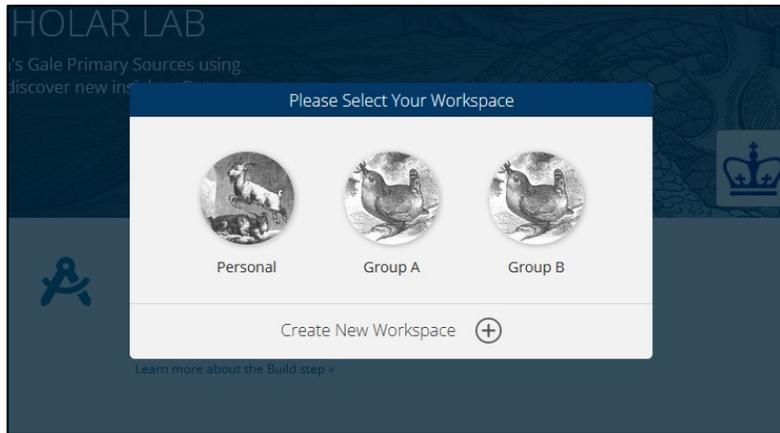
**WHY**

47% of survey respondents (Spring 2021) identified collaboration tools as Extremely or Very Important. Integrating the Lab in the classroom increases usage and the likelihood of renewal.
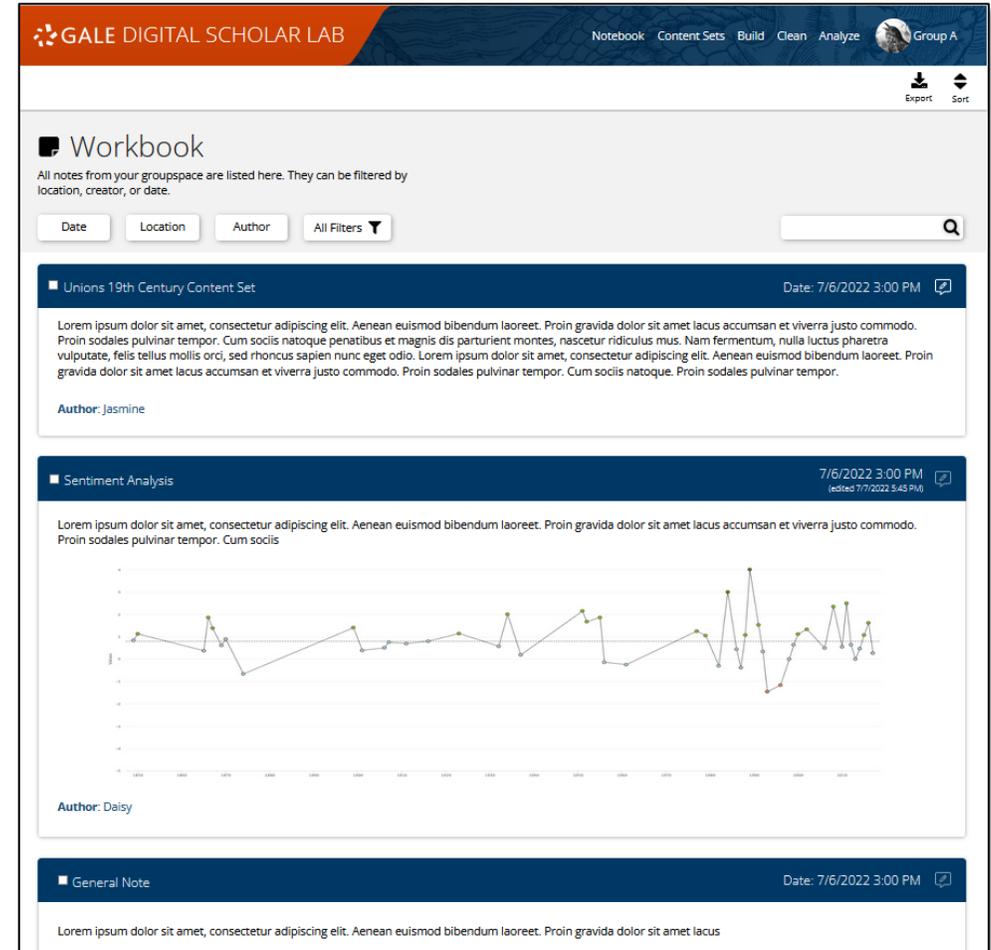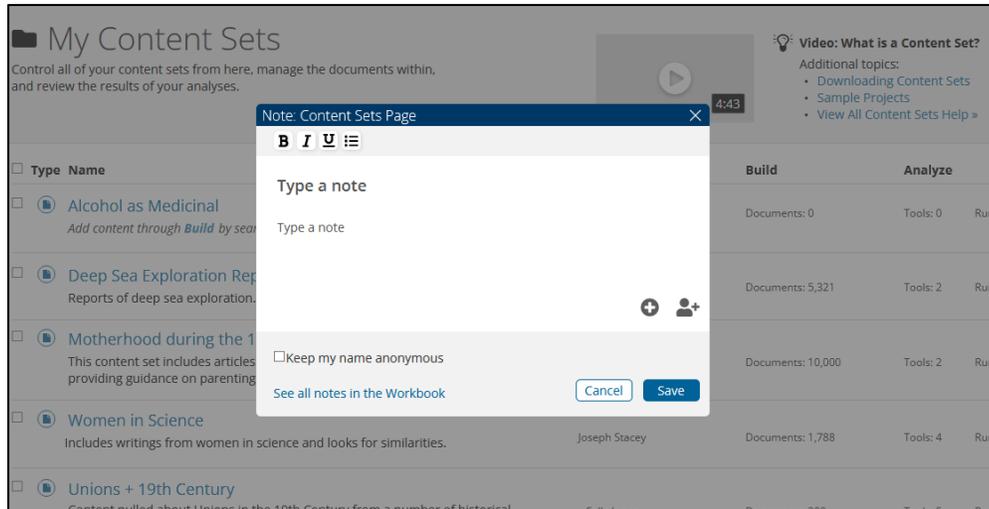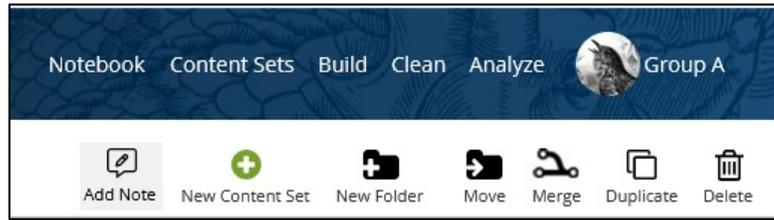
**HOW**

Working closely with our users to define requirements, build and validate prototypes and undertake extensive usability testing to ensure the we are meeting user needs.

Gale, here for **everyone.**

# Group Workspace

# Group Collaboration Notebook

# Thank you!