

# **Reputation in Computer Science on a per Subarea Basis**

---

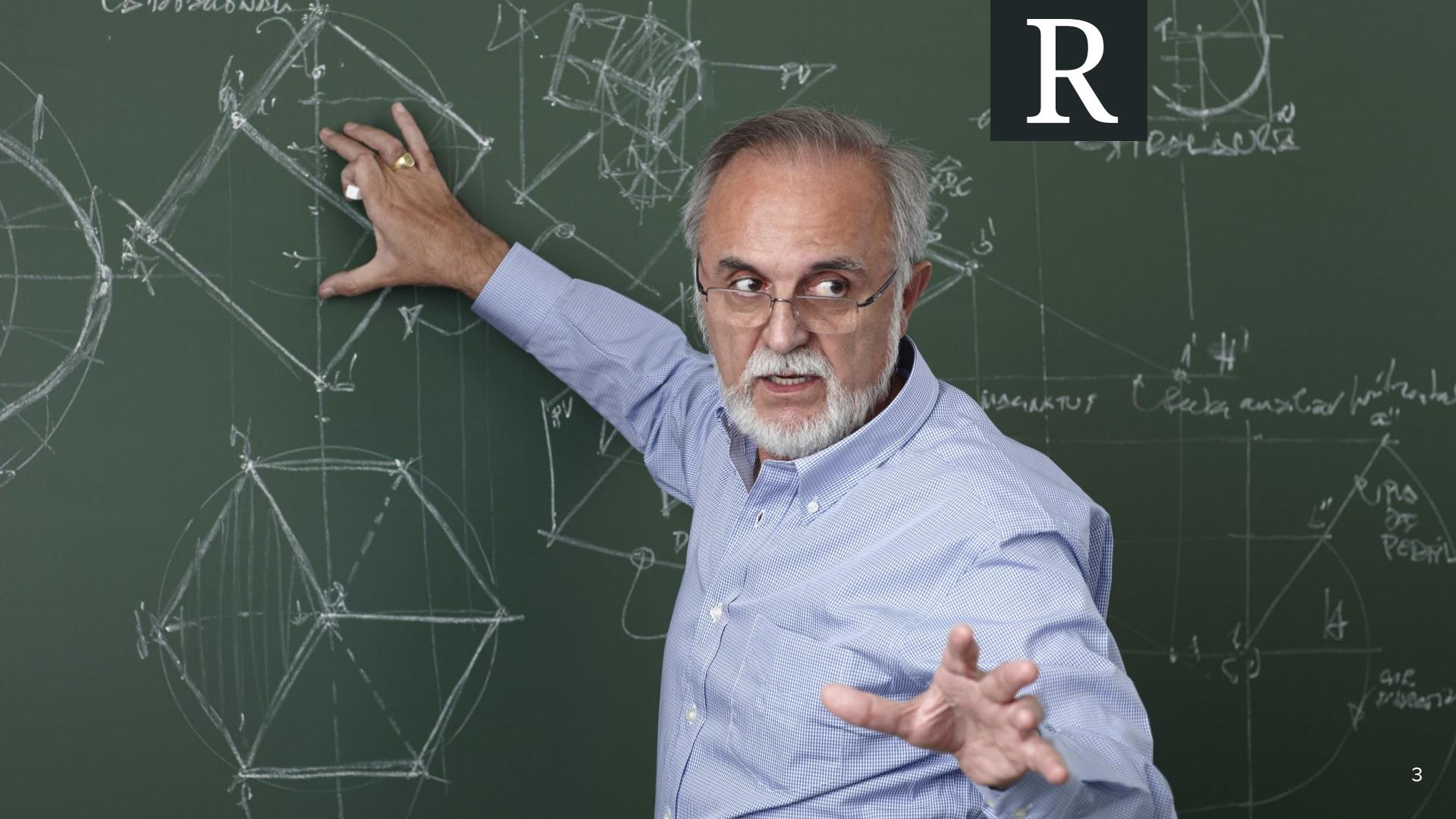
Alberto Hideki Ueda  
Berthier Ribeiro-Neto and Nivio Ziviani

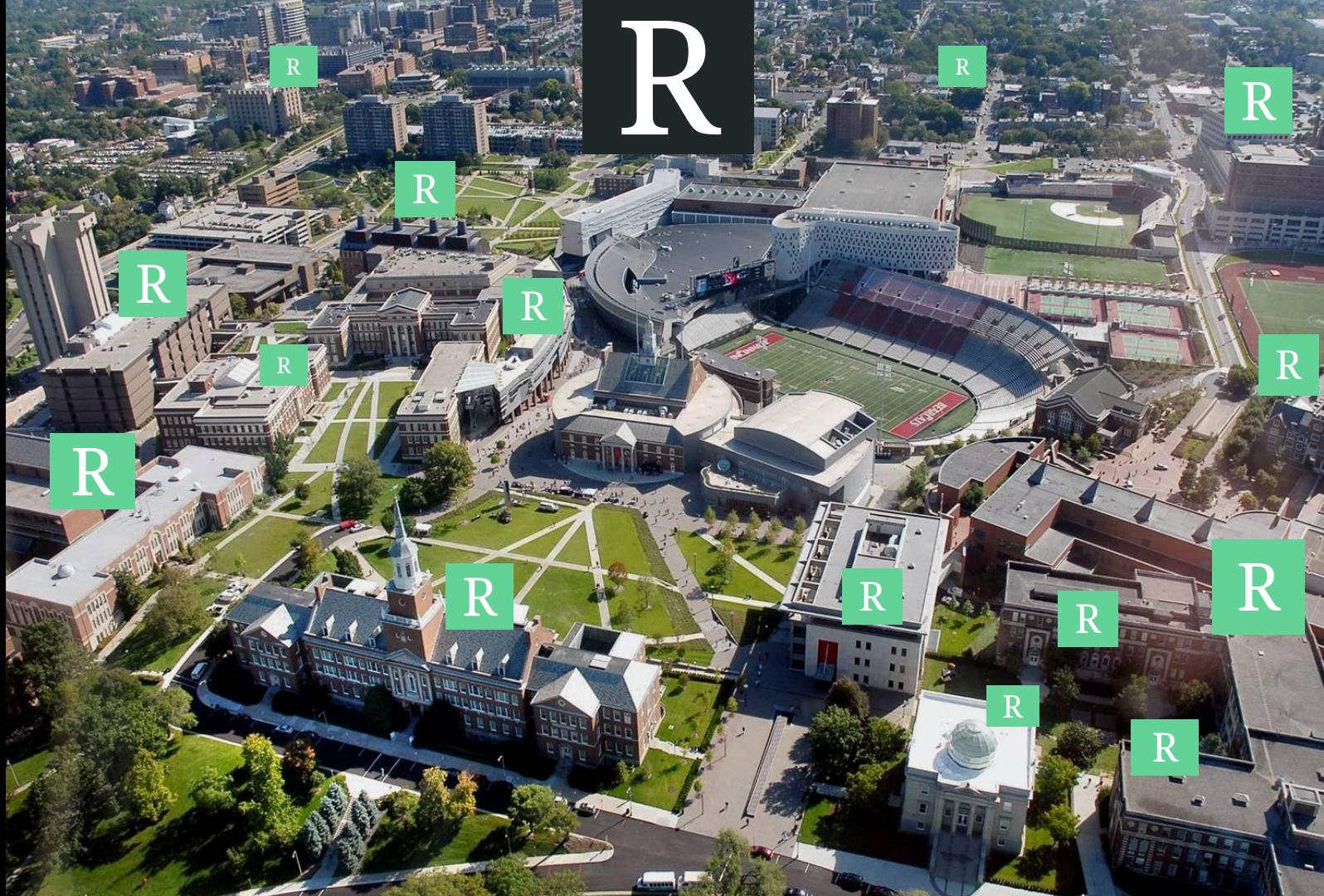
July 14, 2017  
[ueda@dcc.ufmg.br](mailto:ueda@dcc.ufmg.br)

# R



R





# Motivation

Assessing the reputation of academic entities is an important task from different perspectives



\$\$\$

# 37 Subareas in CS

Microsoft Academic, 2017

---

## CS Subareas

---

Algorithm	Internet privacy
Artificial intelligence	Knowledge management
Bioinformatics	Machine learning
Cognitive science	Management science
Computational biology	Mathematical optimization
Computational science	Multimedia
Computer architecture	Natural language processing
Computer graphics	Operating system
Computer hardware	Operations research
Computer network	Parallel computing
Computer security	Pattern recognition
Computer vision	Programming language
Data mining	Real-time computing
Data science	Simulation
Database	Speech recognition
Distributed computing	Telecommunications
Embedded system	Theoretical computer science
Human-computer interaction	World Wide Web
Information retrieval	

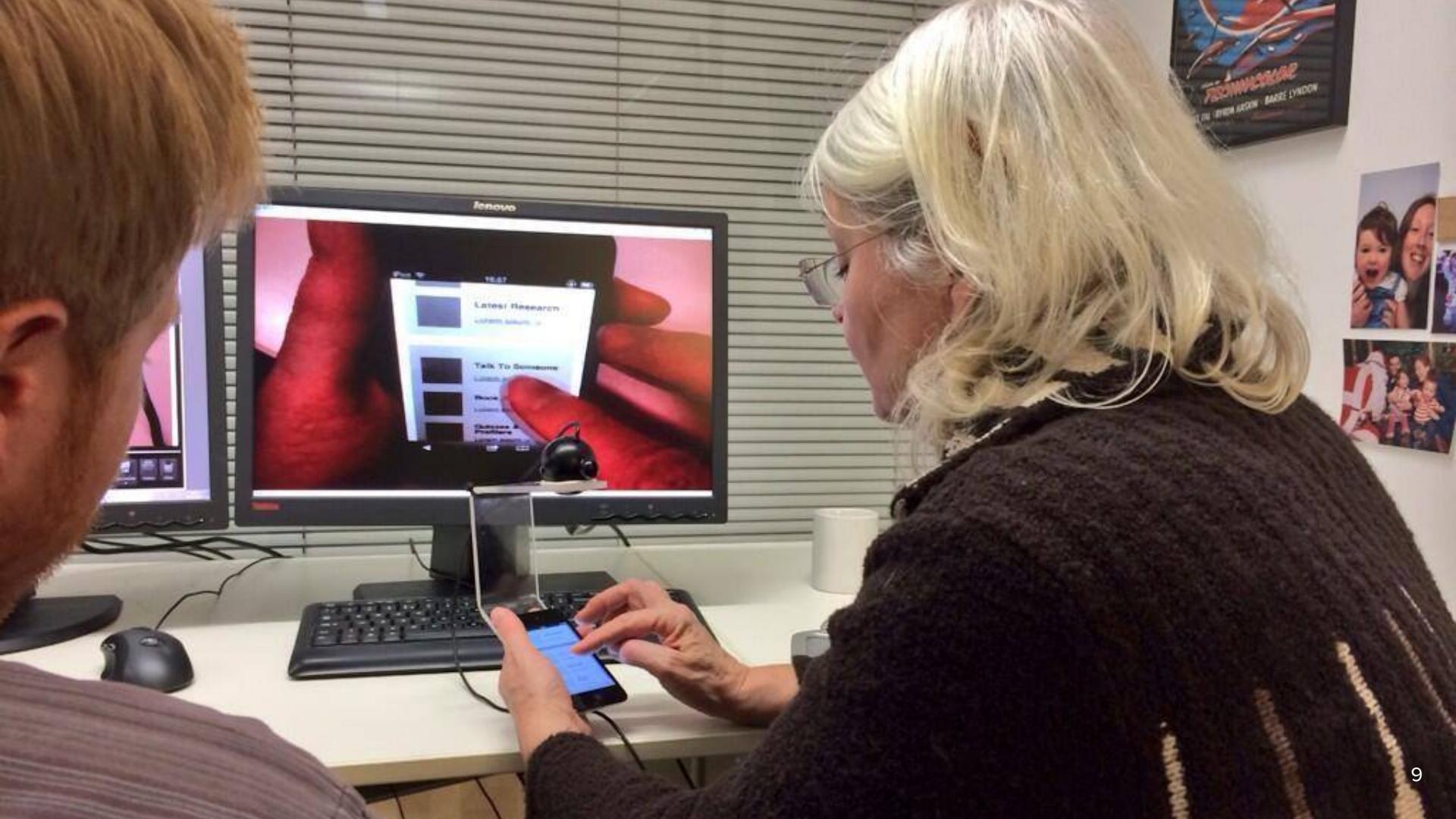
---



Image Processing and Computer Vision  
Average of **5 publications/year** in top venues



Operational Research and Optimization  
**Average of 5 publications/year** in top venues



PROOF OF SPECTRAL THEOREM  
FOR SQUARE MATRICES  
(from P. Lax "Linear Algebra")

Lemma 1: Let  $p \& q$  be 2 polynomials with complex coefficients.  
 Assume  $p \& q$  have no common zero. Then, one can find  
 2 other polynomials  $a \& b$  such that  

$$ap + bq = 1.$$

Proof:  $p \& q$  are given. Consider the set  
 $\mathcal{P} = \{ap + bq \mid a, b \text{ polynomials}\}.$   
 There is at least one non-zero polynomial in  $\mathcal{P}$  of lowest  
 degree. Call  $d = \alpha p + \beta q$  such a polynomial.  
 By the division algorithm, one can write  $p = md + r$   
 where  $m$  and  $r$  are polynomials.  
 Since  $r = p - md = p - m(\alpha p + \beta q)$ , then  $r \in \mathcal{P}$ .  
 But the degree of  $r$  must be strictly less than the degree  
 of  $d$ , which together with  $r \in \mathcal{P}$  implies  $r=0$ . Then,  
 $p=md$ , i.e.  $d$  divides  $p$ . Similarly,  $d$  divides  $q$ .  
 Now we show that  $d$  is a constant. If  $d$  is not  
 of degree 0, then it has at least a root, which is  
 also a root of  $p$  and  $q$ . Since we have assumed that  $p$   
 and  $q$  have no common roots, this is a contradiction.  
 So  $\alpha p + \beta q = d = \text{constant} \Rightarrow \frac{\alpha}{d} p + \frac{\beta}{d} q = 1.$

Lemma 2: Let  $p \& q$  be as in Lemma 1, and let  $A$  be a  
 square matrix with complex entries. Denote by  $N_p$ ,  $N_q$   
 and  $N_{pq}$  the nullspaces of  $p(A)$ ,  $q(A)$  and  $p(A)q(A)$   
 respectively.

# 37 Subareas in CS

Microsoft Academic, 2017

---

## CS Subareas

---

Algorithm	Internet privacy
Artificial intelligence	Knowledge management
Bioinformatics	Machine learning
Cognitive science	Management science
Computational biology	Mathematical optimization
Computational science	Multimedia
Computer architecture	Natural language processing
Computer graphics	Operating system
Computer hardware	Operations research
Computer network	Parallel computing
Computer security	Pattern recognition
Computer vision	Programming language
Data mining	Real-time computing
Data science	Simulation
Database	Speech recognition
Distributed computing	Telecommunications
Embedded system	Theoretical computer science
Human-computer interaction	World Wide Web
Information retrieval	

---

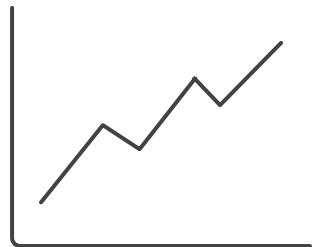
# R



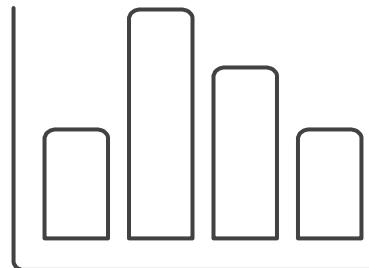


# Reputation in CS on a per Subarea Basis

How to  
**quantify** it?



How does it **vary**  
per CS subarea?



How does it **differ**  
in BR and US?



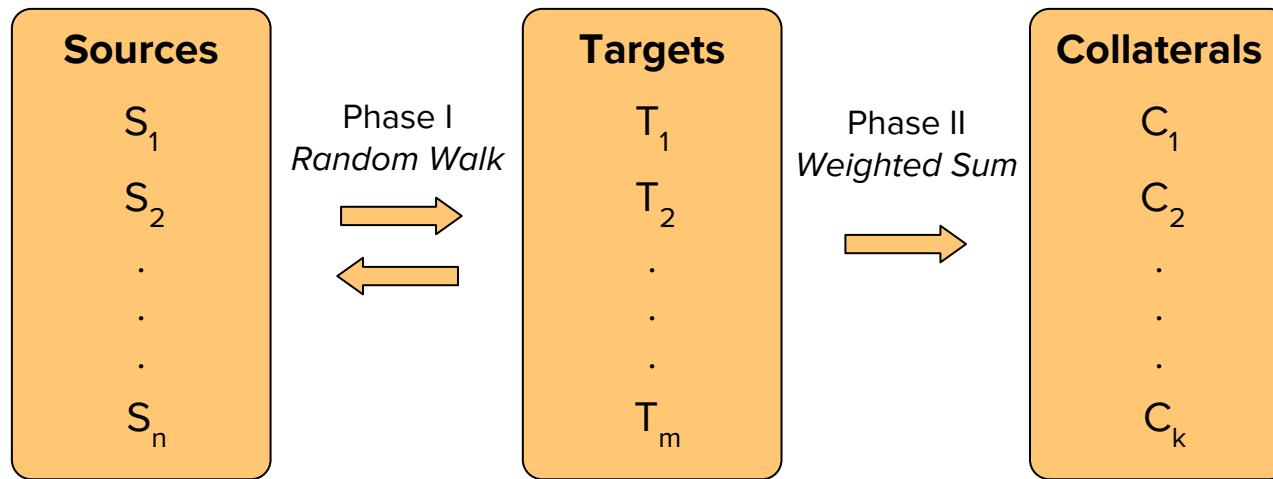
# Our hypothesis

Using P-score on a per subarea basis we can improve rankings of academic entities.

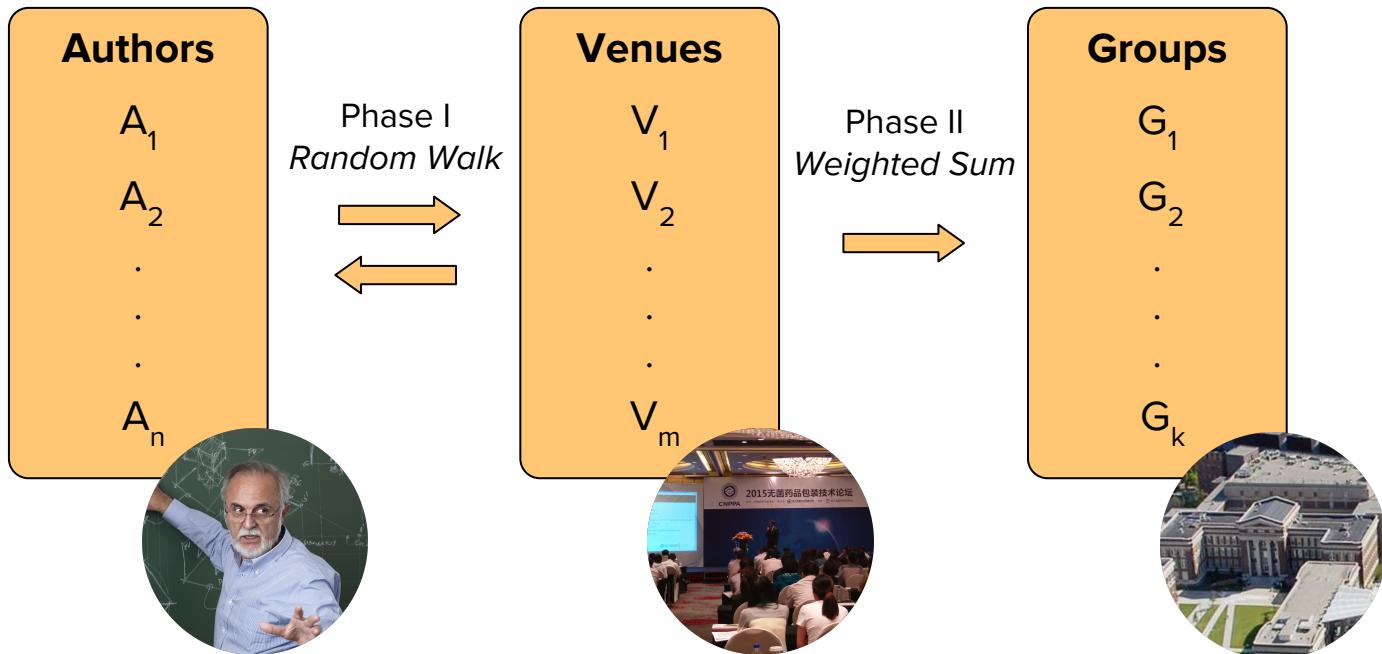
# Reputation Flows



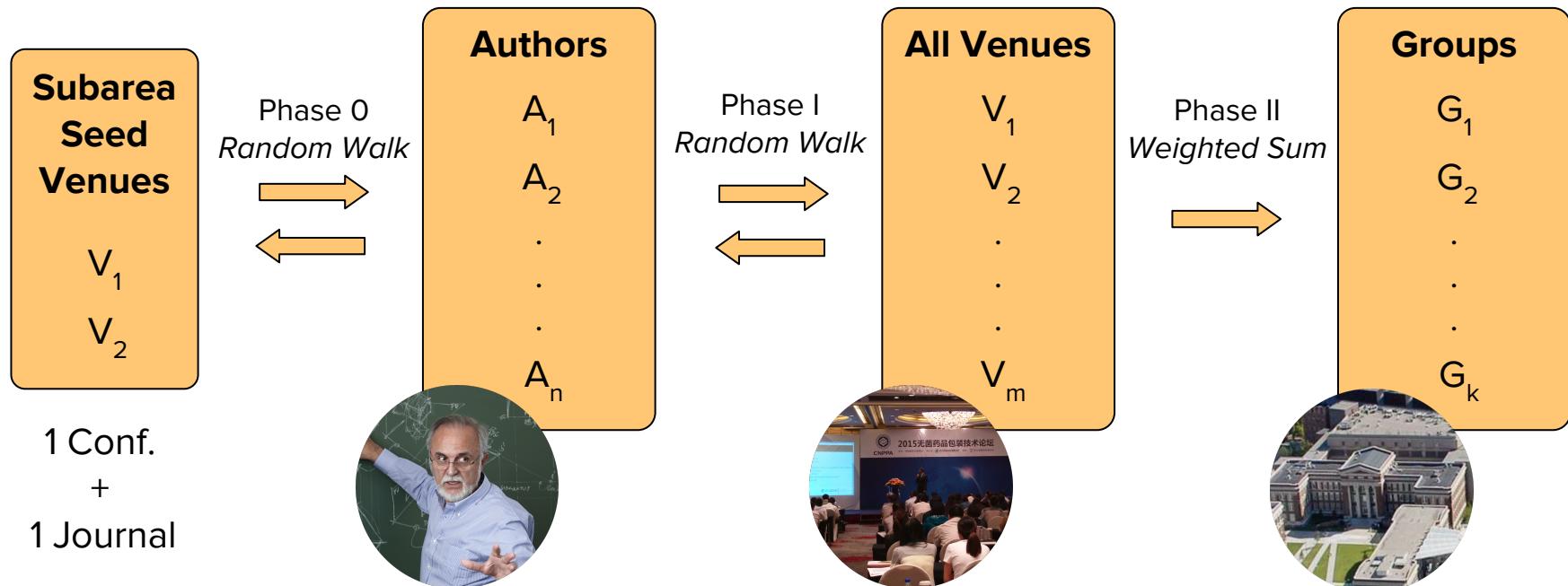
# Reputation Flows



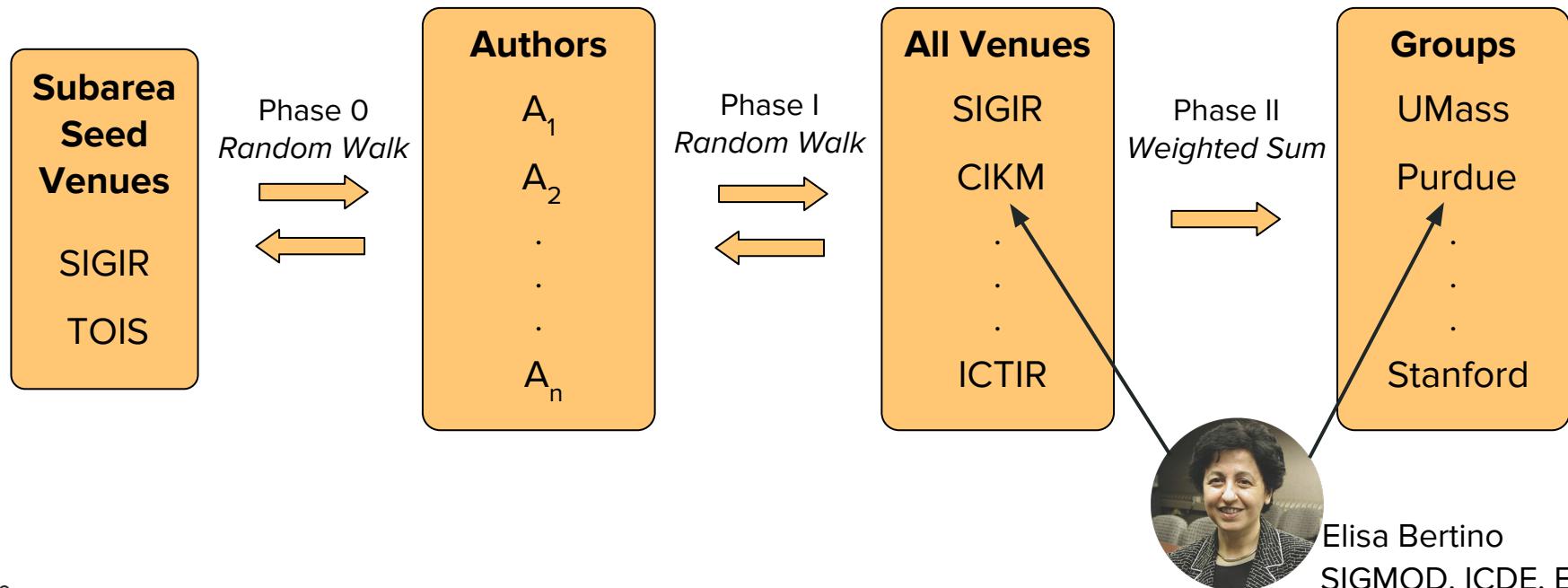
# Reputation Flows



# Standard P-score

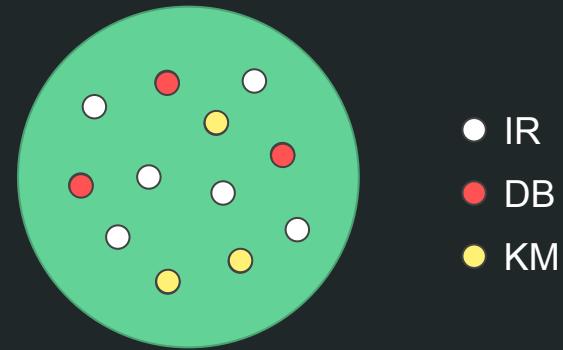


# Standard P-score for IR



# The Encroachment Problem

Publication venues possibly  
cover multiple subareas.



# Normalized P-score for Venues



$$norm\text{-}P\text{-score}(v) = \frac{P\text{-score}(v)}{number\_of\_publications(v)}$$

$$V_s = \{v \in V \mid norm\text{-}P\text{-score}(v) > \alpha_s\}$$

$v_1, v_2, \dots, v_n$



$$\gamma(a, s) = \frac{1}{|\theta(a)|} \sum_{p \in \theta(a)} \begin{cases} \frac{1}{\log_2(y(0)-y(p)+2)} & \text{if } p \in V_s \\ 0 & \text{otherwise} \end{cases}$$

Factor  $\gamma$  of author  $a$  for subarea  $s$

Set of publications of author  $a$

Year of publication of paper  $p$

# Weighted P-score for Graduate Programs

$$P\text{-score}(g) = \sum_{a \in \phi(g)} \sum_{p \in \theta(a)} \frac{P\text{-score}(venue(p))}{number\_of\_authors(p)}$$

$$weighted\text{-}P\text{-score}(g, s) = \sum_{a \in \phi(g)} \gamma(a, s) \times \sum_{p \in \theta(a)} \frac{P\text{-score}(venue(p))}{number\_of\_authors(p)}$$

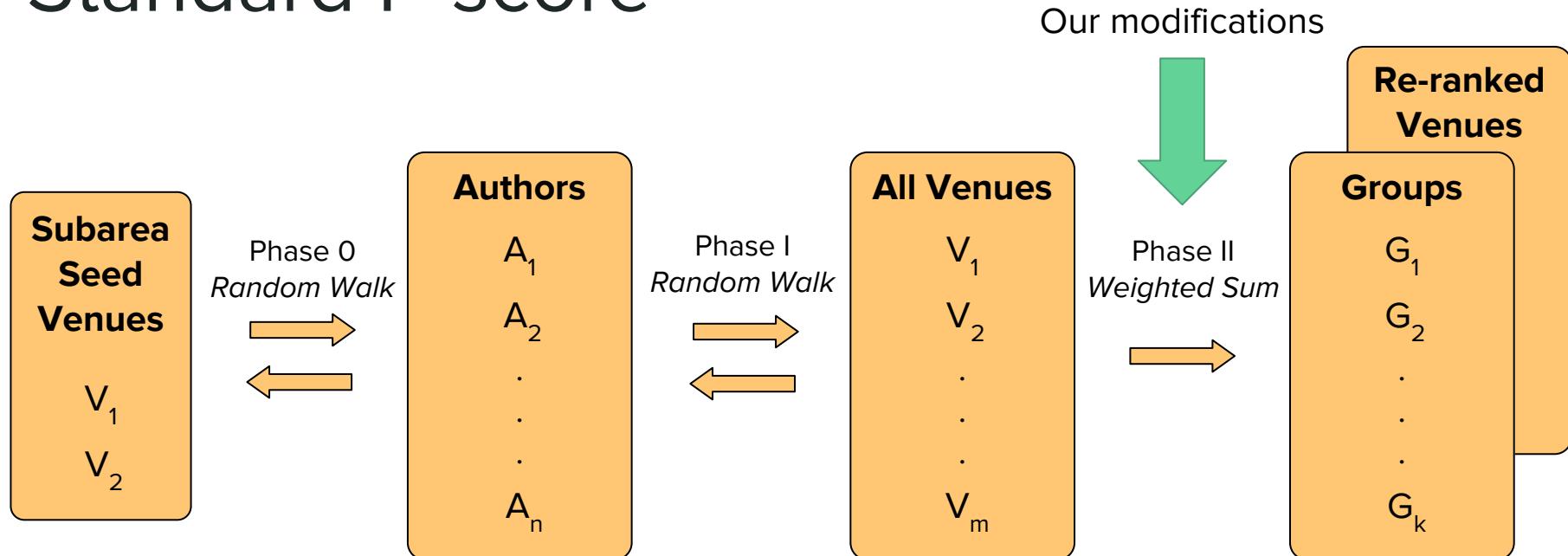
Set of researchers of group  $g$

Factor  $\gamma$  of author  $a$  for subarea  $s$

Set of publications of author  $a$



# Standard P-score



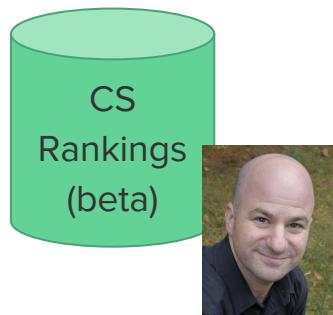
# Experimental Setup

Dataset

Seed Venues

Venues Ground-truth

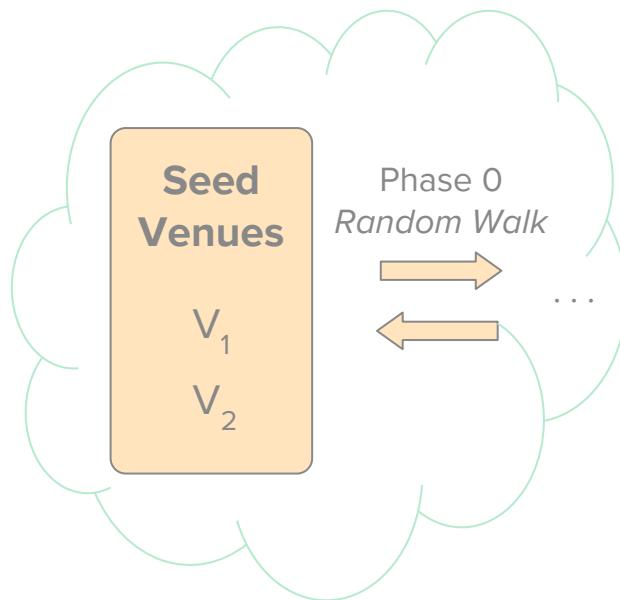
# Dataset



Emery Berger  
UMass

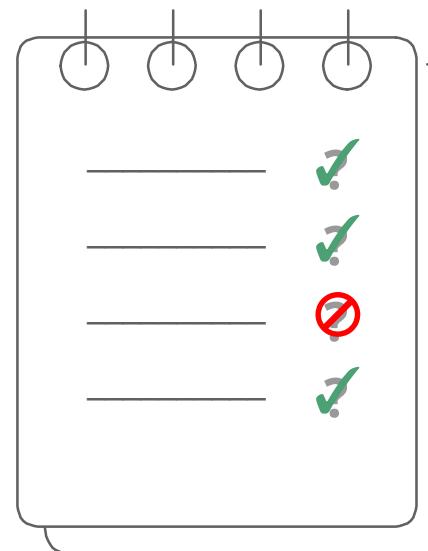
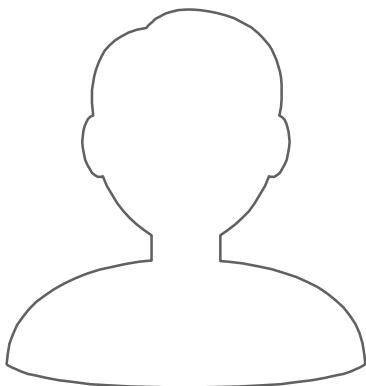
Attribute	Value
Number of Papers	2,931,849
Number of Authors	1,595,771
Number of Venues	5,765
Number of US Depts	126
Avg. number of professors per US Dept	42.4
Number of BR Depts	25
Avg. number of professors per BR Dept	47.8

# Seed Venues

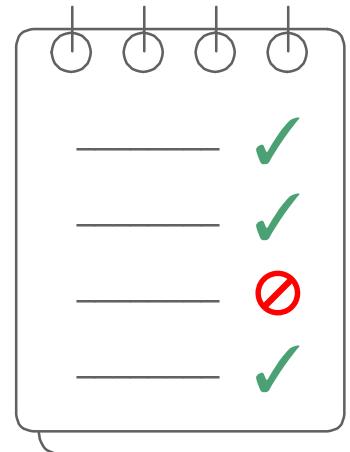
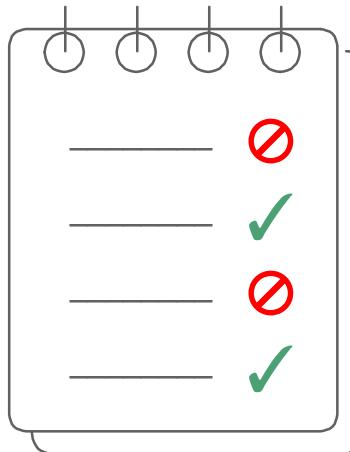
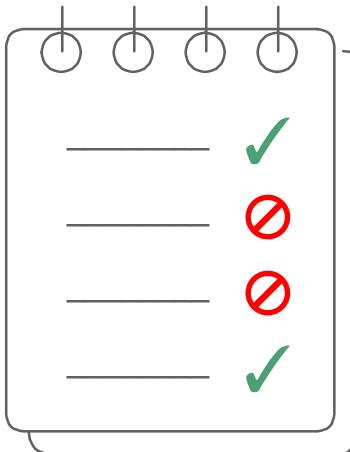
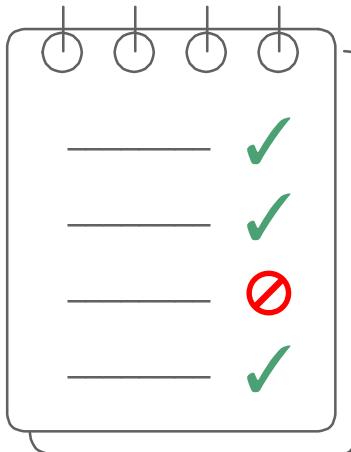


Subarea	Abbreviation	Seed venues
Algorithm	Alg	SODA, Algorithmica
Artificial intelligence	AI	IJCAI, AI
Bioinformatics	Bio	BIBM, Bioinformatics
Computer graphics	CG	SIGGRAPH, TCVG
Computer network	CN	INFOCOM, TON
Computer security	CS	CCS, TISSEC
Computer vision	CV	CVPR, IJCV
Data mining	DM	KDD, SIGKDD
Database	DB	SIGMOD, TODS
Distributed computing	DC	ICDCS, TPDS
Human-computer interaction	HCI	CHI, TOCHI
Information Retrieval	IR	SIGIR, TOIS
Machine learning	ML	ICML, JMLR
Natural language processing	NLP	EMNLP, COLING
Operating system	OS	SOSP, SIGOPS
Parallel computing	PC	IPPS, TPDS
Programming language	PL	PLDI, TOPLAS
Speech Recognition	SR	INTERSPEECH, TCOM
Theoretical computer science	TCS	STOC, SIAMCOMP
World Wide Web	WWW	WWW, WS

# Venues Ground-Truth

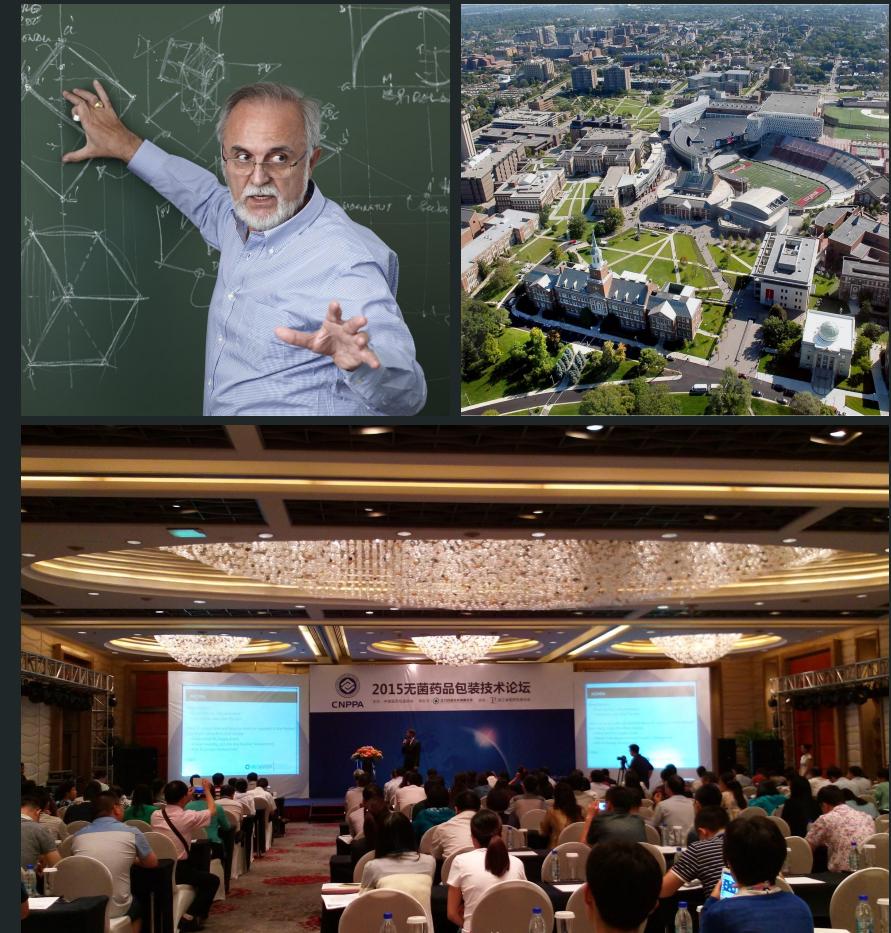


# Venues Ground-Truth

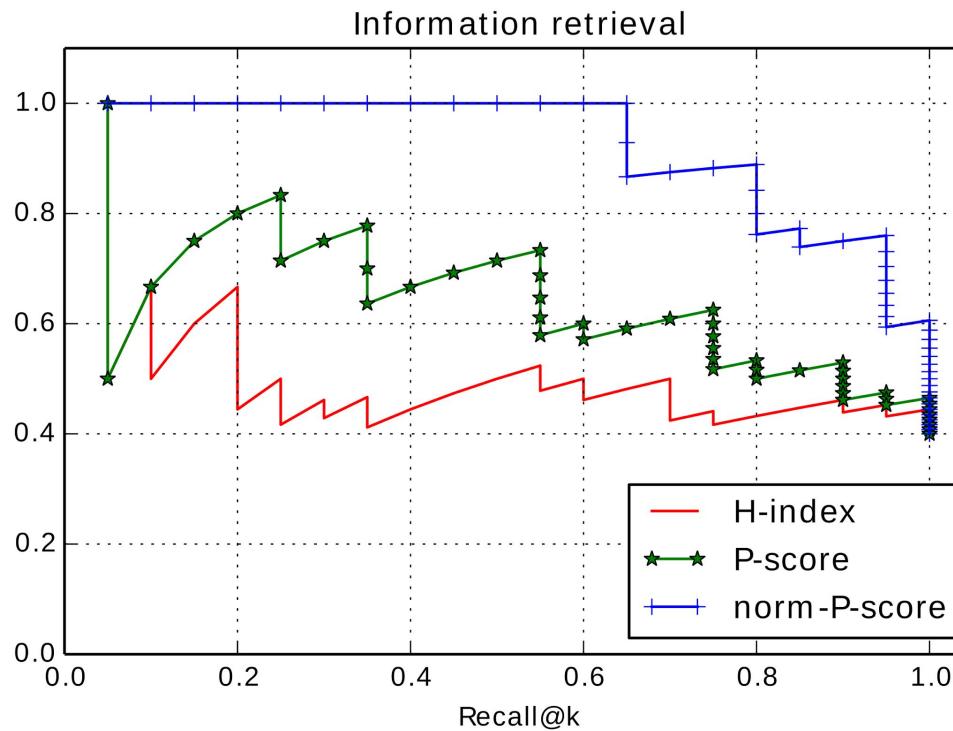


# Experimental Results

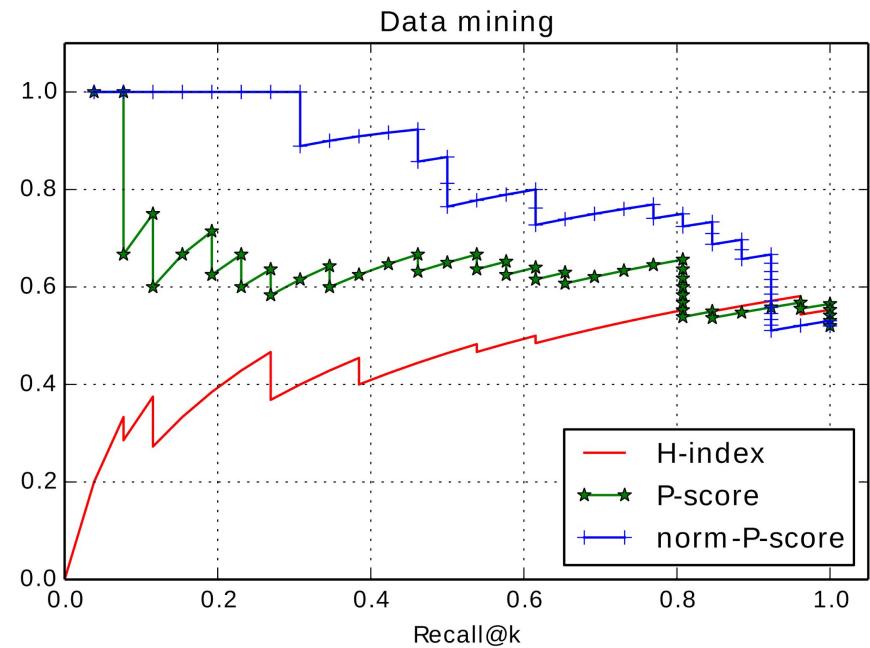
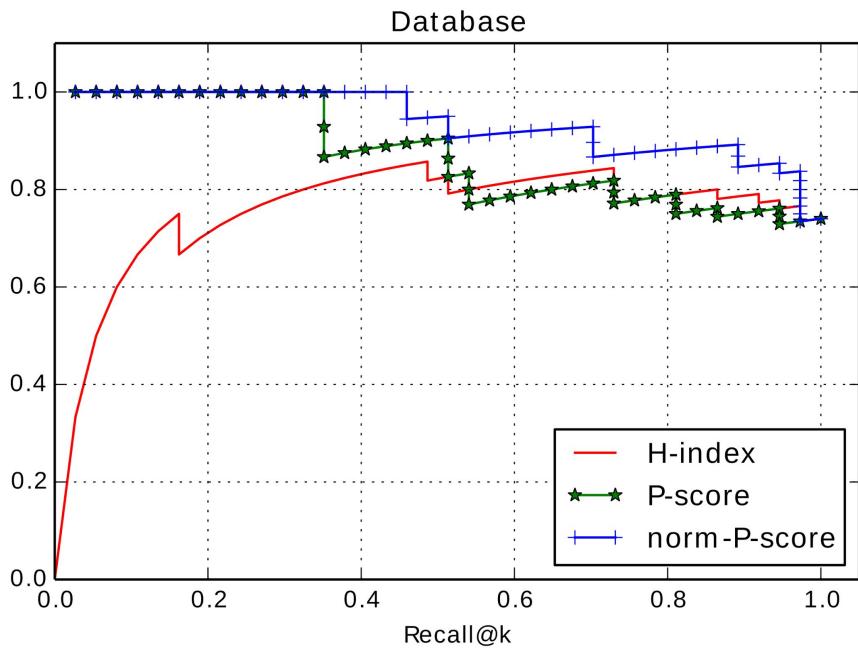
Ranking publication venues and  
graduate programs in BR and US



# Identifying Venues of a Subarea



# Identifying Venues of a Subarea



# Ranking of Venues (IR)

#	Standard P-score	norm-P-score	#	Final Ranking
1	SIGIR (c)	SIGIR (c)	1	SIGIR (c)
2	CIKM	CIKM	2	CIKM
3	TREC	ICTIR	3	TREC
4	ECIR	IIIIX	4	ECIR
5	CLEF	TREC	5	CLEF
6	WWW	SPIRE	6	SIGIR (j)
7	JASIS	ECIR	7	JCDL
8	IPM	ADCS	8	TOIS
9	SIGIR (j)	RAIO	9	IR
10	MM	IR	10	WSDM
11	JCDL	AIRS	11	NTCIR
12	TOIS	NTCIR	12	SPIRE
13	IR	INEX	13	AIRS
14	WSDM	WSDM	14	RAIO
15	NTCIR	TOIS	15	INEX
16	KDD	JCDL	16	IIIIX
17	TKDE	SIGIR (j)	17	ICTIR
18	ACL	TWEB	18	ADCS
19	ICDM	CLEF	19	LA-WEB
20	SPIRE	LA-WEB	20	TWEB

# Ranking of Graduate Programs (IR)

#	University	P-score
1	Carnegie Mellon University	1
2	University of Massachusetts Amherst	0.8082
3	University of Illinois at Urbana-Champaign	0.6735
4	University of Southern California	0.4541
5	Georgia Institute of Technology	0.4341
6	Stanford University	0.3493
7	University of Illinois at Chicago	0.3409
8	Cornell University	0.3344
9	University of California-Berkeley	0.3337
10	Purdue University	0.3120

#	University	weighted-P-score
1	University of Massachusetts Amherst	1
2	University of Illinois at Urbana-Champaign	0.4830
3	Carnegie Mellon University	0.4625
4	University of Delaware	0.2452
5	Purdue University	0.2276
6	Northeastern University	0.1633
7	Lehigh University	0.0964
8	Cornell University	0.0552
9	University of Iowa	0.0494
10	University of Illinois at Chicago	0.0477

# Ranking of Graduate Programs (IR)

---

## # Authors' universities

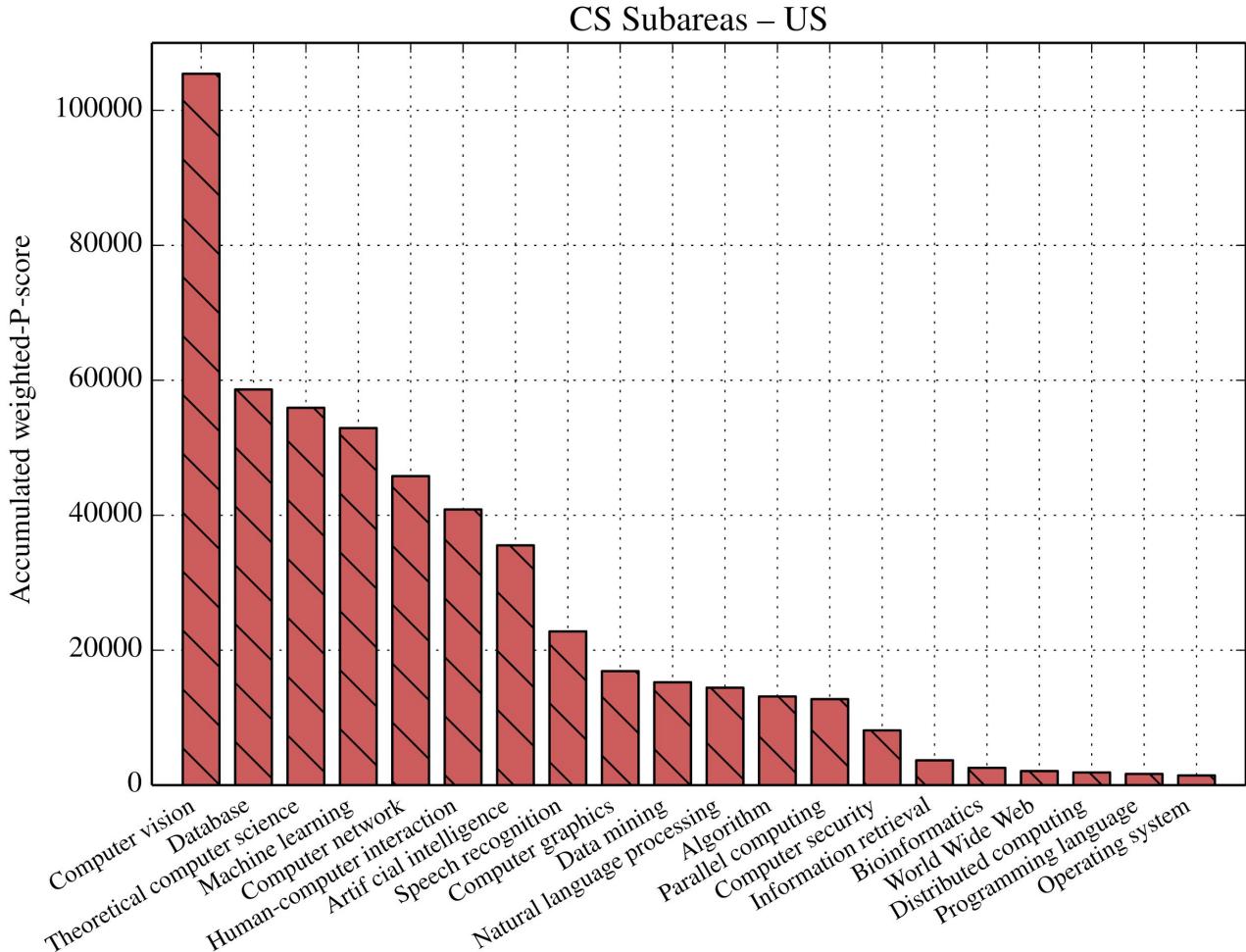
---

- 1 University of Massachusetts Amherst #1
  - 2 University of Massachusetts Amherst #2
  - 3 Carnegie Mellon University #1
  - 4 University of Illinois at Urbana-Champaign #1
  - 5 Purdue University
  - 6 University of Delaware
  - 7 Northeastern University
  - 8 University of Illinois at Urbana-Champaign #2
  - 9 Lehigh University
  - 10 Carnegie Mellon University #2
  - 11 University of Iowa
  - 12 University of Illinois at Chicago
  - 13 Georgia Institute of Technology
  - 14 University of Virginia
  - 15 Carnegie Mellon University #3
  - 16 Texas A&M University
  - 17 Cornell University
  - 18 University of Michigan
  - 19 University of Massachusetts Amherst #3
  - 20 New York University
- 

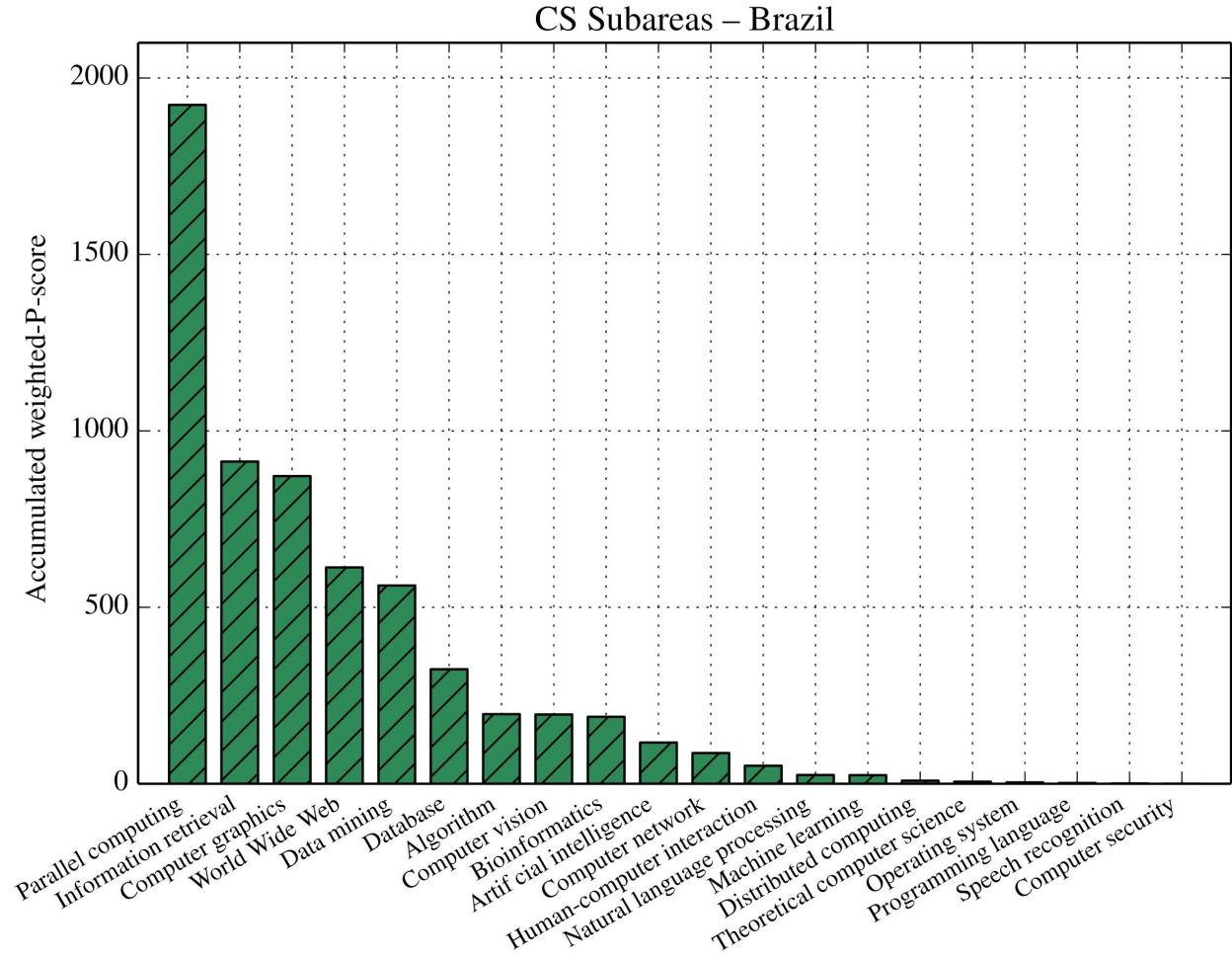
#	University	weighted-P-score
1	University of Massachusetts Amherst	1
2	University of Illinois at Urbana-Champaign	0.4830
3	Carnegie Mellon University	0.4625
4	University of Delaware	0.2452
5	Purdue University	0.2276
6	Northeastern University	0.1633
7	Lehigh University	0.0964
8	Cornell University	0.0552
9	University of Iowa	0.0494
10	University of Illinois at Chicago	0.0477

---

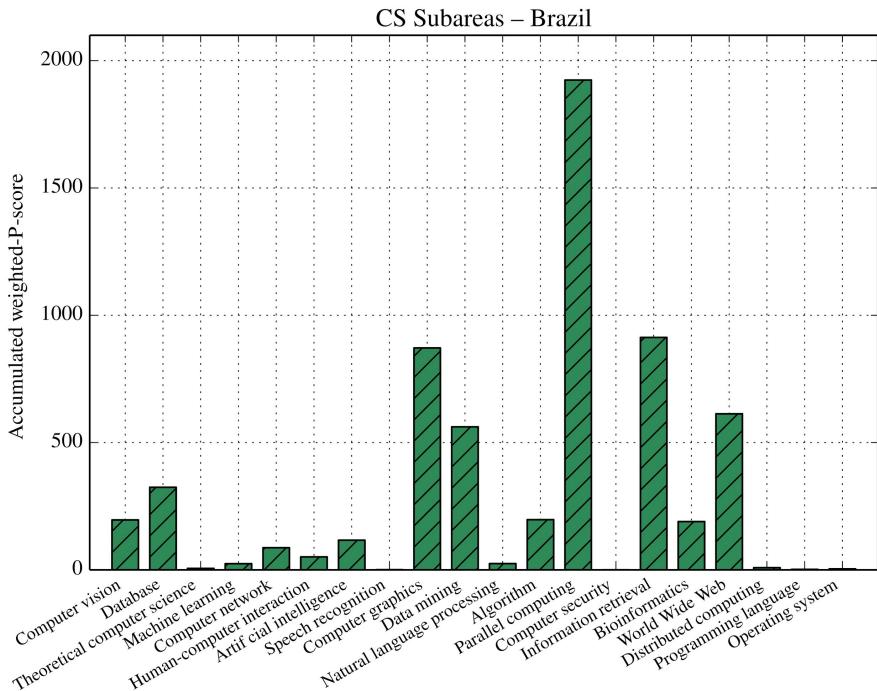
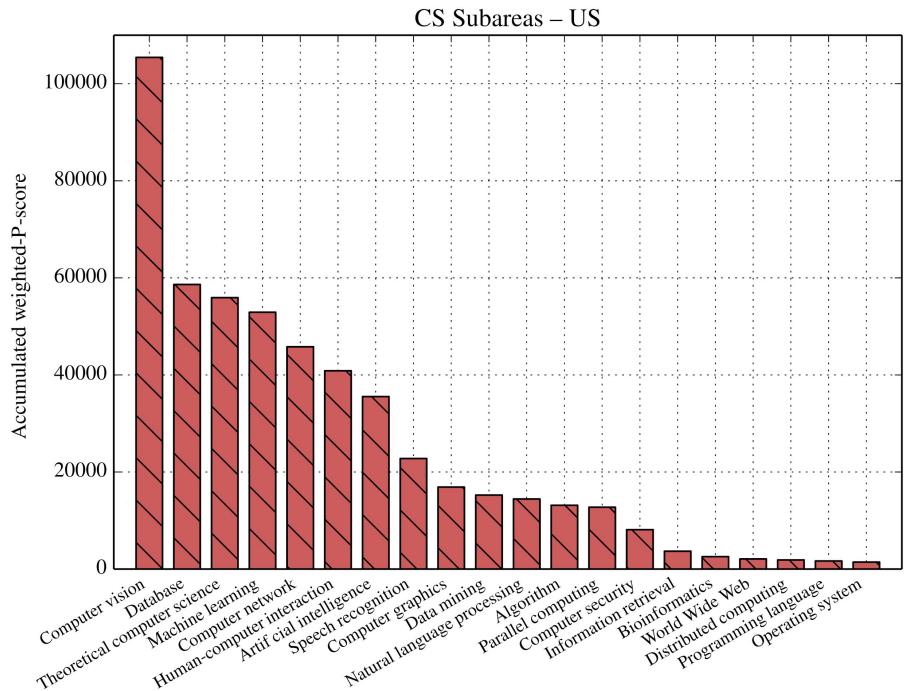
# US Research Directions



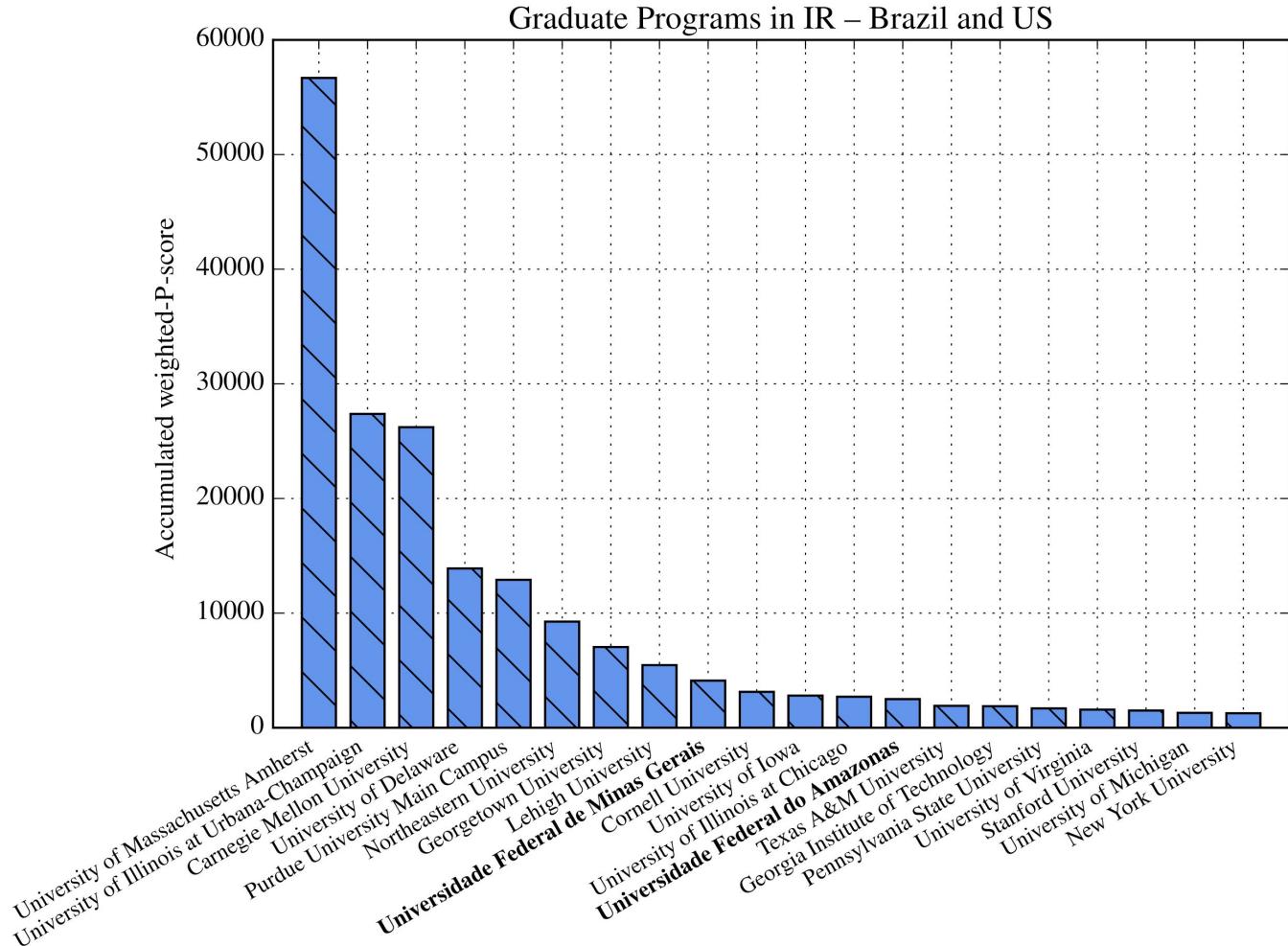
# Brazil Research Directions



# Brazil and US Research Directions

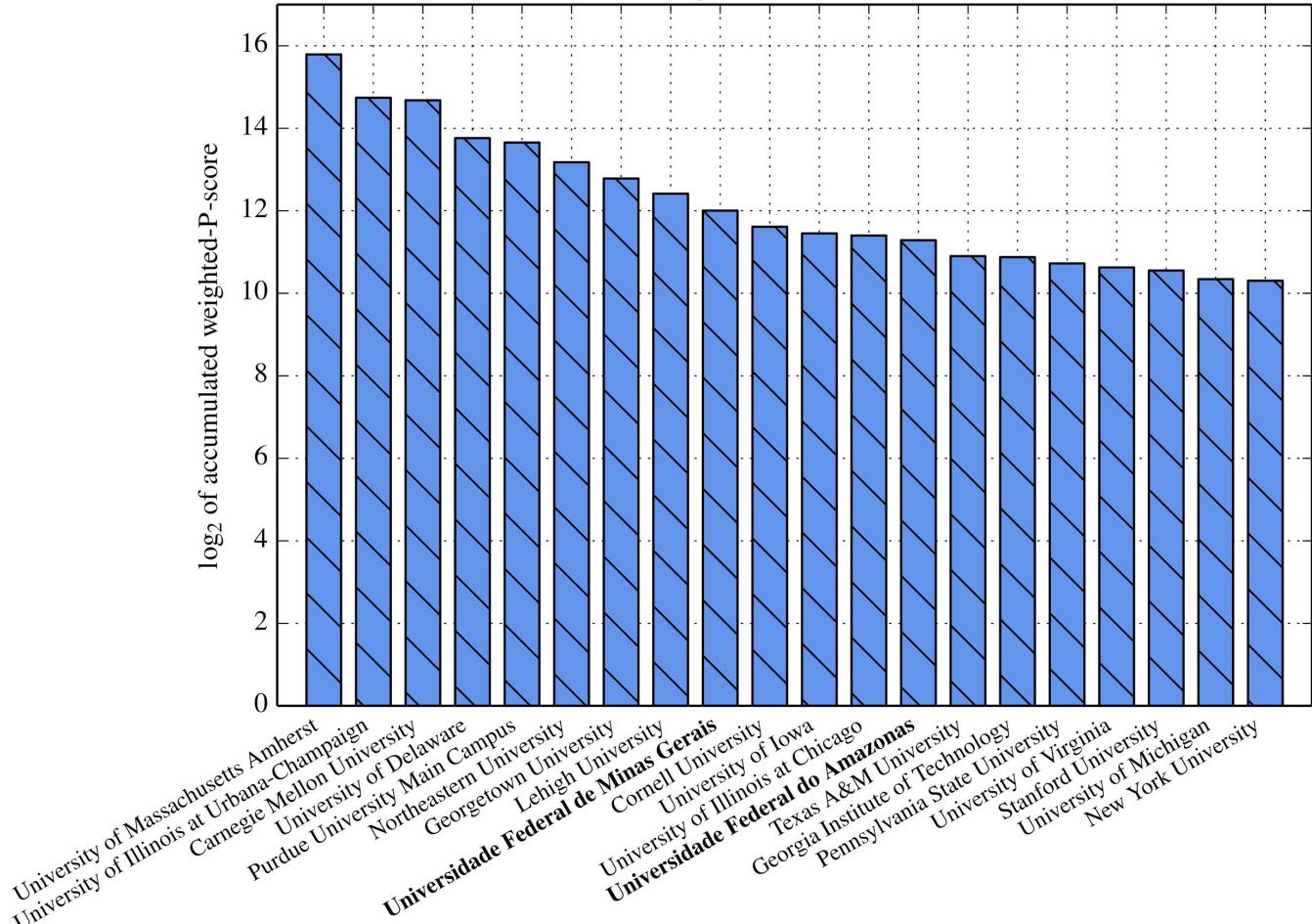


# IR Top Graduate Programs



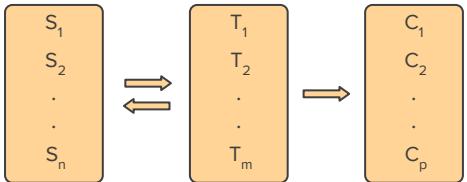
# IR Top Graduate Programs

Graduate Programs in IR – Brazil and US



# Reputation in CS on a per Subarea Basis

How to  
quantify it?



$$\gamma(a, s) = \frac{1}{|\theta(a)|} \sum_{p \in \theta(a)} \begin{cases} \frac{1}{\log_2(y(0)-y(p)+2)} & \text{if } p \in V_s \\ 0 & \text{otherwise} \end{cases}$$

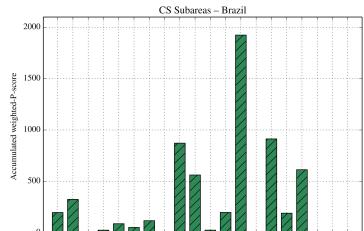
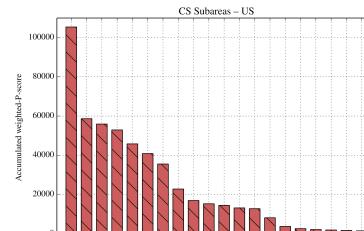
$$\text{weighted-P-score}(g, s) = \sum_{a \in \phi(g)} \gamma(a, s) \times \sum_{p \in \theta(a)} \frac{\text{P-score}(\text{venue}(p))}{\text{number\_of\_authors}(p)}$$

How does it **vary**  
per CS subarea?

#	Standard P-score	norm-P-score	#	Final Ranking
1	SIGIR (e)	SIGIR (e)	1	SIGIR (e)
2	CIKM	CIKM	2	CIKM
3	TREC	TREC	3	TREC
4	ECIR	ECIR	4	ECIR
5	CLEF	CLEF	5	CLEF
6	WWW	WWW	6	SIGIR (j)
7	JASIS	TREC SPIRE	7	JCDL
8	IPM	ECIR ADCS	8	TOIS
9	SIGIR (j)	RIAO IR AIRS	9	IR
10	MM	WSDM	10	WSDM
11	JCDL	NTCIR INEX	11	NTCIR
12	TOIS	WSDM TOIS	12	SPIRE
13	IR	JCDL SIGIR (j)	13	AIRS
14	WSDM	TWEB CLEF	14	IR
15	NTCIR	LA-WEB	15	INEX
16	KDD		16	HIX
17	TKDE		17	ICTIR
18	ACM		18	WWW
19	ICDM		19	LA-WEB
20	SPIRE		20	TWEB

#	University	weighted-P-score
1	University of Massachusetts Amherst	1
2	University of Illinois at Urbana-Champaign	0.4830
3	Carnegie Mellon University	0.4625
4	University of Delaware	0.2452
5	Purdue University	0.2276
6	Northeastern University	0.1633
7	Lehigh University	0.0964
8	Cornell University	0.0552
9	University of Iowa	0.0494
10	University of Illinois at Chicago	0.0477

How do the research  
in CS in BR and US  
differ?



# R



R R R R R R R R R



# Conclusions

Solving the venues **encroachment problem** allows us to improve the ranking of graduate programs in a given subarea.

**P-score** can be used to find core venues of a subarea and also to rank graduate programs on a per subarea basis.

The CS subareas in which **Brazil and US** have major scientific impact are basically disjoint.

# What to do next?

Include academic data of other regions of world (e.g. Europe and Asia), identify research trends, perform a temporal analysis of CS communities, validate the reputation model in other broad areas such as Economics.



# **Reputation in Computer Science on a per Subarea Basis**

---

Alberto Hideki Ueda  
Berthier Ribeiro-Neto and Nivio Ziviani

July 14, 2017  
[ueda@dcc.ufmg.br](mailto:ueda@dcc.ufmg.br)

“We have a responsibility to better inform the public, by providing them with relevant information.”

Moshe Y. Vardi, editor-in-chief of Communications of the ACM 2016, about academic rankings

# Ranking of Venues (DB and DM)

#	Standard P-score	norm-P-score	#	Final Ranking	#	Standard P-score	norm-P-score	#	Final Ranking
1	SIGMOD (c)		1	SIGMOD (c)	1	KDD		1	KDD
2	ICDE		2	ICDE	2	ICDM		2	ICDM
3	VLDB (c)		3	VLDB (c)	3	CIKM		3	CIKM
4	PVLDB		4	PVLDB	4	ICDE		4	ICDE
5	TKDE		5	TKDE	5	ICML		5	ICML
6	DEBU		6	DEBU	6	SDM		6	SDM
7	SIGMOD (j)		7	SIGMOD (j)	7	TKDE		7	TKDE
8	EDBT		8	EDBT	8	WWW		8	PKDD
9	CIKM		9	PODS	9	SIGMOD		9	PAKDD
10	PODS		10	TODS	10	AAAI		10	SIGKDD
11	TODS		11	VLDB (j)	11	PKDD		11	DATAMINE
12	KDD		12	DASFAA	12	NIPS		12	KAIS
13	VLDB (j)		13	SSDBM	13	PAKDD		13	PVLDB
14	WWW		14	ICDT	14	VLDB (c)		14	WSDM
15	ICDM		15	CIDR	15	SIGIR		15	TKDD
16	DASFAA		16	WEBDB	16	SIGKDD		16	VLDB (c)
17	SSDBM		17	SSD	17	DATAMINE		17	RECSYS
18	IS		18	DPD	18	KAIS		18	TIST
19	ICDT		19	COMAD	19	JMLR		19	SSD
20	CIDR		20	TKDD	20	PVLDB		20	SADM

# Ranking of Grad. Programs (DB and DM)

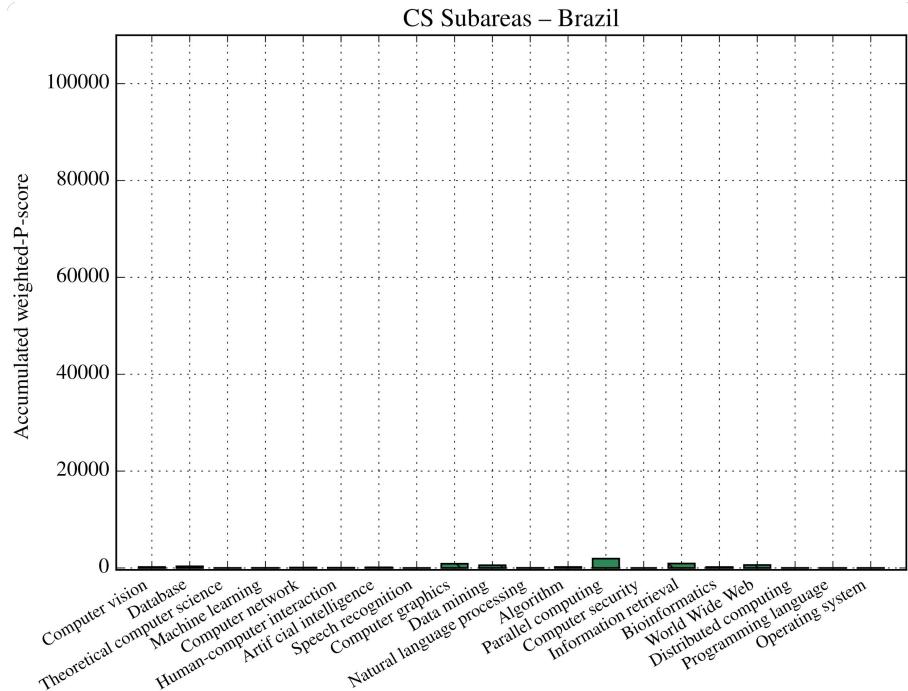
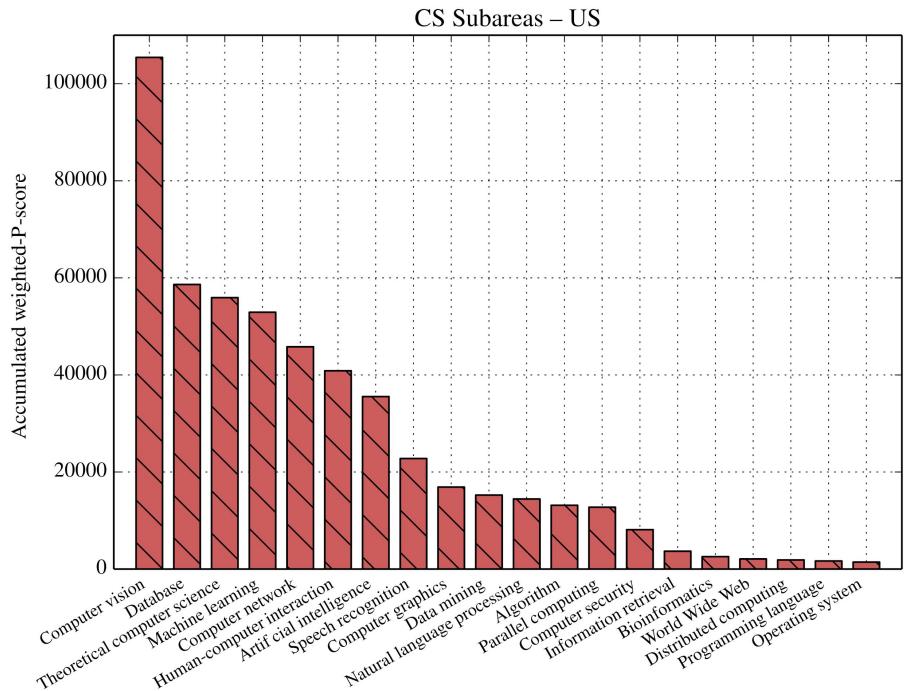
#	University	weighted-P-score
1	University of Wisconsin-Madison	1
2	Stanford University	0.6570
3	University of Illinois at Urbana-Champaign	0.5687
4	Massachusetts Institute of Technology	0.4975
5	Duke University	0.4616
6	University of Massachusetts Amherst	0.4243
7	University of Michigan	0.4195
8	University of California-Irvine	0.4120
9	University of Maryland-College Park	0.4101
10	University of California-Santa Cruz	0.3982

#	University	weighted-P-score
1	University of Illinois at Chicago	1
2	Carnegie Mellon University	0.6857
3	University of Illinois at Urbana-Champaign	0.6344
4	University of Minnesota	0.5350
5	Arizona State University	0.4276
6	University of California-Riverside	0.4212
7	Georgia Institute of Technology	0.3955
8	University of Michigan	0.3275
9	Rensselaer Polytechnic Institute	0.2761
10	University of California-Davis	0.2593

# Standard P-score Suggested Configurations

	Venues	Authors	Groups
Groups	$G^* \rightleftharpoons V \rightarrow V$	$G^* \rightleftharpoons V \rightarrow A$	$G^* \rightleftharpoons V \rightarrow G$
Authors	$A^* \rightleftharpoons V \rightarrow V$	$A^* \rightleftharpoons V \rightarrow A$	$A^* \rightleftharpoons V \rightarrow G$
Venues	$V^* \rightleftharpoons A \therefore A^* \rightleftharpoons V \rightarrow V$	$V^* \rightleftharpoons A \therefore A^* \rightleftharpoons V \rightarrow A$	$V^* \rightleftharpoons A \therefore A^* \rightleftharpoons V \rightarrow G$

# Brazil and US Research Directions



“When a measure becomes a target, it ceases to be a good measure.”

Goodhart's Law,  
Charles Goodhart, Economist, 1975