# Classifying the End of Term Web Archive

**Kathleen Murray & Cathy Hartman**
**University of North Texas, UNT Libraries, 1155 Union Circle #305190, Denton, TX 76203-5017**

## Archive Background

- Captured US government's Web presence:
  - September 2008 - November 2009
  - Transition between George W. Bush & Barack Obama
  - 16 TB
  - 160 million URLs (files/documents)

| Domains | # URLs | Subdomains |
|---|---|---|
| gov | 137,847,822 | 14,339 |
| mil | 3,555,425 | 1,677 |

## Problem

- Absence of descriptive metadata & Website classification thwarts discovery & access
- Standard file format (WARC - ISO 28500)
  - Specifies formats for storage, management, & exchange of data objects
  - Not designed for user access
- Access requires knowledge of a resource's URL (Wayback interface)

## Objective

To enable government information librarians to utilize existing selection practices to identify materials in the EOT Archive
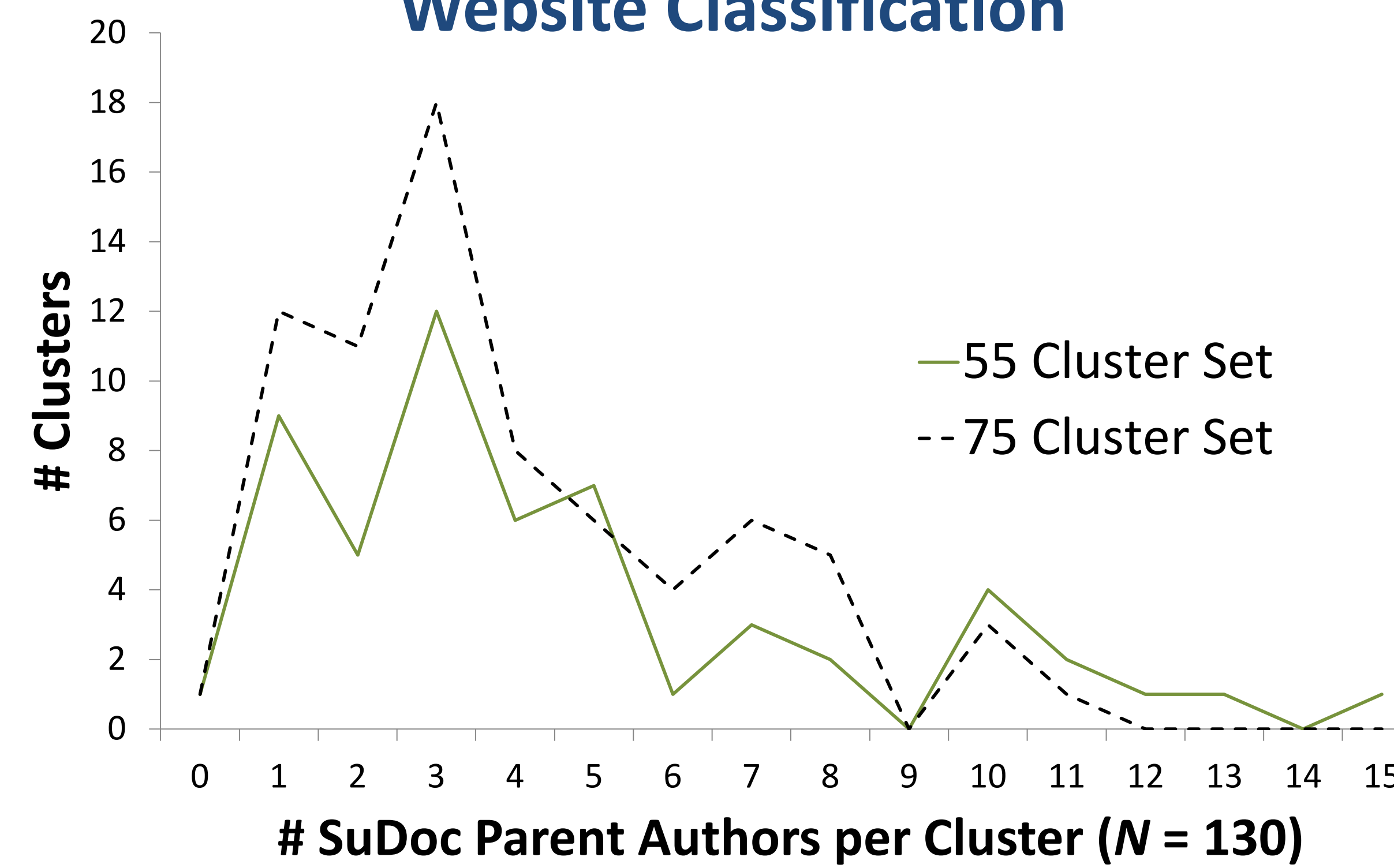
## Methods

- Cluster Identification (1,151 subdomains )
  - Linlog Coordinates with Agglomerative Hierarchical Clustering
  - Two sets of related Websites (55-set & 75-set)
- Website Classification (1,151 subdomains )
  - Government publication classification scheme: Superintendent of Documents (SuDocs) Classification Numbering System
  - 10 Subject Matter Experts (SMEs) & 3 arbitrators
- Cluster Tagging
  - 12 SMEs assigned subject terms to clusters
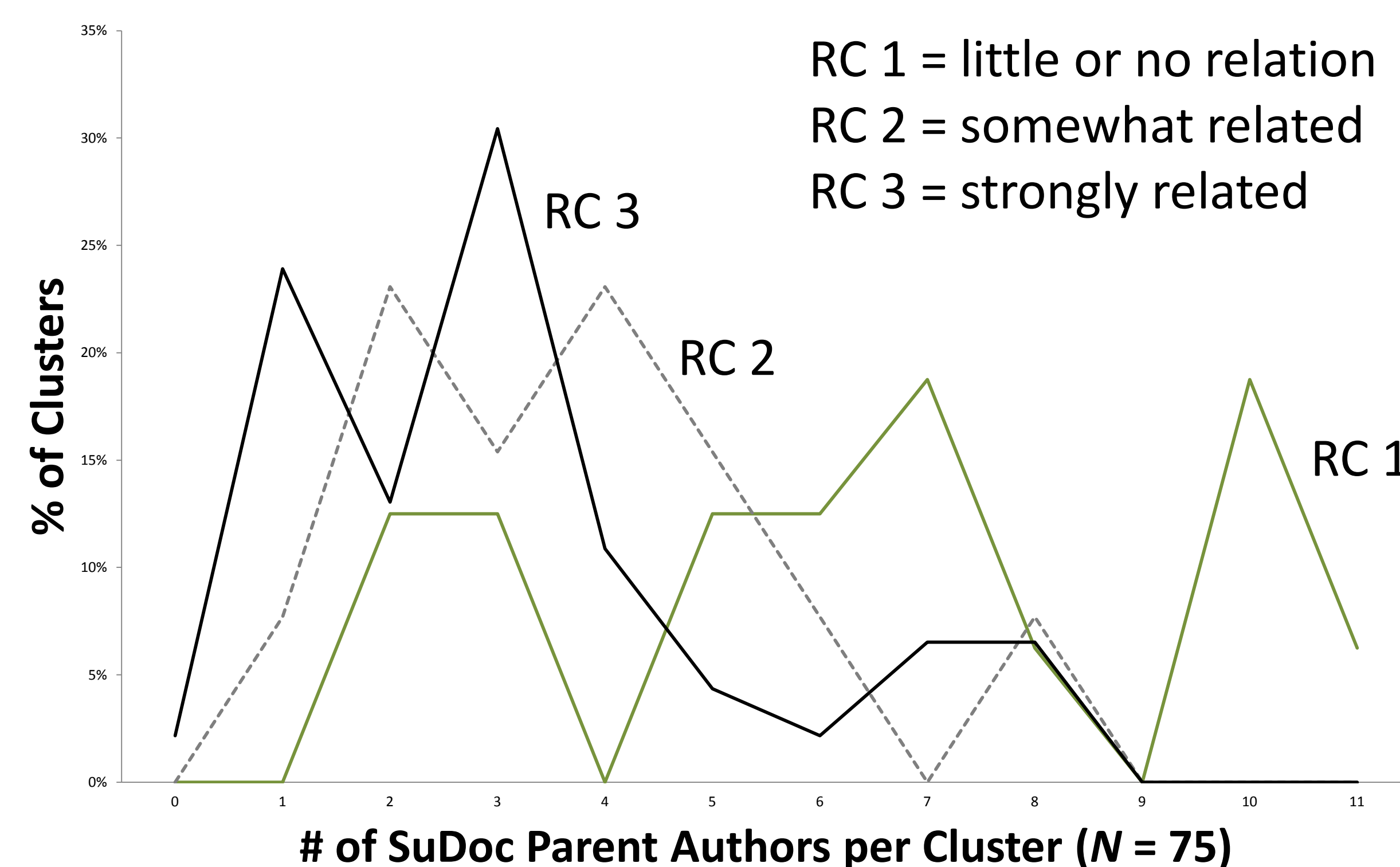
## Website Clusters

| | 55-Cluster Set | 75-Cluster Set |
|---|---|---|
| Identical Clusters | 39 | 39 |
| Unique Clusters | 16 | 36 |

### Website Classification



**# SuDoc Parent Authors per Cluster ($N$ = 130)**

| # Parent Authors | 55-Cluster Set | 75-Cluster Set |
|---|---|---|
| 1 | 16% | 16% |
| ≤ 2 | 27% | 32% |
| ≤ 4 | 60% | 67% |

### Topical Relatedness



RC 1 = little or no relation
RC 2 = somewhat related
RC 3 = strongly related

**# of SuDoc Parent Authors per Cluster ($N$ = 75)**

## Results

- Best clustering result with Linlog Coordinates with Agglomerative Hierarchical Clustering
- Overall: SuDocs Scheme worked well
  - Assigned SuDocs authors to 1,040 subdomains
- Cluster Analysis successfully identified strongly related subject content in the subdomains of 61% of clusters

| Clusters | # | RC 1 | RC 2 | RC 3 |
|---|---|---|---|---|
| Identical | 39 | 18% | 10% | **72%** |
| All | 130 | 21% | 18% | **61%** |
| | | | | |
| Unique in 75-Set | 36 | 22% | 14% | **64%** |
| Unique in 55-Set | 16 | 25% | 31% | **44%** |

- Identical clusters had the highest percentage of topically related subdomains (72%)
- Unique clusters had a substantially higher percentage of topically related subdomains after subdivision (64% v. 44%)

## Conclusions

- Cluster analysis holds promise for organizing Web archives into topically related groupings
- Involving SMEs in limited-scope classification activities may generate meaningful descriptive metadata for resources in focused Web archives
- Using the Web graph
  - How do we leverage the graph for identifying content?
- Describing the collection
  - How can we engage faculty with our Web archives?

## Acknowledgements

UNT Libraries