# CITESEERX DATA: SEMANTICIZING SCHOLARLY PAPERS

Jian Wu, IST, Pennsylvania State University

Chen Liang, IST, Pennsylvania State University

Huaiyu Yang, EECS, Vanderbilt University
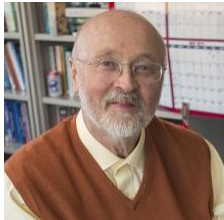
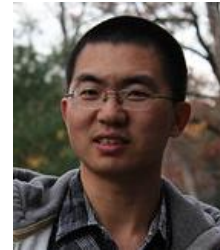C. Lee Giles, IST & CSE Pennsylvania State University

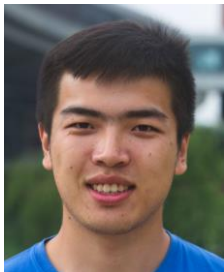The International Workshop on Scholarly Big Data (SBD 2016)
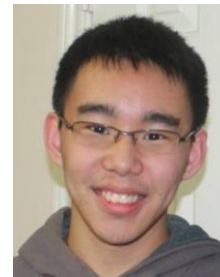
# Self-Introduction



Dr. C. Lee Giles
David Reese Professor
PI and Director of CiteSeerX

Dr. Jian Wu
Postdoctoral scholar
Tech leader of CiteSeerX

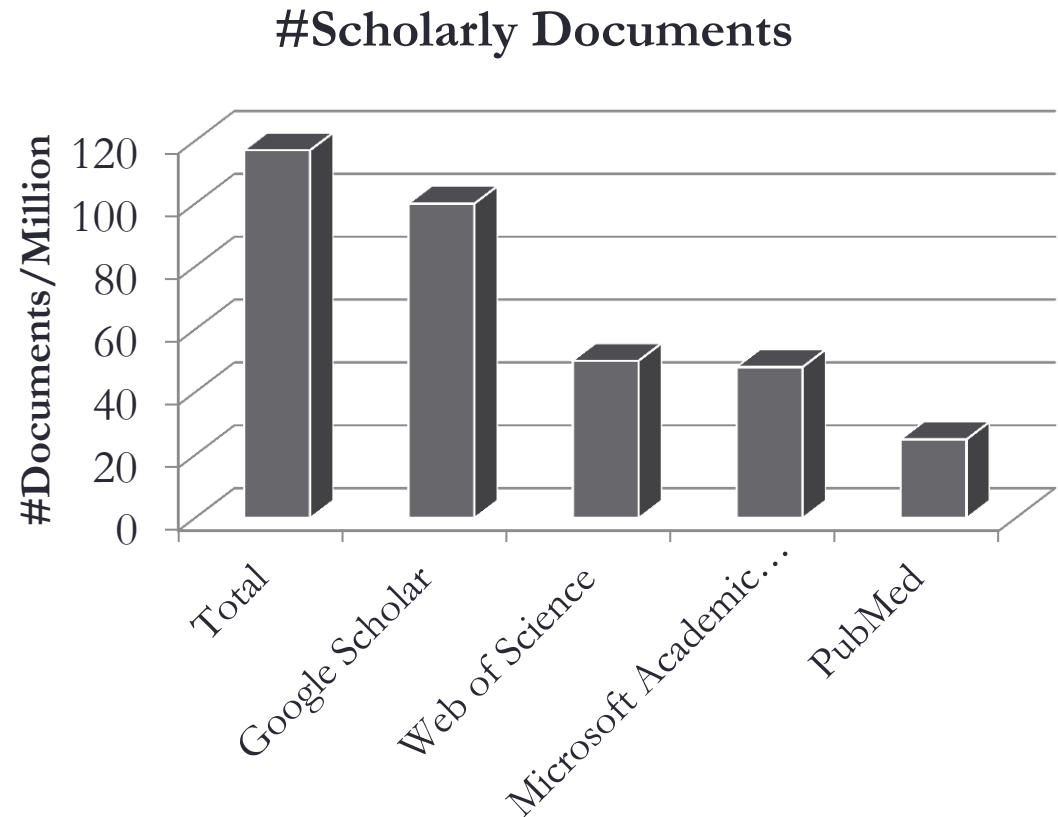Chen Liang
PhD student
Pennsylvania State University

Huaiyu Yang
Undergraduate student
Vanderbilt University

# Outline

- Scholarly Big Data and the Uniqueness of CiteSeerX Data
- Data Acquisition and Extraction
- Data Products
  - Raw Data
  - Production Database
  - Production Repository
- Data Management and Access
- Semantic Entity Extraction From Academic Papers

# Scholarly Data as Big Data

- "Volume"
  - About 120 million scholarly documents on the Web – 120TB or more [1]
  - Growing at a rate of >1 million annually
  - *English only – factor of 2 more with other languages*
  - Compare:
    NASA Earth Exchange Downscaled Climate Projections dataset (17TB)

**#Scholarly Documents**



[1] Khabsa and Giles (2014, PLoS ONE)

# Scholarly Big Data Features

- "Variety"
  - Unstructured: document text
  - Structured: title, author, citation, etc - metadata
  - Semi-structured: tables, figures, algorithms, etc.
  - Rich in facts and knowledge
  - Related data
    - Social networks, slides, course material, data "inside" papers
- "Velocity"
  - Scholarly Data is expected to be available in real time
- On the whole, scholarly Data can be considered an important instance of big data.

# Digital Library Search Engine (DLSE)

- Crawl-based vs. submission-based DLSEs

|  | **Crawl-based** | **Submission-based** |
|---|---|---|
| Data Source | Internet | Author upload |
| Metadata Source (majority) | Automatically Extracted | Author input + Automatically Extracted |
| Data Quality | varies | high |
| Human Labor (relatively) | Low | High |
| Accessibility | Open (or partially) | Subscription |

- Crawl-based DLSEs are important sources of scholarly data for *research* tasks such as citation recommendation, author name disambiguation, ontologies, document classification, and Science of Science
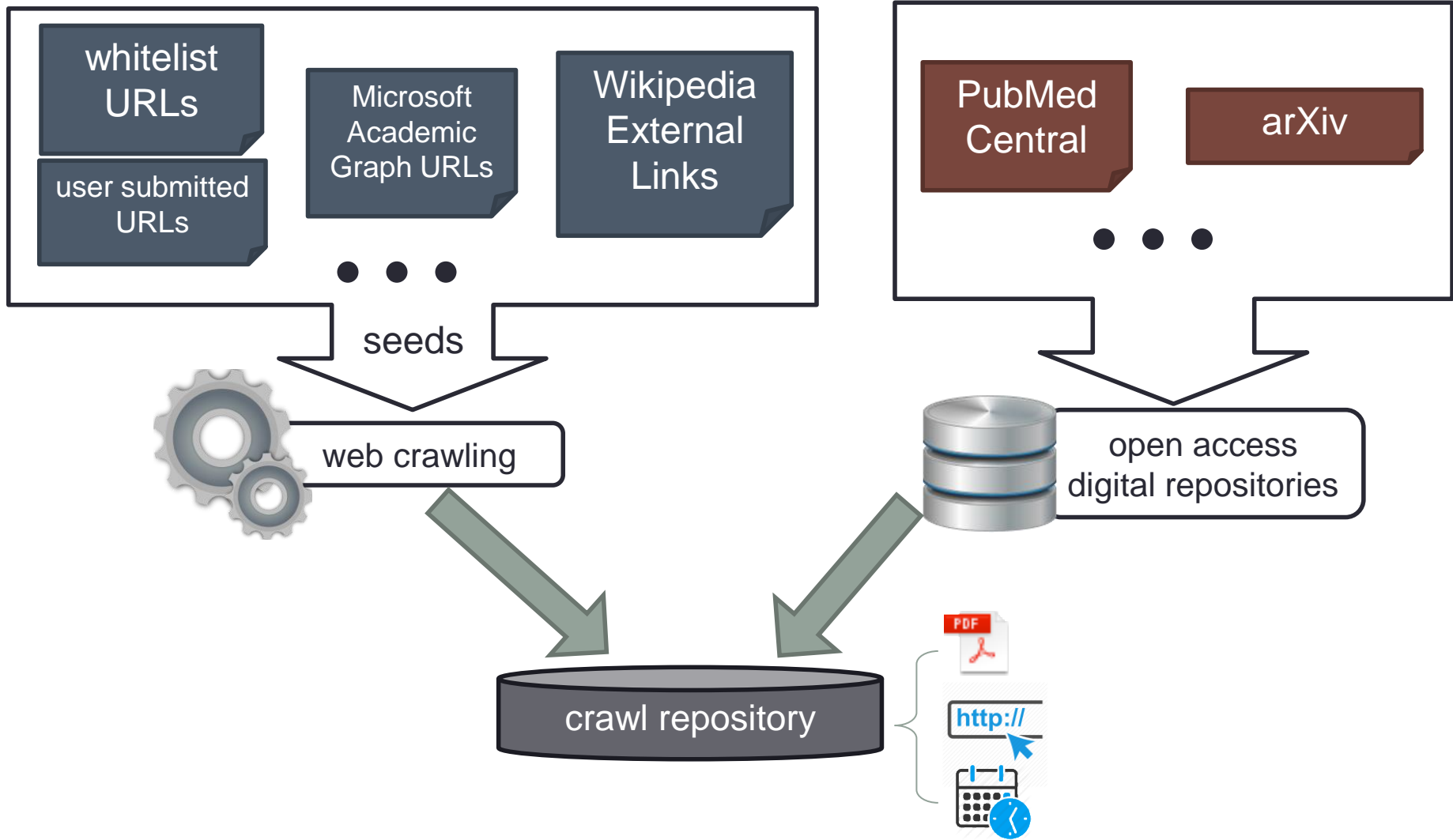
# The Uniqueness of CiteSeerX Data
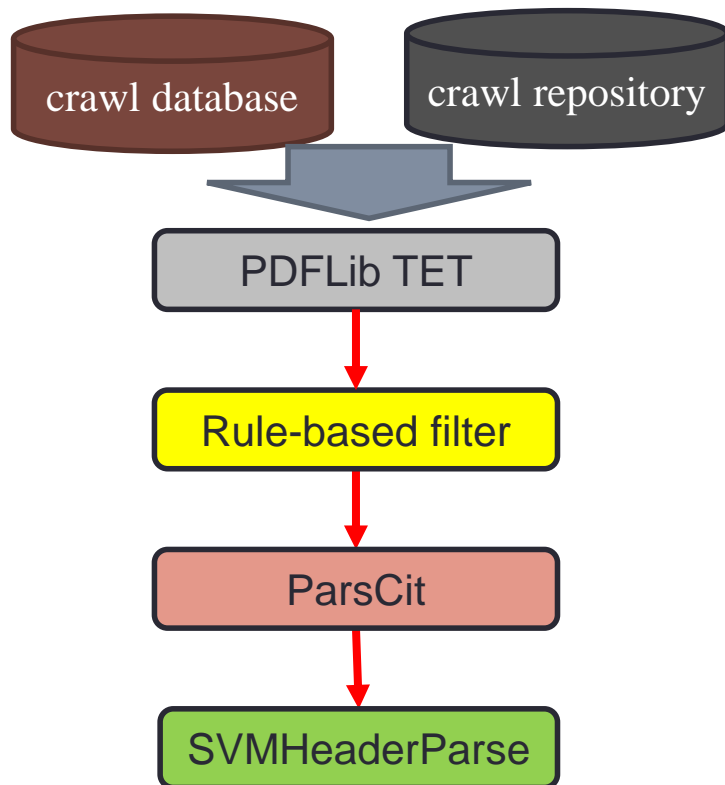
- Open-access Scholarly Data sets

| Datasets | DBLP | MAG* | CiteSeerX |
|---|---|---|---|
| Documents | 5 million | 100 million | 7 million |
| Header | y | y | y |
| Citations | n | y | y |
| URLs | y (publishers) | y (open + publishers) | y (open) |
| Full text | n | n | y |
| Disambiguated author names | n | n | y |

* MAG: Microsoft Academic Graph

# Data Acquisition

# Metadata Extraction



crawl database

crawl repository

PDFLib TET

Rule-based filter

ParsCit

SVMHeaderParse

Currently

crawl database

crawl repository

PDFMEF

PDFBOX/Xpdf

ML-based Filter

GROBID
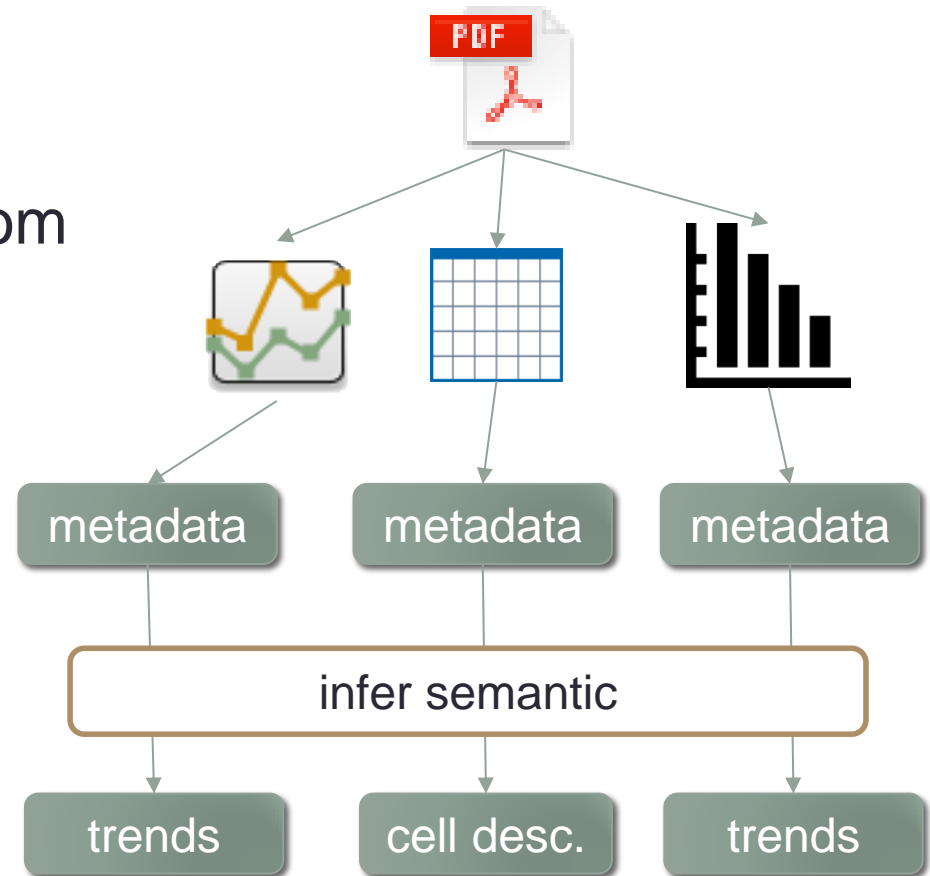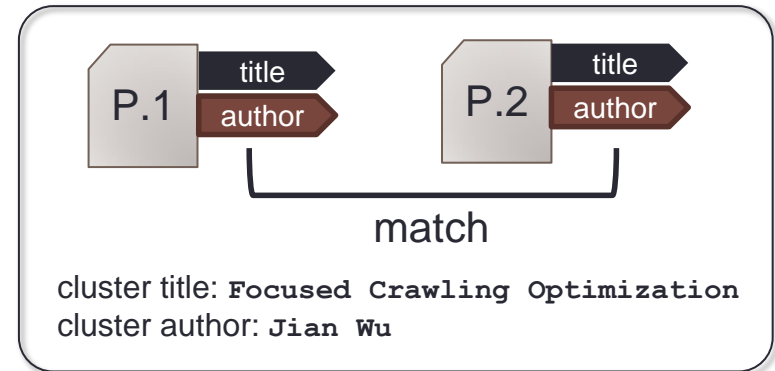
ParsCit

Under test

# Figures/Table/Barchart Extraction

- Data: CiteSeerX papers
- Extraction:
  - Extract figures + tables from papers
  - Extract metadata from figures + tables
- Large scale experiment
  - 6.7 Million papers in 14 days with 8 processes
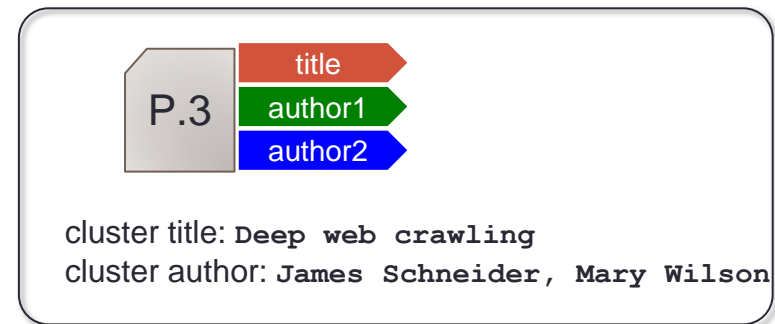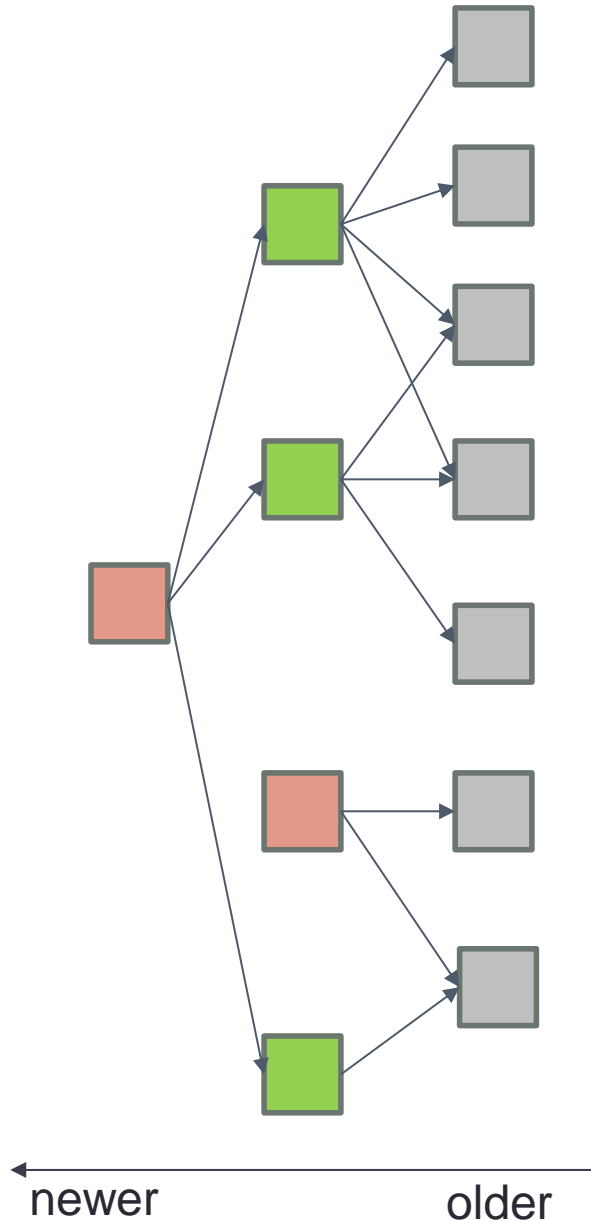
# Ingestion

- Ingestion feeds data and metadata to the *production* retrieval system
  - Ingestion clusters near-duplicate documents
  - Ingestion generate the citation graph (next slide)
  - Relational database
  - File system
  - Apache Solr
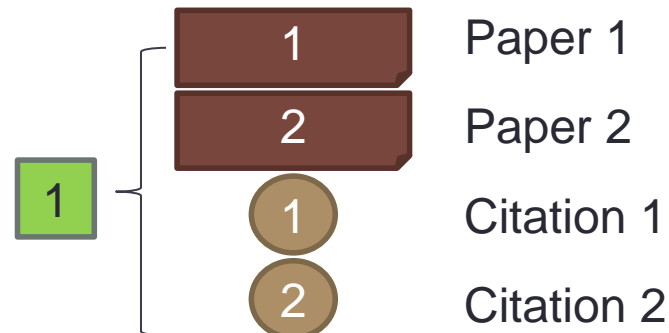


paper cluster 1



paper cluster 2

**1** Type 1 node: clusters with both in-degrees and out-degrees, containing papers, may contain citations

**2** Type 2 node (root): clusters with zero in-degree and non-zero out-degrees, only containing papers, i.e., papers that are not cited yet.

**3** Type 3 node (leaf): clusters with non-zero in-degree and zero out-degrees, only containing citation records, i.e., records without full text papers.

**Characteristics**:
- Directed
- No cycles: old papers cannot not cite new papers

1 — Paper 1
2 — Paper 2
1 — Citation 1
2 — Citation 2

newer                older

# Name Disambiguation

- Challenging due to name variations and entity ambiguity
- Task 1: distinguish different entities with the same surface name
- Task 2: resolve same entities with different surface names

Michael J. Jordan

Michael I. Jordan

Michael Jordan

**?**

Michael W. Jordan (footballer)

Michael Jordan (mycologist)

C L Giles

Lee Giles

C Lee Giles

Clyde Lee Giles

# User Correction

Web-crawling reliability (2004)

by Viv Cothey

**Venue:**   JOURNAL OF THE AMERICAN SOCIETY FOR INFORMATION SCIENCE AND
TECHNOLOGY

**Citations:** 10 - 1 self

Save to List
Add to Collection
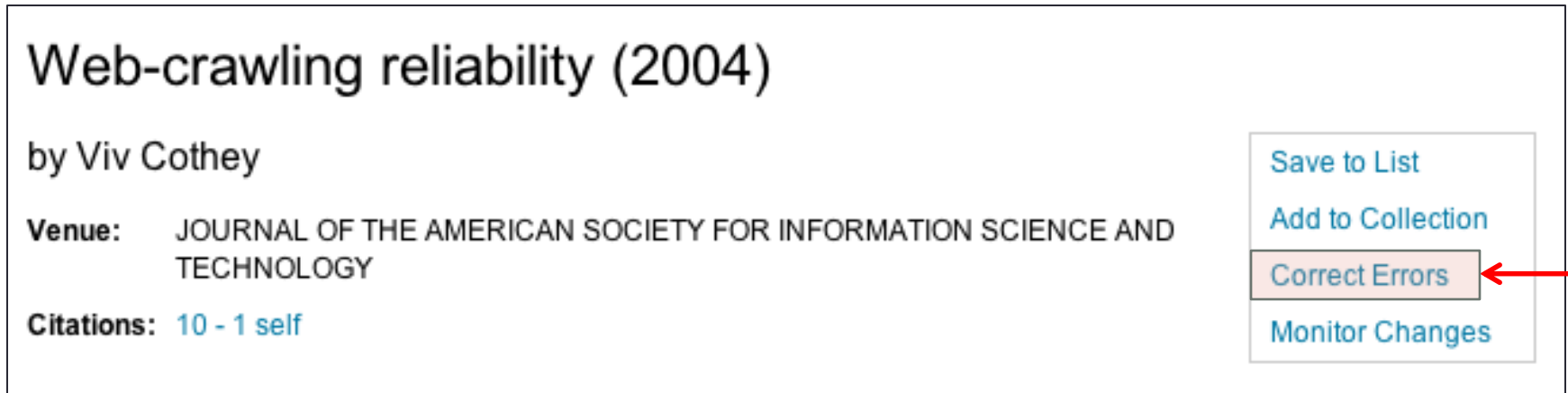Correct Errors
Monitor Changes

Figure: user-correction link on a paper summary page.

- Users can change almost all metadata fields
- New values are effective immediately after changes are submitted
- Metadata can be changed multiple times
- Version control
- About 1 million user corrections since 2008.

# Data Products

- ## Raw Data
  - ### Crawl repository
    - 24TB PDFs
  - ### Crawl database
    - 26 million <span style="color:red">document URLs</span>
    - 2.5 million <span style="color:blue">parent URLs</span>
    - 16GB

Document Collection of CiteSeerX

other page

PDF ← document URL

homepage ← parent URL

26 million

1.9 million

| | 0 | 5 | 10 | 15 | 20 | 25 | 30 |

2015, 2014, 2013, 2012, 2011, 2010, 2009, 2008

■ Indexed  ■ Ingested  ■ Crawled

# Data Products

- Crawl website   http://csxcrawlweb01.ist.psu.edu/

## CiteSeerX Crawler

Donate | Home | Submit | Query

- Document History Statistics
- User Submission Statistics (Last 30 Days)

submit a URL to crawl

### Country Ranking

1. Country Ranking by Number of Documents

Country ranking by number of docs

### Institution Level Domain Ranking

1. Domain Ranking by Number of Documents     [Cached Version]
2. Domain Ranking by Number of Citations     [Cached Version]
3. Domain Ranking by Citation Number Per Document     [Cached Version]

Domain ranking by number of crawled docs

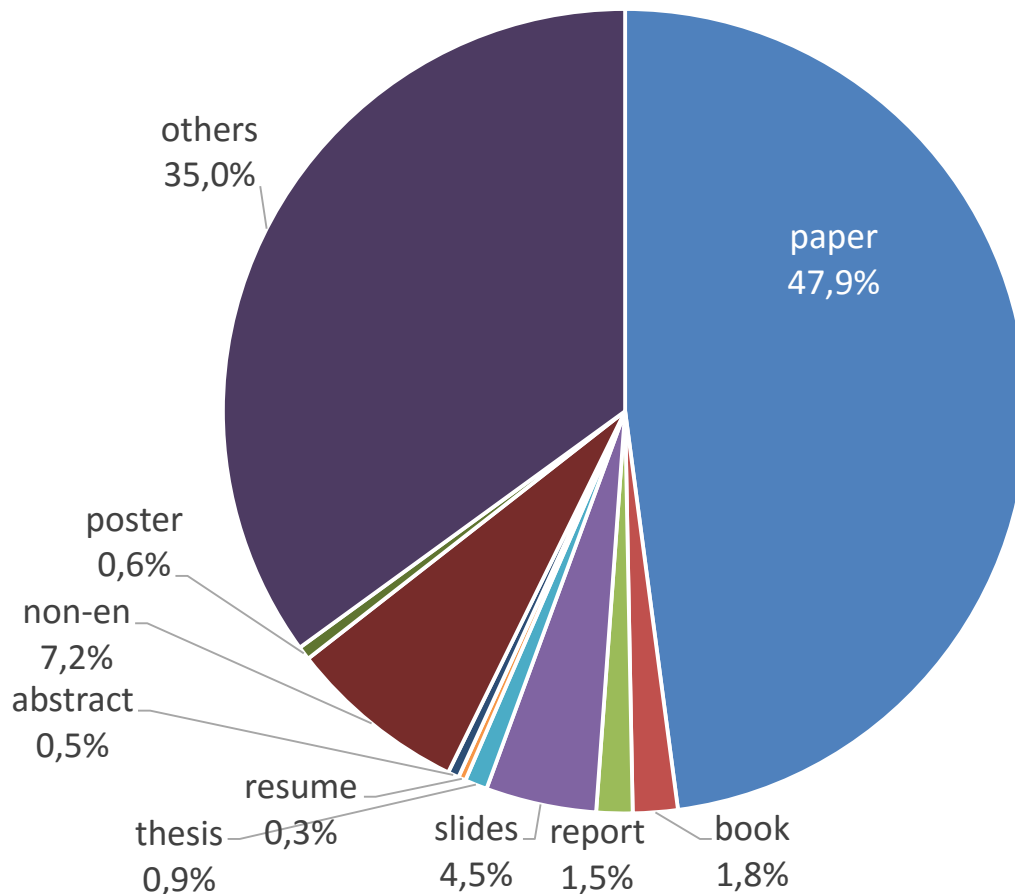### Top (Country) Level Domain Ranking

1. Top Level Domain Ranking by Number of Documents     [Cached Version]
2. Top Level Domain Ranking by Number of Citations     [Cached Version]
3. Top Level Domain Ranking by Citation Number Per Document     [Cached Version]

# What Documents Have We Crawled



- Manually label 1000 randomly selected crawled documents

- *Crawl repository* can be used for documents classification experiments to improve web crawling

- *Crawl database* can be used to generate whitelists and schedule crawl jobs

# Production Databases

- **`citeseerx`**
  - metadata directly extracted from papers

- **`csx_citegraph`**
  - paper clusters
  - citation graph

| database.table | description | rows |
|---|---|---|
| `citeseerx.papers` | header metadata | 6.8 million |
| `citeseerx.authors` | author metadata | 20.6 million |
| `citeseerx.cannames` | authors (disambiguated) | 1.2 million |
| `citeseerx.citations` | references | 150.2 million |
| `citeseerx.citationContext` | citation context | 131.9 million |
| `csx_citegraph.clusters` | citation graph (nodes) | 45.7 million |
| `csx_citegraph.citegraph` | citation graph (edges) | 112.5 million |

* Data are collected at the beginning of 2016.

# What Does Citation Graph Look Like



Suitable for large scale graph analysis

In-degree and out-degree distribution of CiteSeerX Citation Graph. Plots made by SNAP. Data are collected at the beginning of 2016.

# Production Repository

- 7 million academic documents (beginning of 2016)
- 9TB
  - PDF
  - XML (metadata)
  - body text
  - reference text
  - full text
  - version metadata files

- Classification Accuracy

| | | |
|---|---|---|
| paper | 83.0% | |
| others | 7.5% | |
| report | 4.5% | |
| thesis | 2.6% | |
| slides | 0.8% | academic documents |
| book | 0.7% | 92.1% |
| abstract | 0.3% | |
| non-en | 0.3% | |
| poster | 0.2% | |
| resume | 0% | |

# Production Repository

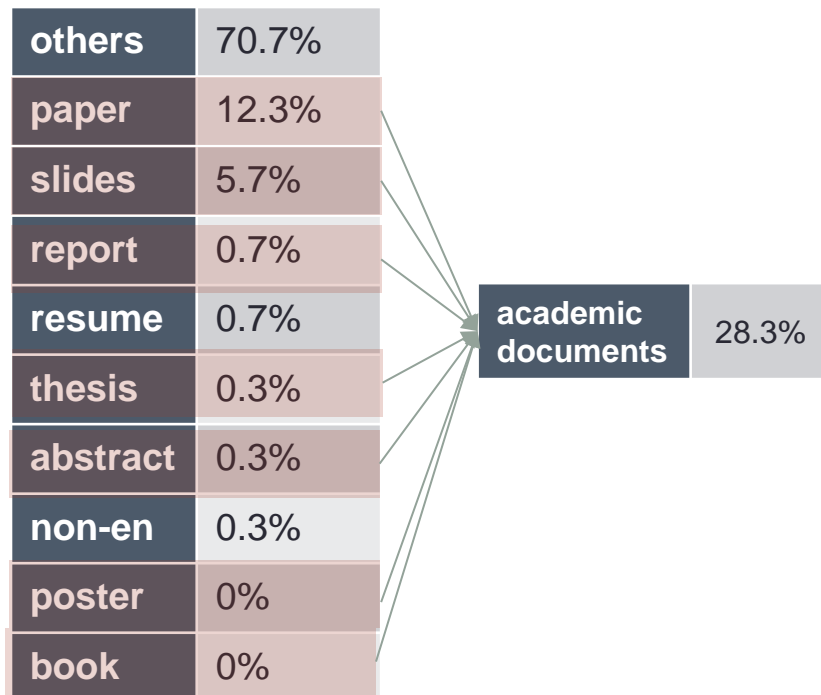- ## False Negatives
  - Documents mis-classified as non-academic documents

| | |
|---|---|
| **others** | 70.7% |
| **paper** | 12.3% |
| **slides** | 5.7% |
| **report** | 0.7% |
| **resume** | 0.7% |
| **thesis** | 0.3% |
| **abstract** | 0.3% |
| **non-en** | 0.3% |
| **poster** | 0% |
| **book** | 0% |

**academic documents** 28.3%

- ## Improving Classification Accuracy
  - Classifier based on Machine Learning and Structural features (Caragea et al. 2014 WSC; Caragea et al. 2016 IAAI)
  - Accuracy > 90%

# Estimate Near-duplication Rate

- Directly evaluating de-duplication is non-trivial.
- Infer and derive the near-duplication rate indirectly from two samples
  - Sample A: 100 clusters, $S = 2$, 200 documents
  - Sample B: 100 clusters, $S > 2$, 430 documents
  - Ground truth: manually extract titles, authors, years, and venues
  - Metrics:
    - Sample A: true duplication rate
    - Sample B: partial duplication rate

| Sample | S | NC | %True | D-ratio |
|--------|-----|-----|-------|---------|
| A | 2 | 100 | 84% | 1.16 |
| B | >2 | 100 | 70% | 2.26 |

S: Cluster size
NC: Number of clusters in a sample
%True: Percentage of true clusters in a sample

$$\text{D-ratio} = \frac{\text{Number of distinct documents}}{\text{NC}}$$

# Near-duplication Rate of CiteSeerX Data

| Cluster Sizes | 1 | 2 | 3 | 4 | >4 |
|---|---|---|---|---|---|
| NC (million) | 5.08 | 0.45 | 0.10 | 0.03 | 0.03 |
| Percentage | 92.8% | 7.91% | 1.76% | 0.53% | 0.53% |

Total number of distinct documents = 5.08+0.45x1.16+0.16x2.26 $\simeq$ **5.96**

Near-duplication rate = (1 – 5.96/6.70) x 100% = 11%

Number of clusters = 5.08+0.45+0.10+0.03+0.03=5.69 < **5.96**

Improve de-duplication accuracy:
- Cleansing metadata: GROBID [1]
- Alternative algorithms: e.g., *simhash* [2]

[1] Jian Wu, Jason Killian, Huaiyu Yang, Kyle Williams, Sagnik Ray Choudhury, Suppawong Tuarob, Cornelia Caragea, and C. Lee Giles. "PDFMEF: A Multi-Entity Knowledge Extraction Framework for Scholarly Documents and Semantic Search." In: Proceedings of The 8th International Conference on Knowledge Capture (K-CAP 2015), Palisades, NY, USA
[2] Kyle Williams, Jian Wu, and C. Lee Giles. "SimSeerX: A Similar Document Search Engine." In:The 14th ACM Symposium on Document Engineering (DocEng 2014), Fort Collins, CO, USA

# Data Management and Access

- Master database: 2x replication VMs hosted in a local private cloud; 2x copies of database dumps

- Search index: Apache Solr 4.9 replicated on a pair of twin VMs. Successfully indexed data on *SolrCloud*

- Production Repository: 2x sync'ed virtual servers; 2x snapshots; accessed via a RESTful API

- Public accessibility: Amazon S3, updated every 2-3 months

- Please contact us if you are interested in using CiteSeerX data

Include Citations

Advanced Search

**Cite**
**Seer**
**X** =7M

Most Cited: Documents , Citations , Authors , Venue Impact Rating

Powered by: Solr

About CiteSeerX        Submit and Index Documents        Privacy Policy        Help        Data        Source        Contact Us

# Semantic Scholarly Entity Extraction

- Motivation
  - Traditional search
    - Indexing metadata
    - Itemizing results
  - Intelligent Semantic Search
    - Answer questions
    - Recommendation
    - Summarization
    - Comparison

| Structural entities | Semantic entities |
|---|---|
| Title | People |
| Authors | Locations |
| Year | Concepts |
| Venue | Tools |
| Figures | Methods |
| Tables | Datasets |

# Scholarly Semantic Entities

- A Scholarly Semantic Entity (SSE) is a semantic entity that appears and/or is described in an academic document that delivers *domain specific knowledge* including a concept, a tool, a method, or a dataset.

- Examples:

  - IPv6 (concept)
  - NLTK (tool)
  - Conditional random field (method)
  - WebKB (dataset)

- Keyphrases in general constitute a subset of SSEs, but SSEs include a broader range of words and phrases.

- Entity linking can resolve a fraction of SSEs, e.g., using Wikifier (UIUC), but there are more to be discovered.

- Few research articles on extracting SSEs.

# Entity Linking Experiments

- 24859 papers randomly selected from CiteSeerX repository
- UIUC Wikifier [1,2]
- 21300 are successfully processed
- Outputs: Wikipedia terms + link score ($S$)
- Empirical cut-off of $S=0.8$ to remove less meaningful terms and single character symbols



Examples of high frequency terms: *Algorithm, Cell (biology), Matrix (mathematics), Protein, United States, Energy, Temperature, One half, Need To, Theorem*

[1] X. Cheng and D. Roth. Relational inference for wikication. In EMNLP, 2013.
[2] L.-A. Ratinov, D. Roth, D. Downey, and M. Anderson. Local and global algorithms for disambiguation to wikipedia. In ACL, 2011.

# On-going Work on Extracting SSEs

- Knowledge base independent
- Applying lexical semantic tools such as NLTK and Stanford CoreNLP tools. Will try Google SyntaxNet
- Supervised Machine Learning
- Focusing on Computer and Information Sciences and Engineering (CISE) papers, e.g., WWW, VLDB, ACL conferences/journals

- Examples of Tagged SSEs
  - Digital Library Search Engine
  - DB Entity Model
  - XML Beans
  - XML Query Language
  - Microsoft SQL Server
  - WCF
  - Loosely Type XML object
  - LINQ Query Translator
  - XML Schema Types
  - HUB4

# Future Work

- CiteSeerX Data
  - Scale-up to 30 million academic documents
  - Improve metadata quality
  - More open access entities, e.g., figures+tables
  - Integrate extraction, ingestion, and indexing; goal: process 1 million docs in 2 days

- SSE Extraction
  - Increase labeled sample sized and quality
  - Develop more efficient features
  - Start with basic ML models
  - Make it scalable

# Summary

- CiteSeerX **_actively_** crawls researcher homepages on the web for scholarly papers, formerly in computer science
  - Converts PDF to text
  - Automatically extracts OAI metadata and <span style="color:red">other data</span>
  - Automatic citation indexing, links to cited documents, creation of document page, author disambiguation
  - Software **open source** – can be used to build other such tools
  - **All** data shared

- 7 M documents
- 150 M citations
- 21 M authors
  - 1.2 M disambiguated
- 3 M hits per day on average
- 1 M page views/month
- 200k documents added monthly
- 150 million documents downloaded annually
- 1 M individual users
- ~40 TB