

# Computational Analysis and Visualization of Electronic Records Collections

Visualization and Data Analysis Group  
Texas Advanced Computing Center, UT Austin  
Research Forum SAA, 2009

# Research Motivation

Digital tools for archivists and archives' users to:

*Make sense, explore, appraise, describe, discover, manage, and preserve electronic records collections of varied structures and formats*

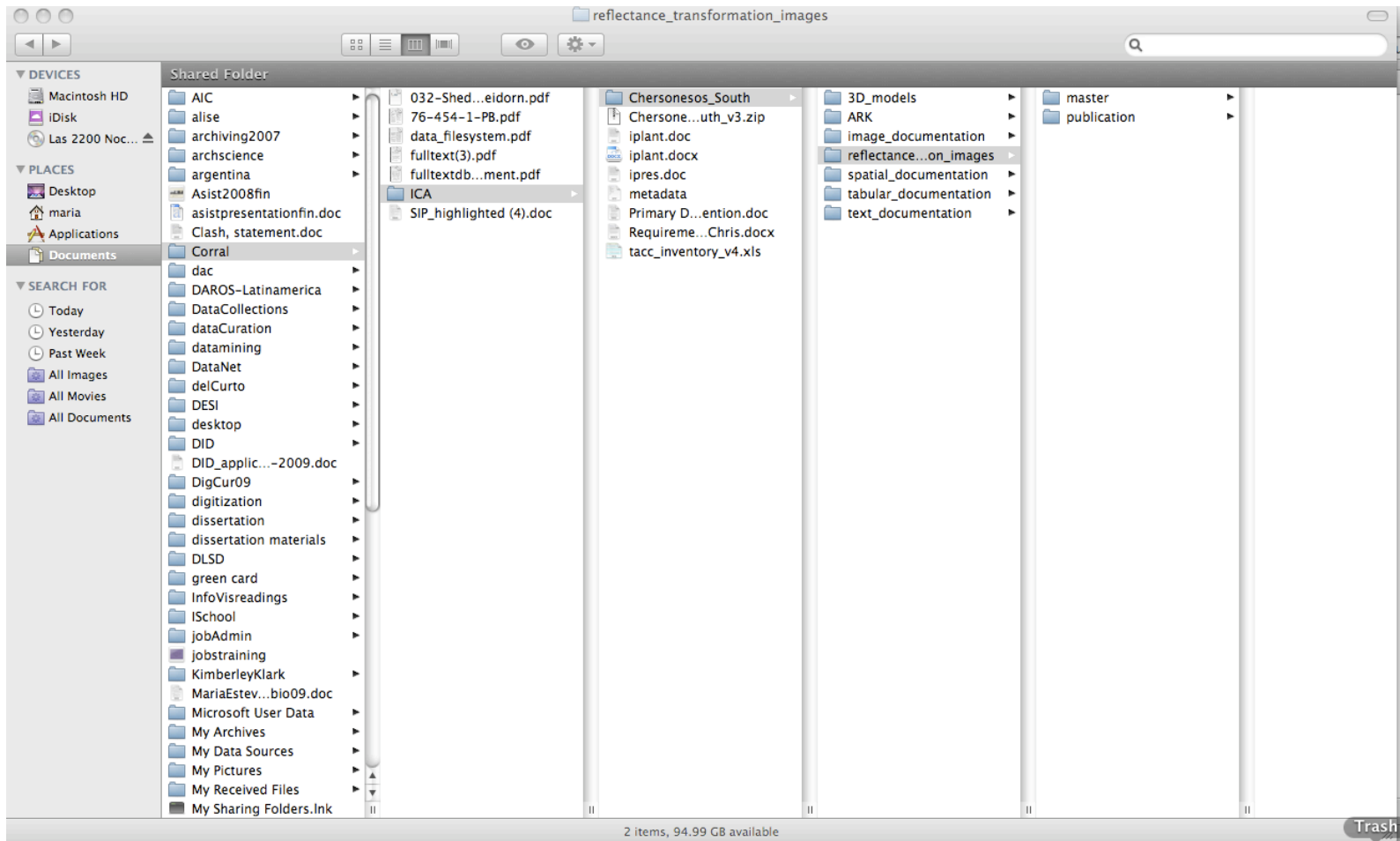
# Challenges

- Large archival data
- Minimal metadata
- What metadata can be extracted from electronic records collections?
- What can be inferred from the structure and the content of electronic records records?
- Abstract representations of collections: visual literacy development

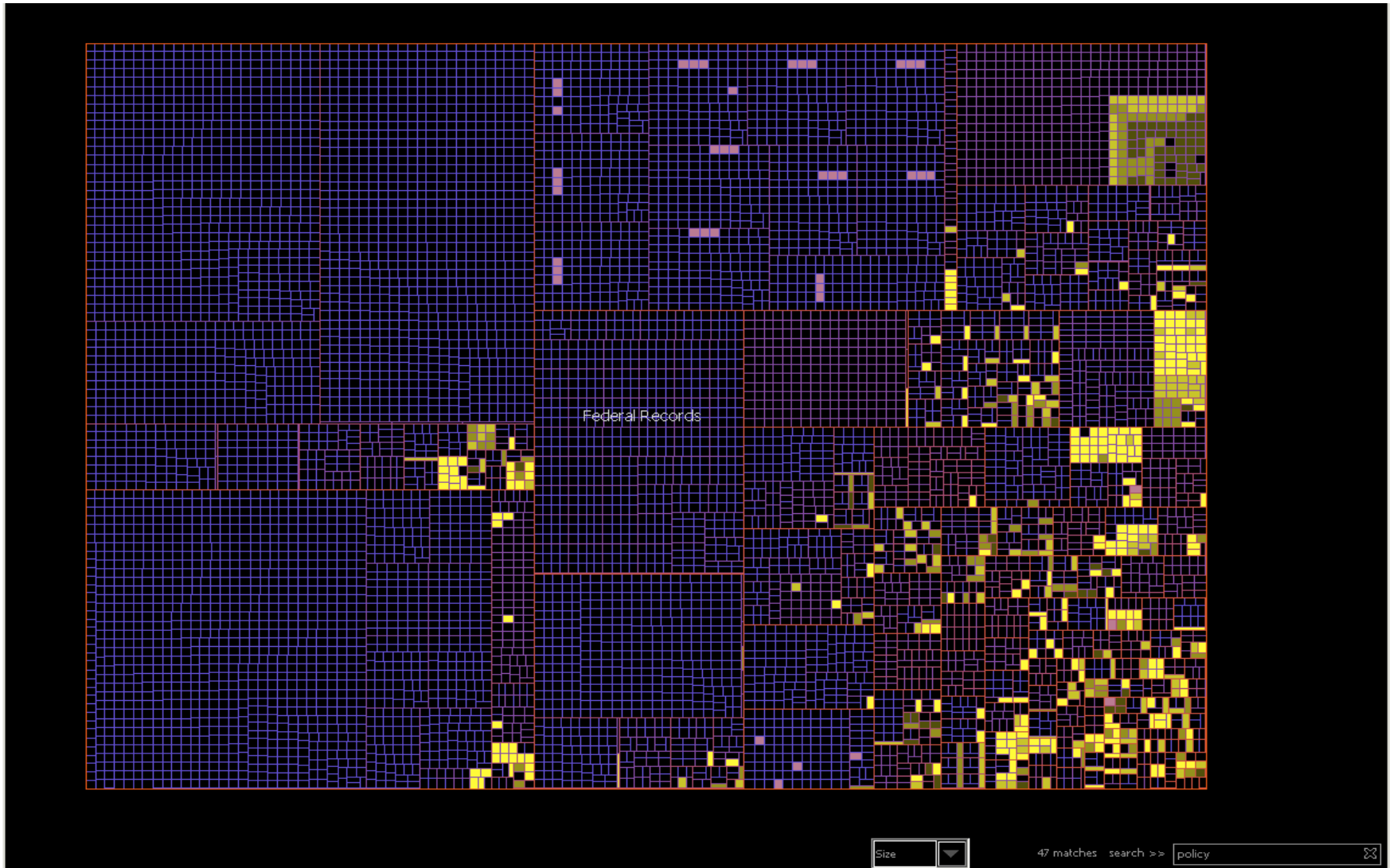
# I. Treemaps

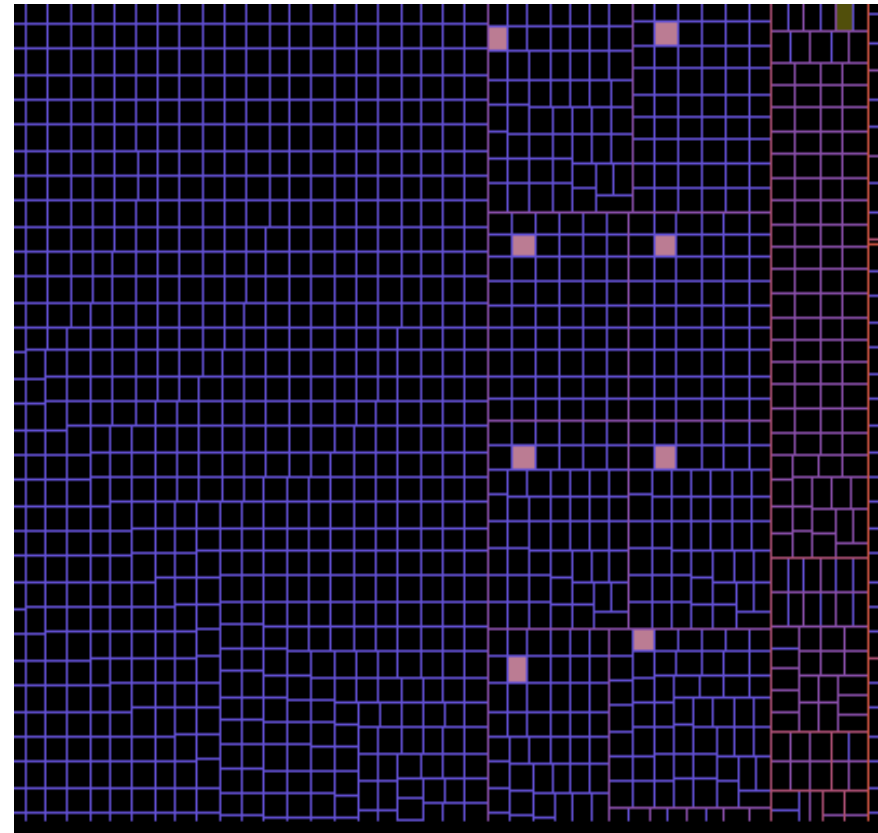
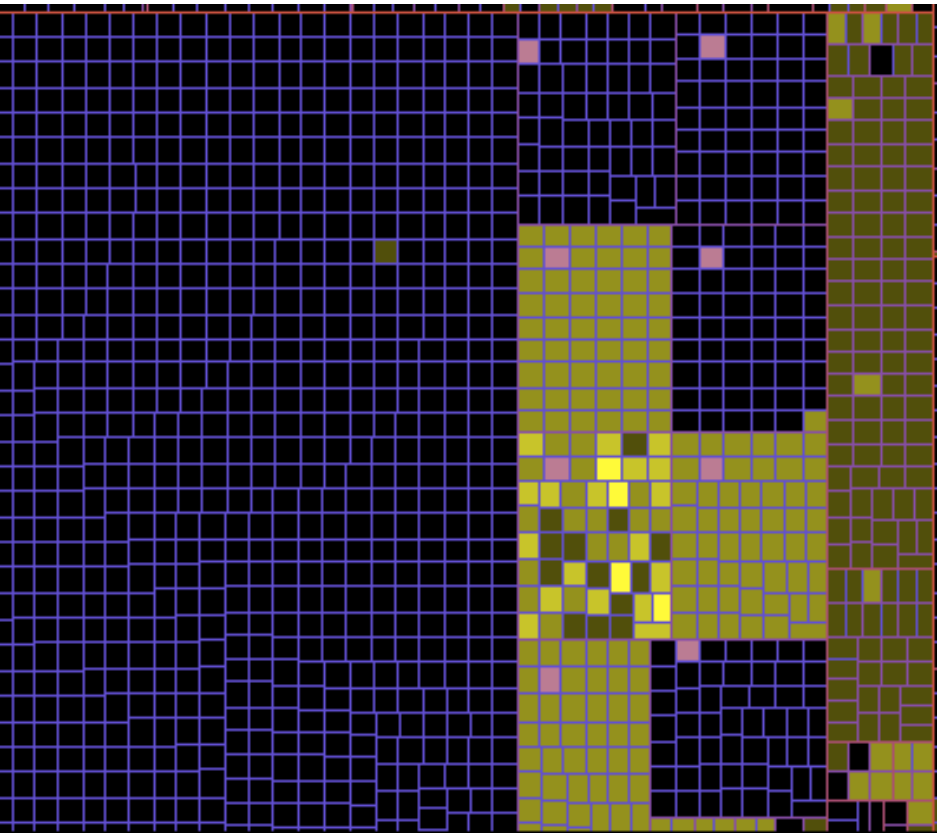
- Developed by Ben Shneiderman in the 1990's
- Adapted to visually scan large collections and to focus on smaller parts
- Collection properties are extracted and stored in a database
  - Structural, descriptive, technical
- Rendered through a treemap visualization application

# Directories view



# NARA test-bed collections in the Transcontinental Persistent Archives Prototype





## Records Group: Records of the Bureau of Census

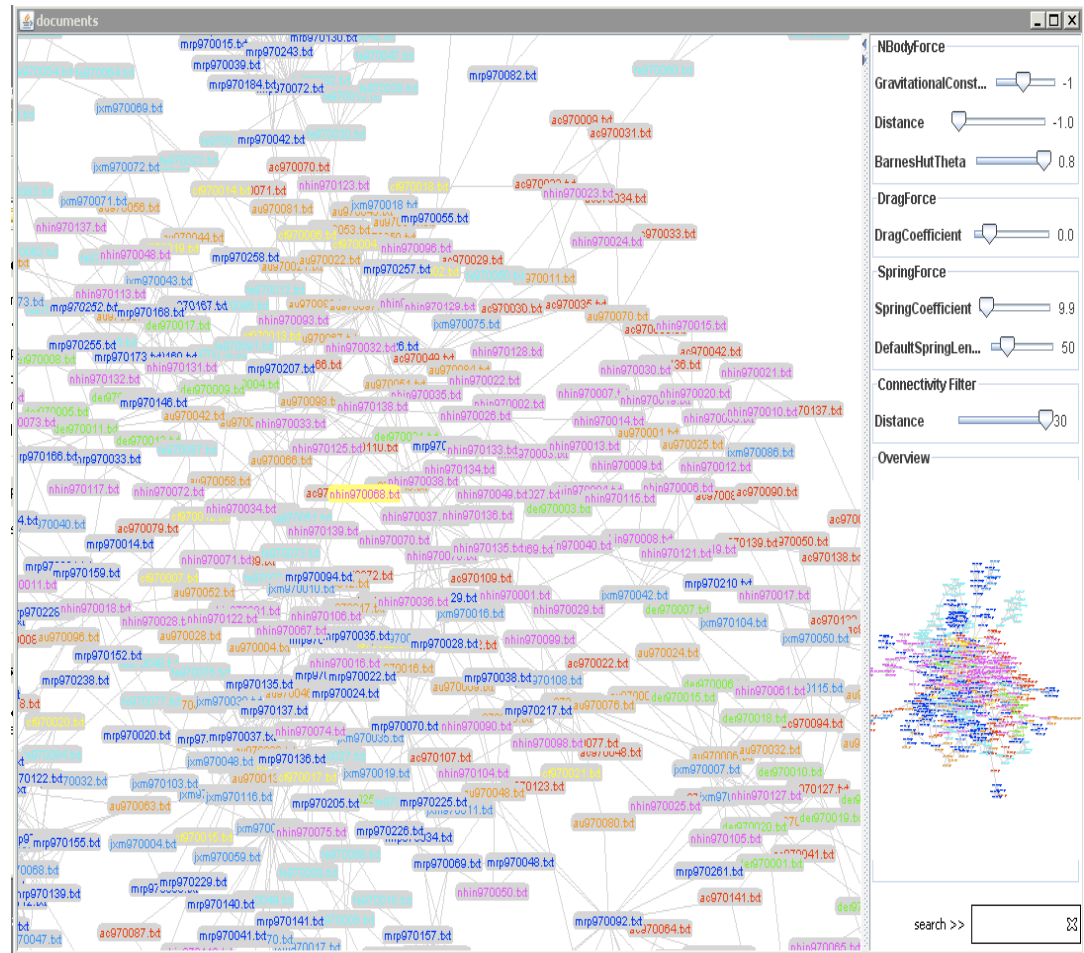
Smallest partition: one directory

Pink fill: searched keyword (2004) present in the name of the directory

Degrees of yellow-green-brown: from more to less number of files present in each directory / from more to less number of different file extensions present in each directory

# II. Paragraph alignment visualization

- Identify related records
  - Belonging to same activities, projects, transactions, events, etc.
  - Content based relationships





A	B	C	D		E
B					
C			C1	D1	D2
			C2	D3	D4
D					
E					

Compute relationships through bioinformatics-inspired method called “paragraph alignment”

- Unstructured collection of electronic records of different authors, topics, sizes with not much of an a-priori organization

# Stories

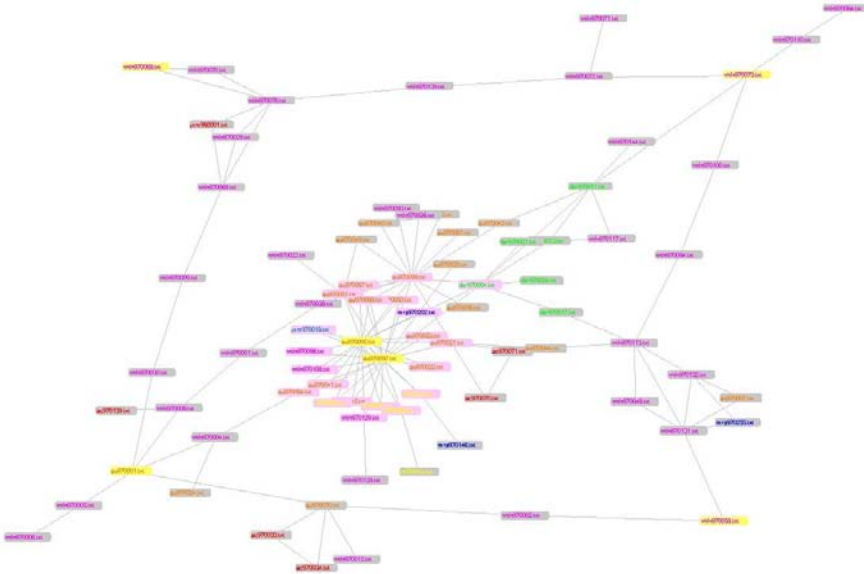
The image displays a network diagram of text files, overlaid on a background of several overlapping document windows. The diagram consists of nodes representing text files, connected by lines. The nodes are color-coded: yellow (e.g., au970001.txt, au970034.txt), purple (e.g., nhin970076.txt, nhin970029.txt), orange (e.g., nhin970008.txt, nhin970004.txt), and red (e.g., ac970034.txt, ac970033.txt). The background windows show snippets of text, including headers like 'nhin970058.txt', 'nhin970100.txt', 'nhin970139.txt', and 'nhin970151.txt'. The text in these windows appears to be a mix of administrative notes, meeting minutes, and correspondence, mentioning dates, names, and organizational details.

# Discovering context

- Explore and discover relationships between records and their authors

- Context

- Evidenced of cooperative writing
- Work-processes



# Conclusions

- Preliminary research
- Usability, display, interoperability and direct adaptation to archival tasks need to be resolved
  - EAD, JHOVE, PREMIS, METS
- Results stored in a database,
  - can be combined in myriad ways and with other tools to make abstractions, synthesis, and new discoveries
- Find useful/new visual metaphors