# Working around Walled Gardens

The Princeton Prosody Archive as Workflow
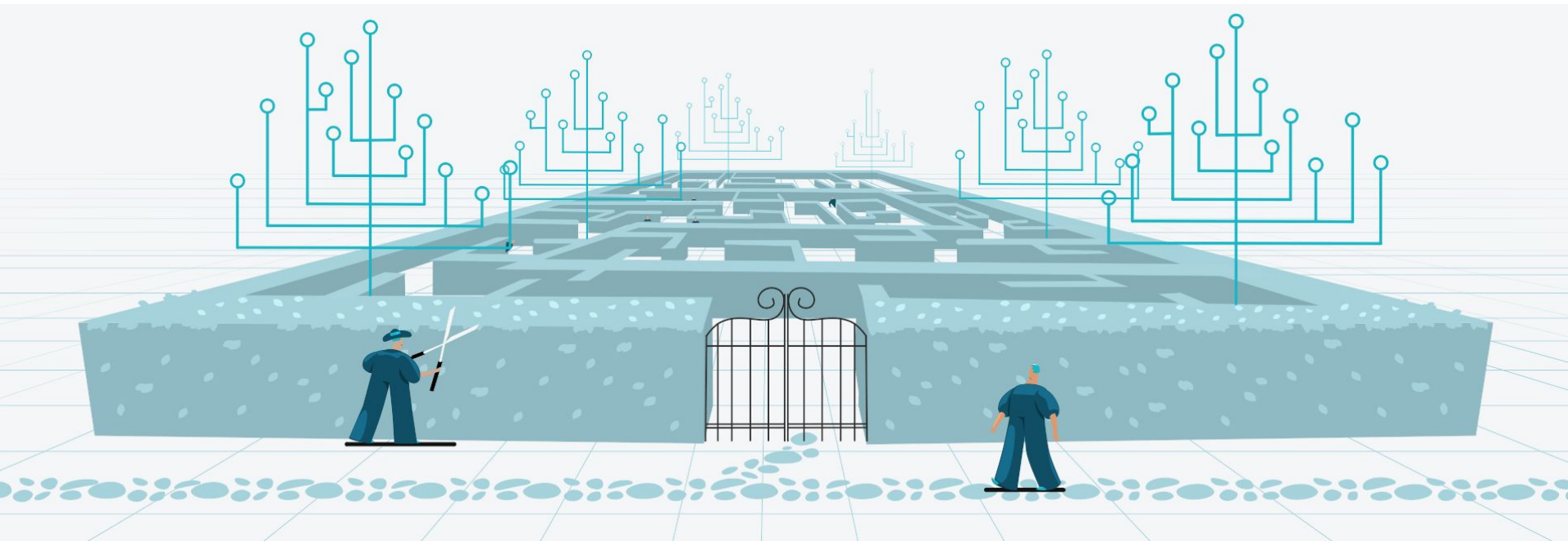
Mary Naydan, Rebecca Sutton Koeser, Meredith Martin

DARIAH Annual Event  |  June 20, 2024

THE CENTER
FOR DIGITAL
HUMANITIES
@PRINCETON

# Overview

1. **What are Walled Gardens?**
   a. The Digital Research Landscape: How We Got Here
   b. Limitations of Existing "Workflows"
   c. Additional Challenges with TDM
   d. Collaborative DH Projects
2. **Case Study: The Princeton Prosody Archive**
   a. Overview
   b. Data Sources
   c. Technical Architecture as Workflow
   d. "Dead Ends" & Roadblocks
   e. The Path through the Hedge Maze
3. **Futures Beyond Walled Gardens**

THE CENTER
FOR DIGITAL
HUMANITIES
@PRINCETON

# What are Walled Gardens?
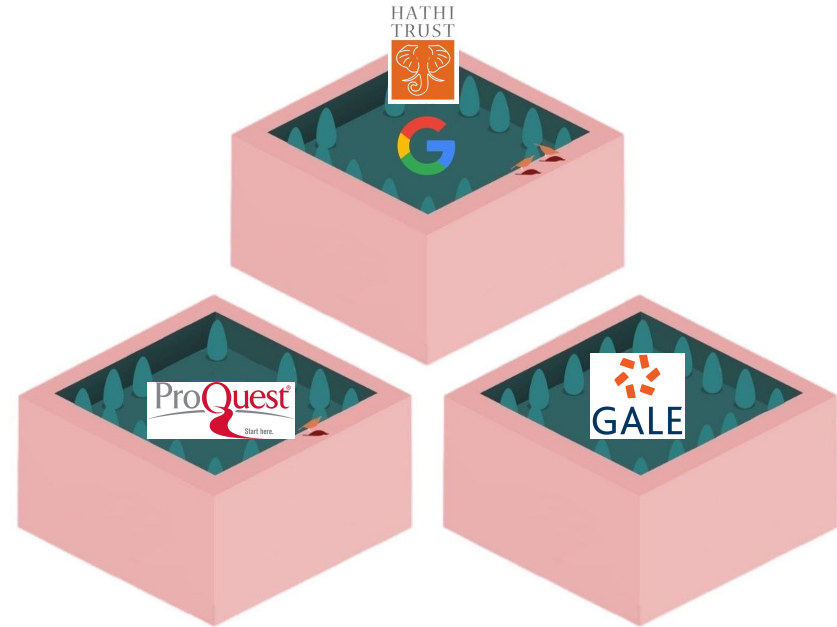
***Since* 1995–**

A **closed platform**, **walled garden**, or **closed ecosystem** is a software system wherein the carrier or service provider has <span style="color:pink">control</span> over applications, content, and/or media, and <span style="color:pink">restricts</span> convenient access to non-approved applicants or content. This is in contrast to an open platform, wherein consumers generally have unrestricted access to applications and content.

– *Wikipedia*, accessed 2024

THE CENTER
FOR DIGITAL
HUMANITIES
@PRINCETON

# The Digital Research Landscape: How We Got Here

- 2008 Google Books project & lawsuits → "propertization and privatization of culture" (Frosio, 2011)

- "Uneven digitization" of texts privilege works from the US and UK (Risam, 2019)

- Prohibited, costly, or unclear access for TDM since 2015 (McCracken & Raub, 2022)



THE CENTER
FOR DIGITAL
HUMANITIES
@PRINCETON

# The Digital Research Landscape: How We Got Here

"[Researchers] seek materials that are the most relevant to their theses, which likely come from multiple vendors. While acceptance of TDM by vendors is growing, one of the current problems with vendor-permitted TDM is an **uncooperative, siloed approach**."



Peter McCracken and Emma Raub, "Licensing Challenges Associated with Text and Data Mining," *JLSC* (2023)

# Limitations of Existing "Workflows"



Peter McCracken and Emma Raub, "Licensing Challenges Associated with Text and Data Mining," *JLSC* (2023)
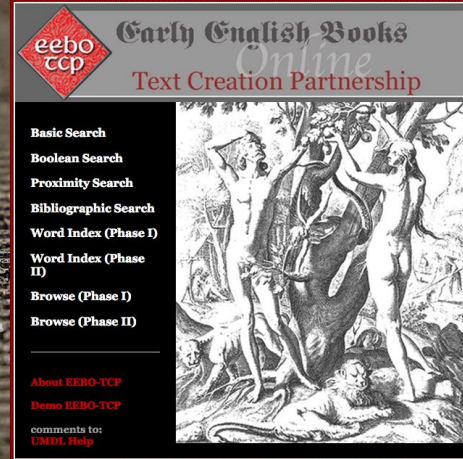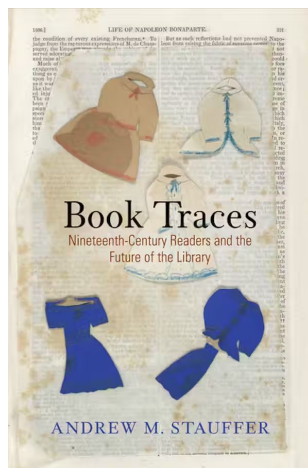
# Additional Challenges of TDM

- Inhibitors to TDM research:
  - Strong copyright laws, privacy, and anti-hacking laws that vary country to country
  - Researchers' lack of legal knowledge
  - Private corporation and government grip on data, which impedes access/sharing & disincentivizes them from improving data as technology improves
- Effects:
  - Lessens cross-border/international collaboration
  - Places onus on researcher to transform data into suitable format for TDM
  - Leads US-based researchers to avoid certain projects more frequently than EU-based researchers due to litigious nature of US, despite fair use clause

Patricia Aufderheide, Brandon Butler, and Kimberly Anastacio, "The Chilling Effects of Obstacles to Accessing, Using, and Sharing In-Copyright Data for Quantitative Research," *Information & Culture* (2024)

# Collaborative DH Projects

- Typically focus on one data source
- Often reinvent what libraries could be doing
- Digital editions and text corpora are both workarounds to the walled garden
- Front-end UX focus
- Where's the data?

# Case Study: The Princeton Prosody Archive

# The Princeton Prosody Archive - Overview



- Full-text searchable database of ~7,000 English-language works from 1559 to 1928 about the study of versification and pronunciation

# The Princeton Prosody Archive - Overview



- Includes "how-to" poetry handbooks, textbooks, grammar books, dictionaries, scholarly articles, elocution manuals, pronunciation guides…
- Curated into 6 collections

THE CENTER
FOR DIGITAL
HUMANITIES
@PRINCETON

# Data Sources

# Technical Architecture as Workflow

# "Dead Ends" & Roadblocks

# "Dead Ends" & Roadblocks



non-compete w/ Google!

# "Dead Ends" & Roadblocks

# "Dead Ends" & Roadblocks

# "Dead Ends" & Roadblocks

# The Path through the Hedge Maze

- Leverage the power of your institution
- Devise Memoranda of Understandings (MOUs)
- Cultivate relationships with librarians and vendors
- Prepare for change (in contacts & data)
- Create a culture of documentation
- Team continuity is key
- Be persistent advocating for your research needs

# Futures Beyond Walled Gardens

# Futures Beyond Walled Gardens

- Move toward Open Source options like Internet Archive, which just made ECCO available
- Generalizable code like the PPA that cuts across walled gardens
  - Partnership/Fellowship with HathiTrust so that others can create "PPAs" with different content?
- Cross-institutional leverage
- Cultural Data Collectives movement

# References

- Giancarlo F. Frosio, "Google Books Rejected: Taking the Orphans to the Digital Public Library of Alexandria," *Santa Clara Computer and High Technology Law Journal* 28, no. 1 (Nov. 2011): 81-141.
- Roopika Risam, *New Digital Worlds: Postcolonial Digital Humanities in Theory, Praxis, and Pedagogy* (Northwestern UP, 2019).
- Patricia McCracken and Emma Raub, "Licensing Challenges Associated With Text and Data Mining: How Do We Get Our Patrons What They Need?", *Journal of Librarianship and Scholarly Communication* 11 no. 1 (2023): 1-14, https://doi.org/10.31274/jlsc.15530.
- Patricia Aufderheide, Brandon Butler, and Kimberly Anastacio, "The Chilling Effects of Obstacles to Accessing, Using, and Sharing In-Copyright Data for Quantitative Research," *Information & Culture* 59, no. 1 (2024): 44-65, https://doi.org/10.7560/IC59103.
- "About the Collection," HathiTrust, https://www.hathitrust.org/the-collection/.
- Ted Underwood, Patrick Kimutis, and Jessica Witte, "NovelTM Datasets for English-Language Fiction, 1700-2009," *Journal of Cultural Analytics* 5, no. 2 (2020): 1-30, https://doi.org/10.22148/001c.13147.
- John A. Walsh et al., "'The Library is Open!': Open Data and an Open API for the HathiTrust Digital Library," *CHR2023: Computational Humanities Research Conference* (December 6–8, 2023, Paris, France), https://ceur-ws.org/Vol-3558/paper7875.pdf.

THE CENTER
FOR DIGITAL
HUMANITIES
@PRINCETON