

# What drives the productivity of scientific labor?

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



# network data science



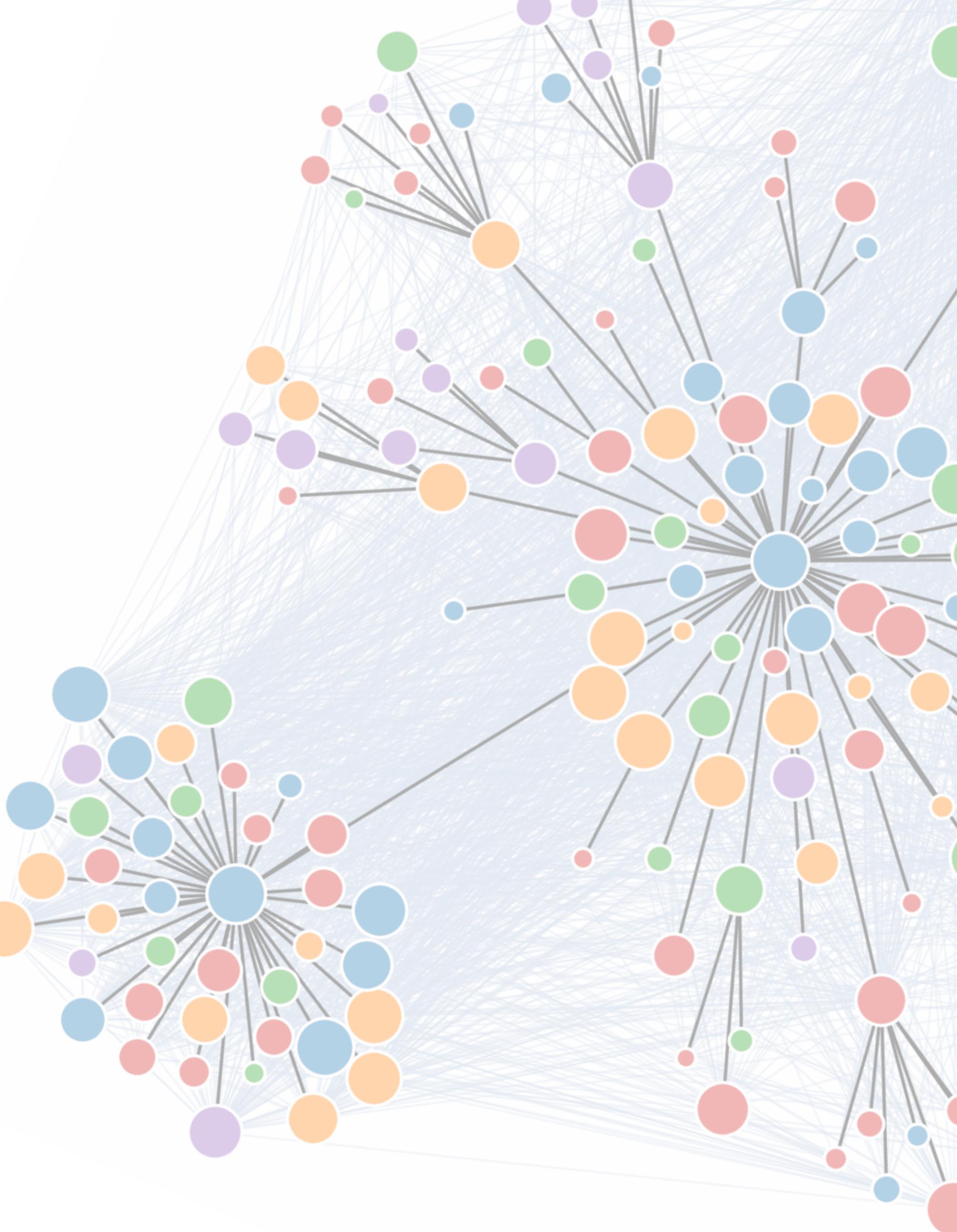
University of Colorado **Boulder**

## Network Analysis and Modeling

This graduate-level course will examine modern techniques for analyzing and modeling the structure and dynamics of complex networks. The focus will be on statistical algorithms and methods, and both lectures and assignments will emphasize model interpretability and understanding the processes that generate real data. Applications will be drawn from computational biology and computational social science. No biological or social science training is required. (Note: this is not a scientific computing course, but there will be plenty of computing for science.)

*Full lectures notes online (~150 pages in PDF)*

<https://aaronclauset.github.io/courses/5352/>



icon.colorado.edu

# ICON

# Colorado Index of Complex Networks

A comprehensive index of research-quality network data sets

## What is ICON?

The Colorado Index of Complex Networks (ICON) is a comprehensive index of research-quality network data sets from all domains of network science, including social, web, information, biological, ecological, connectome, transportation, and technological networks.

Each network record in the index is annotated with and searchable or browsable by its graph properties, description, size, etc., and many records include links to multiple networks. The contents of ICON are curated by volunteer experts from Prof. Aaron Clauset's research group at the University of Colorado Boulder.

Click [NETWORKS](#) to view the networks in the index.

## Network editors & packages

[NetworkX \[python\]](#)

[igraph \[python, R, c++\]](#)

[graph-tool \[python, c++\]](#)

[GraphLab \[python, c++\]](#)

[UCI-Net](#)

[NodeXL](#)

[Gephi](#)

[Pajek](#)

[Network Workbench](#)

[Cytoscape](#)

[yEd graph editor](#)

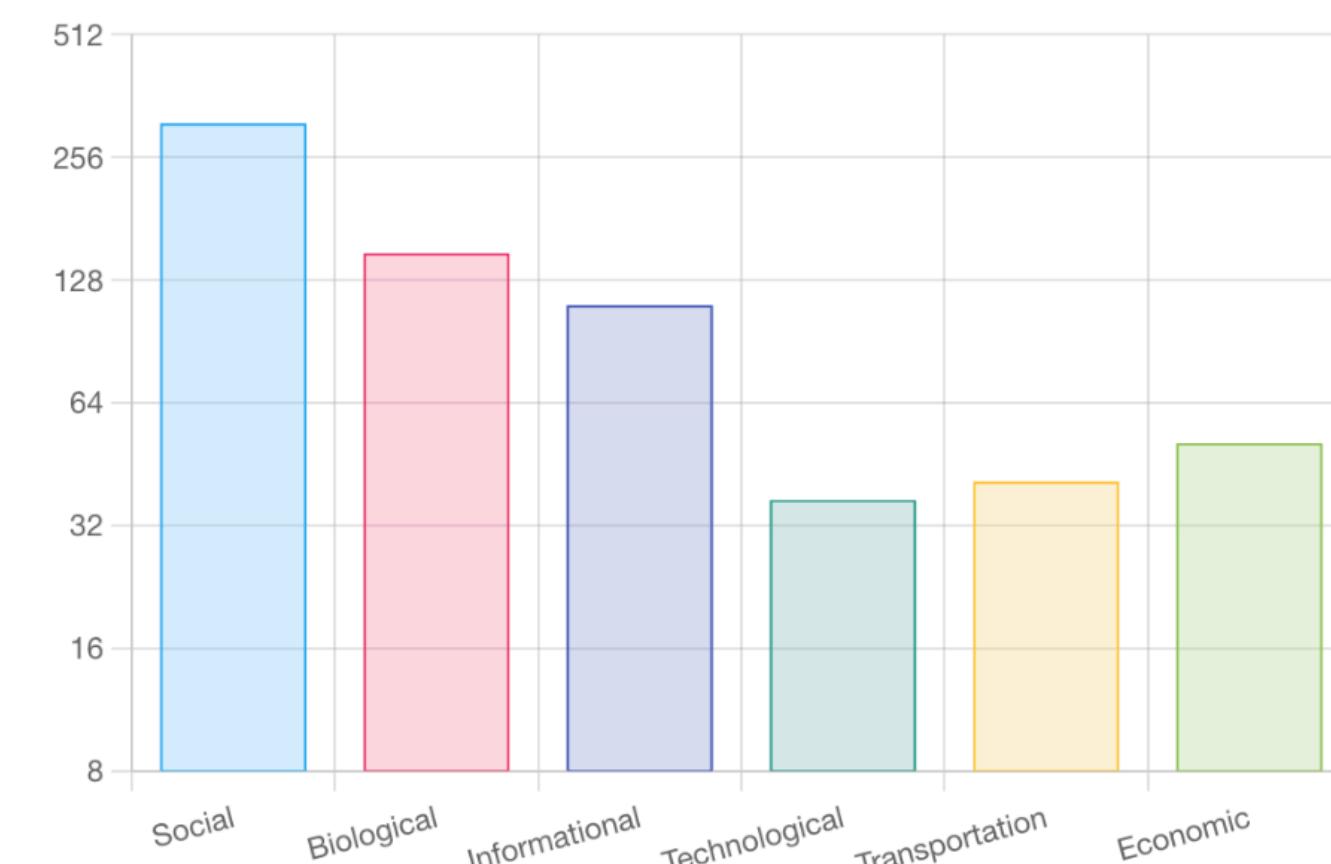
[Graphviz](#)

## Network data sets

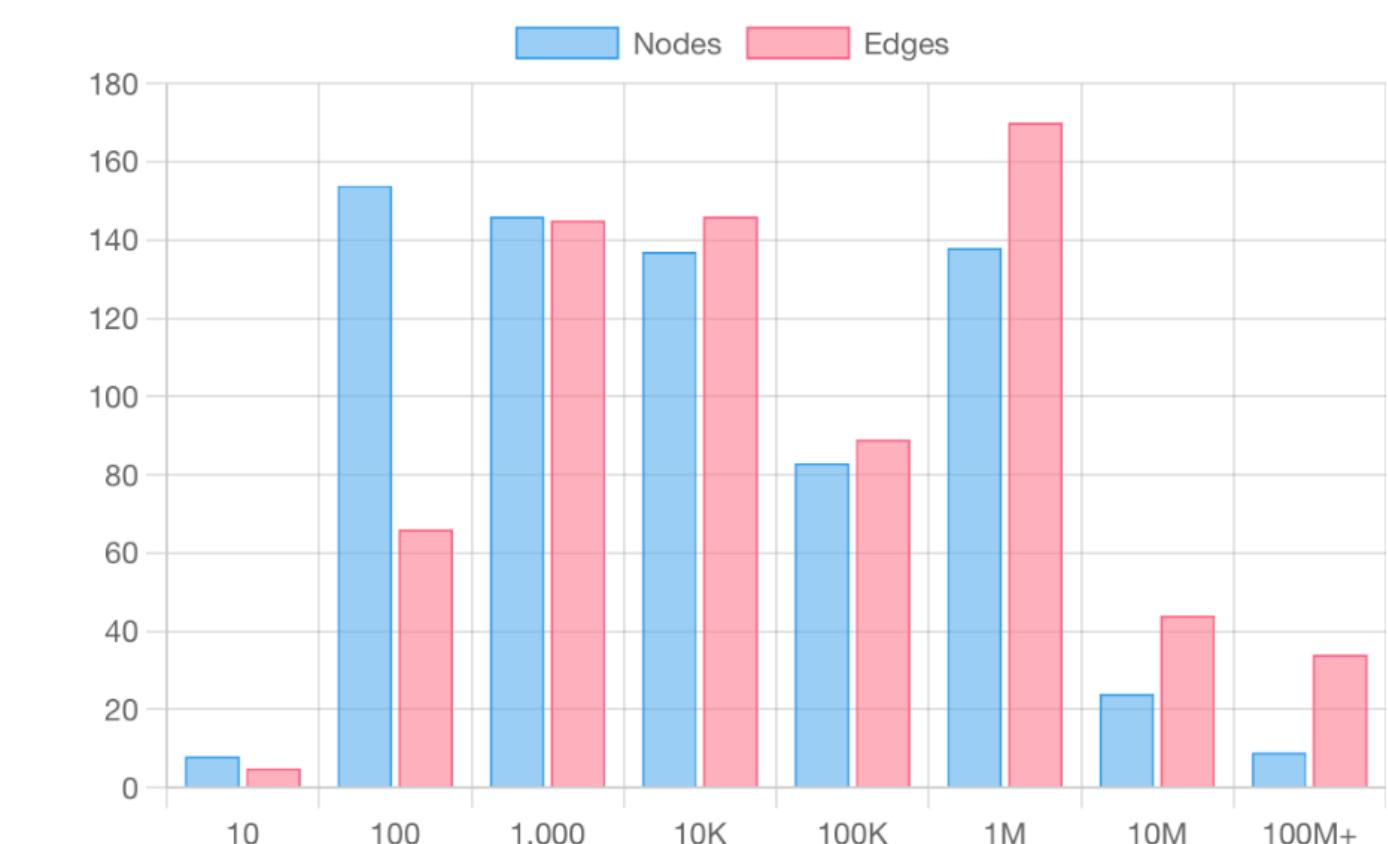
[Colorado Index of Complex Networks](#)

[icon.colorado.edu](#)

Network Domains



Network Nodes and Edges Distribution



# network effects in scientific labor

▶ networks mediate most scientific activities:

scientific training, hiring, collaboration, teaching, attention, peer review, etc.

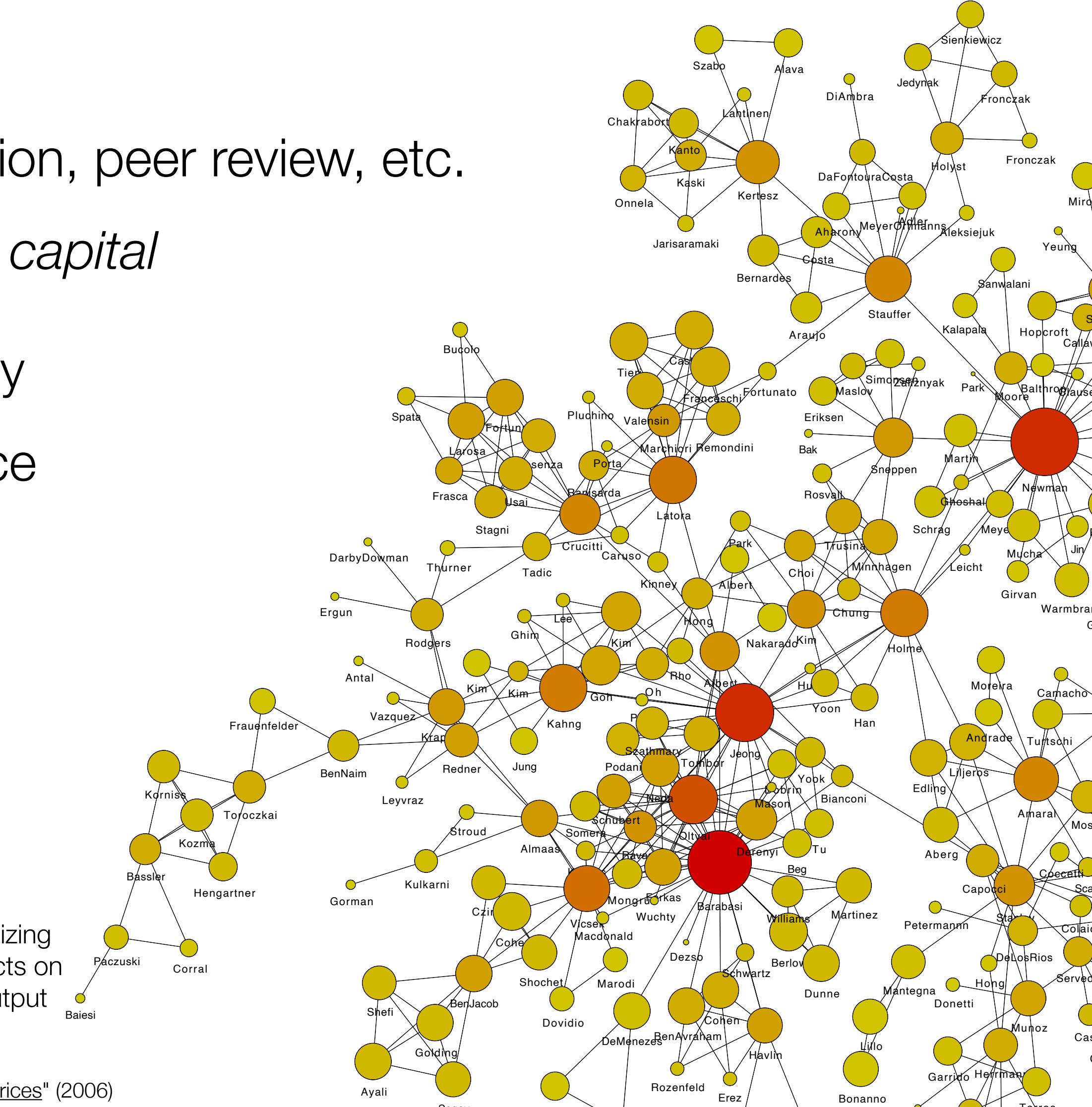
networks act like a form of *unequally distributed social capital*

- a productive collaborator → increases your productivity
- a prominent collaborator → increases your prominence



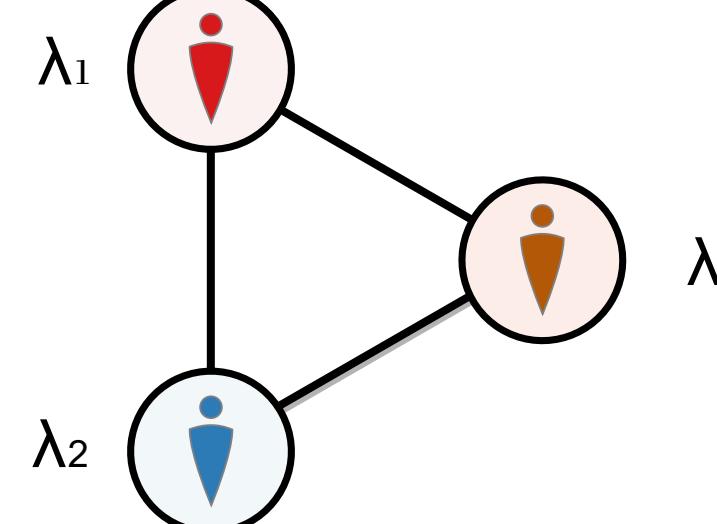
how much does who you work with impact  
your productivity and prominence?

this requires marginalizing  
out the network's effects on  
a node's scholarly output

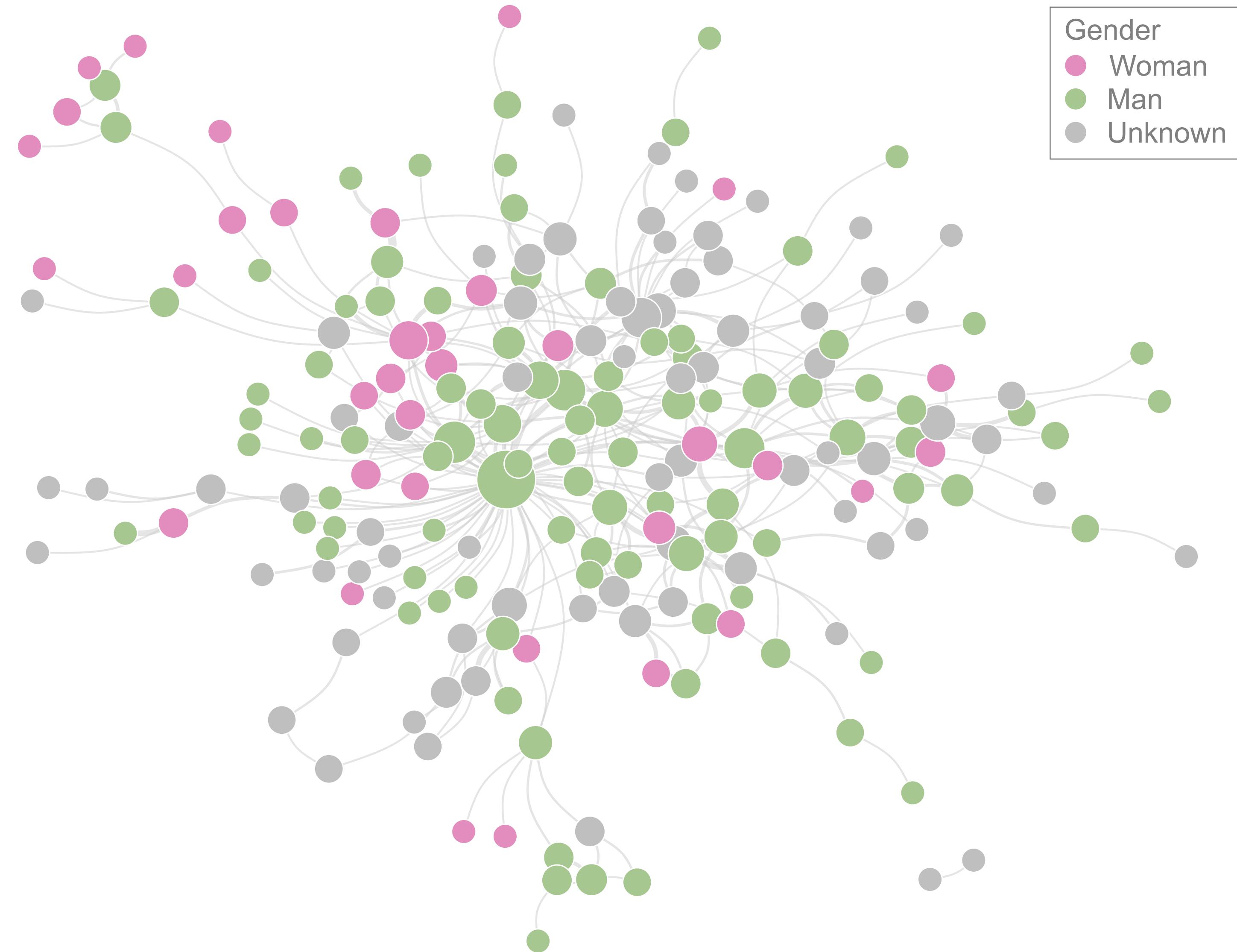


# untangling the network's effects

Coauthorship network

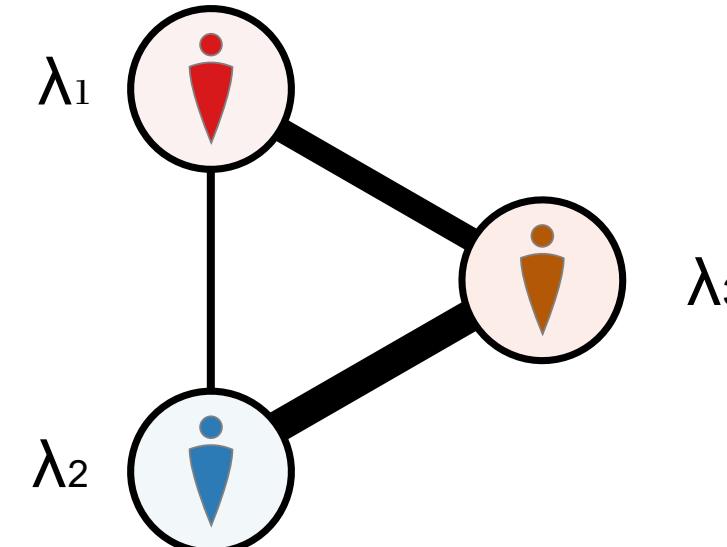


how to estimate  
individual productivity,  
net of coauthors' own  
individual productivity?  
for example...



# untangling the network's effects

Coauthorship network

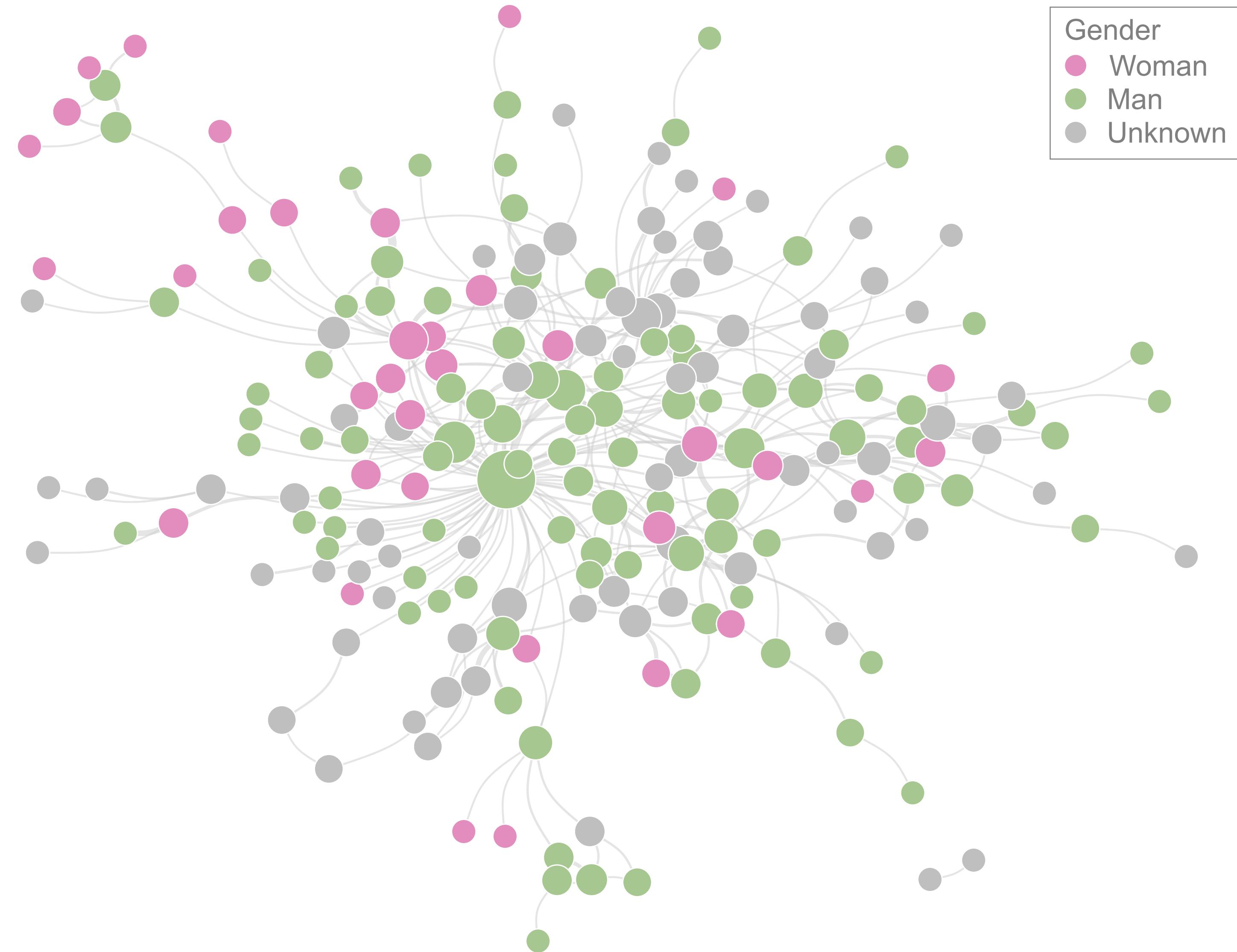


Pairwise productivity

| Author pair | Papers | Time |
|-------------|--------|------|
| $i$ $j$     | 1      | 2    |
| $i$ $k$     | 3      | 2    |
| $j$ $k$     | 4      | 3    |

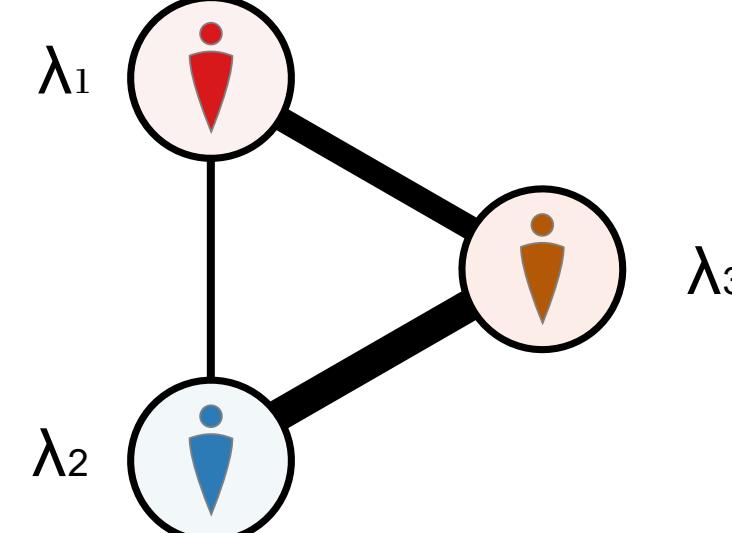
$(i, j)$        $N_{ij}$        $t_{ij}$

who is the most  
individually productive?



# untangling the network's effects

Coauthorship network



Pairwise productivity

| Author pair | Papers   | Time     |
|-------------|----------|----------|
| $i \cdot j$ | 1        | 2        |
| $i \cdot k$ | 3        | 2        |
| $j \cdot k$ | 4        | 3        |
| $(i, j)$    | $N_{ij}$ | $t_{ij}$ |

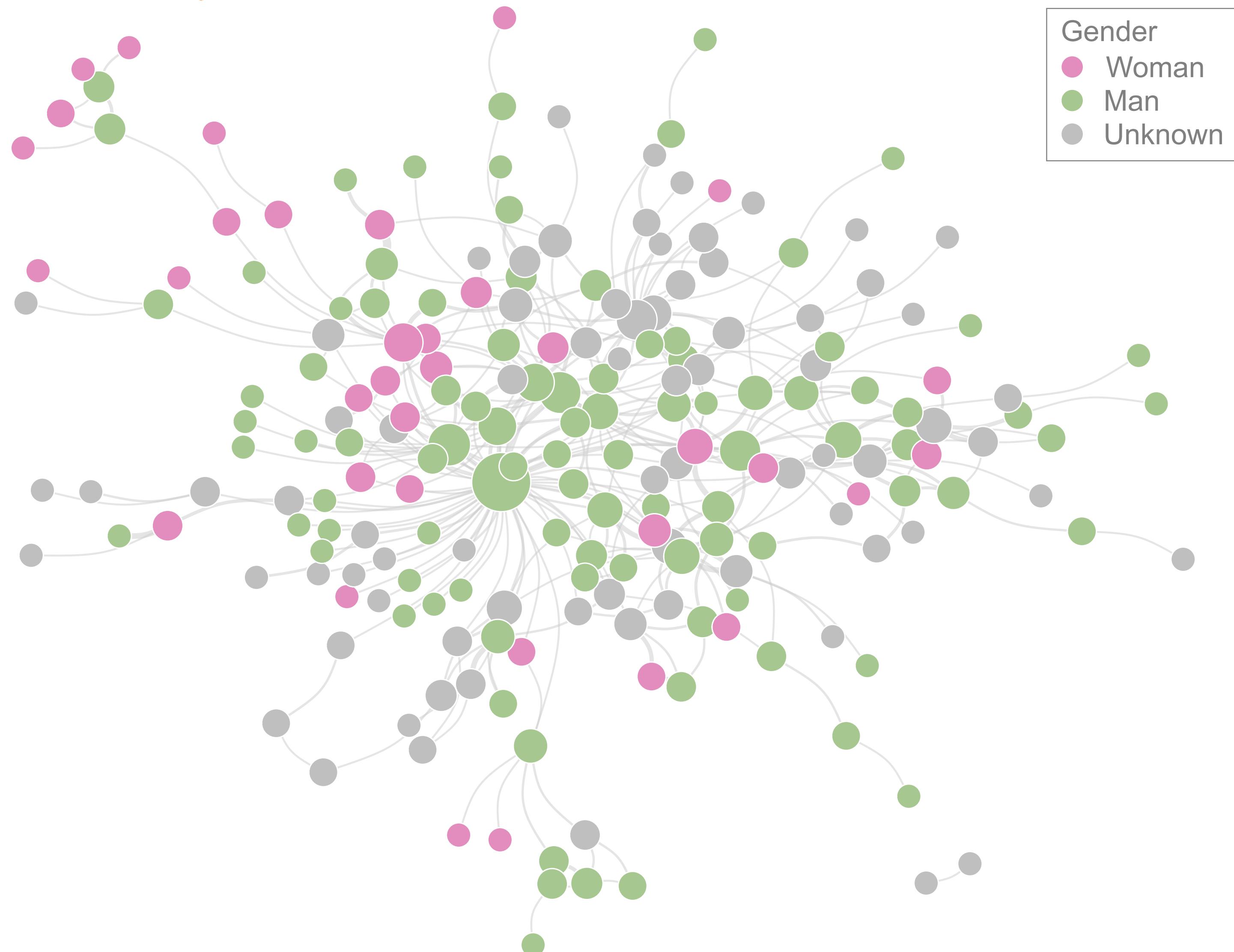
probabilistic model of *pairwise* productivity

number  $(i, j)$ -coauthored papers

$$\Pr(N_{ij}, t_{ij} | \lambda_i, \lambda_j)$$

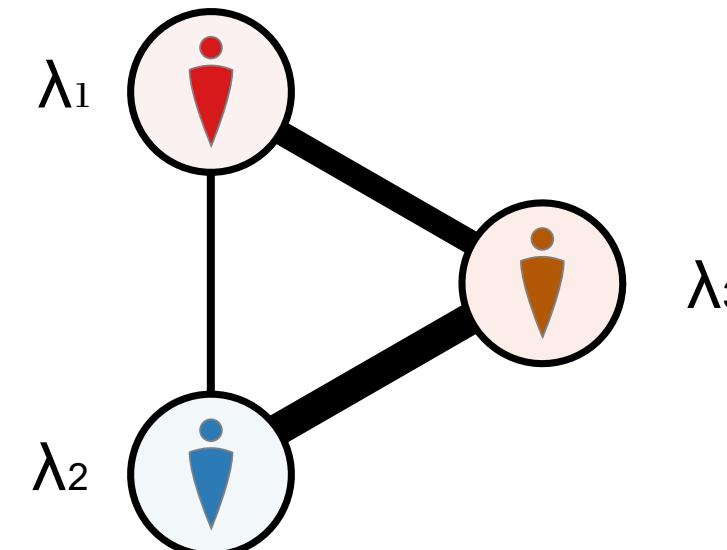
given time period

individual productivities



# untangling the network's effects

Coauthorship network



Pairwise productivity

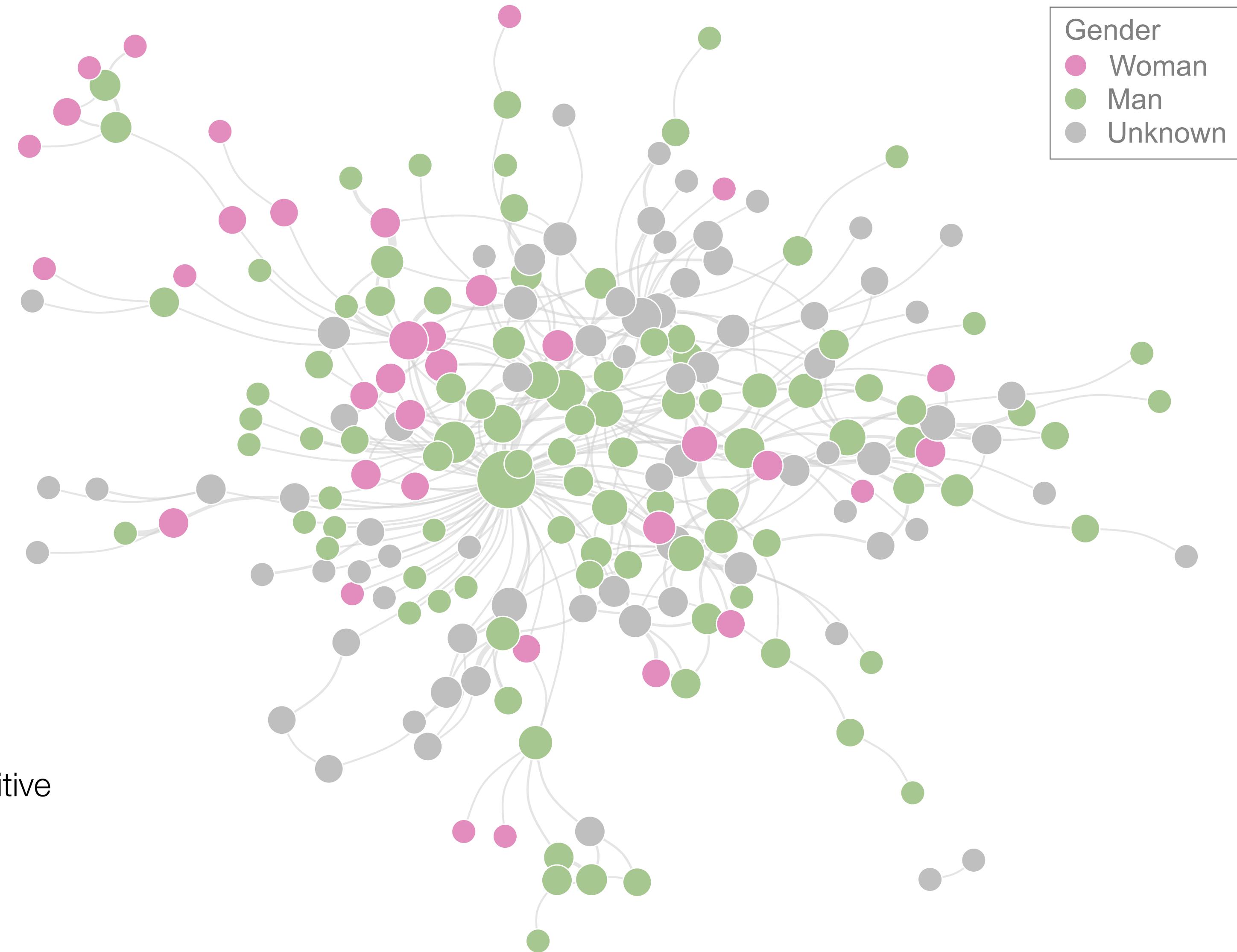
| Author pair | Papers | Time |
|-------------|--------|------|
| • •         | 1      | 2    |
| • •         | 3      | 2    |
| • •         | 4      | 3    |

$(i, j)$        $N_{ij}$        $t_{ij}$

probabilistic model of *pairwise* productivity

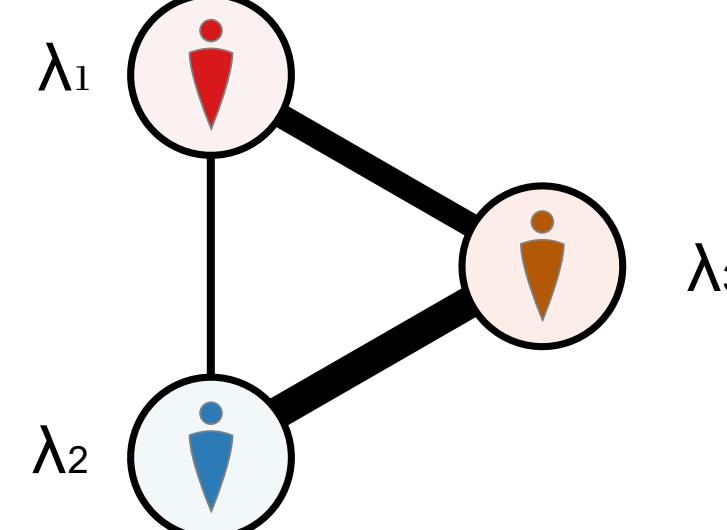
$$\Pr(N_{ij}, t_{ij} | \lambda_i, \lambda_j) = \text{Poisson}([\lambda_i + \lambda_j]t_{ij})$$

productivity is additive



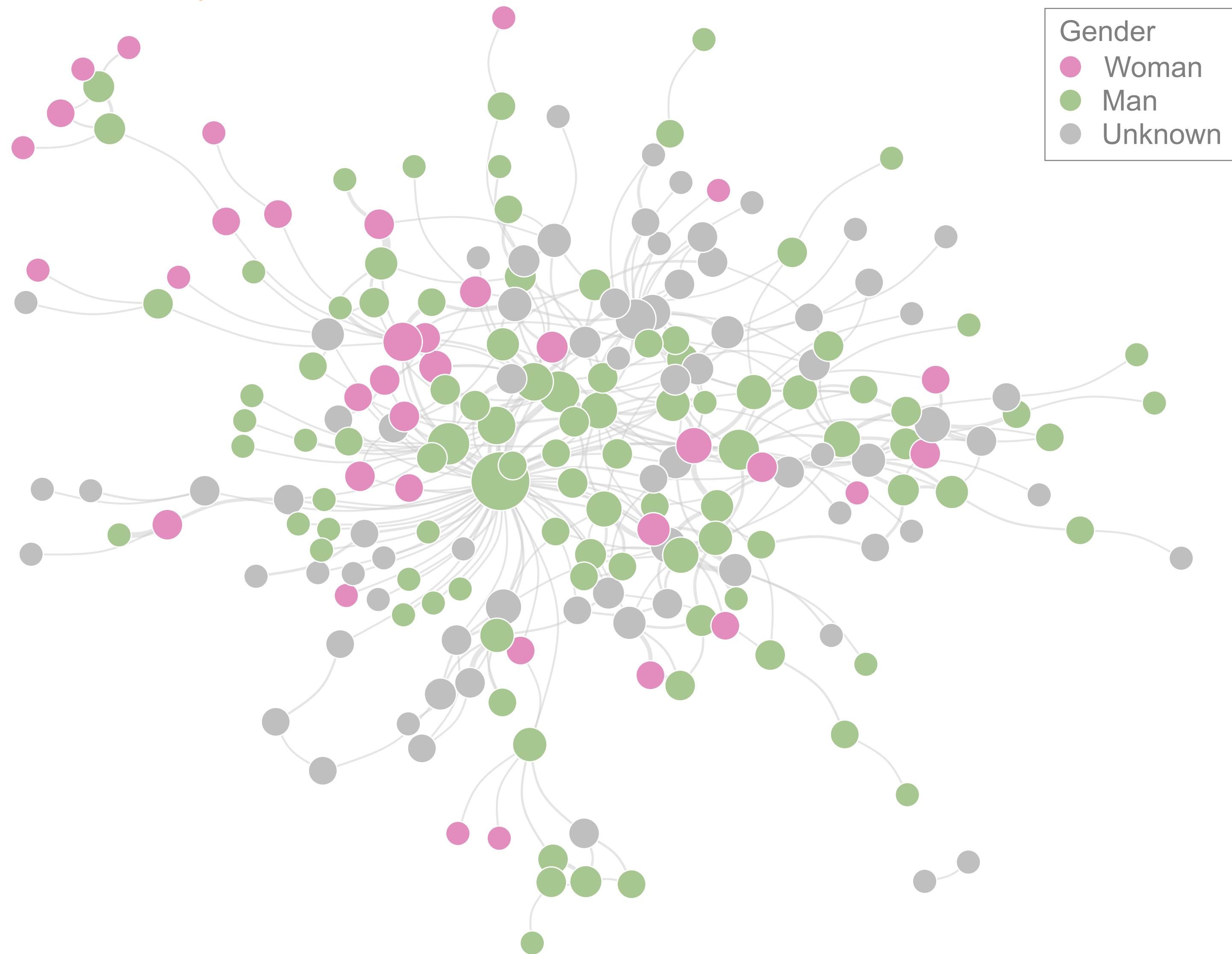
# untangling the network's effects

Coauthorship network



Individual researcher metrics

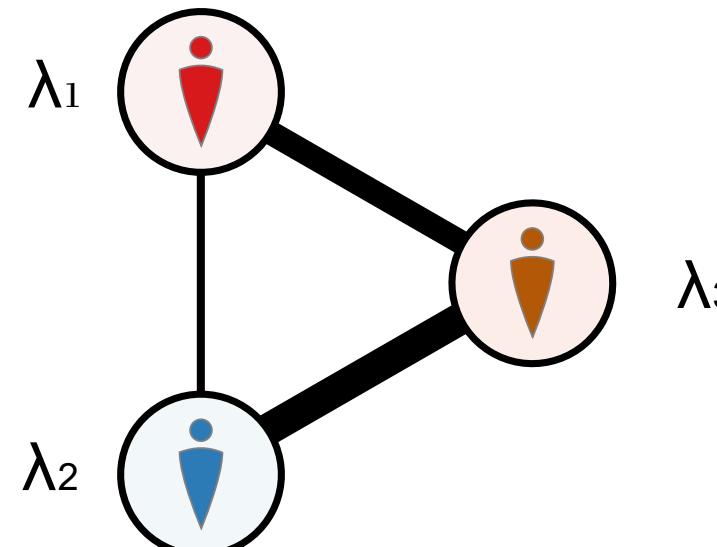
| Author | Papers | Estimate |
|--------|--------|----------|
| $i$    | 4      |          |
| $j$    | 5      |          |
| $k$    | 7      |          |



$$\Pr(N_{ij}, t_{ij} | \lambda_i, \lambda_j) = \text{Poisson}([\lambda_i + \lambda_j]t_{ij})$$

# untangling the network's effects

Coauthorship network



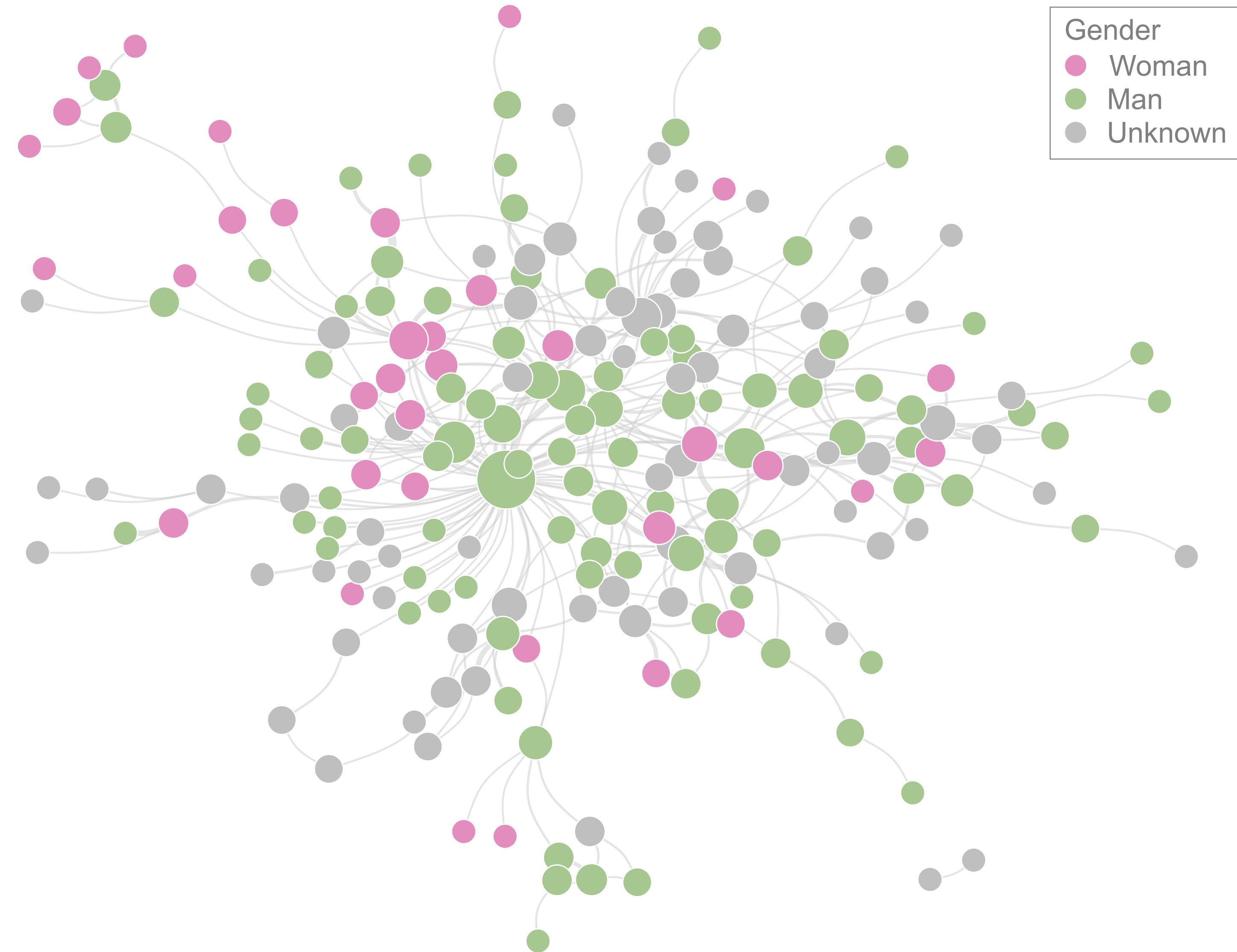
Individual researcher metrics

| Author | Papers | Estimate           |
|--------|--------|--------------------|
| i      | 4      | $\lambda_1 = 0.33$ |
| j      | 5      | $\lambda_2 = 0.17$ |
| k      | 7      | $\lambda_3 = 1.17$ |

$N_i$

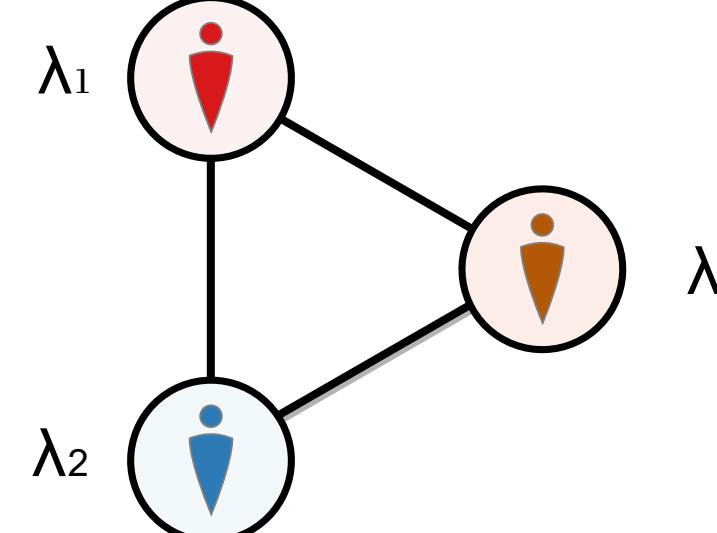
who is the most  
individually productive?

$$\Pr(N_{ij}, t_{ij} | \lambda_i, \lambda_j) = \text{Poisson}([\lambda_i + \lambda_j]t_{ij})$$

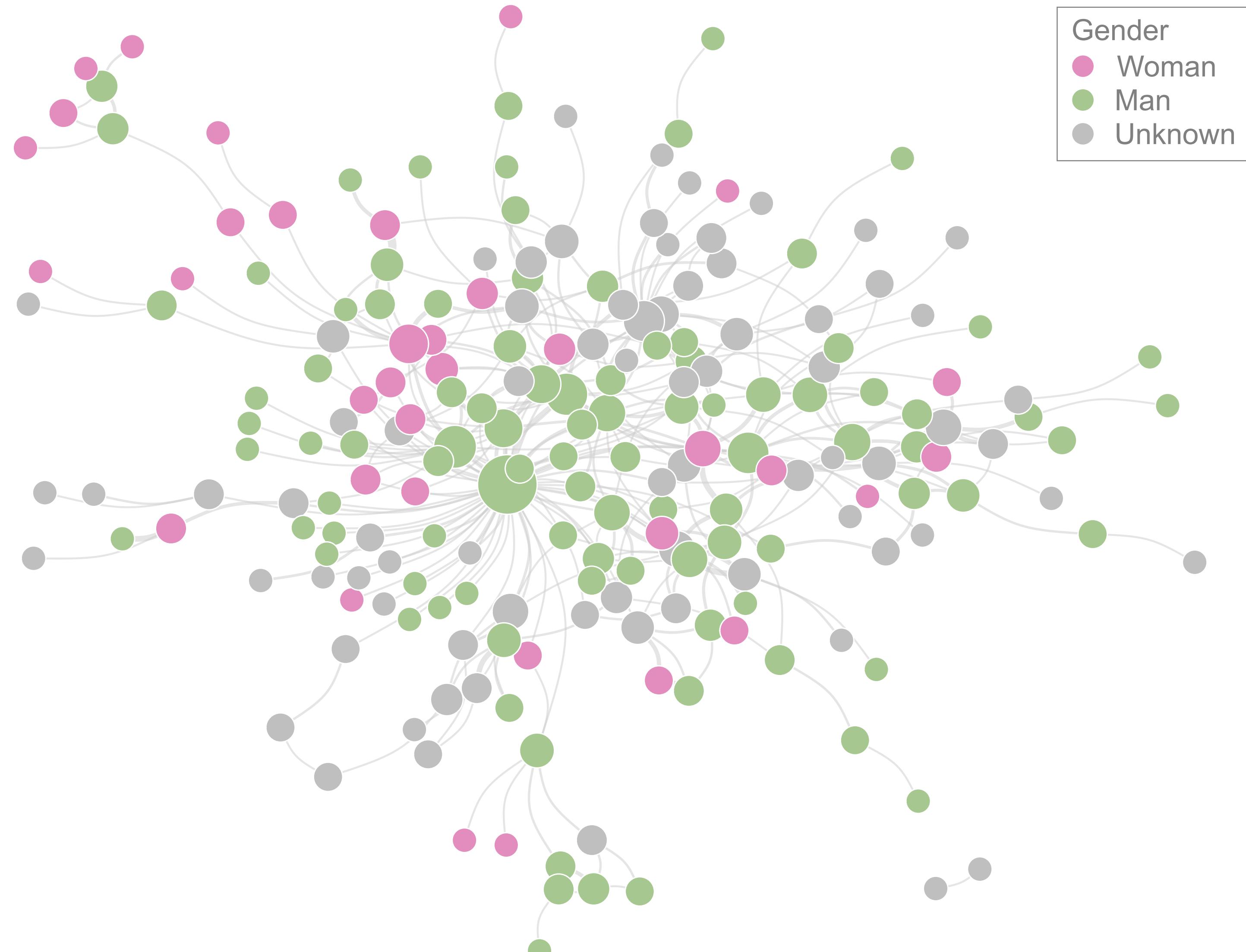


# untangling the network's effects

Coauthorship network

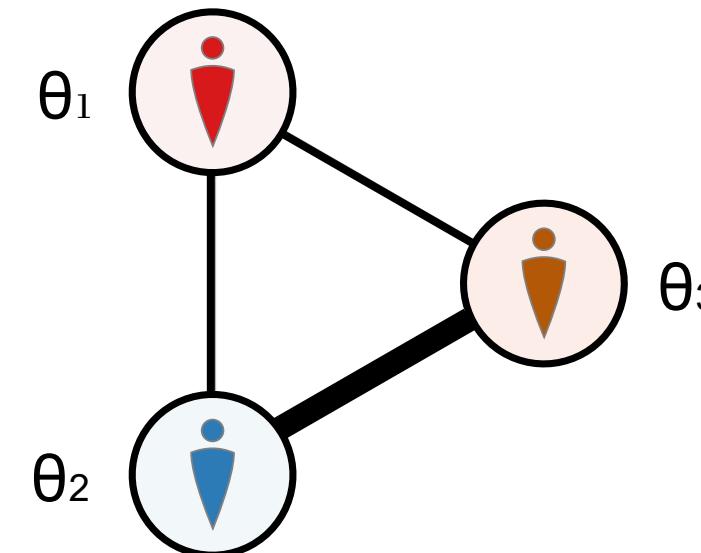


how to estimate  
individual prominence,  
defined as number of  
"high impact" papers?  
for example...



# untangling the network's effects

Coauthorship network

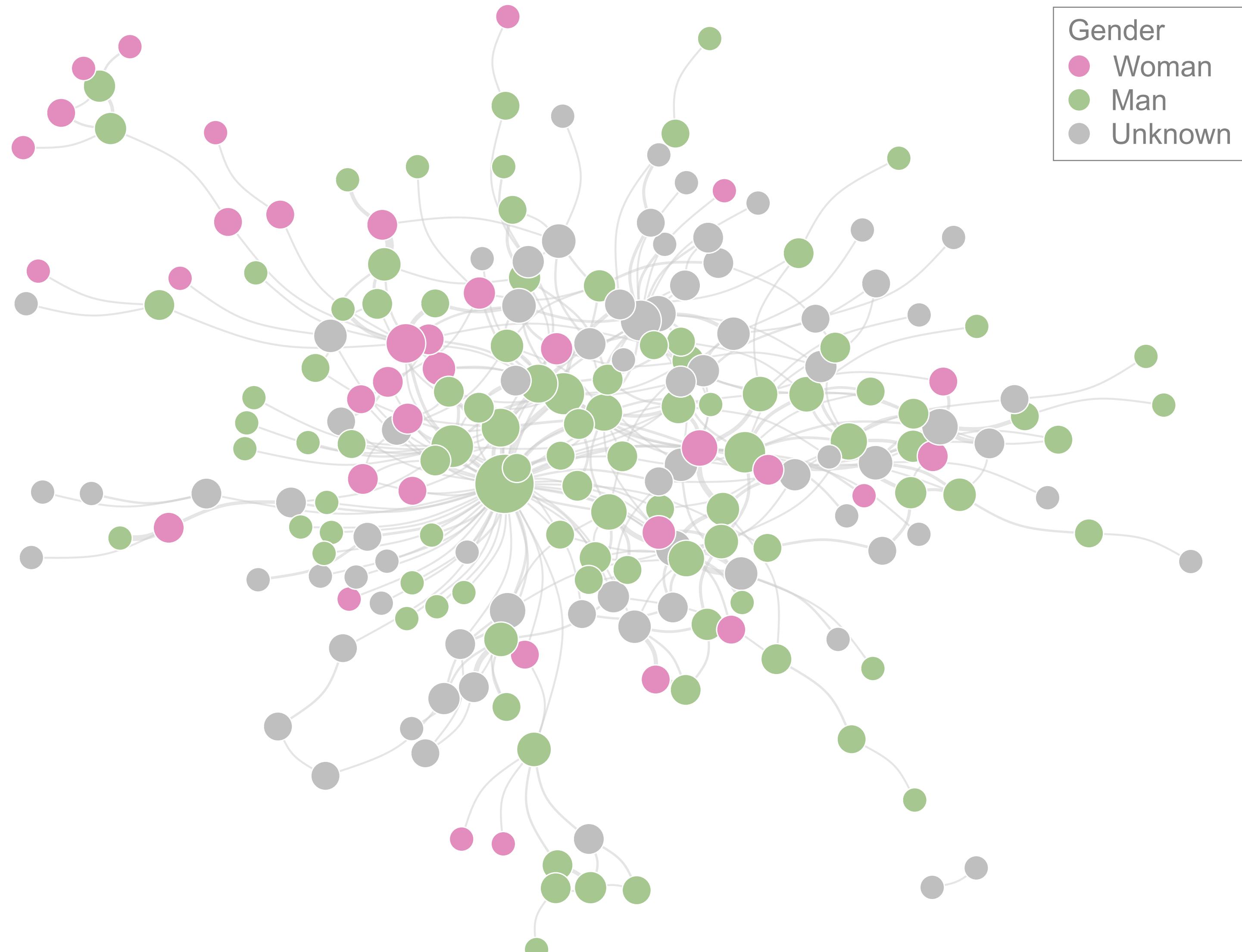


Pairwise impact

| Author pair | Papers | Hit papers |
|-------------|--------|------------|
| • ⚡         | 1      | 1          |
| • ⚡         | 3      | 1          |
| ⚡ ⚡         | 4      | 3          |

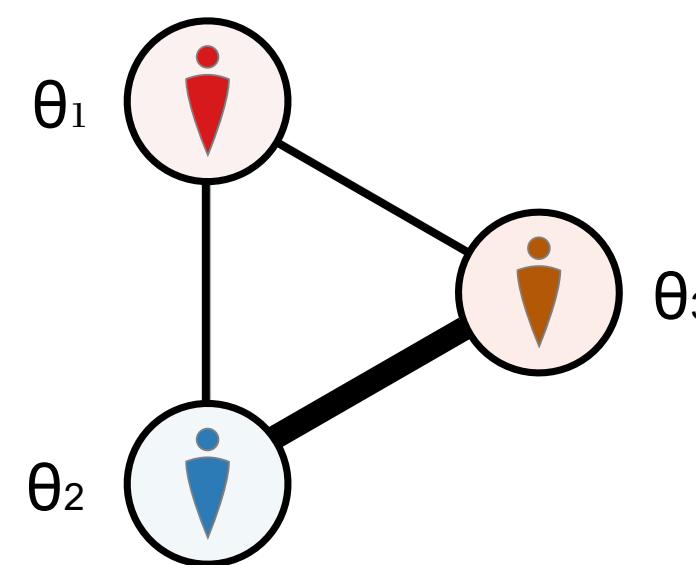
$(i, j)$        $N_{ij}$        $m_{ij}$

who is the most  
individually prominent?



# untangling the network's effects

Coauthorship network



Pairwise impact

| Author pair | Papers   | Hit papers |
|-------------|----------|------------|
| $(i, j)$    | $N_{ij}$ | $m_{ij}$   |
| ●, ●        | 1        | 1          |
| ●, ○        | 3        | 1          |
| ○, ○        | 4        | 3          |

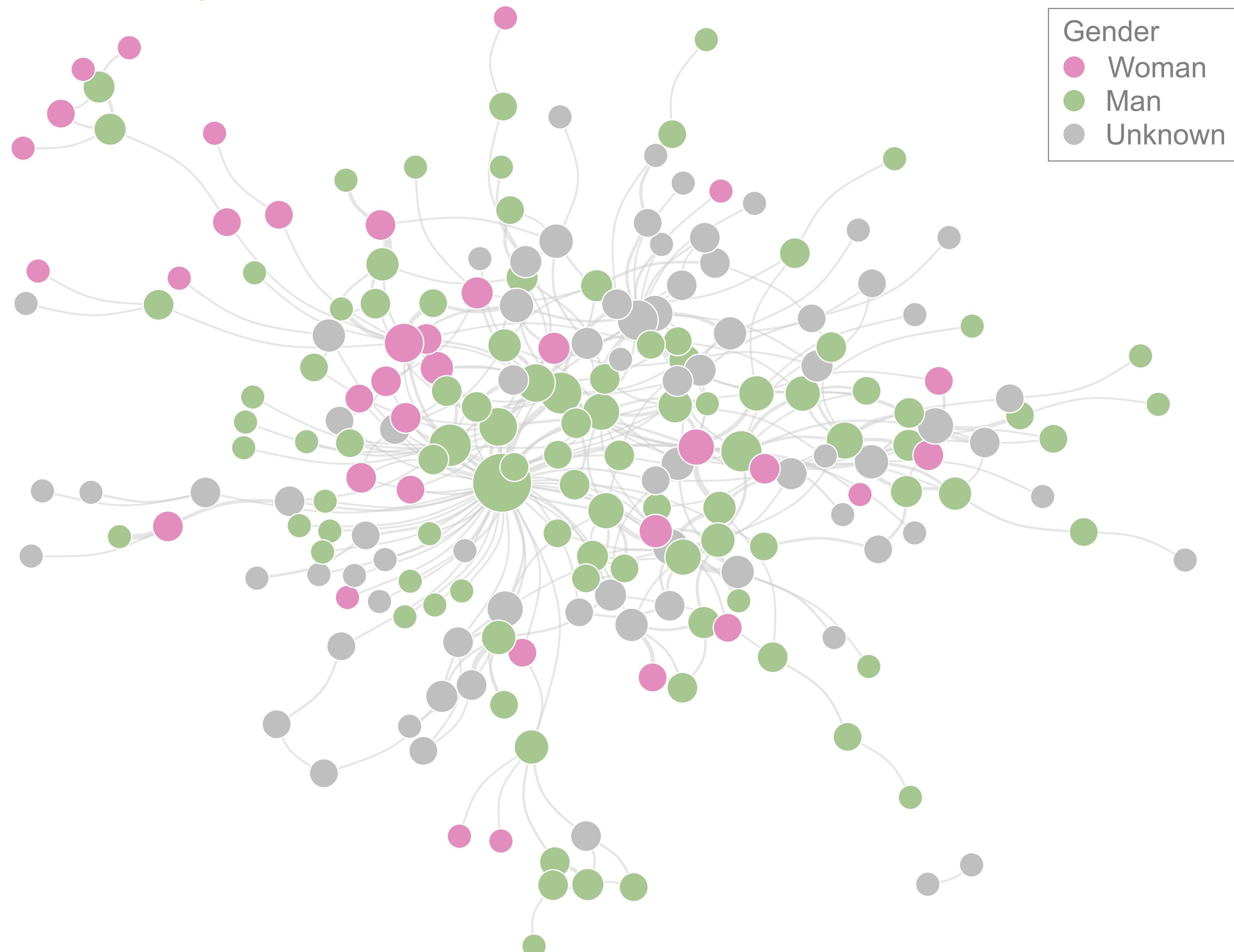
probabilistic model of *pairwise* prominence

number  $(i, j)$ -coauthored papers

$$\Pr(N_{ij}, m_{ij} | \theta_i, \theta_j) =$$

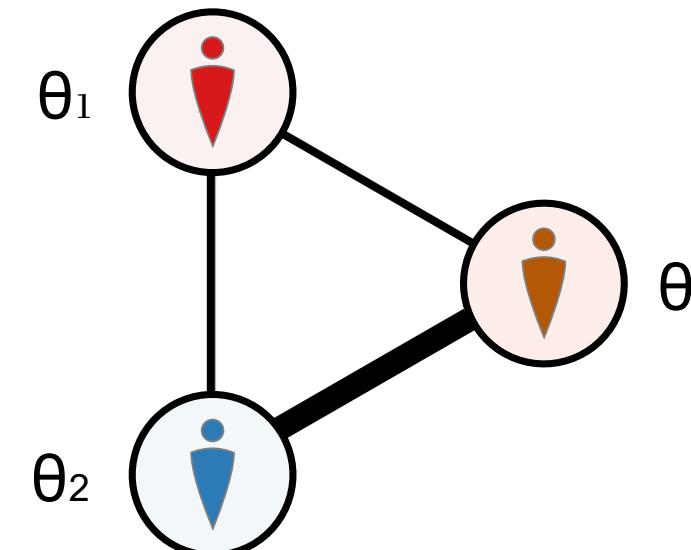
individual prominences

number of "high impact" papers



# untangling the network's effects

Coauthorship network



Pairwise impact

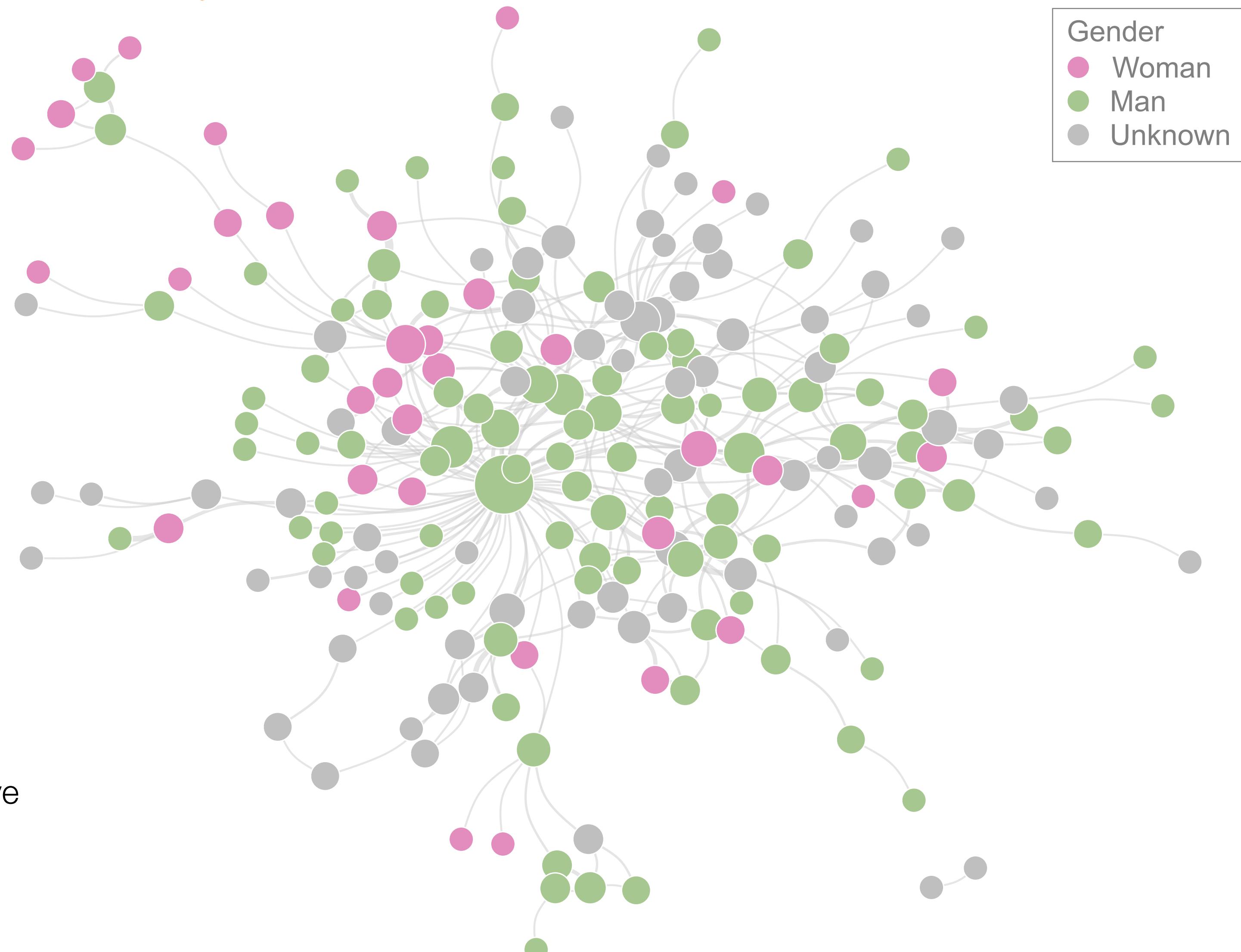
| Author pair | Papers | Hit papers |
|-------------|--------|------------|
| • ↓         | 1      | 1          |
| • ↓         | 3      | 1          |
| ↓ •         | 4      | 3          |

$(i, j)$        $N_{ij}$        $m_{ij}$

probabilistic model of *pairwise* prominence

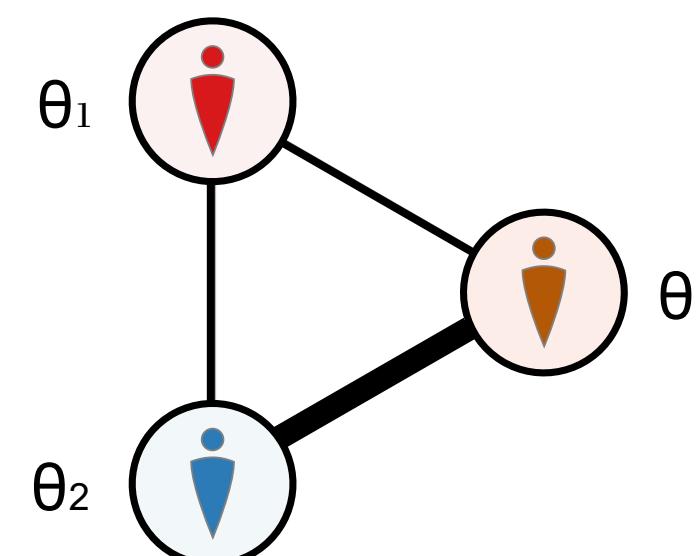
$$\Pr(N_{ij}, m_{ij} | \theta_i, \theta_j) = \text{Binomial}(N_{ij}, [\theta_i + \theta_j])$$

↑  
prominence is additive

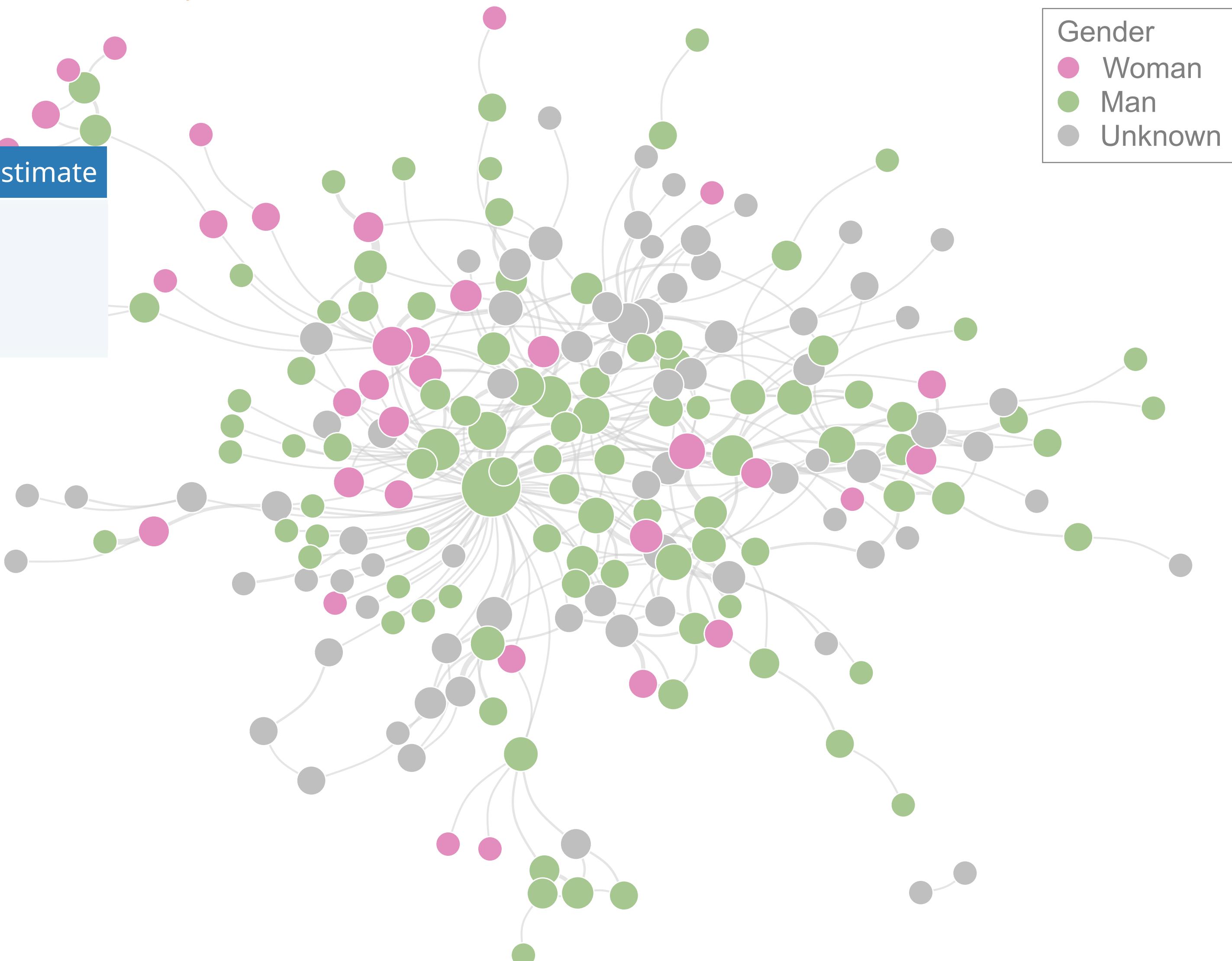


# untangling the network's effects

Coauthorship network



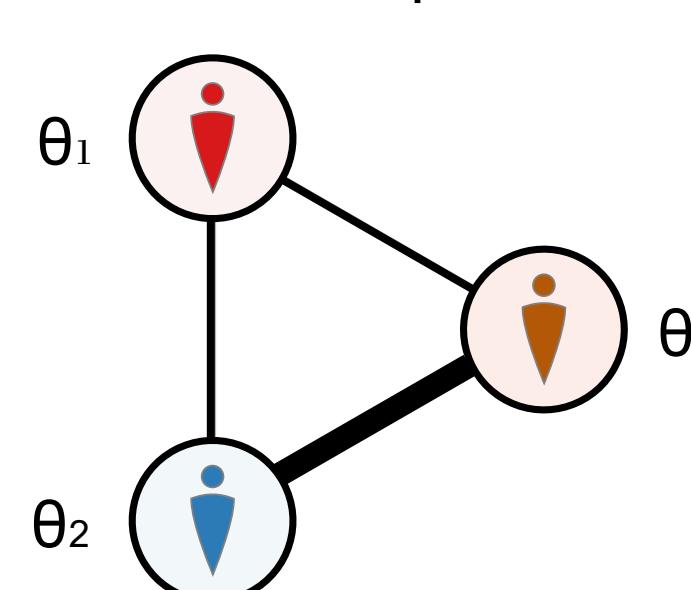
| Author | Papers | Hit papers | Estimate           | Estimate |
|--------|--------|------------|--------------------|----------|
| $i$    | $N_i$  | $m_i$      |                    |          |
| ●      | 4      | 2          | $\lambda_1 = 0.33$ |          |
| ●      | 5      | 4          | $\lambda_2 = 0.17$ |          |
| ●      | 7      | 4          | $\lambda_3 = 1.17$ |          |



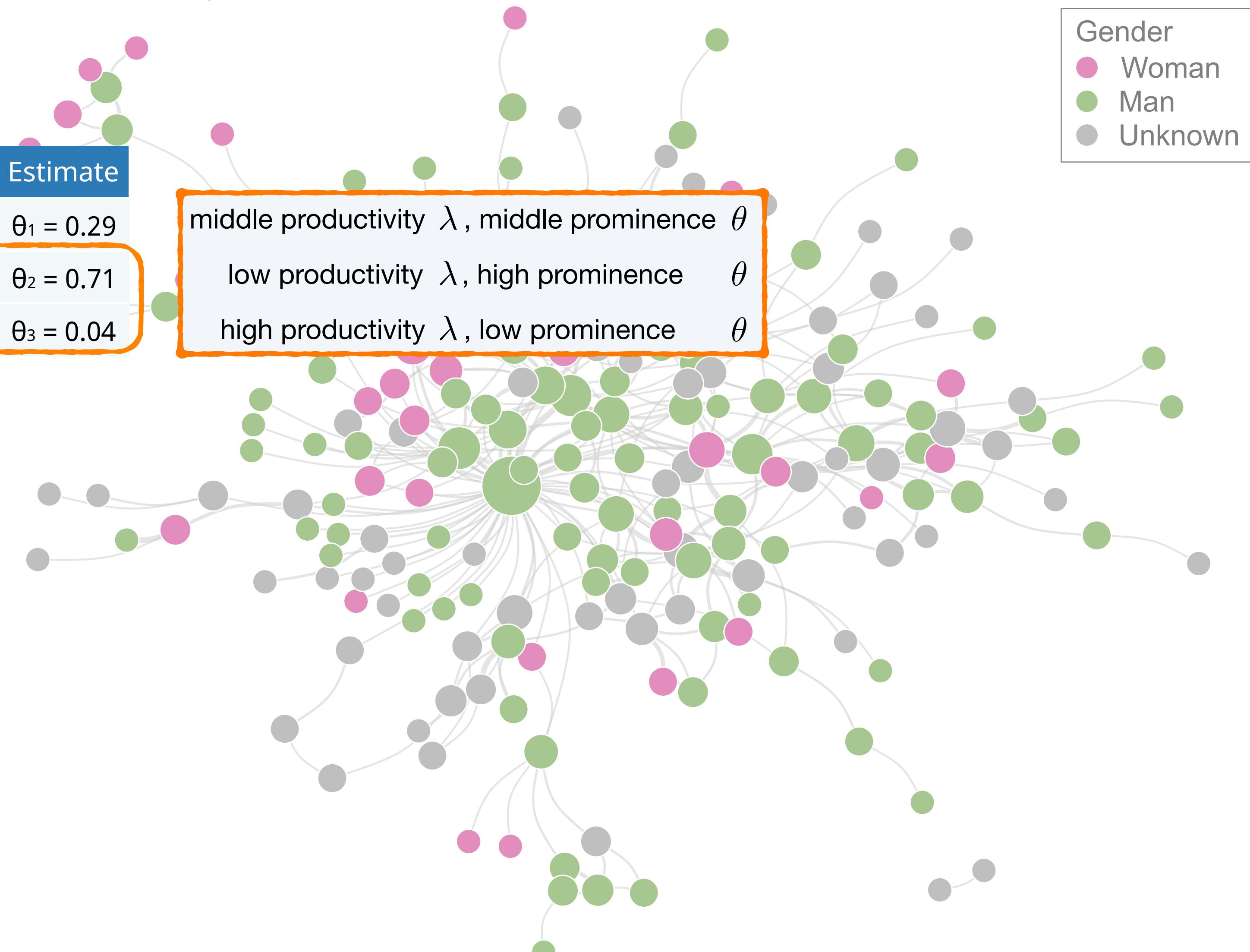
$$\Pr(N_{ij}, m_{ij} | \theta_i, \theta_j) = \text{Binomial}(N_{ij}, [\theta_i + \theta_j])$$

# untangling the network's effects

Coauthorship network



| Author     | Papers | Hit papers | Estimate           | Estimate          |
|------------|--------|------------|--------------------|-------------------|
| $i$        | $N_i$  | $m_i$      |                    |                   |
| ● (red)    | 4      | 2          | $\lambda_1 = 0.33$ | $\theta_1 = 0.29$ |
| ● (blue)   | 5      | 4          | $\lambda_2 = 0.17$ | $\theta_2 = 0.71$ |
| ● (orange) | 7      | 4          | $\lambda_3 = 1.17$ | $\theta_3 = 0.04$ |

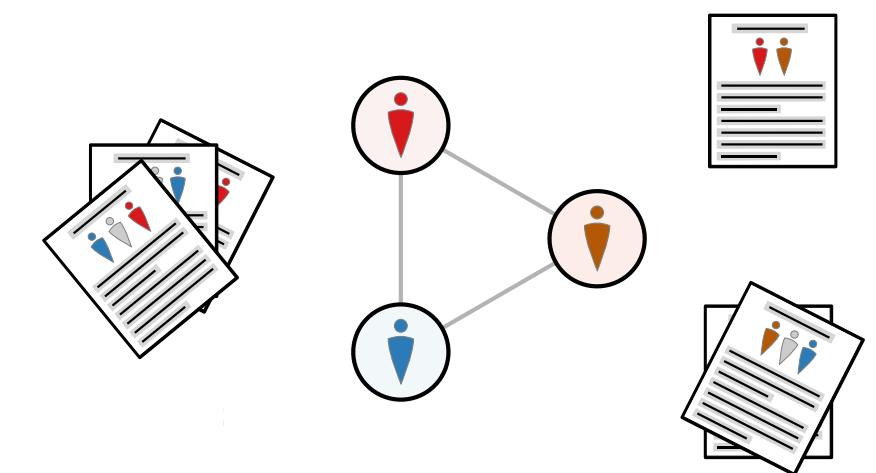


# untangling the network's effects

Microsoft Academic Graph (MAG) 1950–2019



- 198,202 mid-career researchers, with 10+ papers by 15th year of publishing history
- spanning 6 STEM fields (biology, chemistry, CS, math, medicine, physics)
- analyze only first-last author pairs (mitigates middle-author effects)
- 'high impact' paper = in upper 8% for given year-field

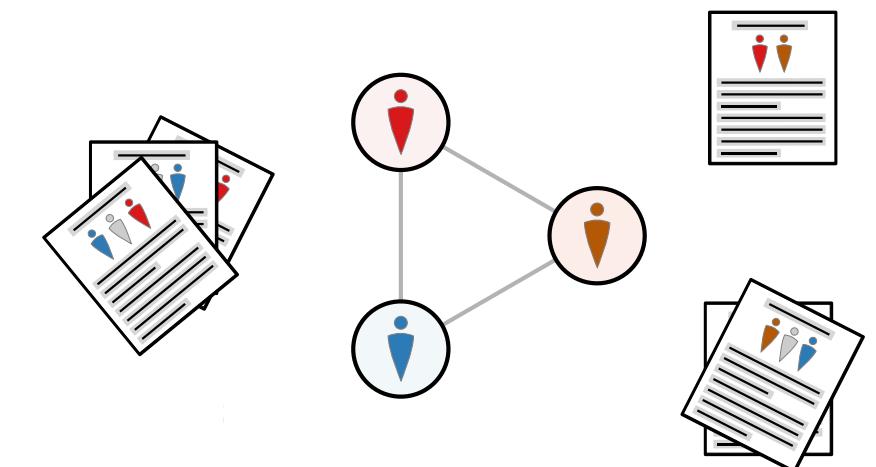


# untangling the network's effects

Microsoft Academic Graph (MAG) 1950–2019



- 198,202 mid-career researchers, with 10+ papers by 15th year of publishing history
- spanning 6 STEM fields (biology, chemistry, CS, math, medicine, physics)
- analyze only first-last author pairs (mitigates middle-author effects)
- 'high impact' paper = in upper 8% for given year-field

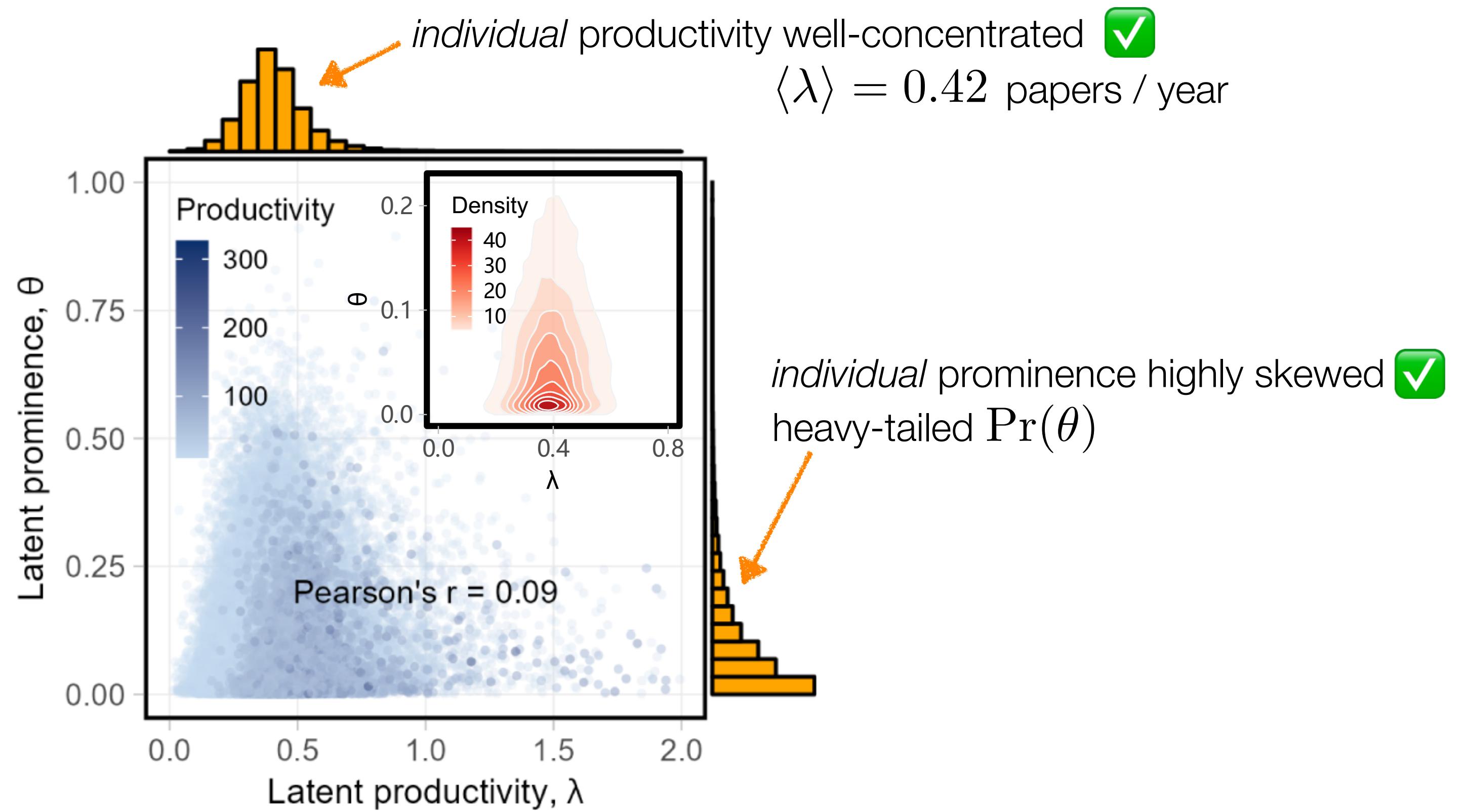


Estimate individual  $(\lambda_i, \theta_i)$  parameters for each scientist

- estimate using network 1950– $T$ , for variable  $T$  in [1975,2017]
- bootstrapped convex optimization, then record mean  $\lambda_i$  and  $\theta_i$
- investigate how  $(\lambda_i, \theta_i)$  covary with gender, prestige, etc.

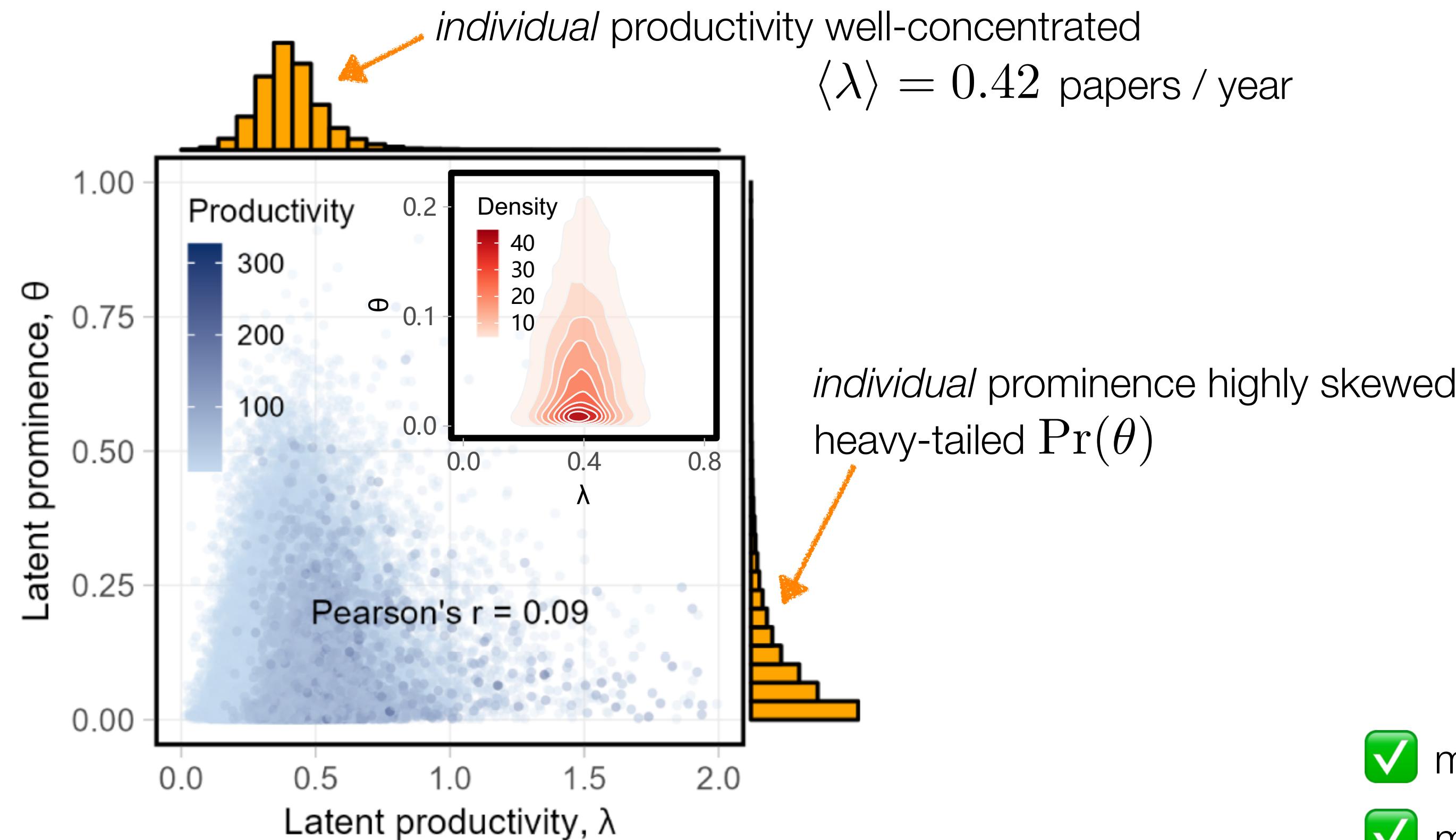
# model checking

▶ applied to 198,202 mid-career STEM researchers 1975-2017



# model checking

▶ applied to 198,202 mid-career STEM researchers 1975-2017

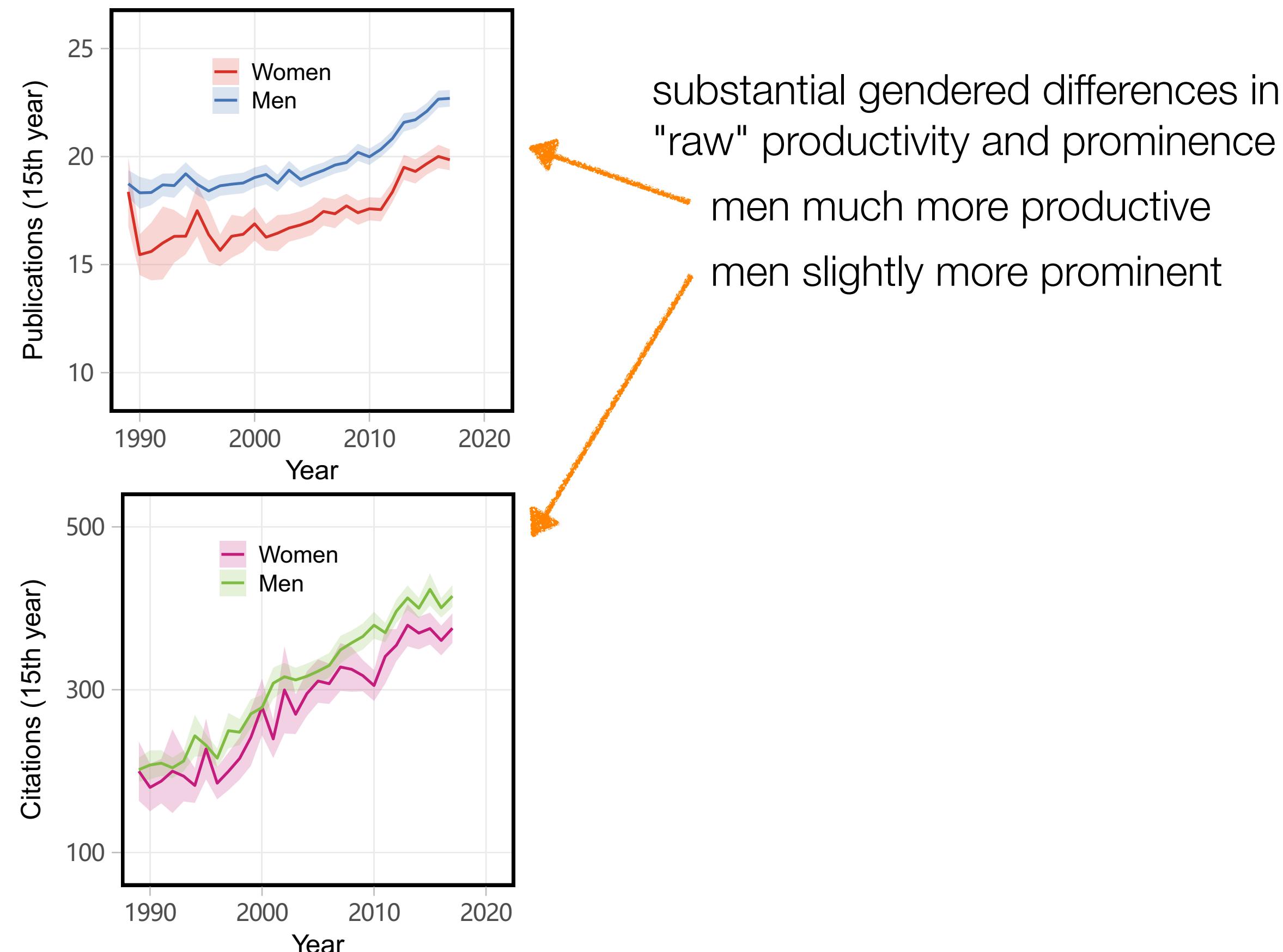


|                          | Prestige | Papers | Citations | $\lambda$ | $\theta$ | High $\lambda$ coauthors | High $\theta$ coauthors |
|--------------------------|----------|--------|-----------|-----------|----------|--------------------------|-------------------------|
| Prestige                 |          | 0.06   | 0.15      | 0.02      | 0.15     | 0.04                     | 0.13                    |
| Papers                   | 0.06     |        | 0.4       | 0.21      | -0.02    | 0.7                      | 0.44                    |
| Citations                | 0.15     | 0.4    |           | 0.12      | 0.38     | 0.27                     | 0.49                    |
| $\lambda$                | 0.02     | 0.21   | 0.12      |           | 0.15     | 0.31                     | 0.14                    |
| $\theta$                 | 0.15     | -0.02  | 0.38      | 0.15      |          | 0.06                     | 0.25                    |
| High $\lambda$ coauthors | 0.04     | 0.7    | 0.27      | 0.31      | 0.06     |                          | 0.43                    |
| High $\theta$ coauthors  | 0.13     | 0.44   | 0.49      | 0.14      | 0.25     | 0.43                     |                         |

- ✓ my paper count = highly correlated with high  $\lambda$  coauthors
- ✓ my citations = well correlated with high  $\lambda$  &  $\theta$  coauthors

# gender vs. productivity & prominence

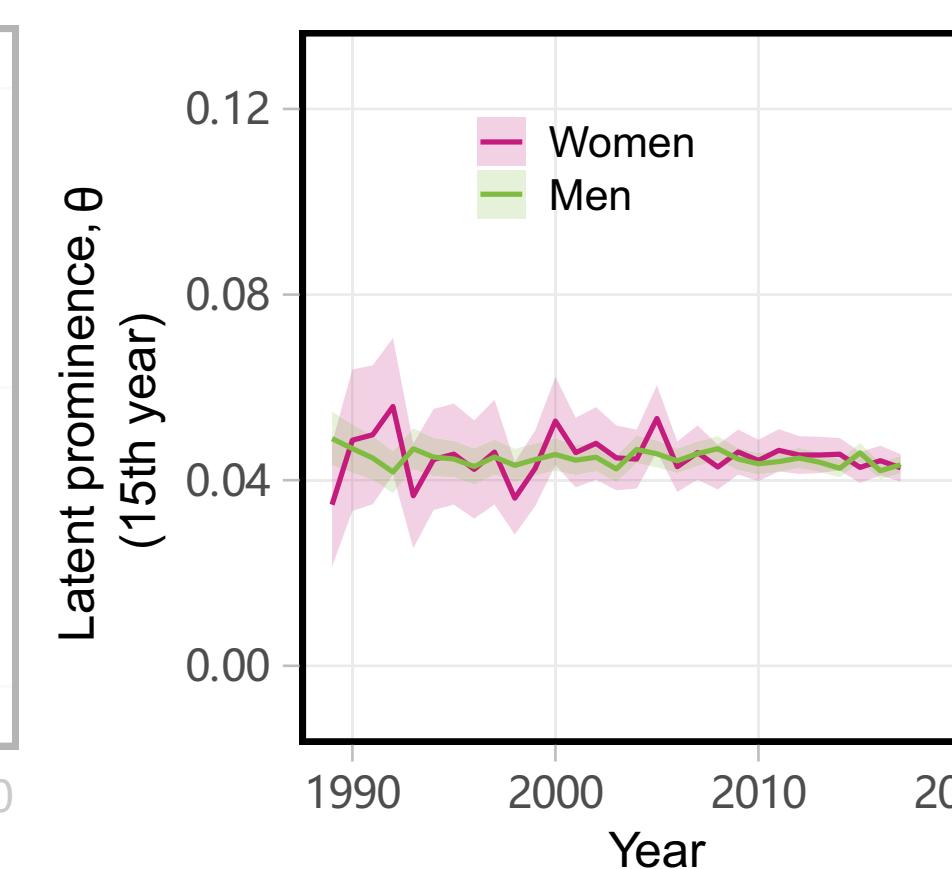
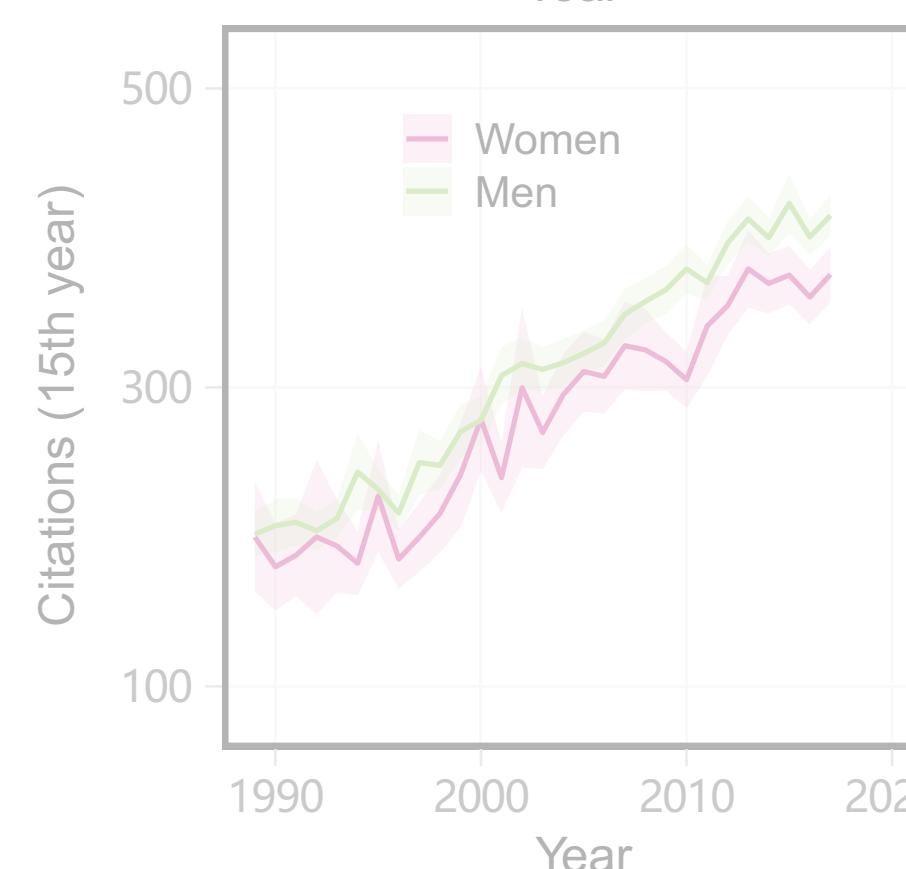
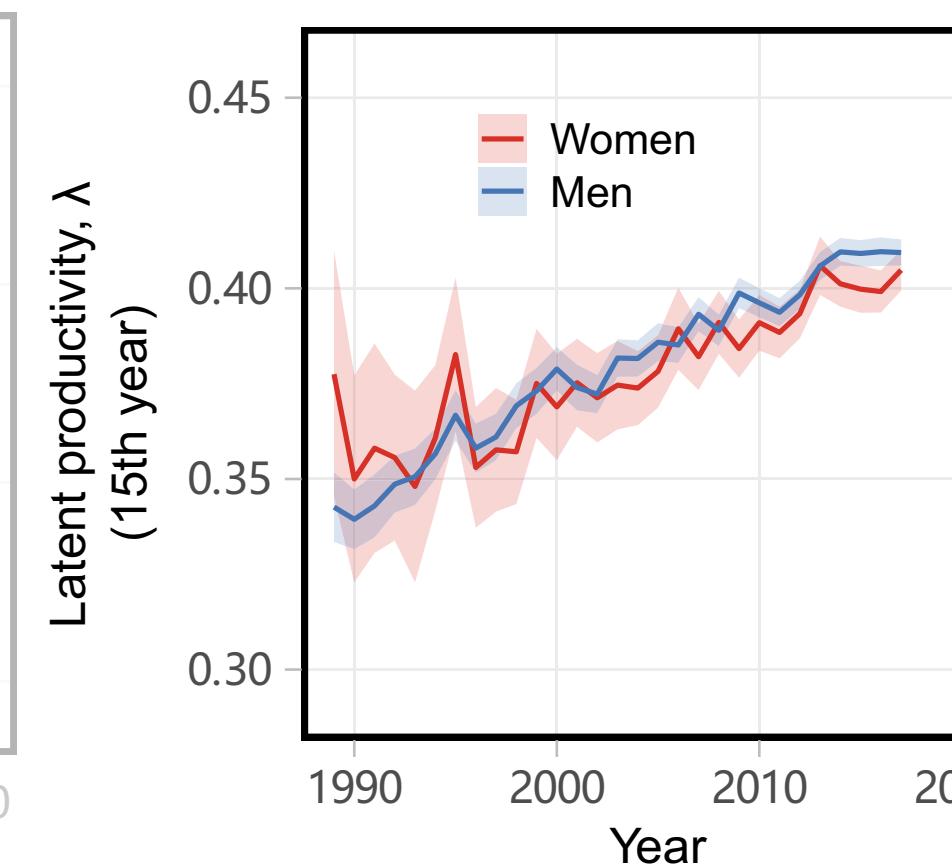
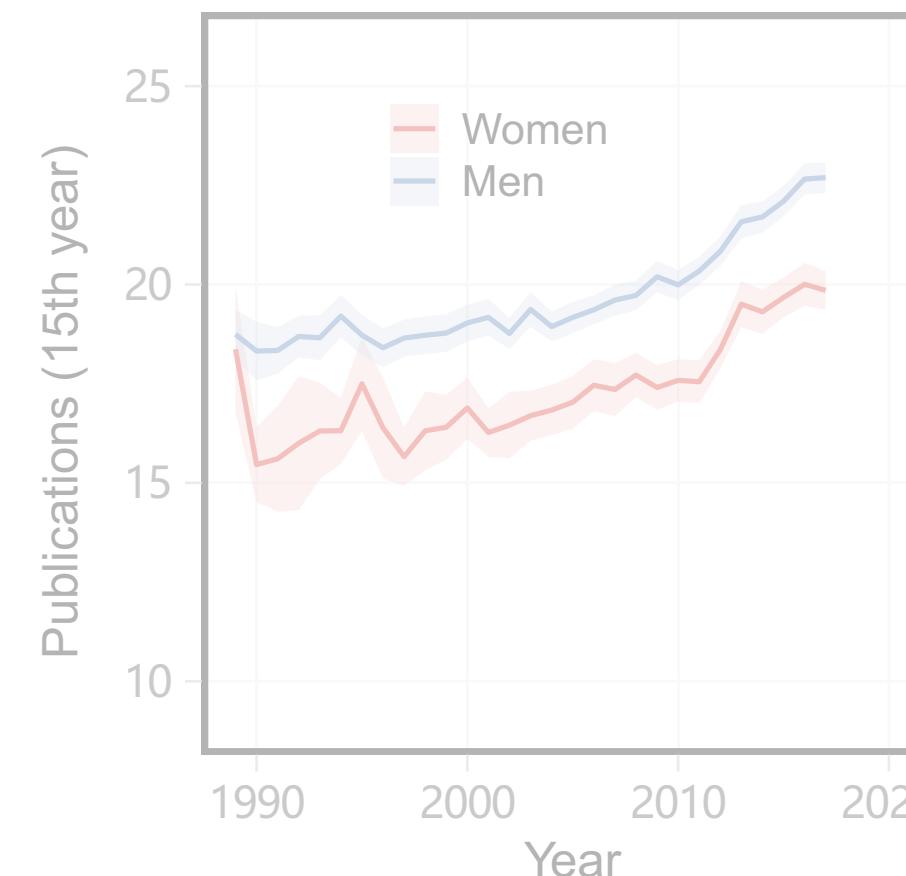
- ▶ past work : men publish more papers than women & receive more citations
  - compare  $(\lambda_i, \theta_i)$  over time, for men and women



shaded areas are 95% confidence intervals

# gender vs. productivity & prominence

- ▶ past work : men publish more papers than women & receive more citations
- compare  $(\lambda_i, \theta_i)$  over time, for men and women → their networks are different



but: latent productivity and prominence  
is *not* gendered

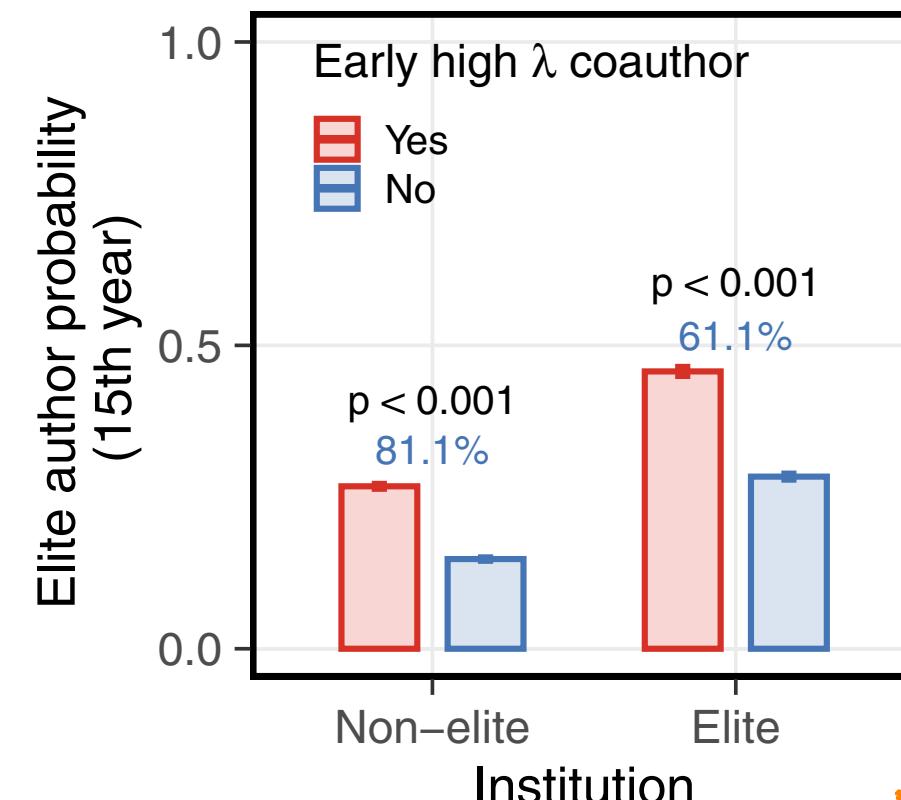
- ▶ size and composition of collaboration networks is gendered
- ▶ latent productivities increase steadily
- ▶ latent prominence stable over time

\*not causal, but implies effects of known gendered causal factors on productivity (eg, parenthood) may operate by reshaping collaborating networks

# effects of elite collaborators

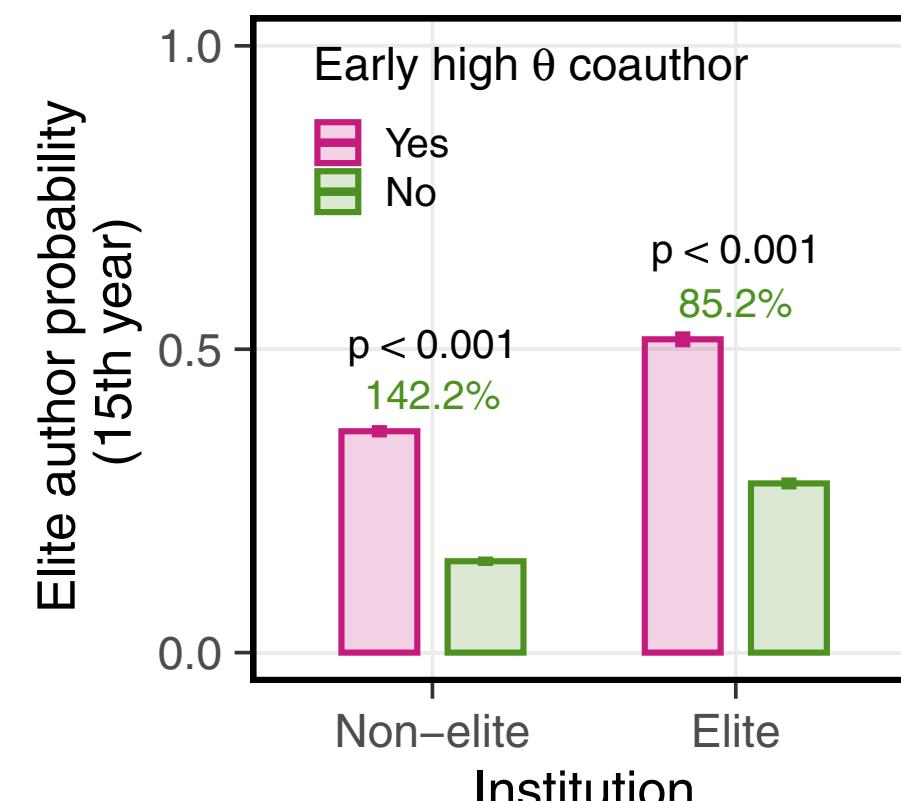
▶ how much does an early-career collaboration with an elite senior researcher influence you?

- elite senior researches with high  $\lambda$  or high  $\theta$



early high- $\lambda$  or early high- $\theta$  collaborator substantially increases likelihood of high prominence in mid-career  
▶ much more common at elite institutions  
▶ effect appears at non-elite institutions too ✓

Pr. of high- $\lambda$  early collab = 0.18 (elite) & 0.15 (non-elite)  
Pr. of high- $\theta$  early collab = 0.14 (elite) & 0.07 (non-elite)

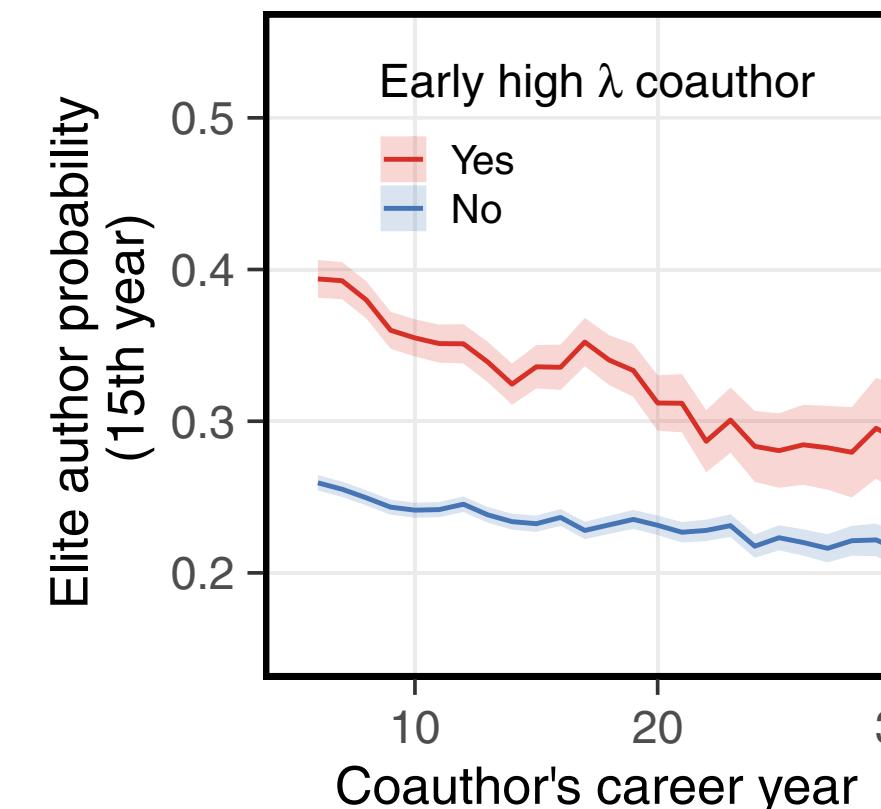
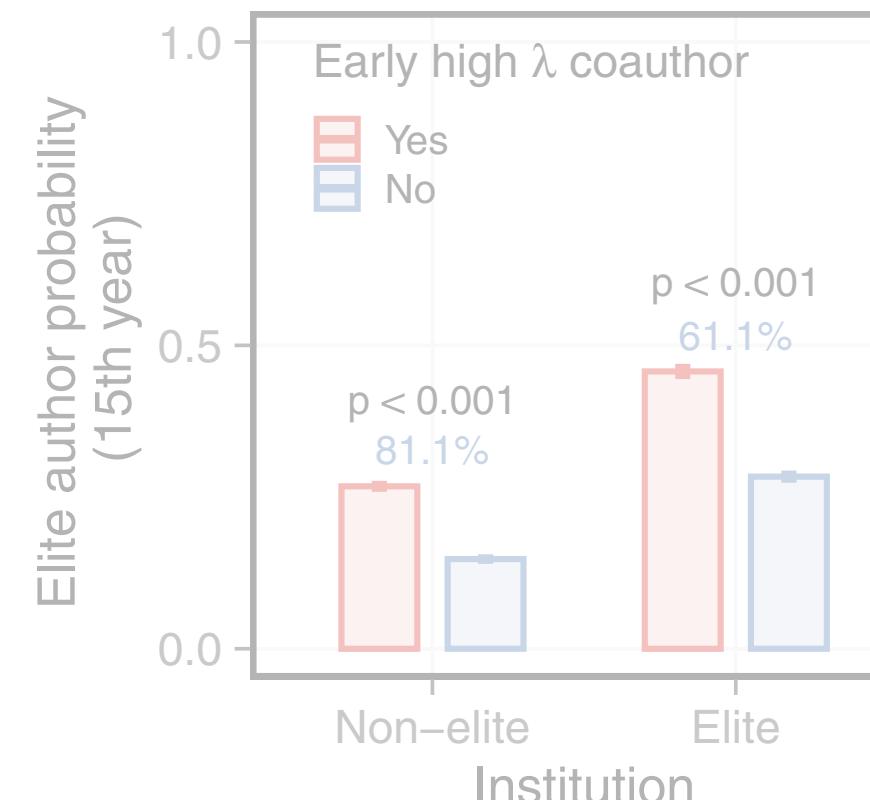


elite institutions = top 10 by z-score of high impact papers  
early-career collaboration = within first 5 years of publishing history  
"elite author" = upper 5% of citations among authors in a given field-year  
shaded areas are 95% confidence intervals

# effects of elite collaborators

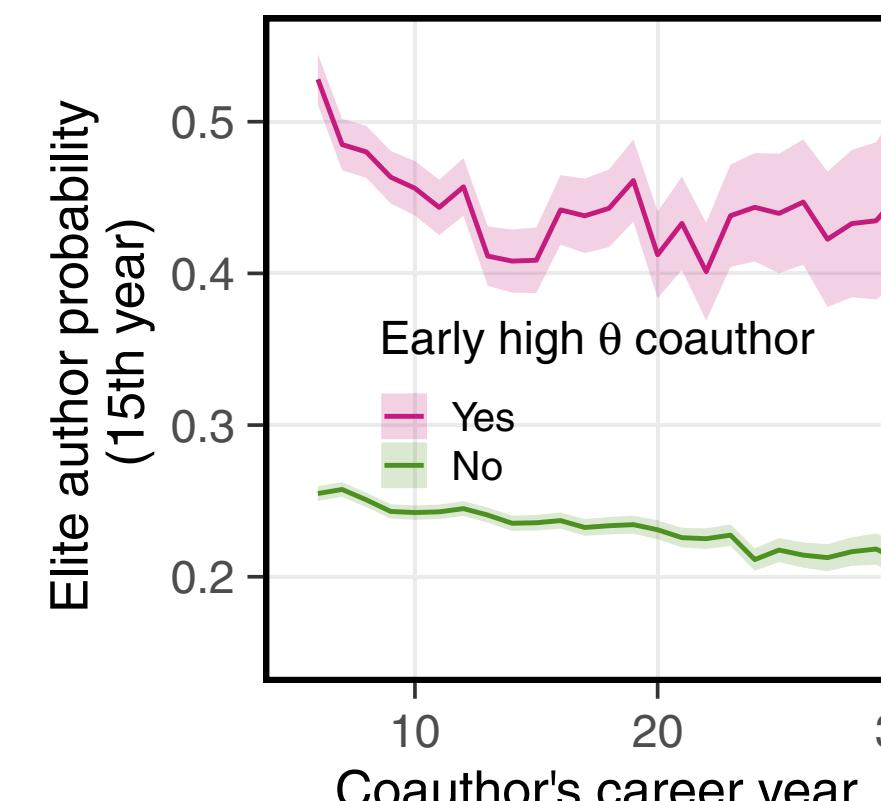
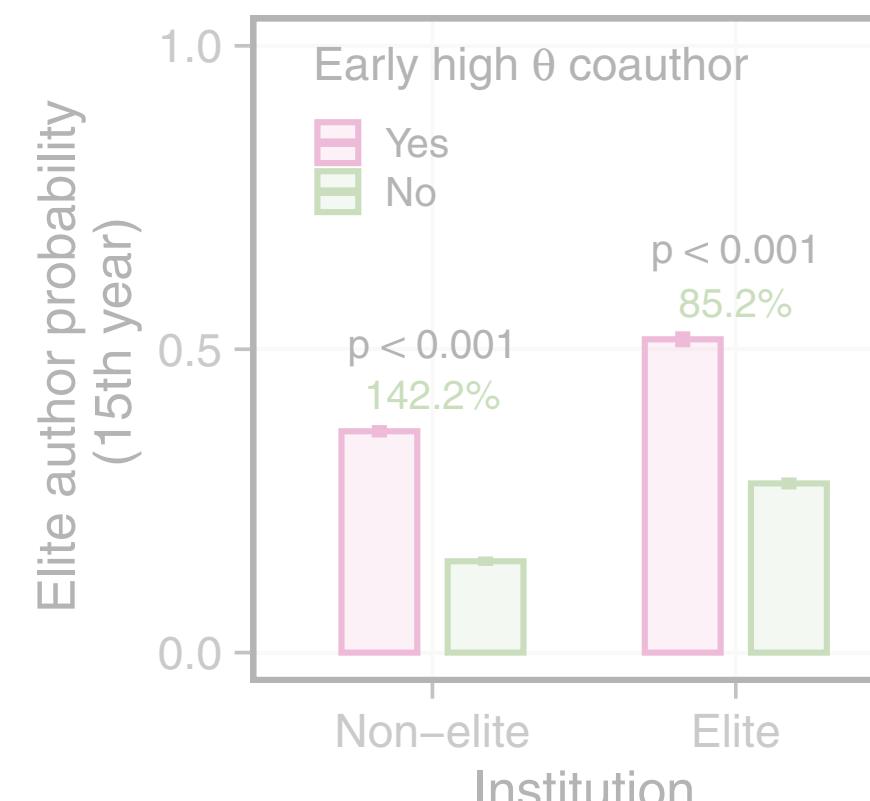
▶ how much does an early-career collaboration with an elite senior researcher influence you?

- elite senior researches with high  $\lambda$  or high  $\theta$



the 'benefits' are substantial regardless of coauthors career age

- ▶ slight decrease for most senior coauthors
- ▶ collaboration networks act like a partially transferrable form of *social capital* in science ✓

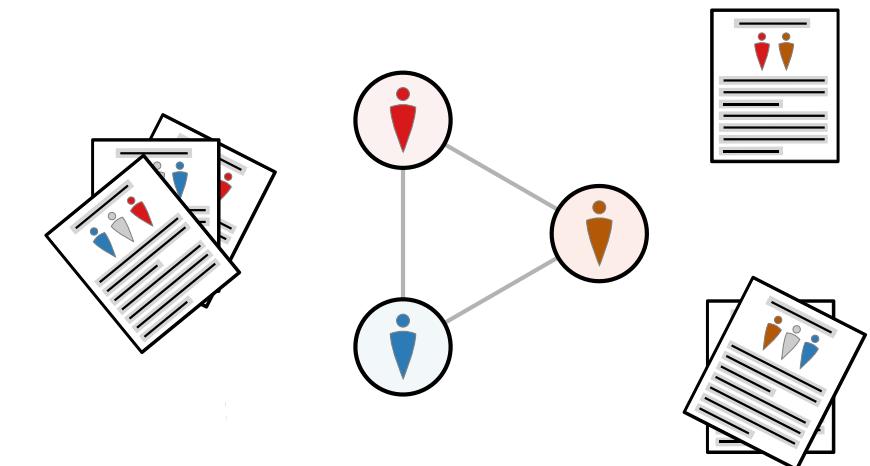


elite institutions = top 10 by z-score of high impact papers  
early-career collaboration = within first 5 years of publishing history  
"elite author" = upper 5% of citations among authors in a given field-year  
shaded areas are 95% confidence intervals

# how important is who you work with?

networks act like unequally distributed social capital in science

- *they mediate our scientific attention, evaluation, and collaboration*



differences in collaboration networks can explain

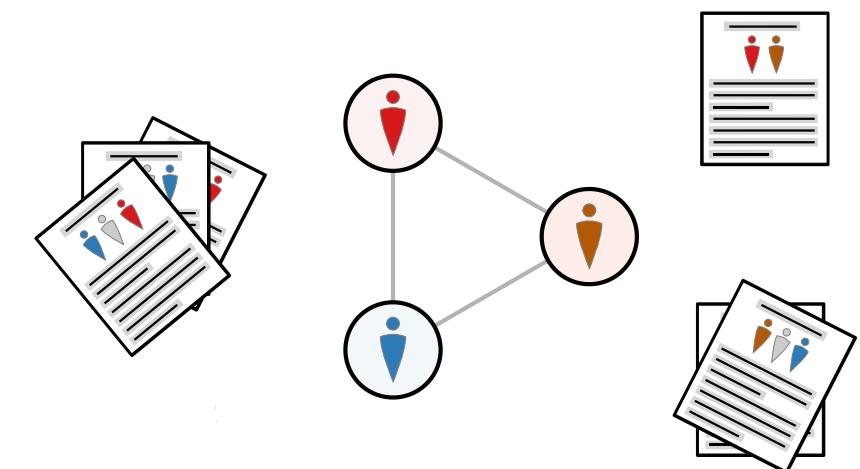
- gendered differences in productivity & prominence →
  - no difference in *individual* productivity & prominence ( $\lambda_i, \theta_i$ )
  - but men do publish more and are more prominent, implying the difference is due to their networks
  - hence, should not compare *unadjusted* measures of productivity & prominence (the network confounds)



# how important is who you work with?

networks act like unequally distributed social capital in science

- *they mediate our scientific attention, evaluation, and collaboration*



differences in collaboration networks can explain

- gendered differences in productivity & prominence
- early-career productivity & prominence

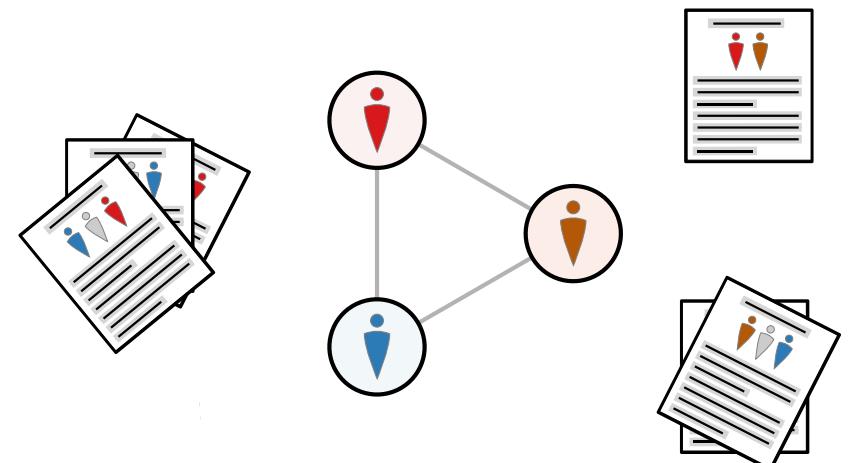
- ➔ • elite senior coauthors bequeath some of their networks to junior coauthors — inter-generational transfers
- this effect is independent of institutional prestige
- but junior scholars have greater access to elite coauthors at elite institutions



# how important is who you work with?

networks act like unequally distributed social capital in science

- *they mediate our scientific attention, evaluation, and collaboration*



differences in collaboration networks can explain

- gendered differences in productivity & prominence
- early-career productivity & prominence
- what else?

can we intervene in these networks to mitigate inequalities? 🤔

- funds for new collaborations, eg, after parenthood?
- early-career fellowships to work with elite senior coauthors? ➔

**Edge interventions can mitigate demographic and prestige disparities in the Computer Science coauthorship network**

Kate Barnes  
kathryn.barnes@colorado.edu  
University of Colorado Boulder  
Boulder, Colorado, USA

Nayera Hasan  
nhasan1@haverford.edu  
Haverford College  
Haverford, Pennsylvania, USA

Sorelle Friedler  
sorelle@cs.haverford.edu  
Haverford College  
Haverford, Pennsylvania, USA

Mia Ellis-Einhorn  
melliseinh@gmail.com  
Haverford College  
Haverford, Pennsylvania, USA

Mohammad Fanous  
mfanous@haverford.edu  
Haverford College  
Haverford, Pennsylvania, USA

Blair D. Sullivan  
sullivan@cs.utah.edu  
The University of Utah  
Salt Lake City, Utah, USA

Carolina Chávez-Ruelas  
carolina.chavezruelas@colorado.edu  
University of Colorado Boulder  
Boulder, Colorado, USA

Aaron Clauset  
aaron.clauset@colorado.edu  
University of Colorado Boulder  
Boulder, Colorado, USA

study limitations:  
• none of these analyses are causal, although they do suggest specific mechanisms that can be tested  
• no data on race/ethnicity in our bibliographic analyses, although literature suggests under-represented minorities may have similar or larger network differences as women  
• we focus only on STEM fields, which tend to have strong collaboration norms; unclear what results might be for fields with different collaboration norms

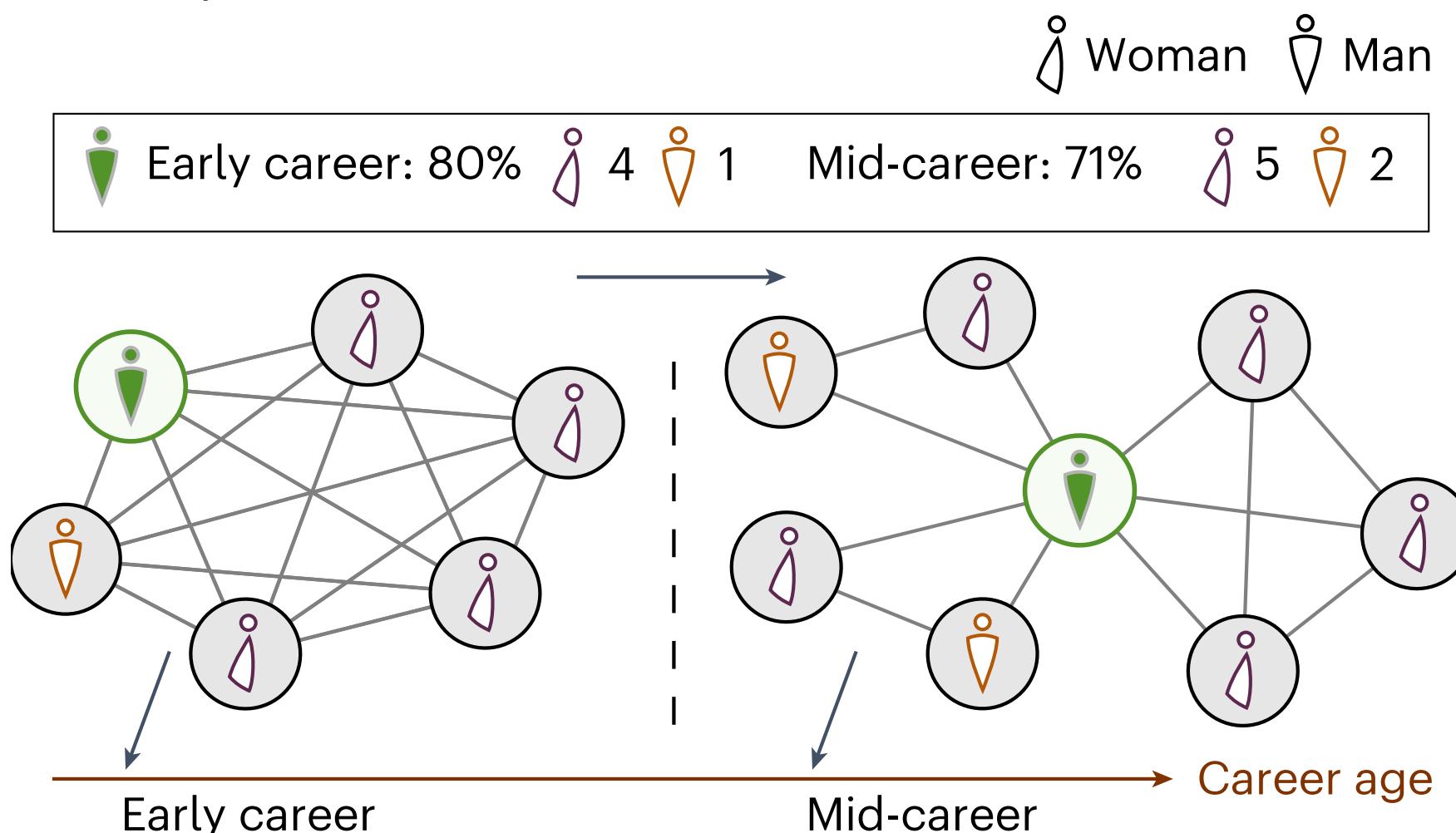
Preprint, arxiv:2506.04435 (2025)

# gender and racial diversity socialization

▶ how does the *diversity* of early-career teams influence later-career collaboration networks?

homophily : women prefer women, men prefer men

*what if preferences can be socialized?*



man who trains with high density of women → man who advises high density of early-career women

Article

<https://doi.org/10.1038/s43588-025-0122-1>

## Gender and racial diversity socialization in science

Received: 4 December 2024

Weihua Li <sup>1,2,3,4,5,6</sup>, Hongwei Zheng <sup>7</sup>✉, Jennie E. Brand<sup>8</sup> &

Accepted: 20 March 2025

Aaron Clauset <sup>9,10,11</sup>✉

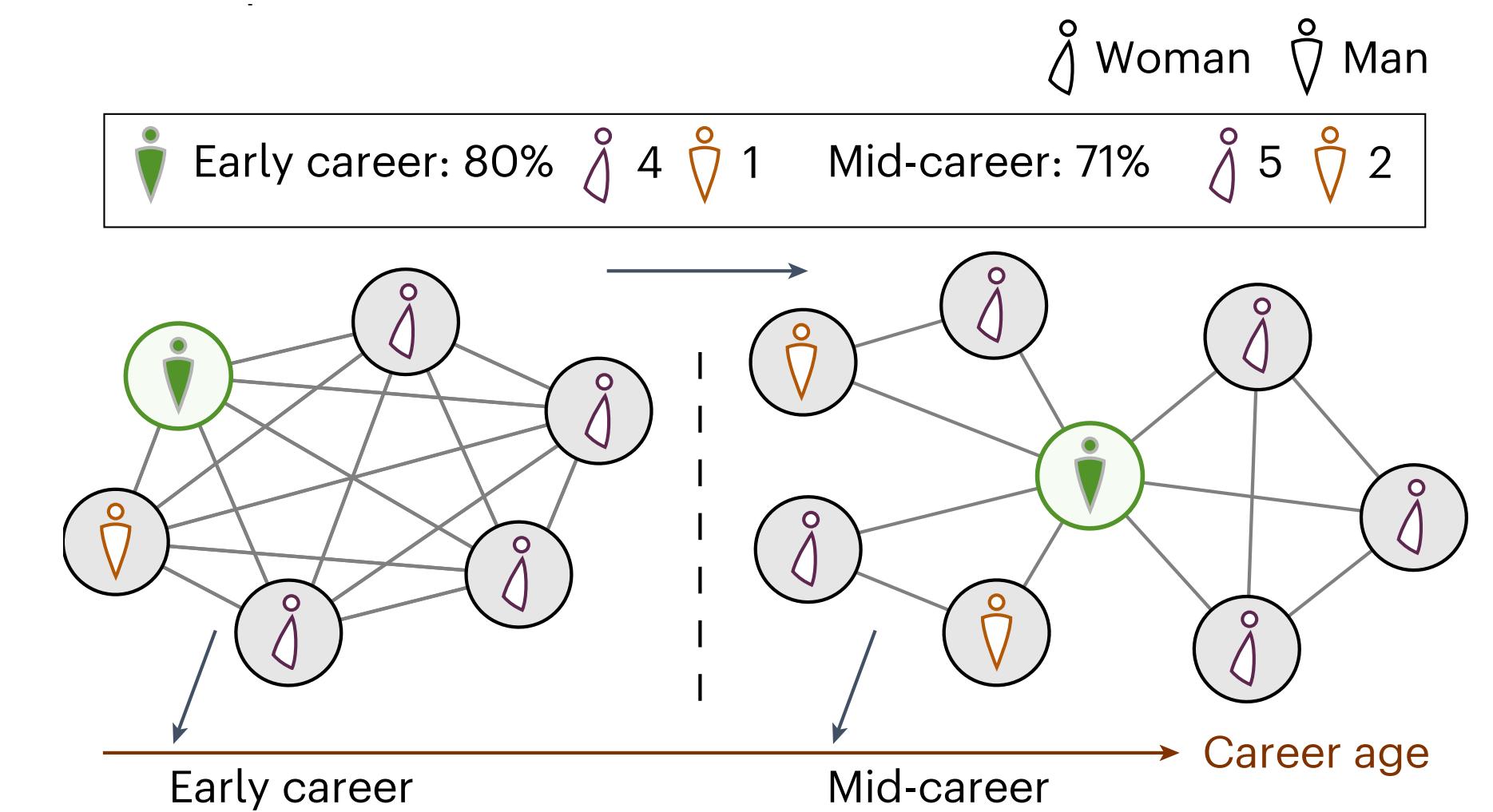
Nature Computational Science (2025)

# gender and racial diversity socialization

▶ how does the *diversity* of early-career teams influence later-career collaboration networks?

homophily : women prefer women, men prefer men

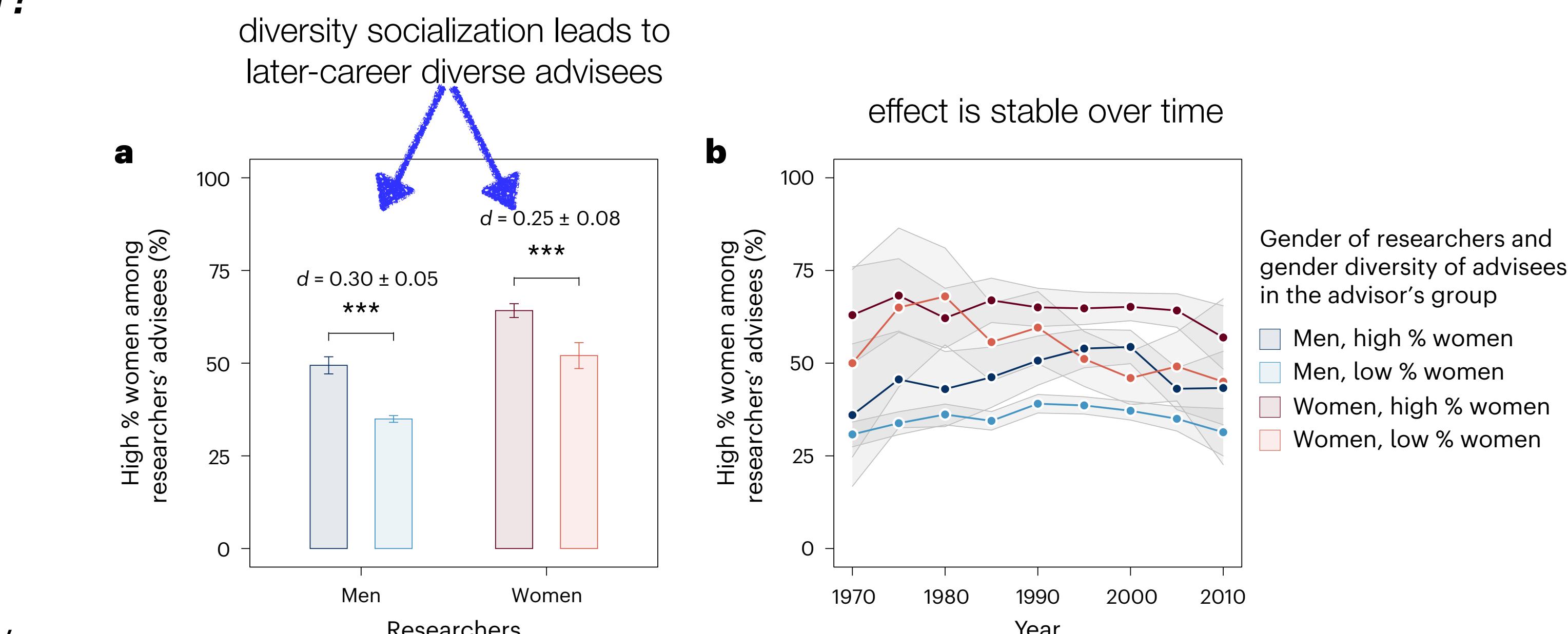
*what if preferences can be socialized?*



man who trains with high density of women



man who advises high density of early-career women



homophily is not destiny – diversity is a learnable preference

Article

<https://doi.org/10.1038/s43588-025-0122-1>

## Gender and racial diversity socialization in science

Received: 4 December 2024

Weihua Li <sup>1,2,3,4,5,6</sup>, Hongwei Zheng <sup>7</sup>✉, Jennie E. Brand <sup>8</sup> &

Accepted: 20 March 2025

Aaron Clauset <sup>9,10,11</sup>✉

Nature Computational Science (2025)

# the scientific ecosystem

no scientist is an island → productivity is driven by environmental & network effects

- where you train *doesn't* seem to matter (conditioned on getting a faculty job)

- where you work *does* matter : prestige → available labor → more scientific output



- prestige shapes your working environment – like a slope underneath the whole ecosystem
- no results on the *value* of different scientific outputs, only their total volume
- many unobserved confounds: quality of job talk, fundability of ideas, quality of writing, etc.

# the scientific ecosystem

no scientist is an island → productivity is driven by environmental & network effects

- where you train *doesn't* seem to matter (conditioned on getting a faculty job)
- where you work *does* matter : prestige → available labor → more scientific output



productivity & prominence are network effects → who you work with shapes your output

- net of their coauthors, men & women *equal* in their *individual* productivity & prominence
- early-career researchers can *inherit* collaboration networks → like generational wealth transfer

- prestige shapes who you *can* work with (environment)
- no results about scientific *value*, only the volume and the attention outputs receive
- comparing individual researchers is unfair without accounting for different coauthor network effects

# the scientific ecosystem

no scientist is an island → productivity is driven by environmental & network effects

- where you train *doesn't* seem to matter (conditioned on getting a faculty job)
- where you work *does* matter : prestige → available labor → more scientific output



productivity & prominence are network effects → who you work with shapes your output

- net of their coauthors, men & women *equal* in their *individual* productivity & prominence
- early-career researchers can *inherit* collaboration networks → like generational wealth transfer

ecosystem metaphor is rich → what other environmental or network effects?

- how should we intervene in an ecosystem?
- what might accelerate scientific discovery?

- The undervaluing of elite women in physics
- Gender and racial diversity socialization in science
- Epistemic inequality in the diffusion of scientific ideas

# references & collaborators

## Productivity, prominence, and the effects of academic environment

Samuel F. Way<sup>a,1</sup>, Allison C. Morgan<sup>a</sup>, Daniel B. Larremore<sup>a,b,2</sup>, and Aaron Clauset<sup>a,b,c,1,2</sup>

<sup>a</sup>Department of Computer Science, University of Colorado, Boulder, CO, USA; <sup>b</sup>BioFrontiers Institute, University of Colorado, Boulder, CO, USA; <sup>c</sup>Santa Fe Institute, Santa Fe, NM, USA

PNAS 116(22), 10729–10733 (2019)

## Quantifying hierarchy and dynamics in US faculty hiring and retention

<https://doi.org/10.1038/s41586-022-05222-x> K. Hunter Wapman<sup>1</sup>✉, Sam Zhang<sup>2</sup>, Aaron Clauset<sup>1,3,4</sup> & Daniel B. Larremore<sup>1,3</sup>✉

Nature 610, 120–127 (2022)

## Labor advantages drive the greater productivity of faculty at elite universities

Sam Zhang<sup>1\*</sup>, K. Hunter Wapman<sup>2</sup>, Daniel B. Larremore<sup>2,3</sup>, Aaron Clauset<sup>2,3,4\*</sup>

Science Advances 8, eabq7056 (2022)

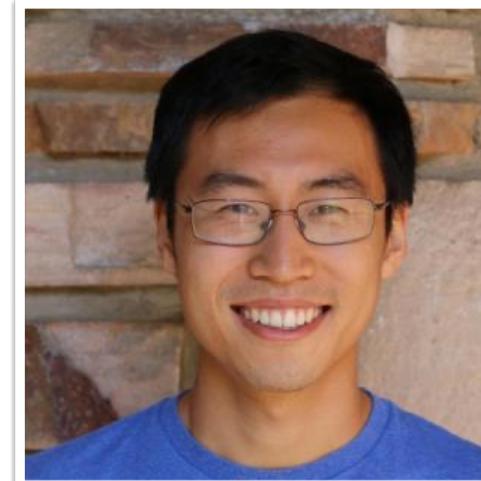
## Untangling the network effects of productivity and prominence among scientists

Weihua Li<sup>1,2,3,4</sup>✉, Sam Zhang<sup>1</sup>✉, Zhiming Zheng<sup>1,2,3,4</sup>, Skyler J. Cranmer<sup>6</sup> & Aaron Clauset<sup>1,7,8,9</sup>

Nature Communications 13, 4907 (2022)



Dr. Samuel F Way  
(now: Spotify)



Sam Zhang  
(Colorado)



Dr. K. Hunter Wapman  
(Colorado)



Prof. Weihua Li  
(Beihang)



Prof. Zhiming Zhang  
(Beihang)



Dr. Allison Morgan  
(now: Code for America)



Prof. Skyler Cranmer  
(Ohio State)



Prof. Daniel Larremore  
(Colorado)

Funding:

