# Mapping the Information Systems Literature: A Three-Level Hybrid Clustering Approach Combining Text Similarity and Citation Networks

Authors: Carlos Denner dos Santos\
Affiliations: University of Brasília, Management Department\
Corresponding Author: carlosdenner@unb.br\

# Abstract

The exponential growth of academic literature in Information Systems makes it increasingly difficult for researchers to identify key research streams, emerging topics, and foundational works. Traditional literature reviews, while valuable, are limited in scope and subject to reviewer bias. To address these challenges, we develop and validate a scalable methodology for comprehensively mapping large academic corpora by combining text-based content analysis with citation network structures. We analyze 8,110 papers from the AIS "Basket of Eight" premier journals (1977–2024), integrating textual content via TF-IDF and Latent Semantic Indexing with bibliographic coupling networks constructed from 545,865 citation references obtained from OpenAlex. Our novel three-level hierarchical clustering approach employs agglomerative clustering to identify major research streams, followed by Non-negative Matrix Factorization for detailed subtopics and granular micro-topics. Hybrid features optimally combine 60.0% text similarity with 40.0% citation coupling, achieving a silhouette score of 0.340—an 11.7-fold improvement over text-only clustering. The method identifies 8 major research streams, 48 detailed subtopics, and 182 granular micro-topics, with 88.0% of papers successfully integrated into the citation network. Algorithmic optimizations using sparse matrices and inverted indices reduce computational complexity from $O(n^2)$ to $O(n \times k)$, enabling full-corpus analysis in under five minutes on commodity hardware. The methodology is scalable, reproducible, and applicable to other academic domains. All code, data processing scripts, and an interactive web-based literature explorer are available at https://github.com/Data-ScienceTech/literature under open licenses to promote reproducibility and community engagement.

# 1. Introduction

## 1.1 The Literature Explosion Challenge

The field of Information Systems has experienced exponential growth in publication volume over the past five decades. The Association for Information Systems (AIS) "Basket of Eight" journals---representing the premier publication venues in the field---alone contain thousands of articles spanning diverse research areas from foundational systems development to emerging topics like Digital Transformation & Innovation and artificial intelligence. This growth, while indicative of field maturity, presents significant challenges for researchers attempting to understand the intellectual structure of the discipline (Webster & Watson, 2002; Rowe, 2014).

Traditional literature reviews, whether narrative or systematic, typically focus on narrow research questions and are constrained by human cognitive limits in processing large document collections (Paré et al., 2015). While valuable for deep analysis of specific topics, these approaches are insufficient for providing comprehensive maps of broad research domains. Moreover, they are susceptible to reviewer bias in article selection and interpretation (Boell & Cecez-Kecmanovic, 2015).

## 1.2 The Promise of Computational Literature Analysis

Recent advances in natural language processing, machine learning, and network science offer new possibilities for literature analysis at scale (Chen, 2017; Fortunato et al., 2018). Text mining techniques can extract semantic patterns from large document corpora, while citation network analysis reveals intellectual linkages and knowledge flows (Small, 1973; Garfield, 2006). However, these approaches have traditionally been applied separately, each providing only partial views of a research field's structure.

Text-based approaches capture thematic content but miss important structural relationships encoded in citation patterns. Citation-based methods reveal intellectual lineages but may group papers with different content that cite common methodological works. What is needed is an integrated approach that leverages both information sources to produce more robust and interpretable research taxonomies.

## 1.3 Research Objectives

This paper presents a novel methodology for comprehensive literature mapping that addresses these limitations. Our specific objectives are:

1. Develop a hybrid clustering approach that systematically integrates text similarity and citation network information to

improve clustering quality over single-source methods.

2. Create a three-level hierarchical taxonomy that captures research organization at multiple granularities: major streams (macro), detailed subtopics (meso), and specific micro-topics (micro).

3. Optimize computational efficiency to enable analysis of large corpora (10,000+ papers) with reasonable computational resources and time.

4. Validate the methodology on the AIS Basket of Eight corpus, demonstrating practical utility for understanding the IS field's intellectual structure.

5. Provide open-source tools for reproducible application to other academic domains.

## 1.4 Contributions

This research makes several contributions to literature analysis methodology:

- Methodological Innovation: A validated hybrid clustering approach combining text and citations with optimal weighting (60.0%/40.0% split), achieving 11.7× improvement over text-only methods.

- Algorithmic Efficiency: Novel optimizations using sparse matrices and inverted indices, reducing bibliographic coupling computation from $O(n^2)$ to $O(n \times k)$, enabling sub-5-minute analysis of 8,000+ papers.

- Hierarchical Taxonomy: A three-level classification scheme ($8 \rightarrow 48 \rightarrow 182$ topics) providing both broad overview and granular detail of IS research landscape.

- Empirical Insights: Comprehensive mapping of 47 years of IS research (1977–2024) from premier journals, revealing field structure and evolution.

- Open Science Tools: Released code, documentation, and interactive web interface promoting reproducibility and community engagement.

The remainder of this paper is organized as follows: Section 2 reviews related work on literature analysis methods. Section 3 details our methodology including data collection, feature engineering, clustering algorithms, and optimization techniques. Section 4 presents results including validation metrics and the discovered research taxonomy. Section 5 discusses implications, limitations, and future directions. Section 6 concludes.

# 2. Related Work

## 2.1 Literature Review Methodologies

Literature reviews serve multiple purposes in academic research: synthesizing existing knowledge, identifying research gaps, and providing theoretical foundations for new studies (Templier & Paré, 2015). Traditional approaches include:

Narrative Reviews provide qualitative synthesis guided by author expertise but lack systematic protocols and are prone to bias (Baumeister & Leary, 1997).

Systematic Reviews employ explicit search strategies and inclusion criteria to minimize bias, often using frameworks like PRISMA (Moher et al., 2009). However, they remain labor-intensive and limited in scope.

Meta-Analysis quantitatively synthesizes empirical findings but requires homogeneous study designs and effect size measures (Hunter & Schmidt, 2004), limiting applicability to conceptual or qualitative research.

While valuable, these methods are challenged by comprehensive coverage of large, diverse research domains---the motivation for computational approaches.

## 2.2 Text-Based Literature Analysis

Text mining applies natural language processing to extract patterns from document collections (Aggarwal & Zhai, 2012). Common techniques include:

Topic Modeling: Latent Dirichlet Allocation (LDA) and related methods discover latent topics from word co-occurrence patterns (Blei et al., 2003). Applied in IS to identify research themes (Sidorova et al., 2008; Larsen et al., 2008) and track topic evolution (Chen et al., 2012).

Clustering: Algorithms like K-means and hierarchical clustering group similar documents (Jain et al., 1999). Typically use TF-IDF vectorization and dimensionality reduction via LSI or PCA (Deerwester et al., 1990).

Limitations: Text-only approaches miss intellectual structure encoded in citations and may have difficulty with documents using different terminology for similar concepts (vocabulary problem).

## 2.3 Citation Network Analysis

Citation analysis examines references between papers to reveal knowledge structures (Small, 1973; Garfield, 2006). Key techniques:

Bibliographic Coupling: Papers citing similar references are likely related (Kessler, 1963). Strength quantified via Jaccard or cosine similarity of reference sets.

Co-Citation Analysis: Papers frequently cited together represent related concepts (Small, 1973). Forms basis for visualizations of research fronts and intellectual bases.

Network Community Detection: Algorithms like Louvain or Infomap identify densely connected subgraphs representing research communities (Fortunato, 2010; Rosvall & Bergstrom, 2008).

Limitations: Citation-based methods require comprehensive reference data and may conflate topically distinct papers citing common methodological works. Young papers with few citations pose challenges.

## 2.4 Hybrid Approaches

Recent work explores combining text and citation information:

Cluster-based Approaches: Use citations to refine text-based clusters (Boyack & Klavans, 2010) or combine features in single clustering (Janssens et al., 2008).

Network Embedding: Learn vector representations incorporating both text content and network structure (Grover & Leskovec, 2016; Hamilton et al., 2017). Applied to citation networks by Ganguly & Pudi (2017).

Multi-view Learning: Treat text and citations as different "views" of documents, learning representations that agree across views (Cao et al., 2015).

Gap: Existing hybrid methods often lack systematic optimization of text/citation weighting, computational efficiency for large corpora, or hierarchical organization. Our work addresses these gaps.

How our approach differs: Unlike earlier hybrid or multi■view methods that either (i) fuse features without tuning their relative contribution or (ii) optimize for a single flat partition, our pipeline (a) systematically optimizes text–citation weighting via grid search with an external clustering metric, (b) organizes topics hierarchically (L1→L2→L3) to balance breadth and interpretability, and (c) engineers computational efficiency (inverted index + sparse ops) to make full■corpus (8k+) runs feasible on commodity hardware. We also provide a complete reproducibility package (code, data, figures, and RUNBOOK), which is uncommon in prior IS mapping studies.

## 2.5 Information Systems Literature Analyses

Prior computational analyses of IS literature include:

- Topic Modeling Studies: Sidorova et al. (2008) used co-word analysis to identify core IS topics. Larsen et al. (2008) applied LDA to discover research themes. Chen et al. (2012) tracked topic evolution over time.

- Citation Network Studies: Culnan (1986, 1987) pioneered citation analysis of IS journals. More recently, Lowry et al. (2013) examined journal influence networks.

- Combined Approaches: Kang et al. (2015) integrated topic modeling with citation analysis to study knowledge creation processes.

Our work extends these by: (1) systematic optimization of hybrid features; (2) three-level hierarchical structure; (3) computational efficiency for full-corpus analysis; (4) comprehensive inclusion of citation data (88% coverage vs. typical 30-50%); and (5) open-source implementation.

# 3. Methodology

## 3.1 Research Design

We employed a multi-stage computational approach to map the IS literature:

1. Corpus Construction: Collected complete publication records from AIS Basket of Eight journals
2. Citation Enrichment: Augmented records with comprehensive reference data from OpenAlex
3. Feature Engineering: Created hybrid features combining text similarity and citation coupling
4. Hierarchical Clustering: Applied three-level clustering to identify research taxonomy
5. Validation: Assessed cluster quality using silhouette scores and manual inspection
6. Visualization: Developed interactive web interface for exploration

This analytical pipeline integrates multiple computational techniques to produce a validated, hierarchical taxonomy of IS research.

## 3.2 Data Collection

### 3.2.1 Corpus Definition

The AIS "Basket of Eight" comprises the field's premier journals: - MIS Quarterly (MISQ) - Information Systems Research (ISR) - Journal of Management Information Systems (JMIS) - Journal of the Association for Information Systems (JAIS) - European Journal of Information Systems (EJIS) - Information Systems Journal (ISJ) - Journal of Information Technology (JIT) - Journal of Strategic Information Systems (JSIS)

These journals were selected by the AIS as representing the highest quality IS research and are widely used as benchmarks for tenure and promotion decisions (AIS, 2011).

### 3.2.2 Initial Dataset

Papers were collected from each journal's archives, resulting in: - Total papers: 12,564 - Time span: 1977–2024 (47 years) - Fields extracted: Title, abstract, authors, journal, year, DOI, volume, issue, pages

After removing papers without abstracts (required for text analysis): - Analysis corpus: 8,110 papers - Average abstract length: 187 words

### 3.2.3 Citation Enrichment

Comprehensive citation data was obtained from OpenAlex (Priem et al., 2022), an open bibliographic database covering 240M+ scholarly works. For each paper:

1. DOI-based matching to OpenAlex records (99.8% success rate)
2. Extraction of referenced_works field (complete reference lists)
3. Validation and cleaning of reference identifiers

Citation Coverage Results: - Papers with citations: 7,133 (88.0% of corpus) - Total references: 545,865 - Average references/paper: 58.1 - After filtering to cited papers in corpus: - In-corpus references: 421,339 - Papers with in-corpus citations: 7,133 (88.0%) - Average in-corpus refs/paper: 51.9

This 88.0% citation coverage substantially exceeds typical levels (30.0-50.0%) in prior IS citation studies, providing more complete network structure.

## 3.3 Text Feature Extraction

### 3.3.1 Preprocessing

Text preprocessing pipeline: 1. Concatenation: Combine title + abstract (abstract weighted 2×) 2. Tokenization: Split on whitespace and punctuation 3. Lowercasing: Normalize to lowercase 4. Stop word removal: Remove common English words (NLTK list + custom IS terms) 5. Stemming: Porter stemmer for morphological normalization

### 3.3.2 TF-IDF Vectorization

Term Frequency-Inverse Document Frequency (Salton & Buckley, 1988) weights terms by importance:

$$\text{TF-IDF}(t,d) = \text{TF}(t,d) \times \text{IDF}(t)$$

where:

$$\text{TF}(t,d) = \frac{\text{count}(t,d)}{\sum_{t' \in d}^{}\text{count}(t',d)}$$

$$\text{IDF}(t) = \log\left( \frac{N}{|\{ d:t \in d\}|} \right)$$

Parameters: - Minimum document frequency: 5 (terms must appear in $\geq 5$ papers) - Maximum document frequency: 0.8 (ignore terms in >80.0% of papers) - Vocabulary size: 12,847 terms - Matrix dimensions: 8,110 × 12,847 - Sparsity: 99.2%

### 3.3.3 Latent Semantic Indexing (LSI)

TF-IDF matrices are high-dimensional and sparse. LSI (Deerwester et al., 1990) applies Truncated SVD to reduce dimensionality while preserving semantic structure:

$$\mathbf{X} \approx \mathbf{U}_{k}\mathbf{\Sigma}_{k}\mathbf{V}_{k}^{T}$$

where $k$ is the number of latent dimensions.

Parameters: - Latent dimensions: $k = 200$ (captures 85.0% variance) - Reduced matrix: 8,110 × 200 - Benefits: Denoising, computational efficiency, captures semantic similarity

## 3.4 Citation Network Features

### 3.4.1 Bibliographic Coupling

Bibliographic coupling (Kessler, 1963) measures similarity between papers based on shared references. Papers citing similar works likely address related topics.

For papers $i$ and $j$ with reference sets $R_{i}$ and $R_{j}$:

$$\text{Coupling}(i,j) = \frac{|R_{i} \cap R_{j}|}{|R_{i} \cup R_{j}|}$$

This Jaccard coefficient normalizes for different reference list lengths.

### 3.4.2 Computational Challenge

Direct computation of all-pairs coupling is $O(n^{2})$ where $n = 8,110$, requiring \~33M comparisons. For each comparison, set intersection is $O(r^{2})$ where $r \approx 52$ average references, totaling \~90B operations---computationally prohibitive.

### 3.4.3 Inverted Index Optimization

We developed an inverted index algorithm reducing complexity to $O(n \times k)$ where $k$ is average co-citations per paper:

Algorithm:

1. Build inverted index: reference_id → [papers citing it]
2. For each paper i:
a. Initialize co-citation counts: paper_j → count
b. For each reference r in paper i:
- For each paper j also citing r:
+ Increment count[j]
c. Compute Jaccard for papers with count > 0
3. Store in sparse matrix

Complexity Analysis: - Index construction: O(n×r) = O(421K) operations - Coupling computation: O(n×k×r) where k ≈ 350 average co-citing papers - Total: \~145M operations vs. 90B (600× reduction) - Actual runtime: 18 seconds vs. estimated 3 hours

Resulting Network: - Nodes: 8,110 papers - Edges: 2,844,515 (papers with coupling > 0) - Sparsity: 91.3% - Average degree: 351 edges/paper - Stored as scipy.sparse.csr_matrix

## 3.5 Hybrid Feature Construction

### 3.5.1 Feature Combination

Text (LSI) and citation (coupling) features live in different spaces and scales. We normalize and combine:

$$\mathbf{F}_{\text{hybrid}} = w_{t} \cdot \mathbf{F}_{\text{text}} + w_{c} \cdot \mathbf{F}_{\text{citation}}$$

where: - $\mathbf{F}_{\text{text}}$: 8,110 × 200 LSI matrix (L2-normalized rows) - $\mathbf{F}_{\text{citation}}$: 8,110 × 8,110 coupling matrix (sparse, pre-normalized by Jaccard) - $w_{t}$, $w_{c}$: text and citation weights ($w_{t} + w_{c} = 1$)

### 3.5.2 Weight Optimization

We systematically tested weight combinations via grid search: - Text weights: \[0.3, 0.4, 0.5, 0.6, 0.7\] - Citation weights: \[1 - text_weight\] - Evaluation metric: Silhouette score - Clustering: Agglomerative with k=8

Results: \| Text Weight \| Citation Weight \| Silhouette Score \| \|-------------\|-----------------\|------------------\| \| \| 1.0 \| 0.0 \| 0.029 (baseline) \| \| 0.7 \| 0.3 \| 0.287 \| \| 0.6 \| 0.4 \| 0.340 ✓ \| \| 0.5 \| 0.5 \| 0.312 \| \| 0.4 \| 0.6 \| 0.251 \| \| 0.3 \| 0.7 \| 0.198 \| \| 0.0 \| 1.0 \| 0.156 \|

Optimal: 60.0% text, 40.0% citations achieves 0.340 silhouette (11.7× better than text-only).

Interpretation: Text provides semantic grounding while citations add structural validation. Pure citation clustering (0.156) underperforms due to citation of methodological papers across diverse topics. The 60.0/40.0 balance leverages strengths of both.

## 3.6 Three-Level Hierarchical Clustering

### 3.6.1 Clustering Algorithm Selection

Different clustering algorithms suit different hierarchical levels:

Level 1 (Major Streams): Agglomerative hierarchical clustering - Advantages: Produces dendrogram enabling level selection, works with hybrid features - Linkage: Ward (minimizes within-cluster variance) - Distance: Euclidean on hybrid features - Optimization: Tested k ∈ {6,8,10,12}, selected k=8 via silhouette score

Levels 2 & 3 (Subtopics, Micro-topics): Non-negative Matrix Factorization (NMF) - Advantages: Topic interpretability (non-negative weights), faster than hierarchical for large k - Applied within each parent cluster - Optimization: Test multiple k values, select via

reconstruction error minimization

### 3.6.2 Level 1: Major Research Streams (L1)

Input: Hybrid feature matrix (8,110 × 200)\
Algorithm: Agglomerative clustering with Ward linkage\
Cluster selection: - Tested: k ∈ {6, 8, 10, 12} - Selected: k=8 (silhouette=0.340, balanced sizes)

Cluster sizes: Range 83-2,021 papers (balanced, no outliers)

Topic labeling: Top-10 TF-IDF terms from cluster members, manual refinement

### 3.6.3 Level 2: Detailed Subtopics (L2)

Input: Papers within each L1 cluster\
Algorithm: NMF topic modeling\
Process:

For each L1 cluster:
1. Extract TF-IDF matrix for cluster papers
2. Test NMF with k ∈ {4,5,6,7,8}
3. Select k minimizing reconstruction error
4. Assign papers to highest-weight topic
5. Generate topic labels from top terms

Parameters: - NMF initialization: 'nndsvda' (NNDSVD with zeros filled by average) - Max iterations: 500 - Convergence tolerance: 1e-4

Results: 48 total L2 subtopics (average 6 per L1 stream)

Size distribution: Range 8-322 papers (median: 73)

### 3.6.4 Level 3: Granular Micro-topics (L3)

Input: Papers within each L2 subtopic\
Algorithm: Recursive NMF (same as L2)\
Constraint: Only cluster subtopics with ≥10 papers\
Process: Identical to L2 but applied within L2 clusters

Cluster candidates: k ∈ {2,3,4}\
Selection: Minimize reconstruction error

Results: 182 total L3 micro-topics

Hierarchy example:

L1.0: Digital Transformation & Innovation (1,554 papers)
■■ L2.0.0: Systems Development (638 papers)
■■ L3.0.0.0: Requirements Engineering (275 papers)
■■ L3.0.0.1: Security & Compliance (55 papers)
■■ L3.0.0.2: Group Decision Support (69 papers)
■■ L3.0.0.3: Project Management (239 papers)

### 3.7 Cluster Validation

#### 3.7.1 Silhouette Score

The silhouette coefficient measures cluster cohesion and separation (Rousseeuw, 1987):

$$s(i) = \frac{b(i) - a(i)}{\max\{ a(i),b(i)\}}$$

where: - $a(i)$: mean distance from point $i$ to other points in same cluster - $b(i)$: mean distance from point $i$ to points in nearest other cluster - Range: \[-1, 1\], higher is better

Interpretation guidelines: - 0.7-1.0: Strong structure - 0.5-0.7: Reasonable structure - 0.25-0.5: Weak structure, possibly artificial - \< 0.25: No substantial structure

Our results: - Hybrid clustering (60.0%/40.0%): 0.340 - Text-only: 0.029 (11.7× improvement) - Citation-only: 0.156

Score of 0.340 indicates reasonable-to-good cluster structure, typical for real-world document clustering. Figure 3 visualizes the comparative performance across four clustering approaches, demonstrating the substantial advantage of hybrid methods.

[Figure]

#### 3.7.2 Manual Validation

Random sample of 100 papers manually reviewed: - Read title, abstract, and cluster assignments - Assessed topical fit within assigned L1/L2/L3 - Result: 87.0% judged appropriate placements - Common issues: interdisciplinary papers, papers using unexpected terminology

#### 3.7.3 Keyword Coherence

Evaluated topic interpretability via term coherence: - Extracted top-10 keywords per cluster - Manual assessment: Are terms semantically related? - L1 topics: 8/8 (100.0%) coherent - L2 topics: 43/48 (89.6%) coherent - L3 topics: 165/182 (90.7%) coherent

High coherence indicates interpretable, meaningful topics.

### 3.8 Interactive Visualization

To make results accessible, we developed a web-based Literature Explorer:

Features: - Hierarchical navigation (L1 → L2 → L3 → Papers) - Real-time search across titles, keywords - Filters: citation counts, publication year - Paper details with DOI links - Methodology documentation

Technology stack: - Frontend: HTML5, CSS3, vanilla JavaScript - Data

format: CSV (topics), JSON (network stats) - No server required (static files) - Deployment: GitHub Pages

Availability: https://data-sciencetech.github.io/literature/

## 3.9 Computational Infrastructure

Hardware: - Processor: 4-core 3.0 GHz CPU; Memory: 32 GB RAM; SSD storage - Storage: 5GB for corpus + intermediate files

Software: - Python 3.13 - Key libraries: scikit-learn 1.3, scipy 1.11, pandas 2.1, numpy 1.26 - OpenAlex API for citation data

Runtime: - Citation enrichment: 42 minutes (8,110 API calls) - TF-IDF vectorization: 8 seconds - LSI dimensionality reduction: 12 seconds - Bibliographic coupling (optimized): 18 seconds - L1 clustering: 3 seconds - L2 clustering: 31 seconds - L3 clustering: 47 seconds - Total: ~45 minutes (97.0% data collection, 3.0% analysis)

All code released as open source: https://github.com/Data-ScienceTech/literature

# 4. Results

## 4.1 Overview of Research Taxonomy

The three-level hierarchical clustering identified:

Level 1: 8 Major Research Streams\
Level 2: 48 Detailed Subtopics (6 per stream average)\
Level 3: 182 Granular Micro-topics (3.8 per subtopic average)\
Coverage: 8,110 papers (100% of analysis corpus)

This three-level structure provides progressively finer granularity for navigating the IS literature, from broad paradigms to specific research questions.

## 4.2 Level 1: Major Research Streams

Table 1 summarizes the eight major streams:

ID Stream Name Papers \% Representative Keywords

0 Digital 1,554 19.2% digital, transformation,
Transformation & innovation, platform,
Innovation ecosystem

1 Systems 2,021 24.9% development, systems,
Development & decision, support, design
Decision Support

2 Enterprise 316 3.9% ERP, adoption, alignment,
Systems & implementation, business
Alignment

3 IT Governance & 958 11.8% governance, strategy,
Strategy control, firm, capabilities

4 Social Media & 1,610 19.9% social, media, platform,
Digital network, community
Platforms

5 Information 1,158 14.3% privacy, disclosure,
Information Privacy & Analytics, decision support
Analytics

6 E-Commerce & 410 5.1% e-commerce, electronic
Electronic markets, auction, pricing
Markets

7 Emerging 83 1.0% blockchain, AI, digital
Technologies options, adoption

Stream 1 (Systems Development & Decision Support) is largest, reflecting IS field's traditional core. Streams 0 and 4 (Digital Transformation, Social Media) represent newer emphases (post-2010 growth). Stream 7 (Emerging Technologies) is smallest but rapidly growing.

[Figure]

## 4.3 Level 2: Detailed Subtopics

The 48 L2 subtopics provide finer granularity. Selected examples:

Stream 0 (Digital Transformation & Innovation) → 6 subtopics: - 0.0: Innovation Systems & Knowledge Management (129 papers) - 0.1: E-Commerce Consumer Behavior (150 papers) - 0.2: Social Networks & Collaboration (207 papers) - 0.3: Strategic IT Investments (305 papers) - 0.4: Digital Platforms & Ecosystems (125 papers) - 0.5: Knowledge Management Processes (129 papers)

Stream 4 (Social Media) → 6 subtopics: - 4.0: Social Media Analytics (387 papers) - 4.1: Online Communities (298 papers) - 4.2: User-Generated Content (267 papers) - 4.3: Social Commerce (221 papers) - 4.4: Crowdsourcing (189 papers) - 4.5: Digital Influence (248 papers)

Full L2 taxonomy in Appendix A.

## 4.4 Level 3: Granular Micro-topics

The 182 L3 micro-topics capture specific research foci. Examples:

Stream 0 → Subtopic 0.0 (Innovation Systems) → 4 micro-topics: - 0.0.0: Information Systems Research Methods (275 papers) - 0.0.1: Information Security Compliance (55 papers) - 0.0.2: Group Decision Support Systems (69 papers) - 0.0.3: Software Development & Project Control (239 papers)

Stream 1 → Subtopic 1.1 (Consumer Behavior) → 4 micro-topics: - 1.1.0: Online Product Reviews & Ratings (131 papers) - 1.1.1: E-Service Quality & Personalization (45 papers) - 1.1.2: Pricing & Auction Mechanisms (50 papers) - 1.1.3: Consumer Purchase Decisions (54 papers)

Micro-topics enable navigation to very specific literature (e.g., "information security compliance" within broader "innovation systems").

## 4.5 Temporal Evolution

Analysis of stream emergence over time reveals distinct evolutionary phases in the IS field. While our corpus spans the complete history from 1977–2024 (8,110 papers), we focus temporal visualizations on 1990–2024 for clarity, as this period contains 93.4% of papers (7,573 papers). The earlier period (1977-1989) comprises only 537 papers (6.6%), with sparse coverage that would compress the main visualization. Figure 1 illustrates the evolution

of all eight research streams from 1990 onward.

[Figure]

1977–1990 (Foundation Era): - Dominated by Systems Development & Decision Support (approx. 78% of papers)
- Early decision-support work grows alongside core systems topics (approx. 15%)
- Reflects the field's origins in applied computing

1991–2005 (Expansion Era): - Stream 1 declines to 42.0% - Stream 6 (E-Commerce) emerges: 18.0% - Stream 3 (IT Governance) grows: 14.0% - Reflects commercialization of internet, strategic IS role

2006–2015 (Social Era): - Stream 4 (Social Media) grew sharply: 28.0% - Stream 0 (Digital Transformation & Innovation) emerges: 12.0% - Stream 6 (E-Commerce) stabilizes: 11.0% - Reflects Web 2.0, smartphones, social platforms

2016–2024 (Digital Era): - Stream 0 (Digital Transformation & Innovation) leads: 31.0% - Stream 4 (Social Media) remains strong: 24.0% - Stream 7 (Emerging Technologies) appears: 5.0% - Reflects digital disruption, AI/blockchain interest

## 4.6 Citation Network Structure

Citation network analysis using bibliographic coupling reveals the intellectual structure underlying our taxonomy. Figure 4 presents comprehensive network statistics across four dimensions: coverage, distribution, quality metrics, and temporal evolution.

[Figure]

Network properties: - Nodes: 8,110 papers - Edges: 2,844,515 coupling relationships - Density: 8.7% (sparse but well-connected) - Clustering coefficient: 0.52 (strong local clustering) - Average path length: 2.8 (small-world property)

Highly coupled papers (top 5 by total coupling strength): 1. Davis (1989): Technology Acceptance Model - coupled to 1,247 papers 2. Delone & McLean (2003): IS Success Model - 1,156 papers 3. Venkatesh et al. (2003): UTAUT - 1,089 papers 4. Compeau & Higgins (1995): Computer Self-Efficacy - 892 papers 5. Gefen et al. (2003): Trust in E-Commerce - 847 papers

These are foundational theoretical papers cited across many topics---correctly identified as central via network analysis.

Stream coupling patterns: - Intra-stream edges: 2,341,892 (82.0%) - Inter-stream edges: 502,623 (18.0%) - Most connected stream pair: Stream 0 ↔ Stream 4 (87,234 edges) - Least connected: Stream 2 ↔ Stream 7 (1,234 edges)

High intra-stream coupling validates cluster coherence. Inter-stream coupling reveals intellectual bridges (e.g., Digital Transformation & Innovation ↔ social media).

## 4.7 Validation Against Manual Classification

We compared our automated taxonomy to expert-curated classifications where available:

AIS eLibrary Subject Classifications: AIS provides optional subject tags for papers. For 2,847 papers with tags: - Agreement rate: 73.0% (assigned to same/related stream) - Common discrepancies: Papers with multiple subjects assigned to one stream

Journal Special Issues: 487 papers from 34 special issues with explicit themes: - Same-cluster rate: 81.0% (special issue papers clustered together) - Indicates topic coherence

Author Keywords: For 5,234 papers with author keywords: - Keyword overlap: 0.64 average Jaccard between our keywords and authors' - Shows semantic alignment

These comparisons demonstrate reasonable alignment with human expert judgments while offering finer granularity.

## 4.8 Comparison to Prior Classifications

Sidorova et al. (2008) identified 5 core IS topics via co-word analysis: 1. IT & Organizations 2. IS Development 3. IT & Individuals 4. IT & Markets 5. IT & Groups

Our taxonomy at L1 maps to these but provides: - Finer granularity (8 vs. 5 streams) - Explicit separation of newer topics (digital transformation, social media) - Hierarchical detail (48 → 182 subtopics) - Quantitative validation via citation networks

Larsen et al. (2008) used LDA to find 6 research themes in top IS journals. Our results align with 4/6 themes, with additional granularity for emerging areas.

Our approach complements these seminal works by providing updated, more granular, and citation-validated taxonomy.

# 5. Discussion

## 5.1 Key Findings

1. Hybrid clustering substantially outperforms single-source methods

The 11.7× improvement in silhouette score (0.340 vs. 0.029) demonstrates clear advantages of combining text and citations. Text-only clustering produced overly broad, poorly separated clusters. Citation-only clustering conflated topically distinct papers citing common methods. The optimal 60.0%/40.0% weighting balances semantic content with structural validation.

This finding has implications for all bibliometric studies: citation patterns and textual content provide complementary information that should be jointly leveraged.

2. Hierarchical organization effectively captures multi-scale structure

IS research exhibits organization at multiple levels: broad paradigms (L1), research programs (L2), and specific questions (L3). A single-level clustering would miss either broad structure (too many clusters) or fine detail (too few). The three-level hierarchy accommodates both.

Different levels serve different purposes: - L1: Field overview for outsiders, curriculum design - L2: Literature review scope definition - L3: Targeted search for specific topics

3. The IS field shows clear evolutionary patterns

Temporal analysis reveals the field's maturation: - 1970s-80s: Technical focus (development, decision support) - 1990s-2000s: Organizational/strategic turn (governance, e-commerce) - 2010s: Social/digital emphasis (platforms, transformation) - 2020s: Emerging Technologies (AI, blockchain)

These shifts reflect broader technology and societal trends, demonstrating IS field's responsiveness to real-world phenomena.

4. Computational optimization enables practical application

Algorithmic innovations (inverted index, sparse matrices) reduced computation time from estimated hours to minutes, making full-corpus analysis feasible on commodity hardware. This democratizes literature analysis---researchers without specialized infrastructure can apply these methods.

5. Open science tools enhance impact and reproducibility

Releasing code, data, and interactive interface serves multiple goals: - Reproducibility: Others can validate findings - Extension:

Methods applicable to other domains - Accessibility: Non-experts can explore results - Education: Students learn literature analysis techniques

## 5.2 Implications for Research

For Literature Review Practice: - Computational methods complement (not replace) expert reviews - Hybrid approaches provide more robust starting points for manual synthesis - Hierarchical taxonomies guide systematic coverage of research areas

For Field Reflexivity: - Quantitative field mapping reveals blind spots and over-researched areas - Temporal analysis tracks paradigm shifts and emerging topics - Citation networks identify foundational works and intellectual lineages

For Research Planning: - Micro-topics reveal specific gaps in literature - Network centrality identifies influential works for theory development - Cluster boundaries suggest opportunities for synthesis and bridging

For Education: - Taxonomies inform curriculum design (cover major streams) - Hierarchical navigation helps students locate relevant literature - Visual representation aids understanding of field structure

## 5.3 Methodological Contributions

1. Optimal hybrid weighting

Systematic grid search identified 60.0%/40.0% text/citation split as optimal. This weighting likely transfers to other fields with similar characteristics (abstract-rich publications, strong citation norms). Future work should test domain-specificity.

2. Efficient bibliographic coupling

Inverted index algorithm provides 600× speedup over naive implementation. Generalizable to any bibliographic coupling application. Could be further optimized via parallelization.

3. Recursive NMF for hierarchy

Applying NMF within parent clusters (rather than full hierarchical clustering) balances interpretability with computational efficiency. Allows heterogeneous cluster numbers at different levels (unlike k-means trees).

4. Validation framework

Multi-faceted validation (silhouette, manual review, keyword coherence, comparison to expert classifications) provides confidence in results. Template for validating other literature mining studies.

## 5.4 Limitations

1. English-language bias: Analysis limited to English papers in top journals. Excludes non-English IS research (significant in Europe, Asia) and conference proceedings, dissertations, books. Results represent "elite" IS discourse.

2. Citation data availability: Despite 88% coverage (high for IS), 12% of papers lack citations. Disproportionately affects recent papers (citations not yet indexed) and older papers (references not digitized). May underweight emerging topics.

3. Abstract-only text: Analyzed titles/abstracts only (full text unavailable for many papers). Misses methodological details, results, deep concepts. Could miss papers using unexpected terminology in abstracts.

4. Static snapshot: Analysis treats corpus as static (2024 snapshot). Field continues evolving. Taxonomy will require periodic updates. Longitudinal analysis could track topic drift.

5. Parameter sensitivity: Results depend on clustering parameters (number of clusters, text/citation weights). While systematically optimized, some arbitrariness remains. Sensitivity analysis (Appendix B) shows results robust to moderate parameter changes.

6. Interpretability: Automated topic labels (top TF-IDF terms) sometimes require manual refinement for clarity. Human judgment still needed for final taxonomy presentation.

7. Disciplinary boundaries: IS field boundaries are fuzzy (overlaps with computer science, management, economics). Journal-based definition (AIS Basket) provides clear scope but may miss relevant work in adjacent fields.

## 5.5 Future Directions

Short-term extensions:

1. Temporal dynamics: Track topic birth, growth, decline over time using dynamic topic models (Blei & Lafferty, 2006)

2. Author analysis: Identify prolific authors per topic, collaboration networks, author migration across topics

3. Journal profiles: Characterize journals by topic distributions, identify specialization vs. generalization patterns

4. Regional differences: Compare topics across geographical regions (if author affiliation data available)

5. Full-text analysis: For papers with full-text access, analyze methods, results sections separately

Medium-term research:

1. Cross-domain application: Apply methodology to other fields (e.g., computer science, medicine, economics) to test

generalizability

2. Enhanced visualization: Develop network visualizations showing paper-paper relationships, temporal flows

3. Recommendation system: Build paper recommendation engine based on hybrid similarity

4. Automated review generation: Use topic models and citation patterns to generate draft literature reviews

5. Theory mapping: Link papers to theoretical frameworks, track theory usage patterns

Long-term vision:

1. Living literature map: Continuously updated taxonomy incorporating new publications in real-time

2. Multi-modal integration: Incorporate figures, tables, datasets, code repositories beyond text/citations

3. Causal inference: Move beyond correlation (coupling) to identify knowledge flows and influence chains

4. Meta-science applications: Use large-scale literature analysis to study science itself (publication biases, citation cartels, paradigm shifts)

5. AI-assisted synthesis: Combine literature mapping with large language models to generate comprehensive, up-to-date review syntheses

# 6. Conclusion

This research developed and validated a novel three-level hybrid clustering approach for mapping large academic literatures. By combining text similarity (TF-IDF + LSI) with citation network structure (bibliographic coupling), optimally weighted at 60.0%/40.0%, we achieved an 11.7-fold improvement in clustering quality over text-only methods.

Applied to 8,110 papers from AIS Basket of Eight journals (1977–2024), the methodology identified 8 major research streams, 48 detailed subtopics, and 182 granular micro-topics, comprehensively organizing the Information Systems field's intellectual structure. Algorithmic optimizations (inverted index, sparse matrices) enabled full-corpus analysis in under 5 minutes on commodity hardware.

The resulting taxonomy reveals the IS field's evolution from technical origins (systems development, decision support) through organizational/strategic emphasis (IT governance, e-commerce) to current focus on Digital Transformation & Innovation, social platforms, and Emerging Technologies. Citation network analysis identified foundational works and intellectual bridges between research areas.

Beyond empirical findings, this work contributes methodological innovations applicable to any research domain: optimal text/citation weighting, efficient bibliographic coupling algorithms, hierarchical topic organization, and comprehensive validation frameworks. Released as open-source tools with interactive web interface, the methodology and findings support reproducible research and community engagement.

As academic literature continues substantial growth across all fields, computational approaches are essential for synthesis and sensemaking. Hybrid methods combining multiple information sources---text, citations, metadata, usage data---will be increasingly important. This work provides a validated template for such integrative analyses.

The complete taxonomy, interactive explorer, and all analysis code are available at https://github.com/Data-ScienceTech/literature, enabling researchers to explore findings, validate results, and extend methods to new domains. We invite the community to build upon this foundation for advancing literature analysis practice and understanding the structure of scientific knowledge.

# References

Aggarwal, C.C., & Zhai, C. (2012). Mining Text Data. Springer.

AIS (2011). Senior Scholars' Basket of Journals. Association for Information Systems.

Baumeister, R.F., & Leary, M.R. (1997). Writing narrative literature reviews. Review of General Psychology, 1(3), 311-320.

Blei, D.M., Ng, A.Y., & Jordan, M.I. (2003). Latent Dirichlet allocation. Journal of Machine Learning Research, 3, 993-1022.

Blei, D.M., & Lafferty, J.D. (2006). Dynamic topic models. *Proceedings of ICML*, 113-120.

Boell, S.K., & Cecez-Kecmanovic, D. (2015). On being 'systematic' in literature reviews. Journal of Information Technology, 30(2), 161-173.

Boyack, K.W., & Klavans, R. (2010). Co-citation analysis, bibliographic coupling, and direct citation. *Journal of the American Society for Information Science and Technology*, 61(12), 2389-2404.

Cao, Y., et al. (2015). Multi-view learning of graph and signal data. IEEE Transactions on Signal and Information Processing over Networks, 1(2), 79-91.

Chen, C. (2017). Science mapping: A systematic review of the literature. Journal of Data and Information Science, 2(2), 1-40.

Chen, H., et al. (2012). Identifying emerging topics in science and technology. Research Policy, 41(8), 1421-1433.

Culnan, M.J. (1986). The intellectual development of management information systems. Management Science, 32(2), 156-172.

Culnan, M.J. (1987). Mapping the intellectual structure of MIS. *MIS Quarterly*, 11(3), 341-353.

Deerwester, S., et al. (1990). Indexing by latent semantic analysis. Journal of the American Society for Information Science, 41(6), 391-407.

Fortunato, S. (2010). Community detection in graphs. Physics Reports, 486(3-5), 75-174.

Fortunato, S., et al. (2018). Science of science. Science, 359(6379), eaao0185.

Ganguly, S., & Pudi, V. (2017). Paper2vec: Combining graph and text information for scientific paper representation. Proceedings of ECIR, 383-395.

Garfield, E. (2006). The history and meaning of the journal impact factor. JAMA, 295(1), 90-93.

Grover, A., & Leskovec, J. (2016). node2vec: Scalable feature learning for networks. Proceedings of KDD, 855-864.

Hamilton, W., Ying, R., & Leskovec, J. (2017). Inductive representation learning on large graphs. Proceedings of NIPS, 1024-1034.

Hunter, J.E., & Schmidt, F.L. (2004). Methods of Meta-Analysis. Sage.

Jain, A.K., Murty, M.N., & Flynn, P.J. (1999). Data clustering: A review. ACM Computing Surveys, 31(3), 264-323.

Janssens, F., et al. (2008). Hybrid clustering for validation and improvement of subject-classification schemes. *Information Processing & Management*, 44(4), 1502-1523.

Kang, Y.J., et al. (2015). Tracing the evolving focus of ISR. Information Systems Research, 26(4), 691-710.

Kessler, M.M. (1963). Bibliographic coupling between scientific papers. American Documentation, 14(1), 10-25.

Larsen, K.R., et al. (2008). A framework to advance theory with IS research. MIS Quarterly, 32(3), 633-656.

Lowry, P.B., et al. (2013). Influential author networks in the IS field. Journal of AIS, 14(Special Issue), 271-294.

Moher, D., et al. (2009). Preferred reporting items for systematic reviews and meta-analyses. PLoS Medicine, 6(7), e1000097.

Paré, G., et al. (2015). Synthesizing information systems knowledge: A typology of literature reviews. Information & Management, 52(2), 183-199.

Priem, J., et al. (2022). OpenAlex: A fully-open index of scholarly works, authors, venues, institutions, and concepts. ArXiv preprint.

Rosvall, M., & Bergstrom, C.T. (2008). Maps of random walks on complex networks reveal community structure. PNAS, 105(4), 1118-1123.

Rousseeuw, P.J. (1987). Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20, 53-65.

Rowe, F. (2014). What literature review is not. *European Journal of Information Systems*, 23(3), 241-255.

Salton, G., & Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. Information Processing & Management, 24(5), 513-523.

Sidorova, A., et al. (2008). Uncovering the intellectual core of the IS discipline. MIS Quarterly, 32(3), 467-482.

Small, H. (1973). Co-citation in the scientific literature. *Journal of the American Society for Information Science*, 24(4), 265-269.

Templier, M., & Paré, G. (2015). A framework for guiding and evaluating literature reviews. Communications of AIS, 37, 112-137.

Webster, J., & Watson, R.T. (2002). Analyzing the past to prepare for

the future. MIS Quarterly, 26(2), xiii-xxiii.

# Appendix A: Complete Level 2 Taxonomy

This appendix presents the complete listing of all 48 Level 2 subtopics identified through hierarchical clustering of the AIS Basket of Eight corpus. Each row represents one L2 subtopic, showing the top-8 representative keywords extracted via TF-IDF weighting. Subtopics are ordered by their parent Level 1 stream (Streams 0-7), with approximately 6 subtopics per stream.

| L2 ID | Representative Keywords |
|:------|:-----------------------|
| 0.0 | systems, information, research, information systems, development, paper, organizational, approach |
| 0.1 | reviews, product, online, consumers, products, ratings, review, consumer |
| 0.2 | social, media, social media, trust, virtual, network, team, users |
| 0.3 | firms, performance, firm, investments, technology, business, outsourcing, value |
| 0.4 | digital, transformation, Digital Transformation & Innovation, innovation, platform, platforms, infrastructure, digital platforms |
| 0.5 | knowledge, knowledge management, knowledge sharing, management, sharing, process, design, knowledge creation |
| 1.0 | systems, information, decision, information systems, management, support, development, design |
| 1.1 | consumers, pricing, product, consumer, market, price, online, advertising |
| 1.2 | social, media, social media, content, news, users, messages, posts |
| 1.3 | software, development, software development, open source, open, oss, source, vendor |
| 1.4 | auctions, bidders, auction, bidding, combinatorial, combinatorial auctions, bid, bids |
| 1.5 | digital, research, platforms, work, jatsp, study, platform, ai |
| 2.0 | adoption, technology, research, information, organizational, model, innovation, satisfaction |
| 2.1 | trust, sellers, online, buyers, online marketplaces, marketplaces, perceived, community sellers |
| 2.2 | perceived, acceptance, ease use, ease, use, usefulness, usage, technology acceptance |
| 2.3 | erp, implementation, erp implementation, erp systems, planning, planning erp, enterprise resource, resource planning |
| 2.4 | knowledge, project, teams, team, development, tms, knowledge sharing, projects |
| 2.5 | business, alignment, capabilities, performance, strategic, value, firms, firm |
| 3.0 | research, theory, design, systems, paper, information, framework, organizational |
| 3.1 | team, virtual, teams, group, communication, virtual teams, knowledge, social |
| 3.2 | firms, firm, governance, business, alignment, innovation, outsourcing, performance |
| 3.3 | project, control, software, development, projects, isd, systems development, risk |
| 3.4 | users, use, user, model, perceived, satisfaction, online, results |
| 3.5 | systems, information, information systems, university, management, mis, professor, journal |
| 4.0 | systems, information, development, paper, case, software, management, technology |
| 4.1 | jatspno, articlejatsp, jatspno abstract, available articlejatsp, abstract available, abstract, available, jatsp |
| 4.2 | university, management, professor, systems, information, information systems, journal, school |
| 4.3 | authors, journal, isj, issue, research, papers, editors, review |
| 4.4 | social, data, media, online, social media, platform, content, users |
| 4.5 | jatsp, reviewers, journal, reviewed, following individuals, thank reviewers, devoted time, reviewers manuscripts |
| 5.0 | systems, information, information systems, research, development, management, paper, organizational |
| 5.1 | firms, jatsp, outsourcing, market, platform, reviews, firm, consumers |
| 5.2 | privacy, information privacy, information, disclosure, individuals, protection, online, privacy protection |
| 5.3 | dss, decision, support, decision support, support systems, systems dss, design, decision making |

| 5.4 | digital, transformation, Digital Transformation & Innovation, technologies, digital technologies, platforms, social, innovation |
| 5.5 | knowledge, team, virtual, software, teams, members, project, work |
| 6.0 | knowledge, project, coordination, performance, team, software, projects, isd |
| 6.1 | perceived, usage, usefulness, trust, perceived usefulness, use, users, beliefs |
| 6.2 | alignment, business, firms, capability, capabilities, strategic, agility, firm |
| 6.3 | erp, erp systems, enterprise, implementation, smes, risk, planning erp, resource planning |
| 6.4 | systems, information, information systems, research, paper, support, group, groups |
| 6.5 | edi, integration, electronic, interchange, data interchange, electronic data, interchange edi, reengineering |
| 7.0 | usage, perceived, perceived usefulness, usefulness, bda, microcomputer, ease use, acceptance |
| 7.1 | use, theory, information, model, systems, research, information systems, support |
| 7.2 | adoption, beliefs, technology, model, influence, technology adoption, social influence, research |
| 7.3 | digital, digital options, options, debt, options thinking, work processes, work, knowledgebased |
| 7.4 | resistance, user resistance, user, change, resistance change, dispositional, implementation, use usefulness |
| 7.5 | trust, online, consumers, purchase, consumer, agents, perceived, recommendation agents |

Note: Some L2 topics (e.g., 4.1, 4.5) contain journal metadata artifacts (e.g., "jatsp", "reviewers") due to incomplete preprocessing of acknowledgment and front matter sections in source documents. These topics represent a small fraction of the corpus and do not affect the overall taxonomy quality.

## Appendix B: Sensitivity Analysis

To assess the robustness of our hybrid clustering approach, we conducted systematic sensitivity analyses by varying key parameters while holding others constant at their optimal values (text weight = 60.0%, citation weight = 40.0%, LSI dimensions = 200, min_df = 5). Table B.1 presents the results.

Table B.1: Parameter Sensitivity Analysis Results

| Parameter Varied | Value | Text Weight | Citation Weight | LSI Dims | min_df | Silhouette Score | Change from Optimal |
|:-----------------|:------|:------------|:----------------|:---------|:-------|:-----------------|:-------------------|
| Baseline (Optimal) | - | 60.0% | 40.0% | 200 | 5 | 0.340 | - |
| Text/Citation Weight | 70/30 | 70.0% | 30.0% | 200 | 5 | 0.287 | -15.6% |
| Text/Citation Weight | 50/50 | 50.0% | 50.0% | 200 | 5 | 0.312 | -8.2% |
| Text/Citation Weight | 40/60 | 40.0% | 60.0% | 200 | 5 | 0.251 | -26.2% |
| LSI Dimensions | 150 | 60.0% | 40.0% | 150 | 5 | 0.328 | -3.5% |
| LSI Dimensions | 250 | 60.0% | 40.0% | 250 | 5 | 0.335 | -1.5% |
| LSI Dimensions | 300 | 60.0% | 40.0% | 300 | 5 | 0.337 | -0.9% |
| Min Document Freq | 3 | 60.0% | 40.0% | 200 | 3 | 0.332 | -2.4% |
| Min Document Freq | 10 | 60.0% | 40.0% | 200 | 10 | 0.336 | -1.2% |
| Min Document Freq | 15 | 60.0% | 40.0% | 200 | 15 | 0.334 | -1.8% |

Interpretation: The analysis demonstrates that the 60.0%/40.0% text-citation weighting is robust, with alternative weightings showing substantial performance degradation (8.2-26.2% reduction in silhouette score). The clustering quality is relatively insensitive to LSI dimensionality above 150 dimensions (≤3.5% variation) and to minimum document frequency thresholds between 3 and 15 (≤2.4% variation), indicating stability across reasonable parameter ranges. The sharp decline in performance when citation weight exceeds text weight (40/60 configuration: -26.2%) confirms that semantic content provides essential discriminative power, while citations primarily serve a validating role. These findings support

the generalizability of our methodology to similar corpora with comparable characteristics.

# Appendix C: Code and Data Availability

### Reproducibility Statement

All materials necessary to reproduce the analyses presented in this paper are publicly available under open-source licenses to promote transparency, validation, and extension by the research community.

Repository: https://github.com/Data-ScienceTech/literature

Licenses:
- Code and scripts: MIT License (unrestricted use, modification, and distribution)
- Data and documentation: Creative Commons Attribution 4.0 International (CC-BY 4.0)

Pipeline Structure:
- current_pipeline/fetcher/ — Data collection from journal archives
- current_pipeline/enricher/ — Citation enrichment via OpenAlex API
- current_pipeline/analysis/ — Clustering and validation scripts
- tools/ — Utility scripts for export and bibliography generation
- submission/ — Manuscript, appendices, and reproducibility artifacts

Key Analysis Scripts:
- current_pipeline/fetcher/fetch_ais_basket_crossref.py — Fetch papers from AIS Basket journals
- current_pipeline/enricher/enrich_ais_basket_openalex.py — Citation enrichment via OpenAlex API
- current_pipeline/analysis/analyze_ais_basket_coverage.py — Coverage statistics and validation
- current_pipeline/analysis/analyze_enrichment_results.py — Citation network analysis
- Clustering implementation: See core Python files in repository root
- tools/export_l2.py — Generate Appendix A (Level 2 topics)
- tools/export_sensitivity.py — Generate Appendix B (sensitivity analysis)
- tools/build_bib.py — Compile references from citations

Artifacts Produced:
- ais_basket_enriched.csv — Full corpus with citation metadata (8,110 papers)
- hierarchical_clusters_L1_L2_L3.csv — Complete cluster assignments
- cluster_keywords.json — Representative keywords for all topics
- citation_network.npz — Sparse bibliographic coupling matrix
- dashboard/ — Interactive web-based literature explorer (HTML/CSS/JS)

Execution Instructions: See RUNBOOK.md in the repository for step-by-step instructions to replicate the complete analytical pipeline, including environment setup, data collection, clustering execution, and visualization generation. Estimated runtime: ~45 minutes on standard hardware (4-core CPU, 16GB RAM).

Interactive Explorer: The web-based literature explorer is deployed at https://data-sciencetech.github.io/literature/, providing searchable access to the complete taxonomy, cluster memberships, and citation networks.

# Appendix D: Interactive Explorer Guide

\[Screenshots and usage instructions for web interface\]

Word Count: \~12,500 words (excluding references and appendices)

Manuscript Status: DRAFT for open peer review\
Suggested Journal: MIS Quarterly, Information Systems Research, Journal of the AIS, or PLOS ONE (computational track)

Open Peer Review: Submitted for community feedback and collaborative revision