

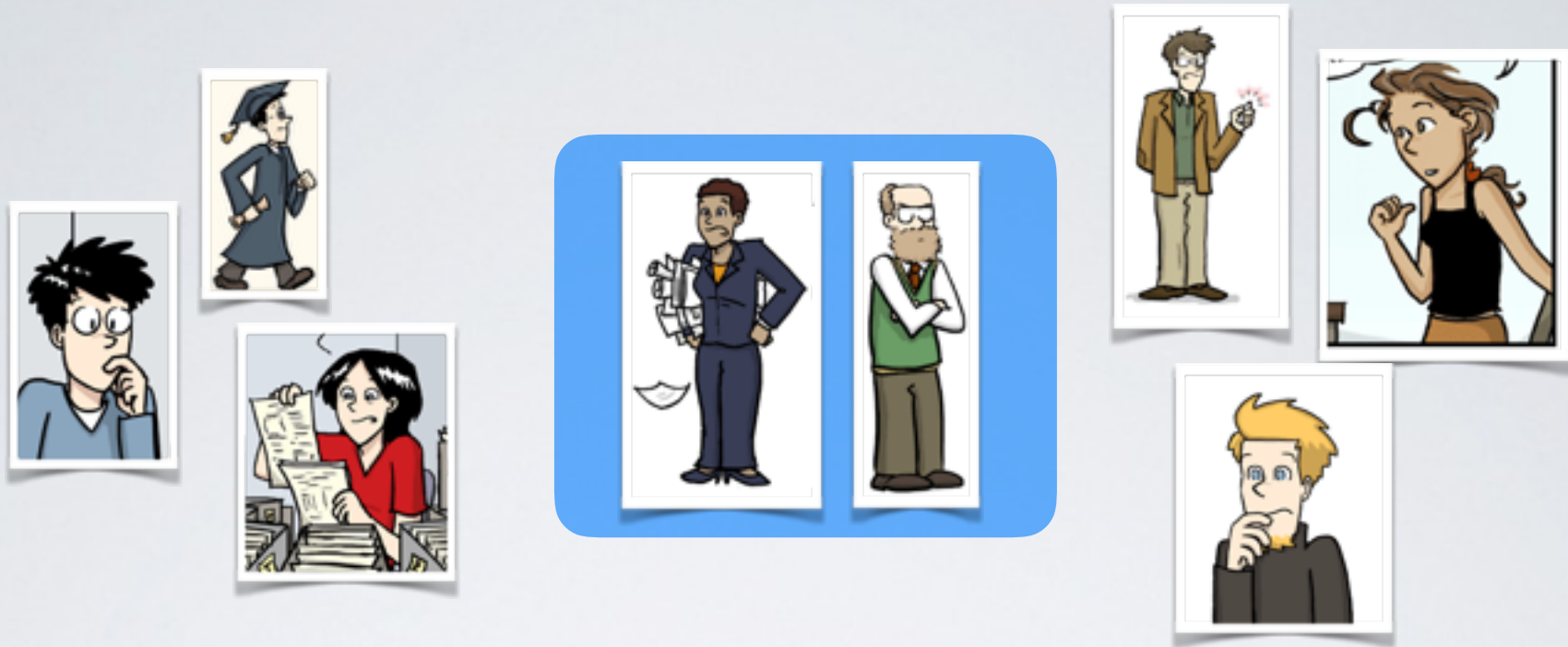


# Network Structure and Dynamics of the Scientific Workforce

Aaron Clauset  
@aaronclauset  
Computer Science Dept. & BioFrontiers Institute  
University of Colorado, Boulder  
External Faculty, Santa Fe Institute

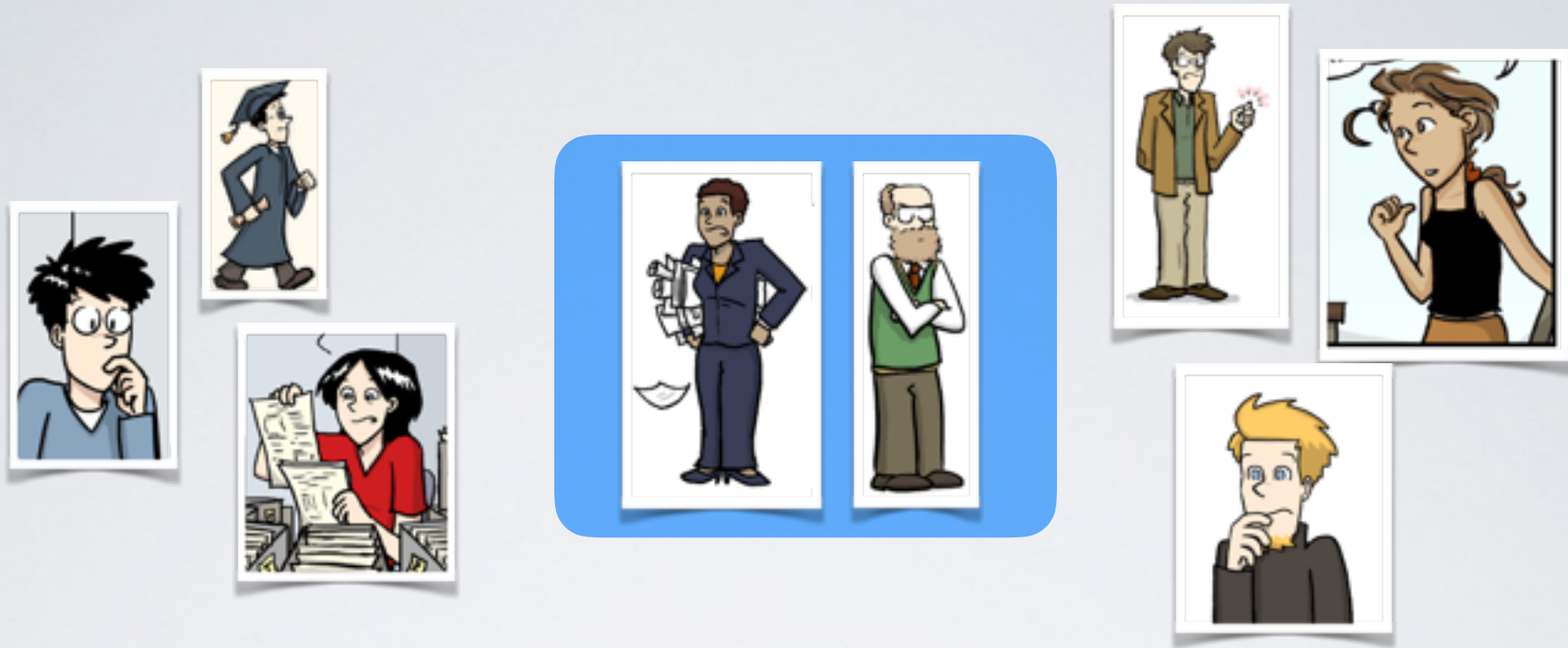


**faculty**



- faculty research agendas shape what discoveries are made
- faculty train students and postdocs
- faculty have long careers
- faculty are numerous

**faculty play a special role in the scientific workforce**



but:

- who hires whose graduates as faculty?
- what does the "system" of faculty production look like?
- what predicts faculty placement?
- where are there inequalities in this system?
- what are their consequences? what drives them?

**who hires whose graduates as faculty?**

**COLLECT**



**ALL THE DATA**





collect all the data

	Computer Science	Business	History
institutions	205	112	144
tenure-track faculty	5032	9336	4556
mean size	25	83	32

$$\sum = 18,924$$


## collect all the data

	Computer Science	Business	History	
institutions	205	112	144	
tenure-track faculty	5032	9336	4556	$\Sigma = 18,924$
mean size	25	83	32	
Full Professors	2400 (48%)	4294 (46%)	2097 (46%)	
Associate Prof.	1772 (35%)	2521 (27%)	1611 (35%)	
Assistant Prof.	860 (17%)	2521 (27%)	848 (19%)	
female	15%	22%	36%	

## collect all the data

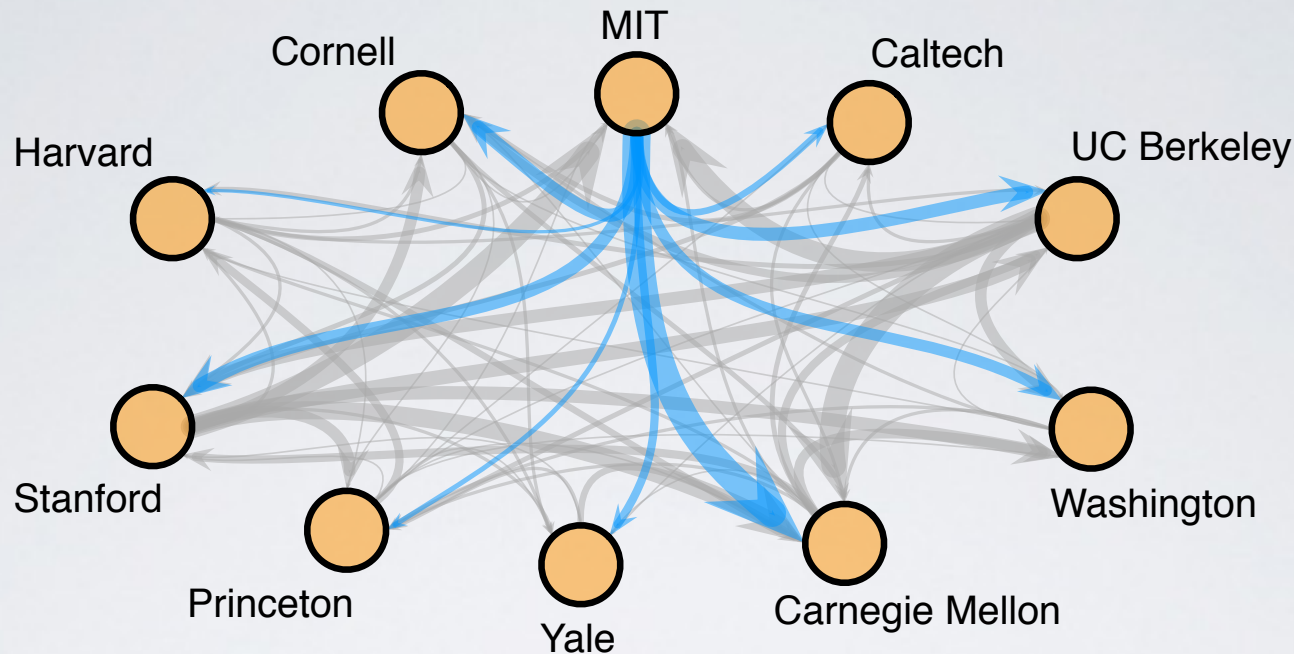
	Computer Science	Business	History	
institutions	205	112	144	
tenure-track faculty	5032	9336	4556	$\Sigma = 18,924$
mean size	25	83	32	
Full Professors	2400 (48%)	4294 (46%)	2097 (46%)	
Associate Prof.	1772 (35%)	2521 (27%)	1611 (35%)	
Assistant Prof.	860 (17%)	2521 (27%)	848 (19%)	
female	15%	22%	36%	
PhDs in-sample	87%	84%	89%	

**nearly closed hiring systems**



## faculty market is a *network*

- vertices are PhD-granting universities
- consumers  $\leftrightarrow$  producers
- $v$  hires from  $u$ , add an edge  $u \rightarrow v$



## faculty market is a *network*

- vertices are PhD-granting universities
- consumers  $\leftrightarrow$  producers
- $v$  hires from  $u$ , add an edge  $u \rightarrow v$



**huge inequalities in faculty production**

## huge inequalities in faculty production

### Gini coefficients (production)

- 0.69, 0.62, 0.72

### 50% of faculty from

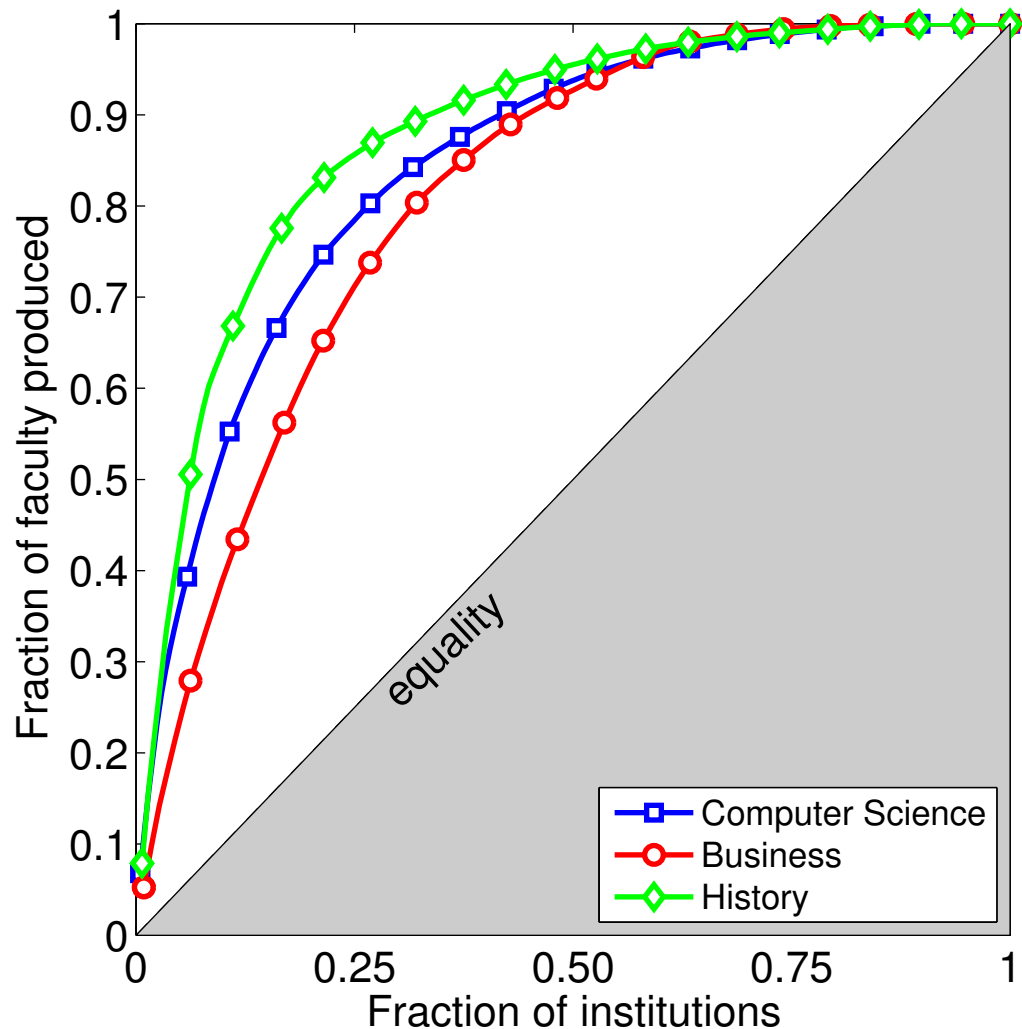
- 18, 16, 8 universities

### net producers $k_{\text{out}}/k_{\text{in}} > 1$

- 24%, 36%, 18%

### 1-10 producers vs.

- 11-20 : 1.6, 2.1, 3.0x more
- 21-30 : 3.1, 2.3, 5.6x more



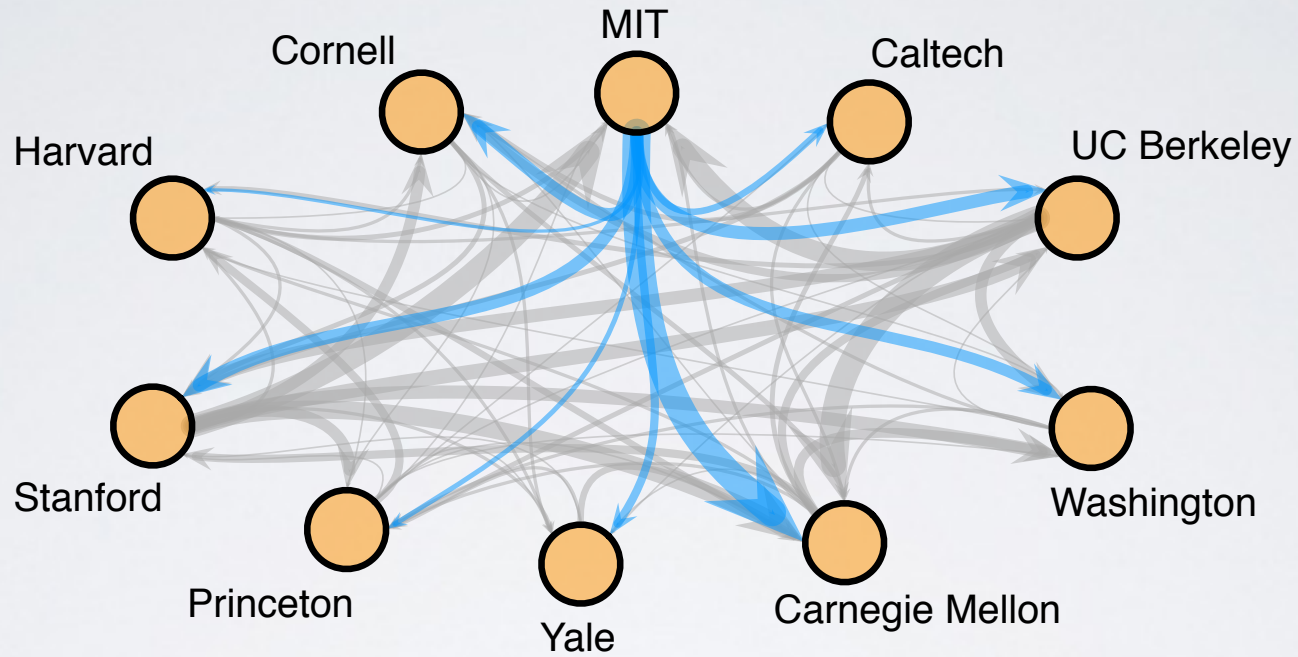
[1] order: CS, Business, History

[2] U.S. Income Gini coefficient = 0.45

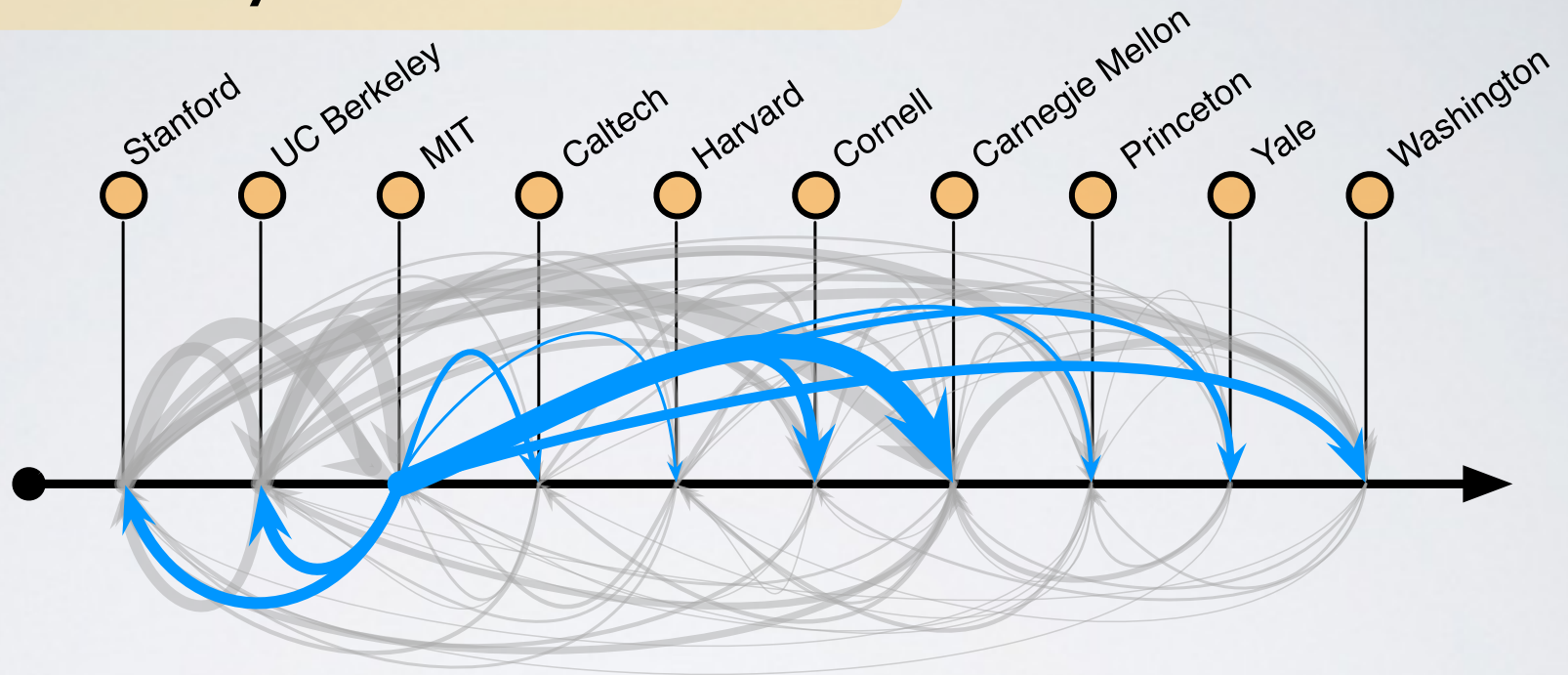
## a prestige hierarchy

- difficult to talk about inequalities in academia without talking about rankings
- let's extract a data-driven ranking from the network

## a prestige hierarchy



## a prestige hierarchy



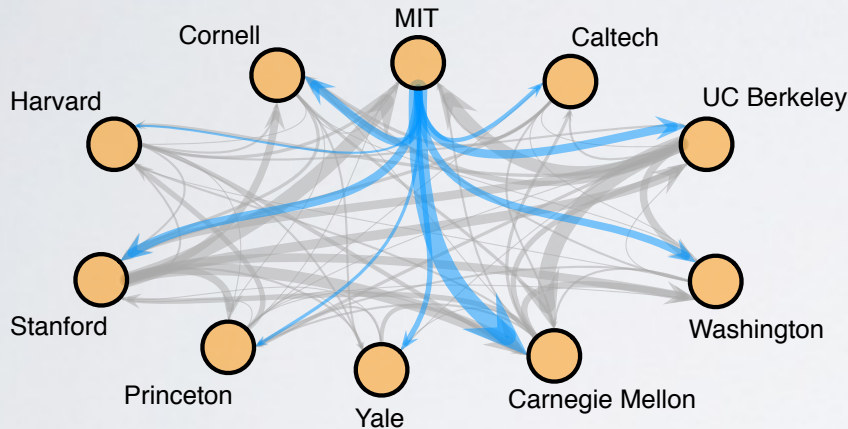
- select permutation (a ranking)  $\pi$  that minimizes the number of "rank violations" : edges  $(u, v)$  where  $\pi_v < \pi_u$
- higher-ranked nodes have greater "placement power"
- equivalent to *minimum feedback arc set* problem (NP-hard)

[1] these "MVR"s have a deep history in social theory for extracting dominance or prestige hierarchies from data, especially in animal behavior

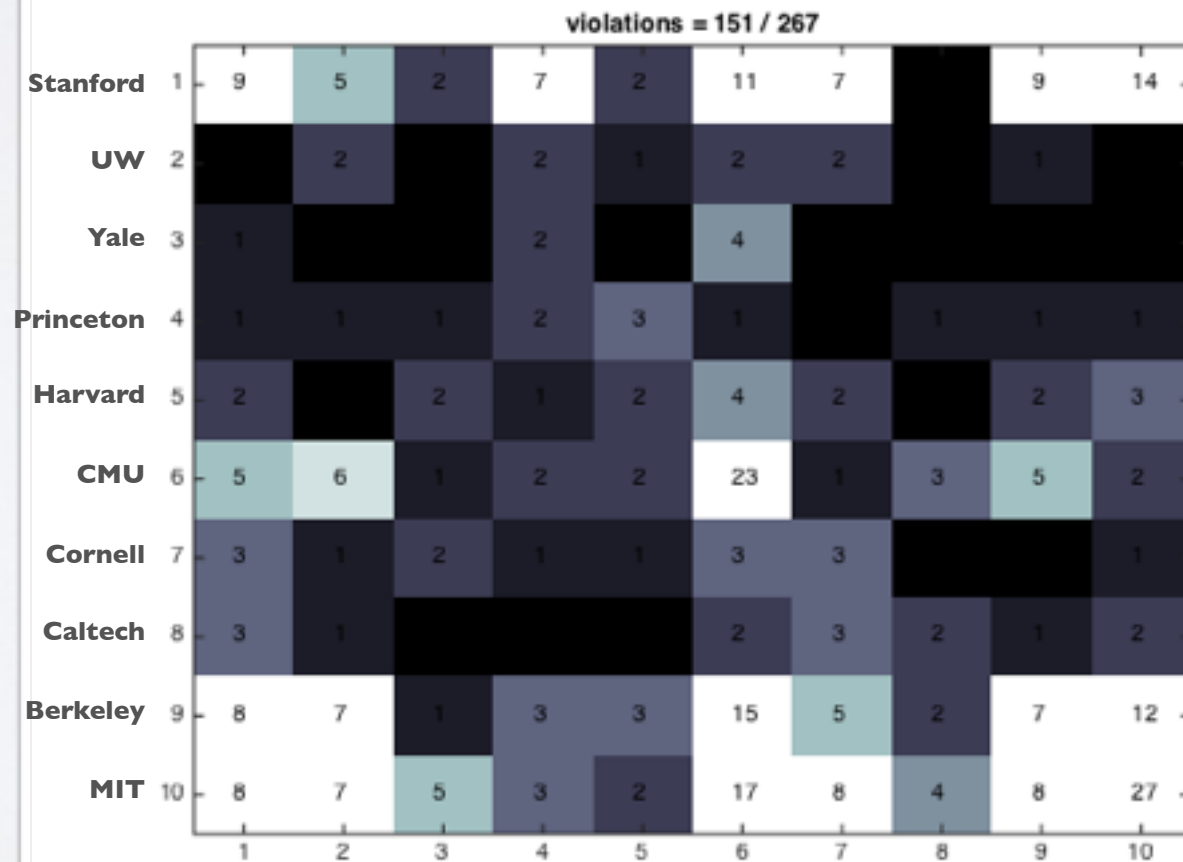
[2] MFAS: find the set of arcs of minimum cardinality whose removal converts a directed graph  $G$  into a directed acyclic graph

[3] there are many equivalent MVRs for our network. we sample these using a zero-temp MCMC, and average across them to obtain  $\langle \pi \rangle$

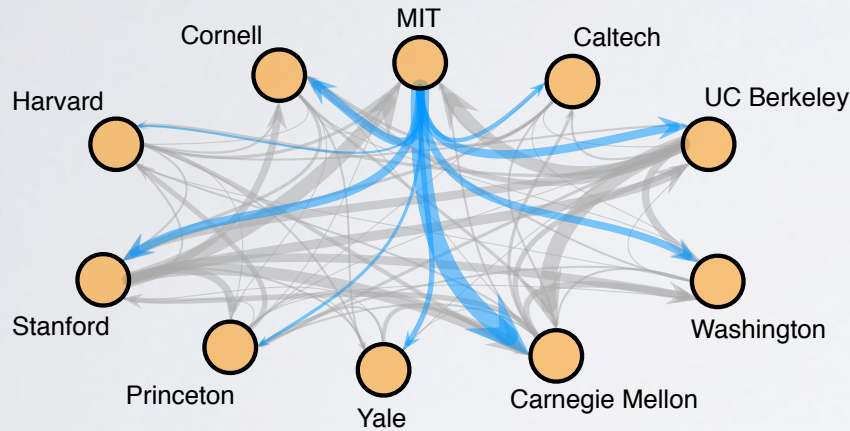
- given an ordering  $\pi$  with  $\psi(\pi, A)$  rank violations on network  $A$
- repeat *ad infinitum*: choose a pair  $(u, v)$ , swap their ranks  $\pi_u \leftrightarrow \pi_v$  to obtain  $\pi'$ , compute  $\psi(\pi', A)$ , accept change if  $\psi(\pi', A) \geq \psi(\pi, A)$
- for instance:



random



- given an ordering  $\pi$  with  $\psi(\pi, A)$  rank violations on network  $A$
- repeat *ad infinitum*: choose a pair  $(u, v)$ , swap their ranks  $\pi_u \leftrightarrow \pi_v$  to obtain  $\pi'$ , compute  $\psi(\pi', A)$ , accept change if  $\psi(\pi', A) \geq \psi(\pi, A)$
- for instance:

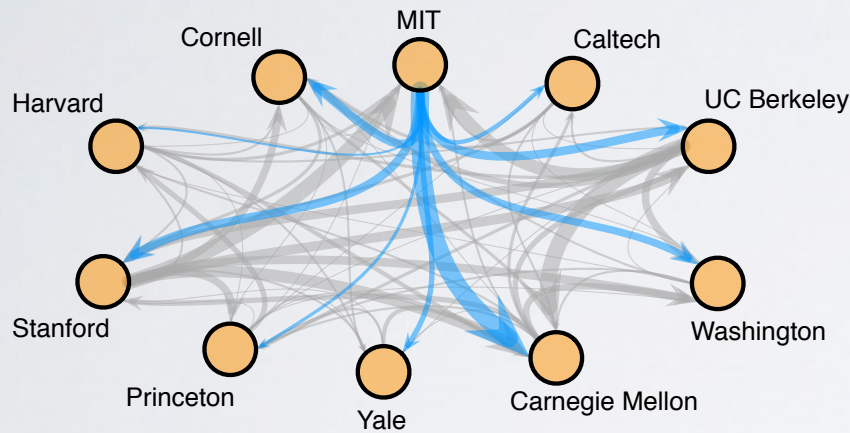


sort by degree

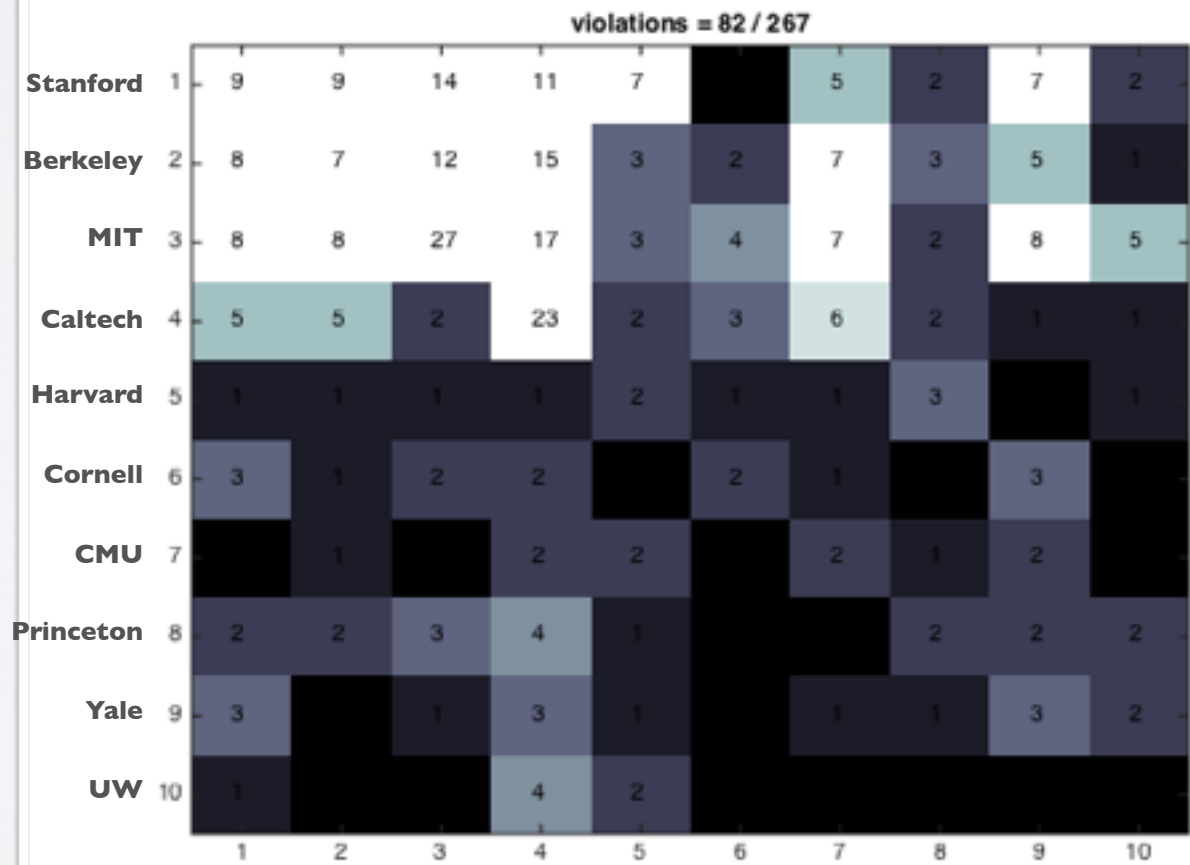
violations = 99 / 267

<b>MIT</b>	1	27	8	8	17	2	8	4	3	7	5
<b>Stanford</b>	2	14	9	9	11	2	7		7	5	2
<b>Berkeley</b>	3	12	8	7	15	3	5	2	3	7	1
<b>CMU</b>	4	2	5	5	23	2	1	3	2	6	1
<b>Harvard</b>	5	3	2	2	4	2	2		1		2
<b>Cornell</b>	6	1	3		3	1	3		1	1	2
<b>Caltech</b>	7	2	3	1	2		3	2		1	
<b>Princeton</b>	8	1	1	1	1	3		1	2	1	1
<b>UW</b>	9			1	2	1	2		2	2	
<b>Yale</b>	10		1		4			2			
		1	2	3	4	5	6	7	8	9	10

- given an ordering  $\pi$  with  $\psi(\pi, A)$  rank violations on network  $A$
- repeat *ad infinitum*: choose a pair  $(u, v)$ , swap their ranks  $\pi_u \leftrightarrow \pi_v$  to obtain  $\pi'$ , compute  $\psi(\pi', A)$ , accept change if  $\psi(\pi', A) \geq \psi(\pi, A)$
- for instance:

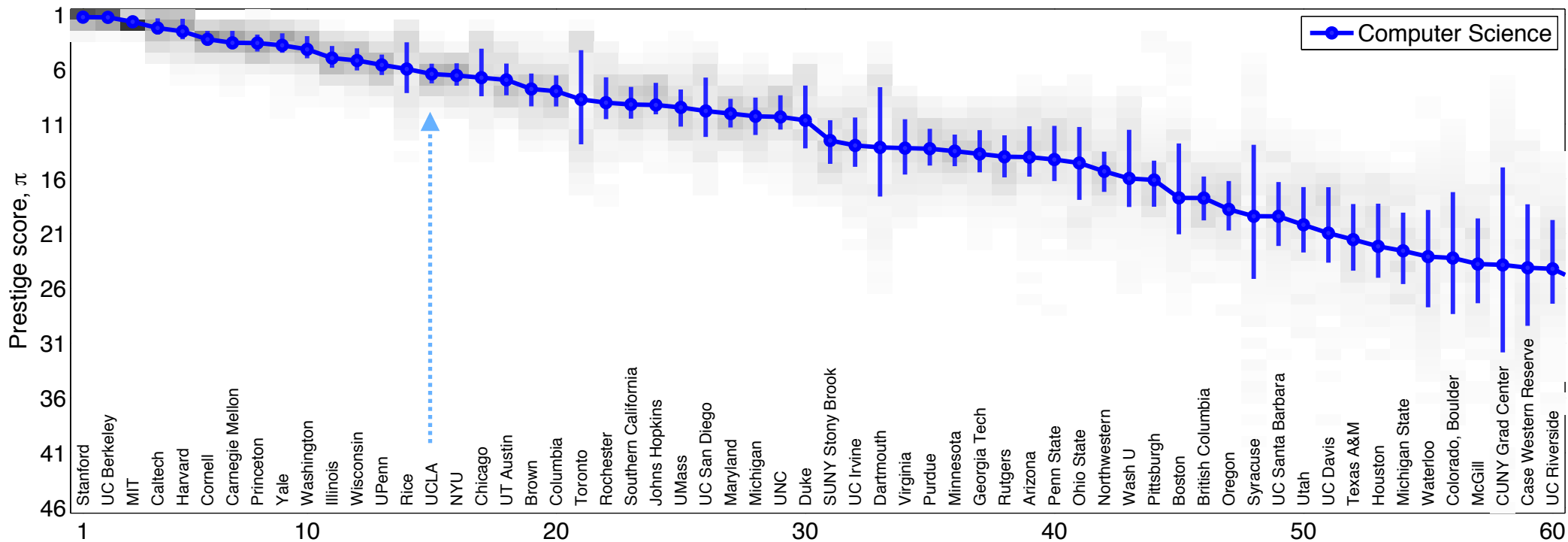


## MVR



## a prestige hierarchy

- what do these prestige hierarchies look like?
- what do they tell us about the structure of faculty hiring?
- what predicts placement?



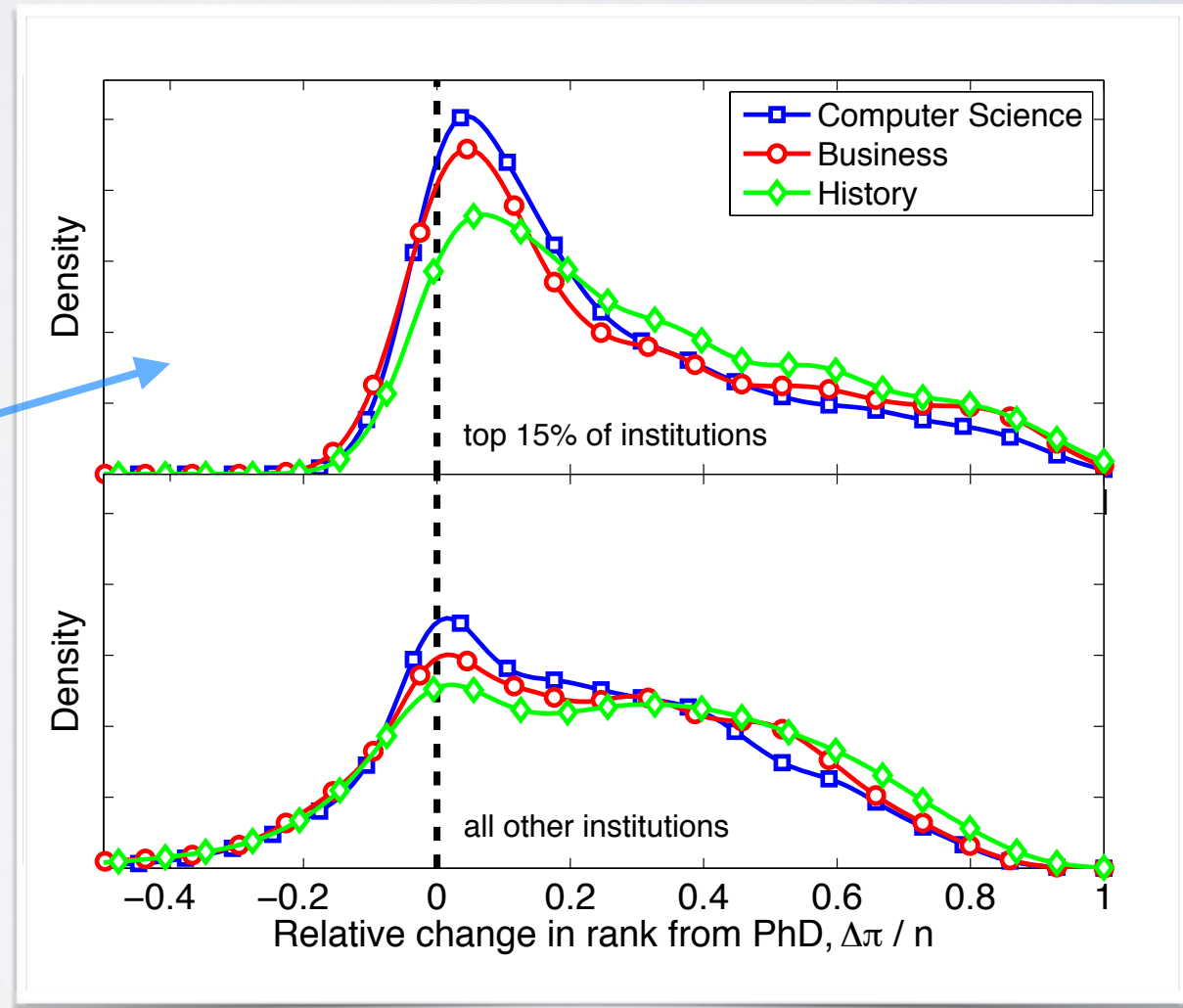
## prestige rankings correlate with USNews and NRC

- here, *prestige*  $\pi$  quantifies *placement power*
- uncertainty increases as prestige decreases
- similar results, but different orderings for Business and History

**most placements are down the hierarchy**

## most placements are down the hierarchy

- down : 88%, 86%, 91%
- up : 12%, 14%, 9%
- $\langle \Delta\pi \rangle = 47, 27, 42$  steps down
- CS: top 15% of departments produce 68% of their own faculty and hire 7% from outside top 25% of departments



## what predicts placement?

- compare 10 single features:

prestige

US News rank

NRC rank

out-degree

in-degree

out/in degree

eigenvector centrality

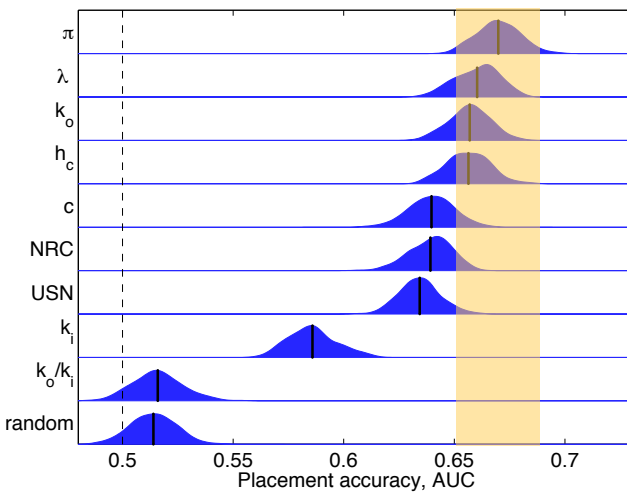
harmonic centrality

closeness centrality

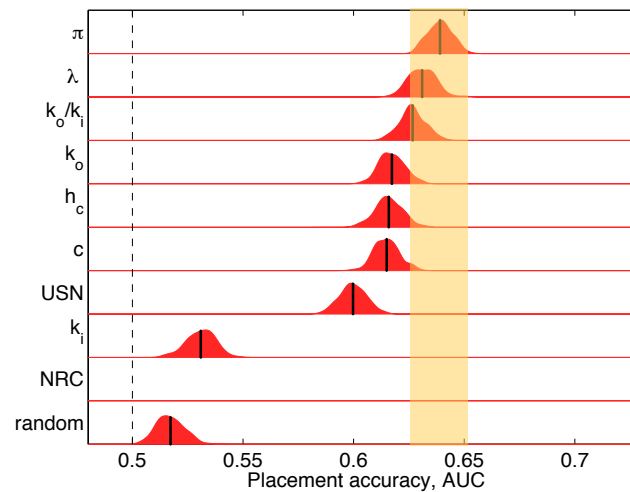
random

## what predicts placement?

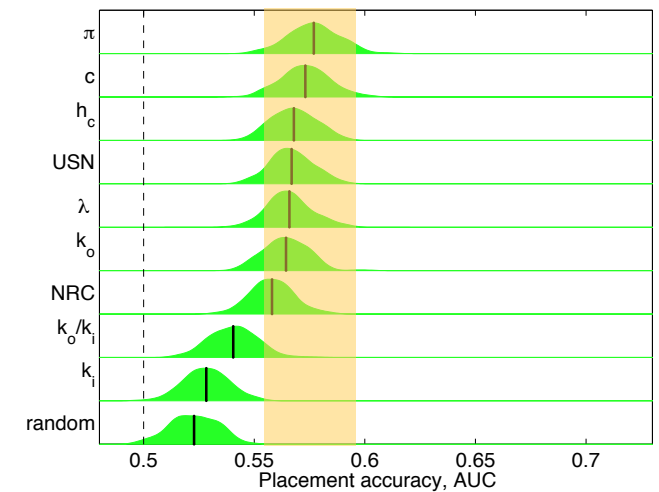
- prestige best *single* predictor in all 3 fields
- order of other features varies by field
- AUCs all below 0.67 = plenty of room for improvement



CS



Business



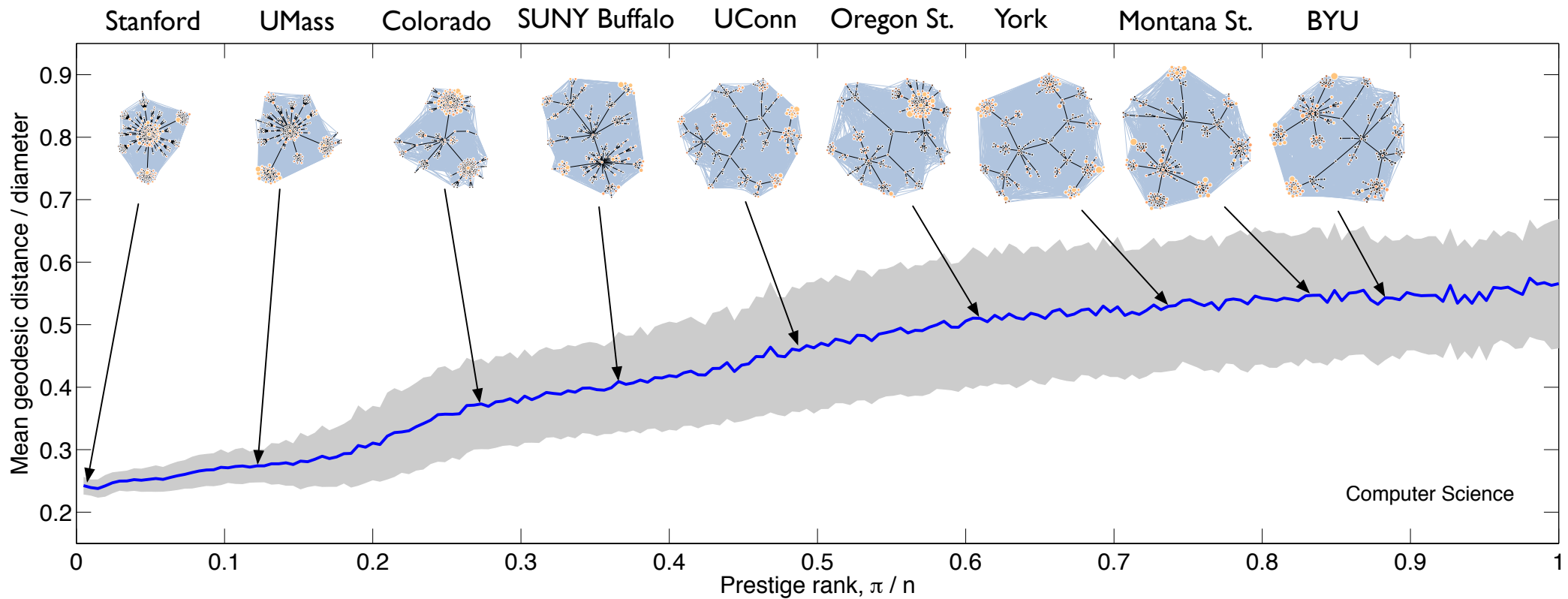
History

## prestige correlates with network position

- core and periphery

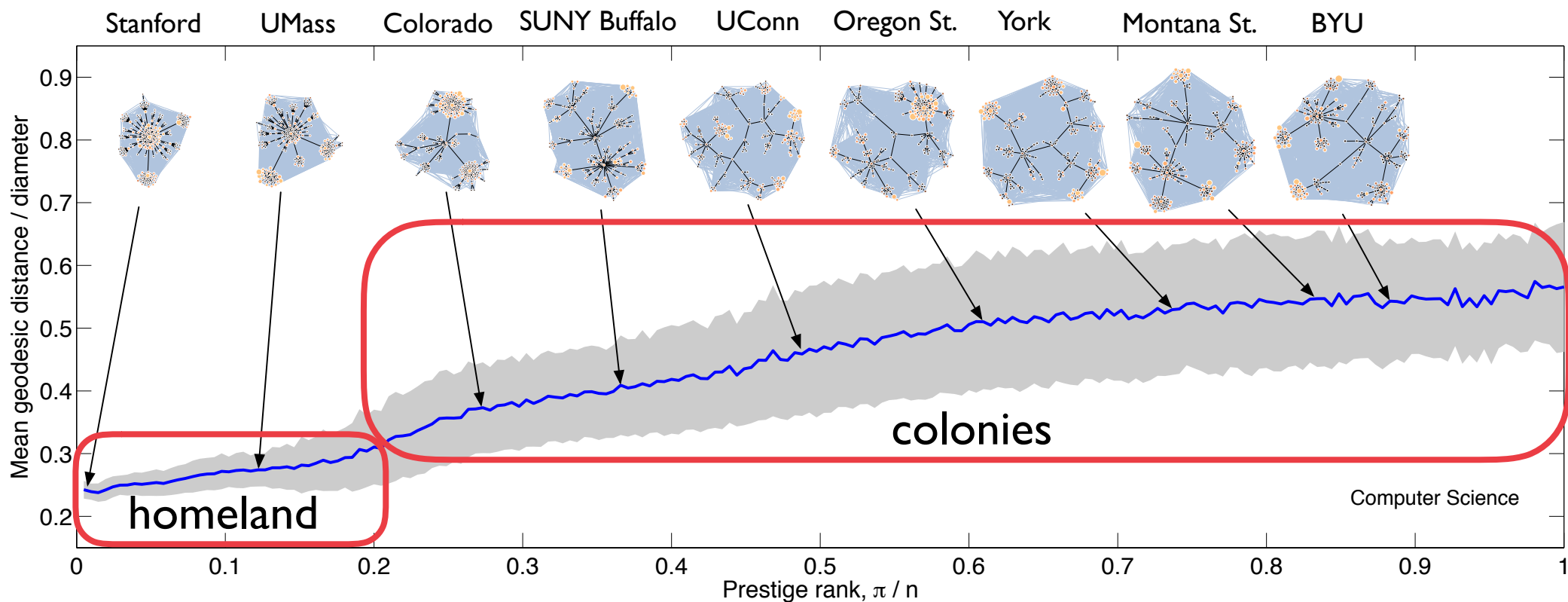
## prestige correlates with network position

- core and periphery



## prestige correlates with network position

- ~~core and periphery~~ homeland and colonies
- prestige is *influence*, via doctoral placement, over research agendas, research communities, and departmental norms across the discipline



## inequality and prestige hierarchies

- prestige is influence, via doctoral placement
- faculty flow out of core, into periphery ("the colonies")
- small fraction stay inside core
- only ~10% of hires flow "upstream"

## inequality and prestige hierarchies

- prestige is influence, via doctoral placement
- faculty flow out of core, into periphery ("the colonies")
- small fraction stay inside core
- only ~10% of hires flow "upstream"

## future work

- how to measure cultural influence of core departments?
- what is different about "upstream" hires?
- what role for other inequalities : gender, ethnicity/race, SES, neighborhood effects, productivity, etc.?

## inequality and prestige hierarchies

- prestige is influence, via doctoral placement
- faculty flow out of core, into periphery ("the colonies")
- small fraction stay inside core
- only ~10% of hires flow "upstream"

## future work

- how to measure cultural influence of core departments?
- what is different about "upstream" hires?
- what role for other inequalities : gender, ethnicity/race, SES, neighborhood effects, productivity, etc.?

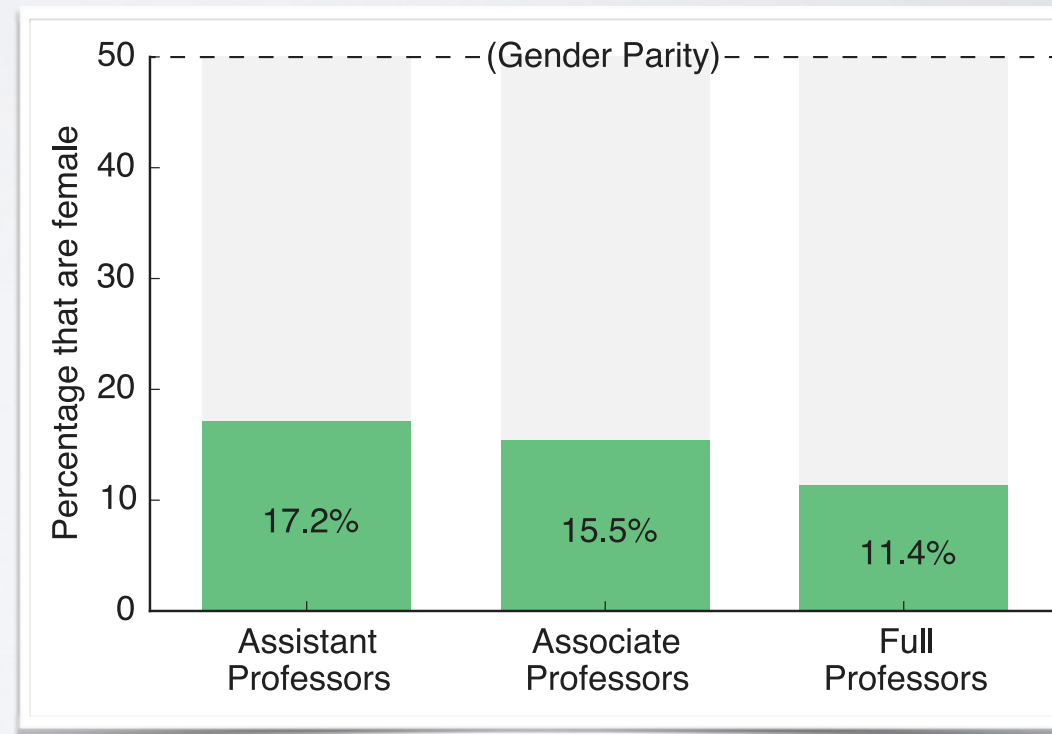
## women are *dramatically* under-represented in computing

18% of bachelors degrees

20% of doctoral degrees

20% of industry positions

**15% of CS faculty positions**



[1] [National Center for Education Statistics](#) (2011 data)

[2] Reliably industry-wide estimates are hard to come by, but see <http://cnet.co/IGZh268>

[3] [Clauset, Arbesman, Larremore, Science Advances](#) 1(1), e1400005 (2015)

**what role does gender play in CS faculty hiring?**

## what role does gender play in CS faculty hiring?

- learn a model of CS faculty hiring for all placements 1970-2010  
via logistic regression (adapted for non-independence of hires)
- consider 6 factors:
  - gender
  - postdoc training
  - changing in geographic region
  - doctoral prestige
  - prestige difference
  - scholarly productivity

## what role does gender play in CS faculty hiring?

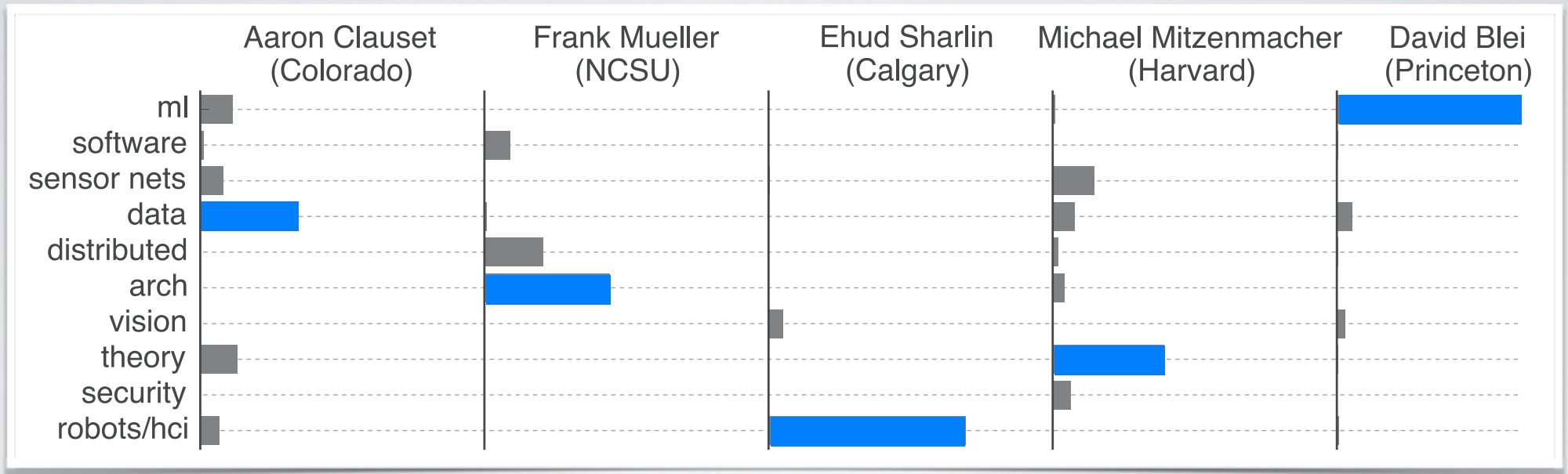
- learn a model of CS faculty hiring for all placements 1970-2010 via logistic regression (adapted for non-independence of hires)
- consider 6 factors:
  - gender
  - postdoc training
  - changing in geographic region
  - doctoral prestige
  - prestige difference
  - scholarly productivity



## faculty productivity data

- publication records (titles, dates, etc.) for 95.1% of sampled faculty
- mean number of pubs prior to first faculty appointment is 11.3  
(but distribution has a heavy tail)
- and, mean varies by subfield

- use a topic model to learn subfields from all pre-hire paper titles:



- for each subfield, we tabulated a distribution over paper counts, weighted by each faculty's inferred emphasis in that field
- for each faculty, we computed a  $z$ -score for their productivity  $\rho_i$  relative to their subfield mixture

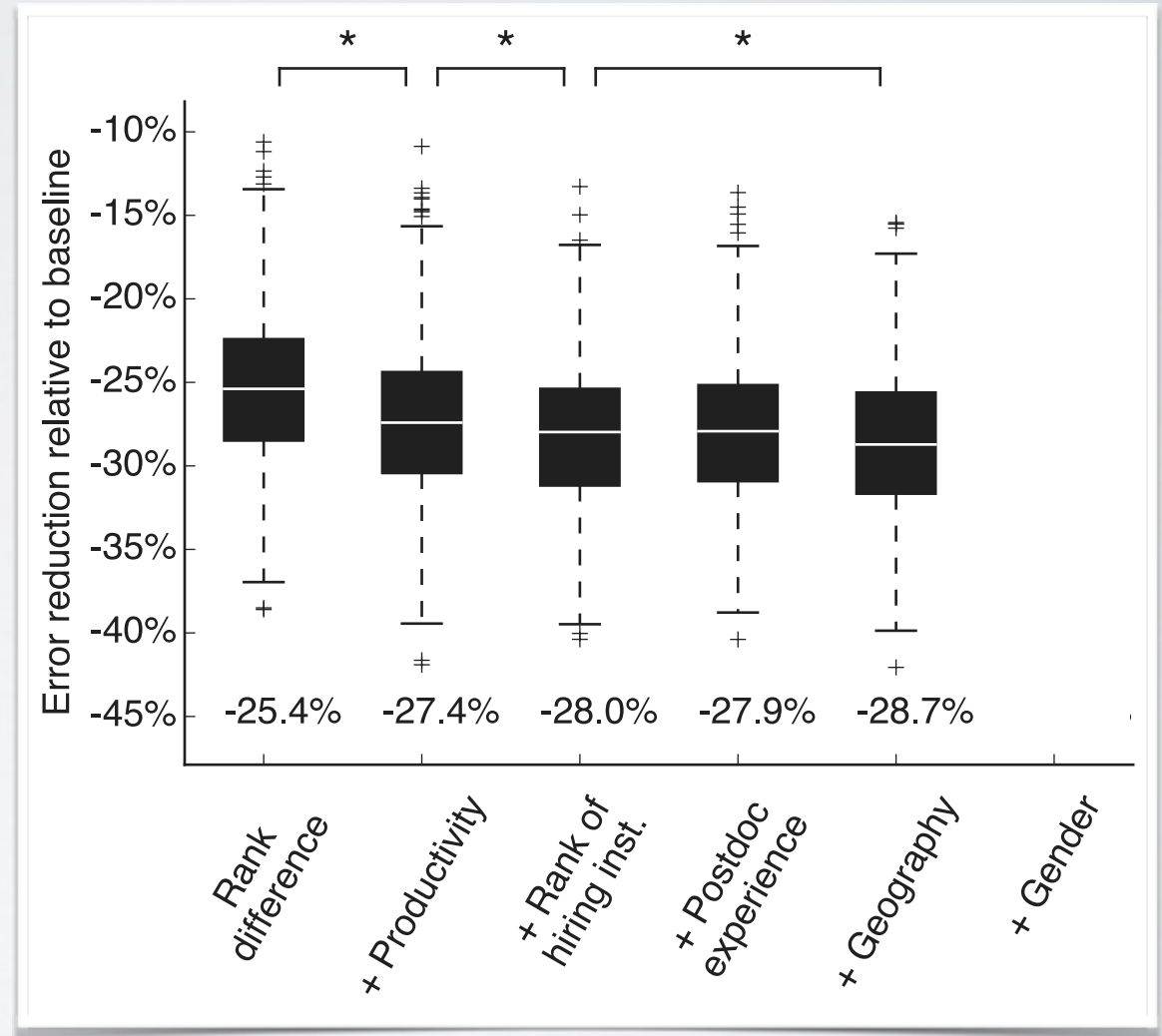
# fitting the model

# fitting the model

- we learn weights  $\vec{w}$  for our covariates by minimizing placement error (with L1 regularization)

$$\text{err} = \frac{1}{m} \sum_{i=1}^m [O(u_i) - M(u_i)]^2 + \lambda \sum_k |\vec{w}_k|$$

- add covariates one-by-one, in greedy fashion, with gender added last

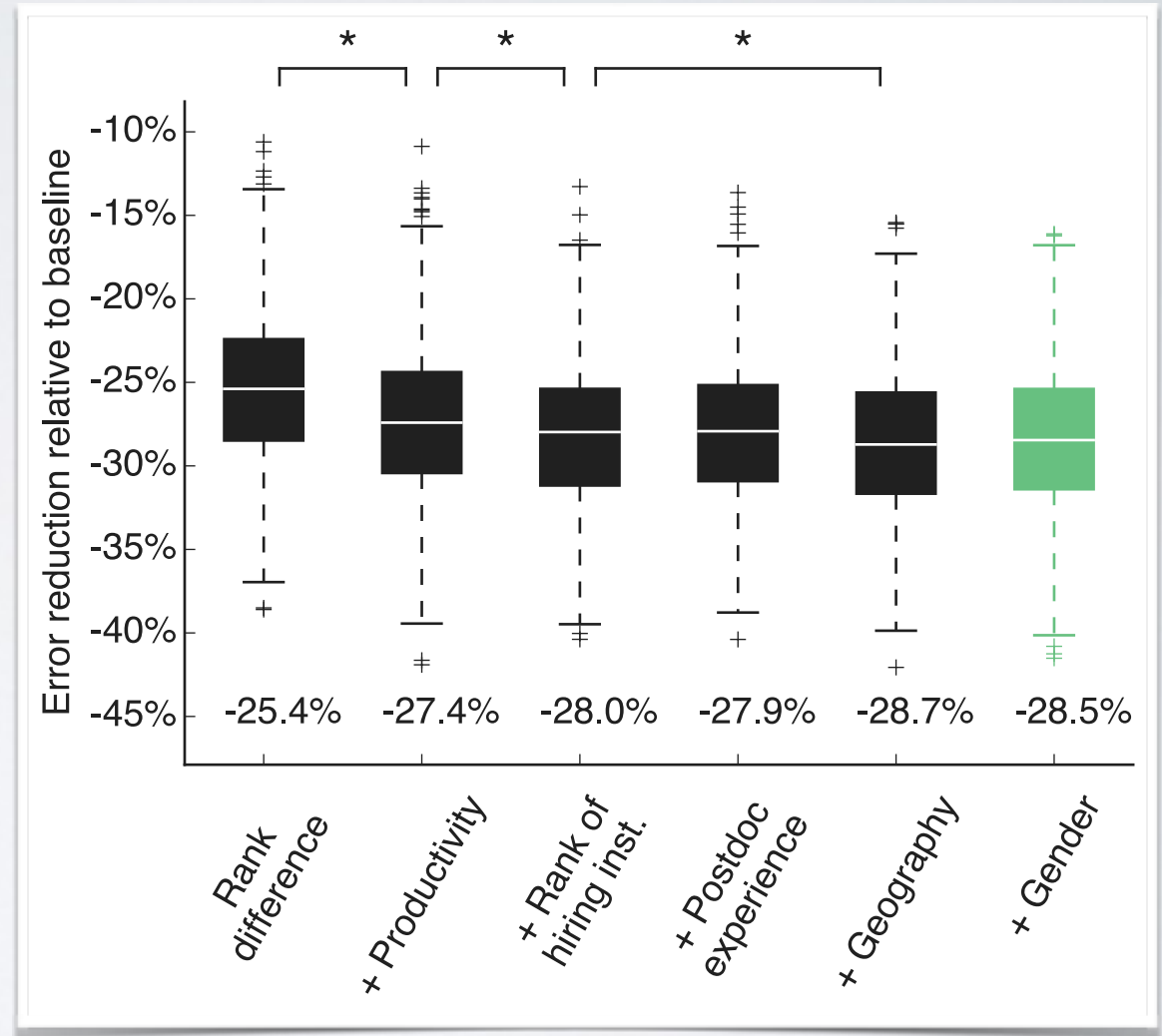


# fitting the model

- we learn weights  $\vec{w}$  for our covariates by minimizing placement error (with L1 regularization)

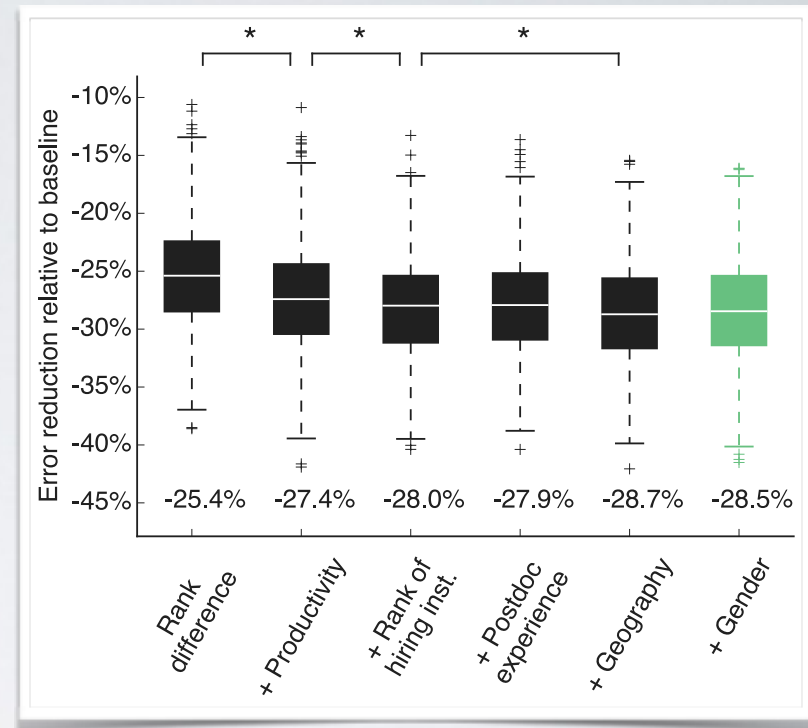
$$\text{err} = \frac{1}{m} \sum_{i=1}^m [O(u_i) - M(u_i)]^2 + \lambda \sum_k |\vec{w}_k|$$

- add covariates one-by-one, in greedy fashion, with gender added last



# system-level results

- adding gender *does not improve* placement accuracy
- three possibilities:
  1. gender is irrelevant to hiring decisions
  2. we modeled gender's effect incorrectly
  3. gender's effect is included in the other variables [see footnotes]

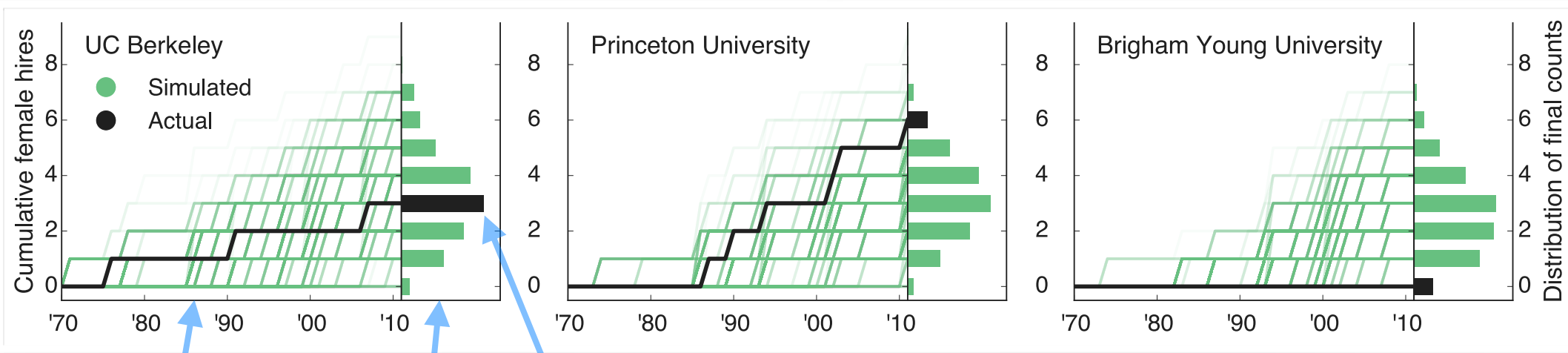


[1] since 2002, men postdoc at equal rates to women (about 28%). prior to 2002, women postdoc'd at greater rates than men (28% vs. 16%)

[2] post-2002, women *with* postdocs are as productive as men *without* postdocs; men *with* postdocs are significantly more productive than women *with* postdocs

# institutions and individuals

- use our learned model to simulate hiring patterns for each institution over 1970-2010
- compare *actual* vs. *expected* number of female hires



1000 simulated histories of hiring

observed count in 2011

simulated distribution of hires

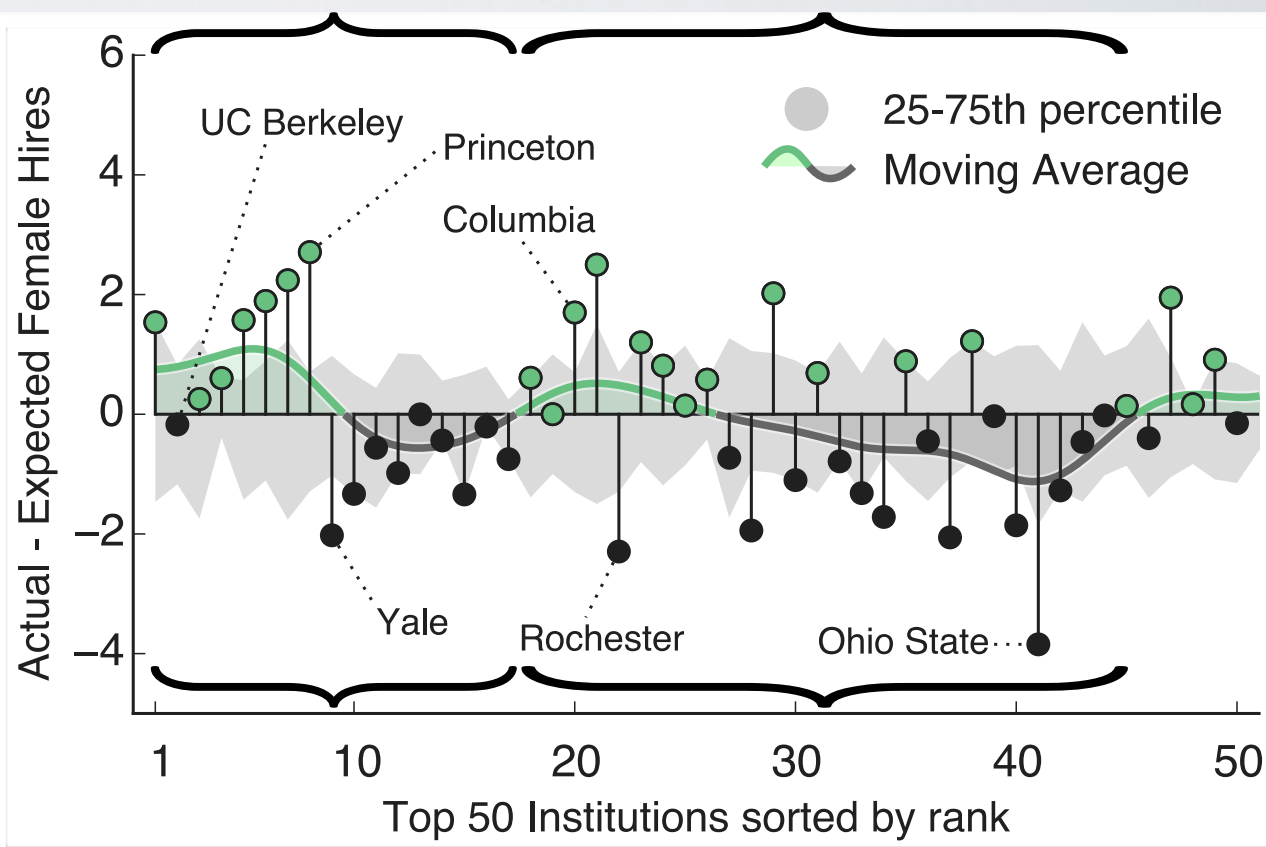
# institutions and individuals

- use our learned model to simulate hiring patterns for each institution over 1970-2010
- compare *actual vs. expected* number of female hires

- for top 50 institutions, *an oscillation:*

an interference effect?  
from non-independence  
of hiring

and, two distinct pools  
of candidates

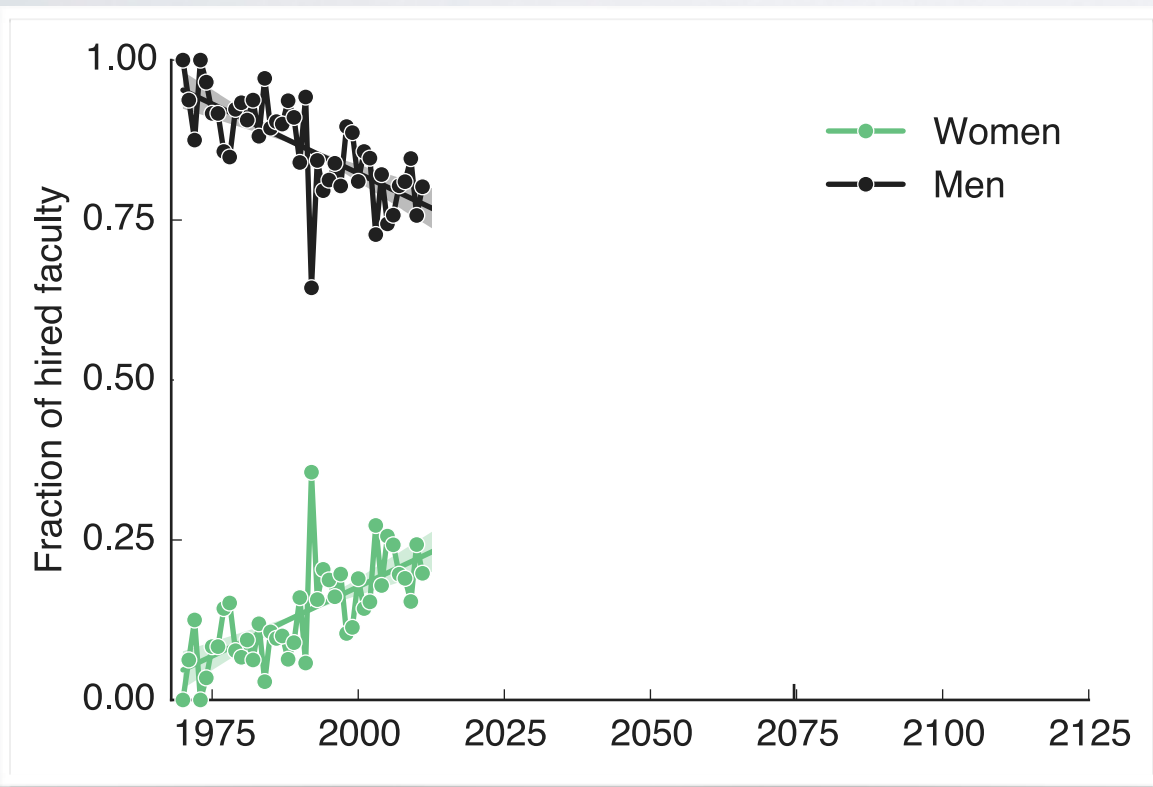


[1] the errors our model makes are interesting, and non-uniform. for instance, men and women tend to exceed the model's expectations at similar rates; but, for under-performing individuals, men tend to fall short of expectations by a wider margin. furthermore, people with postdoc training tend to exceed the model's expectations, but women with postdoc experience tend to exceed expectations by a wider margin than do men

**gender parity**

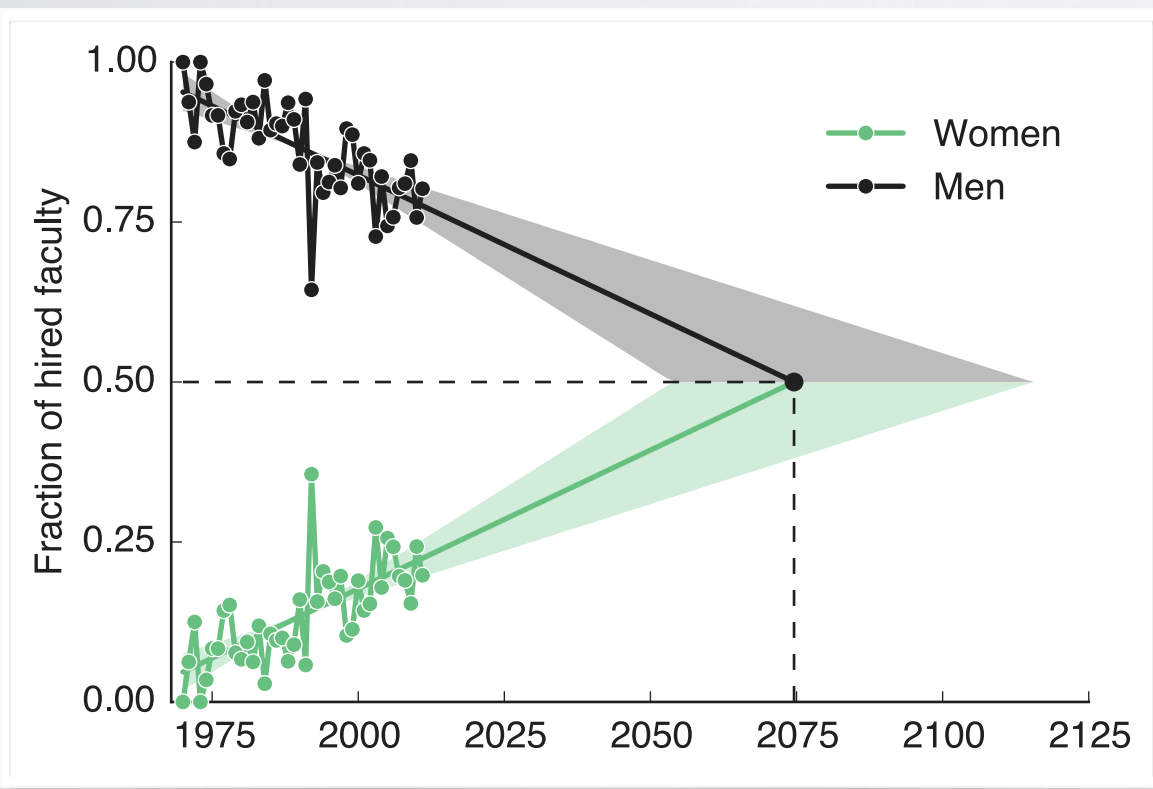
# gender parity

- where is the CS faculty gender ratio going?
- from our 40 years of observable data, the trend is toward parity
- ratio increases by 0.43% per year



# gender parity

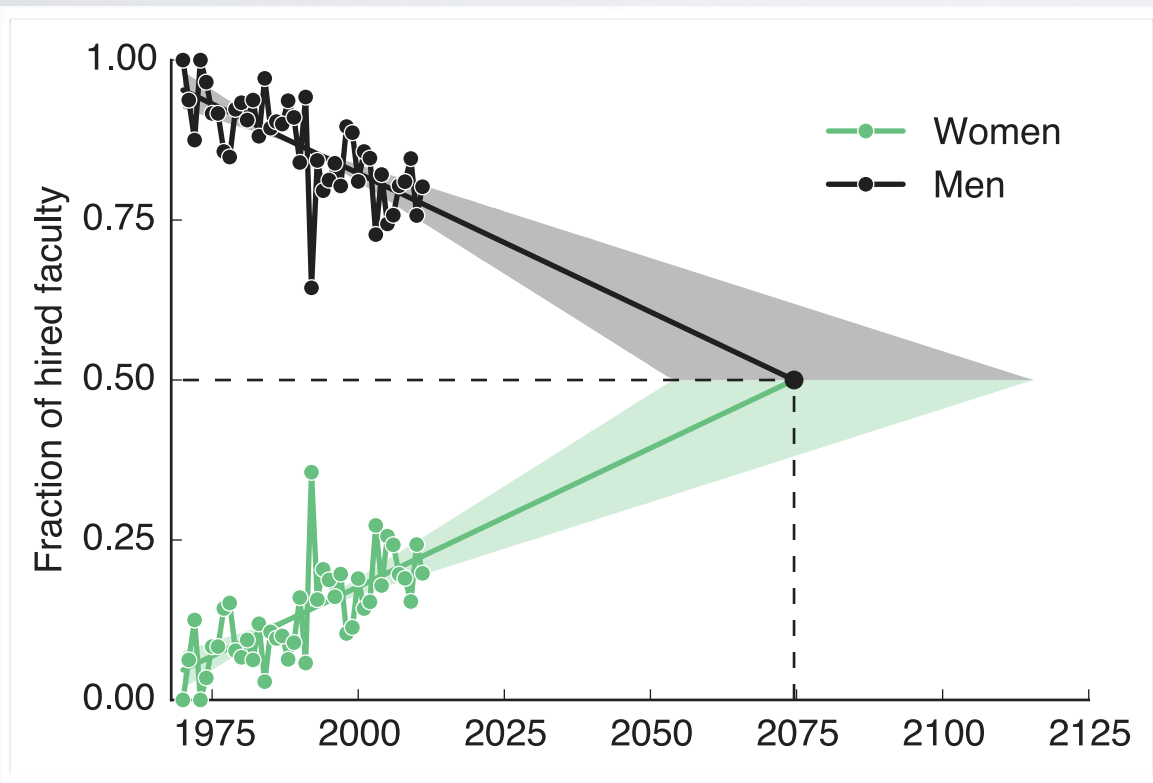
- where is the CS faculty gender ratio going?
- from our 40 years of observable data, the trend is toward parity
- ratio increases by 0.43% per year
- linear extrapolation: *gender parity in new hires around 2075*



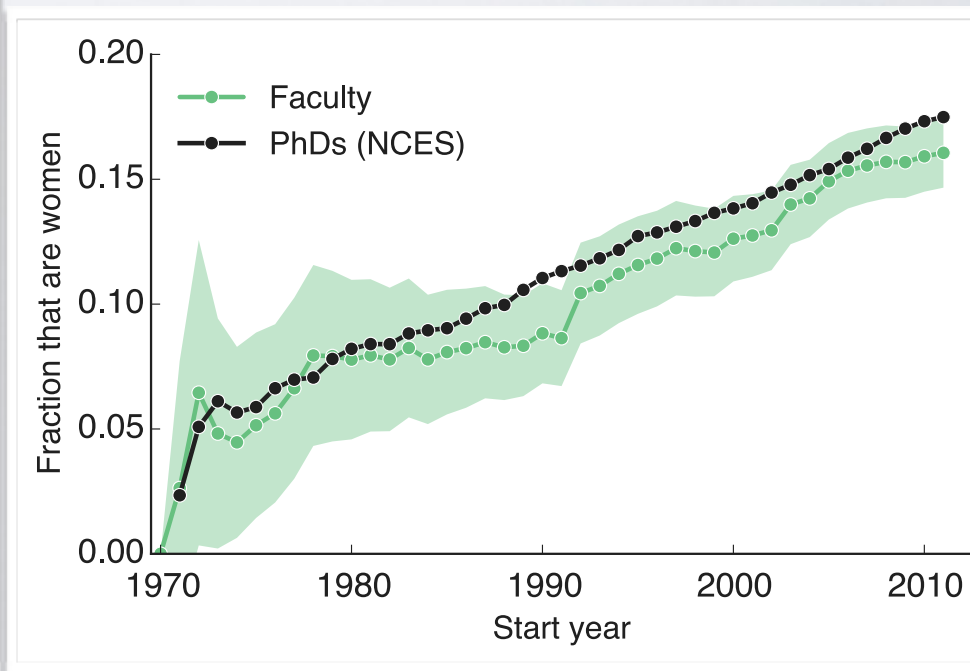
can we change this?!

# gender parity

- where is the CS faculty gender ratio going?
- from our 40 years of observable data, the trend is toward parity
- ratio increases by 0.43% per year
- linear extrapolation: *gender parity in new hires around 2075*



can we change this?!



# the scientific workforce





## the scientific workforce

- common inequalities and hierarchical structure across disciplines
- prestige is cultural influence, via doctoral placement
- prestige predicts\* placement
- dominant core-periphery structure (*homeland vs. colonies*)
- other fields? interdisciplinary work? upstream hires?



## the scientific workforce

- common inequalities and hierarchical structure across disciplines
- prestige is cultural influence, via doctoral placement
- prestige predicts\* placement
- dominant core-periphery structure (*homeland vs. colonies*)
- other fields? interdisciplinary work? upstream hires?

## gender's role in faculty hiring

- no systematic effect alone
- correlates with productivity, postdocs, geography (which are effects)
- interference effect in hiring
- "core" departments drive gender ratios everywhere
- what about other inequalities? other fields?

[1] women *with* a postdoc are as productive as men *without* a postdoc; since 2002, men/women postdoc at similar rates (28%), implying that most female applicants today will appear less productive than most male applicants



**Sam Way  
(Colorado)**



**Dan Larremore  
(Santa Fe)**



**Sam Arbesman  
(Lux Capital)**

RESEARCH ARTICLE

NETWORK SCIENCES

**Systematic inequality and hierarchy in faculty hiring networks**

Aaron Clauset,<sup>1,2,3\*</sup> Samuel Arbesman,<sup>4</sup> Daniel B. Larremore<sup>5,6</sup>

2015 © The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. Distributed under a Creative Commons Attribution Non-Commercial License 4.0 (CC BY-NC). 10.1126/sciadv.1400005

**Science Advances 1(1), e1400005 (2015)**

**Gender, Productivity, and Prestige in Computer Science Faculty Hiring Networks**

Samuel F. Way,<sup>1,\*</sup> Daniel B. Larremore,<sup>2,†</sup> and Aaron Clauset<sup>1,3,2,‡</sup>

<sup>1</sup>Department of Computer Science, University of Colorado, Boulder CO, 80309 USA

<sup>2</sup>Santa Fe Institute, Santa Fe NM, 87501 USA

<sup>3</sup>BioFrontiers Institute, University of Colorado, Boulder CO, 80303 USA

**WWW (2016) [arxiv:1602.00795]**

**Funding:**

- Kauffman Foundation
- University of Colorado Boulder

# Computer Science

Fraction of faculty by type

0.0 0.2 0.4 0.6 0.8 1.0

