Microsoft Research

Microsoft Academic Graph

Alex Wade
Microsoft Research







- Background
- Bing / Cortana
- The Data
- The Challenge (WSDM Cup)



- Background
 - Project Libra (2005 2008) Wei-Ying Ma, et al.
 - Microsoft Academic Search (2008 2012)
 - → Bing, Cortana (2015)

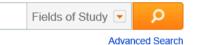
	Bin	g exa	mp	les
--	-----	-------	----	-----

- The Data
- The Challenge (WSDM Cup)

	Papers	Citations	Authors
Libra	~5M	~20M	~1M
MAS	~50M	~350M	~20M
Bing	~125M	~1B	~40M







Keywords **Organizations** Authors **Publications** Conferences **Journals** Top auth **Authors Publications** Conferences **Organizations Journals** Keywords Top p **Authors Publications** Keywords **Organizations** Conferences **Journals** Introdu Indexin Top conferences in information retrieval Informa SIGIR - Research and Development in Inform... ECDL - European Conference on Digital Libr... Term-v **CLEF - Cross-Language Evaluation Forum** CIKM - International Conference on Informa... Lettters TREC - Text REtrieval Conference RIAO - Recherche d'Information Assistee pa... Machin JCDL - ACM/IEEE Joint Conference on Digita... ECIR - European Colloquium on IR Research The Ac DL - Digital Libraries Multimedia Information Retrieval The co See more Automa A vector 3 See more

All Fields of Study

Computer Science

Algorithms & Theory

Artificial Intelligence

Bioinformatics & Computational Biology

Computer Education

Computer Vision

Data Mining

Databases

Human-Computer Interaction

Information Retrieval

Machine Learning & Pattern Recognition

Interaction

Information Retrieval

Machine Learning





Advanced Search

Authors (3516)

W. Bruce Croft

Jamie Callan

James Allan

Cheng-xiang Zhai

Norbert Fuhr

Mark Sanderson

Susan T. Dumais

Justin Zobel

Alistair Moffat

Ryen W. White

Keywords (2238)

Document Retrieval Indexation Information

Need Information

Retrieval System Information Retrieval

Language Model Machine

Learning Query Expansion

Question Answering Relevance

Feedback Search

Engine Test
Collection Text Retrieval

Web Pages Web Search

Academic > Conferences > SIGIR - Research and Development in Information Retrieval







Sort by: Year

Publications (2717)

Improving tweet stream classification by detecting changes in word probability

Kyosuke Nishida, Takahide Hoshide, Ko Fujimura

Conference: Research and Development in Information Retrieval - SIGIR, pp. 971-980, 2012

Combining implicit and explicit topic representations for result diversification

Jiyin He, Vera Hollink, Arjen de Vries

Conference: Research and Development in Information Retrieval - SIGIR, pp. 851-860, 2012

Personalized social query expansion using social bookmarking systems

Mohamed Reda Bouadjenek, Hakim Hacid, Mokrane Bouzeghoub, Johann Daigremont Conference: Research and Development in Information Retrieval - SIGIR, pp. 1113-1114, 2011

Evaluating the Synergic Effect of Collaboration in Information Seeking

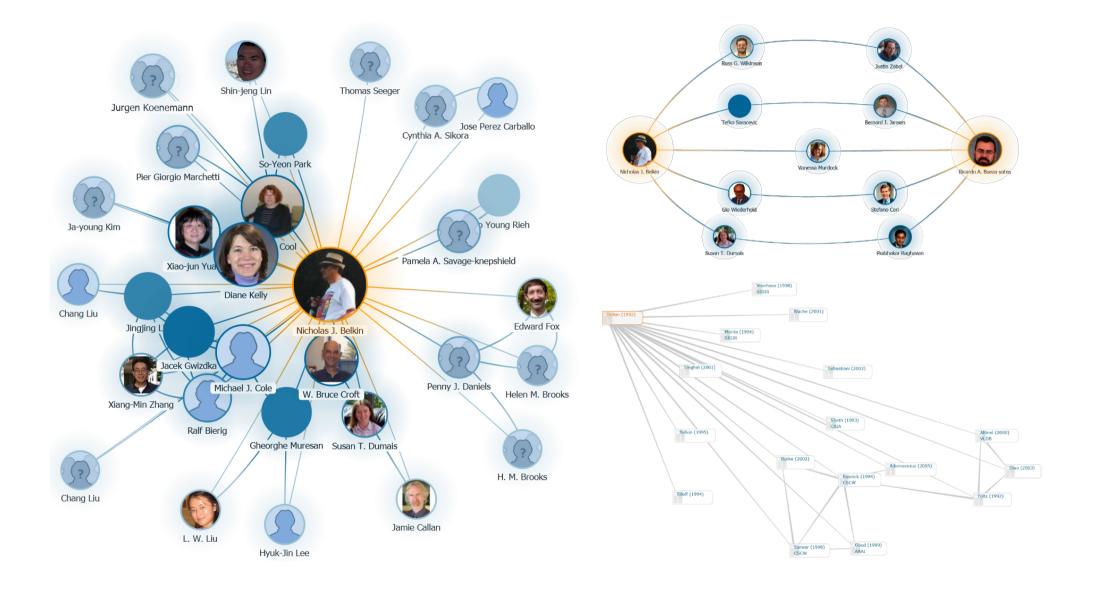
Chirag Shah, Roberto Gonzalez-Ibanez

Conference: Research and Development in Information Retrieval - SIGIR, pp. 913-922, 2011

<u>Jester 2.0</u>: Evaluation of a New Linear Time Collaborative Filtering Algorithm (Citations: 10)

Mark Digiovanni, Hiro Narita, Ken Goldberg

Conference: Research and Development in Information Retrieval - SIGIR, 2010



- Background
 - Project Libra (2005 2008) Wei-Ying Ma, et al.
 - Microsoft Academic Search (2008 2012)
 - → Bing, Cortana (2015)
- Bing examples
- Get the Data
- WSDM Cup

An Overview of Microsoft Academic Service (MAS) and Applications

Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-june (Paul) Hsu, Kuansan Wang Microsoft Research, Redmond, WA 98052, USA

{arsinha, zhihosh, yangsong, haoma, darrine, paulhsu, kuansanw}@microsoft.com

ABSTRACT

In this paper we describe a new release of a Web scale entity graph that serves as the backbone of Microsoft Academic Service (MAS), a major production effort with a broadened scope to the namesake vertical search engine that has been publicly available since 2008 as a research prototype. At the core of MAS is a heterogeneous entity graph comprised of aix types of entities that model the scholarly activities: field of study, author, institution, paper, venue, and event. In addition to obtaining these entities from the publisher feeds as in the previous effort, we in this version include data mining results from the Web index and an in house knowledge base from Bing, from the Web index and an in house knowledge base from Bing, grainor, the new MAS graph sees significant increase in size, with fresh information streaming in automatically following their discoveries by the search engine. In addition, the rich entity relations included in the knowledge base provide additional signals to disambiguate and enrich the entities within and beyond the academic domain. The number of papers indeed by MAS, for instance, has grown from low tens of millions to 83 million while maintaining an above 95% accuracy based on text data set derived from academic activities at Microsoft Research. Based on the data set, we demonstrate disagraph as a feasing securities reactive search and proactive suggestion experience, and a proactive beterogeneous entity recommendation.

Keywords

Academic search; Recommender systems; Entity conflation

1. INTRODUCTION

Recort years have witnessed a paradigm shift in how the knowledge on the Web is made available to the users. The trund is highly visible in the evolution of the Web search engine. The traditional Web search outcomes often serve the users' need at best in a "his ormise" fashion [4, 7]. A multi-year initiative in the industry, called Bing Dalogi in Microsoft [11] and Knowledge Vanit in Google [5], addresses this challenge by using statistical inferences to better organize the Web information and support much richer forms of in-

Copyright is held by the International World Wide Web Conference Committee (IWSC2). IWSC2 reserves the right to provide a hyperlink to the author's site if the Material is used in electronic mediu. WWW 2015 Computation, May 18–22, 2015, Florence, Italy. ACM 978-1-4503-3473-01508. http://dx.doi.org/10.1145/2140908.2742839.

teraction in recognizing and serving the user needs. In addition to reactively retrieving information and answering questions, the model proactively includes additional dialog acts, such as confirmation disambiguation refinement and digression. Coupled with statistical user intent inferences, these acts significantly expedit the process of serving users with the knowledge they need [14]. Our work aims at leveraging this model in addressing the informa-tion needs in areas where the sheer amount of information available through a multitude of channels has exceeded the human canacity in processing them. Although most search engines have provided advanced operators for users to compose elaborated queries to better filter out unwanted materials, their areane syntax has relegated their usages to a negligible rate. A goal of the modern dialog approach to Web search is therefore to utilize advanced technique to enable the search engines to communicate with users in natural language. Because the dialog inferences inevitably require the system to anticipate or predict the needs of the users, another emerg ing trend in the search engine evolution is to extend the prediction behaviors into system initiated notifications. The growing prevalence of mobile personal assistants serve as a natural vehicle to deliver proactive notifications, potentially preempting the needs of user initiated search for information [10].

In this paper, we present two applications in the area of academic publications to demonstrate the potentials of the emerging search paradiem. The first application, described in Section 3.1, illustrates a natural language powered interactive search experience. By lever aging the relationships among the entities in the academic domain the natural language processor is able to harvest the syntactic and semantic cues for parsing and predicting user queries. The second application, described in Section 3.2, demonstrates how a recomnendation system can take advantage of the relationships acros different types of entities to offer heterogeneous suggestions. Noting that the statistical techniques underlying these two applications are by no means perfect, we further decide to make the data used by the two applications publicly available so that the community car jointly attack the challenging unsolved problems. The data set is an update to the corpus previously released for research purposes [2] and will be described in details in Section 2. The two application also exemplify a commonly encountered scenario in which the results presented to the users should be properly ranked. The ranking algorithms and the measurements for determining the ranking or-der remain actively research topics. Given the surge in the count of academic entities and observable limitations of citation count based mpact metrics, the problem of defining meaningful impact metrics of academic entities (e.g. papers, authors, conferences) is gaining substantial interest among the researchers [8, 3]. We hope this open corpus can contribute not only to advance information technologies for other innovative applications but also trigger a new horizon of

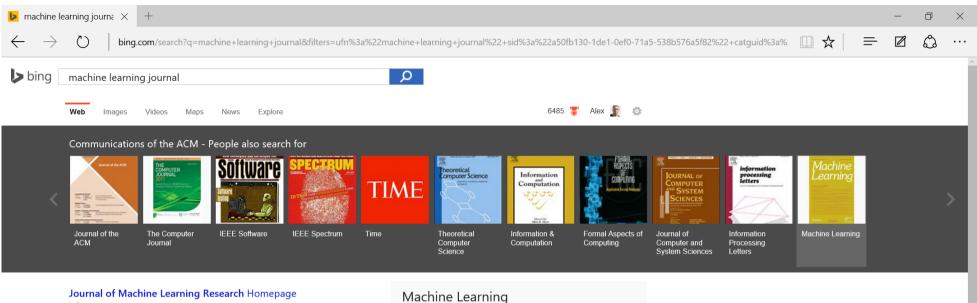
243

- Background
 - Project Libra (2005 2008) Wei-Ying Ma, et al.
 - Microsoft Academic Search (2008 2012)
 - \rightarrow Bing, Cortana (2015)

	Papers	Citations	Authors
Libra	~5M		
MAS	~50M	~350M	~20M
Bing			

- Bing examples
- The Data
- The Challenge (WSDM Cup)





Journal of Machine Learning Research Homepage

Journal of Machine Learning Research The Journal of Machine Learning Research (JMLR) provides an international forum for the electronic and paper publication of high.

Proceedings

Jmlr Volume 14

Lévy ... 703-727, 2013. ..

JMLR: Workshop and Conference

Risk Bounds of Learning Processes for

Proceedings ISSN: 1938-7228. .

Open Source Software Machine Learning Open Source

Software To support the open

Papers

Machine Learning and Large Scale Optimization (Jul 2006 - Oct 2006)

Editorial Board

JMLR Editorial Board Editors-in-Chief Kevin Murphy, Google Bernhard .

Submissions

JMLR seeks previously unpublished papers on machine learning that

See results only from jmlr.org

Machine Learning Journal

www.springer.com/computer/ai/journal/10994 -

Machine Learning is an international forum for research on computational approaches to learning. The journal publishes articles reporting substantive results on a

Machine Learning - Springer

link.springer.com/journal/10994

Machine Learning is an international forum for research on computational approaches to learning. The journal publishes articles reporting substantive results on a .

Machine Learning (journal) - Wikipedia, the free encyclopedia

en.wikipedia.org/wiki/Machine_Learning_(journal) -Machina Learning is a poor reviewed scientific journal nublished since 1086. In 2001

Journal



Machine Learning is a peer-reviewed scientific journal, published since 1986. In 2001, forty editors and members of the editorial board of Machine Learning resigned in order to found the Journal of Machine Learning Research, saying that in th... + en.wikipedia.org

Data from: Wikipedia

Feedback

Related searches

Machine Learning Journal Ranking

Machine Learning Journal Impact Factor

JMLR

JMLR Healthcare

JMLR Impact Factor

Machine Learning Paper

Journal of Learning

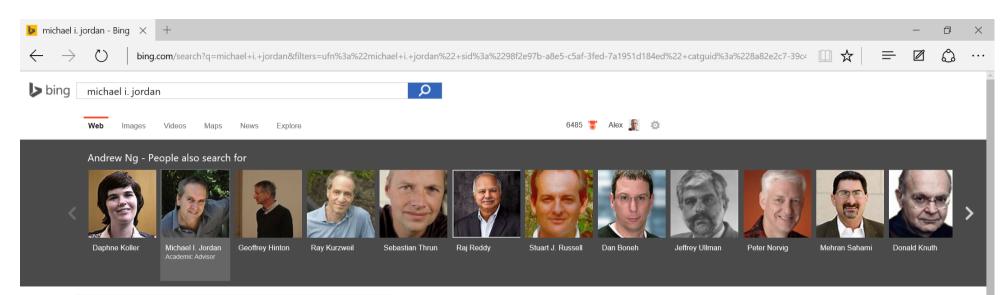
Top Machine Learning Conferences

Twitter

Data Community DC - 2 days ago

Tonight. DC Machine Learning Journal Club discusses journal study

 $http://www.bing.com/search?q=theoretical+computer+science+journal\&filters=ufn\%3a\%22theoretical+computer+science+journal\%2; streaming, \#deeplearning ow.ly/QKN3M_lines_$



Connection to Andrew Ng

Michael I. Jordan and Andrew Ng both wrote Latent dirichlet allocation.

Michael I. Jordan's Home Page

www.cs.berkeley.edu/~jordan ▼

Michael I. Jordan Penong Chen Distinguished Professor Department of EECS Department of Statistics University of California, Berkeley. Emails: EECS Address: Publications · Yuchen Zhang · Martin Wainwright · Andre Wibisono · Ashia Wilson

Michael Jordan - Wikipedia, the free encyclopedia

en.wikipedia.org/wiki/Michael_Jordan >

Michael Jeffrey Jordan (born February 17, 1963), also known by his initials, MJ, is an American former professional basketball player. He is also an entrepreneur, and ... Early years · High school · College · Professional career · Olympic career

Images of michael i. jordan

bing.com/images













Michael I. Jordan - Wikipedia, the free encyclopedia

en.wikipedia.org/wiki/Michael_I._Jordan >

Michael Irwin Jordan (born 1956) is an American scientist, Professor at the University of

Michael I. Jordan

Scientist



Michael Irwin Jordan is an American scientist, Professor at the University of California, Berkeley and leading researcher in machine learning and artificial intelligence. Jordan was born in Ponchatoula, Louisiana, to a working-class fa... +

W Wikipedia

Born: Feb 25, 1956 (age 59) · Louisiana, United States, with Territories

Education: Arizona State University · University of California, San Diego

Academic advisor: David Rumelhart

Written works: Probabilistic Grammars and Hierarchical Dirichlet Processes · Hierarchical Dirichlet Processes · Latent dirichlet allocation

Timeline

1988: He was a professor at MIT from 1988-1998.

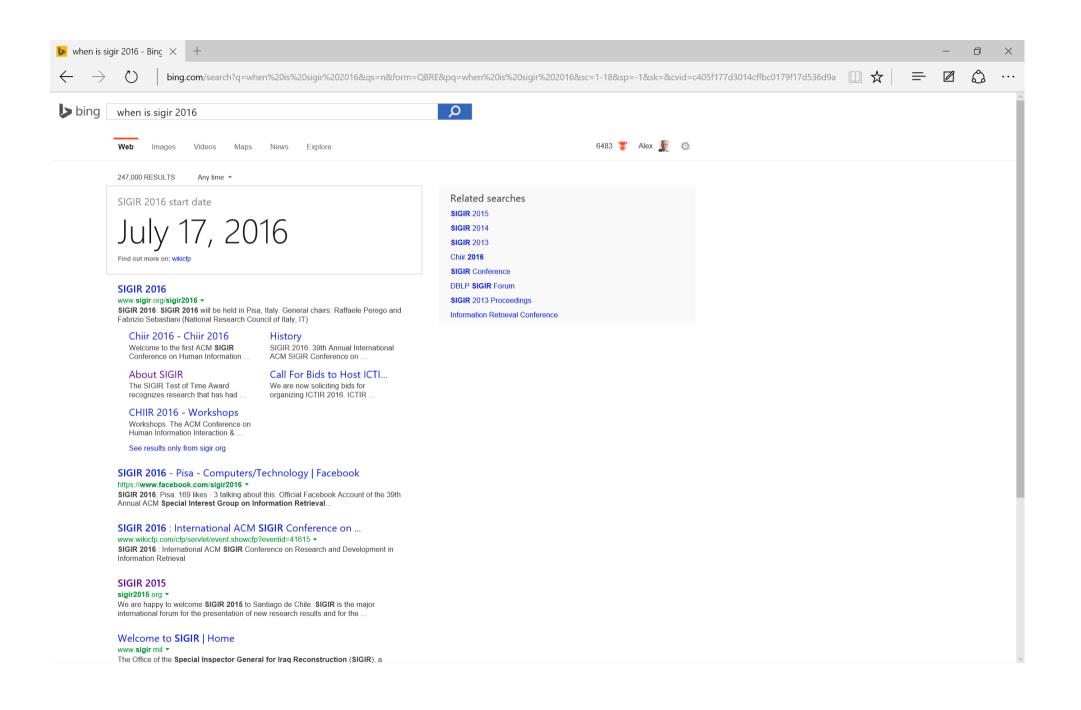
1991: Adaptive mixtures of local experts written by Michael I. Jordan was first published in 1991.

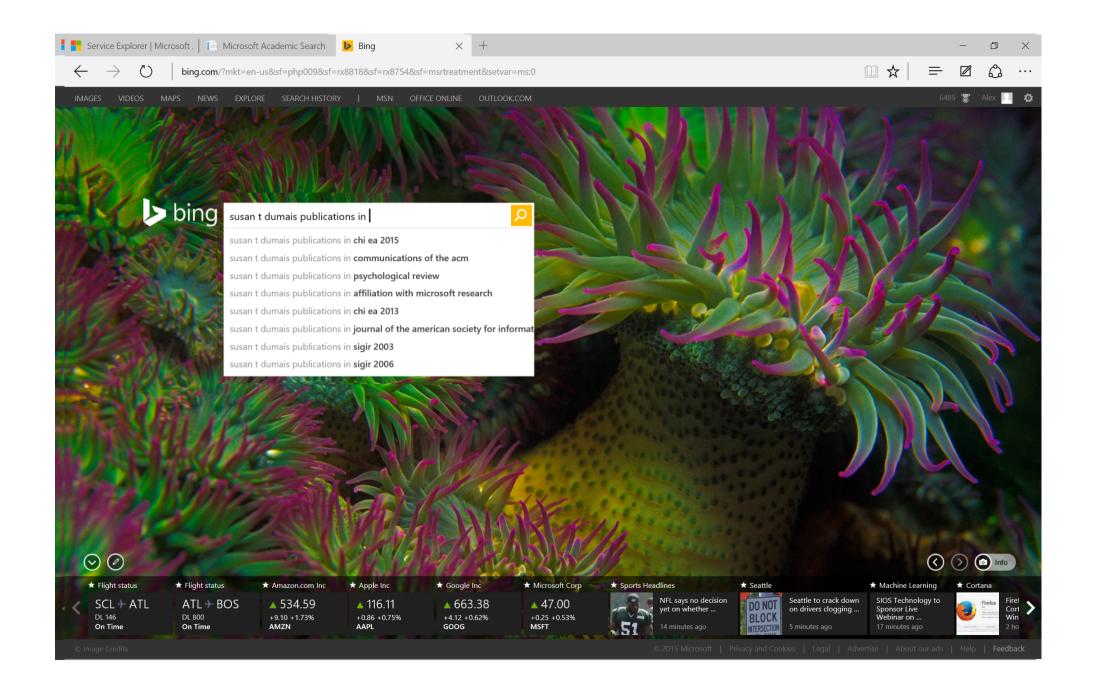
2003: Latent dirichlet allocation written by Michael I. Jordan was first published on March 01, 2003.

See more 💌

Data from: Wikipedia \cdot Microsoft \cdot Berkeley \cdot Freebase

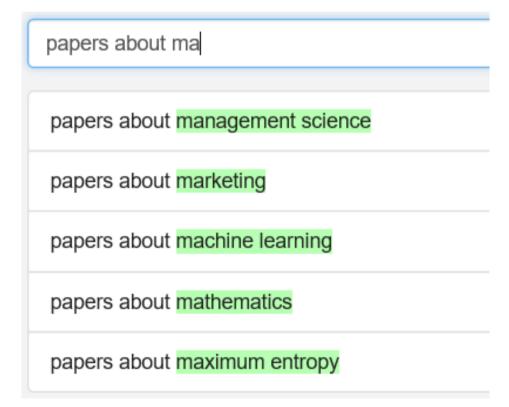
Feedback

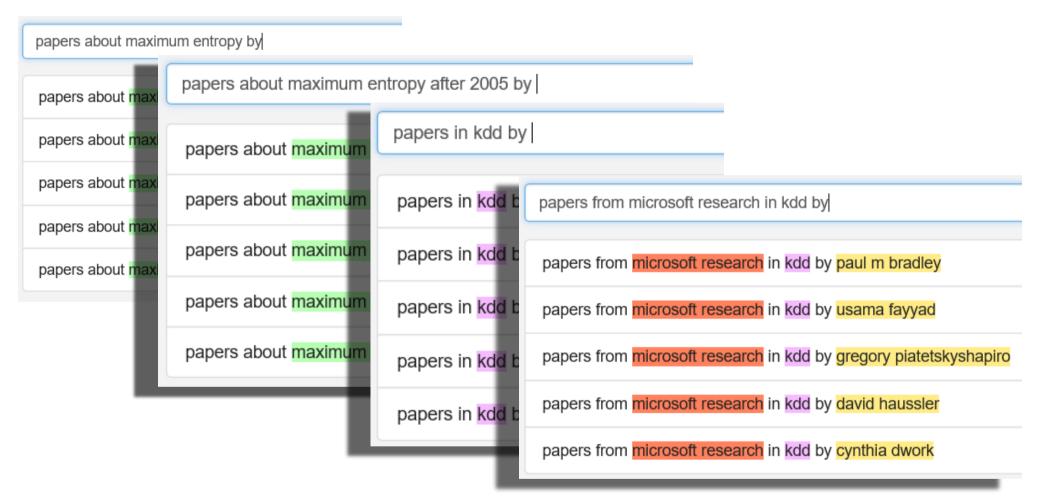




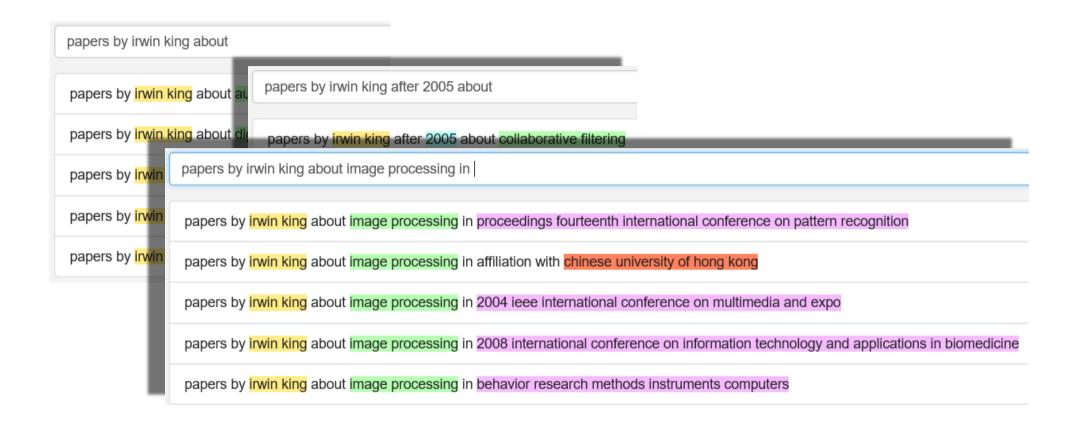
NL Query Completion/Recommendation

- How to complete never foreseen academic queries?
- How to rank completion suggestions?
- How to avoid making completions leading no search results?





Finding reviewers/co-PI made easier



Understanding a colleague made easier

- Background
 - Project Libra (2005 2008) Wei-Ying Ma, et al.
 - Microsoft Academic Search (2008 2012)
 - \rightarrow Bing, Cortana (2015)

	Papers	Citations	Authors
Libra	~5M		
MAS	~50M	~350M	~20M
Bing	~125M		

- Bing examples
- The Data

The Challenge (WSDM Cup)

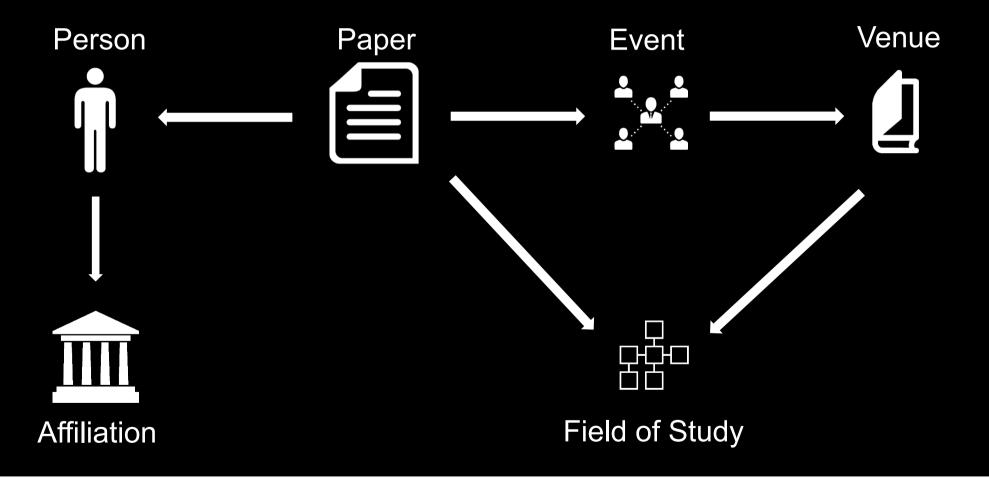
	Papers	Citations	Authors
Microsoft Academic	~40M	~250M	~20M
Microsoft Academic Graph	~100M	~700M	~25M

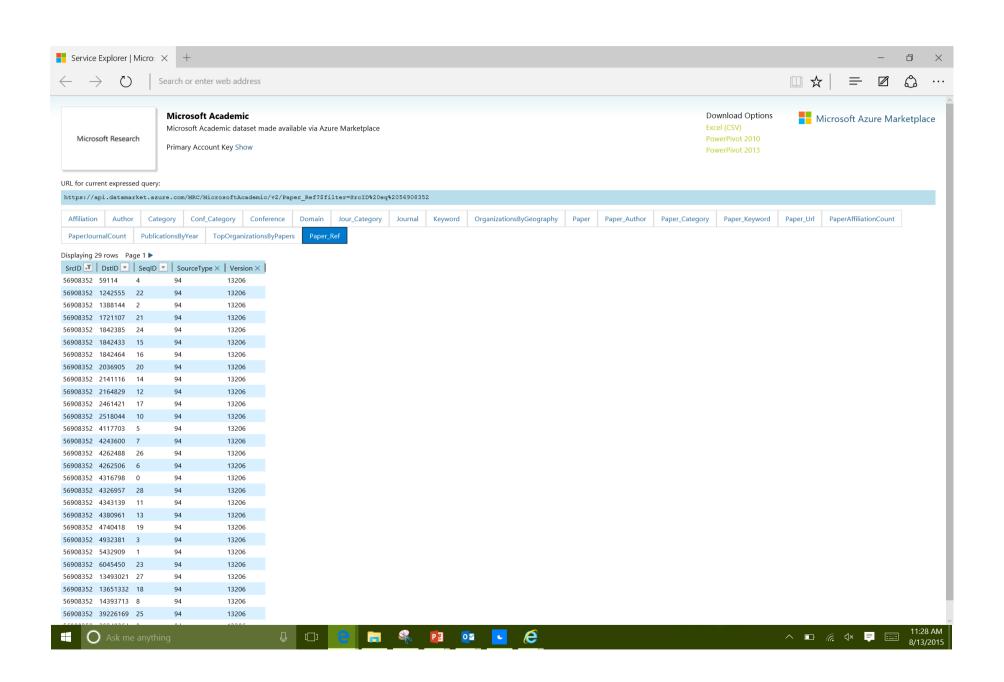


Data Releases

- Microsoft Academic (2012) https://datamarket.azure.com/dataset/mrc/microsoftacademic
 - Azure Datamarket
 - REST API
- Microsoft Academic Graph (2015) http://aka.ms/academicgraph
 - Increased coverage
 - Azure blob storage
 - No API (yet!)
 - Improved topic classification
 - More current (→ early 2015)

	Papers	Citations	Authors
Microsoft Academic	~40M	~250M	~20M
Microsoft Academic Graph	~100M	~700M	~25M





Microsoft Research

Search Microsoft Research

Our research Connections Careers About us

All Downloads Events Groups News People Projects Publications Videos

Microsoft Academic Graph

The Microsoft Academic Graph is a heterogeneous graph containing scientific publication records, citation relationships between those publications, as well as authors, institutions, journals and conference "venues" and fields of study. This data is available as a set of zipped text files stored in Microsoft Azure blob storage and available via HTTP. The file size is ~37GB.

We encourage researchers working with this data to apply for an Azure for Research Award. Details are available on the Azure for Research award submission process at http://azure4research.com. Please include the hashtag #academicgraph in your submission title for easier tracking.

To download the data you need to first agree to the terms of use.

☐ I agree to abide by the terms of use for the Microsoft Academic Graph. Get the data!

Please note: in order to emphasize one important technical challenge that is common in web-scale data collection and aggregation, the data released here have undergone only rudimentary processing, for example in areas of author and paper conflation/deduplication. This noisy yet realistic dataset can provide additional avenues for research in the big data arena.

We kindly request that any published research that makes use of this data cites our data paper listed below.

Publications

 Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, and Kuansan Wang, An Overview of Microsoft Academic Service (MAS) and Applications, WWW – World Wide Web Consortium (W3C), 18 May 2015.



http://aka.ms/academicgraph

- Background
 - Project Libra (2005 2008) Wei-Ying Ma, et al.
 - Microsoft Academic Search (2008 2012)
 - → Bing, Cortana (2015)

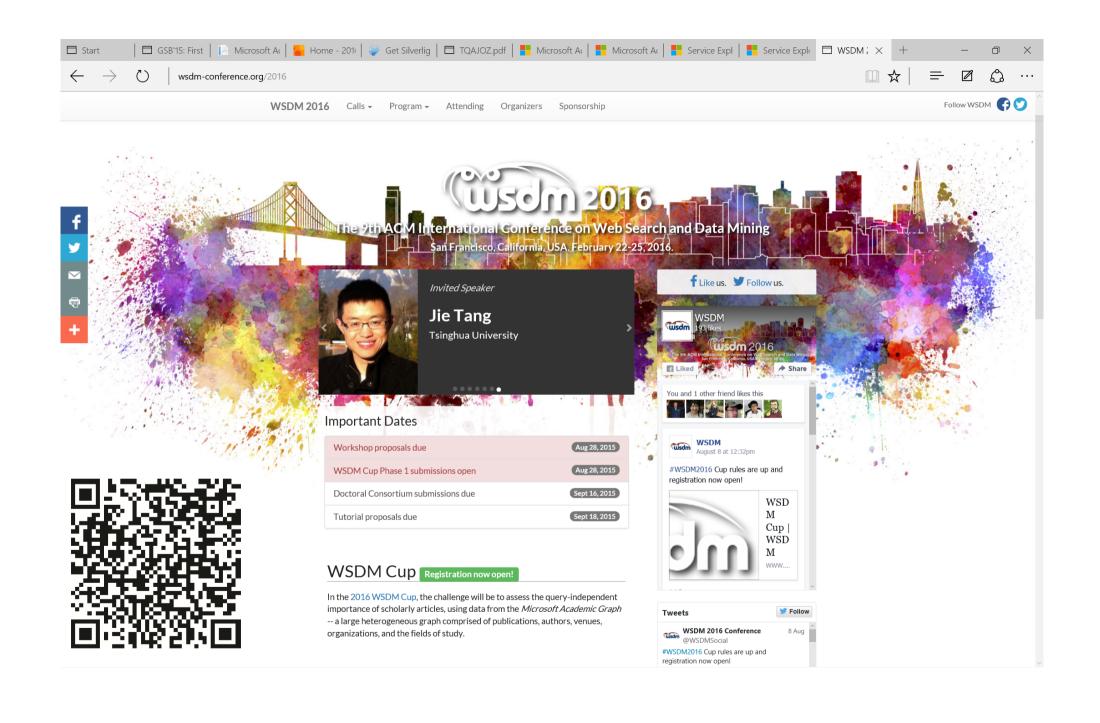
	Papers	Citations	Authors
Libra	~5M		
MAS	~50M	~350M	~20M
Bing	~125M		

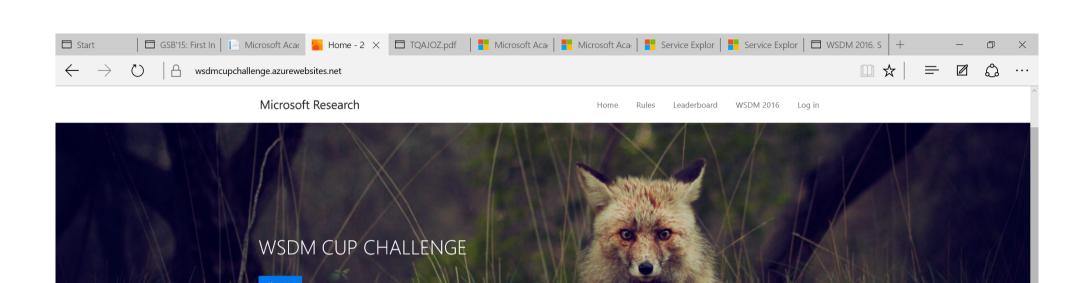
- Bing examples
- The Data

•	The	Challe	enge	(WSDM	Cup)	
---	-----	--------	------	-------	-----	---	--

	Papers	Citations	Authors
Microsoft Academic	~40M	~250M	~20M
Microsoft Academic Graph	~100M	~700M	~25M







The Graph

The Microsoft Academic Graph is a heterogeneous graph containing scientific publication records, citation relationships between publications, as well as authors, institutions, journal and conference "venues," and fields of study.

The Data

This data is available as a set of zipped text files stored in Microsoft Azure blob storage and available via HTTP. The file size (zipped) is ~37GB. The data provided for the challenge may be access and downloaded here.

The Challenge

The goal of the Ranker Challenge is to assess the query-independent importance of scholarly articles, using data from the Microsoft Academic Graph--a large heterogeneous graph comprised of publications, authors, venues, organizations, and the fields of study.



Azure For Research

We encourage researchers working with this data to apply for an Azure for Research Award. Details are available on the Azure for Research award submission process at http://azure4research.com. Please include the hashtag #academicgraph in your submission title for easier tracking.

Microsoft Research

Questions?

Alex Wade

Microsoft Research





