Leveraging Institutional Data For Author Name Disambiguation



Michael Bales, Paul Albert, Jie Lin, and Stephen Johnson

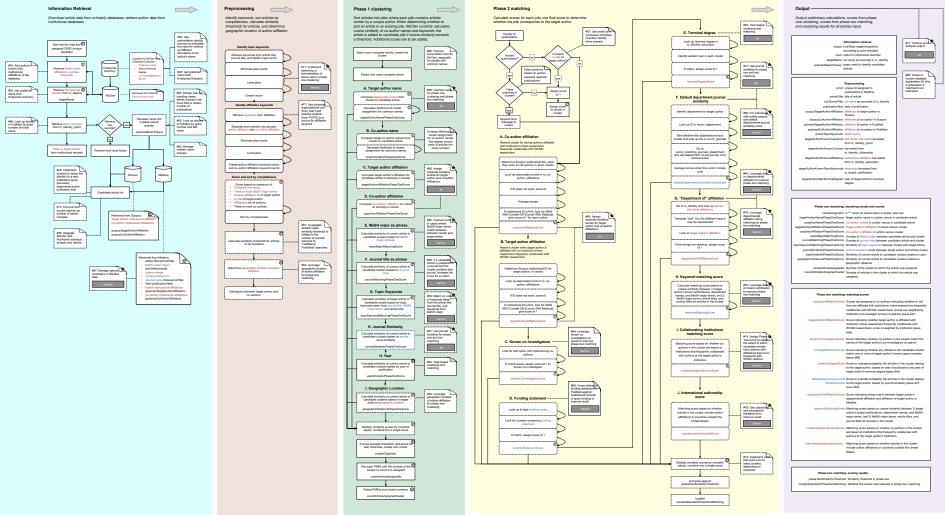
WE WOMEN TO A SOM IN IN TRANSPORT ASSUME

ReCiter Architecture and Data Processing Operations

Last updated June 25th, 2015

Michael Bales, Paul Albert, Jie Lin, and Stephen Johnson, Weill Cornell Medical College



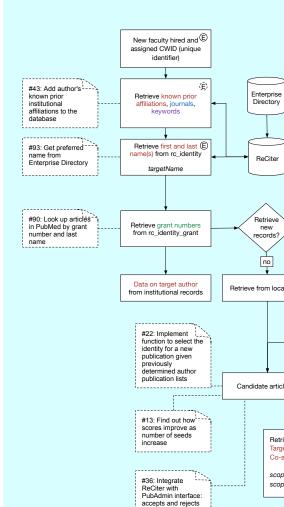


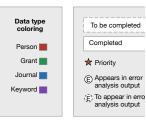


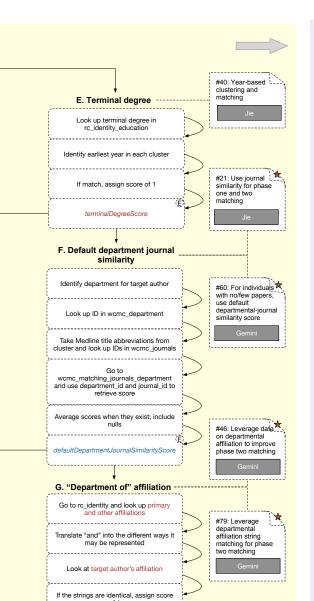


Information Retrieval

Download article data from scholarly databases







Output preliminary calculations, scores from phase one clustering, scores from phase two matching, and clustering results for all articles input. #71: Improve error Information retrieval status: true/false negative/positive, analysis output according to gold standard cwid: author's institutional identifier targetName: full name as recorded in rc_identity pubmedSearchQuery: query used to identify candidate #95: Output a human-readable explanation for why Preprocessing a publication is matched to an pmid: unique ID assigned to individual publications in Medline articleTitle: title of article fullJournalTitle: full name as recorded in rc identity publicationYear: year of publication scopusTargetAuthorAffiliation: affiliation of target author in Scopus scopusCoAuthorAffiliation: affiliation of co-author in Scopus pubmedTargetAuthorAffiliation: affiliation of author in PubMed pubmedCoAuthorAffiliation: affiliation of co-author in PubMed articleTopicKeywords: MeSH terms targetAuthorKnownCoauthors: last name, first initial harvested from rc_identity_grant targetAuthorKnownCountry: harvested from rc_identity_citizenship targetAuthorKnownAffiliations: institutional affiliation harvested from rc_identity_education targetAuthorKnownTopicKeywords: keywords harvested from rc_board_certification targetAuthorYearTerminalDegree: year of target author's terminal degree Phase one clustering: clustering results and scores clusterOriginator: A "*" when an article starts a cluster; else null targetAuthorNamePhaseOneScore: Target author name in cluster versus in candidate article coauthorNamePhaseOneScore: Co-author names in cluster versus in candidate article targetAuthorAffiliationPhaseOneScore: Target author's affiliation in article versus cluster coauthorAffiliationPhaseOneScore: Co-author's affiliation in article versus cluster meshMajorMatchingScore: Overlap of MeSH major between candidate article and cluster journalMatchingPhaseOneScore: Overlap of journal titles between candidate article and cluster topicKeywordMatchingPhaseOneScore: Similarity of topic keywords between cluster and target article

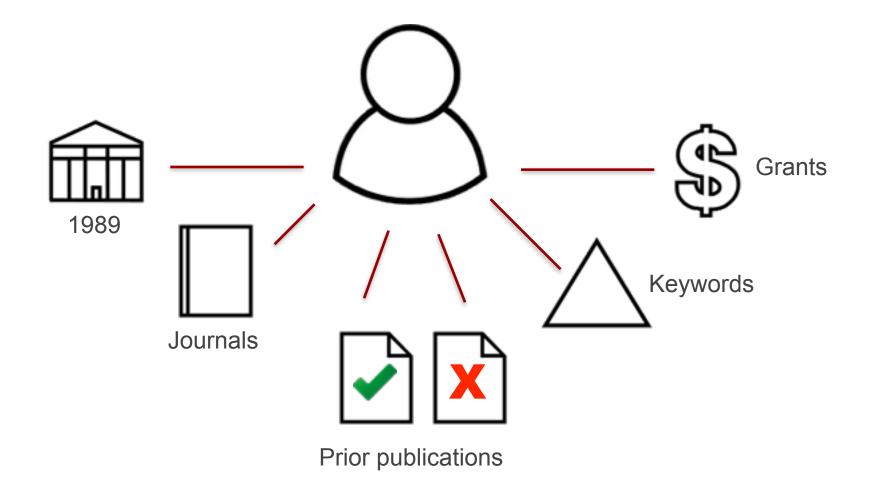
ReCiter Phases

- Information retrieval
- Preprocessing
- Phase one clustering
- Phase two matching
- Output

Use cases

- 1. Identify publications authored by new member of WCMC community
- 2. Assign newly identified publications

Evidence



ReCiter Evidence Types - Target Author

- First name, last name, middle initial
- Board certifications and clinical expertise
- Year of terminal degree
- Institutions where degrees earned
- E-mail address
- Known prior affiliations
- Geographic location of affiliation
- Grants & co-investigators

ReCiter Evidence Types, continued

- Target author prior publications
 - Journals
 - MeSH major topics
 - Co-authors
 - Names
 - Institutional affiliations
- WCMC collaborating institutions

ReCiter - Selected Evidence for Authors

Data	Source	Example	Phase 1	Phase 2
Known publication	Publications management	12923412	Yes	Yes
Job title	WOOFA	Professor of Medicine	No	Yes
Primary department	WOOFA	Medicine	No	Yes
Appointment period	WOOFA	1978 - current	No	Yes
Board certifications	Physicians profile	Cardiovascular disease	No	Yes
Citizenship	WOOFA	United States	No	Yes
Degree	WOOFA	Doctoral, 1971	No	Yes
Alma mater	WOOFA	Yale University, 1971	No	Yes
Grant	Coeus	5 U01 HL54495-10 EWOF	No	Yes



Contents lists available at ScienceDirect

Journal of Biomedical Informatics

journal homepage: www.elsevier.com/locate/yjbin



Automatic generation of investigator bibliographies for institutional research networking systems



Stephen B. Johnson a.*, Michael E. Bales b, Daniel Dine b,c, Suzanne Bakken b,c, Paul J. Albert d, Chunhua Weng b,c

ARTICLE INFO

Article history: Received 13 December 2013 Accepted 20 March 2014 Available online 30 March 2014

Keywords:
Authorship
Bibliography as topic
MEDLINE
Natural language processing
Pattern recognition
Automated

ABSTRACT

Objective: Publications are a key data source for investigator profiles and research networking systems. We developed ReCiter, an algorithm that automatically extracts bibliographies from PubMed using institutional information about the target investigators.

Methods: ReCiter executes a broad query against PubMed, groups the results into clusters that appear to constitute distinct author identities and selects the cluster that best matches the target investigator. Using information about investigators from one of our institutions, we compared ReCiter results to queries based on author name and institution and to citations extracted manually from the Scopus database. Five judges created a gold standard using citations of a random sample of 200 investigators.

Results: About half of the 10,471 potential investigators had no matching citations in PubMed, and about 45% had fewer than 70 citations. Interrater agreement (Fleiss' kappa) for the gold standard was 0.81. Scopus achieved the best recall (sensitivity) of 0.81, while name-based queries had 0.78 and ReCiter had 0.69. ReCiter attained the best precision (positive predictive value) of 0.93 while Scopus had 0.85 and name-based queries had 0.31.

Discussion: ReCiter accesses the most current citation data, uses limited computational resources and minimizes manual entry by investigators. Generation of bibliographies using named-based queries will not yield high accuracy. Proprietary databases can perform well but requite manual effort. Automated generation with higher recall is possible but requires additional knowledge about investigators.

© 2014 Elsevier Inc. All rights reserved.

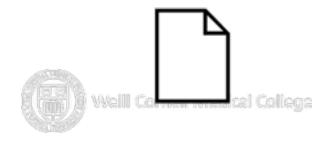


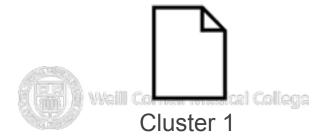
Department of Public Health, Weill Cornell Medical College, New York, United States

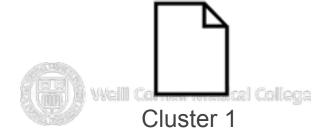
b Department of Biomedical Informatics, Columbia University, New York, United States

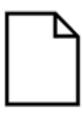
⁶The Irving Institute for Clinical and Translational Research, Columbia University, New York, United States

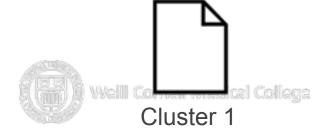
^d Samuel J. Wood Library, Weill Cornell Medical College, New York, United States



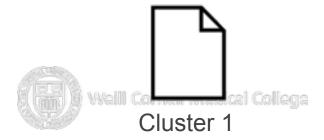


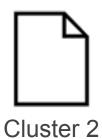




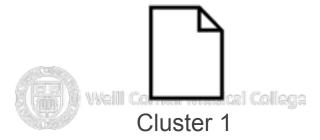












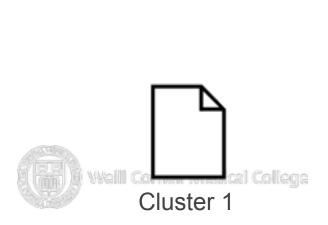


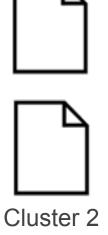


Cluster 3

Ana Santos-Carvallo

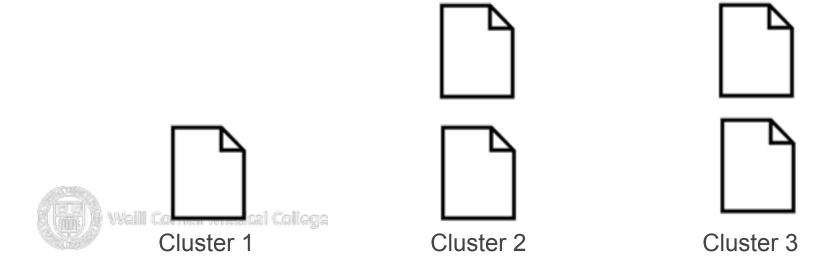
Santos A[AU] OR
Santos-Carvallo A[AU]

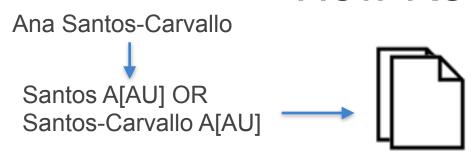


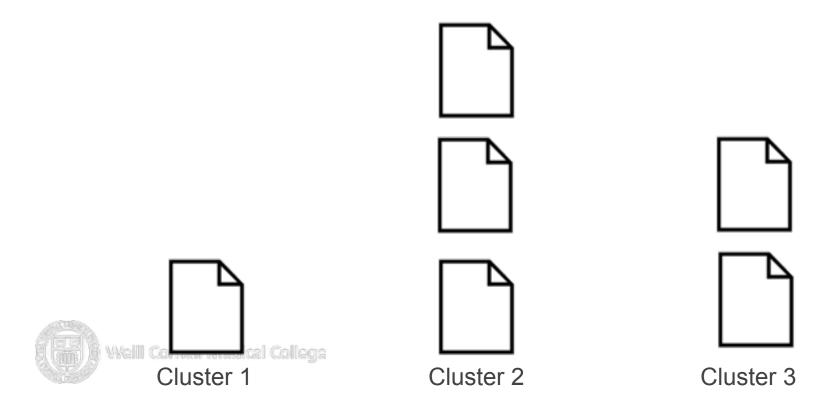


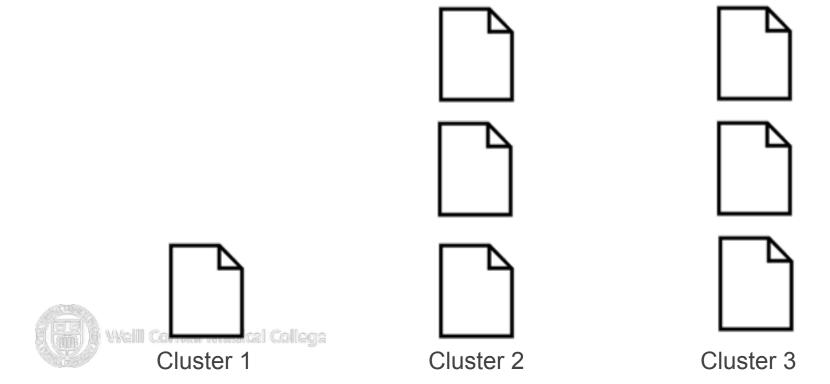


Cluster 3









Ana Santos-Carvallo Santos A[AU] OR Santos-Carvallo A[AU] Cluster 1 Cluster 3 Cluster 2

Ana Santos-Carvallo Santos A[AU] OR Santos-Carvallo A[AU] Cluster 1 Cluster 3 Cluster 2

Person: Jonathan W. Weinsaft



Board Certifications: Cardiovascular Disease, Internal Medicine

Article cluster #1







JACC. Cardiovascular Imaging

Article cluster #2









Person: Jonathan W. Weinsaft



Board Certifications: Cardiovascular Disease, Internal Medicine

Article cluster #1



Duke Cardiovascular Magnetic Resonance Center



Division of Cardiology, New York Methodist Hospital



JACC. Cardiovascular Imaging

Article cluster #2



Dept. of Ophthalmology, Mayo Clinic



Tim Dept. of Ophthalmology and Visual Sciences ...





Person: Jonathan W. Weinsaft



Board Certifications: Cardiovascular Disease, Internal Medicine

Score = 0.27

Article cluster #1







JACC. Cardiovascular Imaging

Article cluster #2



Dept. of Ophthalmology, Mayo Clinic



חבוד Dept. of Ophthalmology and Visual Sciences ...





Person: Jonathan W. Weinsaft



Board Certifications: Cardiovascular Disease, Internal Medicine

Score = 0.27

Article cluster #1







JACC. Cardiovascular Imaging

Article cluster #2



Dept. of Ophthalmology, Mayo Clinic



חבוד Dept. of Ophthalmology and Visual Sciences ...





Person: Jonathan W. Weinsaft



Board Certifications: Cardiovascular Disease, Internal Medicine

Score = 0.27

Article cluster #1



Duke Cardiovascular Magnetic Resonance Center



Division of Cardiology, New York Methodist Hospital



JACC. Cardiovascular Imaging

Article cluster #2



Dept. of Ophthalmology, Mayo Clinic



חבוד Dept. of Ophthalmology and Visual Sciences ...



Cornea

Score = 80.0

Person: Shahin Rafii



Grant co-investigators: Bi-Sen Ding, Zhongwei Cao

Article cluster #1 includes:



Ding BS



Cao Z



Kao DI

Article cluster #2 includes:



Kaira K



Yasuda M



Palefsky JM



Person: Shahin Rafii



Grant co-investigators: Bi-Sen Ding, Zhongwei Cao

Article cluster #1 includes:



Ding BS



Cao Z



Kao DI

Article cluster #2 includes:



Kaira K



Yasuda M



Palefsky JM



Person: Shahin Rafii



Grant co-investigators: Bi-Sen Ding, Zhongwei Cao

Score = 0.45

Article cluster #1 includes:



Ding BS



Cao Z



Kao DI

Article cluster #2 includes:



Kaira K

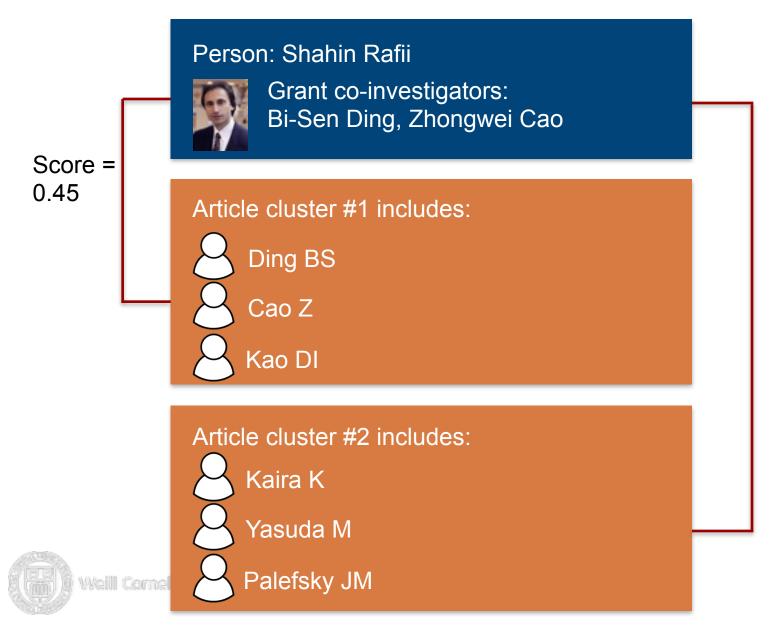


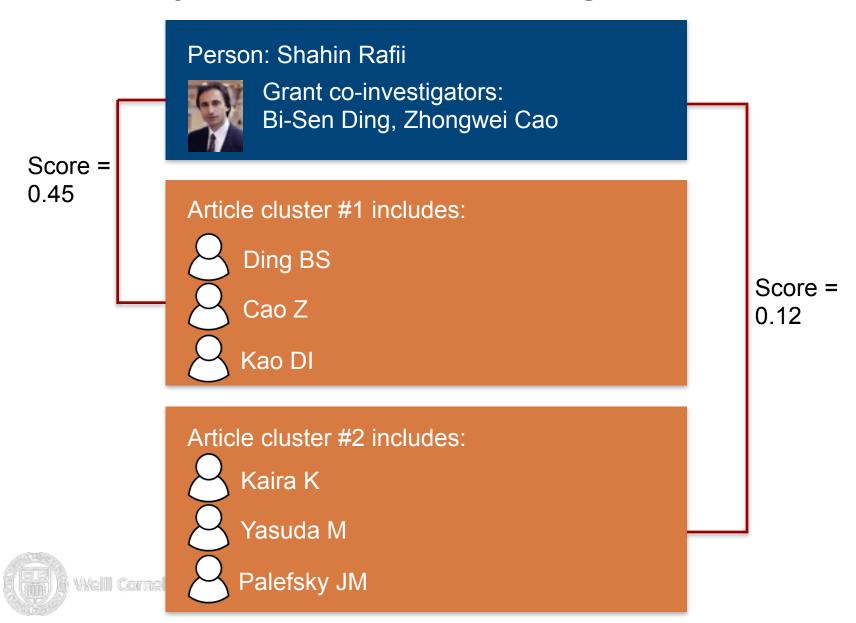
Yasuda M



Palefsky JM



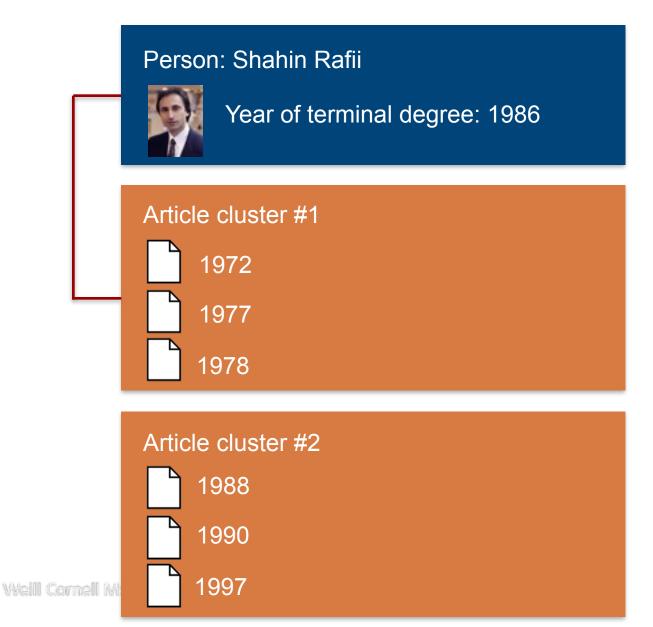


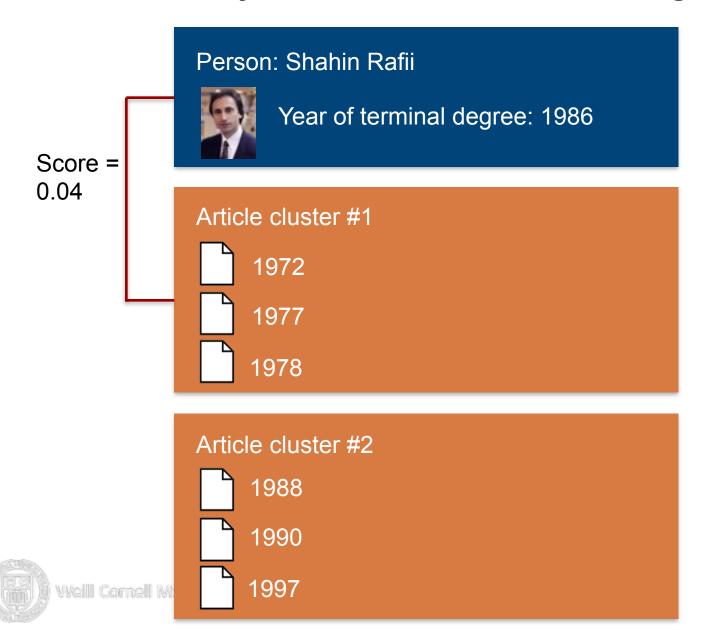


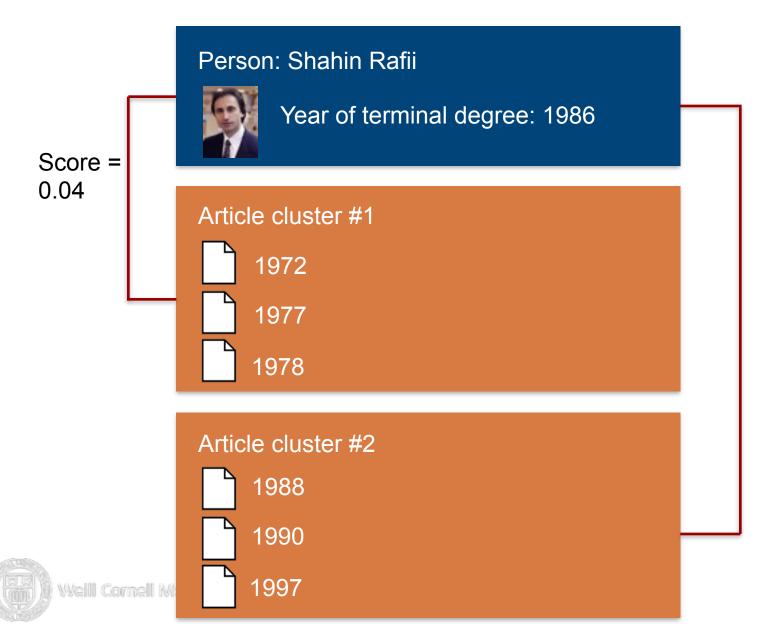
Similarity Score: Year of Terminal Degree

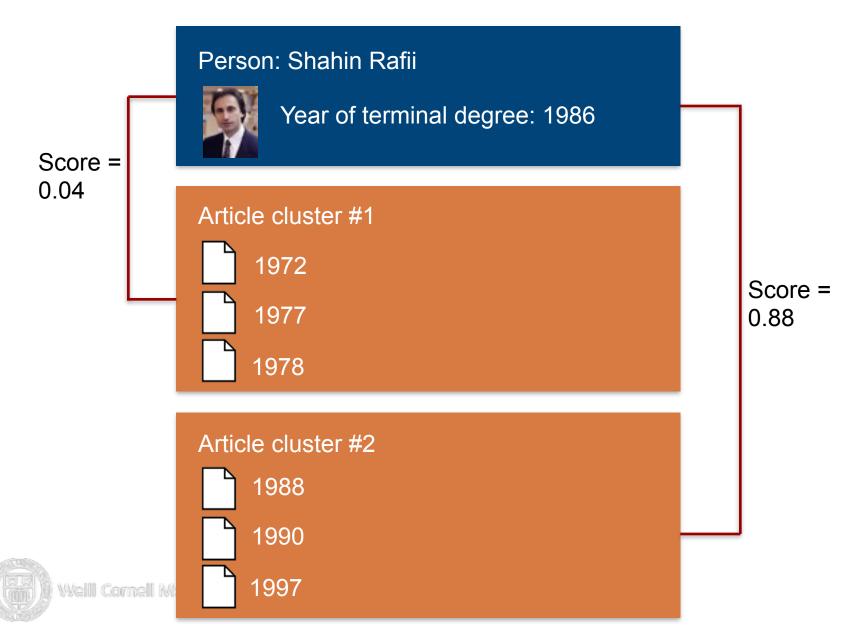
Person: Shahin Rafii Year of terminal degree: 1986 Article cluster #1 1972 1977 1978 Article cluster #2 1988 1990 1997





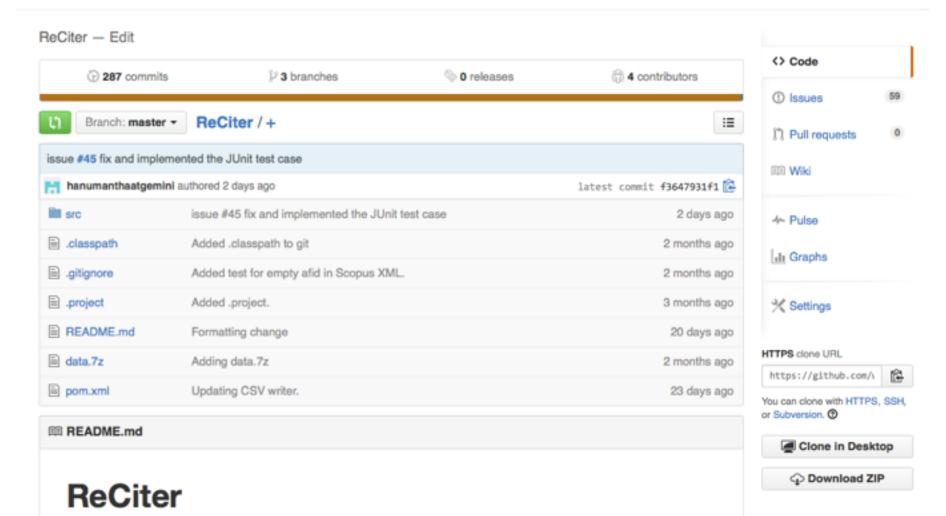






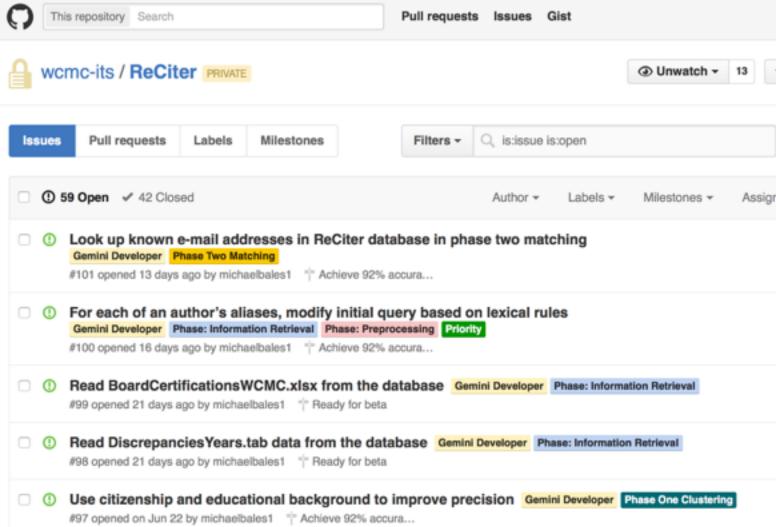




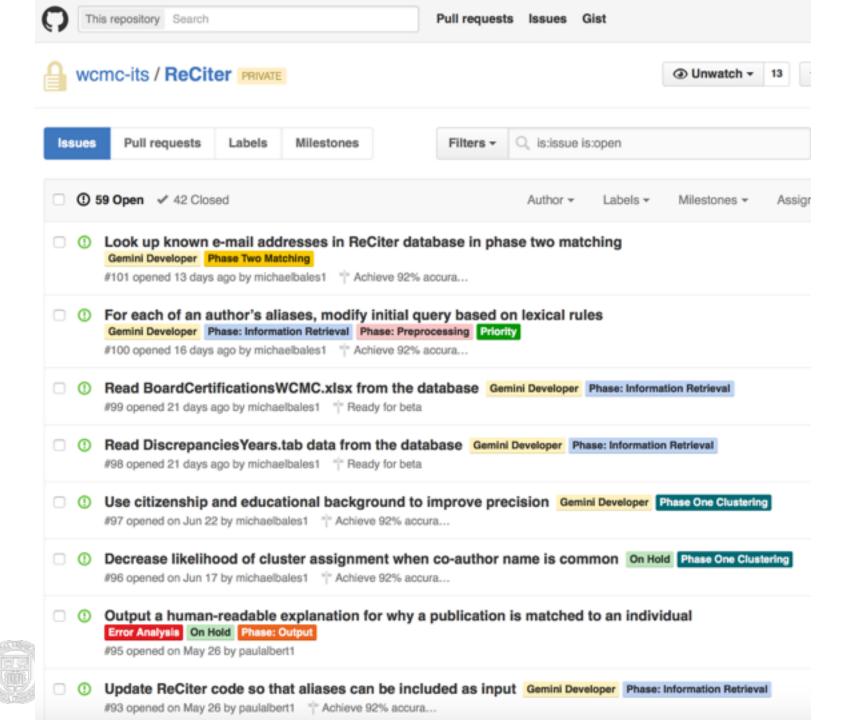


ReCiter wiki

The wiki includes descriptions of files used for computation, an overview of error analysis, a log of performance, and use cases, among other informational material on the project.





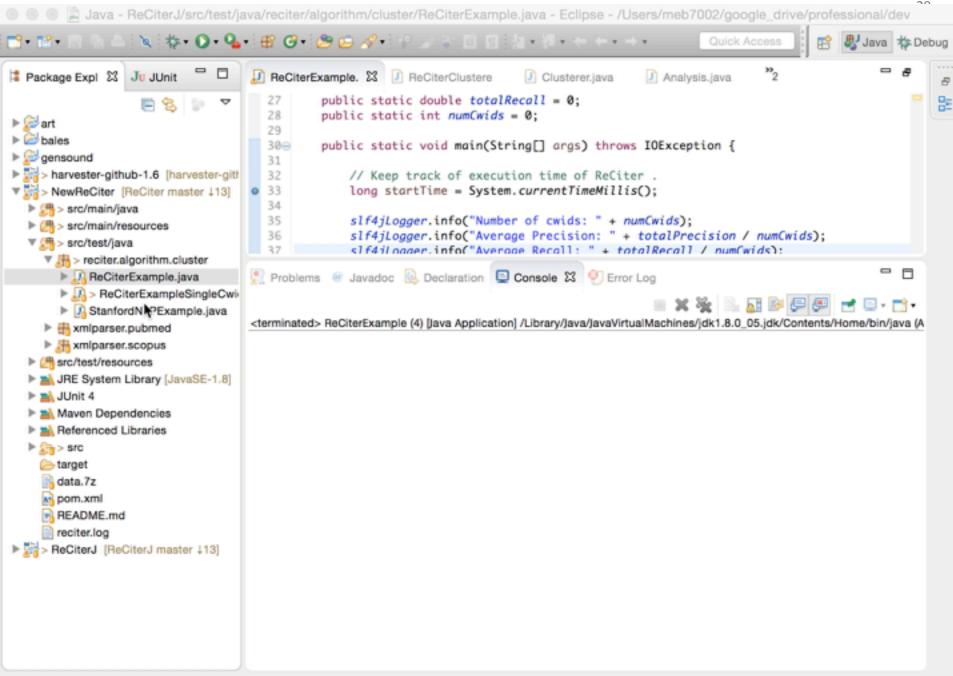


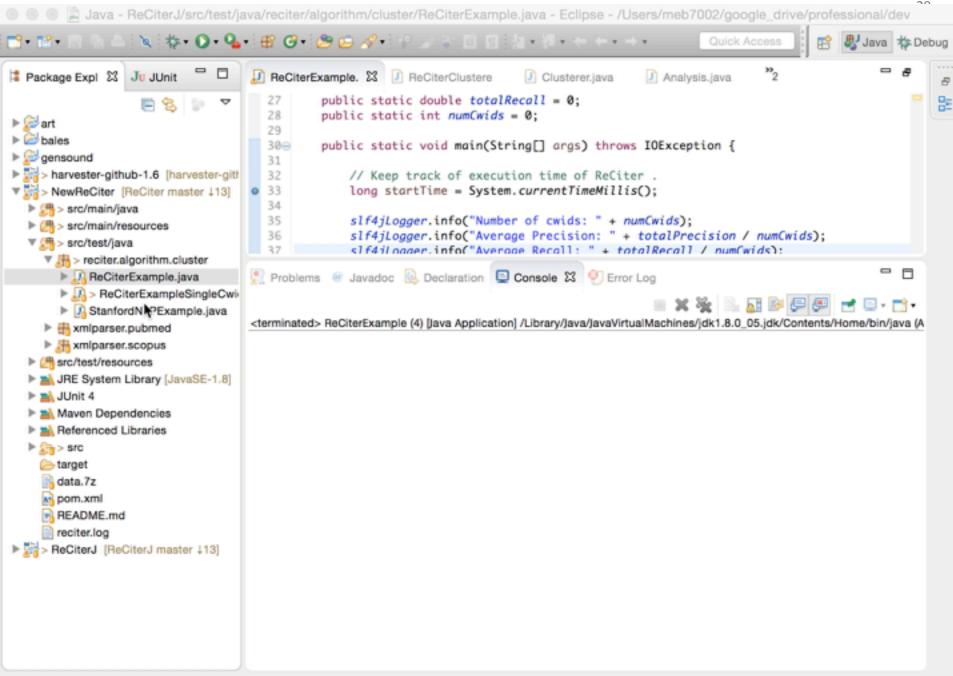
Scope

Туре	Count	Priority
Active faculty	5,500	High
Active students	1,000	High
Postdocs and fellows	400	Medium
Research and staff associates	800	Medium
Alumni	5,000	Medium
Non-WCMC faculty in Graduate School	< 100	Medium
Members of the CTSC including those from CU, MSKCC, NYP, Hunter	> 10,000	Low
Inactive/historical academics	> 10,000	Low

Next steps

- Machine learning
- Open source
- You can help





ReCiter Output (selected fields)

Article ID	Target author	Cluster	Articles in cluster	Cluster ultimately selected	Reference standard status
25313356	aas2004	1	4	No	True Negative
24605052	aas2004	1	4	No	True Negative
19389401	aas2004	1	4	No	True Negative
24767105	aas2004	1	4	No	True Negative
23332979	aas2004	2	2	Yes	False Positive
20489570	aas2004	2	2	Yes	True Positive

ReCiter Team

Name		E-mail
	Paul Albert	paa2013@med.cornell.edu
ı	Michael Bales	meb7002@med.cornell.edu
	Jie Lin	jie265@gmail.com
Balu	ı Mudhavathu	bam3002@med.cornell.edu
Han	umantha Rao	hat3001@med.cornell.edu

HOW MAJOR SYSTEMS TRACKING PUBLICATIONS ARE CONNECTED

