

# Large-scale structures in networks: Hidden communities and latent hierarchies

Daniel Larremore

Assistant Professor

Dept. of Computer Science  
& BioFrontiers Institute

daniel.larremore@colorado.edu  
@danlarremore



University of Colorado **Boulder**

# Note in the online version of slides:

Here you'll find lots of information, but of course, there's so much great work out there that is *not* covered in these slides.

If you see mistakes in what follows, please write! [daniel.larremore@colorado.edu](mailto:daniel.larremore@colorado.edu)

**PDF** of slides available → <http://LarremoreLab.github.io>

## **Goals** for this talk:

1. **Why** do we look for large-scale structure? 🤔
2. **How** do we find communities and hierarchies? 🙄
3. **Where** can we read more details? 📚

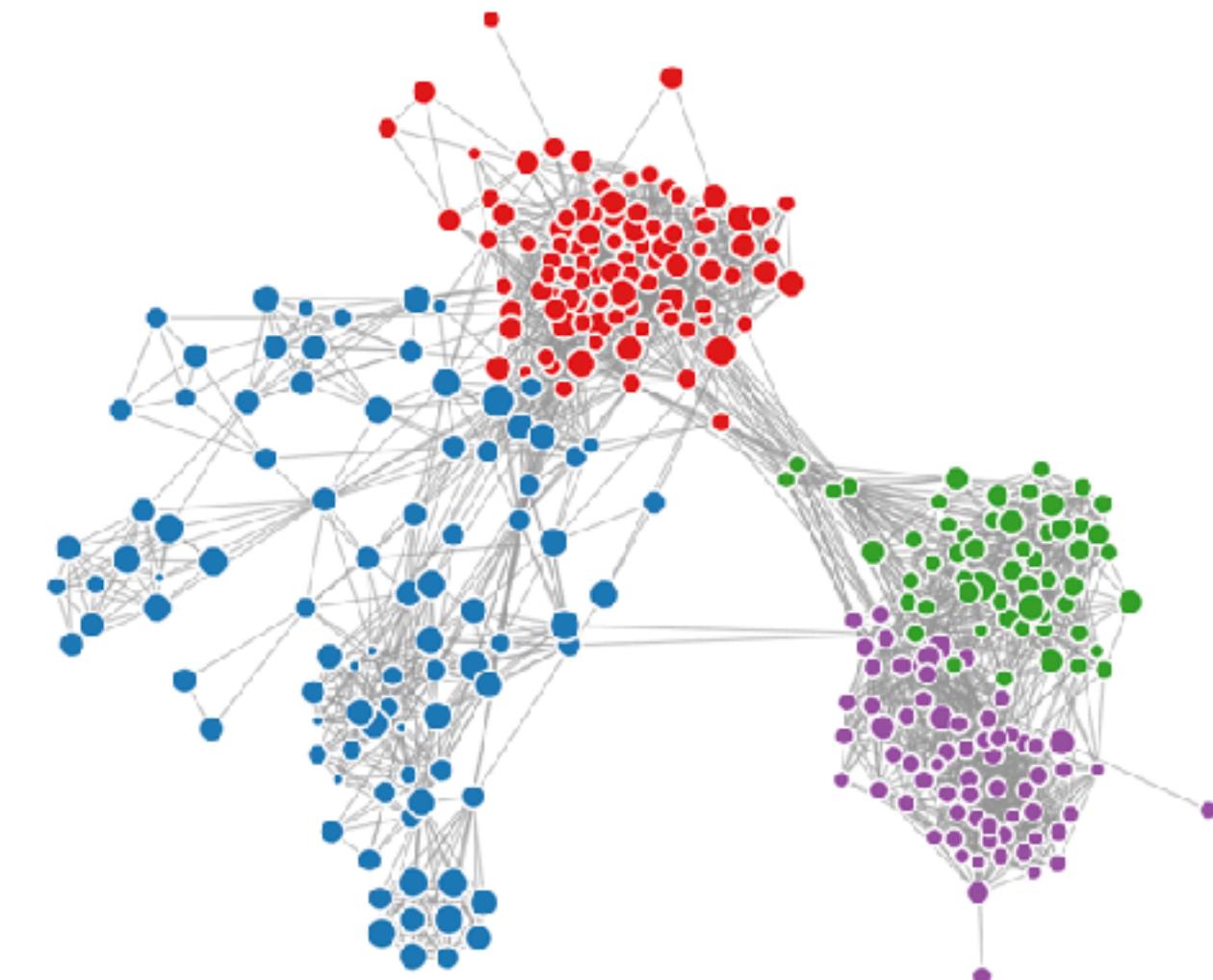
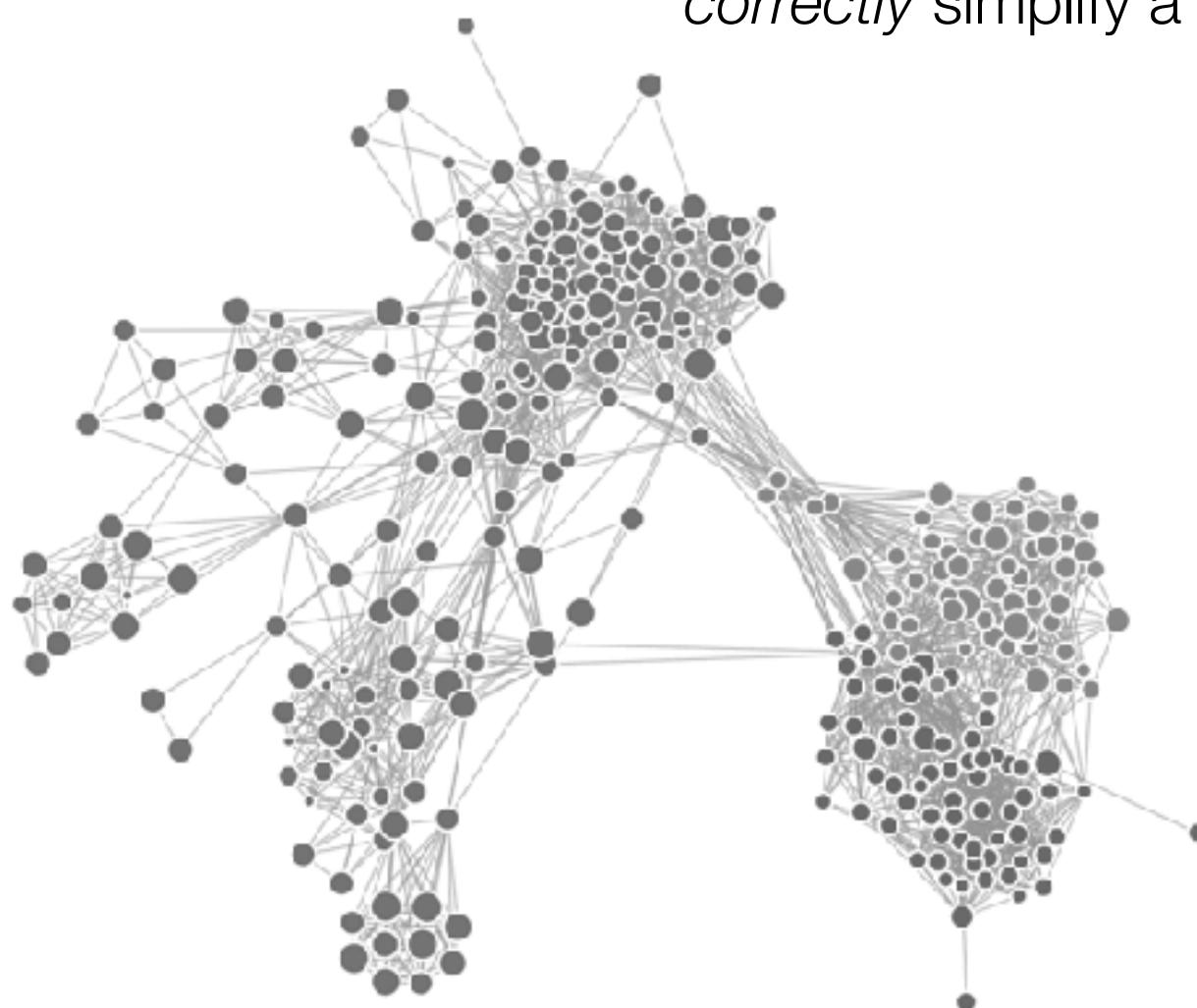
Simplicity is a great virtue but it requires hard work to achieve it and education to appreciate it. And to make matters worse: complexity sells better.

E. W. Dijkstra

We can interpret this in two ways:

**The Cynic:** Pictures of networks can be *really cool* but our goal is to do good science, not make pretty pictures.

**The Scientist:** The most beautiful science is when we *correctly* simplify a complex system.



# What do we mean by “large-scale structure” ?

**Structure is what makes data different from noise.**  
It's what makes a network different from a random graph.

Networks are often too large and complex to be adequately summarized by a few scalars, like the number of nodes, the number of edges, or the mean degree.

However, they are also often too large and complex to be analyzed *without* some kind of simplification!

Therefore, understanding what the network means requires that we identify key structures.

Searching for large-scale structures in a network reflects a belief that in all the complexity there are patterns that make the network less complicated.

We define these large-scale structures—models, really—to compress complex networks.

# Goal: understanding, not a list of parts and dimensions



Finding large-scale structures  
is the same as anything else:

We want a simplified model of  
something very complicated.

We want to know what the  
important pieces are,  
and how they fit together.

# Many uses for models of large-scale structure

## Treat the network like a system:

**Extrapolation.** Make predictions for as-yet unseen nodes (in “space” or time).

**Interpolation.** Identify missing links.

**Generalization.** Nodes of this type are like others of the same type.

## Treat the network like an artifact:

**Mechanisms.** How did this network arise? What rules governed its assembly?

**Explanations.** Coarse-graining or compression.

## Treat the network like a means to an end; an intermediate data structure:

**Useful division.** Need groups so that we can assign treatments in an A/B test.

**Simplification.** Downstream regression model needs ranks or groups.

intuition: compare this list with the list you would write for regression

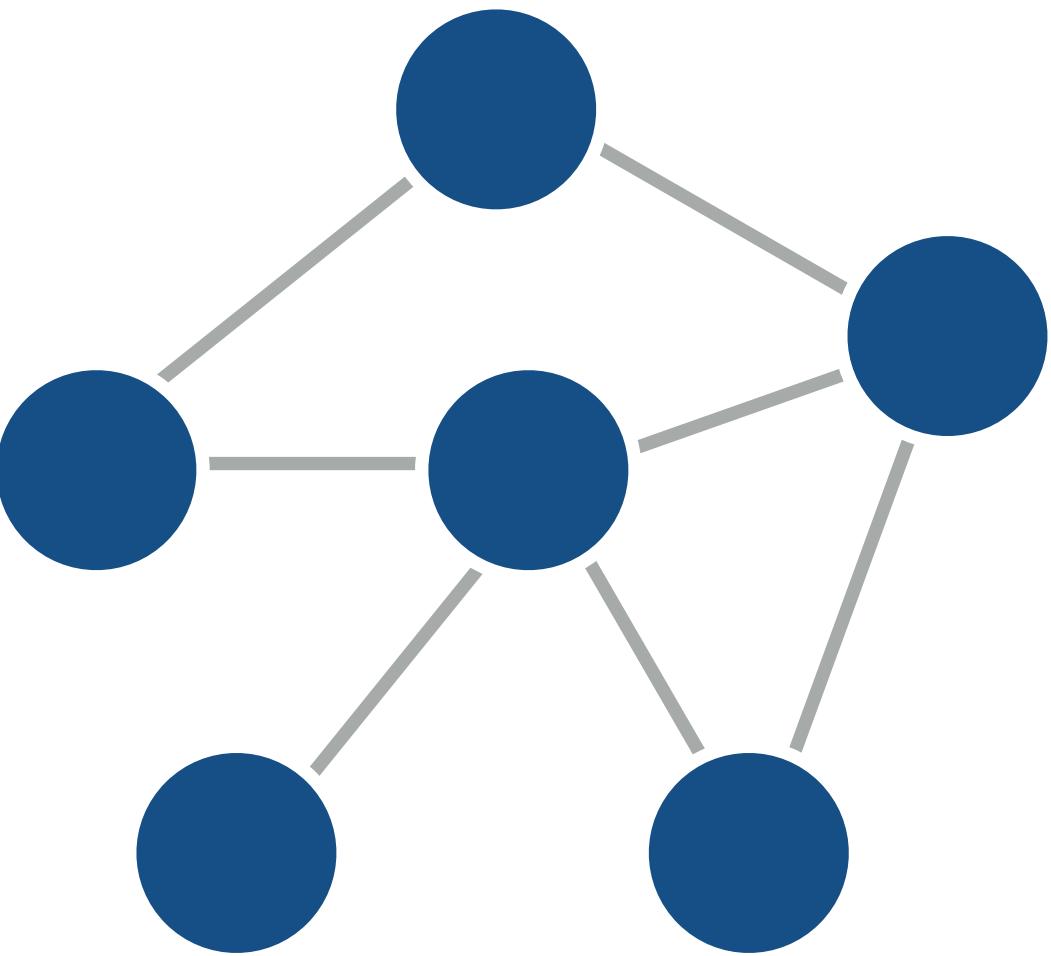
# Community structure



# Homophily & assortative mixing

*like* links with *like*

Assortativity coefficient  $r$  measures extent of homophily.

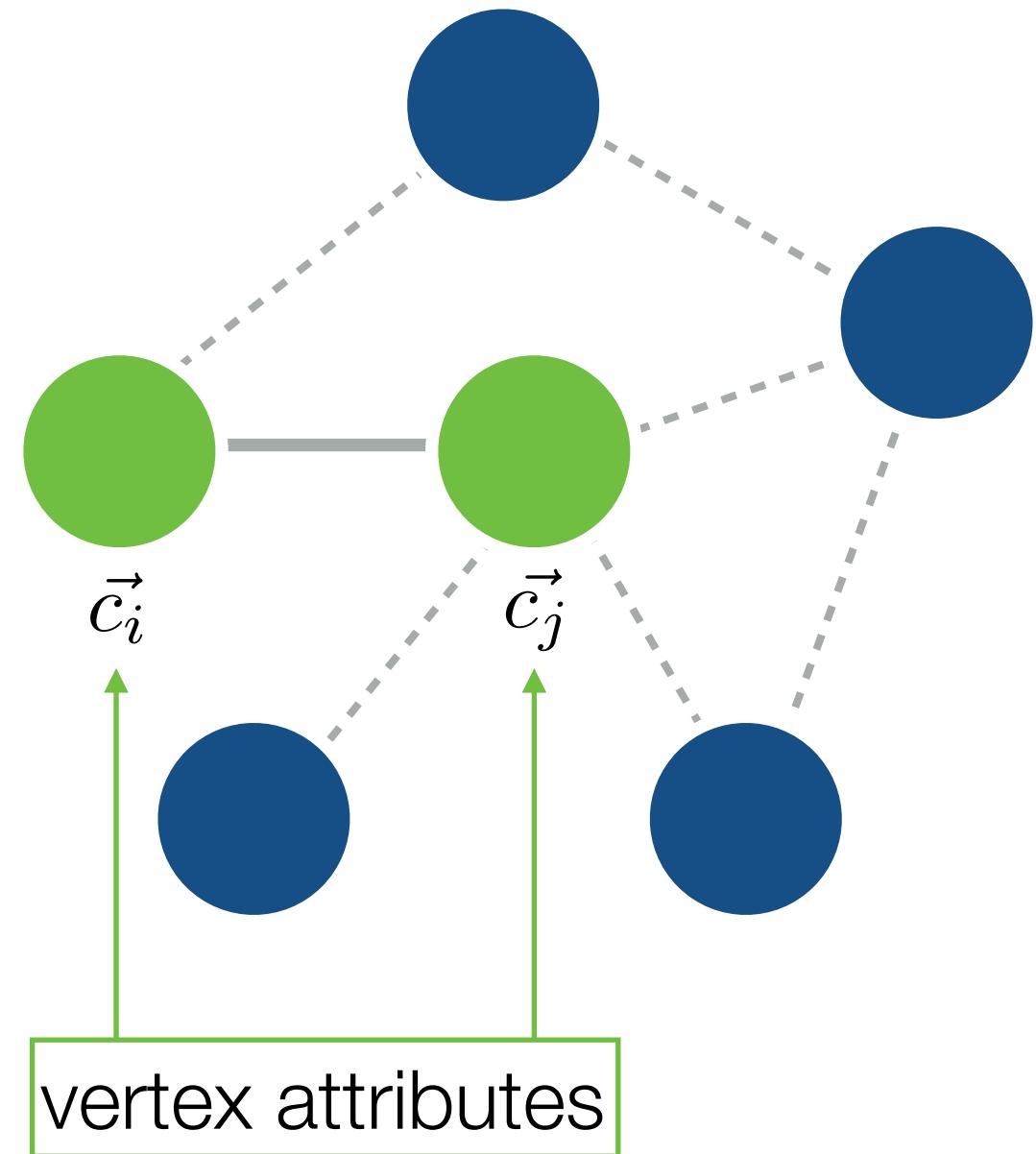


# Homophily & assortative mixing

*like* links with *like*

Assortativity coefficient  $r$  measures extent of homophily.

Three types:  
scalar attributes  
vertex degrees  
categorical variables



# Homophily & assortative mixing

*like* links with *like*

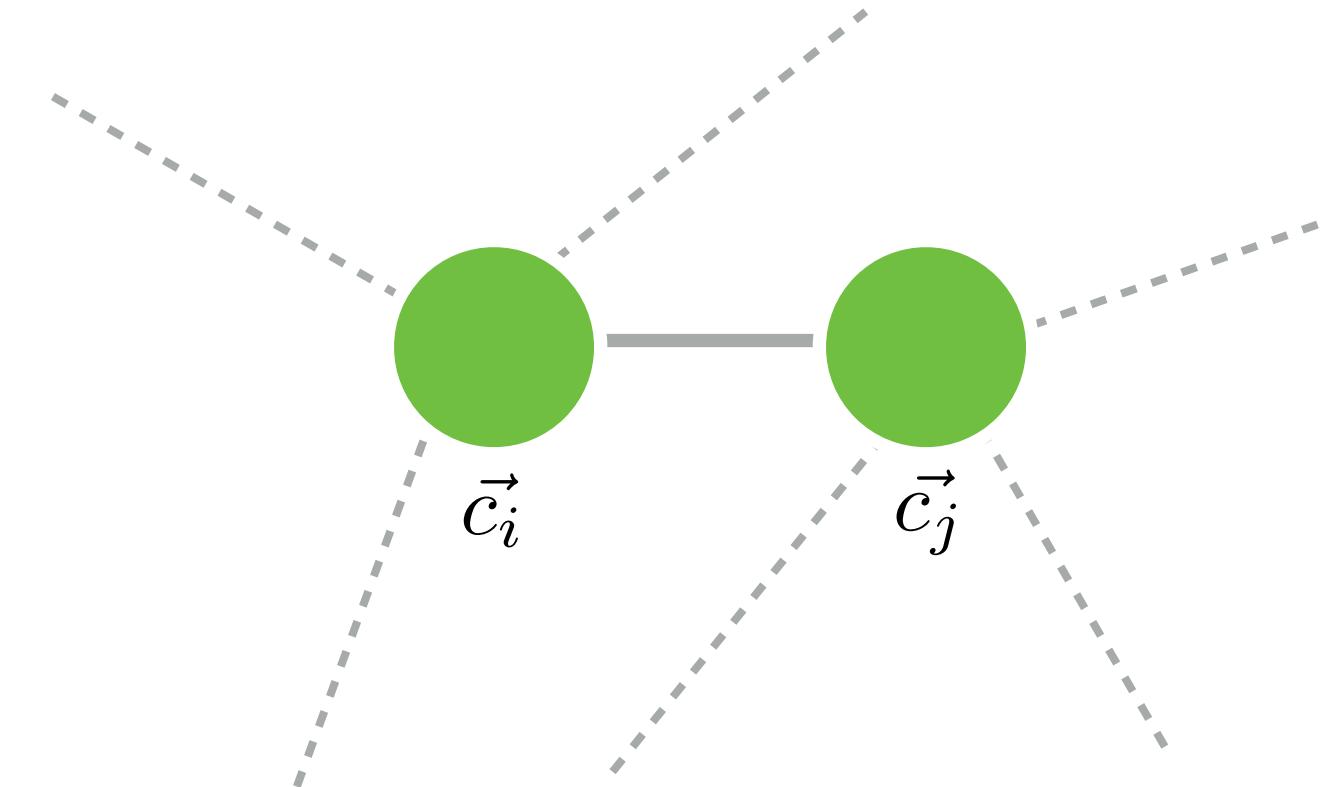
Assortativity coefficient  $r$  measures extent of homophily.

scalar attributes:

compute a Pearson correlation over edges.

start with the mean of  $c$  across ties:

$$\mu = \frac{1}{2m} \sum_i \sum_j A_{ij} c_i = \frac{1}{2m} \sum_i k_i c_i$$



# Homophily & assortative mixing

*like* links with *like*

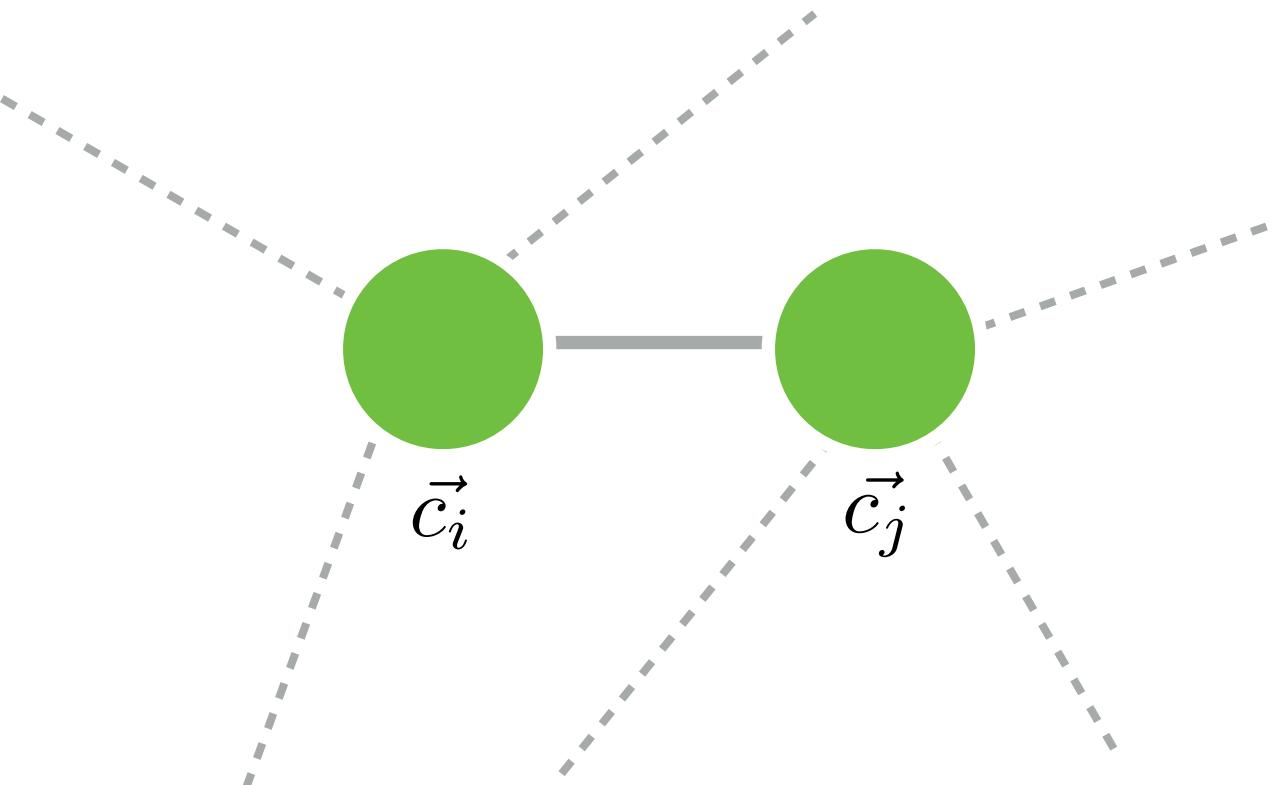
Assortativity coefficient  $r$  measures extent of homophily.

compute covariance:

$$\text{cov}(c_i, c_j) = \frac{\sum_{ij} A_{ij}(c_i - \mu)(c_j - \mu)}{\sum_{ij} A_{ij}}$$

$$= \frac{1}{2m} \sum_{ij} A_{ij} c_i c_j - \mu^2$$

$$= \frac{1}{2m} \sum_{ij} \left( A_{ij} - \frac{k_i k_j}{2m} \right) c_i c_j$$



substitute  $\mu = \frac{1}{2m} \sum_i k_i c_i$

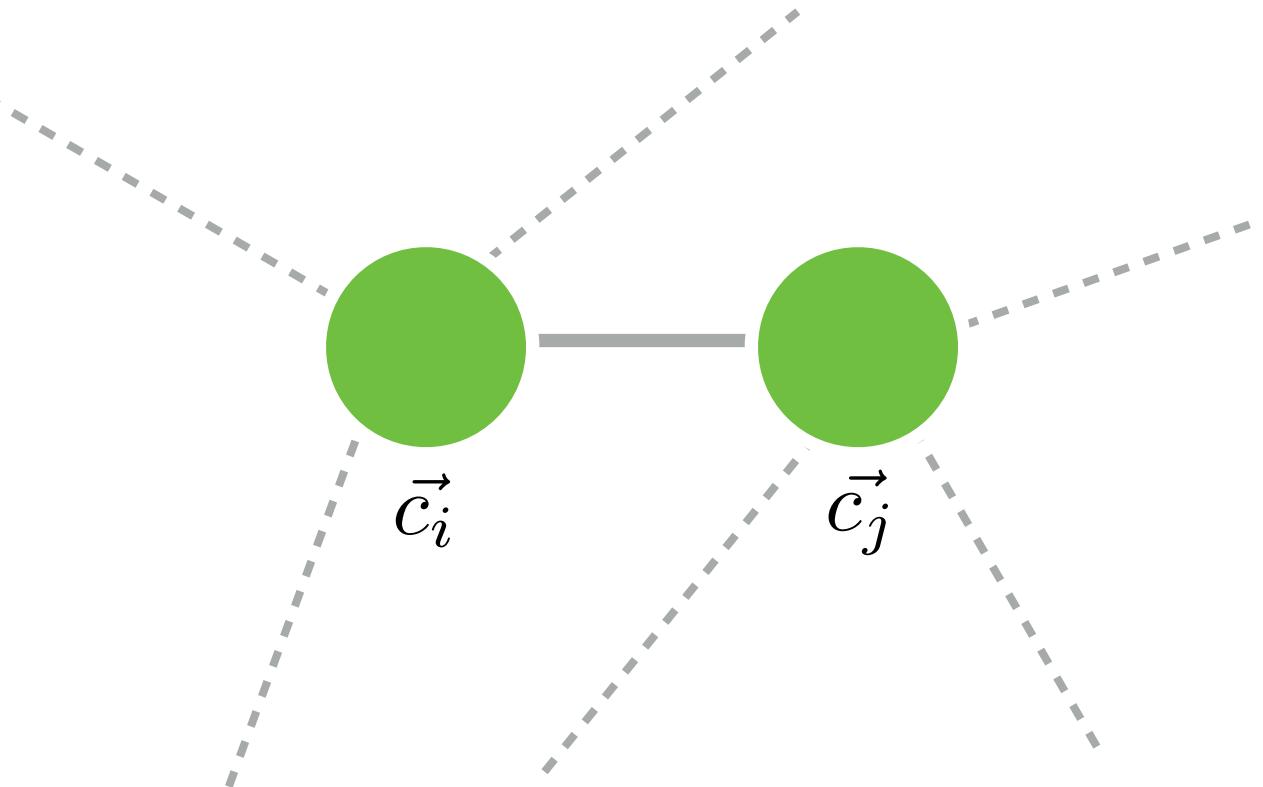
# Homophily & assortative mixing

*like* links with *like*

Assortativity coefficient (scalar).

$$r = \frac{\text{cov}(c_i, c_j)}{\text{var}(c_i, c_j)}$$

$$r = \frac{\sum_{ij} (A_{ij} - k_i k_j / 2m) c_i c_j}{\sum_{ij} k_i \delta_{ij} - k_i k_j / 2m}$$



“it’s that easy!” 😎

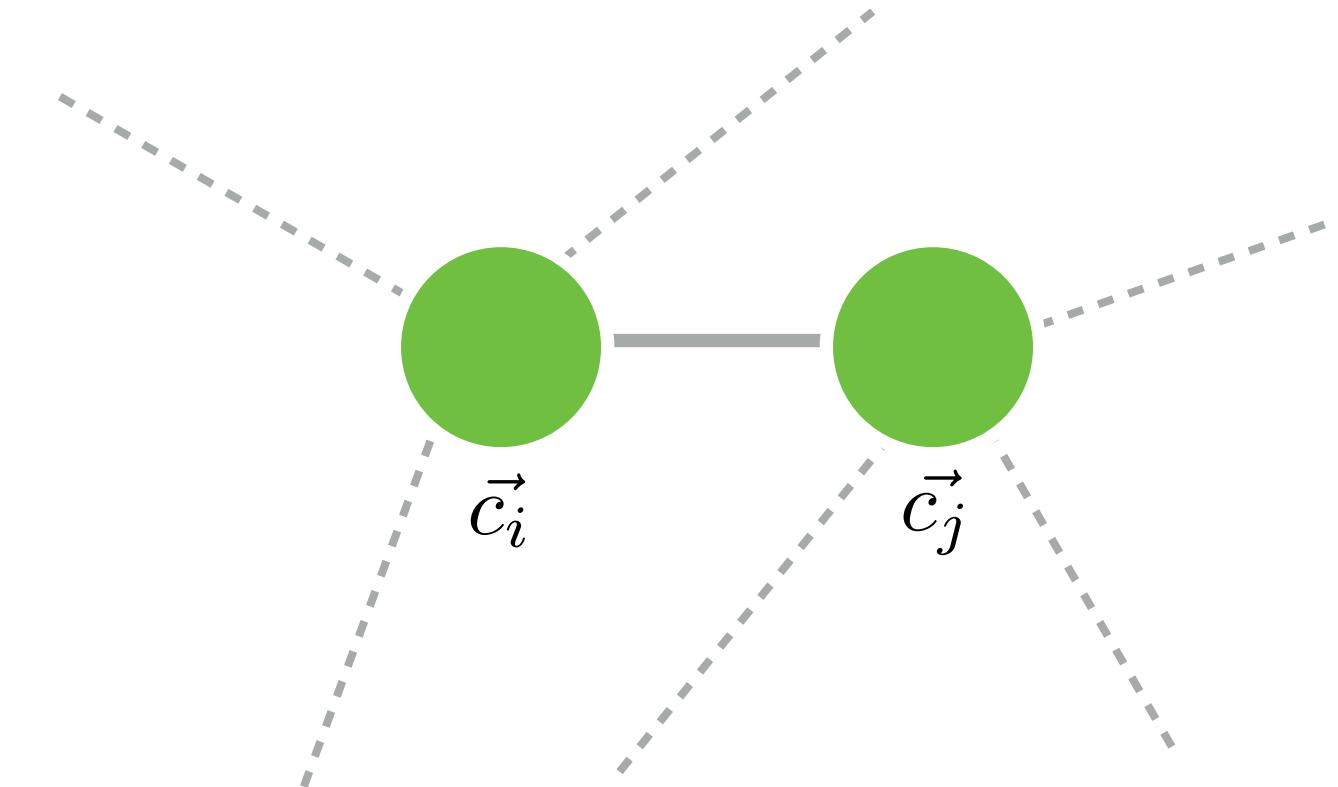
$$-1 \leq r \leq 1$$

# Homophily & assortative mixing

*like* links with *like*

degree: just another scalar.

(very well studied!)



# Homophily & assortative mixing

*like* links with *like*

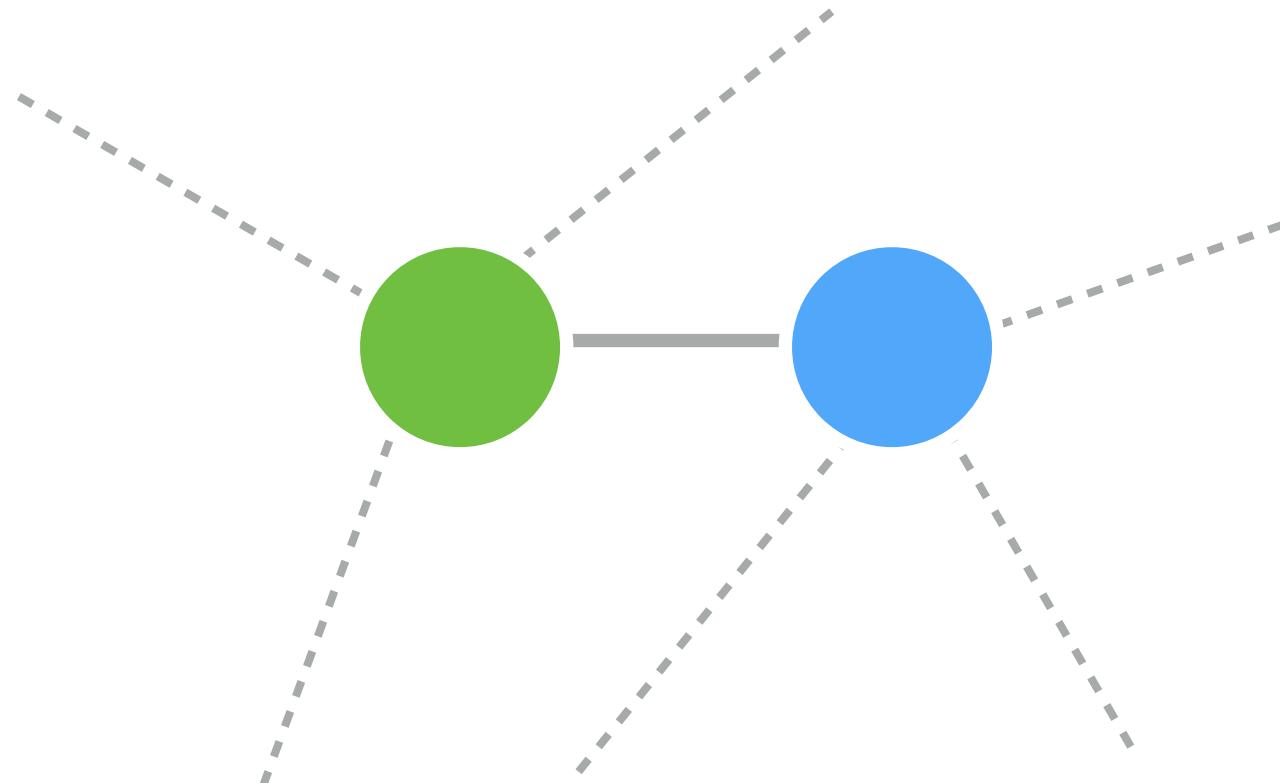
categories: not just another scalar.

Imagine that  $c$  are labels, categories.

let  $e_{rs}$  be the fraction of edges  
between nodes of type  $r$  and  $s$ .

$$\sum_{rs} e_{rs} = 1 \quad \sum_r e_{rs} = a_r \quad \sum_s e_{rs} = b_s$$

$$\text{then, } r = \frac{\sum_r e_{rr} - \sum_r a_r b_r}{1 - \sum_r a_r b_r}$$



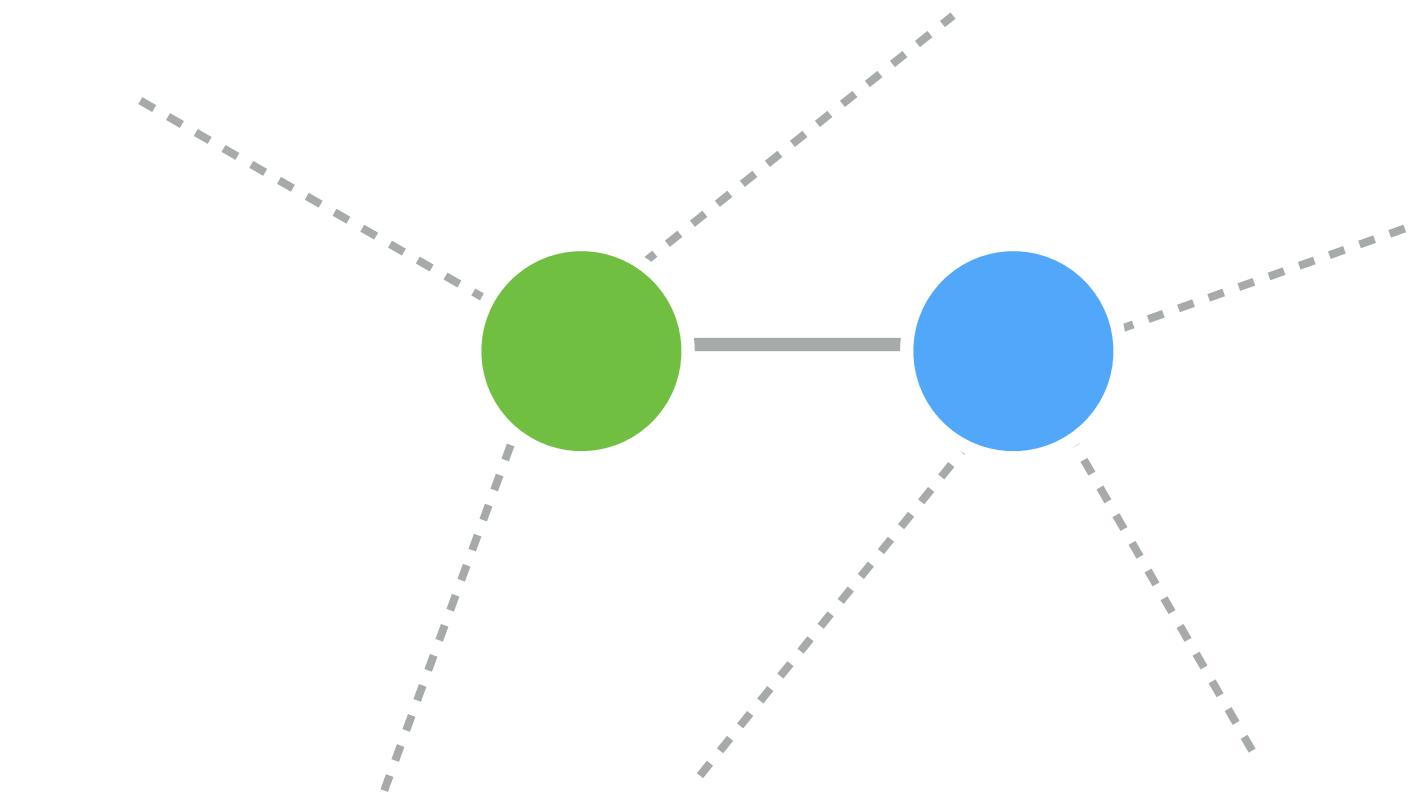
consider: if there were *only* edges  
between nodes of the same type...

# Homophily & assortative mixing

*like* links with *like*

We write the correlation of categories across edges this way, and call it  $Q$ .

Principle: How many more edges are there between nodes with the same label, than we'd expect at random?



$$\sum_{ij} \left( \begin{array}{c} \text{\# actual edges} \\ i \leftrightarrow j \end{array} - \begin{array}{c} \text{\# edges if there were} \\ \text{no correlations} \\ i \leftrightarrow j \end{array} \right) \text{only if } i, j \text{ have the same label}$$

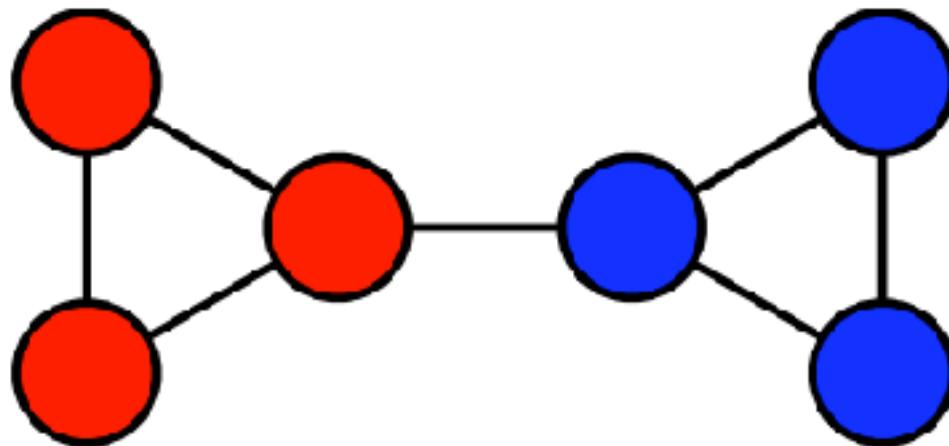
$$Q = \frac{1}{2m} \sum_{ij} \left( A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{b_i, b_j}$$

# Practice makes the master

$$Q = \frac{1}{2m} \sum_{ij} \left( A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{b_i, b_j}$$

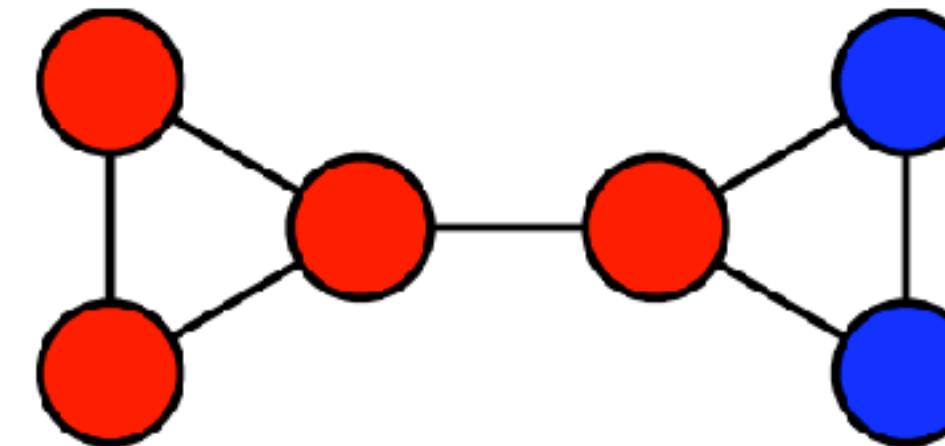
$$Q = \sum_r e_{rr} - a_r^2$$

Equivalent form:  
←  $e_{rs}$  is the fraction of edges connecting labels  $r$  and  $s$



labeling 1		red	blue
red	3/7	1/14	
blue	1/14	3/7	

$$Q_1 = 5/14 = 0.357$$



labeling 2		red	blue
red	4/7	2/14	
blue	2/14	1/7	

$$Q = 6/49 = 0.122$$

# Modularity

Modularity is easily *the* most popular method for community detection. But why?

$$Q = \frac{1}{2m} \sum_{ij} \left( A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{b_i, b_j}$$

Why is this more powerful than simply a measure of correlation over node labels?

## Community structure in social and biological networks

[M Girvan, MEJ Newman - Proceedings of the national ...](#), 2002 - National Acad Sciences

A number of recent studies have focused on the statistical properties of networked systems such as social networks and the Worldwide Web. Researchers have concentrated particularly on a few properties that seem to be common to many networks: the small-world property, power-law degree distributions, and network transitivity. In this article, we highlight another property that is found in many networks, the property of community structure, in which network nodes are joined together in tightly knit groups, between which there are only ...

☆ 99 Cited by 12341  Related articles All 40 versions

## Finding and evaluating community structure in networks

[MEJ Newman, M Girvan - Physical review E, 2004 - APS](#)

We propose and study a set of algorithms for discovering community structure in networks—natural divisions of network nodes into densely connected subgroups. Our algorithms all share two definitive features: first, they involve iterative removal of edges from the network to split it into communities, the edges removed being identified using any one of a number of possible “betweenness” measures, and second, these measures are, crucially, recalculated after each removal. We also propose a measure for the strength of the community structure ...

☆ 99 Cited by 11238  Related articles All 38 versions

# Key: let's reverse our thinking of what Q does

Don't use Q to compute correlation of some given labels.

Instead, **experiment with the labels** and see how you can **maximize Q!**

Now, we have a computer science problem:  
**how do you search the space of partitions?**

(This space is really big!)

How would you do it? 🤔

# People like modularity. Why?

$$Q = \frac{1}{2m} \sum_{ij} \left( A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{b_i, b_j}$$

- Intuitive
- Works for weighted and unweighted networks.
- Corresponds to our social network ideas of what (cohesive) communities are.
  - Automatically choose  $k$ , the number of groups.
  - Rapid approximate solutions.
  - Follows the usual methods trajectory: idea, demonstration, optimization.
- Fun customizations:
  - Resolution parameter to “zoom in” and “zoom out.”
  - Find the clusters. Then cluster the clusters. Then cluster those clusters...
  - Directed. Bipartite.

$$Q = \frac{1}{m} \sum_{ij} \left( A_{ij} - \frac{k_i^{\text{out}} k_j^{\text{in}}}{m} \right) \delta_{b_i, b_j}$$

modularity for directed networks

$$Q = \frac{1}{2m} \sum_{ij} \left( A_{ij} - \gamma \frac{k_i k_j}{2m} \right) \delta_{b_i, b_j}$$

modularity with a resolution parameter

# Why aren't we done here?

Physicists like to minimize things because rocks fall.

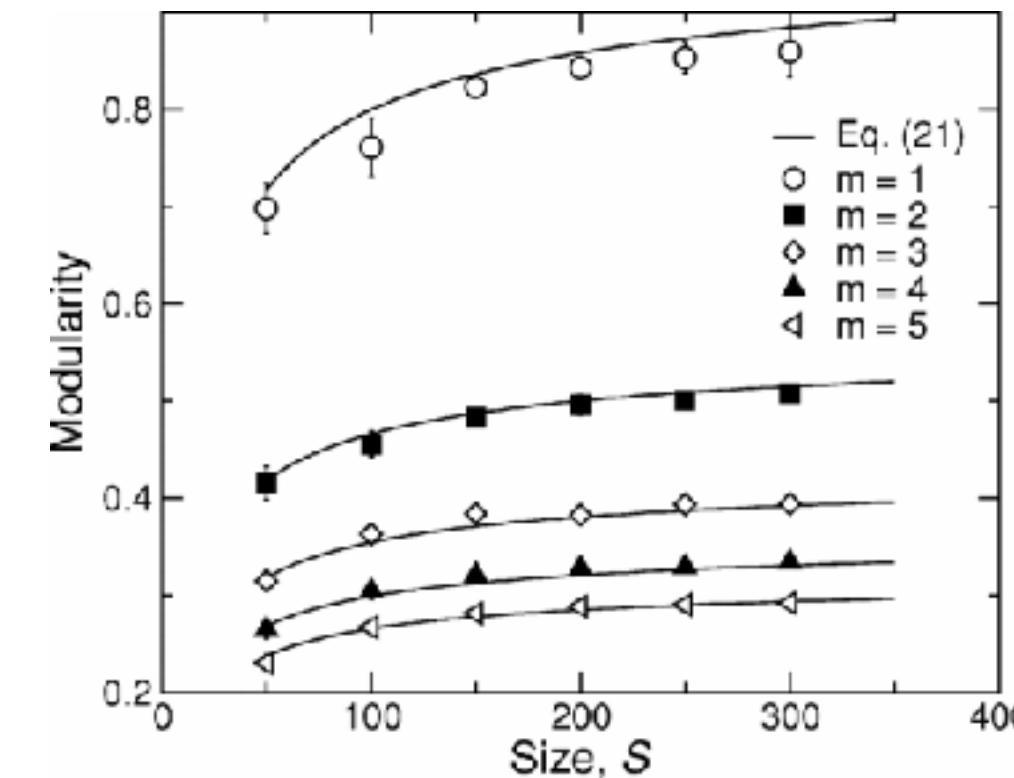
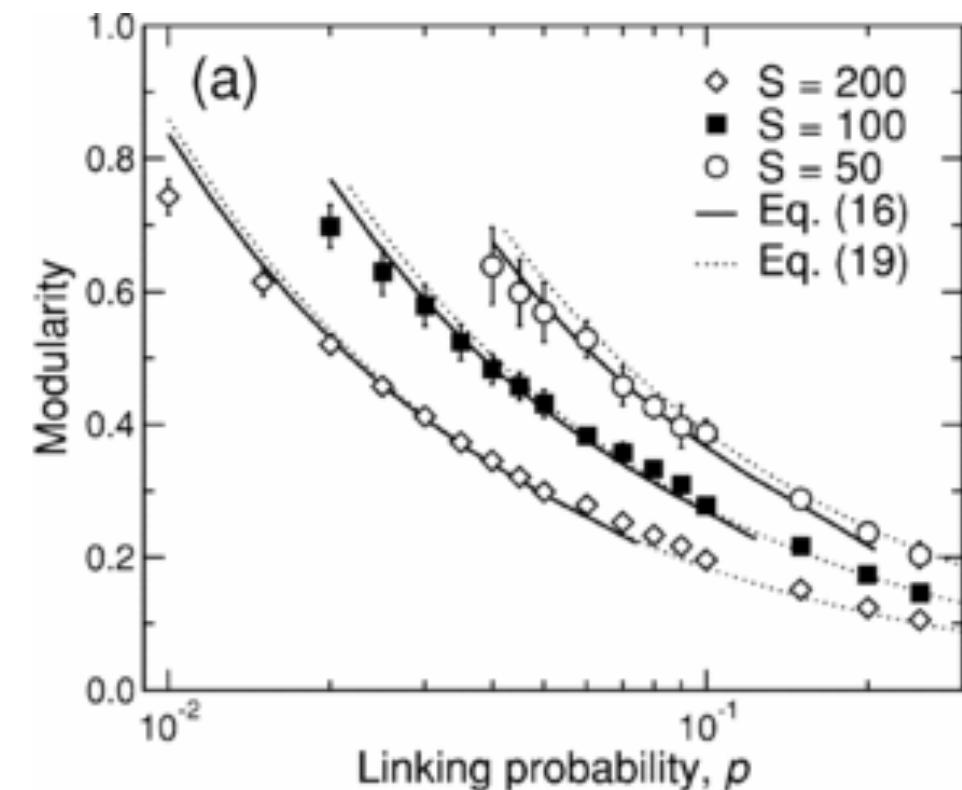
— Cris Moore

We can always maximize  $Q$  to find a partition, but is it meaningful?

# Fooled by “structure” in totally random networks

As it turns out, you can find high-modularity partitions in random networks.

*Structure is what makes data different from noise.  
It's what makes a network different from a random graph.*

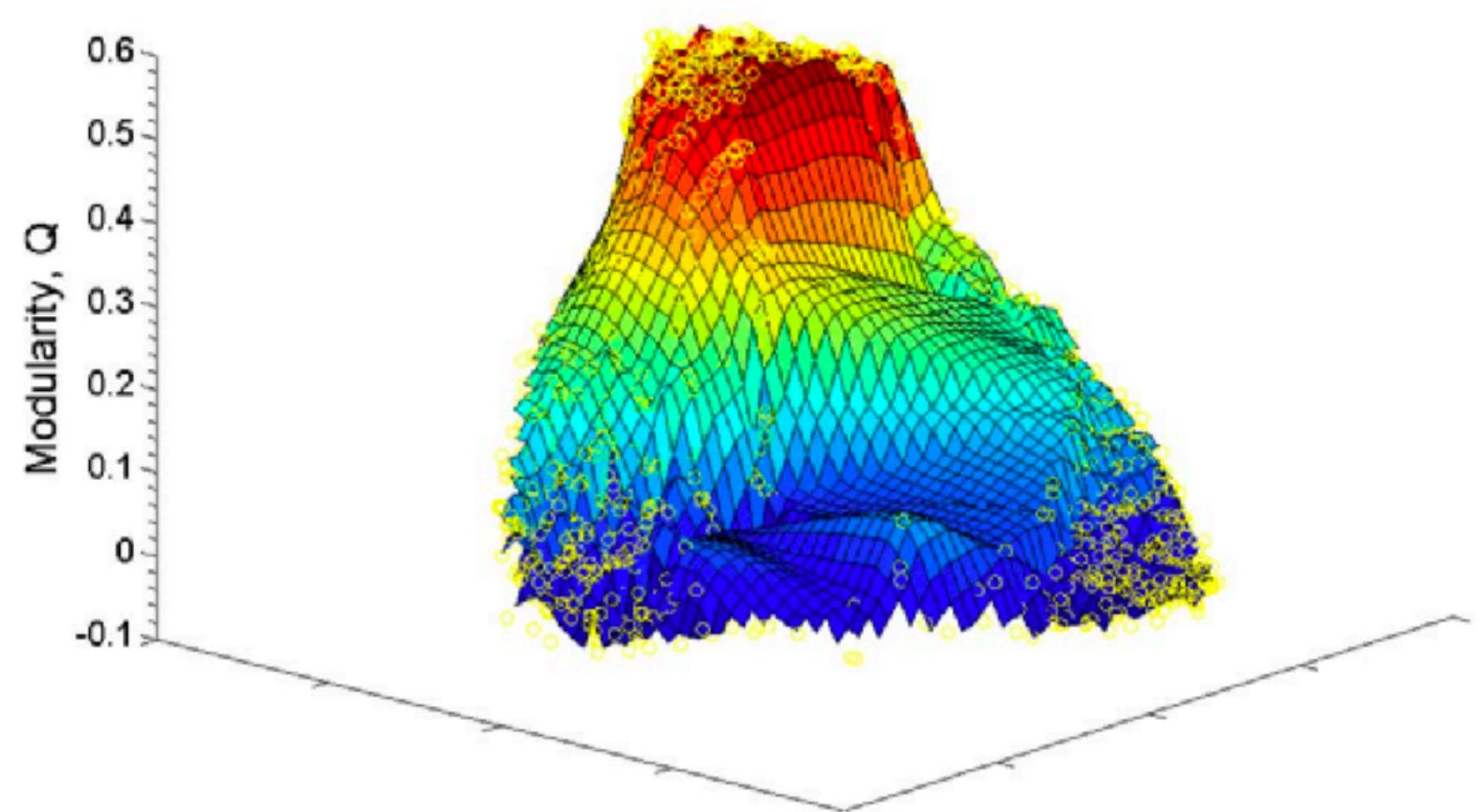
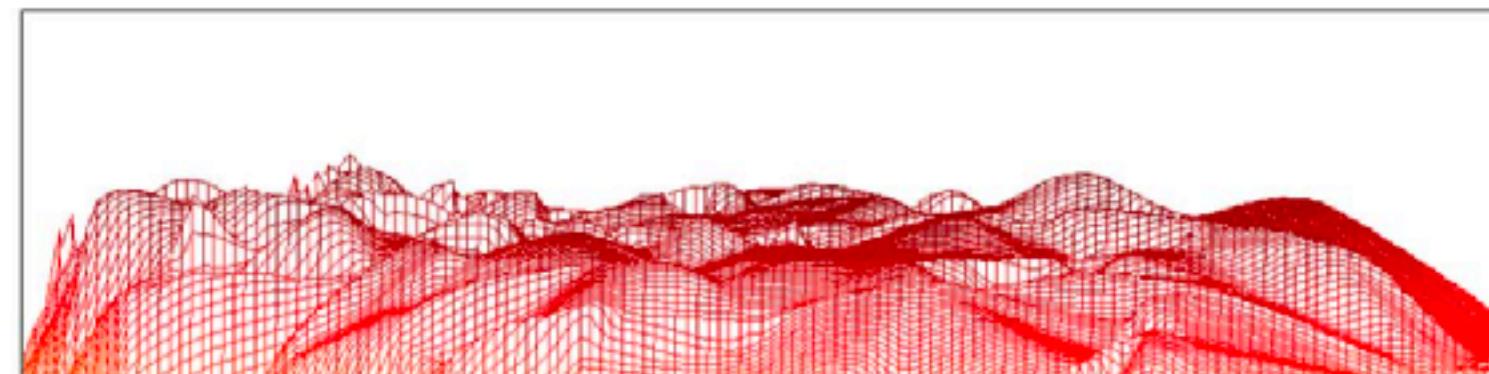


We prefer that our methods fail gracefully, and tell us when they fail. (like  $R^2$ )  
[alternative perspective: maybe you want to find clusters in randomness?]

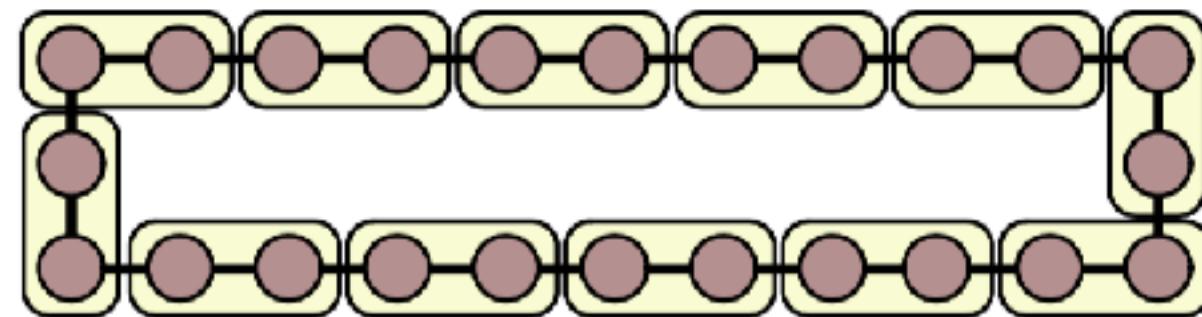
# Modularity: degeneracy and strange behavior

Lots of different but nearly-as-good partitions.

The optimization landscape is *degenerate*.



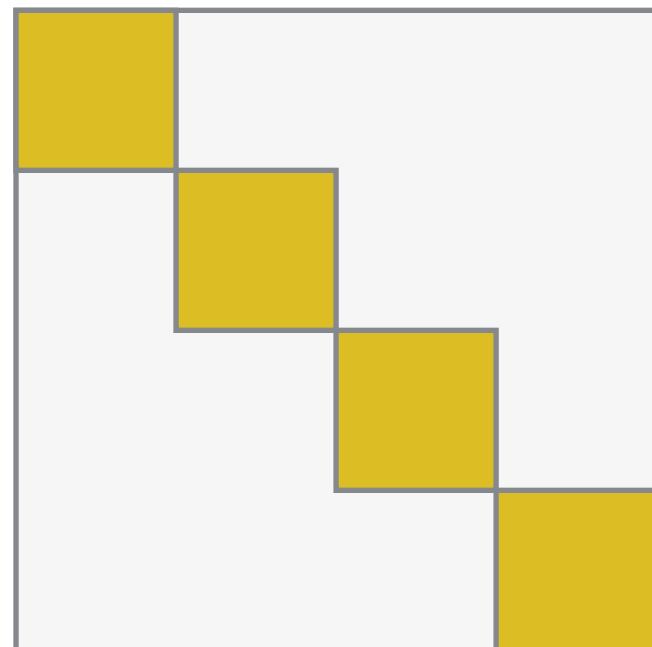
Unintuitive behavior: find the communities in a chain of cliques and you get pairs of cliques...not the cliques themselves!



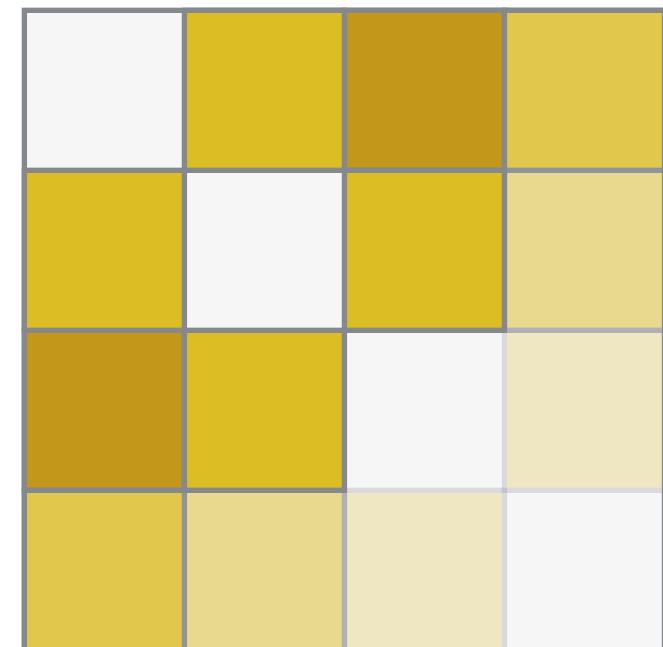
# $Q$ is restricted to *assortative* community structure

The zoo of possible structures is diverse and interesting!

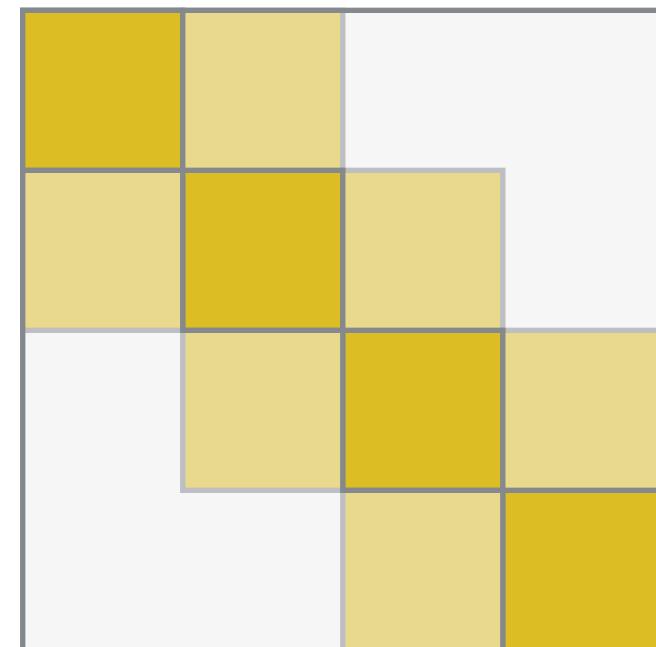
Build intuition: what do these networks look like?



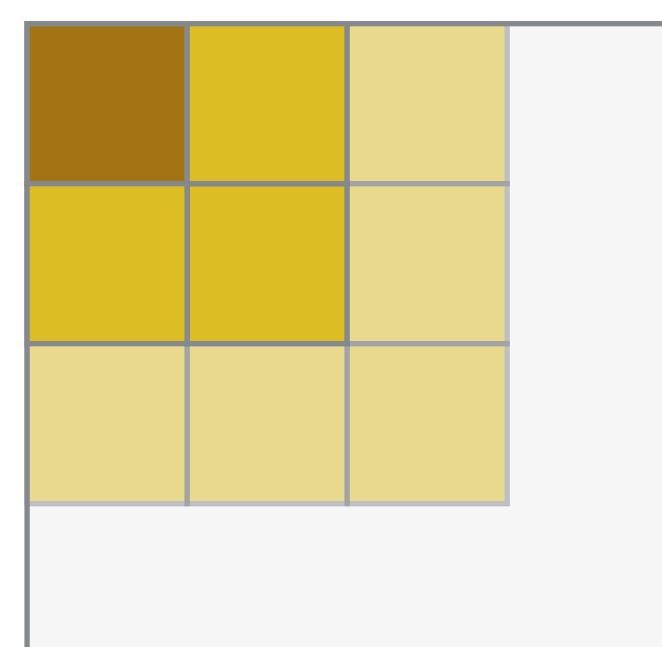
Assortative



Disassortative



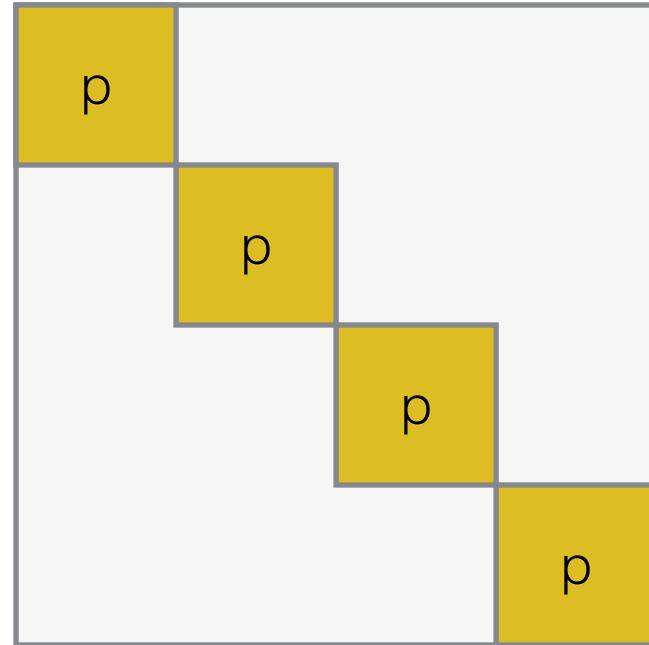
Ordered



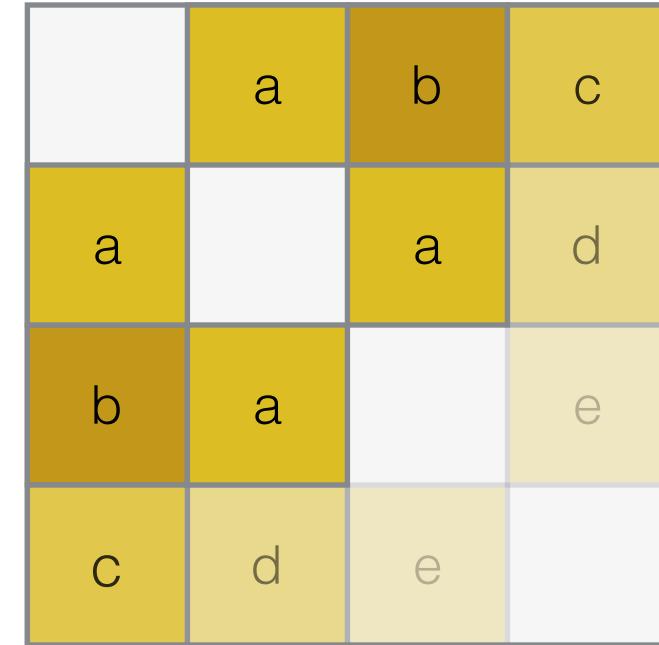
Core-periphery

# Beyond assortativity: block models

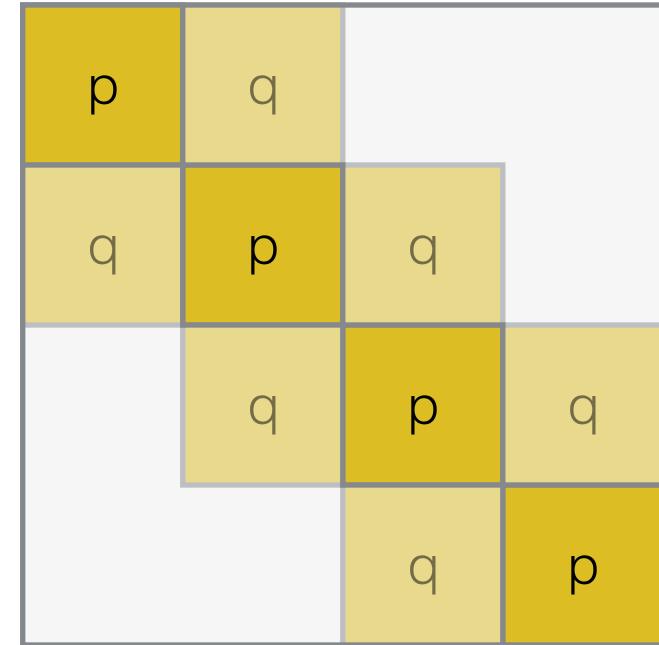
What do these have in common?



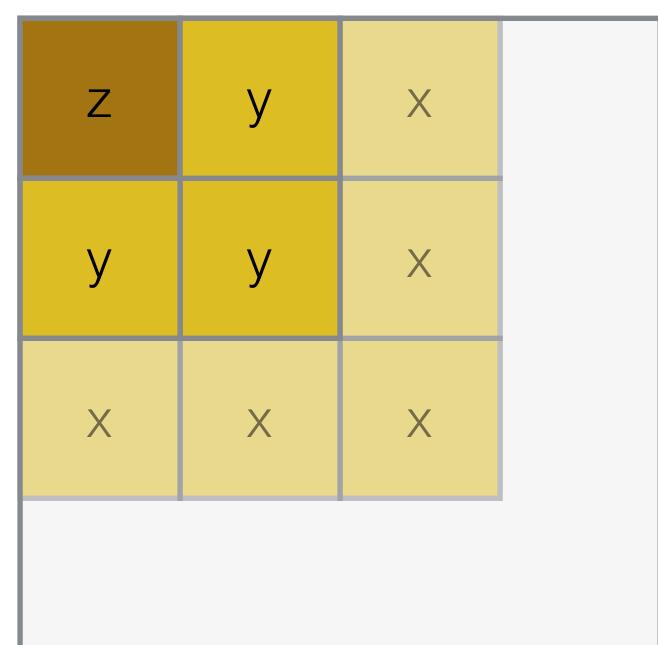
Assortative



Disassortative



Ordered



Core-periphery

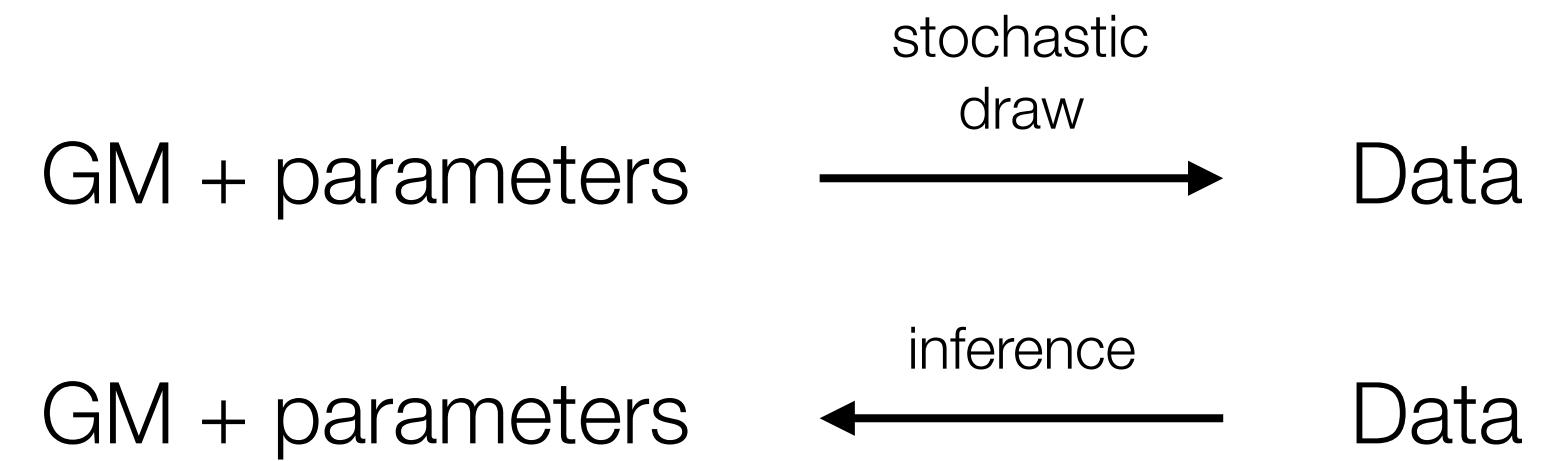
Nodes are in groups with other nodes that connect to other groups *in similar ways*.

**Key idea:** all nodes in a group are *stochastically equivalent*.

# Generative model approach

*Generate the structure you wish to infer.*

We like generative models because they open the door to inference:



In other words: let's write down a recipe for generating block structure.



# The stochastic block model

GM + parameters



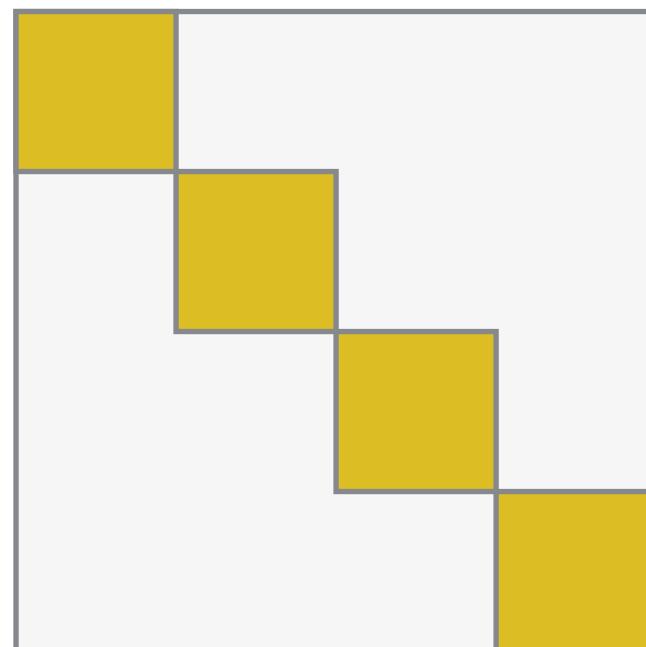
Assign each node to one of  $B$  blocks.

$$b_i$$

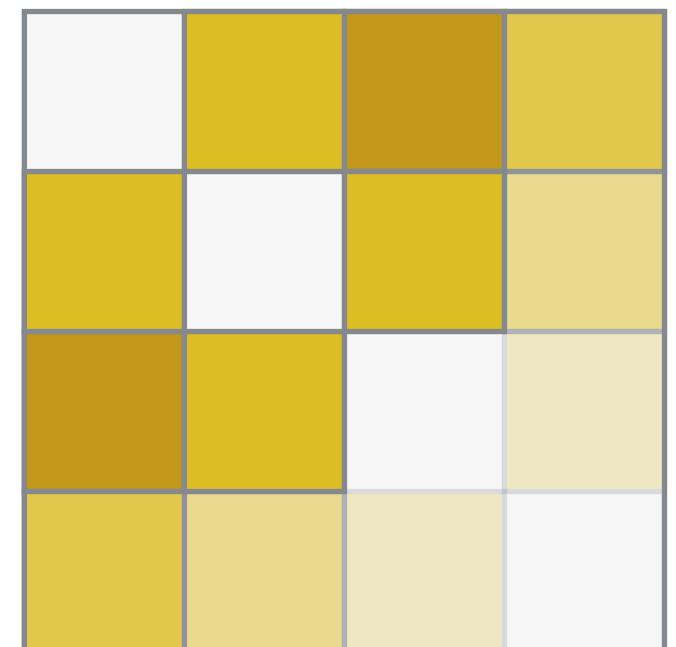
Let the probability that two nodes connect depend *only* on their blocks:

$$\Pr(A_{ij}|b_i, b_j) = \omega_{b_i, b_j}$$

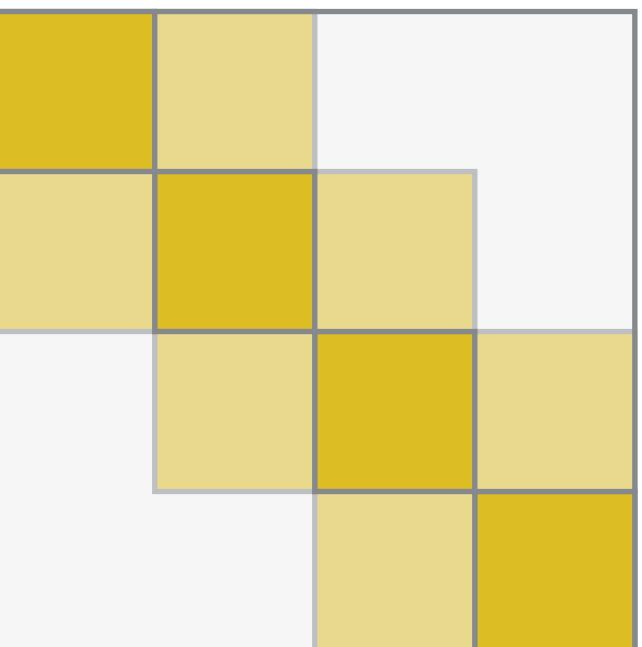
Then we can choose the matrix  $\omega$  to have whatever structure we want!



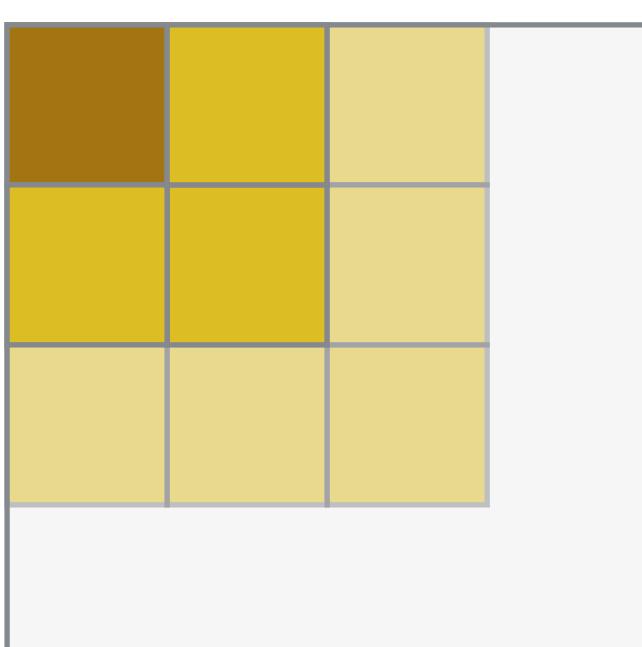
Assortative



Disassortative



Ordered



Core-periphery

# SBM inference

GM + parameters



Data

no more math on slides 😭

but the derivations are beautifully described in:

Karrer, Newman. Stochastic blockmodels and community structure in networks.  
Phys. Rev. E 83, 016107 (2011).

Recommended reading! ..... ↑


example matrix of parameters,  $B=4$

Summary:

1. Write down the SBM *likelihood function* for a fixed number of blocks  $B$ .
2. Maximize the likelihood with respect to matrix parameters.
3. Search over divisions into  $B$  blocks to find the best blocks.

There's a subtlety here, which I haven't written out, called *degree correction*. In practice, we also take into account the exact degree sequence. This allows us to find community structure while controlling for variability in the nodes' degrees.

# The problem with parameterized models...

You have to choose their parameters!

How should we choose  $B$ , the number of blocks?

Hint: we can't simply maximize the likelihood over all choices of  $B$ :

Why? If we place each node in its own community, we can get Likelihood=1.  
[Actually, this wouldn't model the data at all: it would *memorize* it.]

We need a way to penalize the complexity of the model. Any ideas?

# Description length & Occam's razor

The Description Length of a message is:

# bits required to send the compressed message + # bits in encoding scheme.

Occam's razor: among all possible explanation for a phenomenon, choose the simplest one. Therefore, choose the model with **Minimum Description Length** (MDL).

The stochastic block model also has a Description Length:

$$\Sigma = \boxed{\mathcal{S}} + \boxed{\mathcal{L}}$$

**description length** = entropy of data, given the model (fit SBM) + entropy of model

Consider the original problem:

what happens to this equation when I increase the number of blocks B?

# MDL criterion suggests an algorithm:

Fit the SBM with 1 block and record the Description Length.

Fit the SBM with 2 blocks and record the Description Length.

...

when the Description Length starts to increase, go back one step and stop.\*

Bonus: what happens if I *try to trick you* and give you a *random network with no blocks*?

MDL approach will tell you: your network is a random network with *one* block.

\*Actually, use something clever, like Golden Ratio / Fibonacci search

Press et al. Numerical Recipes: The Art of Scientific Computing, (Cambridge University Press, Cambridge, England, 2007), 3rd ed.

# So how does the search part work?

Markov-chain Monte Carlo:

Wander from one partition to another partition by proposing to take a node from one group and move it to a new group.

If this move *increases* the likelihood score, then keep the move.

If this move *decreases* the likelihood score, then maybe keep it, depending on how bad it is.

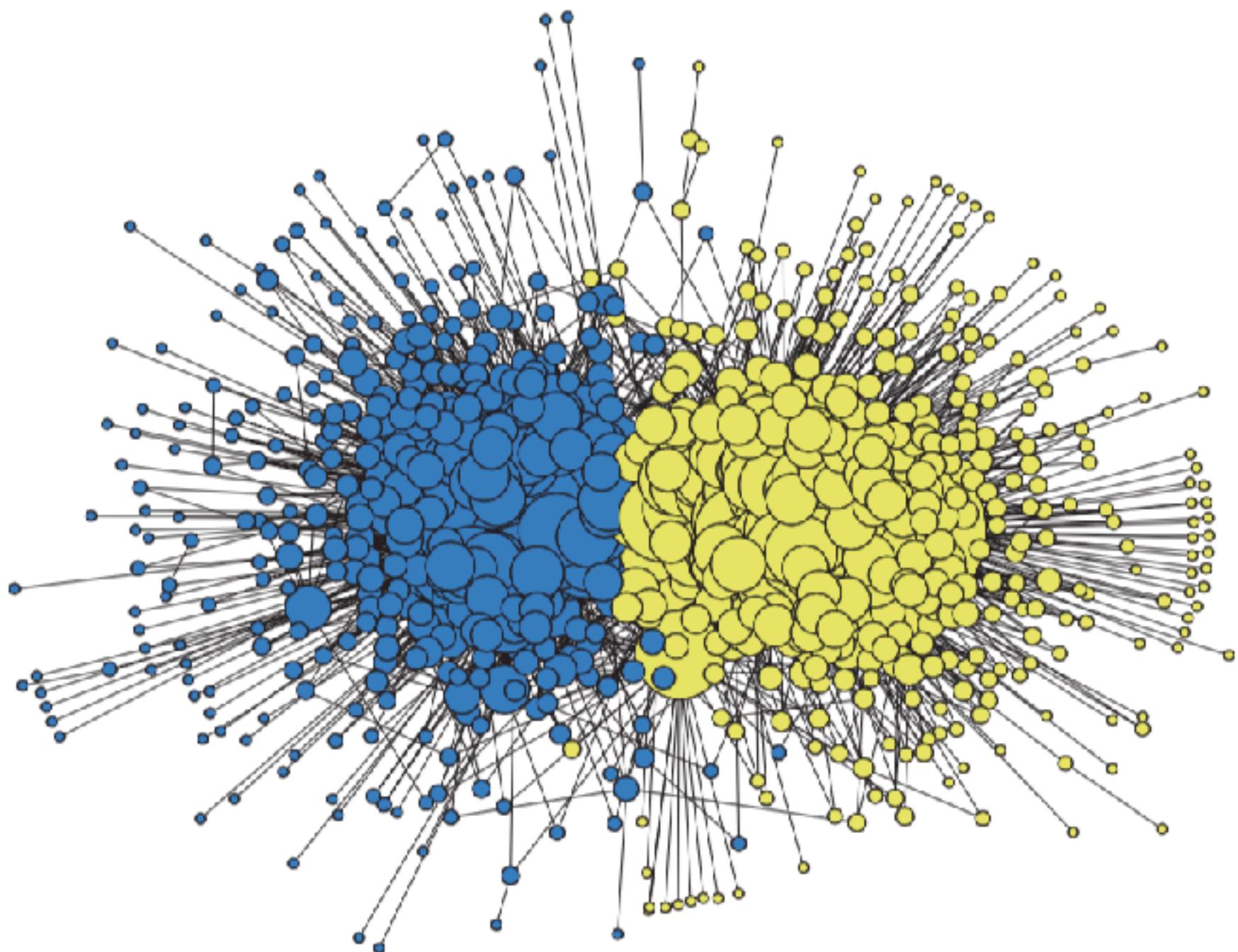
Thorough details in the documentation for graph-tool. <https://graph-tool.skewed.de>

# Does it work?

Adamic & Glance mapped the link structure of USA political blogs in 2004.

Karrer & Newman used this network as a testbed for community detection using the SBM.

What does this say about the process that may be generating (or pruning?) the links in this network?

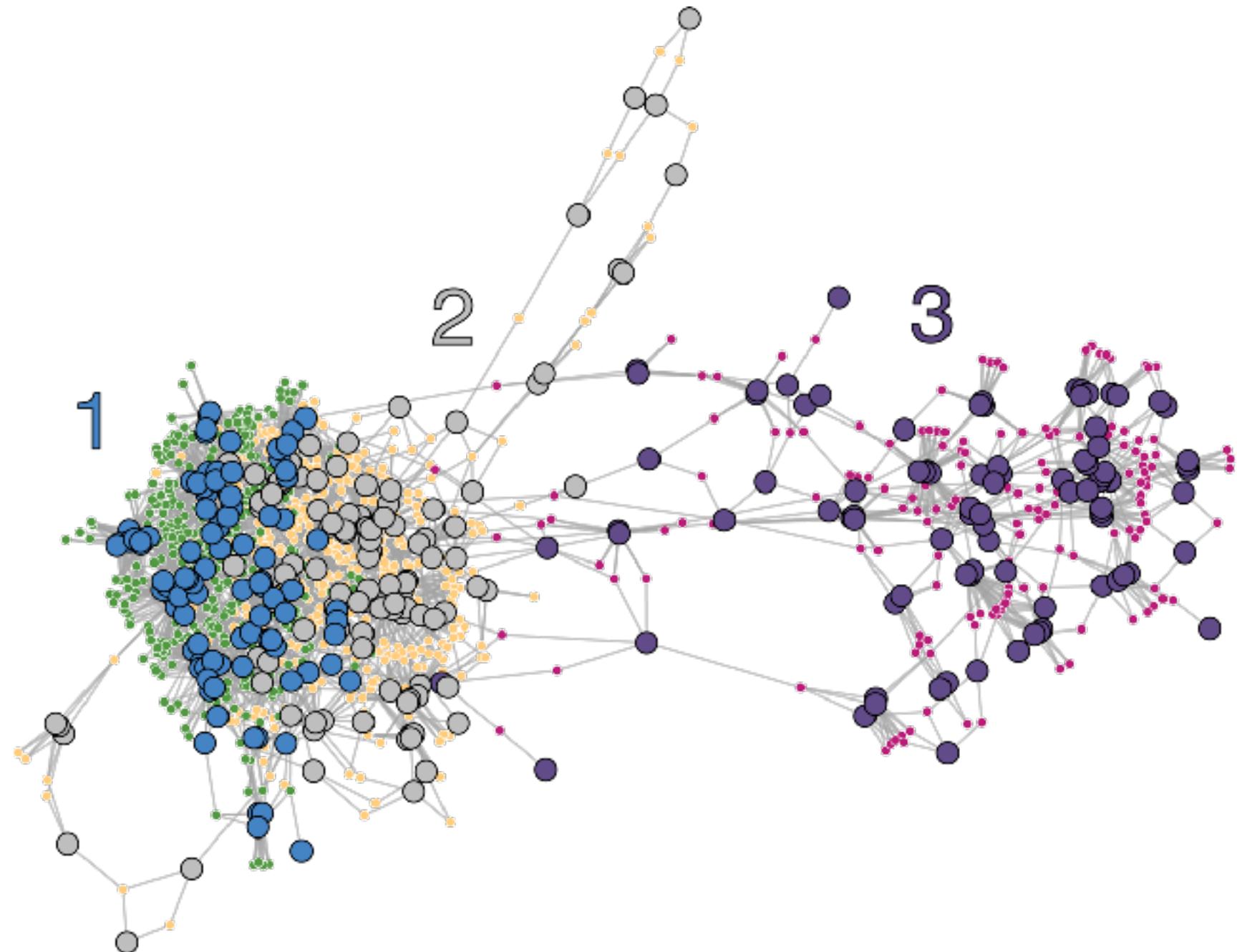


# Does it work?

In bipartite networks, we know the major split in the data already.

Methodologically, we found that exploiting this split improved speed and quality of the partitions we found.

Scientifically, this opened new directions to analyze (and understand) evolutionary constraints on malaria parasites.



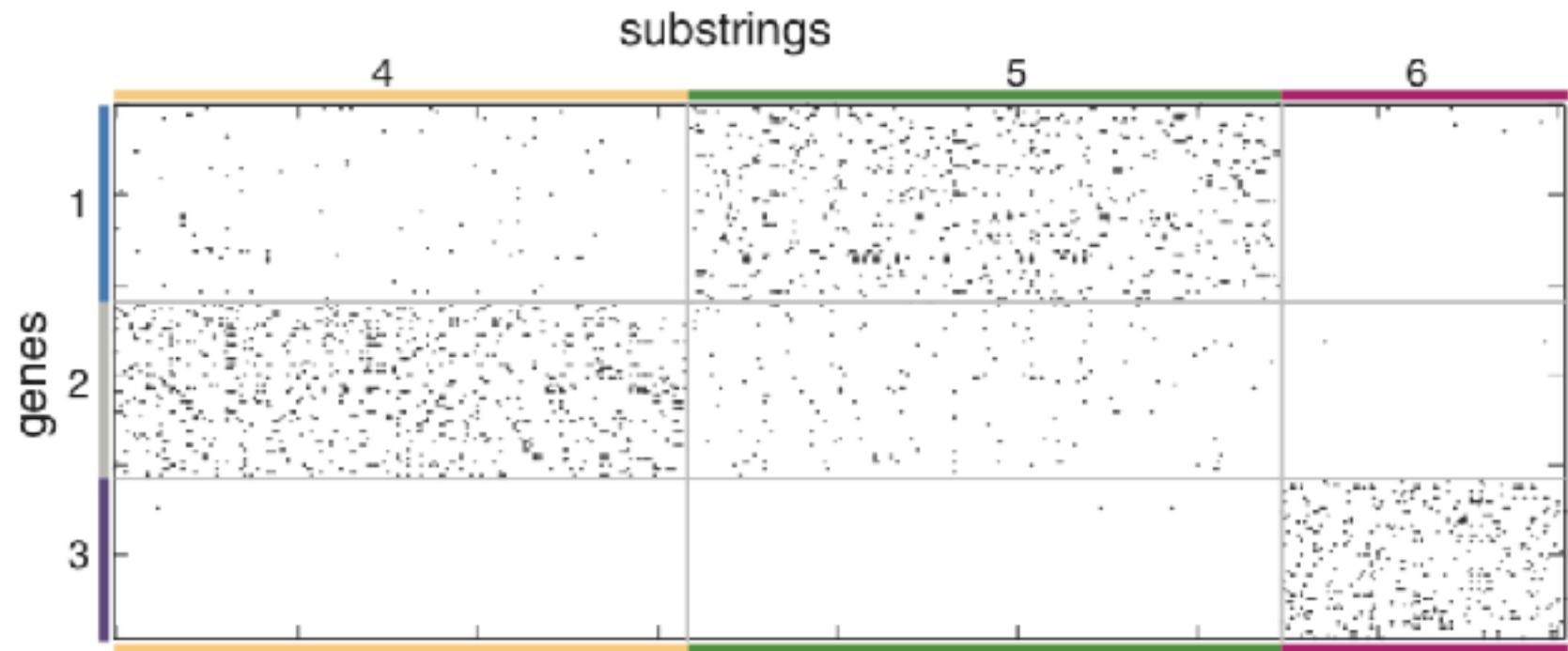
Genes & substrings,  
malaria immune evasion

# Does it work?

In bipartite networks, we *know* the major split in the data already.

Methodologically, we found that exploiting this split improved speed and quality of the partitions we found.

Scientifically, this opened new directions to analyze (and understand) evolutionary constraints on malaria parasites.



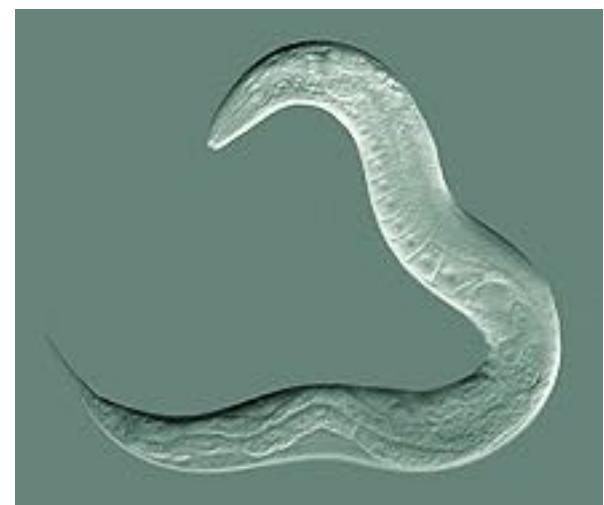
Genes & substrings,  
malaria immune evasion

Larremore, Clauset, Buckee, *PLoS Comp Biol*, 2013.

Larremore, Clauset, Jacobs, *Physical Review E*, 2014.

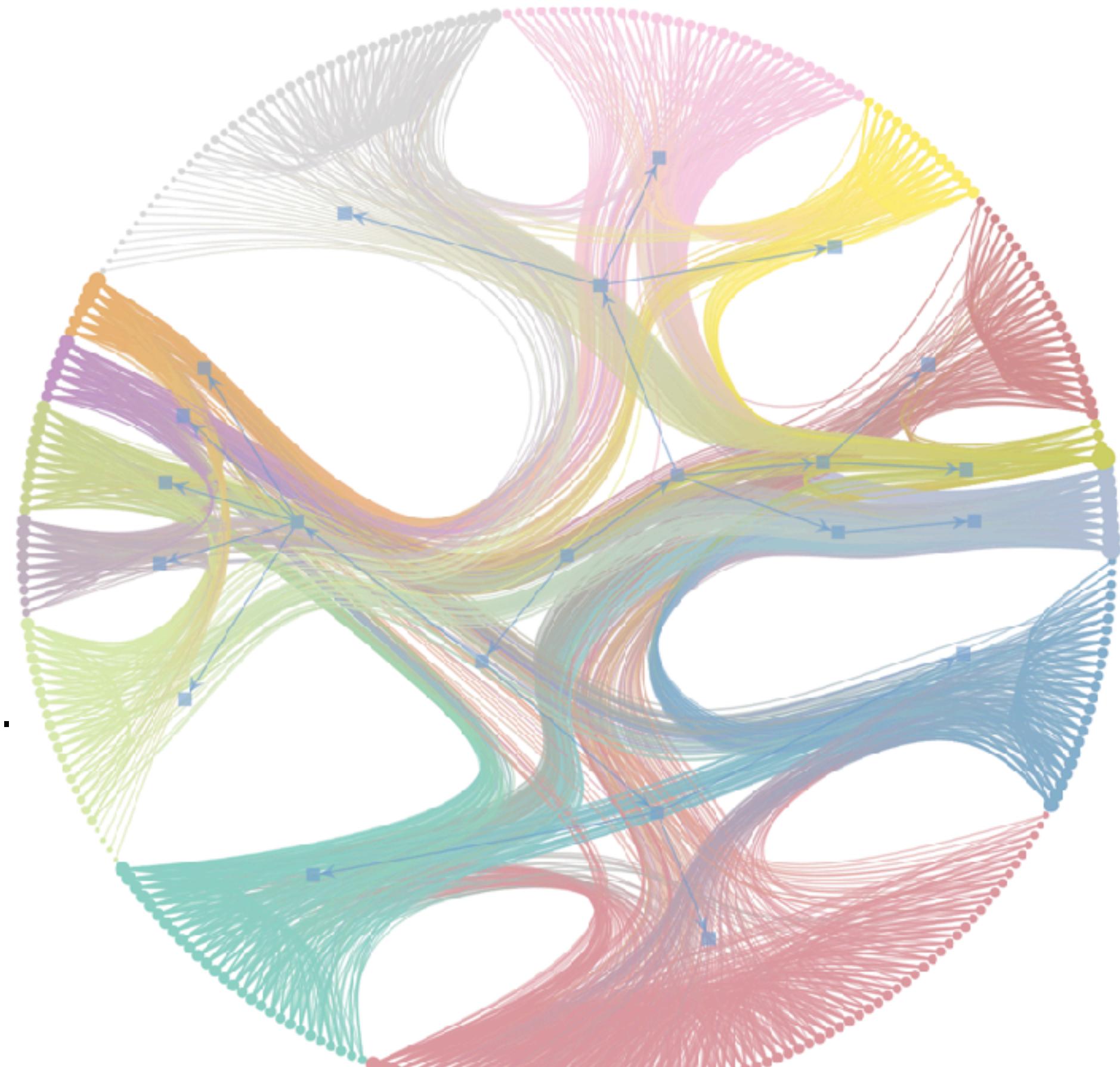
# Does it work?

*C. elegans* neuronal network.



297 neurons, completely mapped.  
The neurons do not fire action  
potentials, and do not express  
any voltage-gated ion channels.

Note the different layout...



# A good alternative? Cross-validation via link prediction

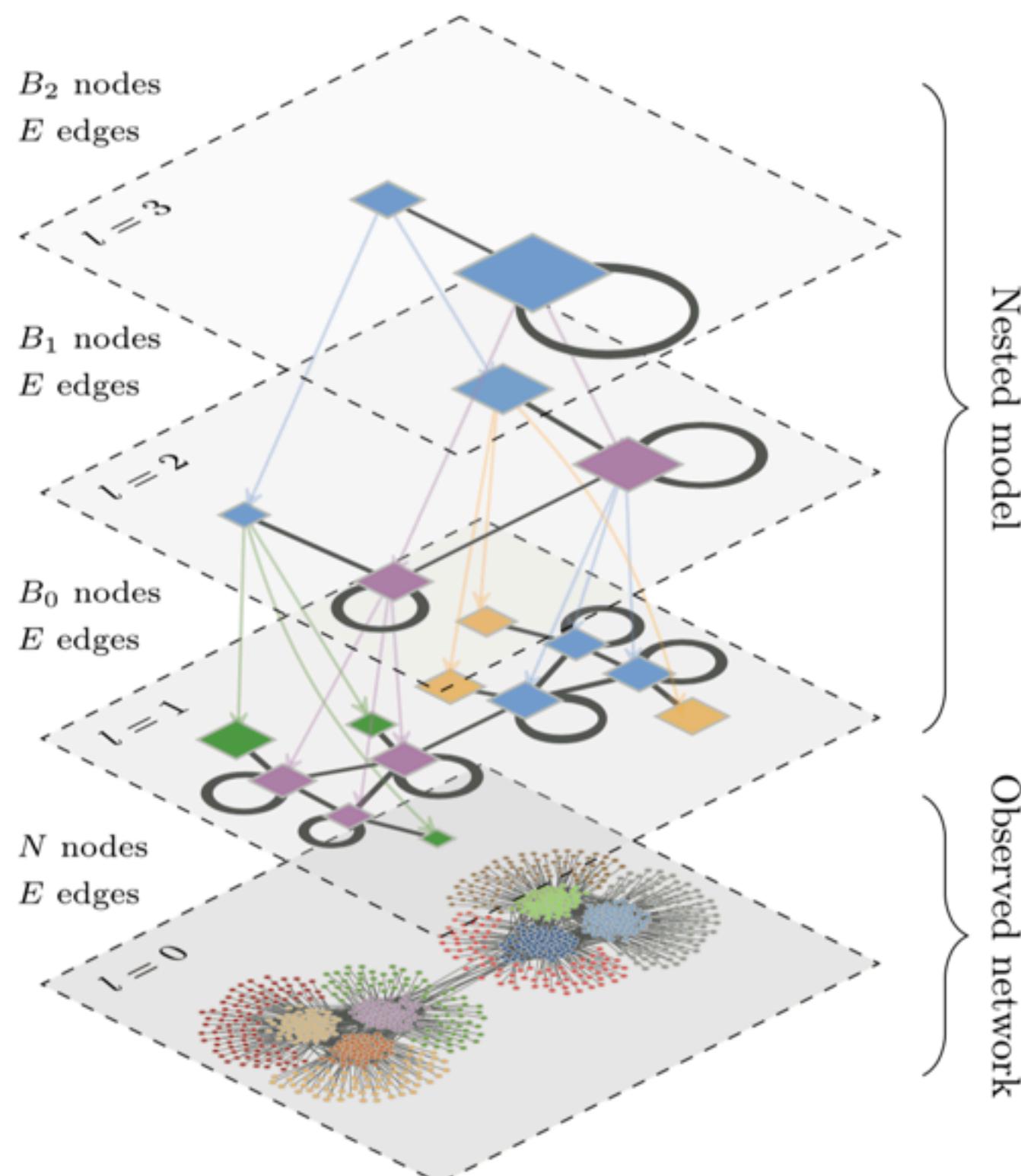
Select  $B$  by choosing the model that **makes the best predictions**.

## **Perform $k$ -fold cross validation:**

1. Divide the edges of the network into  $k$  groups, called folds.
2. Hide one of the folds (the “test set”)
3. Fit each SBM to the remaining  $k-1$  folds (the “training set”), varying  $B$ .
4. Test the ability of the fitted models to predict the hidden test data.
5. Switch which fold is “test” and which are “training” and repeat.

Choose the  $B$  with the highest performance on link prediction over all  $k$  folds.

# Advanced topic 1: hierarchical communities



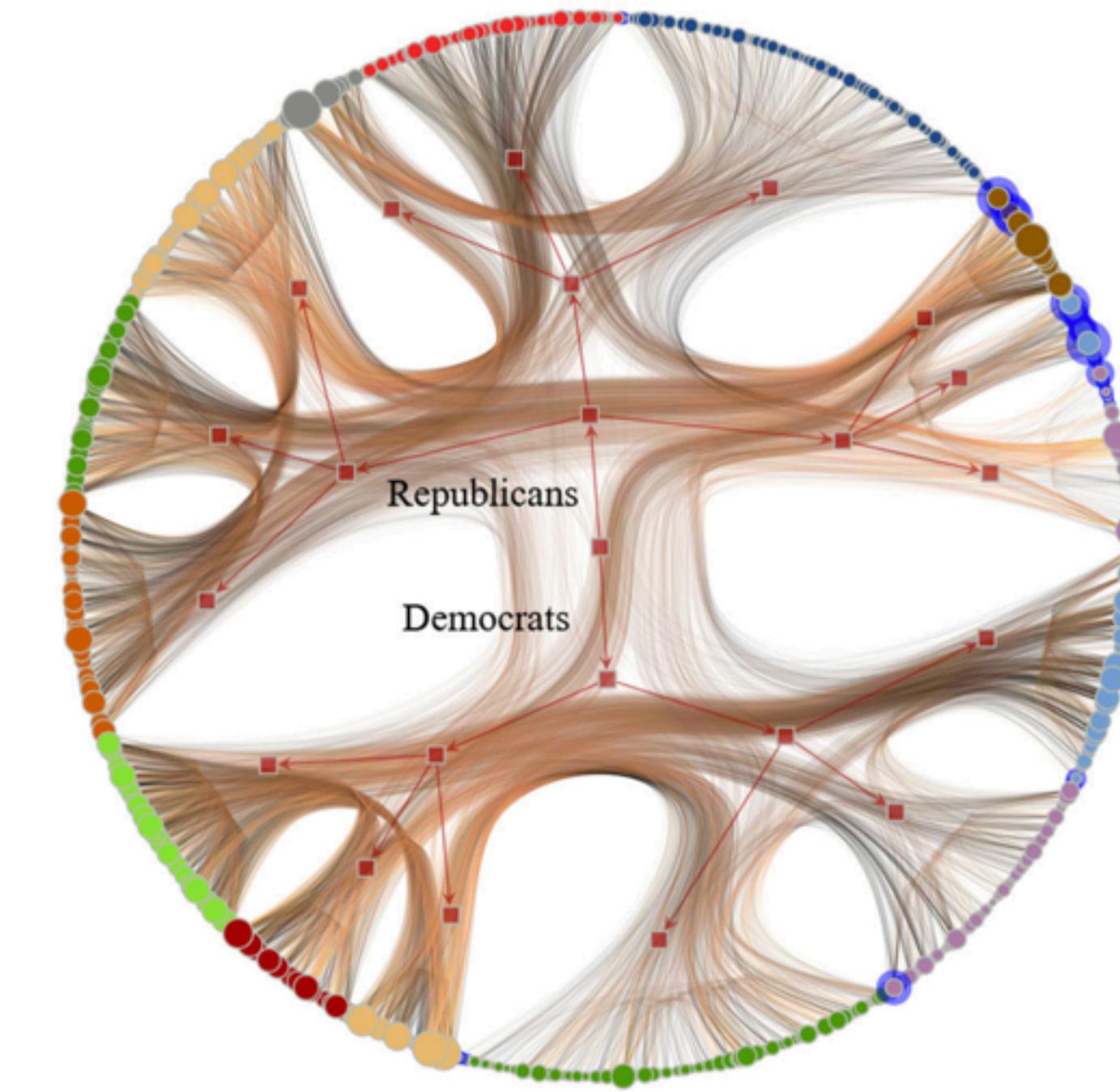
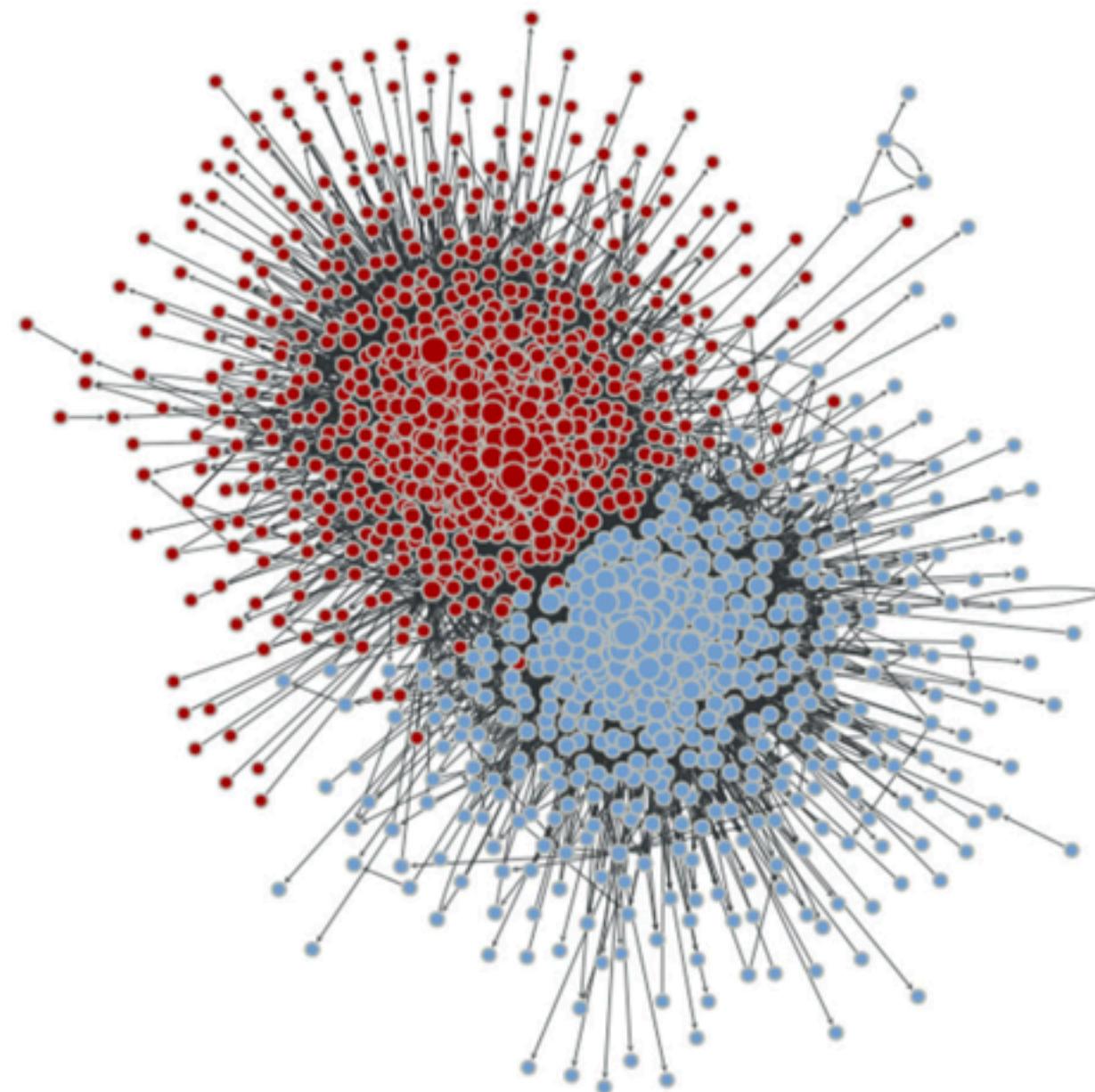
Don't minimize the description length  
of the data and the model...

Model the model as well. Why?

If we compress the model, we can  
afford a bigger model, but a lower  
overall cost.

Except now, the description length  
includes two models. Or three? Or?

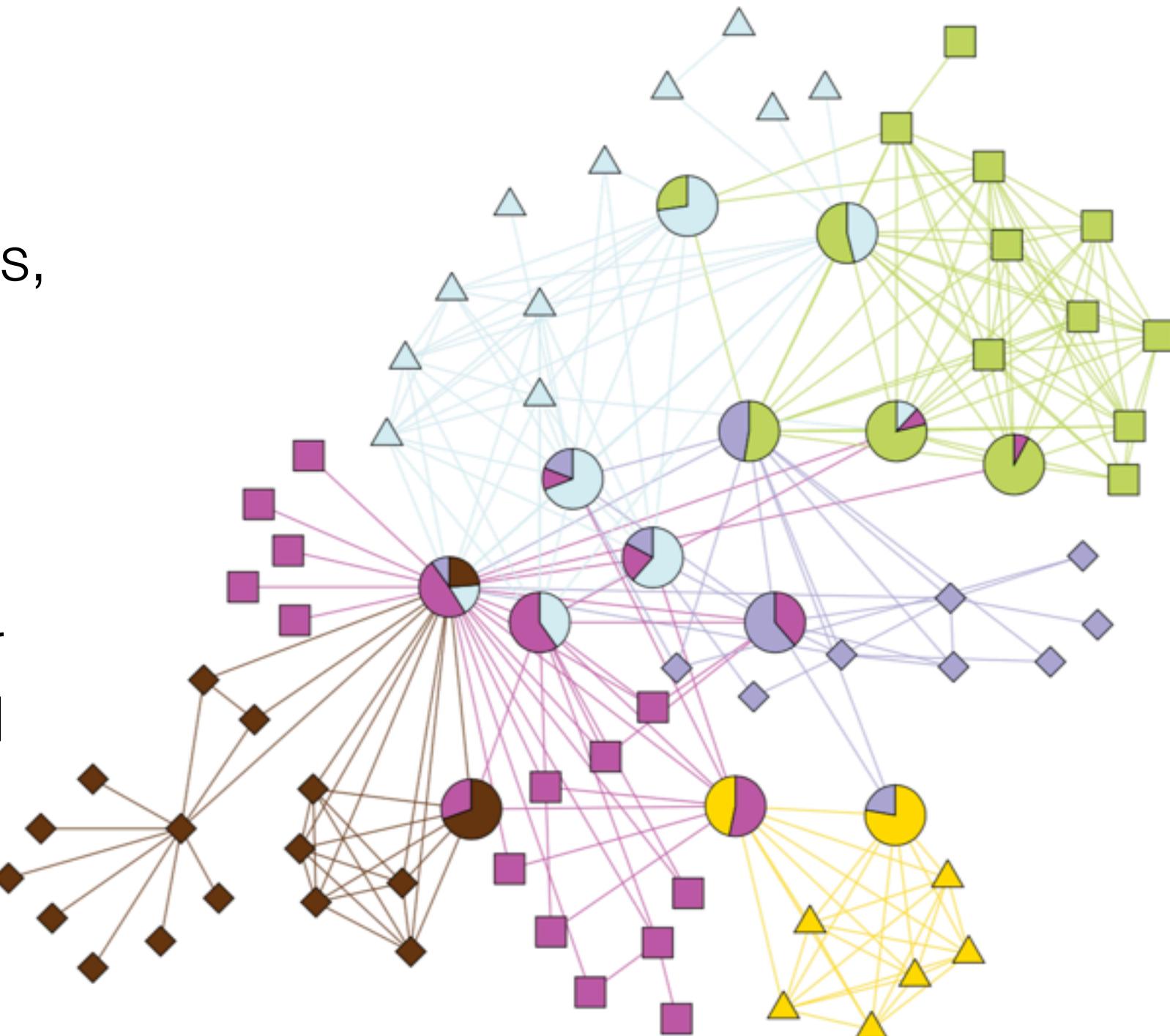
# Advanced topic 1: hierarchical communities



# Advanced topic 2: mixed-membership

Nodes are often pulled between communities. (Or in real social systems, individuals belong to multiple groups.)

“Mixed membership” models allow for that, by assigning *links* to groups, and assigning nodes to groups based on their links.



# Advanced Topic 3: multilayer networks

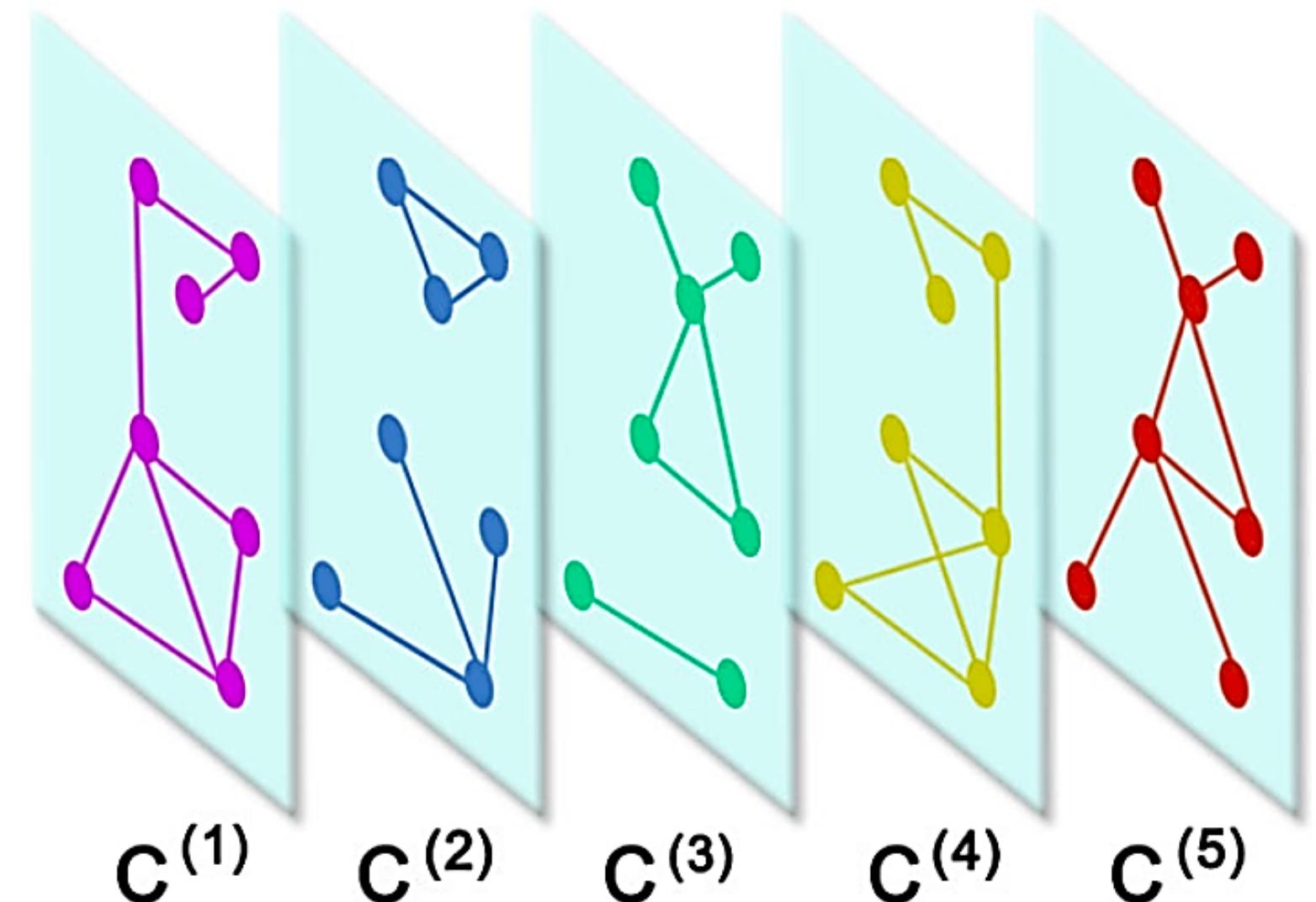
In single-layer networks:  
nodes and edges

In multilayer networks:  
nodes, edges, and layers

**edges**: different types of relationships

**layers**: each layer contains all edges of one type

**nodes**: same nodes in each layer



# Multilayer network: air travel



Ryanair

Lufthansa

Vueling

British airways

Aggregate



traditional: booking with airline



disrupted: booking with kayak, expedia, etc

# Multilayer network: community structure?

## three key approaches:

1. **Non-generative**: modularity maximization; vary inter-layer strength.

Mucha et al *Science* 2010. <http://science.sciencemag.org/content/328/5980/876>

2. **Generative**: SBM for each layer, but jointly model layers whenever their structures are sufficiently similar.

Peixoto, T. P. *Phys. Rev. E* 92, 042807–15 (2015).

3. **Generative**: SBM for each layer, and model all layers simultaneously with same community structure, but allow relationships between groups to vary.

De Bacco Power Larremore Moore. *Phys. Rev. E* 95, 1981–10 (2017).

1 or 2 are preferred if nodes appear/disappear over time.

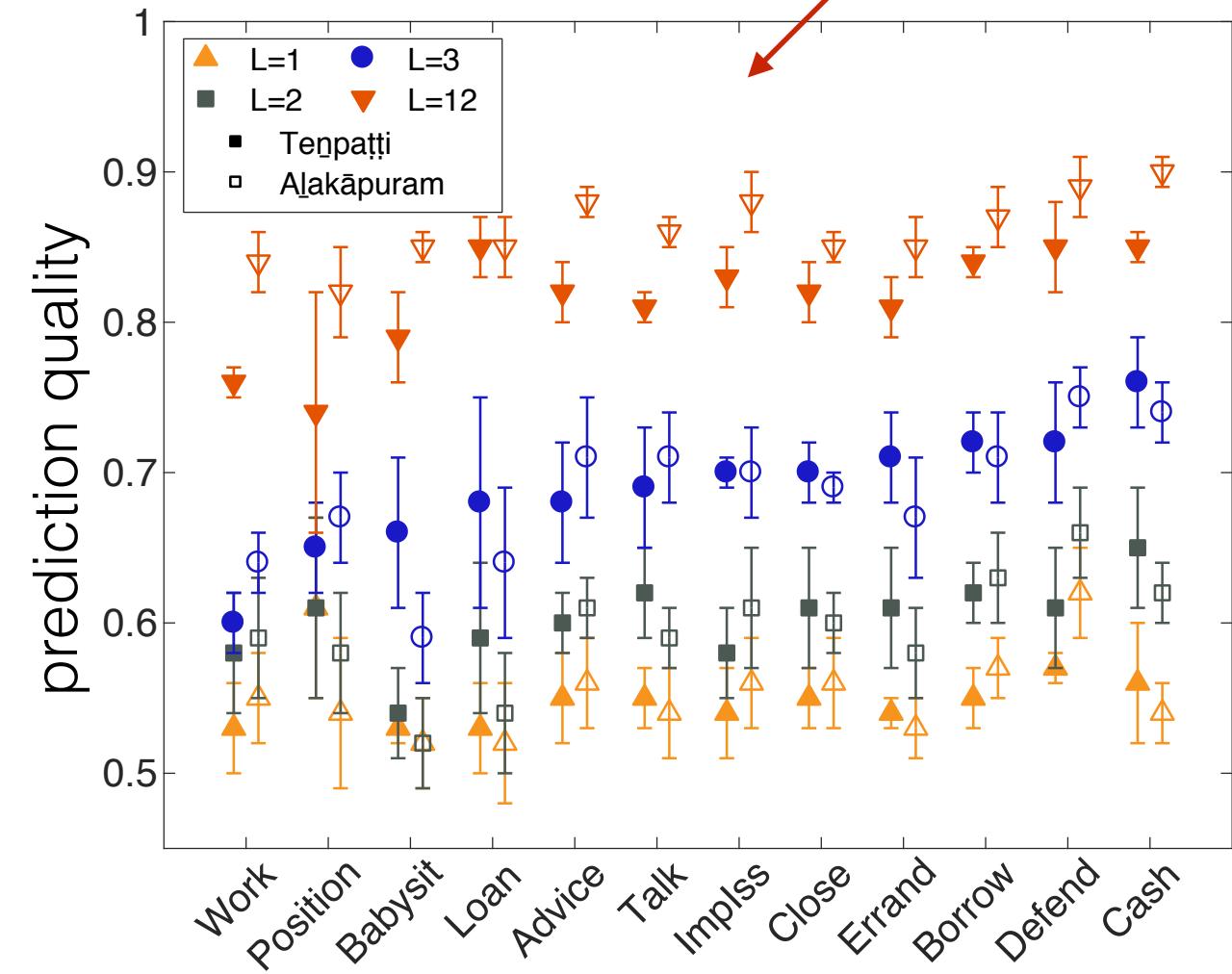
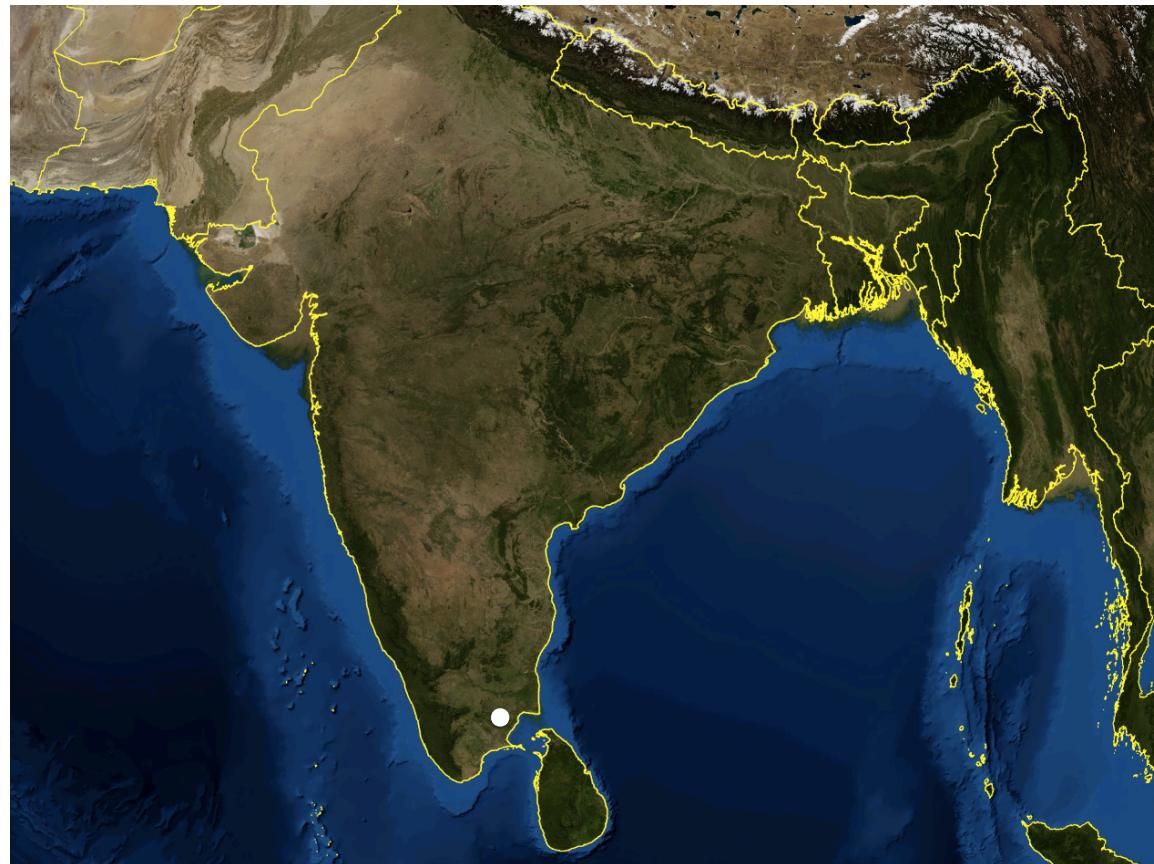
2,3 are preferred to solve the *layer interdependence problem*

# Layer interdependence

more layers = better performance  
(layer structure generated by same social mech.)

Are layers structurally similar? Complementary? Neither?

“Learn” a SBM from  $m$  layers; try to predict links of  $m+1$ .

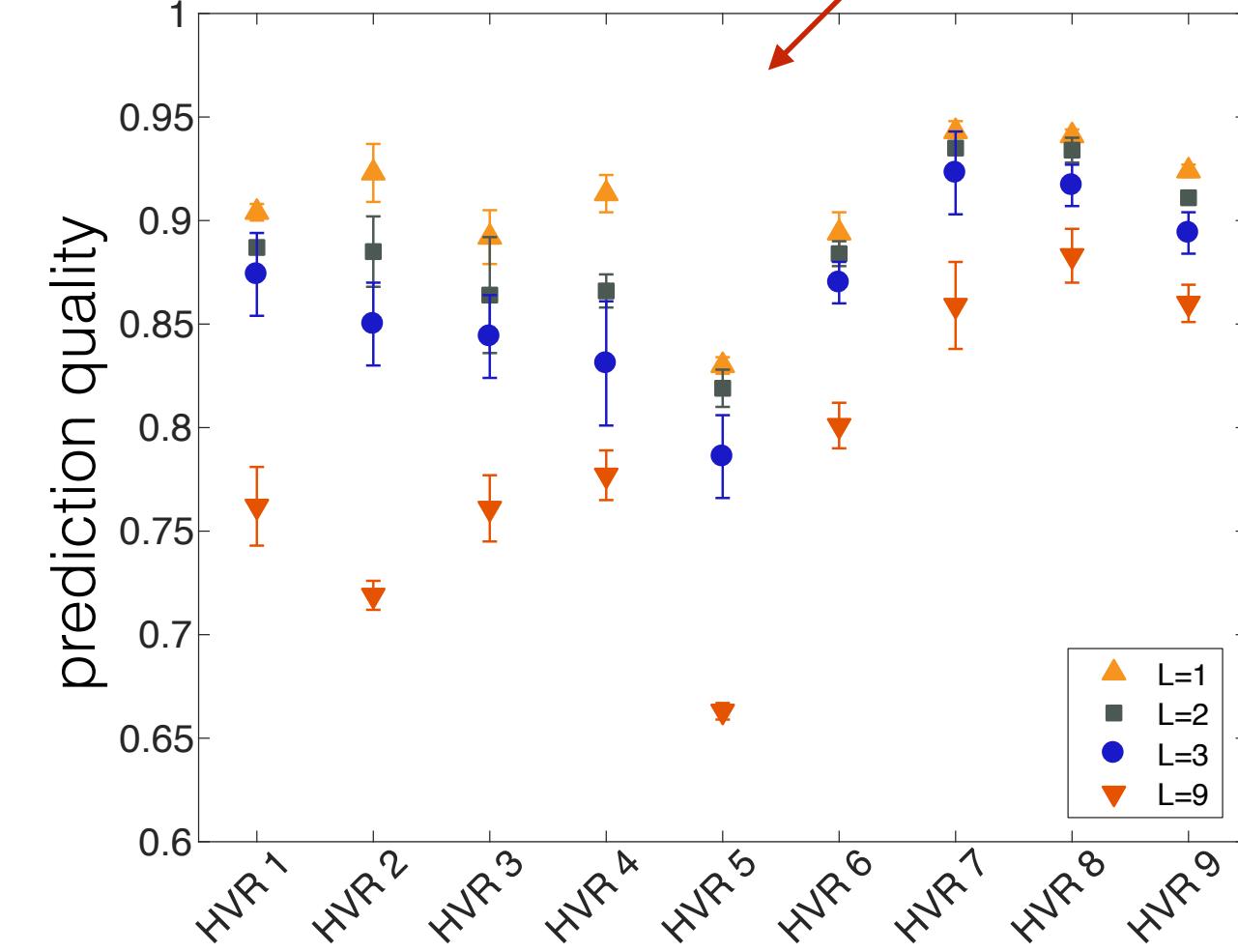


12 layer social support network across 2 villages in South India.

# Layer interdependence - malaria



more layers = worse performance  
(layer structure generated by different biol. mech.)



cannot predict the structure of one region in the immune-evasion genes  
by using other regions; layers are unrelated!

# Advanced topic 4: metadata+communities

## **What are metadata?**

How well do metadata explain the network structure? “BESTest”

Peel\*, Larremore\*, Clauset. *Science Advances* 3(5) e1602548. (2017) <http://advances.sciencemag.org/content/3/5/e1602548.full>

How do metadata relate to network structure? “neoSBM”

Peel\*, Larremore\*, Clauset. *Science Advances* 3(5) e1602548. (2017) <http://advances.sciencemag.org/content/3/5/e1602548.full>

Can we use metadata to guide community detection? “metadata assisted SBM”

Newman, Clauset. *Nature communications* 7 (2016). <https://www.nature.com/ncomms/2016/160616/ncomms11863/full/ncomms11863.html>

Can we find patterns in the metadata itself? Apply multilayer SBM

Peixoto. Hierarchical Block Structures and High-Resolution Model Selection in Large Networks. *Phys. Rev. X* 4, 011047 (2014).

# Blockmodel entropy significance test

How well do the metadata explain the network?

randomly assigned metadata  
→ model gives no explanation, high  $H$

metadata correlated with communities  
→ model gives good explanation, low  $H$

1. Divide the network  $G$  into groups according to metadata labels  $M$ .
2. Fit the maximum likelihood parameters of an *a posteriori* SBM and compute the entropy  $H(G,M)$  of the corresponding ensemble.
3. Compare the entropy of this SBM ensemble to distribution of entropies from SBMs partitioned using shuffled metadata  $\underline{M}$ .

$$\text{p-value} = \Pr[H(G,\{\underline{M}\})] \leq H(G,M)]$$

# Multiple network layers; multiple metadata attributes

Network	Status	Gender	Office	Practice	Law School
Friendship	$< 10^{-6}$	0.034	$< 10^{-6}$	0.033	0.134
Cowork	$< 10^{-3}$	0.094	$< 10^{-6}$	$< 10^{-6}$	0.922
Advice	$< 10^{-6}$	0.010	$< 10^{-6}$	$< 10^{-6}$	0.205

model = SBM

Multiple sets of metadata **significantly explain** multiple networks.  
[Should one particular set of metadata be ground truth?]

# BESTest accommodates many models of group structure

Network	Model	
	SBM	DCSBM
Malaria 1	0.566	0.066
Malaria 2	0.064	0.126
Malaria 3	0.536	0.415
Malaria 4	0.588	0.570
Malaria 5	0.382	0.097
Malaria 6	0.275	0.817
Malaria 7	0.020	0.437
Malaria 8	0.464	0.143
Malaria 9	0.115	0.104

metadata = parasite origin

A negative result: parasite origin is irrelevant to genetic substring-sharing.

Malaria parasites *do not* have a strong strain structure, with implications for diversifying selection among parasites.

# neoSBM

Choose between the **SBM partition** and the **metadata partition**.

$$\mathcal{L}_{\text{neoSBM}} = \mathcal{L}_{\substack{\text{SBM} \\ \text{log likelihood}}} + f(\theta)_{\substack{\text{cost} \\ \text{log likelihood}}}$$

neoSBM  
log likelihood

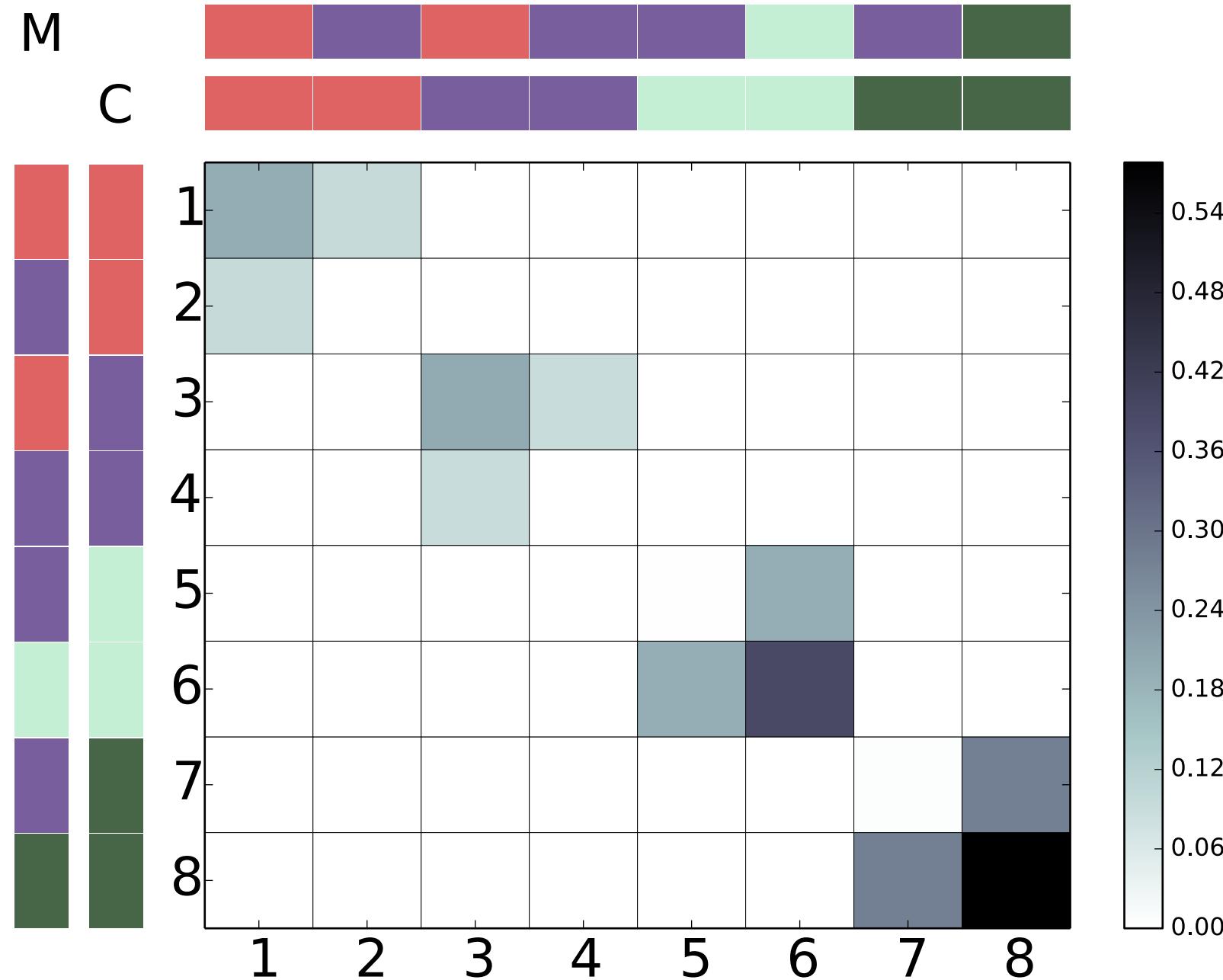
Log likelihood with parameterized prior:

$\theta$  is the parameter of a Bernoulli prior on whether the node is **free to choose its own community** or held **fixed at its metadata label**.

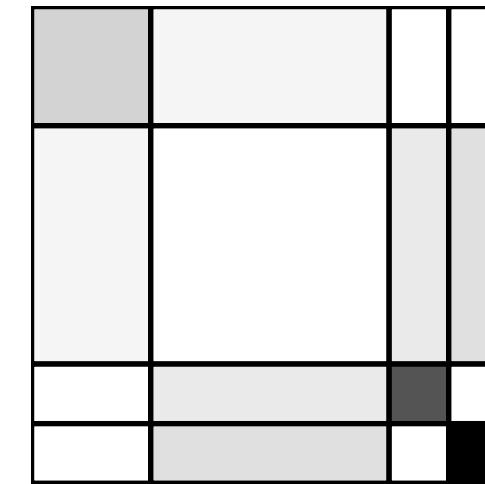
As  $\theta$  increases, the cost of freeing a node decreases.

Varying  $\theta$  in the unit interval **explores the space of partitions** between  $M$  and  $C$ .

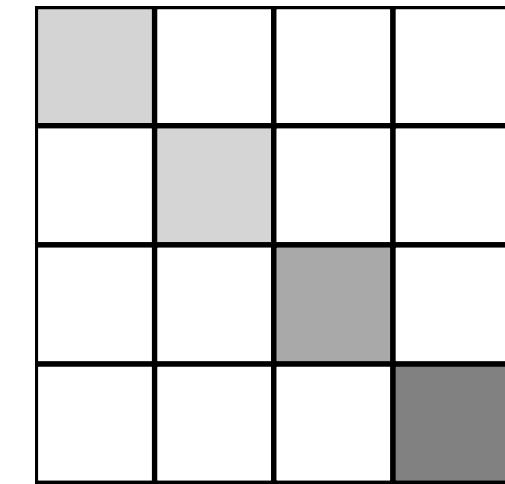
# Plant two different *kinds* of structure in a network



SBM with 8 groups and  
two interesting 4-group partitions:

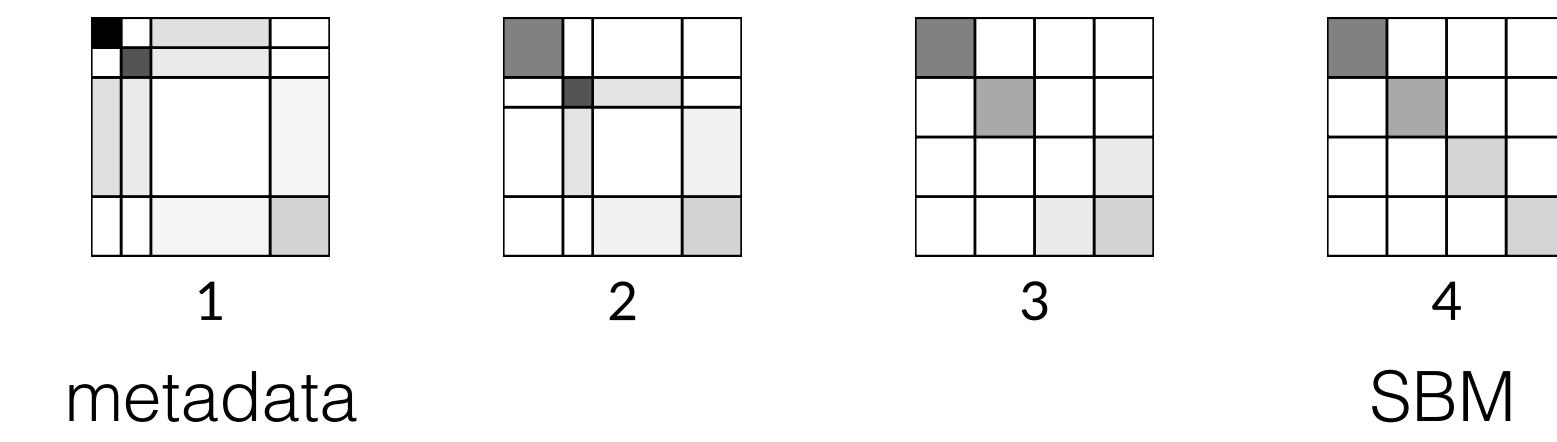
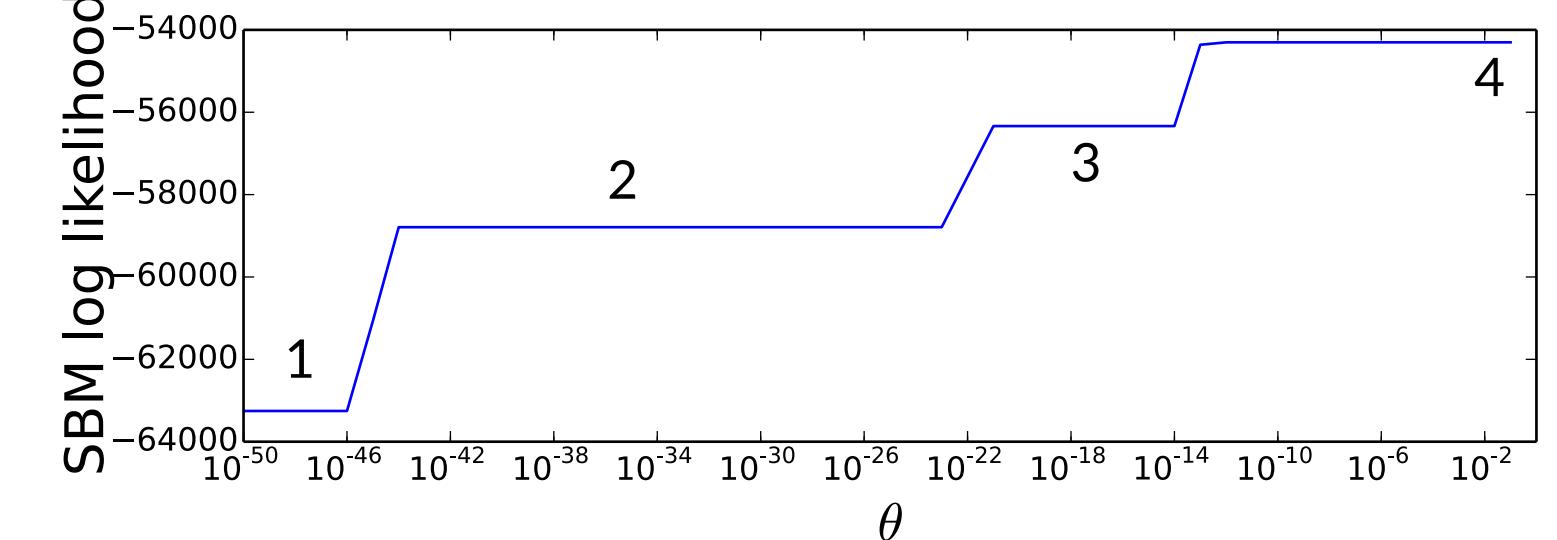
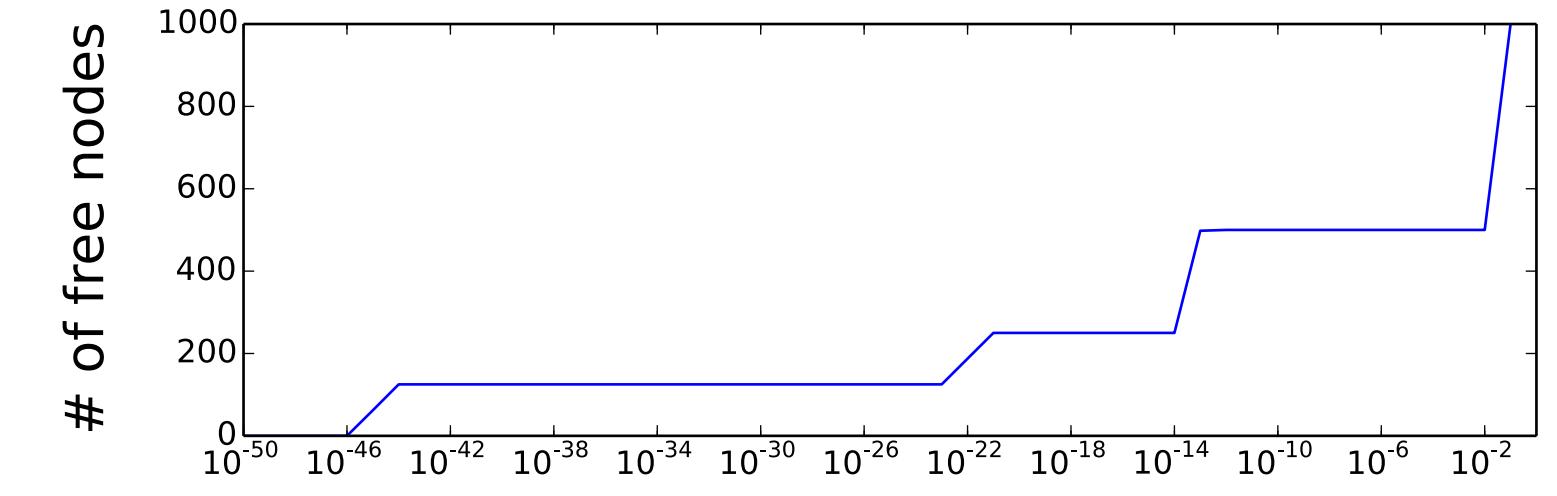
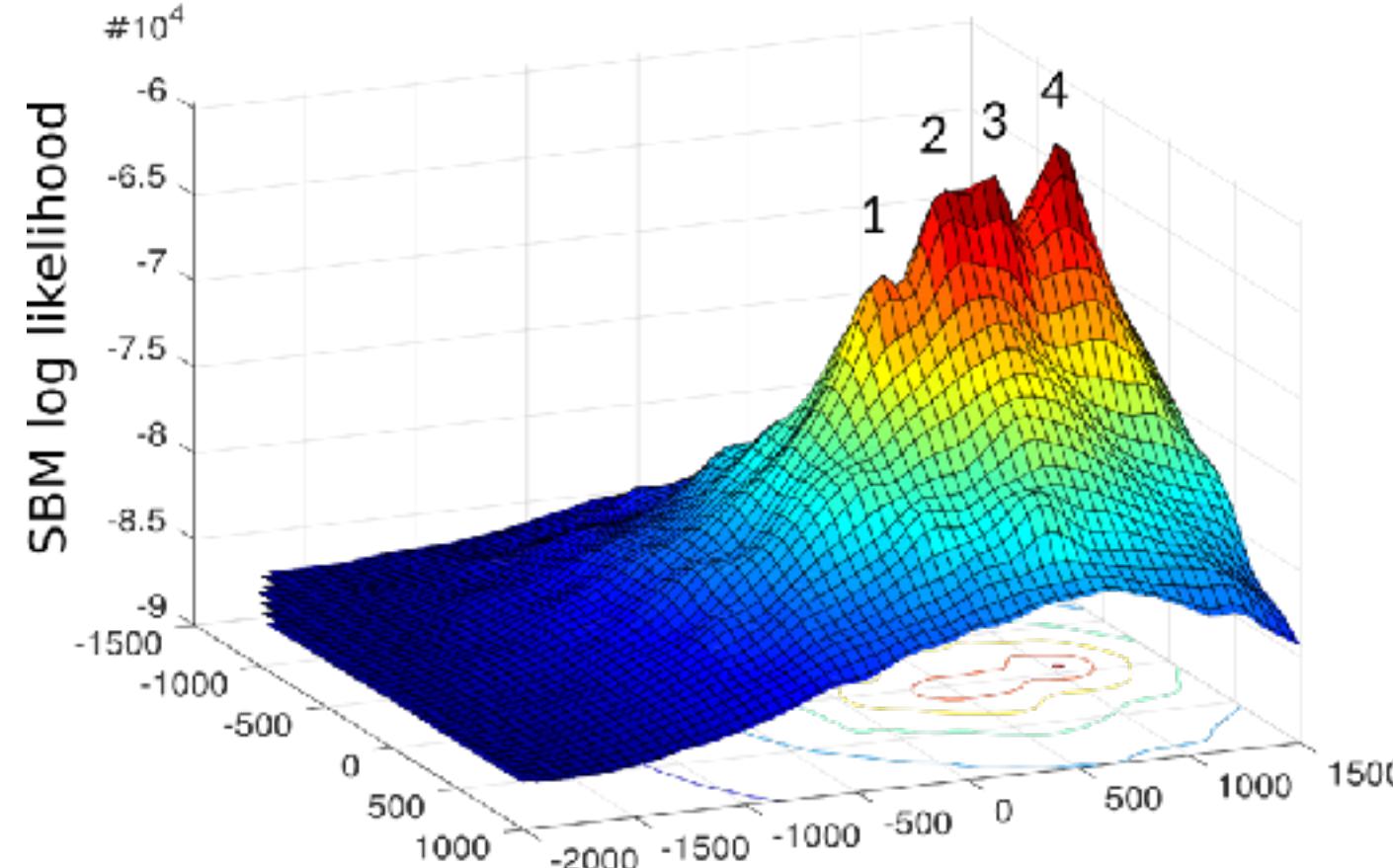


i. core-periphery

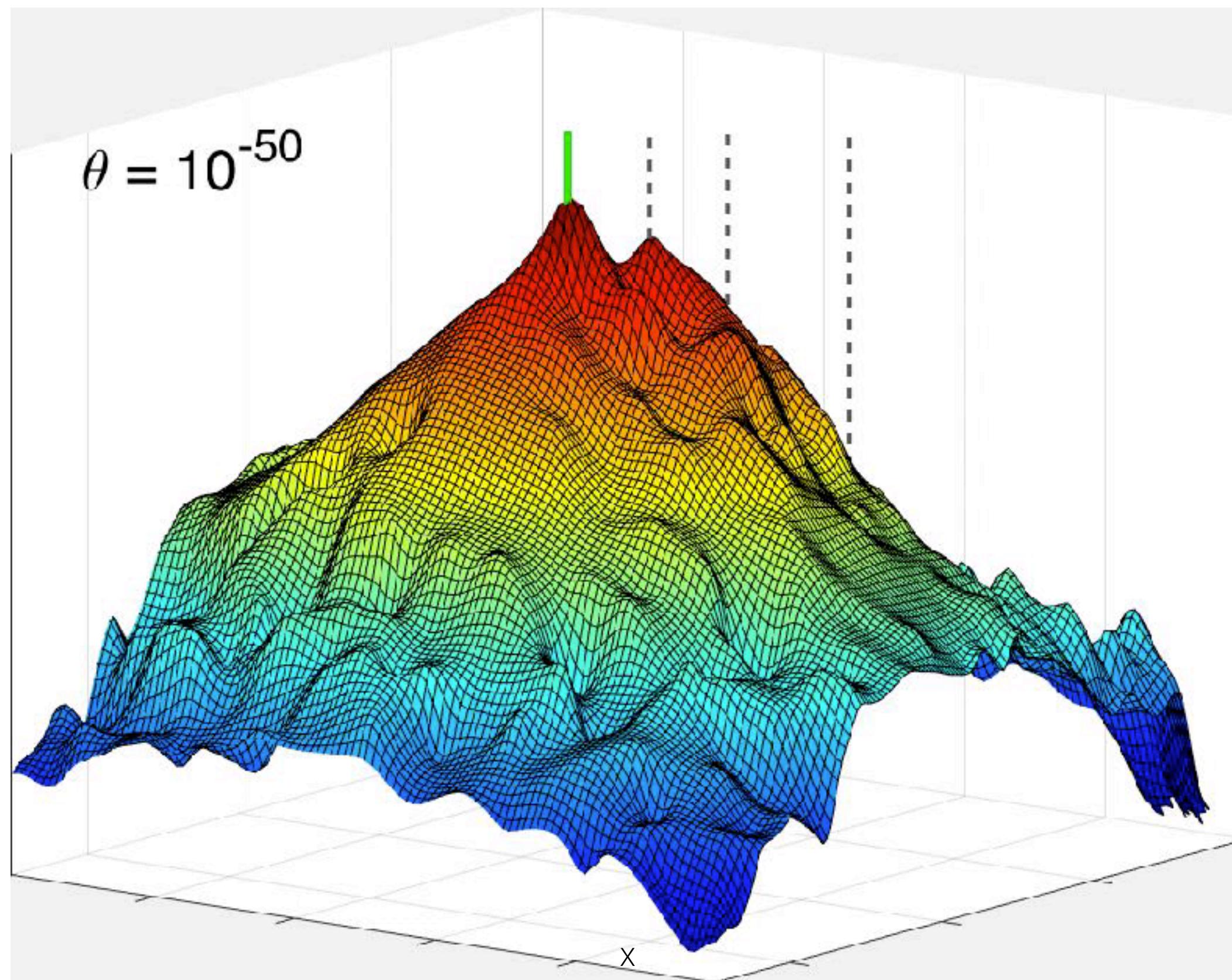


ii. assortative

# The neoSBM identifies four interesting partitions



# The prior parameter changes the likelihood surface



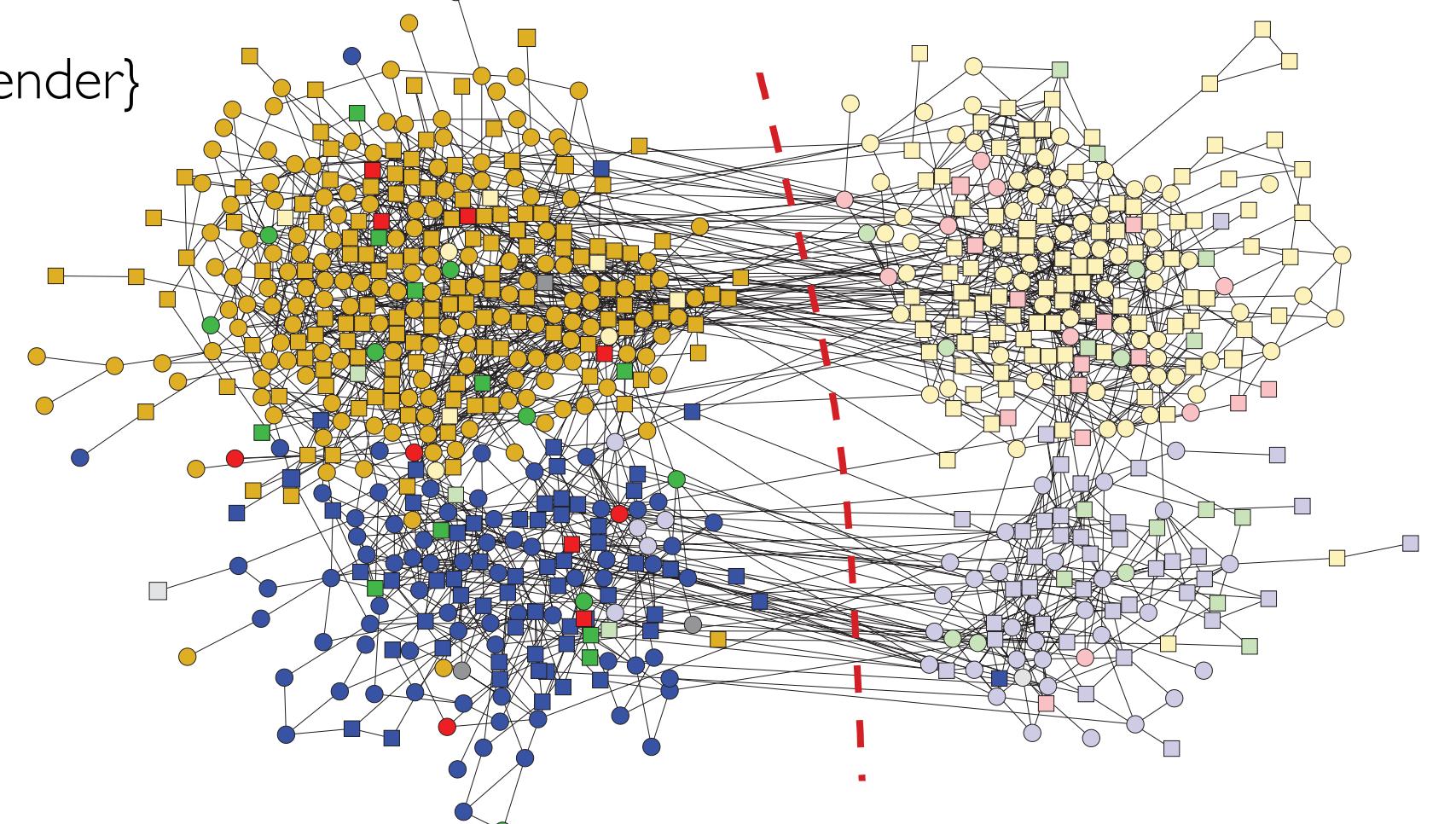
# Metadata-aware SBM

## high school social network



795 students from an American high school + its feeder middle school

- $x = \{\text{grade 7-12, ethnicity, gender}\}$



Male



Female

White

Middle

High

Black

High

Low

Hispanic

Low

High

Other

Low

High

Missing

Low

High

# Metadata-aware SBM

## high school social network



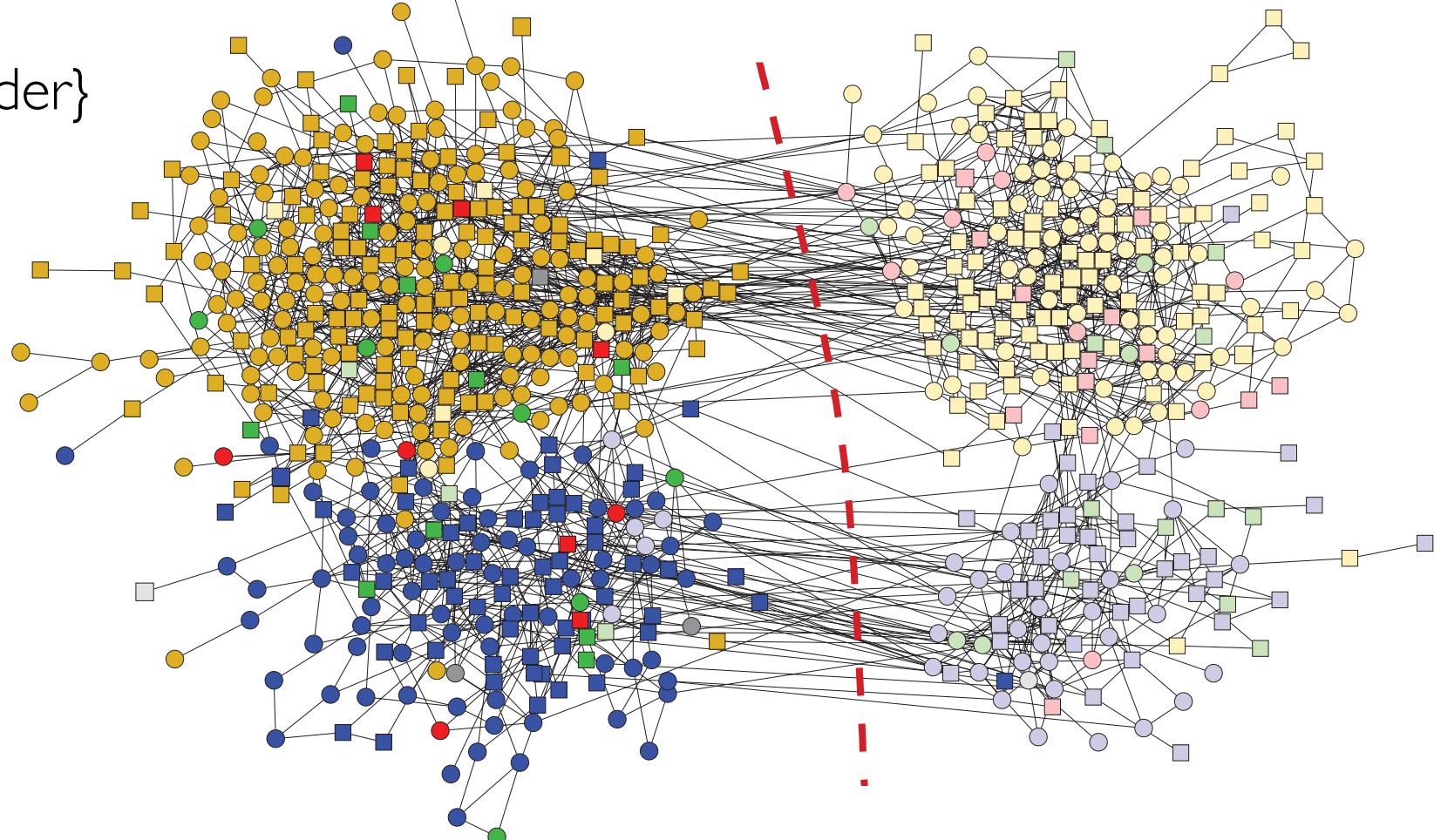
795 students from an American high school + its feeder middle school

- $\mathbf{x} = \{\text{grade 7-12, ethnicity, gender}\}$
- method finds a good partition between high-school and middle-school

$$\text{NMI} = 0.881$$

- without metadata:

$$\text{NMI} \in [0.105, 0.384]$$



	White	Black	Hispanic	Other	Missing	Male	Female
Middle	○	○	○	○	○	○	○
High	○	○	○	○	○	○	○

Newman & Clauset, *Nat. Comms.* 7, 11863 (2016)

Add Health network data, designed by Udry, Bearman & Harris

<http://www.santafe.edu/~aarong/>

# Metadata-aware SBM

## high school social network



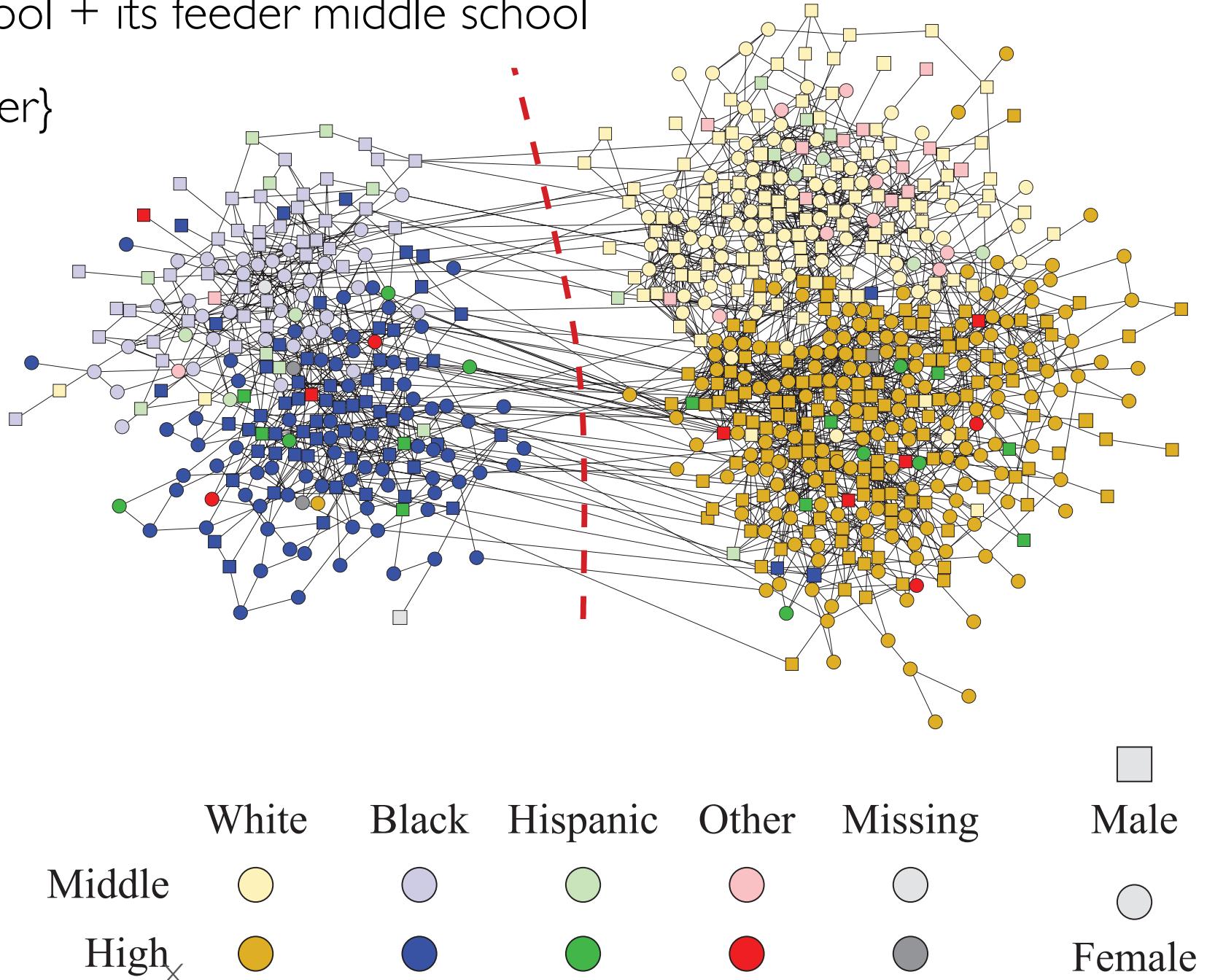
795 students from an American high school + its feeder middle school

- $\mathbf{x} = \{\text{grade 7-12, ethnicity, gender}\}$
- method finds a good partition between blacks and whites (with others scattered among)

NMI = 0.820

- without metadata:

NMI  $\in [0.120, 0.239]$



Newman & Clauset, *Nat. Comms.* 7, 11863 (2016)

Add Health network data, designed by Udry, Bearman & Harris

<http://www.santafe.edu/~aarong/>

# Metadata-aware SBM

## high school social network



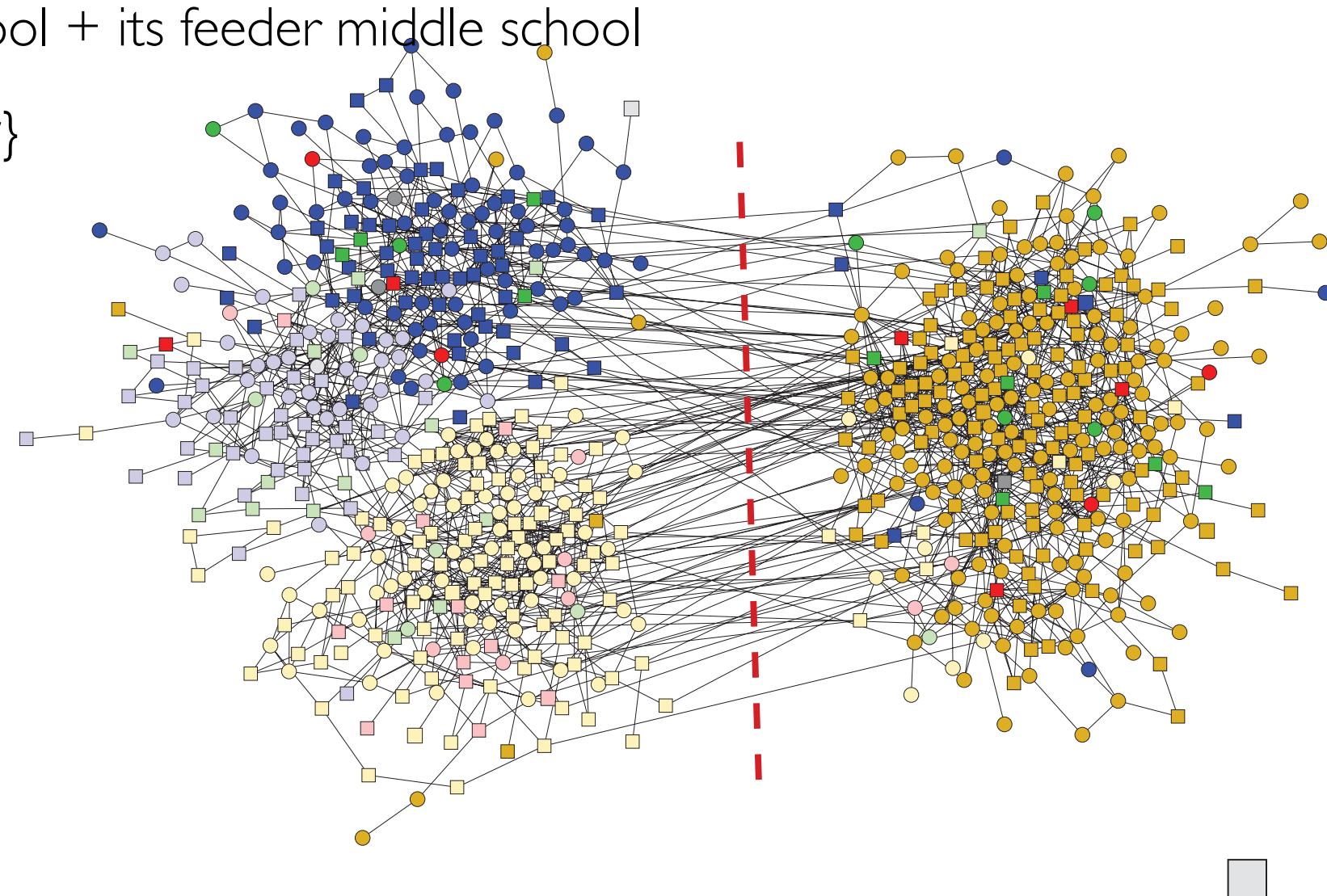
795 students from an American high school + its feeder middle school

- $\mathbf{x} = \{\text{grade 7-12, ethnicity, gender}\}$
- method finds no good partition between males/females.  
instead, chooses a mixture of grade/ethnicity partitions

$$\text{NMI} = 0.003$$

- without metadata:

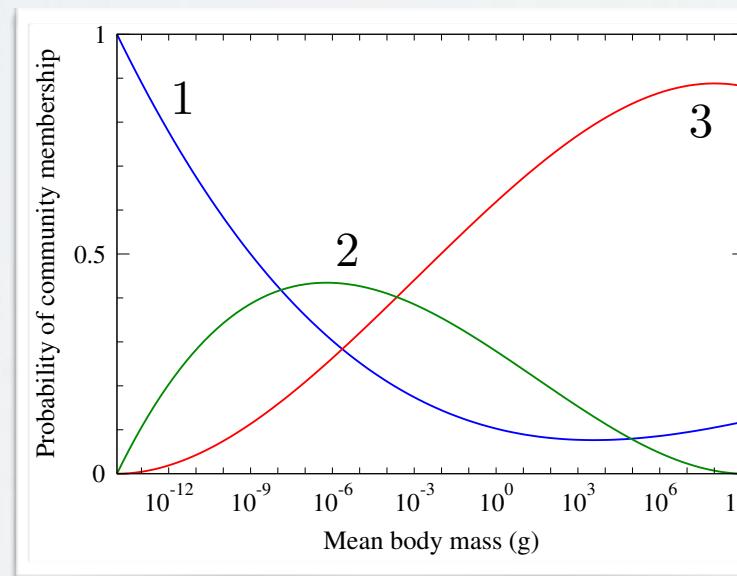
$$\text{NMI} \in [0.000, 0.010]$$



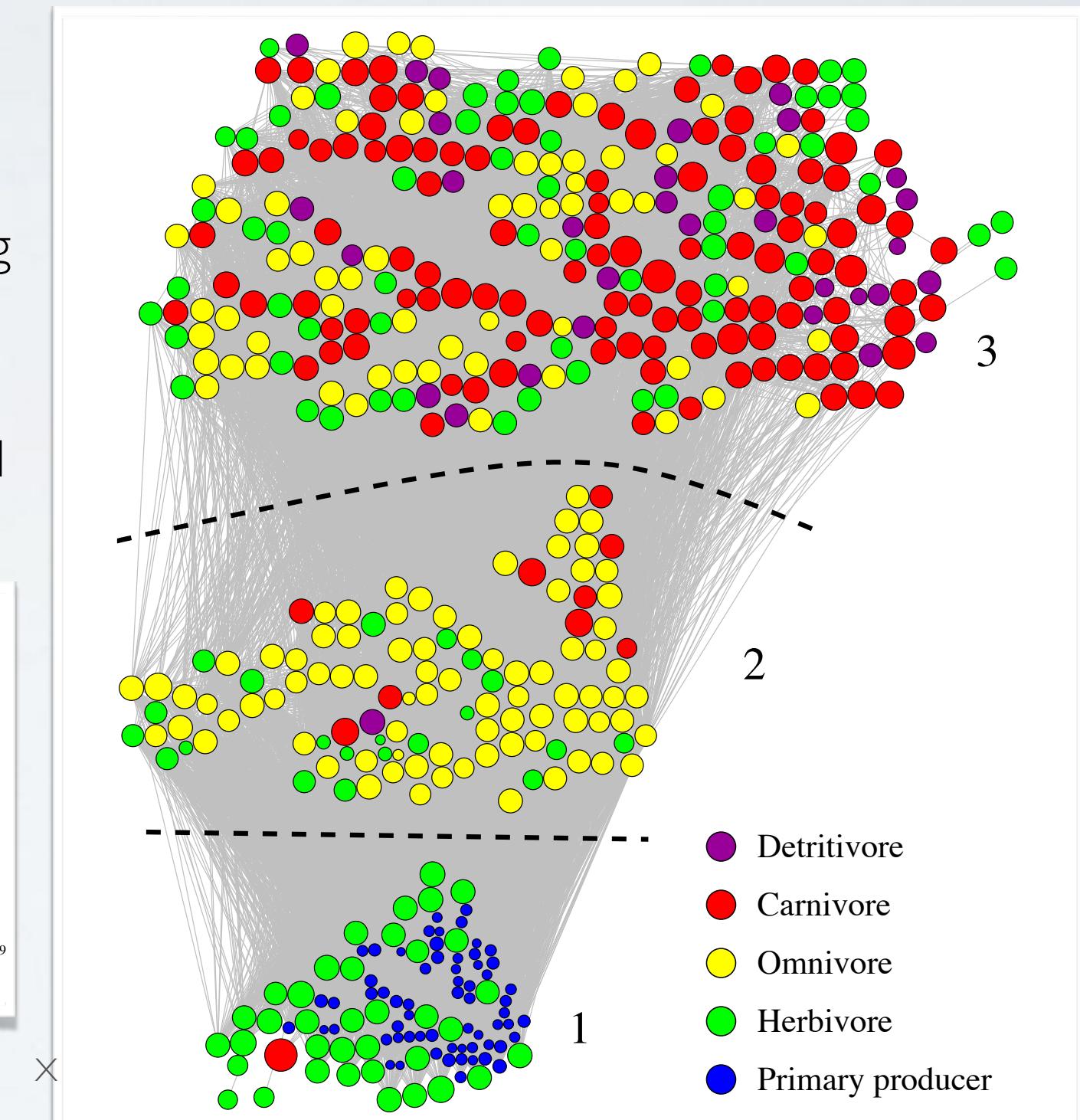
	White	Black	Hispanic	Other	Missing	Male	
Middle	○	○	○	○	○	○	
High	○	○	○	○	○	○	Female

# real-world networks

2. marine food web: predator-prey interactions among 488 species in Weddell Sea in Antarctica
- $x = \{\text{species body mass, feeding mode, oceanic zone}\}$
  - partition recovers known correlation between body mass, trophic level, and ecosystem role:



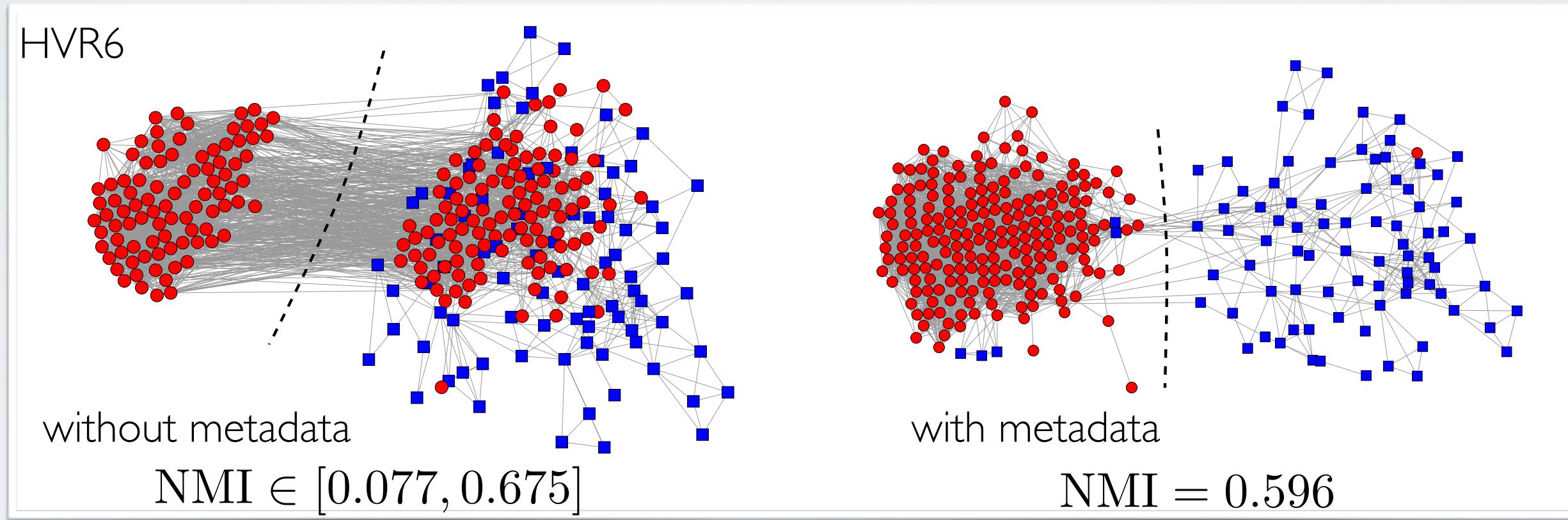
here, we're using a continuous metadata model  
Jacob et al. *Adv. Ecological Res.* 45 (2005)



# real-world networks

3. Malaria gene recombinations: recombination events among 297 var genes

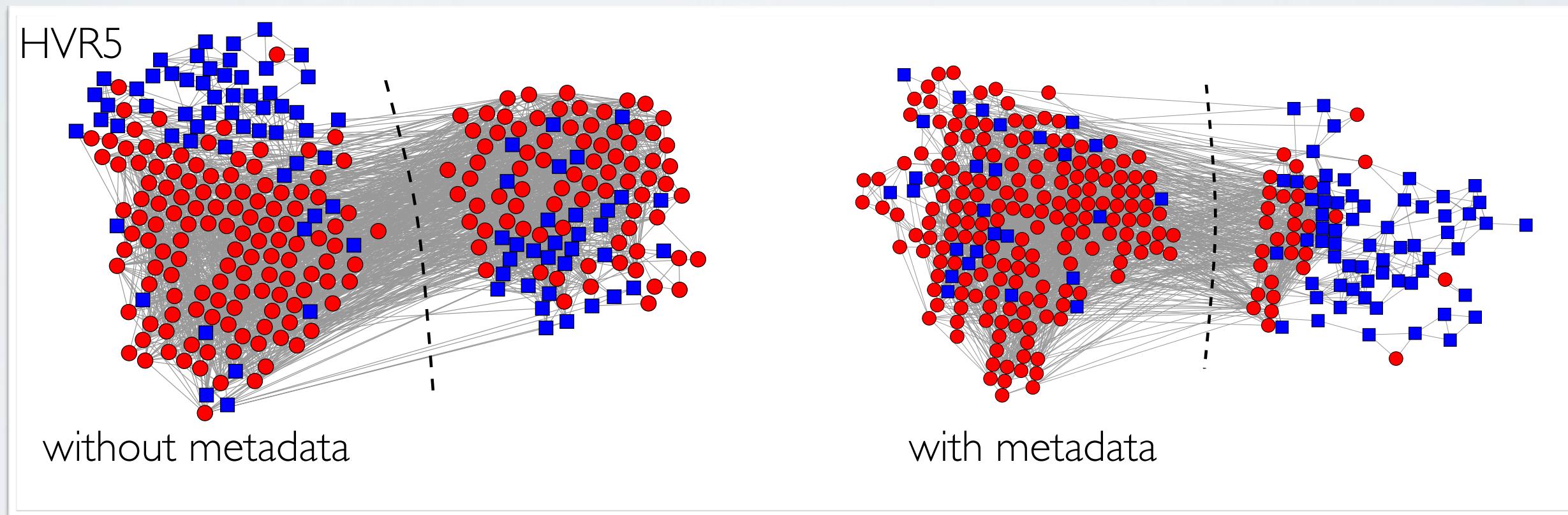
- $x = \{\text{Cys-PoLV labels for HVR6 region}\}$
- with metadata, partition discovers correlation with Cys labels (which are associated with severe disease)



# real-world networks

3. Malaria gene recombinations: recombination events among 297 var genes

- $x = \{\text{Cys-PoLV labels for HVR6 region}\}$
- on adjacent region of gene, we find Cys-PoLV labels correlate with recombinant structure here, too



# Other things to know about 1: “The Louvain Method”

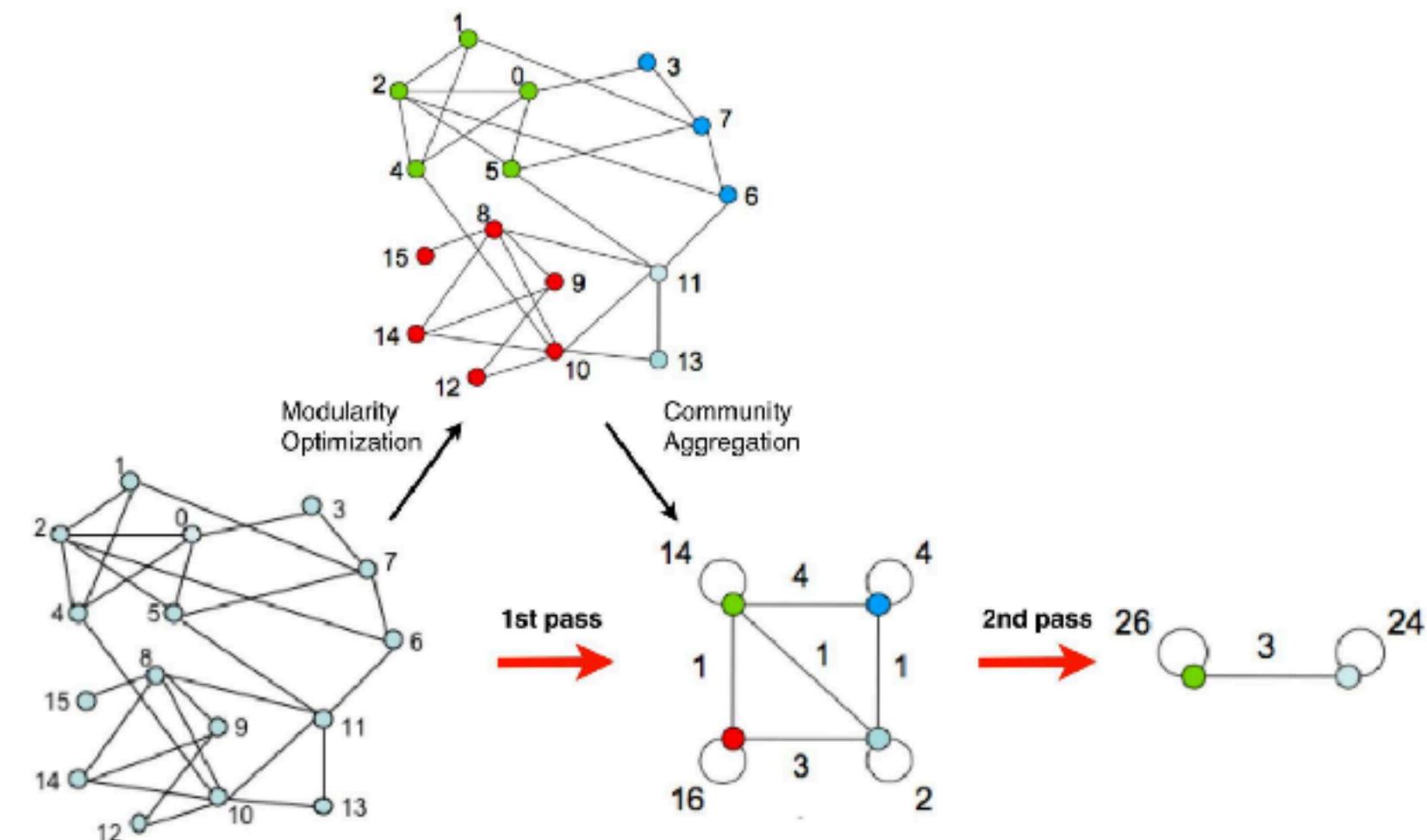
If your network is *really* big. (Millions of nodes, Billions of edges)

Take ClausetNewmanMoore’s approach for greedy Q maximization and find small groups. Run the code again on those groups... And again...

Advantage: fast! big! 

Disadvantage: inherits the assumptions of modularity.  
(clustering vs modeling)

6K citations. People like it!



# Other things to know about 2: InfoMap

Imagine a random walker on a network

A description of her walk can be compressed if the network has regions in which the random walker tends to stay for a long time.

Minimizing the “map equation” over all possible network partitions is the same as finding the best codebook.



<http://www.mapecuation.org/apps/MapDemo.html> 

<http://www.mapecuation.org/code.html>

# Outlook for community detection

**Simply put, we have amazingly powerful tools that did not exist 15 years ago.**

Many are principled, statistically rigorous, and we learn more all the time. Those that aren't statistically rigorous are really, really fast.

**There is no multiple regression for networks.**

“Controlling for C, how important is X in predicting Y?”

**Tradeoffs between general and bespoke methods are still being explored.**

Outside of SBM, Modularity, Louvain, Infomap, it's a wild west.

**Methodologists are keen to be challenged by new problem types.**

New scientific questions inspire new methods.



# The idea of rankings—pervasive!

Assumptions:

1. Competitors have some intrinsic quality (or vector of qualities).
2. Interactions can (stochastically) reveal differences in qualities.
3. Competitions are pair-wise. (Lee Sedol vs. AlphaGo; Astros vs. Dodgers)

In other words: outcomes are generated by a stochastic process, which is some function of the positions of the competitors.



# Systems of dominance

social



mental



physical



financial

**Sam Bennett vs Ryan Johansen**  
Feb 21, 2017 2pd 05:27  
2016-2017 Regular Season

Date / Time	Away / Home Team	Away / Home Player
Feb 21, 2017	Calgary Flames	Sam Bennett
2pd 05:27	Nashville Predators	Ryan Johansen

Your vote:  
You must sign in to vote.  
You can [sign up](#) for free if you do not have an account already.

**Results:**

Sam Bennett	92.9%
Ryan Johansen	6.4%
Draw	1.8%

From 61 votes with an average rating of 5.6

User Name:  User Name:  Remember Me?  
 Password:

Related Links:  
[Sam Bennett](#)  
2016-2017 Regular Season

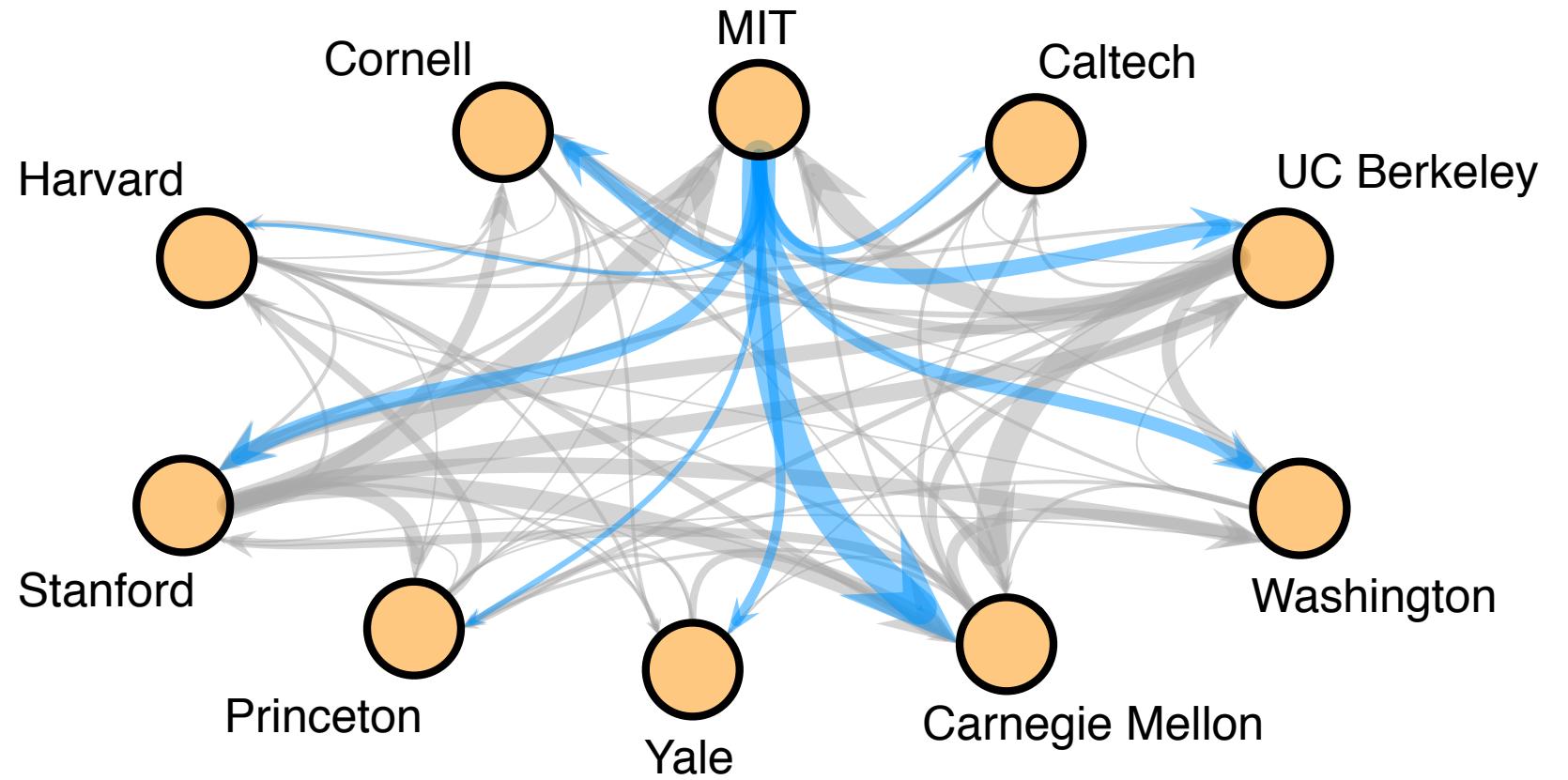
**hockeyfights.com** POWERED BY VIOLENT GENTLEMEN

Video

2nd 12:31 2017  
CDY 4 16  
NSH 1 12

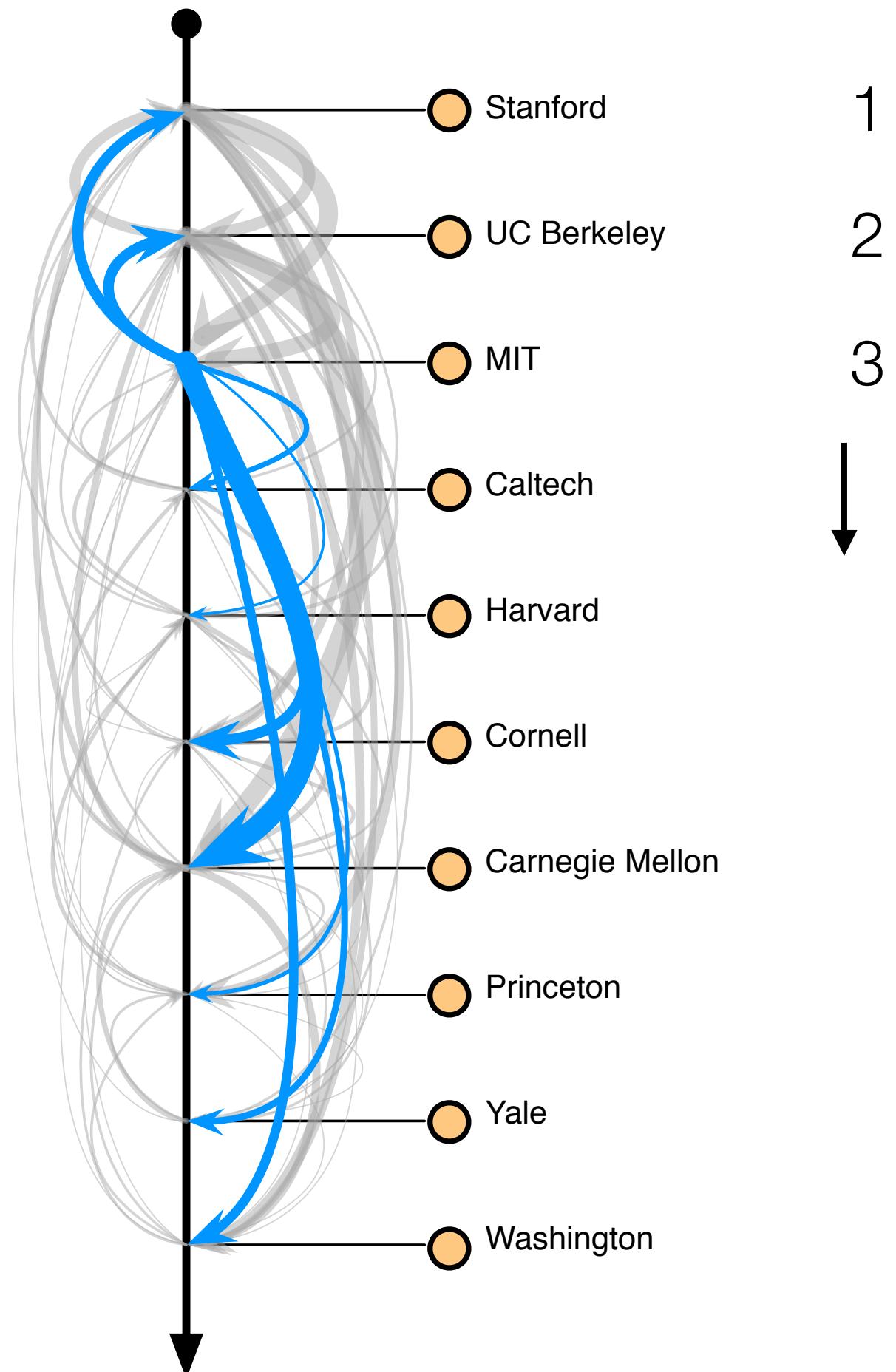
A screenshot of the hockeyfights.com website. It shows a poll for the fight between Sam Bennett and Ryan Johansen on February 21, 2017. The results show Sam Bennett with 92.9% of the votes. There is also a video player showing a fight from the game.

# Systems of endorsement

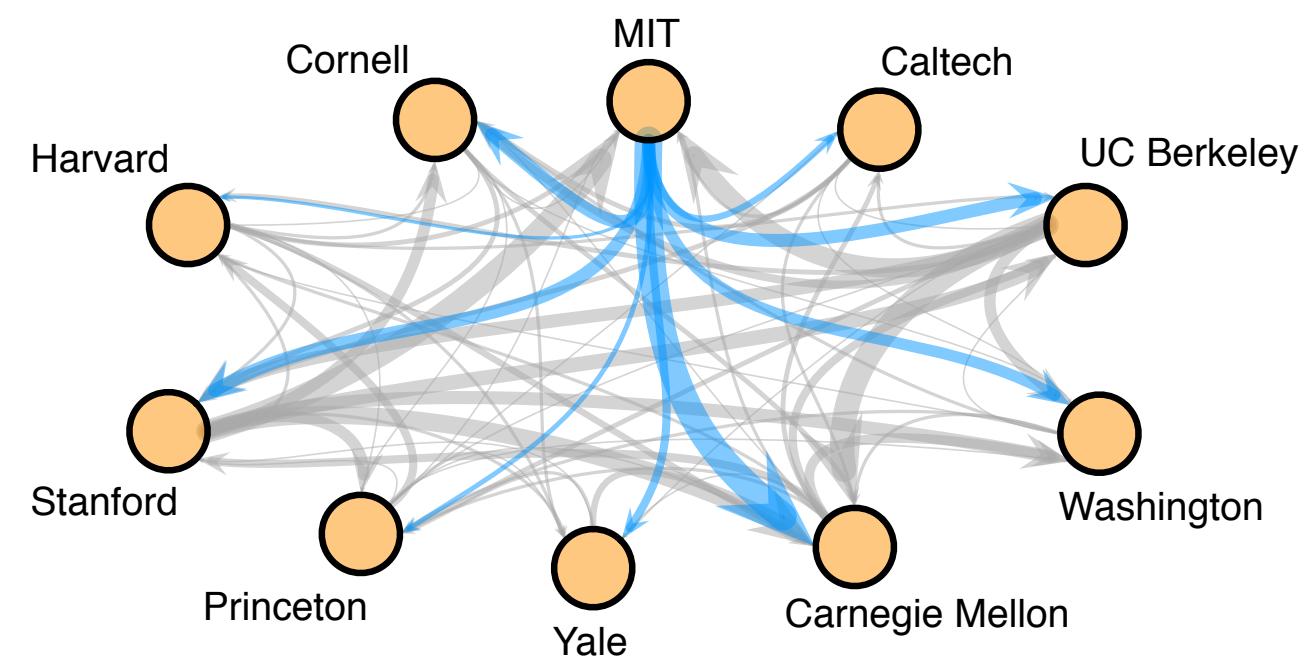


Assumptions:

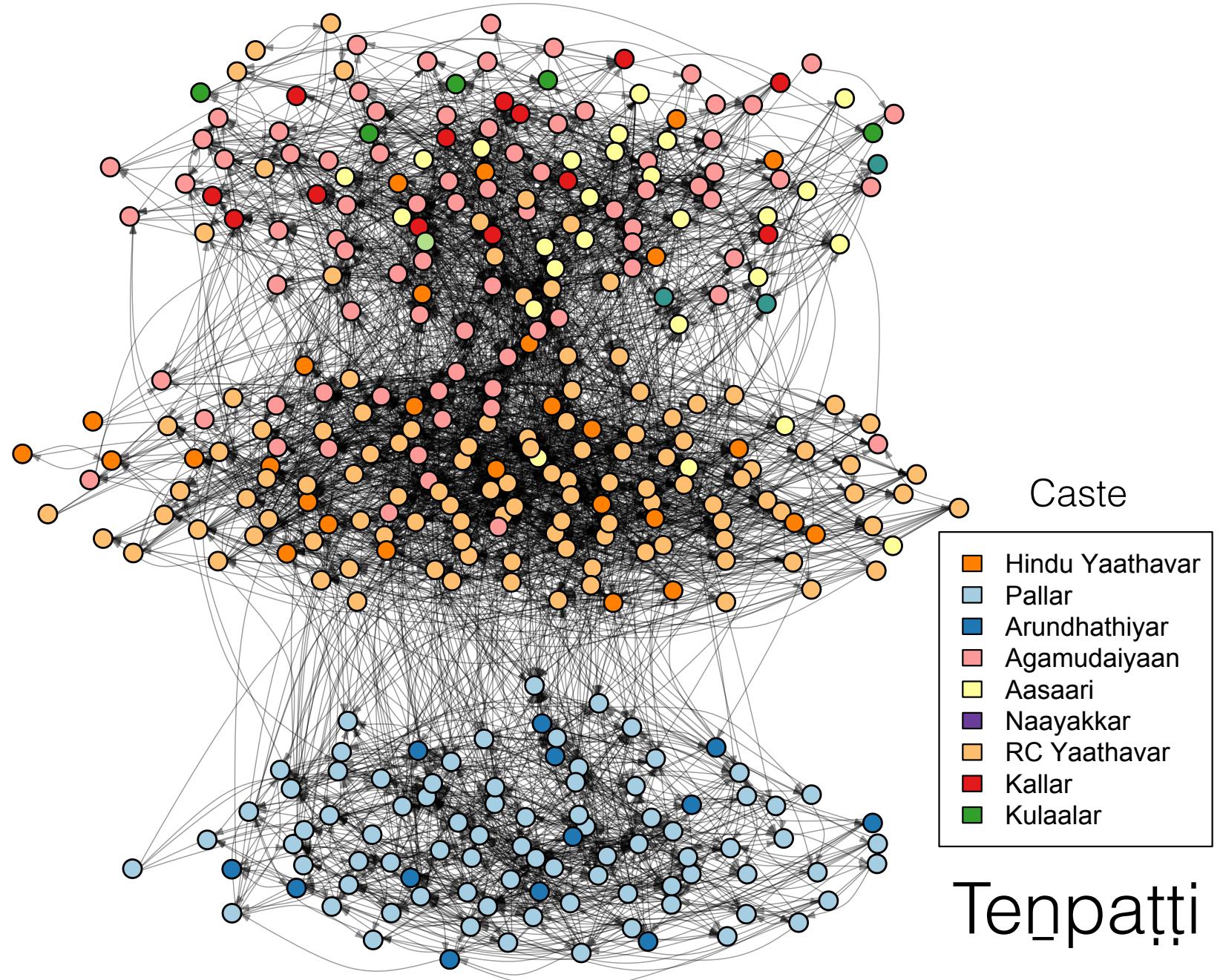
1. Endorsers have some intrinsic quality.
2. Interactions can reveal differences in qualities.
3. Endorsements are pair-wise.



# Systems of endorsement



Latent position can be revealed by dominance or endorsement interactions.



**The setup:** suppose we have a *directed* network.

Its adjacency matrix is  $A$ .

$A_{ij} = A_{i \rightarrow j}$  means  $i$  beat  $j$  or  $i$  was endorsed by  $j$

**The problem:** Rank the nodes.

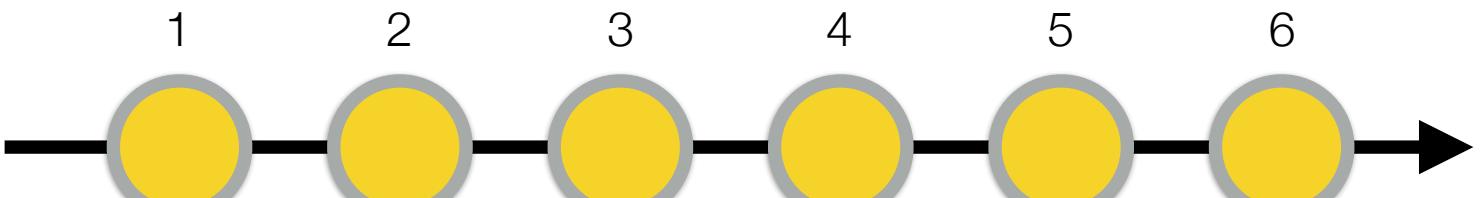
**Alternative view:** there might be no network here. In some cases we're just seeing a network in pairwise comparison data because networks are a convenient data structure.

**Alternative problem:** Which items should be compared next in order to most/best resolve our estimate of the ranks? (sequential tournament design)

# Embeddings vs Orderings

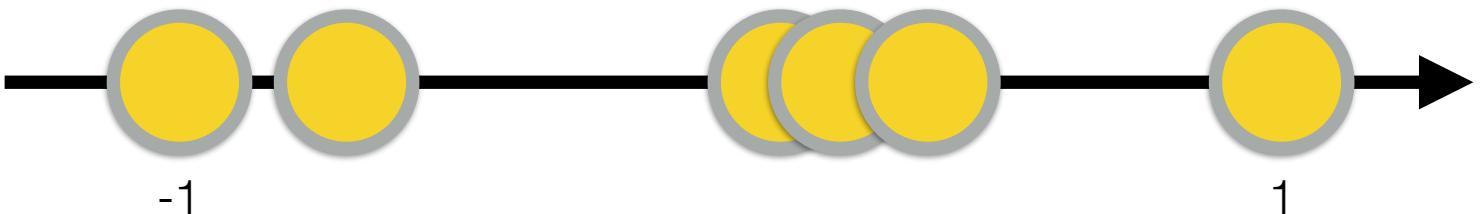
**Ordering** place the nodes in order:

1, 2, 3, ...



**Embedding** assigns a position to each node:

1, 1.2, 7, 20, 21, 21.2, ...



## Which one should I use?

- > Depends on the use case.
- > Is it possible for two nodes to occupy the same rank or position? If so, an embedding is more appropriate. Also better when meaning of 1-rank  $\Delta$  varies.
- > Consider that you can always go from an embedding to an ordering, if you have a rule for breaking ties.

# Win-Loss is not satisfactory: schedule matters

Beating the grandmaster counts for more than beating a novice.

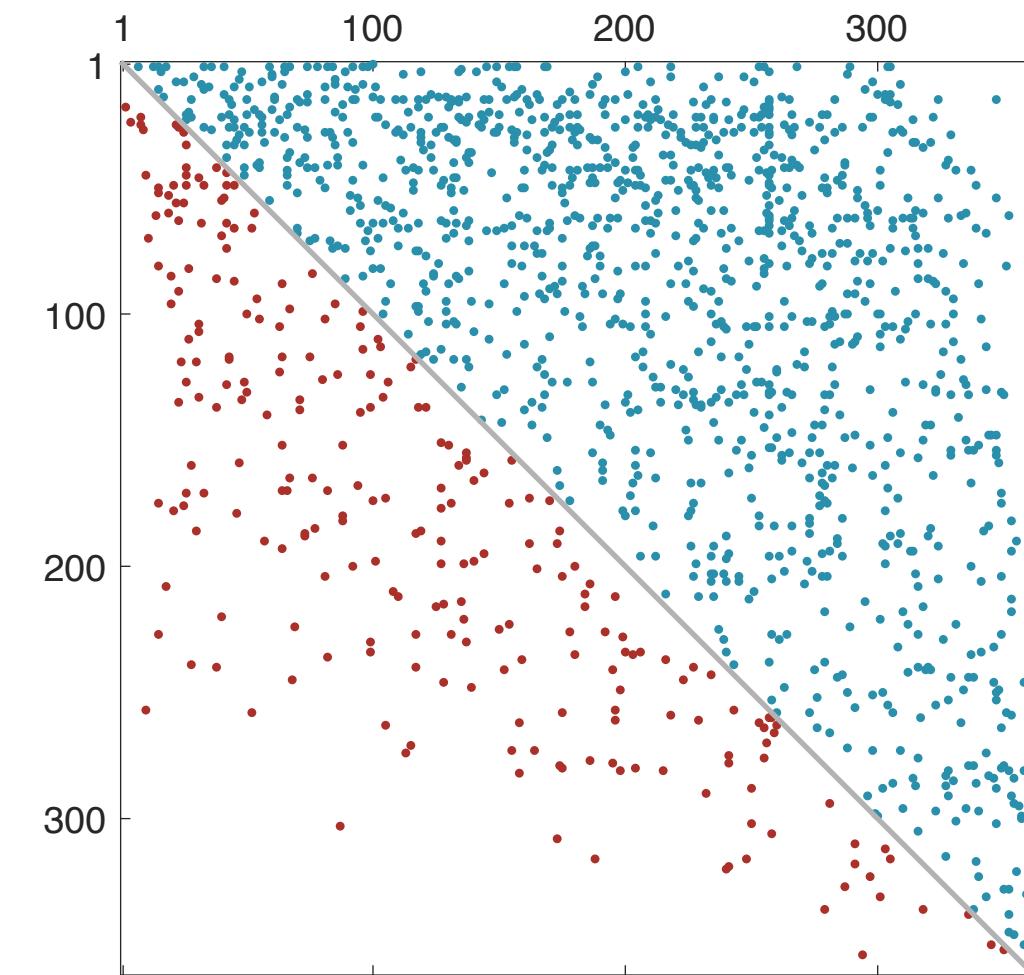
Win and loss tallies don't take this "schedule difficulty" into account. Put differently, win-loss records leave information on the table.

One way to make use of this information:

$i$  beats  $j$  implies  $s_i > s_j$

Therefore if we have a whole list of outcomes, we can try to find a total ordering that breaks as few of these implications as possible.

$A_{ij}$  = number of times that  $i$  beat  $j$ .



minimum violation ranking: sort  $A$ .

# Win-Loss is not satisfactory: schedule matters

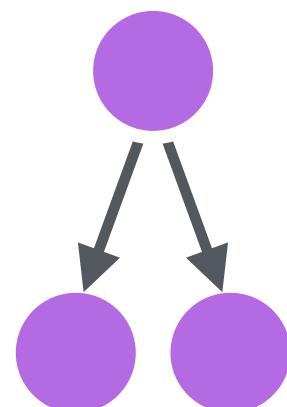
How do we find an ordering that minimizes the number of violations (or upsets) ?

## Recipe (MCMC):

1. Order the nodes randomly.
2. Compute the number of violations. In expectation, this should be 50% of edges.
3. Pick two nodes at random and propose to swap their positions.
4. Compute the number of violations in this scenario.
5. If #violations decreases or stays the same, keep the swap. Otherwise, reject.
6. Repeat until....?

## Notes:

- \* The number of violations is non-increasing over time.
- \* There may be no unique minimum. Consider this scenario:



# Embeddings & Orderings 0: MVR & Agony

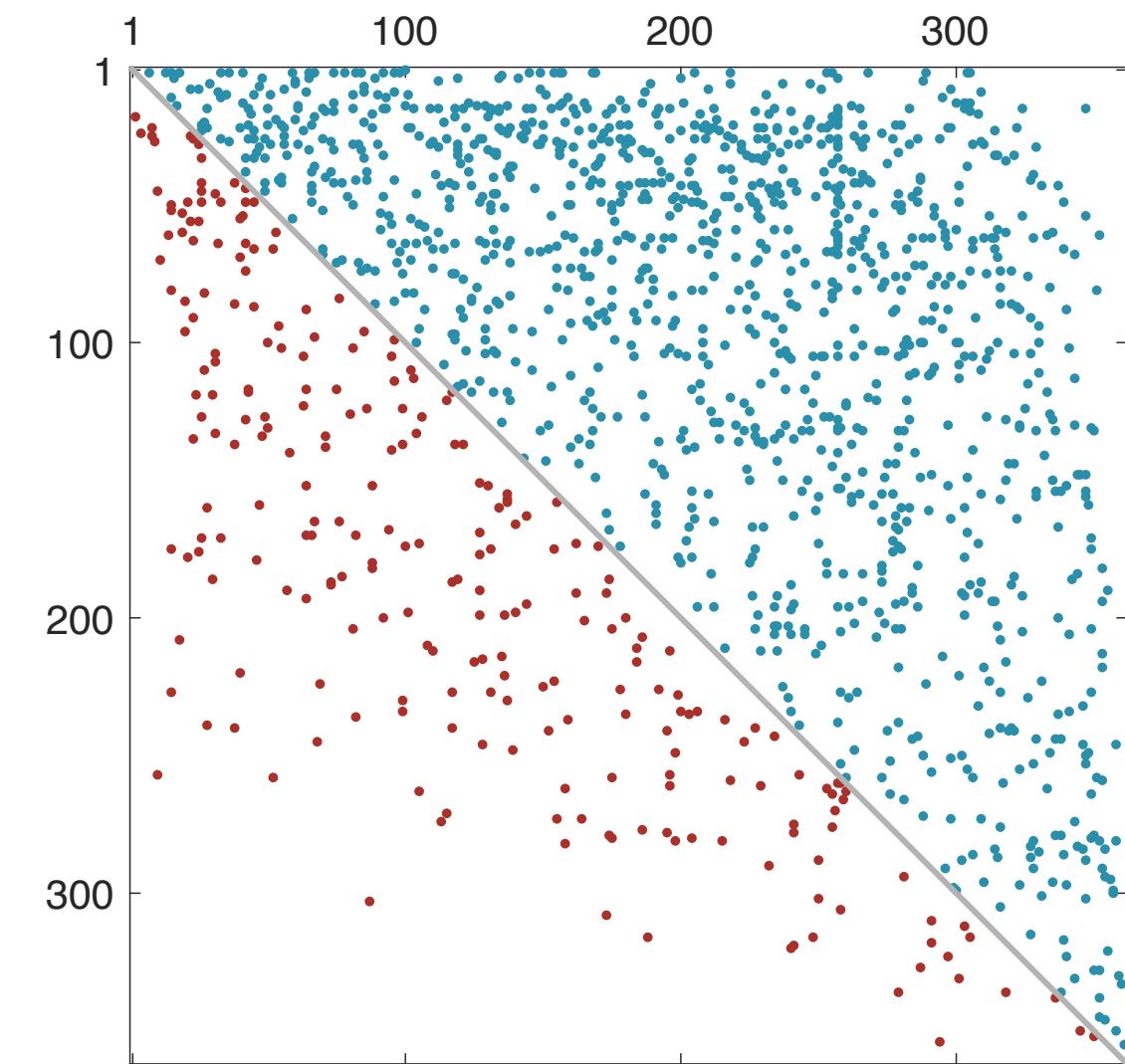
- \* There is no guarantee of a unique minimizing ranking  $s$ .
- \* Space of ordinal rankings has  $n!$  elements, requiring slow search algorithms (e.g. MCMC).
- \* Ordinal. No ties. No interpretability of rank differences.

What if you allowed for **ties** and then ran Minimum Violation Ranking (MVR)? What would happen?

**MVR**: uniform cost (1 per edge).

**Agony**: generic cost function.  
for example, difference in ranks.

What are other premises on which we can base a ranking model?



minimum violation ranking: sort  $A$ .

# Embeddings and Orderings 1: Discrete choice models

Louis Leon Thurstone and Thelma Thurstone

Thelma: Prof. of Education & Psych UNC Chapel Hill. Louis: Worked with Edison.



# Embeddings and Orderings 1: Discrete choice models



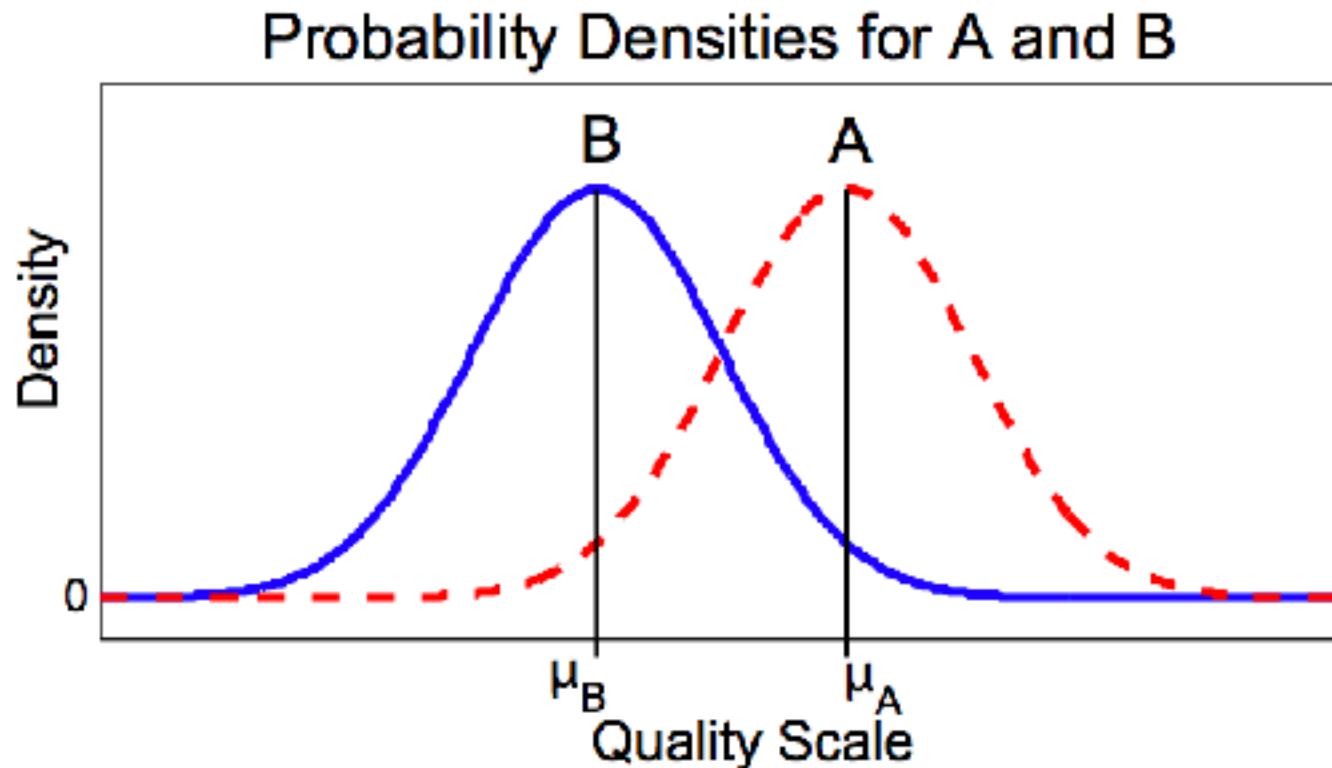
Instead of rating everything from 1 to 10, try *paired comparisons*.

Do you prefer  $i$  or  $j$  ?

Why? Consider: My 3 is not your 3. What is 1 and what is 10?

# Embeddings and Orderings 1: Discrete choice models

**Thurstone:** items have quality distributions. When a person judges whether A is better than B they draw from A's distribution and from B's distribution and see which is higher.



Thurstone modeled these as Gaussians.

$$P(A > B) = P(A - B > 0)$$

Difference of Gaussians is Gaussian.

$$\hat{\mu}_{AB} = \Phi^{-1} \left( \frac{C_{A \rightarrow B}}{C_{A \rightarrow B} + C_{B \rightarrow A}} \right)$$

Where  $\Phi^{-1}(x)$  is the inverse CDF of standard normal, a.k.a. the *probit*.

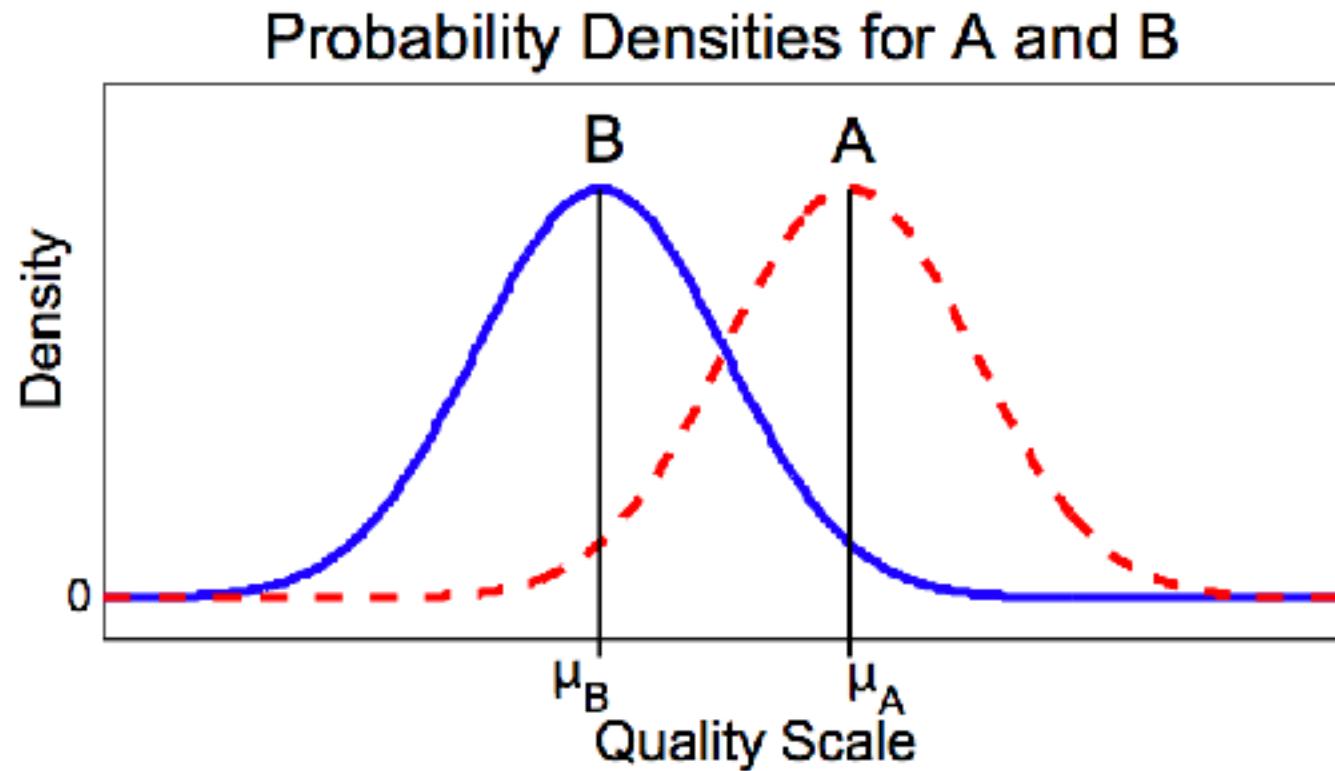
**Powerful idea:** lots of pairwise comparisons = estimates of all the qualities! An embedding!

**Key:** pairwise comparisons = directed network.       $i$  preferred to  $j$  =  $i \rightarrow j$

Finding the qualities of items from pairwise comparisons = Finding embedding of nodes.

# Embeddings and Orderings 1: Discrete choice models

**Bradley-Terry & Luce:** items have quality distributions. When a person judges whether A is better than B they draw from A's and from B's distribution and see which is higher.



BTL

$$P(A > B) = \frac{\pi_A}{\pi_A + \pi_B}$$

Or usually:

$$P(A > B) = \frac{e^{\mu_A/s}}{e^{\mu_A/s} + e^{\mu_B/s}} = \frac{1}{1 + e^{-(\mu_A - \mu_B)/s}}$$

Same idea; different distribution. (*logit* instead of *probit*; *Gumbel* instead of *Gaussian*)

**Powerful idea:** lots of pairwise comparisons = estimates of all the qualities! An embedding!

# BTL avoids non-transitivities (aka rock-paper-scissors)

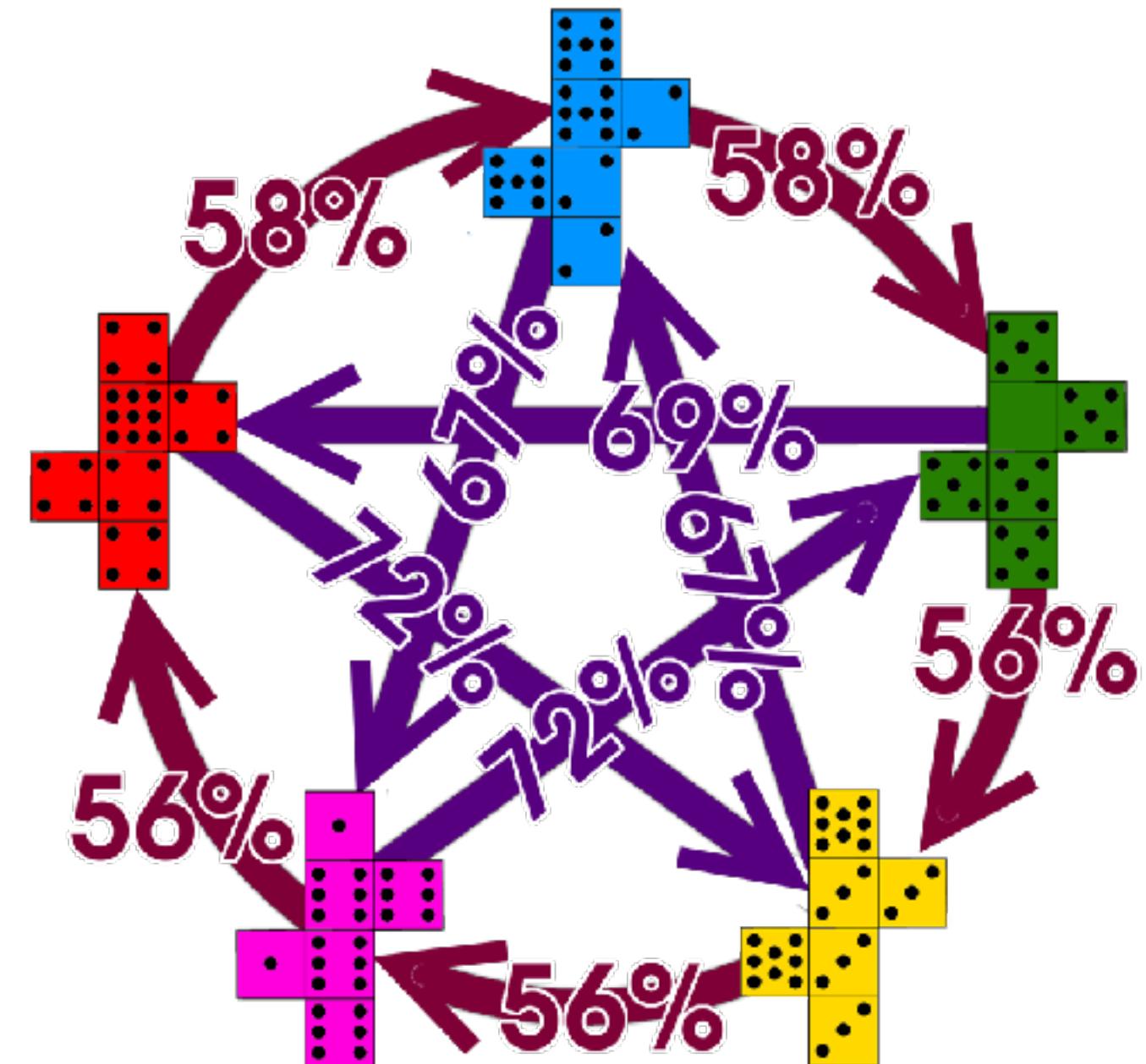
Introducing: **non-transitive dice!**

- 3 (or more) dice {A,B,C}
- faces chosen so that they have the property:
  - A>B more than half the time.
  - B>C more than half the time.
  - C>A more than half the time (?!)

[https://en.wikipedia.org/wiki/Nontransitive\\_dice](https://en.wikipedia.org/wiki/Nontransitive_dice)

A great gift for your favorite nerd's desk!

Go to the makerspace and laserbeam your own!



# Bradley-Terry-Luce

These methods embed items or players in a 1D space.

- Provably avoids non-transitive properties
- Great when lots of data per interaction.

Pairwise ranking is really nice for ordering large sets of preferences too, and this model specifically models the probability that the preference will be for  $i$  over  $j$ .

Iterative algorithms exist. Needs a little regularization so the winningest winners don't fly off to infinity. [why?]

$$P(A \rightarrow B) = \frac{\pi_A}{\pi_A + \pi_B} = \frac{e^{\gamma_A}}{e^{\gamma_A} + e^{\gamma_B}}$$

# Embeddings and Orderings 1: Discrete choice models

**Introductory tutorial (Gupta):**

<http://mayagupta.org/publications/PairedComparisonTutorialTsukidaGupta.pdf>

**Discrete choice today (Ugander):**

<https://web.stanford.edu/~jugander/papers/nips16-pcmc-slides.pdf>

**The textbook (Train):**

<https://eml.berkeley.edu/books/train1201.pdf>

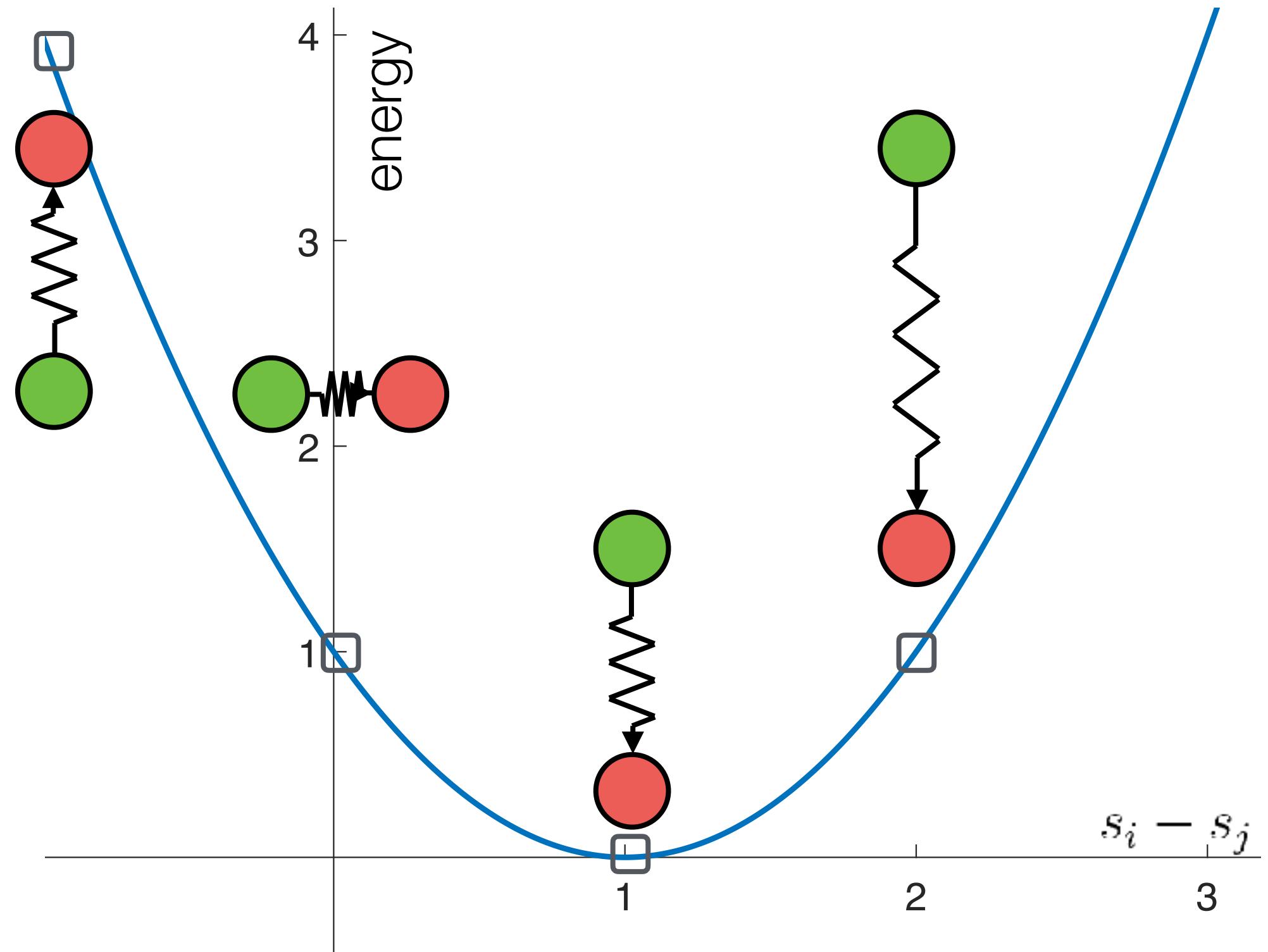
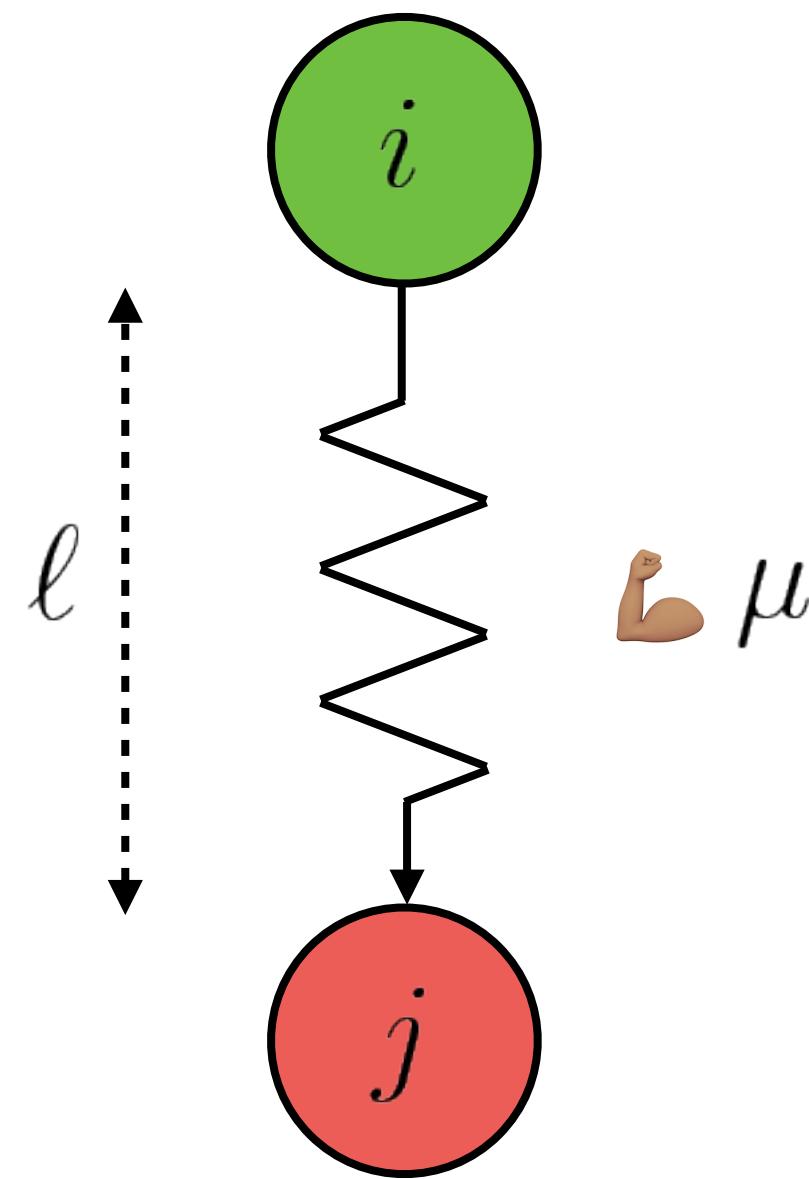
**Nobel Lecture for roadtrip when you're out of podcasts (McFadden)**

<https://www.nobelprize.org/prizes/economic-sciences/2000/mcfadden/facts/>



Embeddings & Orderings 2: SpringRank

Each directed edge = directed spring



# How much energy is this system of springs?

SpringRank Hamiltonian = energy of the system, given the node positions  $s$ .

# Relax and let the springs decide the ranks

$$H(s) = \frac{1}{2} \sum_{i,j=1}^N A_{ij} (s_i - s_j - 1)^2$$

SpringRank Hamiltonian = energy of the system, given the node positions  $s$ .

Because the springs are linear, the potential is quadratic.

The SR Hamiltonian is *convex* in  $s$ .

$$\nabla H(s) = 0$$

The solution is unique...up to an additive constant. (Why?)

# Derivatives work out nicely

$$0 = \frac{\partial H}{\partial s_i} = \sum_j A_{ij}(s_i - s_j - 1) - A_{ji}(s_j - s_i - 1)$$

Rewrite as a linear algebra problem.

$$[D^{\text{out}} + D^{\text{in}} - (A + A^T)] s^* = [D^{\text{out}} - D^{\text{in}}] \mathbf{1}$$

We know *a priori* that the matrix on the left is singular: translational invariance of  $H(s)$ .  
[if  $s$  is a solution, then  $s + k$  is a solution for any constant  $k$ ; eigenvalue 0, eigenvector  $\mathbf{1}$ ]

Notice: the matrix on the left is the *graph Laplacian* of the *undirected* network.

Uniqueness: Set  $s_1=0$ ,  $\min(s)=0$ , or  $\text{mean}(s)=0$ . Or use a pseudoinverse. Or regularize.

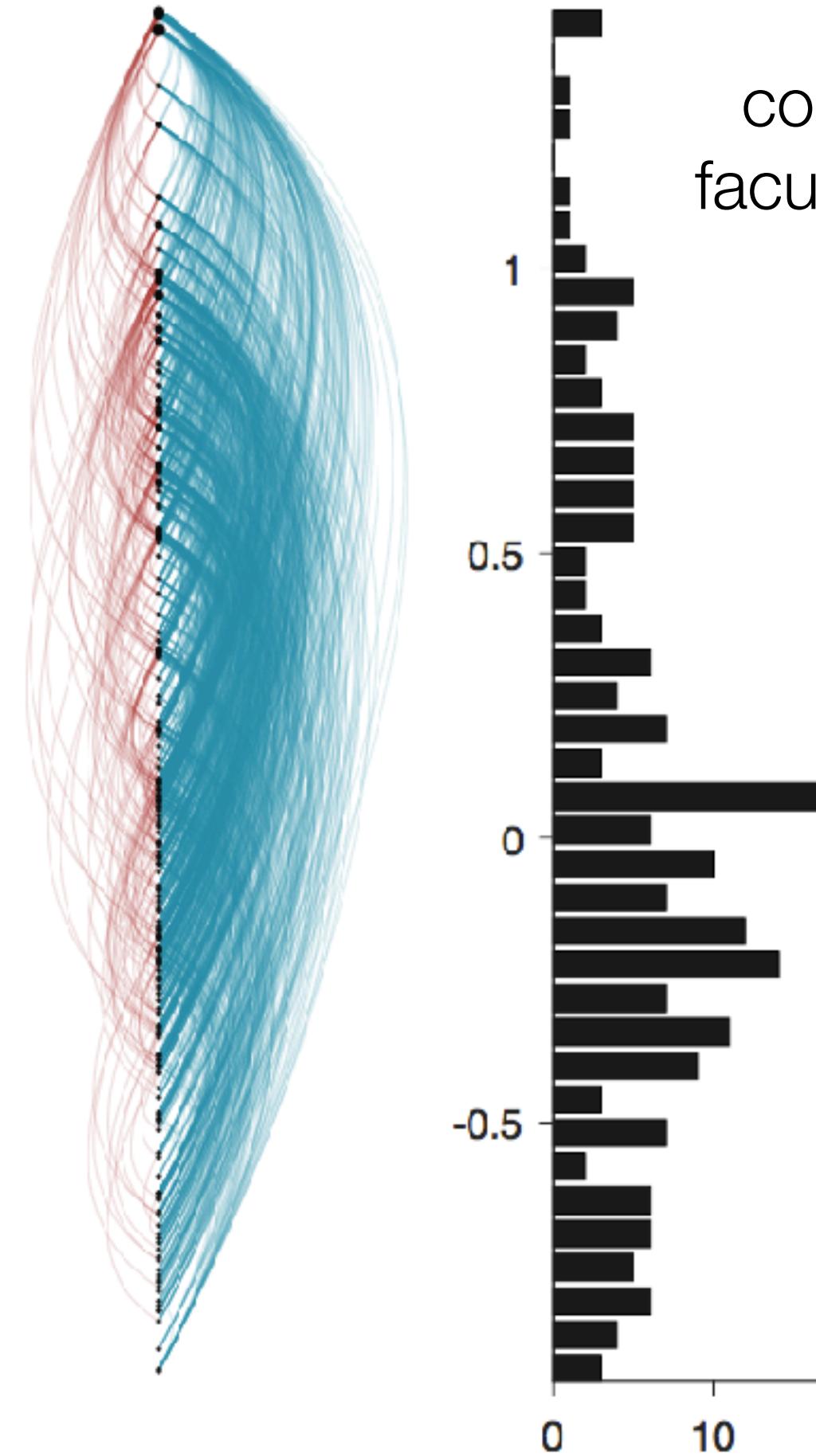
# It works!

Real networks tend to be sparse...  
our linear algebra problem is sparse...  
we can use sparse iterative solvers...  
**millions of edges in seconds.**

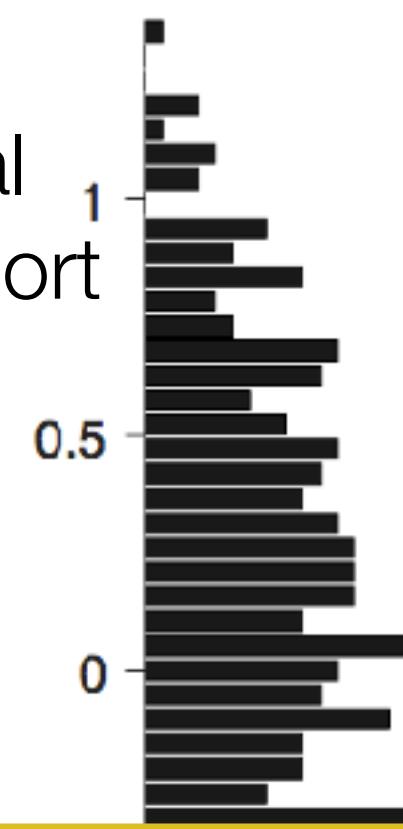
Even better: it's a linear-Laplacian system.

🚀 Near-linear-time (in  $|edges|$ ) solutions.

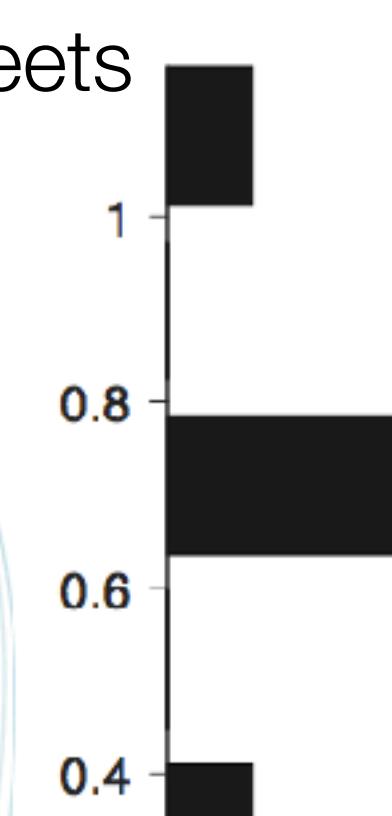
Note that node positions can be clumpy,  
since this is an *embedding*.



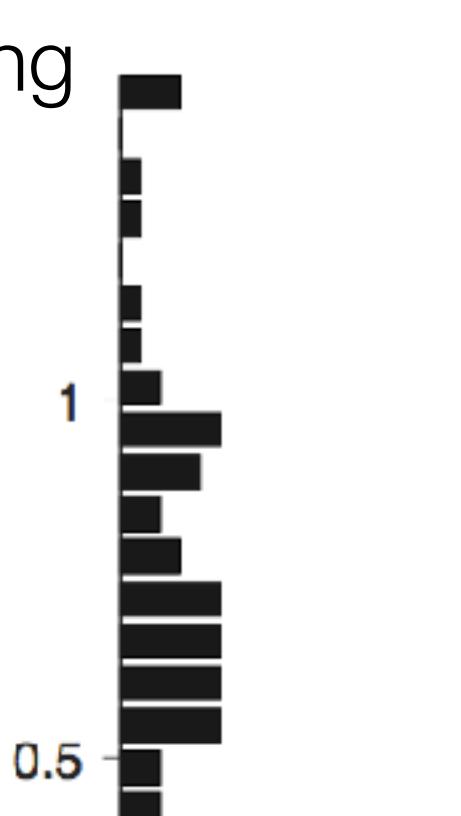
rural  
social  
support



parakeets

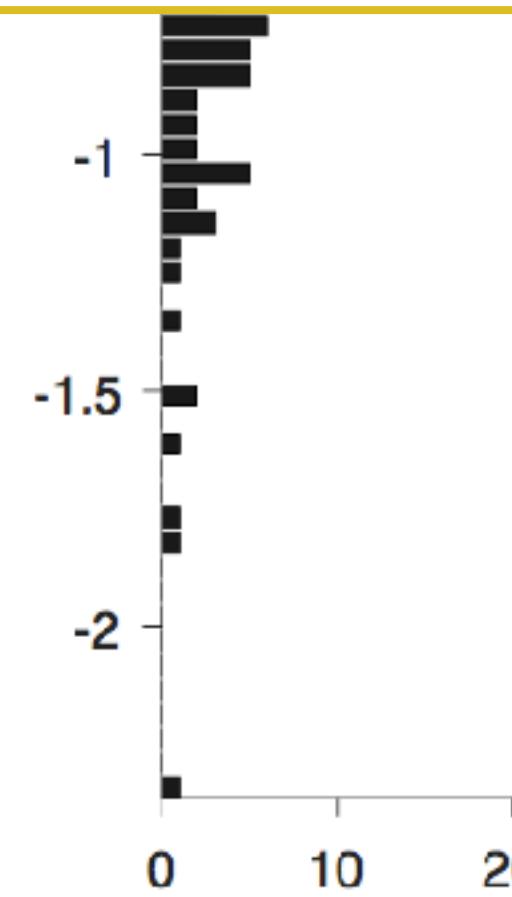


CS hiring

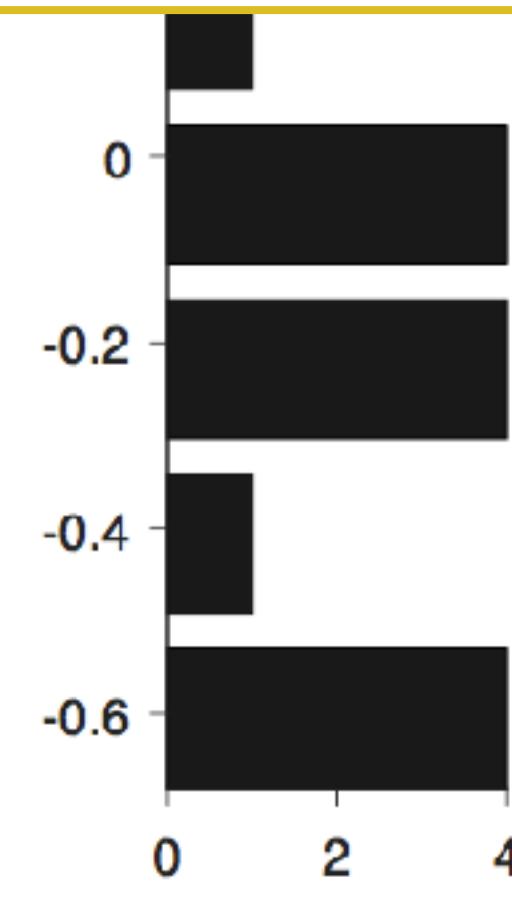
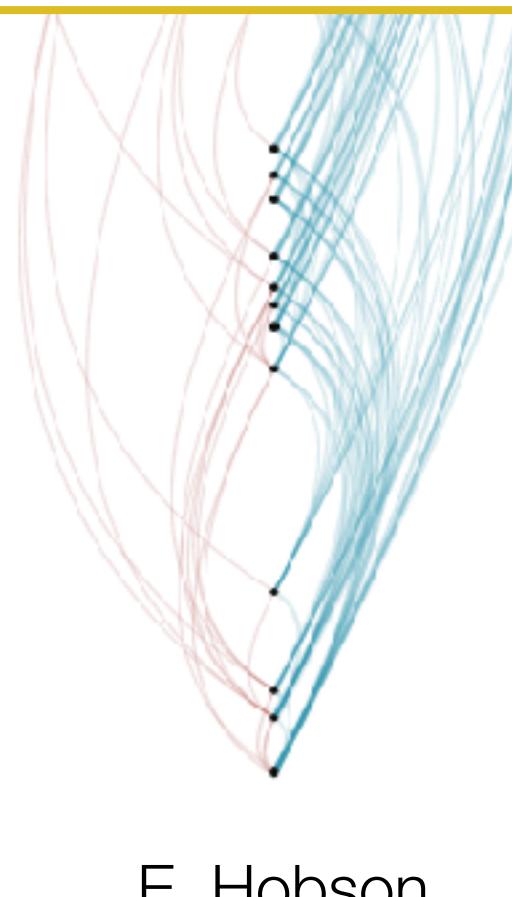


risk: stopping at “ours is faster” + pretty pictures

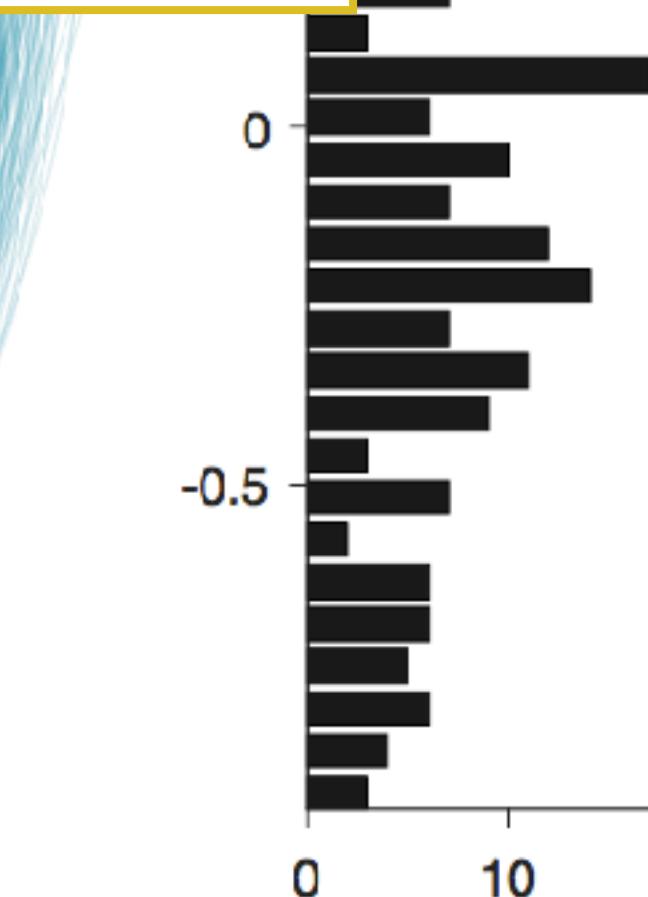
E. Power



E. Hobson



A. Clauset



# Cross validation: train on 80%, predict 20%

In a linear hierarchy the key quantity to predict is *edge direction*, given *edge existence*.

If  $i$  and  $j$  were to face off, who would win?

I'll give you *undirected*( $A$ ), and you predict *directed*( $A$ ).

**Setup:** learn  $s$  from 80% of  $A$ . Then predict edge directions for remaining 20% of  $A$ .

SpringRank predicts edge direction based on the relative direction probabilities:

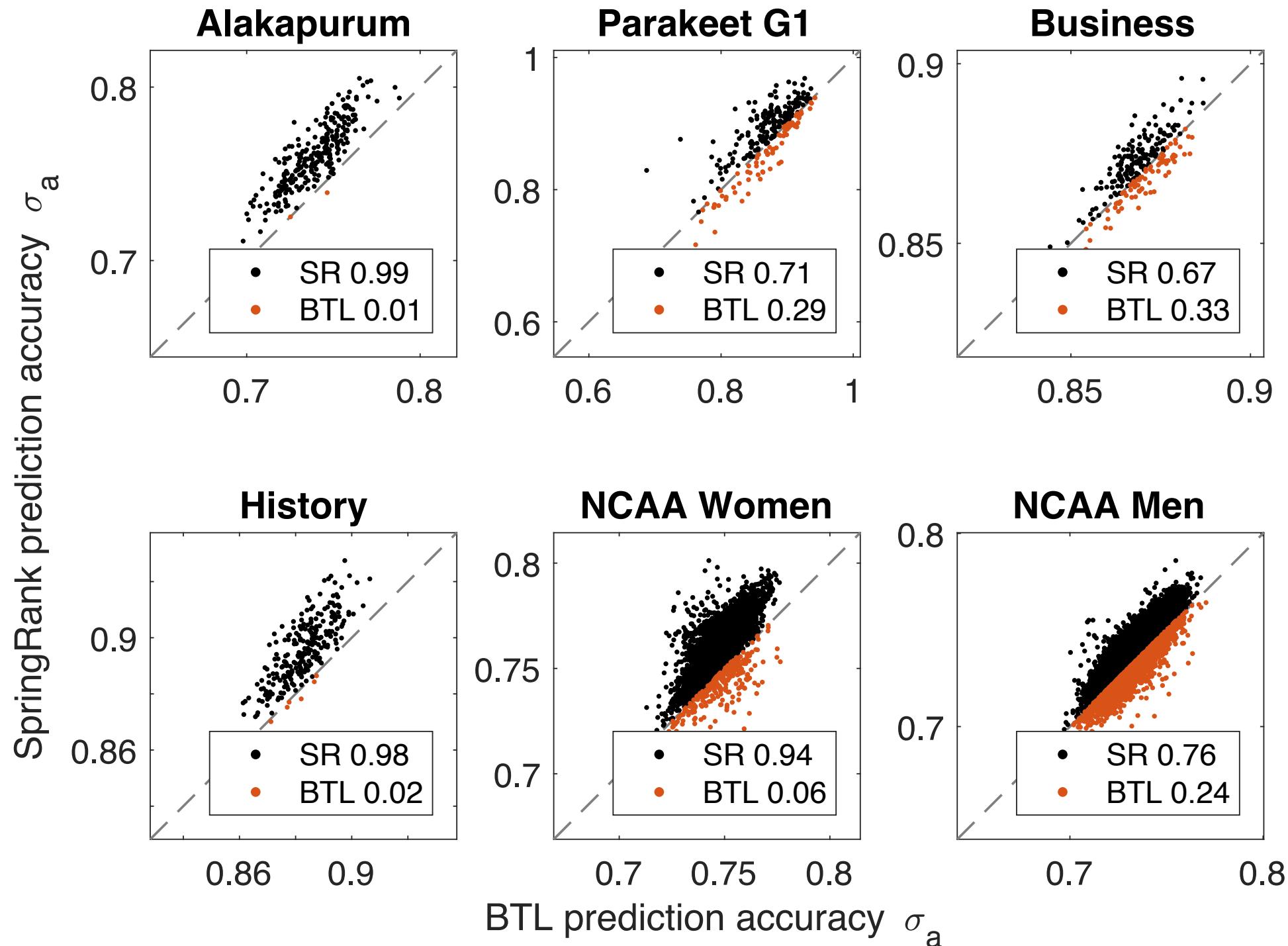
$$P_{ij}(\beta) = \frac{e^{-\beta H_{ij}}}{e^{-\beta H_{ij}} + e^{-\beta H_{ji}}} = \frac{1}{1 + e^{-2\beta(s_i - s_j)}}$$

# Cross validation vs BTL: SR makes better predictions

Accuracy:

$$\sigma_a = 1 - \frac{1}{2M} \sum_{i,j} |A_{ij} - (A_{ij} + A_{ji}) P_{ij}|$$

Goal: maximize the number of correctly predicted edge directions.



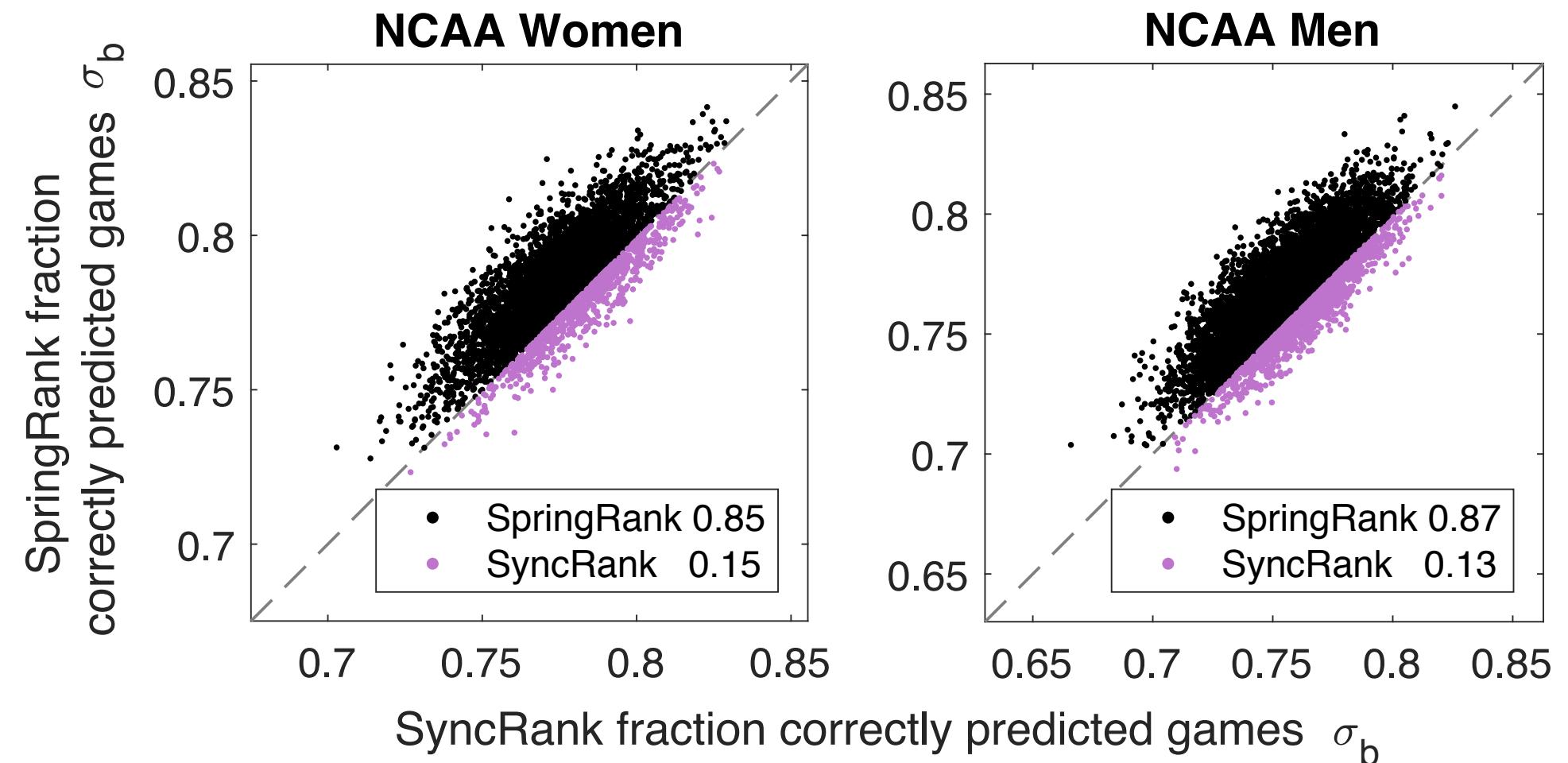
# Cross validation vs SyncRank: SR makes better predictions

“One-bit” Accuracy:

Higher ranked player always wins.

- No probabilistic prediction.
- Bad for gambling.

Goal: maximize the number of correctly predicted edge directions.



Why/when would a model of springs make better predictions than a model of the choices themselves? 🤔

It's unclear *why* we get this result!

Both BTL and SpringRank make logistic predictions about preference.  
Key Idea: SpringRank makes different regularization assumptions.

# Embeddings and Orderings 3: PageRank

**PageRank** defines scalar rank recursively:

*important pages are those that are linked to by important pages.*

- Great at finding the top 3 but limited predictions available using the PageRank scores.

## The PageRank Citation Ranking: Bringing Order to the Web

January 29, 1998

### Abstract

The importance of a Web page is an inherently subjective matter, which depends on the readers interests, knowledge and attitudes. But there is still much that can be said objectively about the relative importance of Web pages. This paper describes PageRank, a method for rating Web pages objectively and mechanically, effectively measuring the human interest and attention devoted to them.

We compare PageRank to an idealized random Web surfer. We show how to efficiently compute PageRank for large numbers of pages. And, we show how to apply PageRank to search and to user navigation.

## The Anatomy of a Large-Scale Hypertextual Web Search Engine

Sergey Brin and Lawrence Page

Computer Science Department,  
Stanford University, Stanford, CA 94305, USA  
scrgcy@cs.stanford.edu and page@cs.stanford.edu

### Abstract

In this paper, we present Google, a prototype of a large-scale search engine which makes heavy use of the structure present in hypertext. Google is designed to crawl and index the Web efficiently and produce much more satisfying search results than existing systems. The prototype with a full

# Embeddings and Orderings 3: PageRank

We imagine a web surfer who choose a starting webpage at random.

From that webpage, she looks at the links on the page, and either

- (a) clicks on a random link or
- (b) stops surfing; when she returns, she starts at a new random page.

What's the probability that she's at a particular page? *That's PageRank.*

$$\pi_{ji} = \frac{A_{ji}}{k_j}$$

$$p_i = \frac{1-d}{N} + d \sum_j p_j \pi_{ji}$$

$$\mathbf{p} = \left( \frac{1-d}{N} \right) \mathbf{1} + d \boldsymbol{\pi}^T \mathbf{p}$$

define a transition matrix

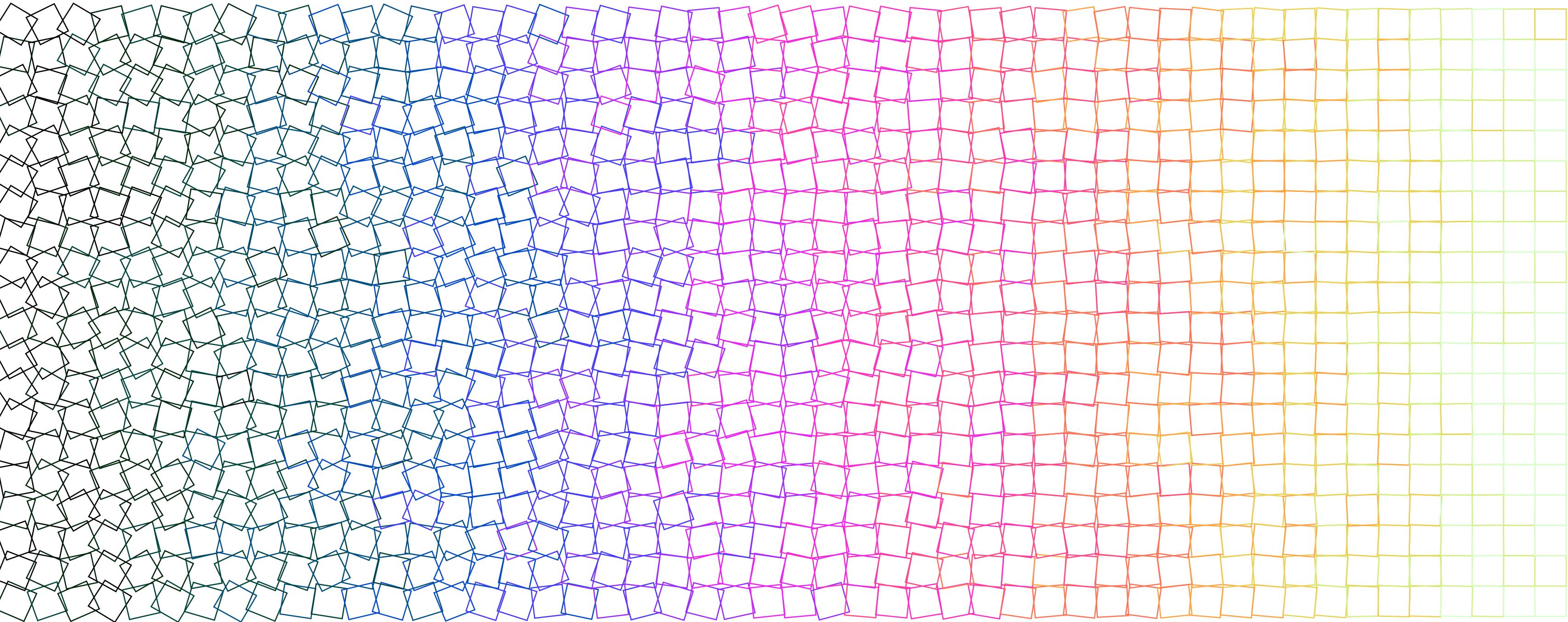
write the equation

matrix-vector form

**Alternative:** stationary distribution of random walk on the network + weak all-to-all links

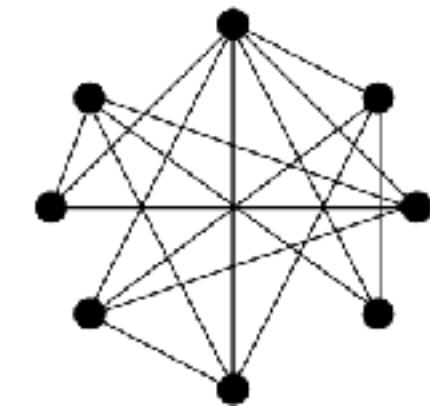
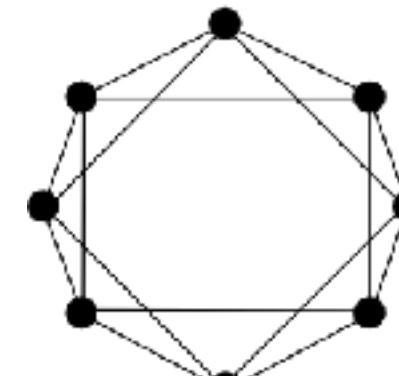
**Jeremy Kun:** <http://www.infinitelooper.com/?v=K3pT0gTaDec&p=n>

# sin/cos: generative models for networks



# Stochastic models, sets, and distributions

- a model is just a recipe:  
choose parameters → make the network
- a **stochastic** (generative) model is also just a recipe:  
choose parameters → draw **a** network
- since a single stochastic generative model can generate many networks, the model itself corresponds to a **set of networks**.
- and since the generative model itself is some combination or composition of random variables, a **random graph model** is a set of possible networks, each with an associated probability, i.e., a distribution.

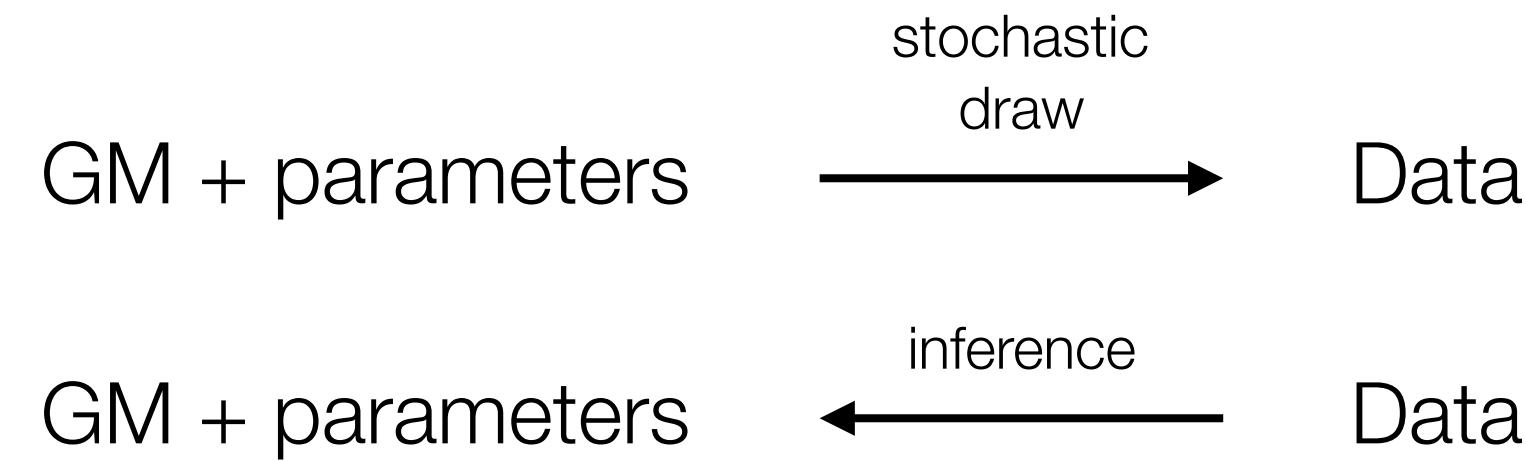


By changing the recipe, we can change the *support* of the distribution and the *probability masses* themselves.

# Generative models for network structure

*Generate the structure you wish to infer.*

We like generative models because they open the door to inference:



In other words: let's write down a model whose ensemble's distribution is not uniform but **highly peaked** around networks with structures that we want to see.

# e.g. The stochastic block model

Assign each node to one of  $B$  blocks.  $b_i$

GM + parameters

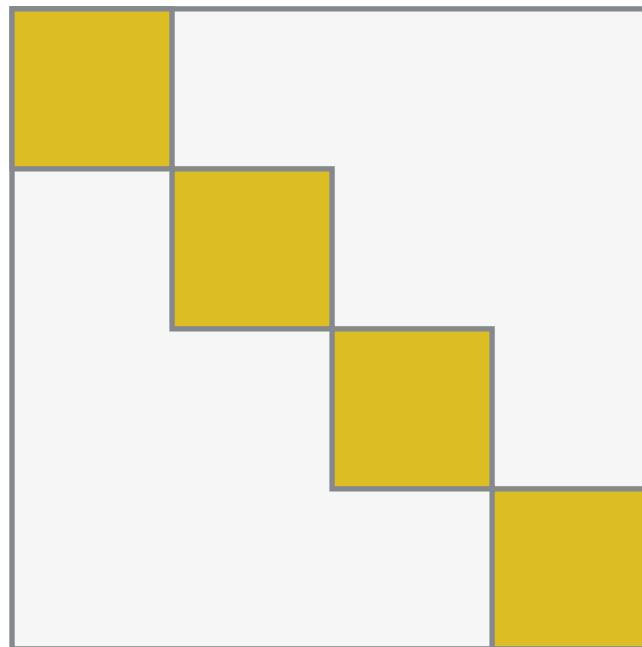
stochastic  
draw

→ Data

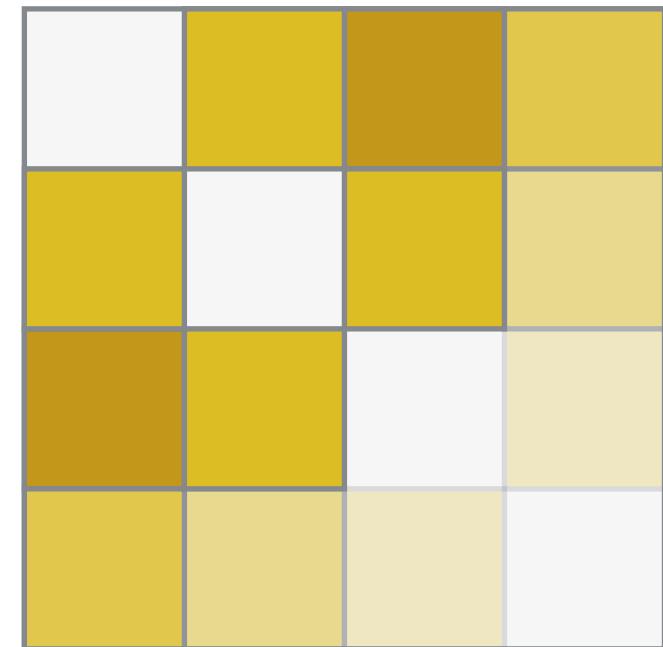
Let the probability that two nodes connect depend *only* on their blocks:

$$\Pr(A_{ij}|b_i, b_j) = \omega_{b_i, b_j}$$

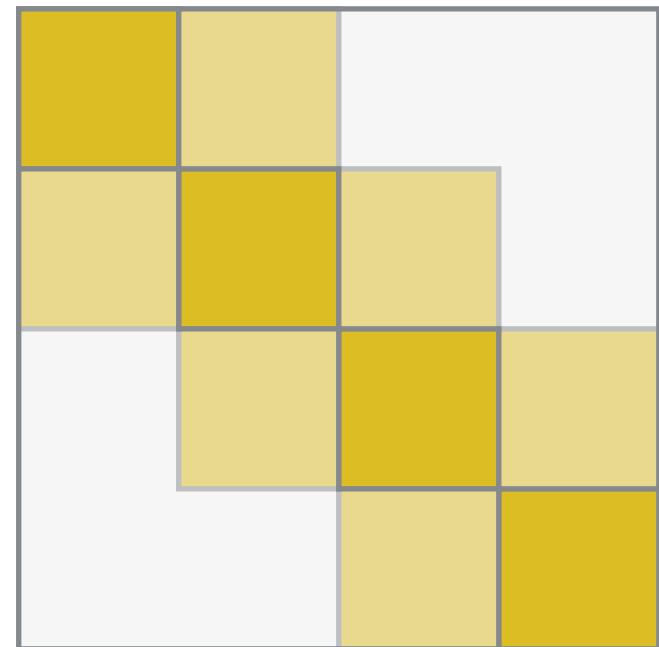
Then we can choose the matrix  $\omega$  to have whatever structure we want!



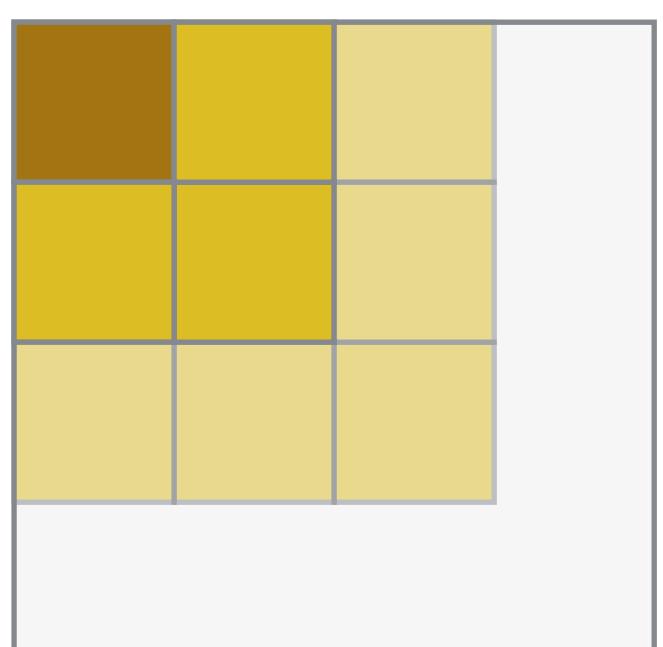
Assortative



Disassortative



Ordered



Core-periphery

# e.g. The stochastic block model

GM + parameters

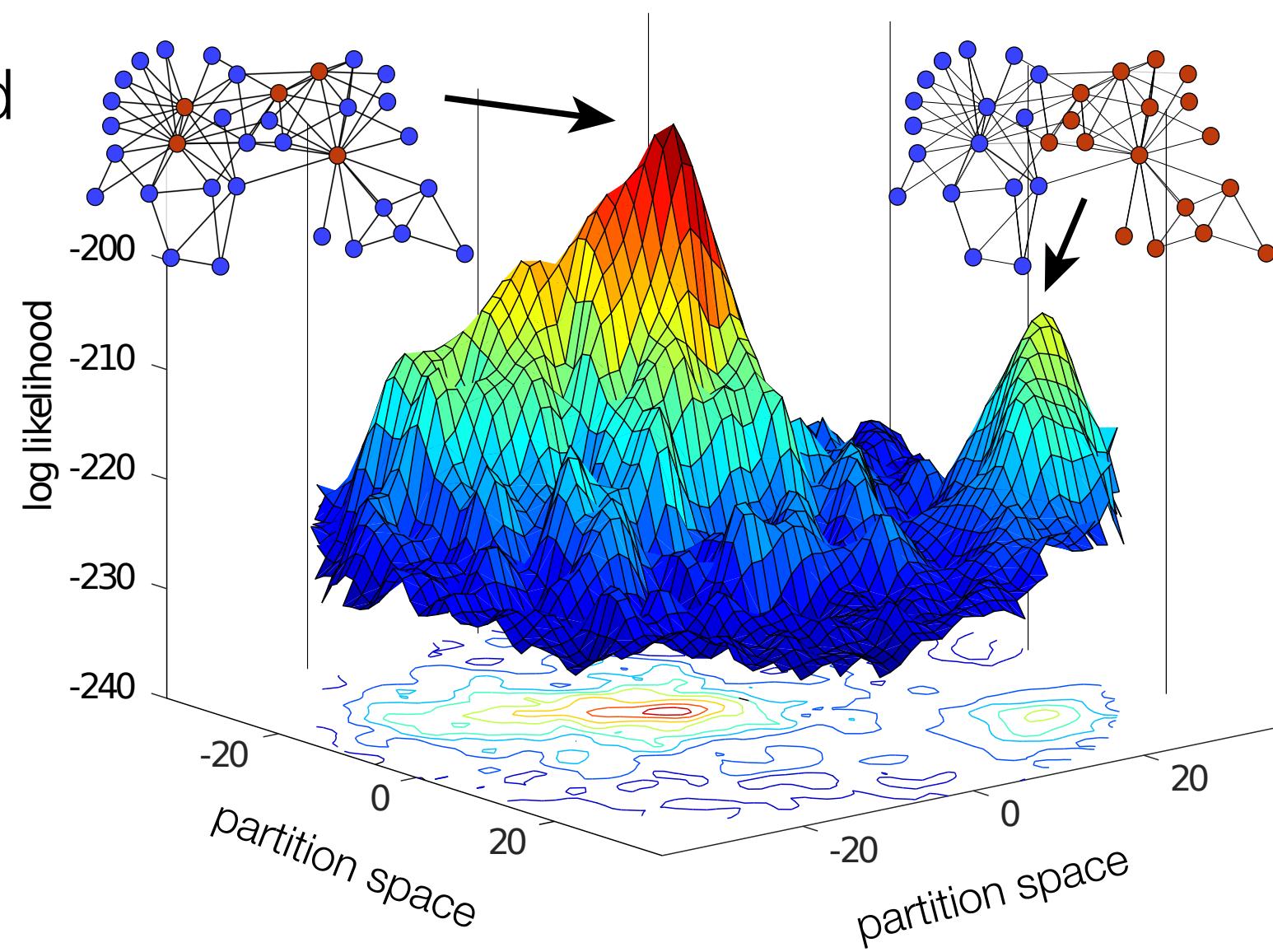
inference

Data

When we run the generative process in reverse (aka inference), we find community structure.

This is nothing more than a statistically principled approach to **fitting a model to data**.

But instead of fitting a line to a scatter of (x,y) data, we're fitting a model for networks with community structure to data.



# Embeddings and Orderings 4: Ball & Newman

## Generative model:

Generate the patterns that you want to identify.

Create N nodes.

Assign each node an integer rank  $r$ , from 1 to N.

IRL, not all friendships are reciprocated 🤦

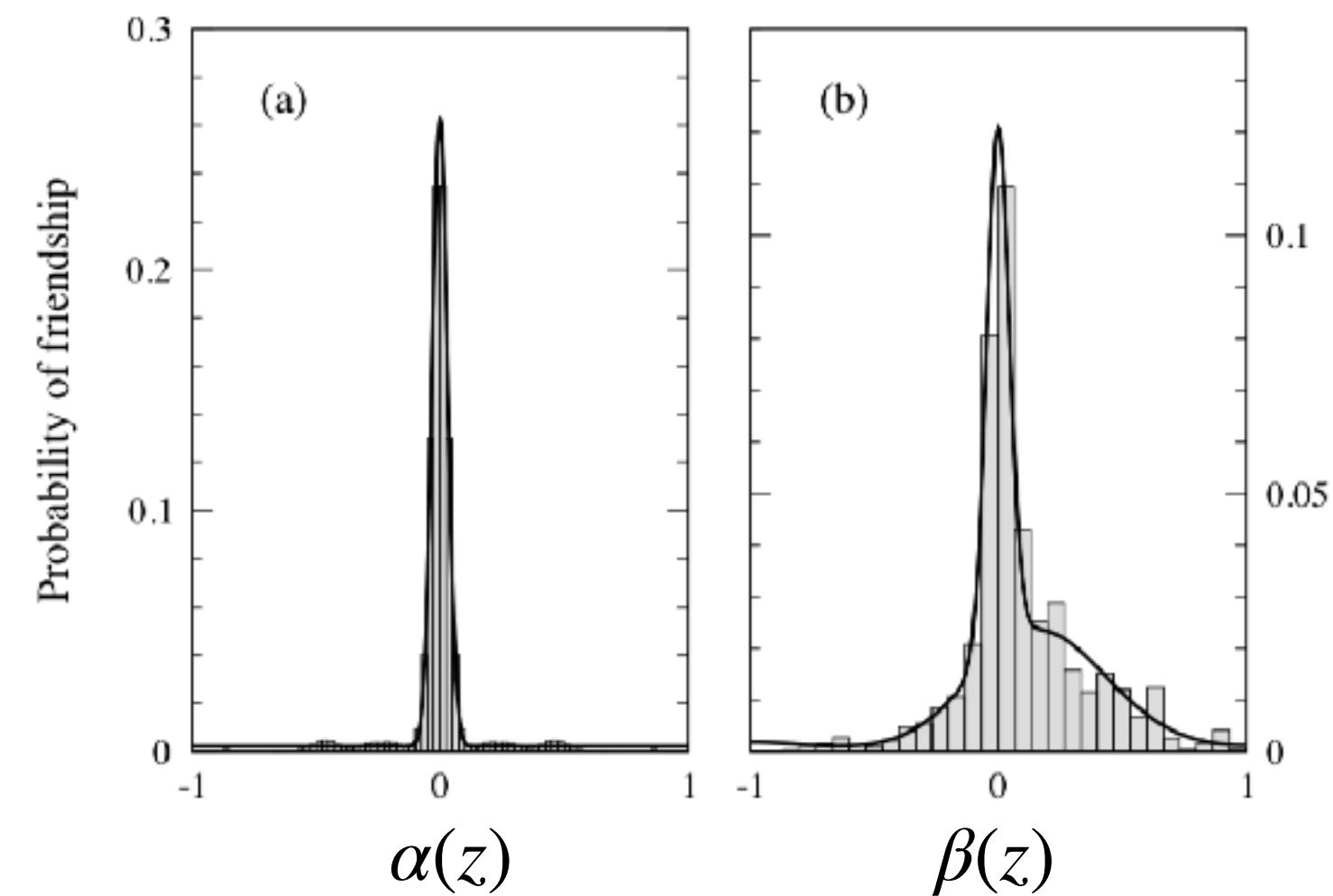
So let's generate undirected AND directed edges:

$$P(i \leftrightarrow j) = \alpha(r_i - r_j)$$

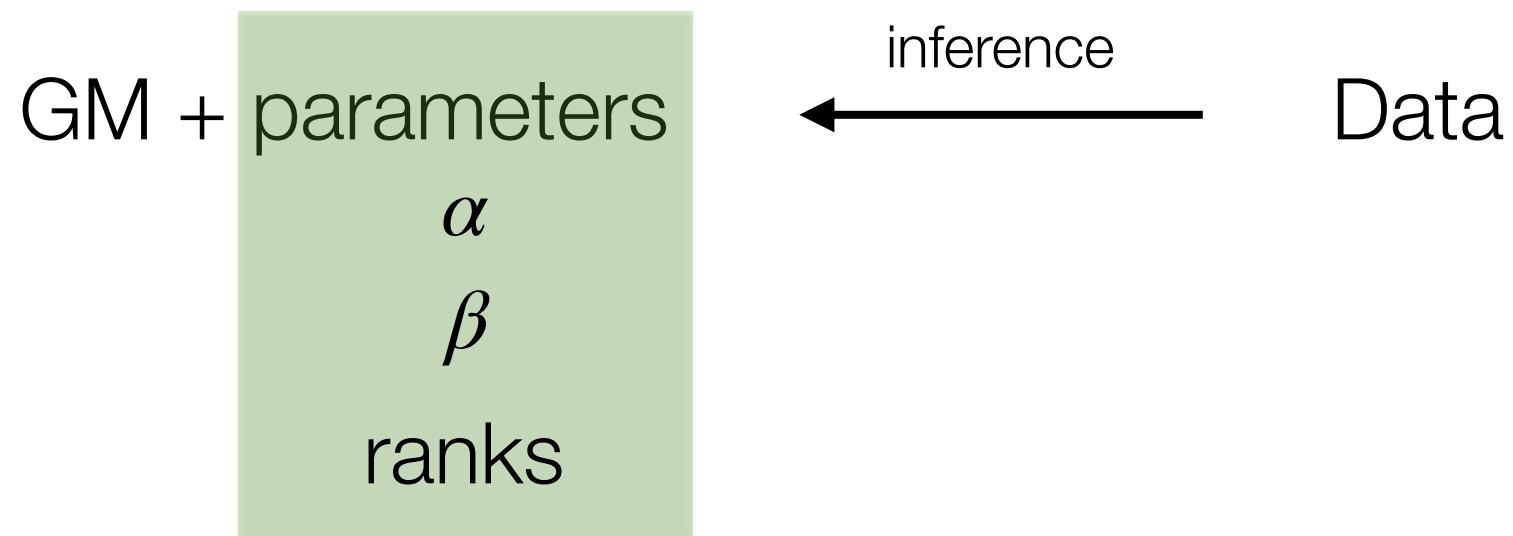
A gaussian centered at 0

$$P(i \rightarrow j) = \beta(r_i - r_j)$$

Fourier cosine series, keeping five terms & squaring to enforce nonnegativity, plus an additional Gaussian peak at the origin.



# Embeddings and Orderings 4: Ball & Newman



Inferred parameters of people's attachment preferences & ranks.

- Identified the need to learn from reciprocated friendships.
- Found that in AddHealth data, teens link to others of *nearby* social status.

12th grade

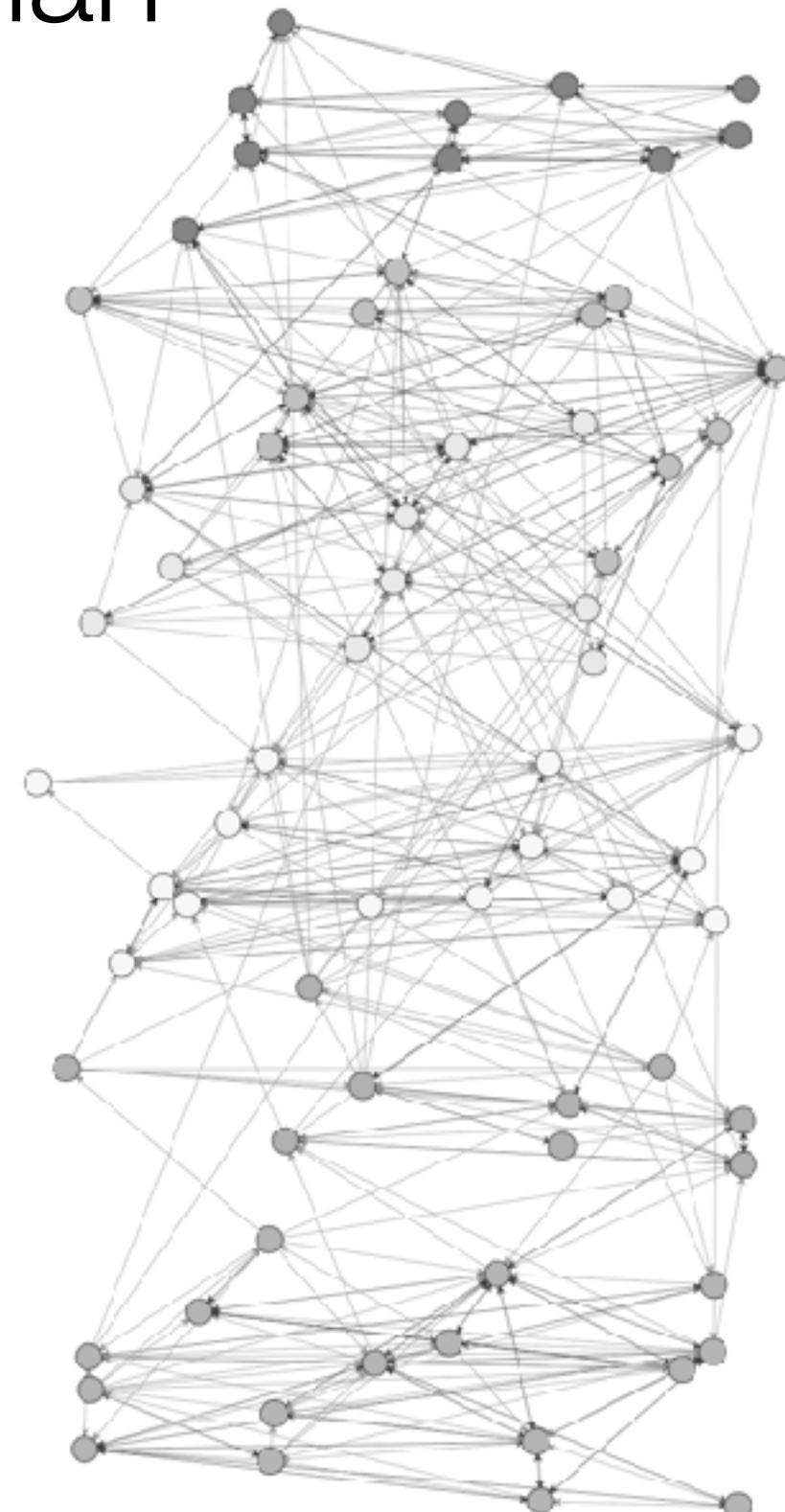
11th grade

10th grade

9th grade

8th grade

7th grade



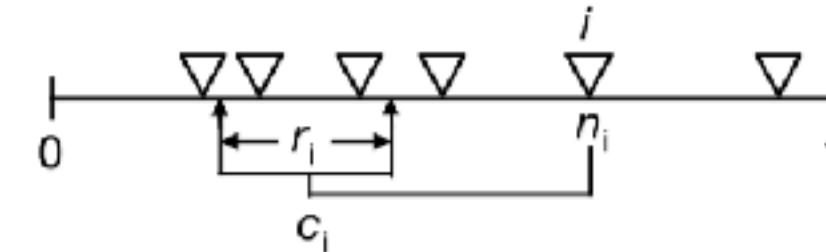
# Embeddings and Orderings 5: Niche Models

**Niche Models** embed species in a latent space based on feeding preferences:

*most species feed from narrow range in a 1-dim. space (~body size).*

- Great for food webs. Inference models v slow for all but small networks.

Want more? Jen Dunne, Cris Moore



**Figure 1** Diagram of the niche model. Each of  $S$  species (for example,  $S = 6$ , each shown as an inverted triangle) is assigned a 'niche value' parameter ( $n_i$ ) drawn uniformly from the interval  $[0,1]$ . Species  $i$  consumes all species falling in a range ( $r_i$ ) that is placed by uniformly drawing the centre of the range ( $c_i$ ) from  $[r_i/2, n_i]$ . This permits looping and cannibalism by allowing up to half of  $r_i$  to include values  $\geq n_i$ . The size of  $r_i$  is assigned by using a beta function to randomly draw values from  $[0,1]$  whose expected value is  $2C$  and then multiplying that value by  $n_i$  [expected  $E(n) = 0.5$ ] to obtain the desired  $C$ . A beta distribution with  $\alpha = 1$  has the form  $f(x|1, \beta) = \beta(1-x)^{\beta-1}$ ,  $0 < x < 1$ , 0 otherwise, and  $E(X) = 1/(1+\beta)$ . In this case,  $x = 1 - (1-y)^{1/\beta}$  is a random variable from the beta distribution if  $y$  is a uniform random variable and  $\beta$  is chosen to obtain the desired expected value. We chose this form because of its simplicity and ease of calculation. The fundamental generality of species  $i$  is measured by  $r_i$ . The number of species falling within  $r_i$  measures realized generality. Occasionally, model-generated webs contain completely disconnected species or trophically identical species. Such species are eliminated and replaced until the web is free of such species. The species with the smallest  $n_i$  has  $r_i = 0$  so that every web has at least one basal species.

# Embeddings and Orderings 6: Centrality?

**Centrality measures** are another form of ordering:

## **Geometric Centralities:**

- Harmonic
- Closeness
- Betweenness

So many varieties exist.

## **Connectivity:**

- Degree
- Eigenvector
- PageRank
- Katz

**Recommendation:** choose the centrality (or method) whose assumptions match what you know about your problem.

</methods>

<applications>

# Many uses for the same techniques. cf regression

## Treat the network like a system:

**Extrapolation.** Make predictions for as-yet unseen nodes (in “space” or time).

**Interpolation.** Identify missing links.

**Generalization.** Nodes of this type are like others of the same type.

## Treat the network like an artifact:

**Mechanisms.** How did this network arise? What rules governed its assembly?

**Explanations.** Coarse-graining or compression.

## Treat the network like a means to an end; an intermediate data structure:

**Useful division.** Need groups so that we can assign treatments in an A/B test.

**Simplification.** Downstream regression model needs ranks or groups.

# Structure and inequality in academic hiring



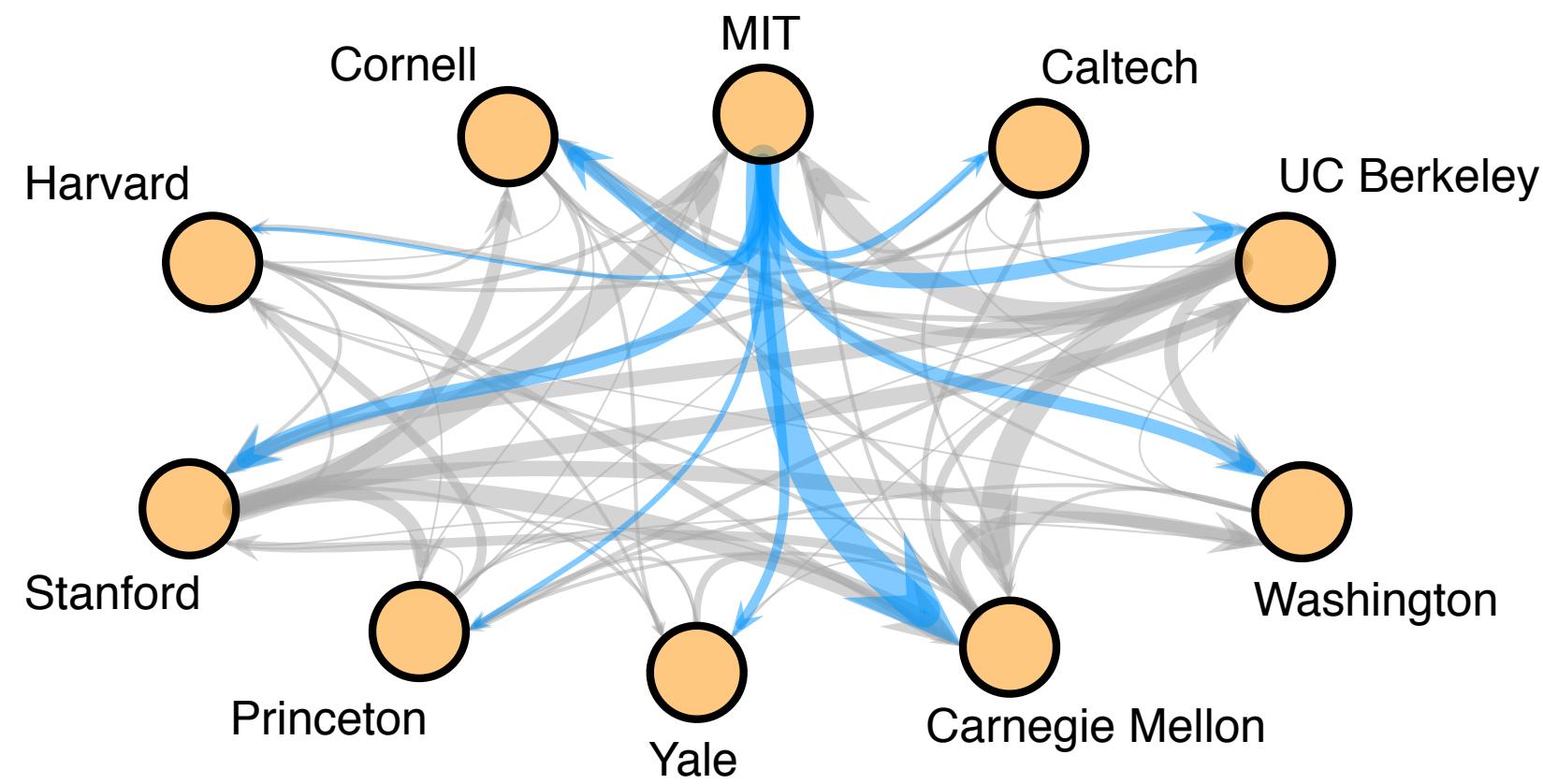
# Collect the data (by hand 😭 )

CVs of all US & Canadian tenure-track faculty in CS, Business, History: 2011-2013.

	Computer Science	Business	History
institutions	205	112	144
tenure-track faculty	5032	9336	4556
mean size	25	83	32
female	15%	22%	36%

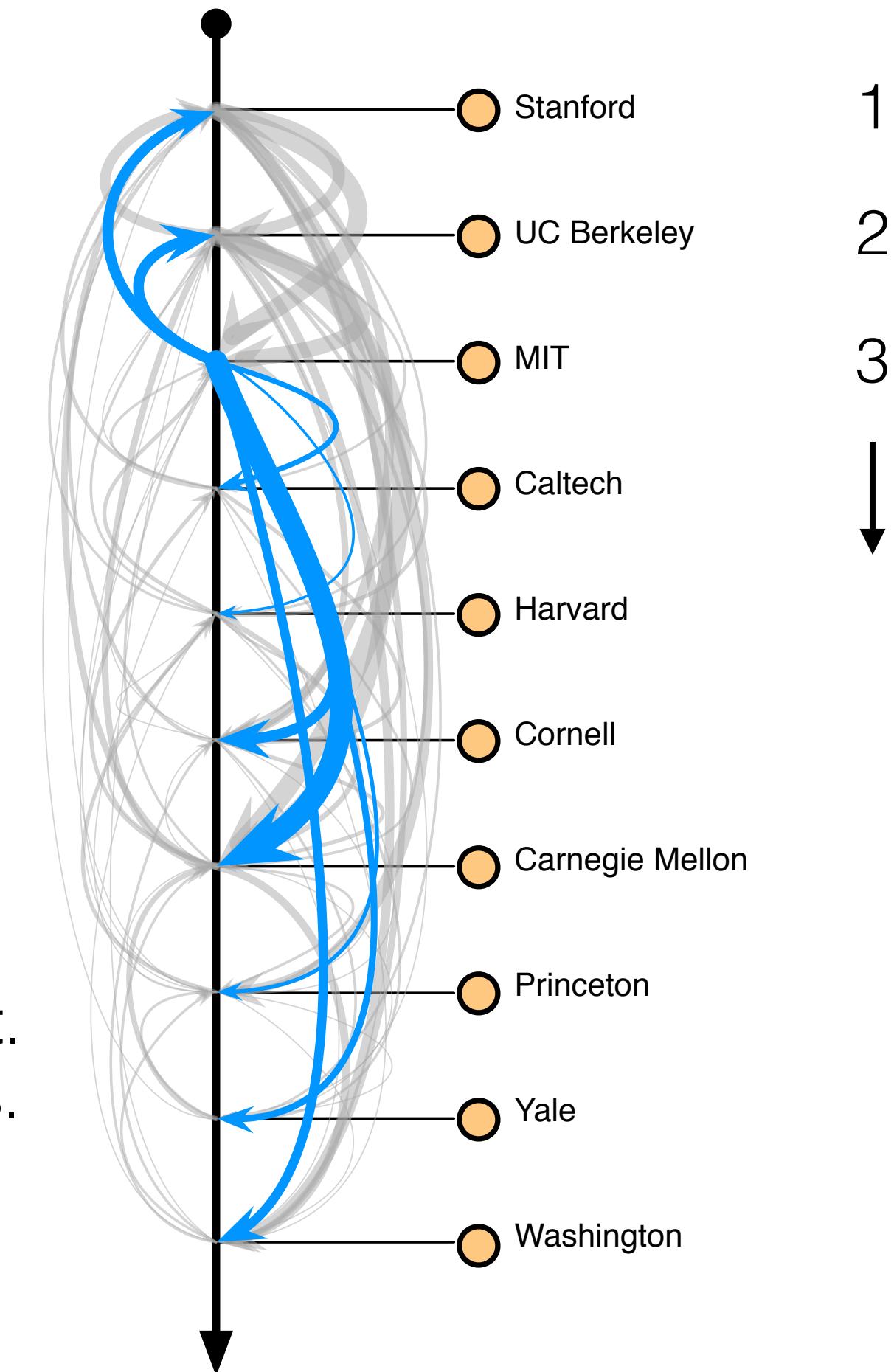
total: **18,924** CVs

# Faculty hiring networks



Premises:

1. Each hiring committee wants to hire the best.
2. Entire network reveals **collective preferences**.



# Faculty hiring networks

**systematic**

90% of hiring movement  
is “down” the hierarchy

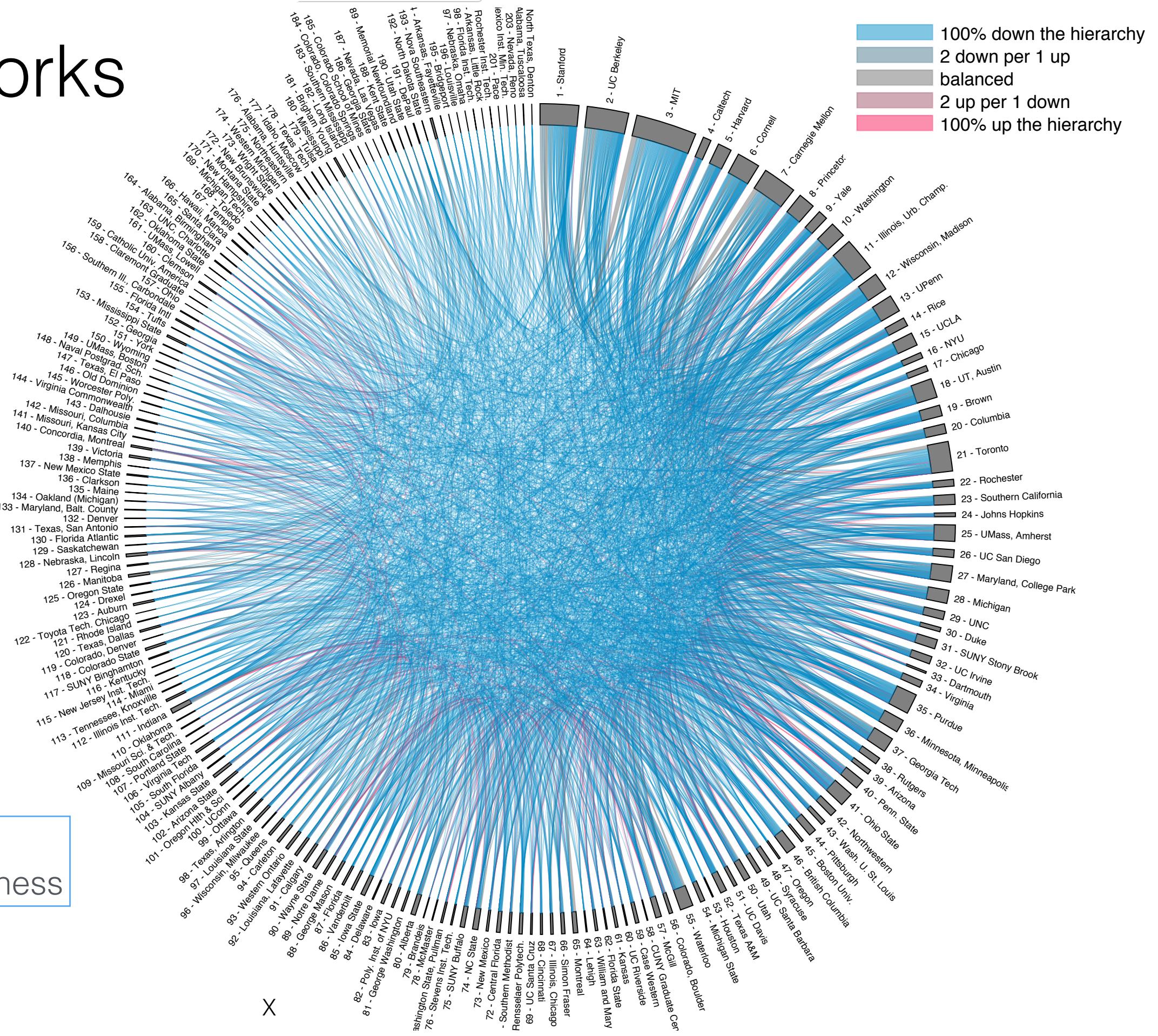
**steep**

< 7% of faculty have PhD  
from lower 75% of universities

**biased**

median change for women  
~3 ranks worse than men

[danlarremore.com/faculty/](http://danlarremore.com/faculty/)  
explore 19,000 hires for History, CS, Business



# What else explains movement in this labor market?

Generative model:

**prestige**  
**productivity**  
**postdoc experience**

**gender**  
**geography**

candidates

Cornell

MIT

Caltech

UW

openings

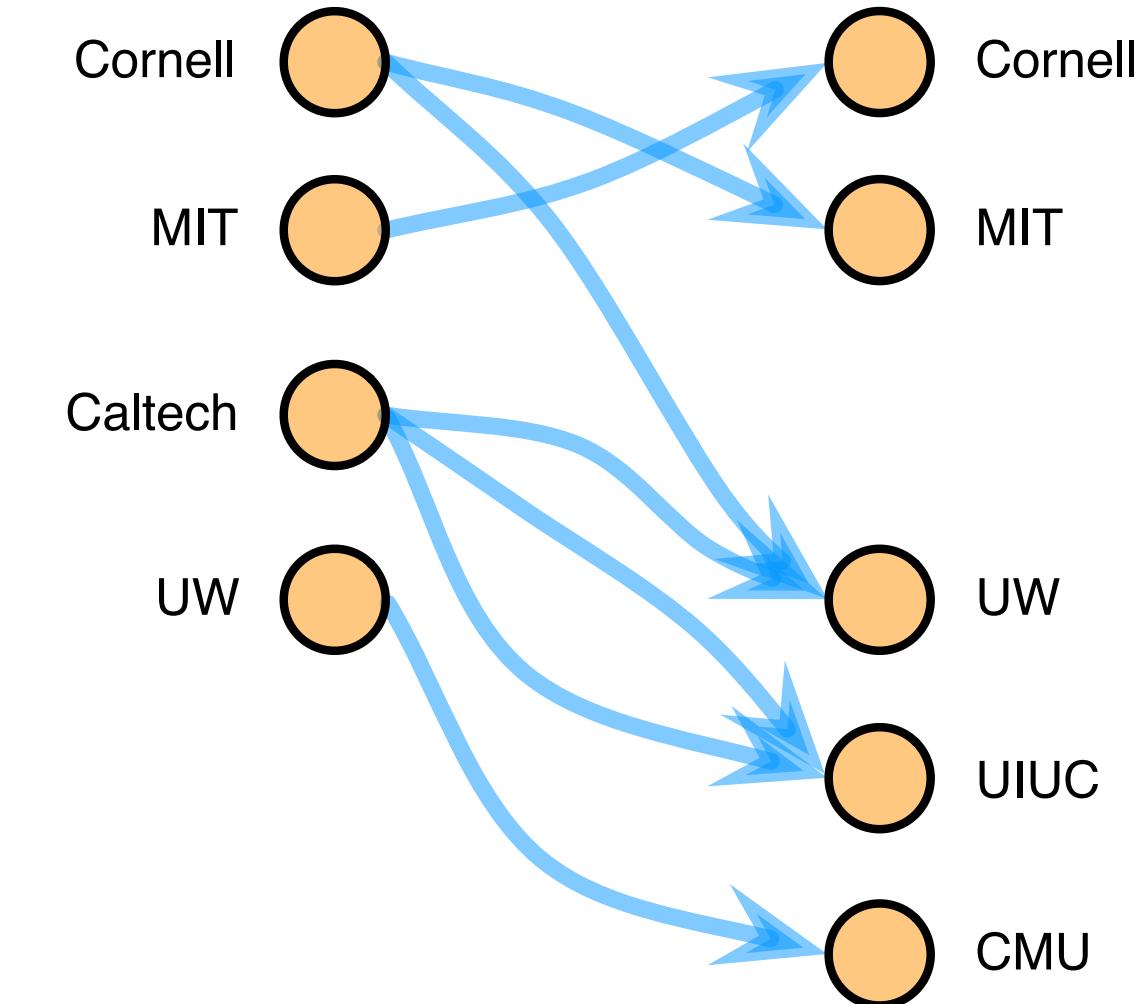
Cornell

MIT

UW

UIUC

CMU

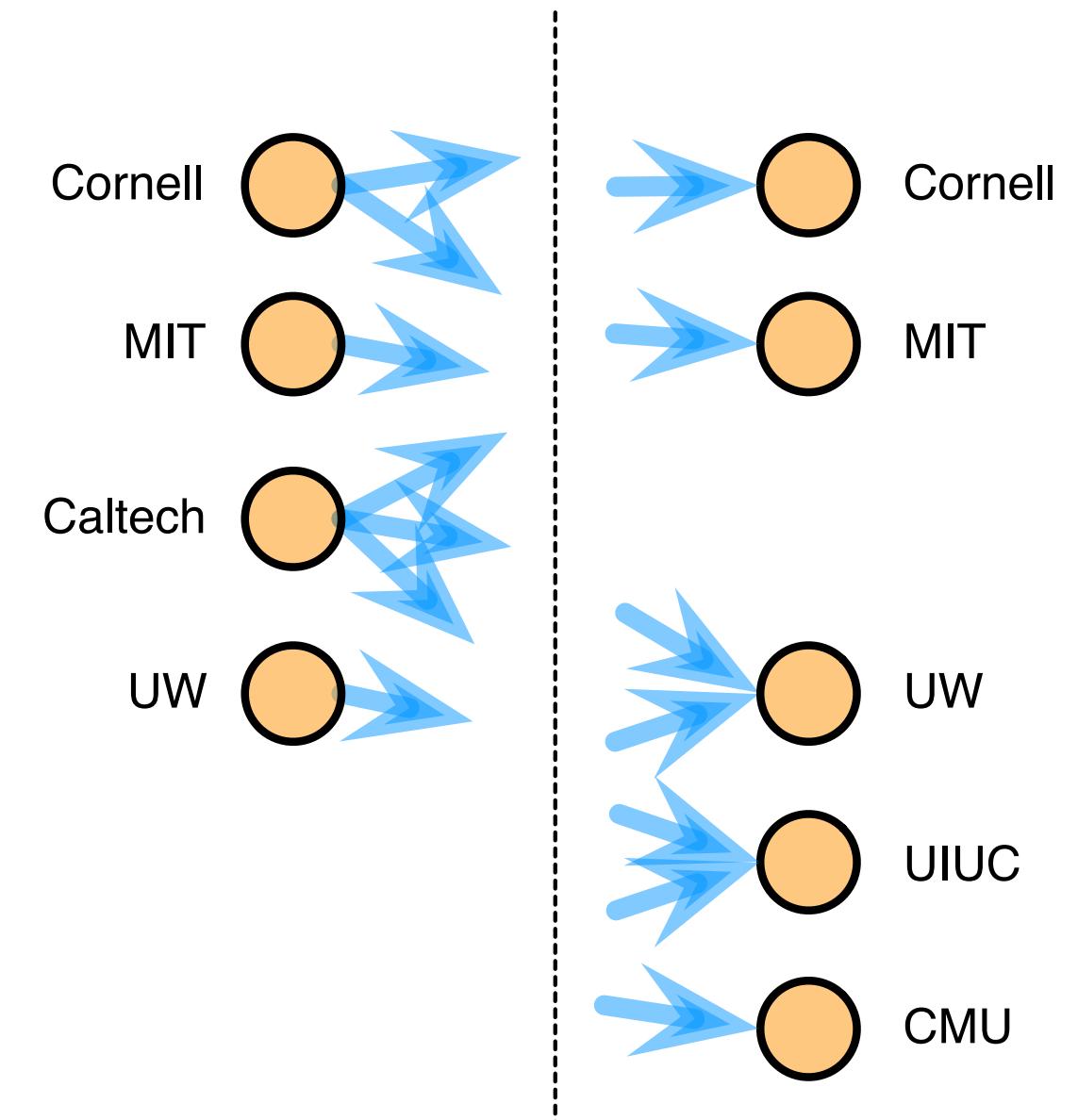


# What else explains movement in this labor market?

Generative model:

**prestige**  
**productivity**  
**postdoc experience**

**gender**  
**geography**



accurately generate the links!

# What else explains movement in this labor market?

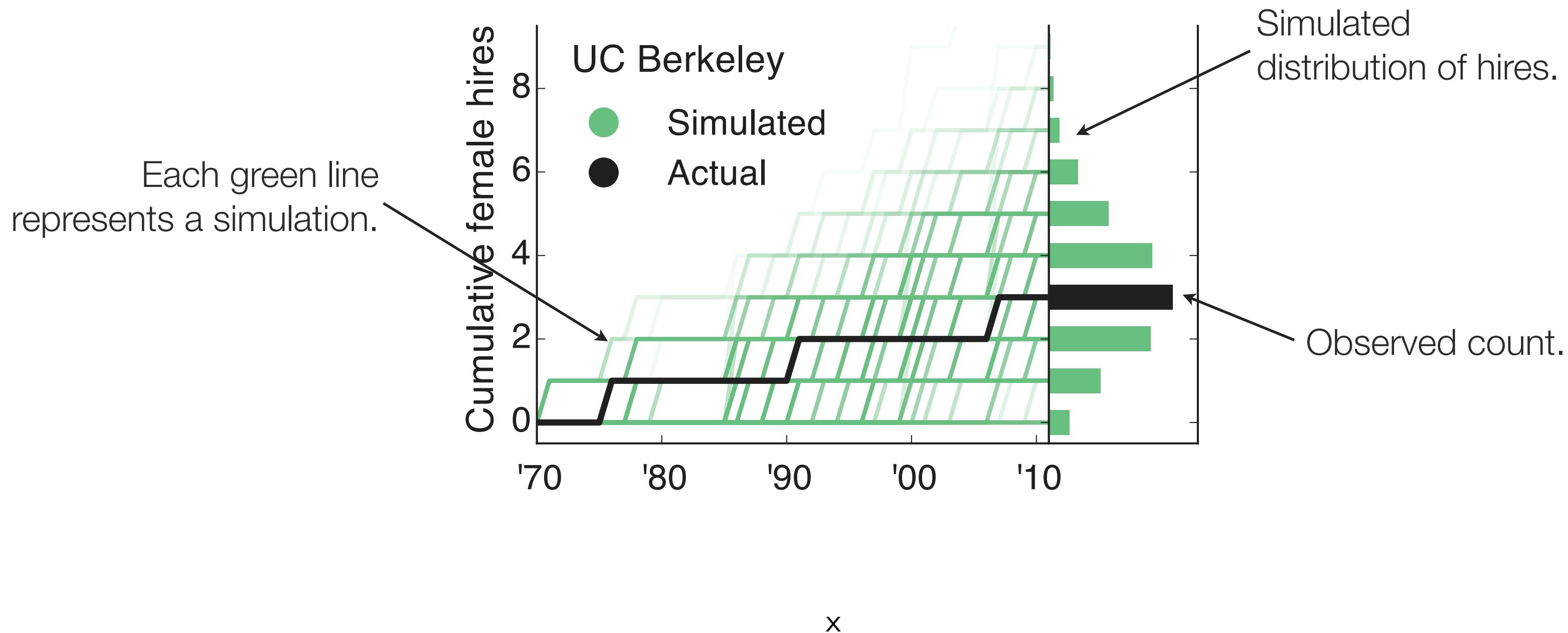
1. **Prestige difference**: Faculty Job vs PhD
2. **Productivity**
3. **Prestige of Faculty Job**
4. **Postdoc experience + geography** (together)
5. **Gender** label *not* significant after other factors included.

a woman on the job market must have published ~1 additional paper  
to be placed the same as an equally qualified man.

# Institution-level results

Using 40 years of actual hiring data, simulate hiring patterns for each institution.

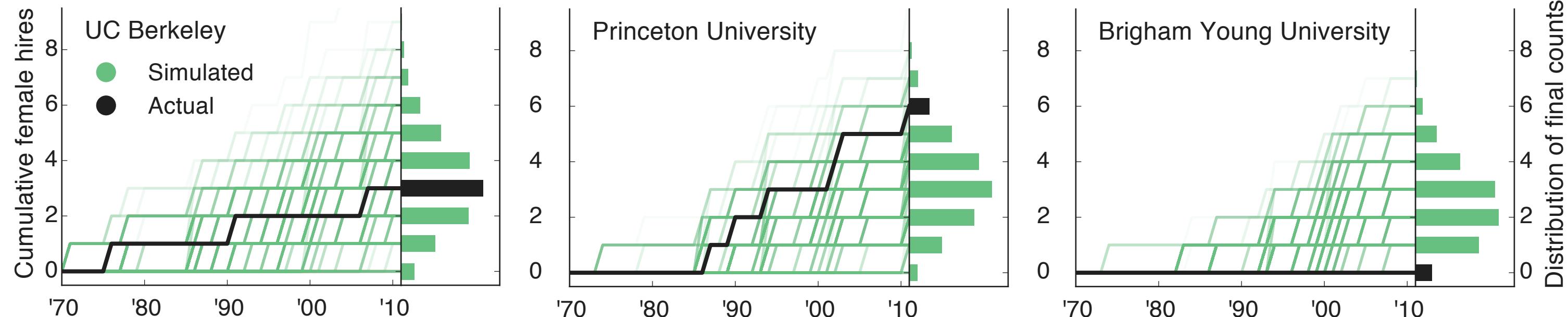
Compare actual vs. expected number of female hires.



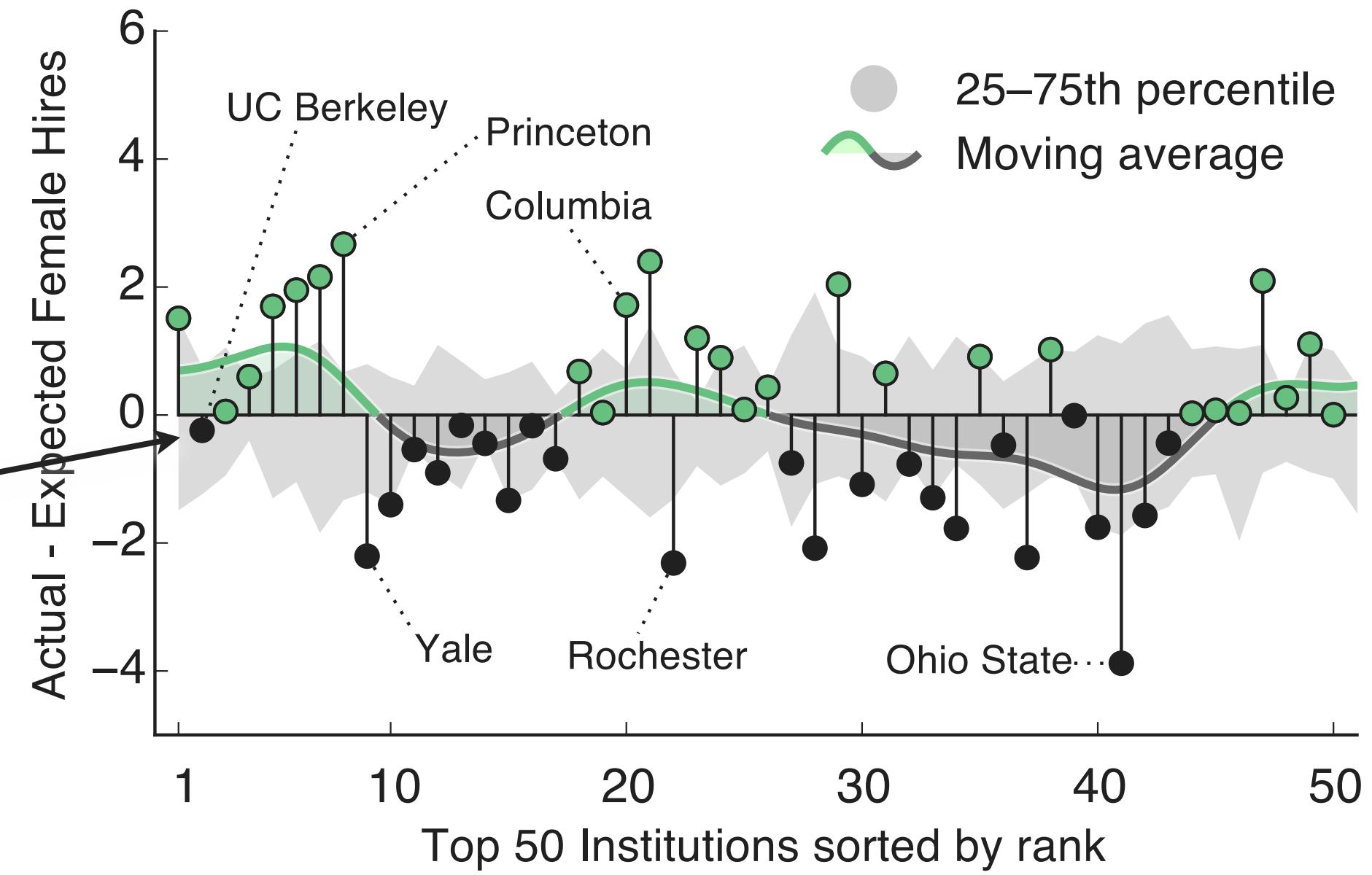
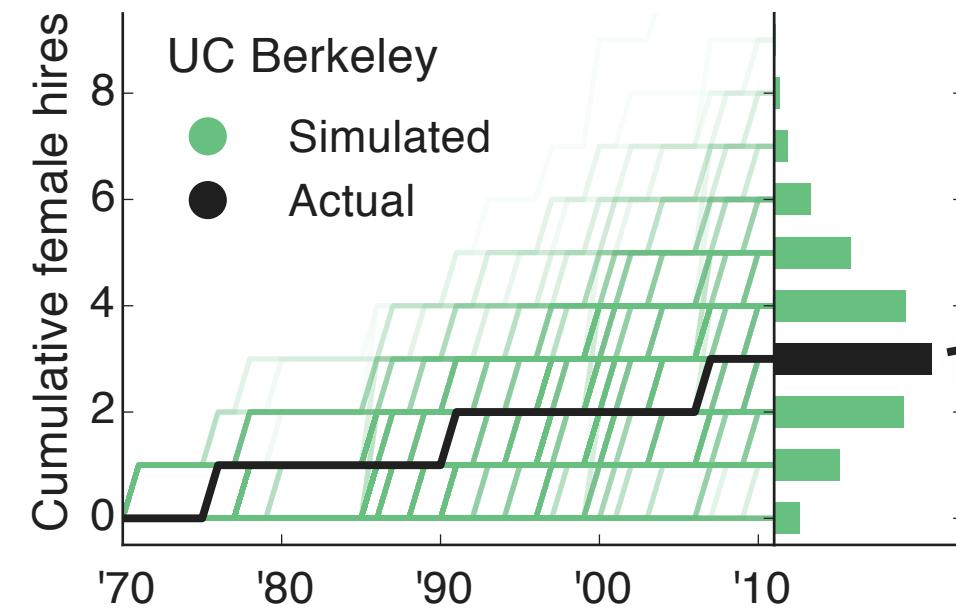
# Institution-level results

Using 40 years of actual hiring data, simulate hiring patterns for each institution.

Compare actual vs. expected number of female hires.



# Institution-level results



# Institution-level results

For the top 50 institutions,  
we see an oscillation.

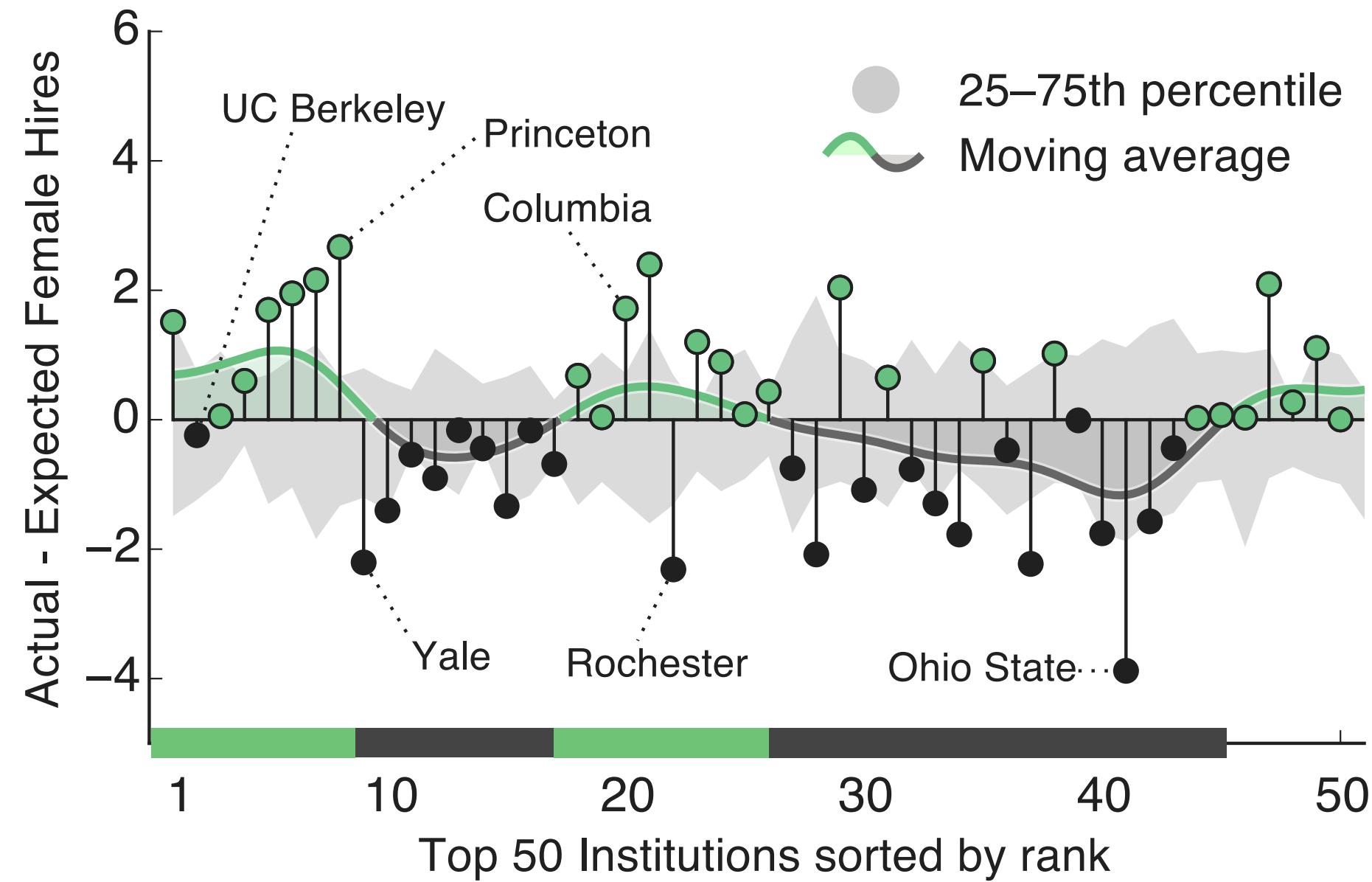


**Why?**

An interference effect?

Two distinct candidate pools?

Is it real?

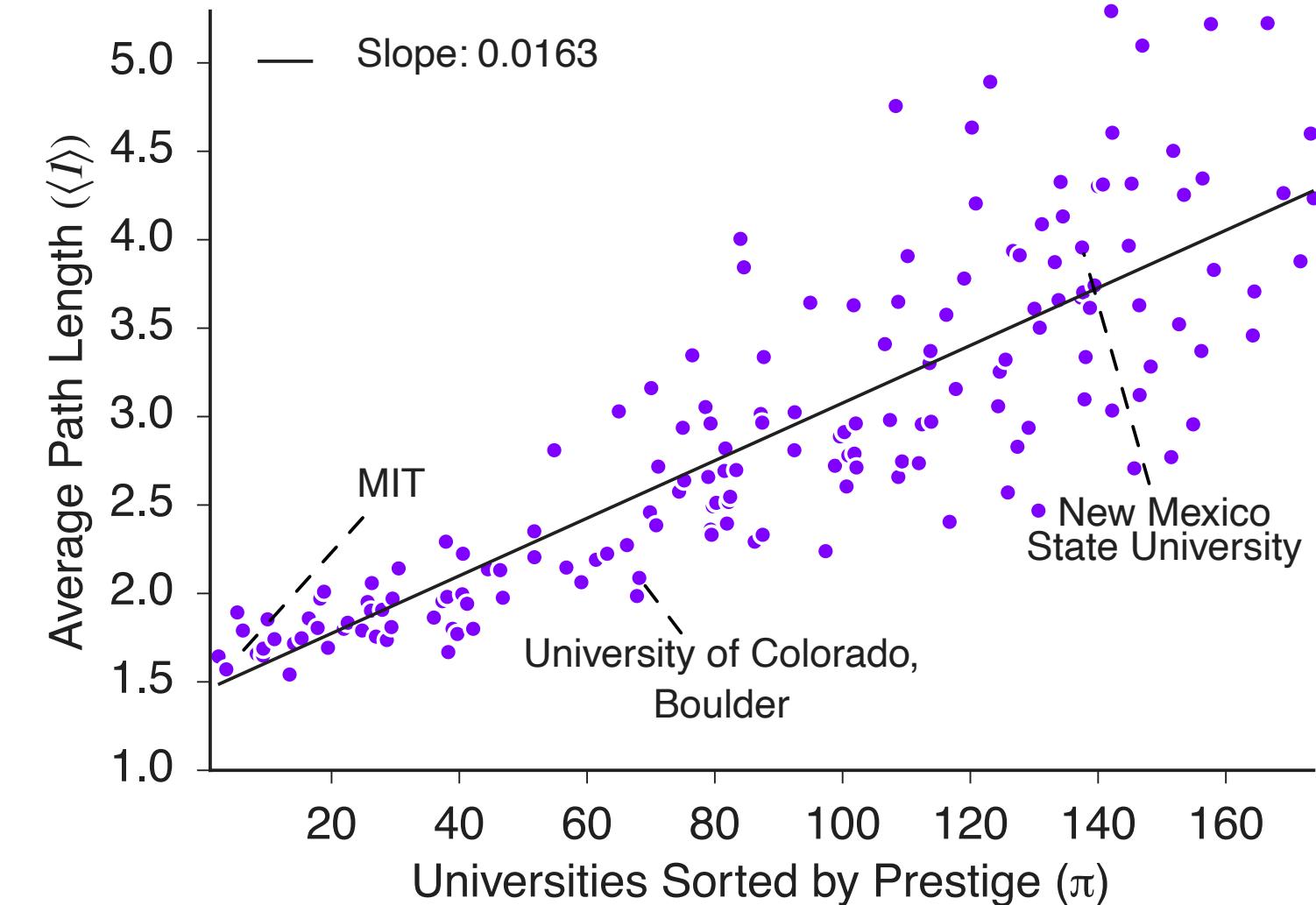


# Does the structure of this network affect *ideas*?

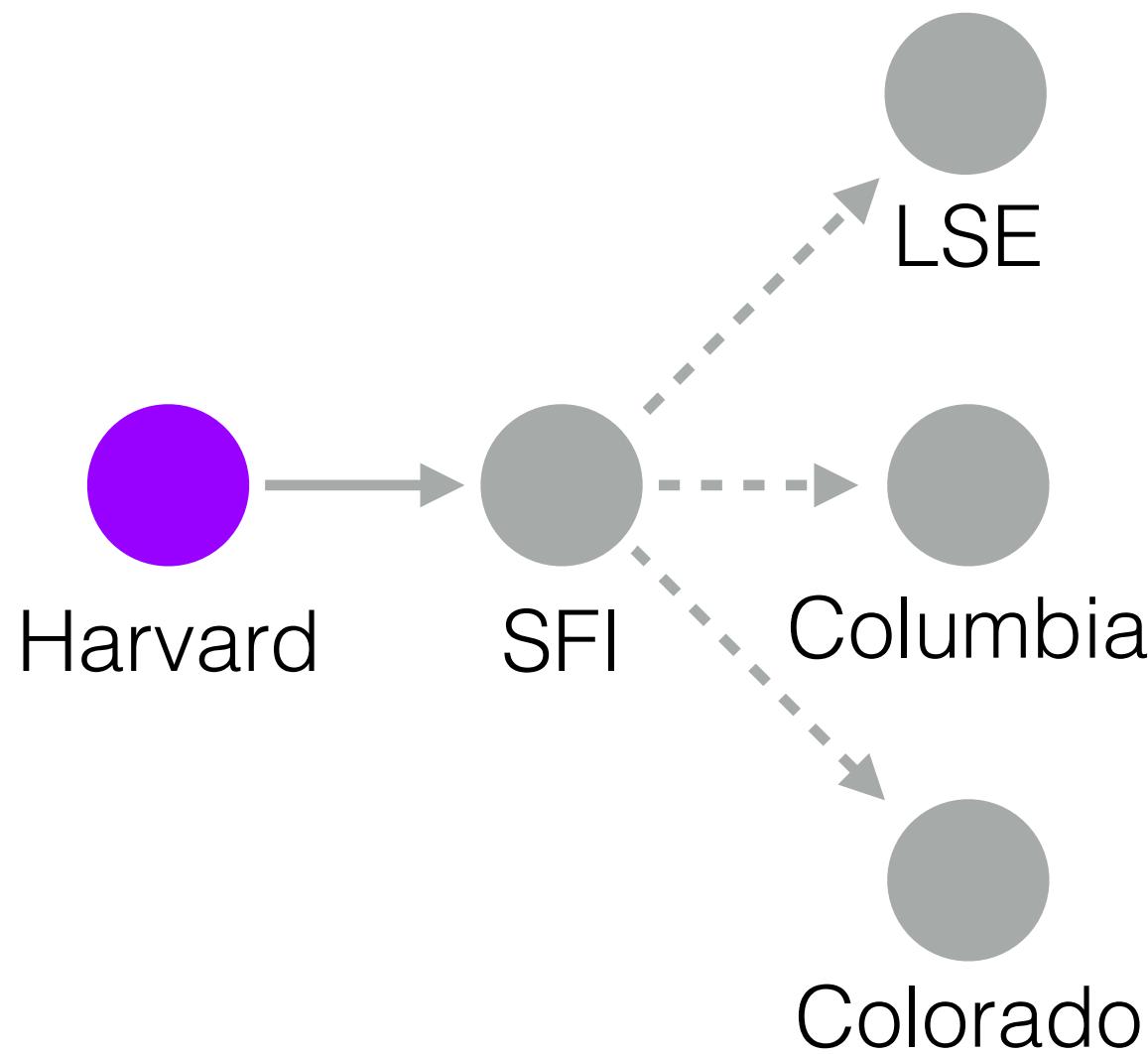
Prestigious institutions are **closer** to all other institutions.

What implications does this have for the **exchange & filtration** of ideas?

Does the prestige hierarchy lead to **epistemic inequality**?



# New hires as vectors for infectious ideas?



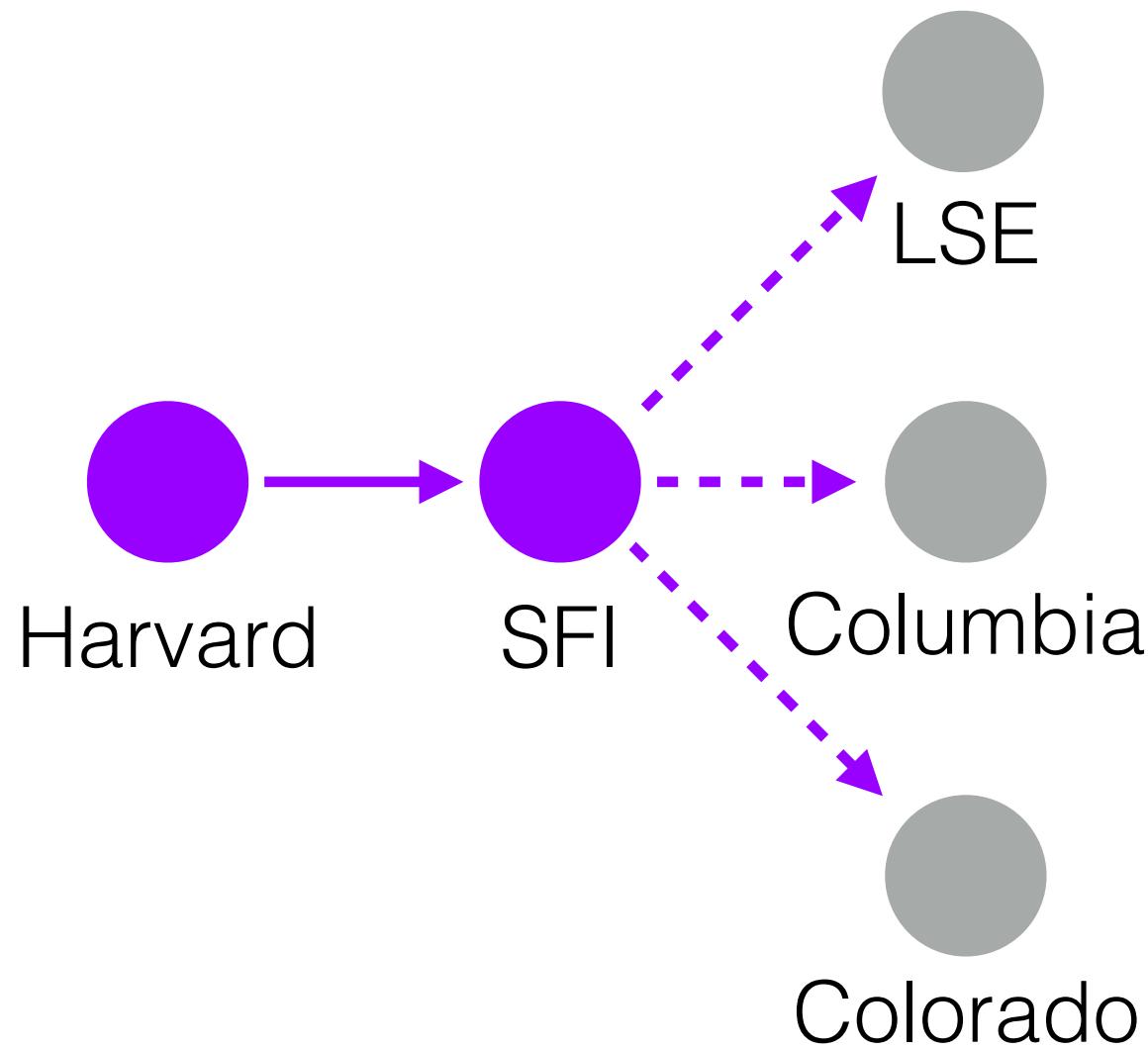
Do new hires *actually* bring ideas with them?  
[or would popular topics get there anyway?]

Are some universities better idea exporters?

Epidemic model: treat the idea as an infection,  
and a new hire as “infectious.”

The probability that a hire transmits the idea to  
an uninfected university:  $p$  [idea quality]

# New hires as vectors for infectious ideas?



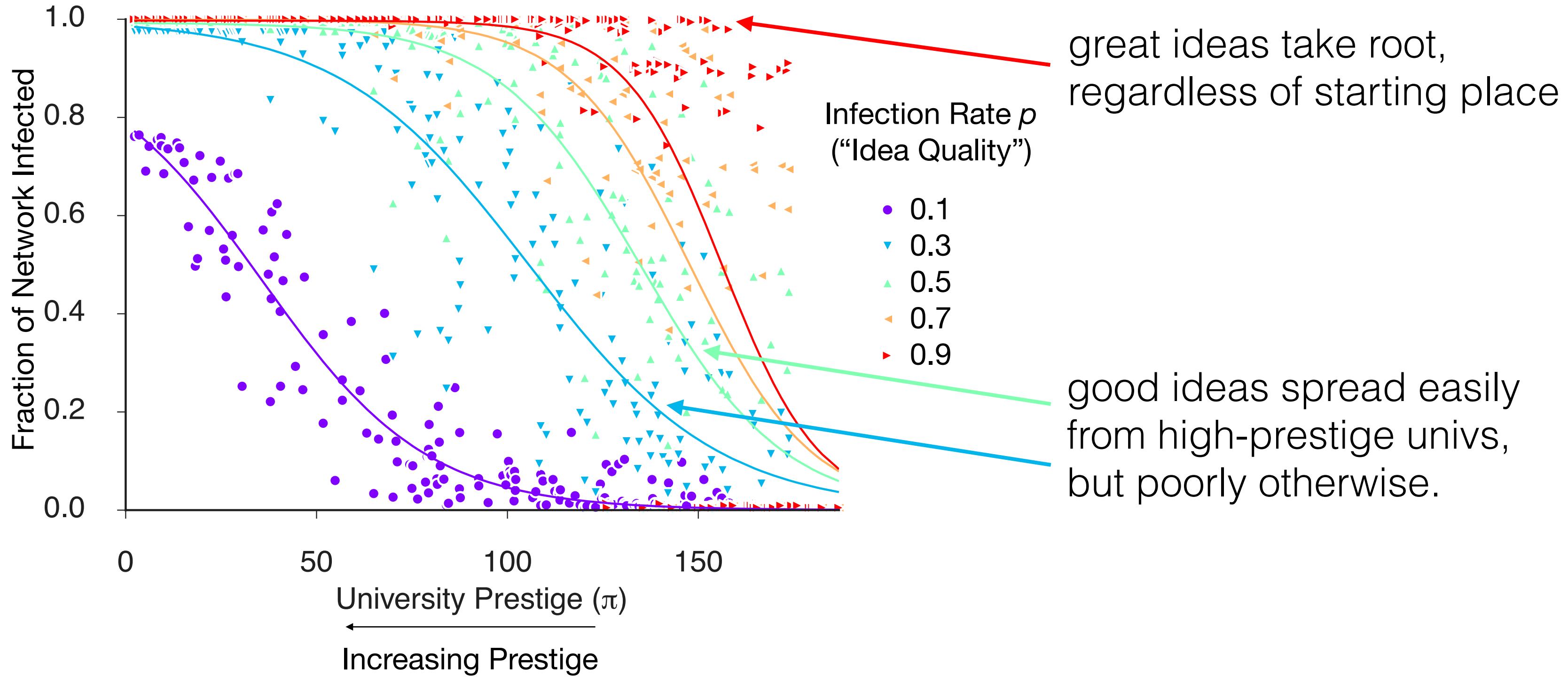
Do new hires *actually* bring ideas with them?  
[or would popular topics get there anyway?]

Are some universities better idea exporters?

Epidemic model: treat the idea as an infection,  
and a new hire as “infectious.”

The probability that a hire transmits the idea to  
an uninfected university:  $p$  [idea quality]

# Network position & the spread of ideas



# Do *real* ideas spread along hiring links?

Analyzed over 200,000 computer science publications and over 2,500 hires.

Flagged publications on topic modeling, incremental computing, deep learning.

Identified faculty who brought [topic] with them when they were hired.

Identified faculty who began working on [topic] only 2+ yrs after being hired.

Compared relative rates of hiring-link spread vs spontaneous spread. (vs random)

Spread of **topic modeling ( $p=0.01$ )** & **incremental computing ( $p=0.01$ )** significantly tied to infection via hiring. Spread of deep learning ( $p=0.2$ ) *not* significantly linked to hiring.



University of Colorado **Boulder**

## Colorado

Sam Way  
Aaron Clauset  
Allison Morgan  
Dimitrios Economou

## Kauffman / Lux

Sam Arbesman



Allie Morgan



Sam Way



Aaron Clauset



Ewing Marion  
**KAUFFMAN**  
Foundation



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

Way, Morgan, Larremore, Clauset. *PNAS* (2019).

### **Prestige drives epistemic inequality in the diffusion of scientific ideas**

Morgan, Economou, Way, Clauset. *EPJ Data Science* (2018).

### **The misleading narrative of the canonical faculty productivity trajectory**

Way, Morgan, Clauset, Larremore. *PNAS* (2017).

### **Data-driven predictions in the science of science**

Clauset, Larremore, Sinatra. *Science* (2017).

### **Gender, productivity, and prestige in computer science faculty hiring networks**

Way, Larremore, Clauset. *Proc. WWW* (2016).

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

Clauset, Arbesman, Larremore. *Science Advances*. (2015).

Conference on Complex Networks  
**COMPLENET '18**

Hosted by Northeastern University Network Science Institute



BOSTON, MA

APRIL 2018

Xindi Wang



 **Northeastern University**  
Network Science Institute

# LEARNING TO PLACE OBJECTS: A NETWORK-BASED APPROACH

Xindi Wang

Onur Varol



Tina Eliassi-Rad



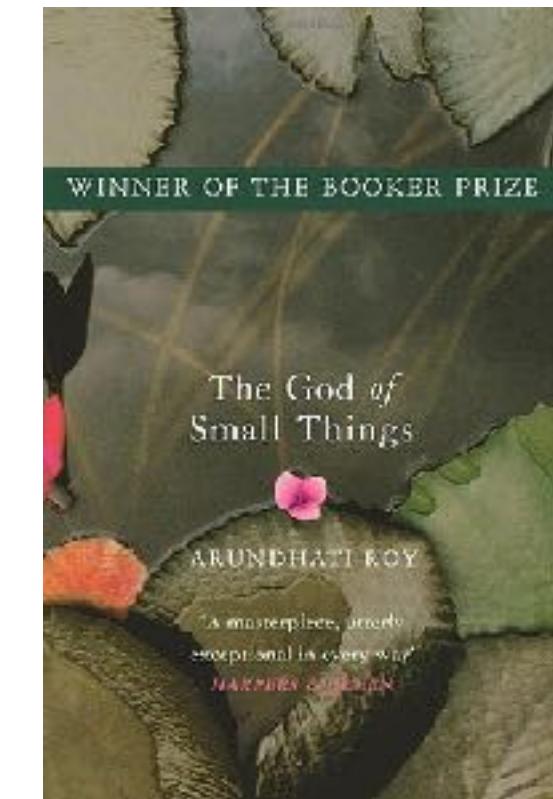
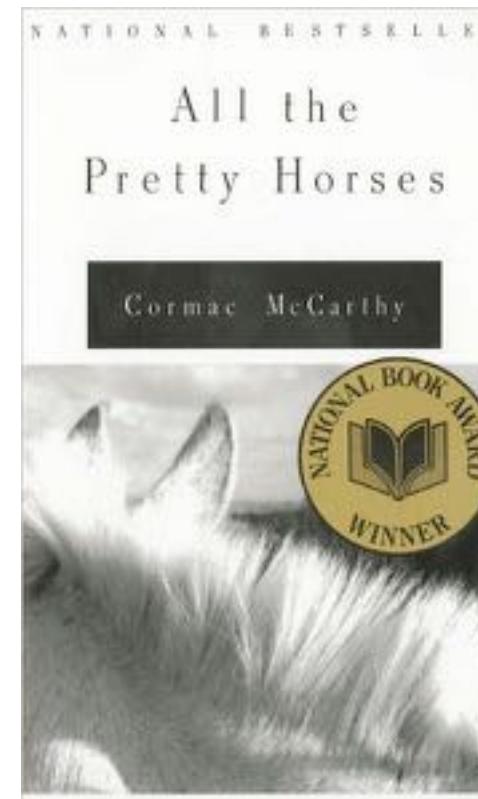
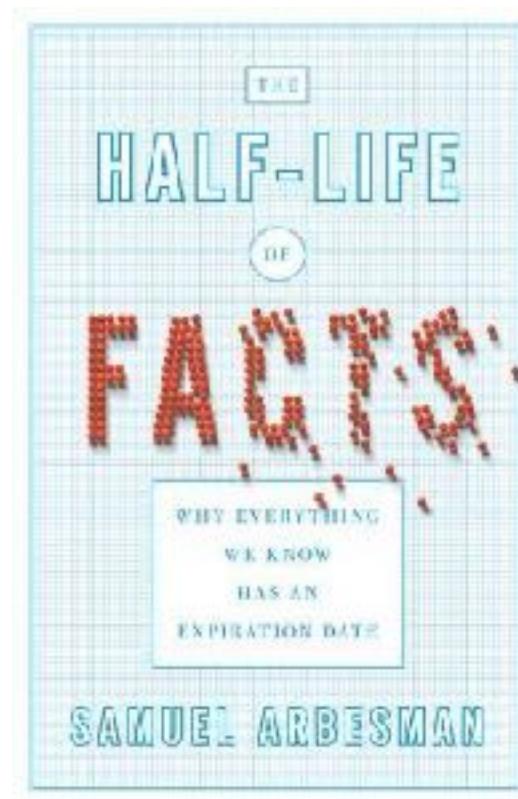
Albert-László Barabási



# Suppose I give you a book. Predict its sales.

Existing data: books and their sales.

1. turn books into feature vectors.

 $\vec{x}_1$  $\vec{x}_2$  $\vec{x}_3$  $\vec{x}_4$ 

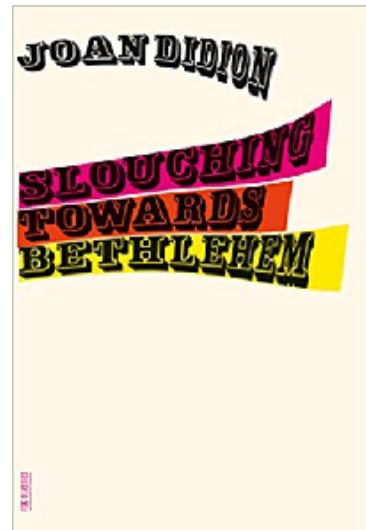
2. Train a model:

$$P(\text{book } i > \text{book } j \mid \vec{x}_i, \vec{x}_j, \theta)$$

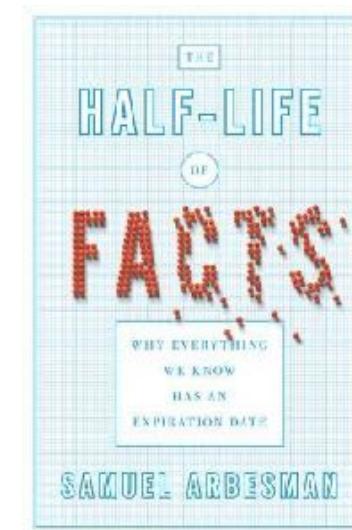
# Suppose I give you a book. Predict its sales.

Existing data: books and their sales.

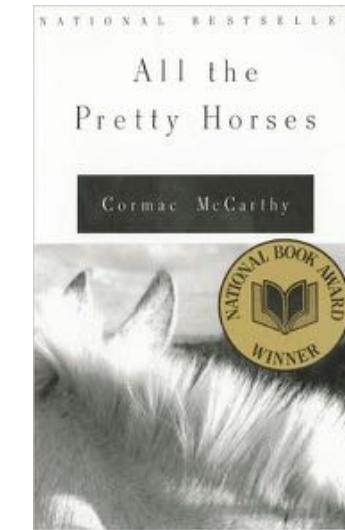
1. turn books into feature vectors.



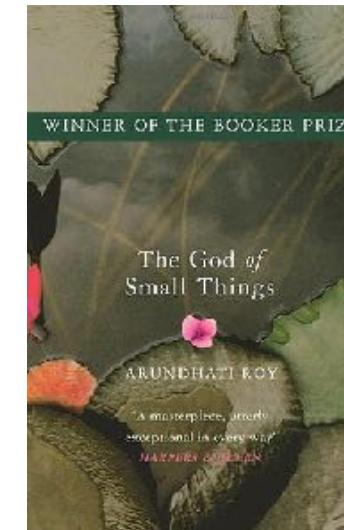
$$\vec{x}_1$$



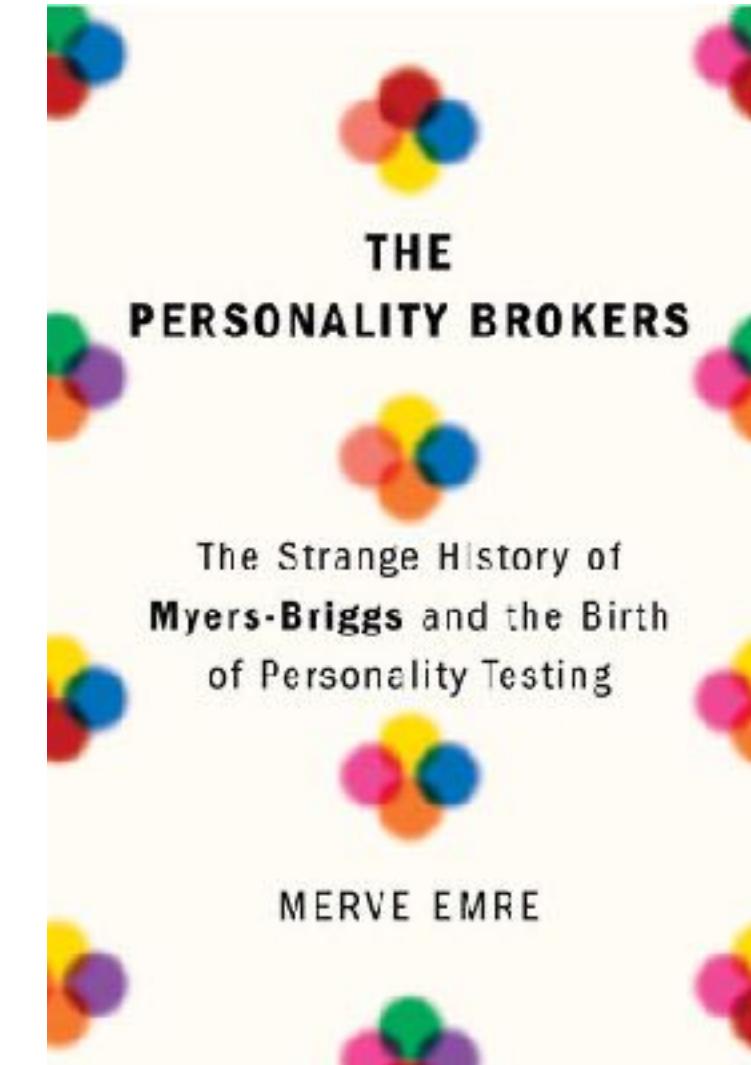
$$\vec{x}_2$$



$$\vec{x}_3$$



$$\vec{x}_4$$



$$\vec{x}_5$$

2. Train a model.

3. Use the model to simulate pairwise competitions.

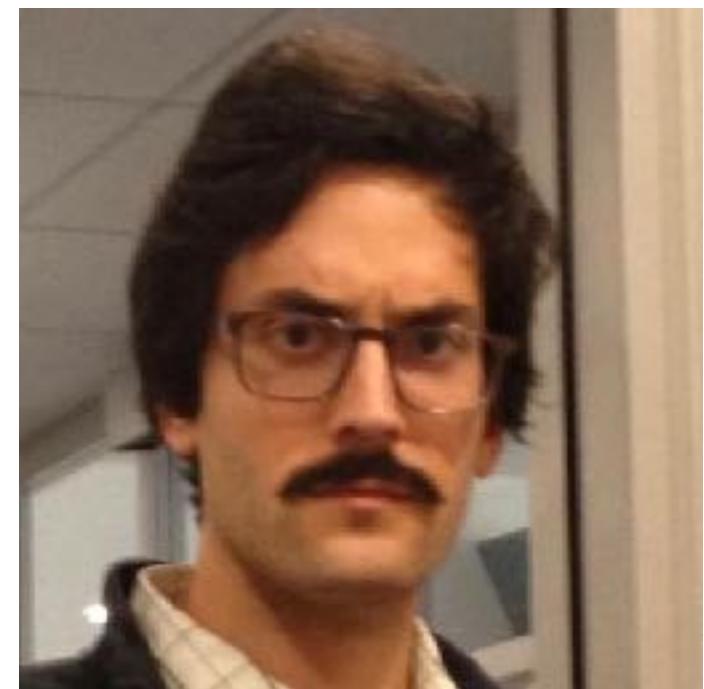
$$P(\text{book } i > \text{book } 5 \mid \vec{x}_i, \vec{x}_5, \theta)$$

4. Use [your favo(u)rite algorithm] to infer  $\text{rank}_5$  from pairwise comparisons.

# Rankings rankings

Area under the receiver-operator curve (AUC)

Method	AUC on Fiction	AUC on Biography
KNN	0.759	0.815
Cohen et al.	0.892	0.871
WTG wave	<b>0.910</b>	<b>0.892</b>
Pairwise + Voting	<b>0.915</b>	<b>0.891</b>
FAS-PIVOT	<b>0.907</b>	<b>0.892</b>
SpringRank	3  <b>0.908</b>	1  <b>0.893</b>



# But actually... many methods performed well!

Now the question is: why do the top four algorithms perform similarly?

What does that tell us about the **structure of the problem** and the **structure of the space** over which we are ranking?

Use the consistency of the ranking results across algorithms to learn about the system itself.

What does it mean for a space or problem or set to be easily ordered or rankable?

# Hierarchies of Desirability in Online Dating



# Hierarchies are encoded in language about courtship

## “She’s out of your league.”

asserts that:

1. Leagues or hierarchies of desirability exist.  
[How can we find them in data?]
2. The relative positions of individuals can be estimated.
3. Positions are predictive of something.  
[What behavior? And how noisy is the prediction?]

# Data: a popular online dating service

## **Characteristics of the online dating service:**

- Free.
- Around 4 million active self-identified heterosexual users.
- Ethnically diverse, urban, youngish.
- Approved and highly restricted data access through collaboration.

## **Characteristics of the data I work with:**

- One month of anonymized and timestamped messaging.
  - Self-identified heterosexual users with genders declared as M or F.
- 

## **How can we use messaging data to answer these questions?**

1. Do desirability hierarchies exist?
2. Are they predictive of behavior?

# Intuition: Something like PageRank?

SCIENCE ADVANCES | RESEARCH ARTICLE

SOCIAL SCIENCES

**Aspirational pursuit of mates in online dating markets**

Elizabeth E. Bruch<sup>1,2\*</sup> and M. E. J. Newman<sup>2,3</sup>

- Bruch & Newman 2018 used PageRank to analyze messaging network.
- Sorted individuals by percentile PageRank scores.
- Aspirational pursuit patterns: people message “up” the sorted network.

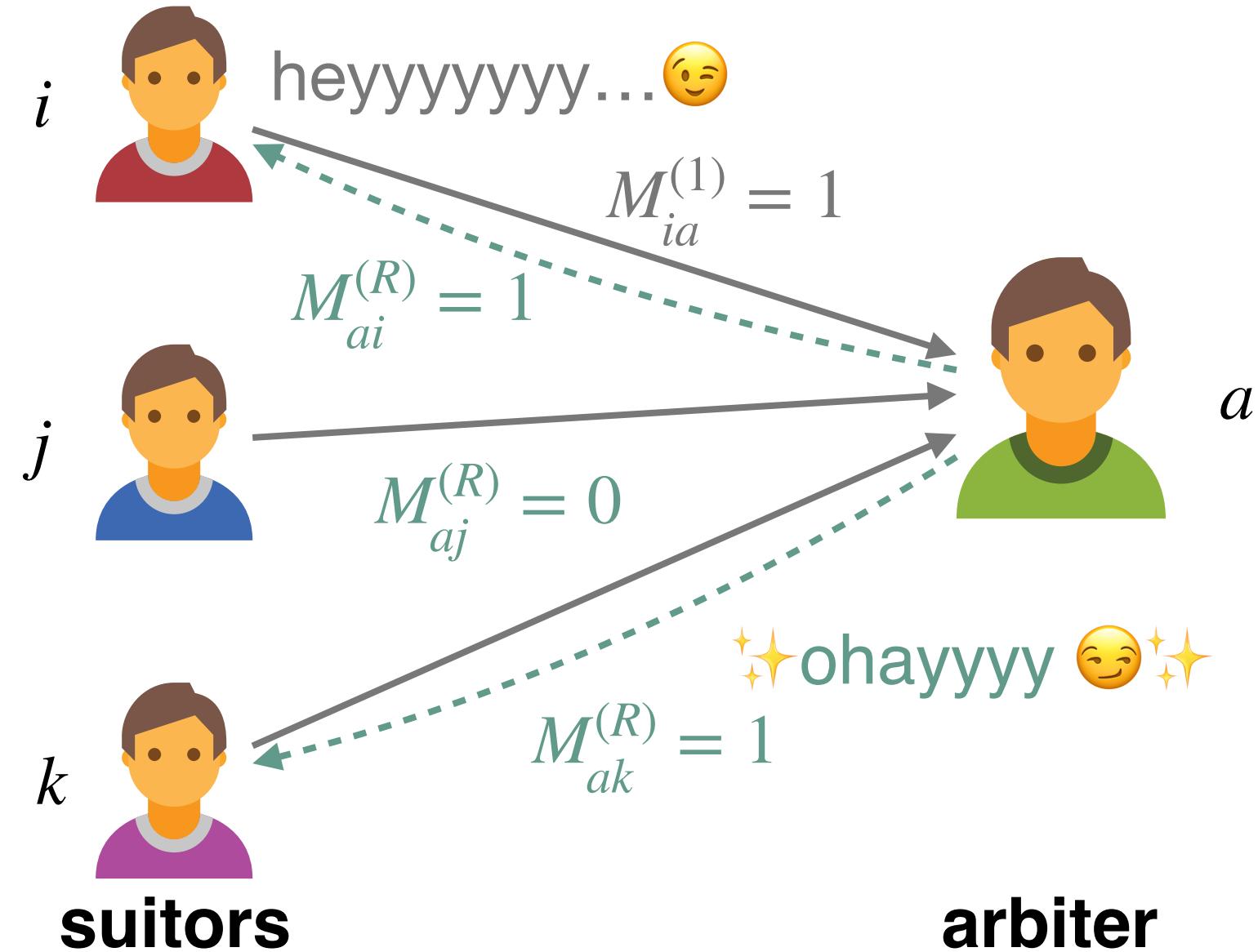
Recommended reading if you are interested in this subject!

Downside: PageRank scores aren't *usefully* interpretable.

[Stationary distrib. of random walk + teleportation, on a network of message passing?]

Can we use the messaging data to find more meaningfully interpretable ranks?

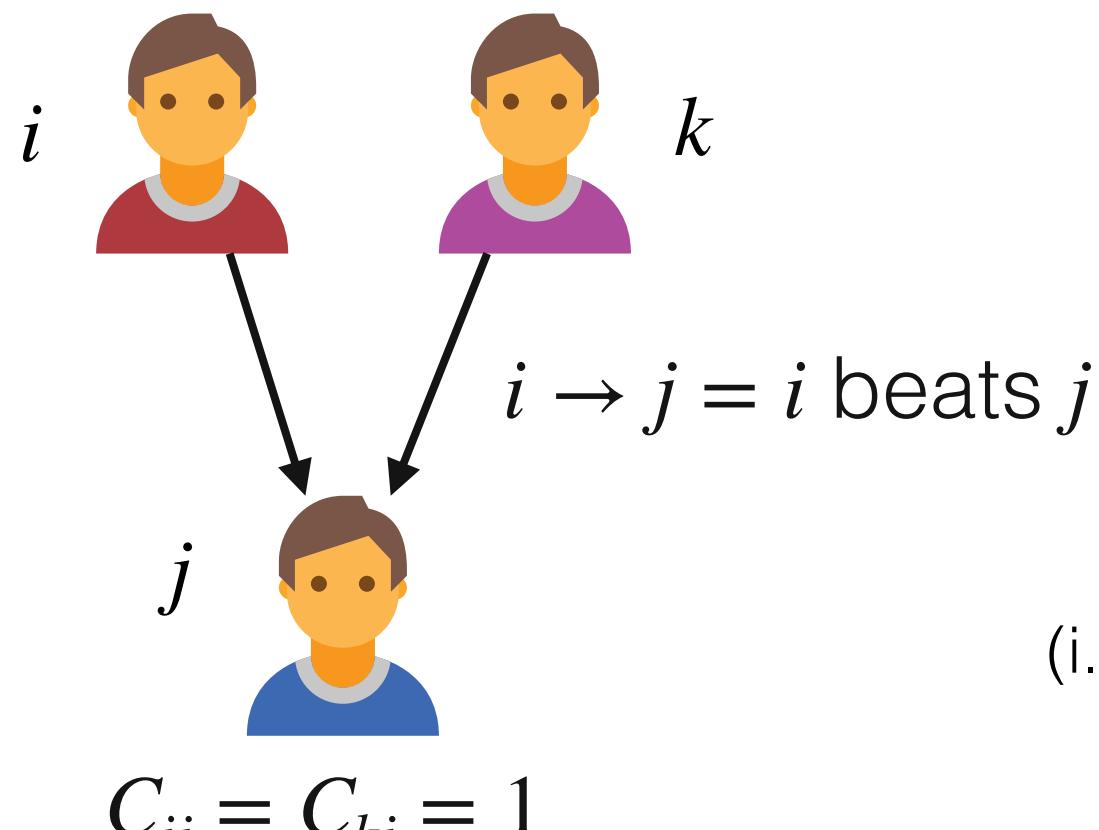
# Insight: messaging is a competition for attention



$M^{(1)}$ : network layer of first messages  
 $M^{(R)}$ : network layer of replies

Who won/lost this competition for attention?

# Hierarchies in the competition for attention



Let  $C_{ij}$  be the number of times that  $i$  outcompetes  $j$ .

$$C_{ij} = \sum M_{ia}^{(1)} M_{ja}^{(1)} M_{ai}^{(R)} \left( 1 - M_{aj}^{(R)} \right)$$

arbiters  
(receivers)

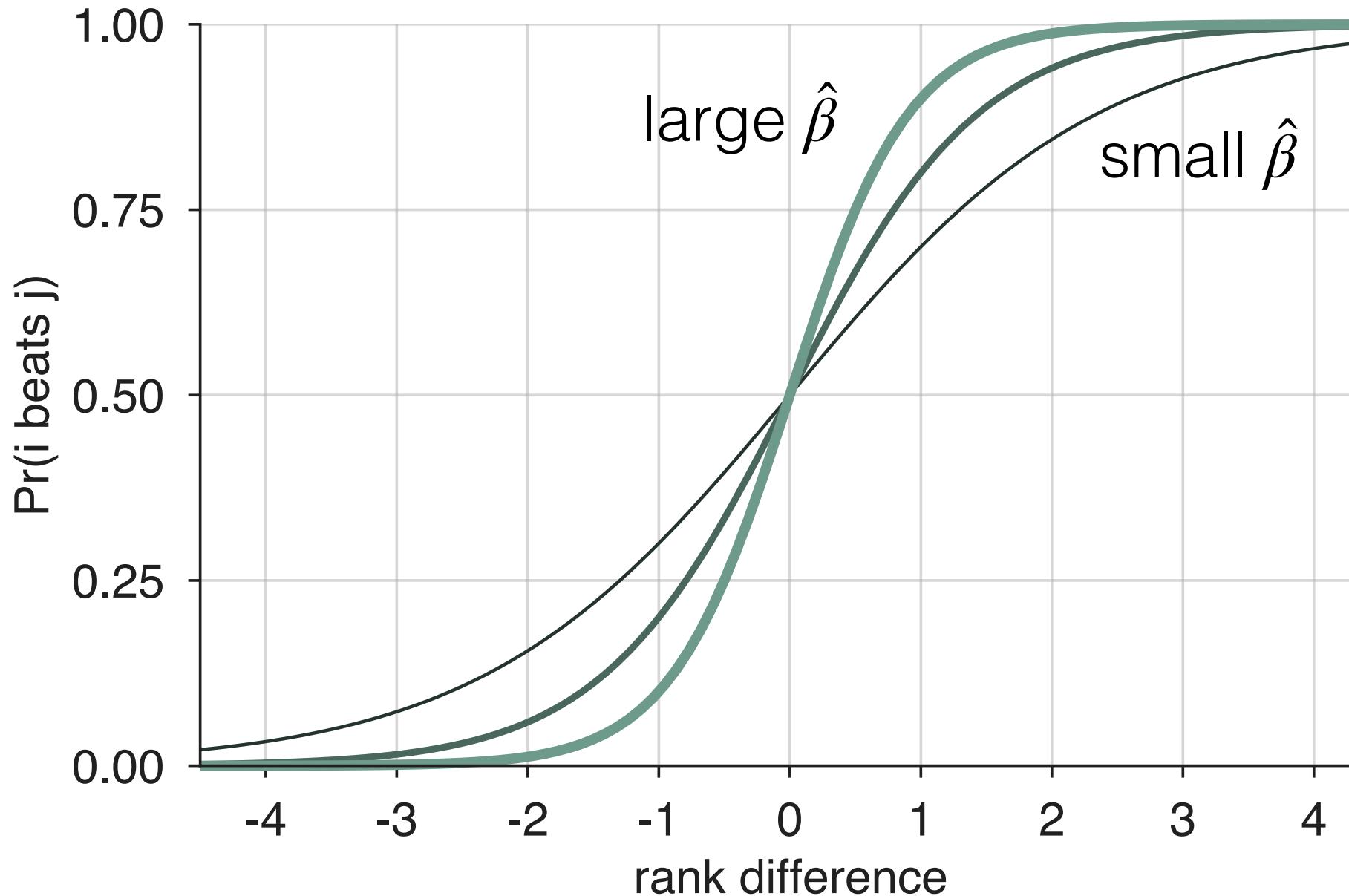
$i \rightarrowtail a$      $i \rightarrowtail a$      $j \rightarrowtail a$

$j \rightarrowtail a$

Each arbiter is presented with a choice set of suitors.  
Collectively, arbiters' choices provide many *partial orderings* of suitors

$C$  is a one-component, directed network representing pairwise comparisons.  
We will use **SpringRank** to find people's latent positions to *predict future behavior*.

# Beta tells us how to interpret rank differences



$$P(i \rightarrow j \mid i \leftrightarrow j) = \frac{1}{1 + e^{-2\hat{\beta}(c_i - c_j)}}$$

$\hat{\beta}$  is the MLE inverse temperature  
of the SpringRank Boltzmann

It tells us the sensitivity/scale for predictions in the ranking space.

# The Depths of Leagues

Dear Editors,

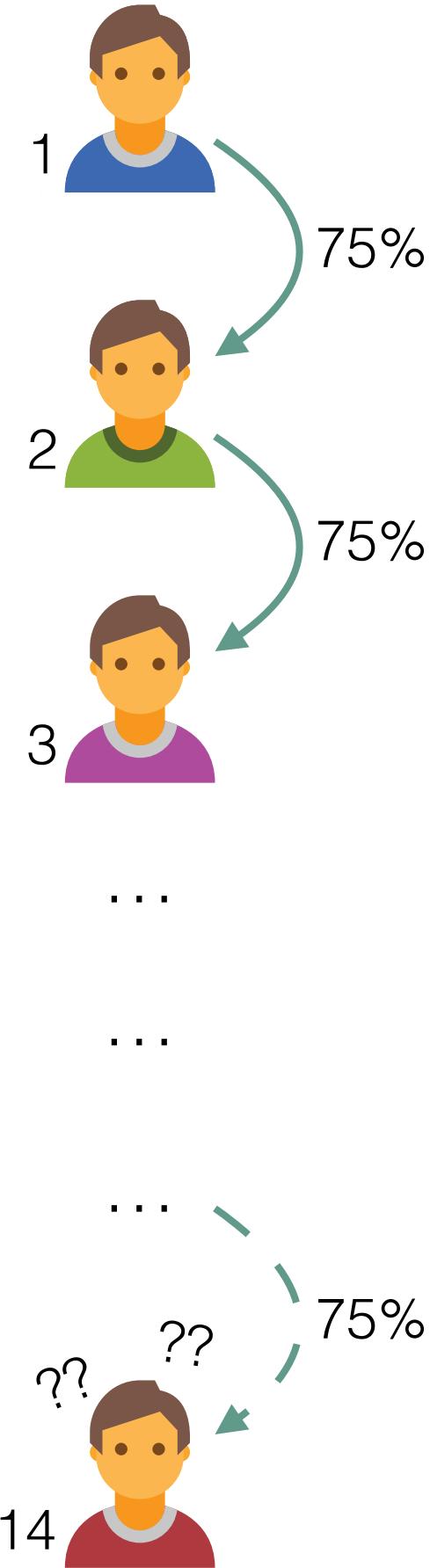
There are some questions I would like to ask. Firstly, how complex is backgammon compared to other games of skill such as chess or bridge?

*Let's start with chess, which has evolved a well-developed rating system over the past 40 years. Chess ratings range from a high of about 2800 to theoretical lows of about 0 (a complete beginner who has just learned the moves). Chess ratings are also designed so that a 200-point rating difference between two players anywhere on the scale means that the higher-rated player has a 70-75% chance of defeating a lower-rated player (discounting draws, which are possible in chess but not in most of the other games we'll consider).*

*Now consider the following experiment:*

- (1) Take the best player in the world (in the case of chess, it's Gary Kasparov). Call him player 1.
- (2) Find someone that the best player beats 70-75% of the time. Call him player 2.
- (3) Call the difference between players 1 and 2 one skill differential.
- (4) Find someone that player 2 can beat 70-75% of the time. Call him player 3. The difference between players 2 and 3 is another skill differential.
- (5) Continue this process until you have taken the chain down to an absolute beginner.
- (6) Count the number of skill differentials involved. This is the **complexity number** of the game.

*In the case of chess, this number is about 14.*



# The Depths of Leagues

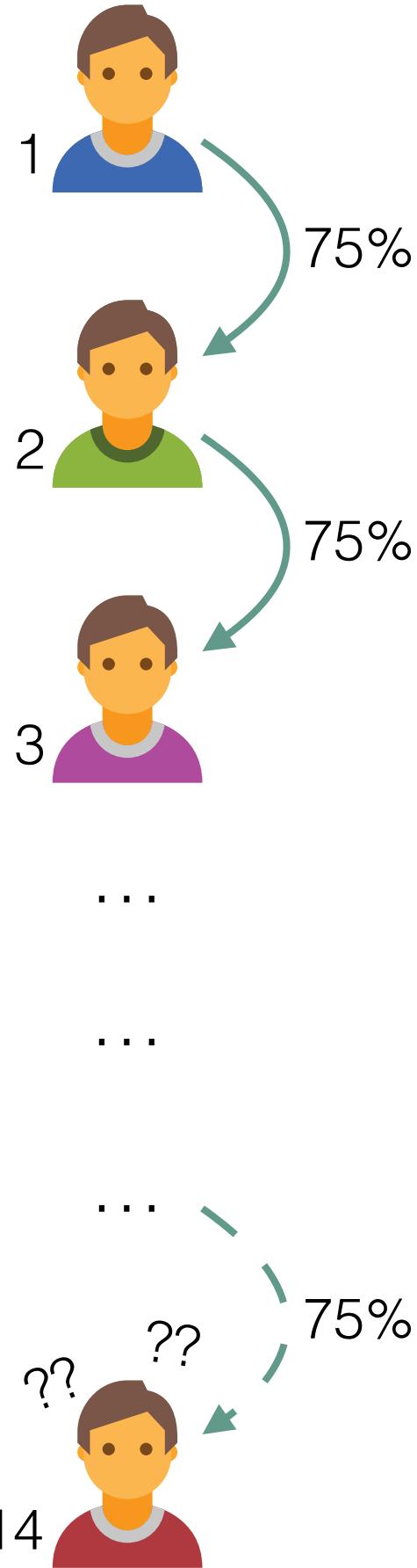
## COMPLEXITY NUMBERS

Go	40
Chess	14
Scrabble	10
Poker	10
Backgammon	8
Checkers	8
Hearts	5
Blackjack	2
Craps	0.001
Lotteries	0.0000001
Roulette	0
<b>Online Dating</b>	<b>???</b>

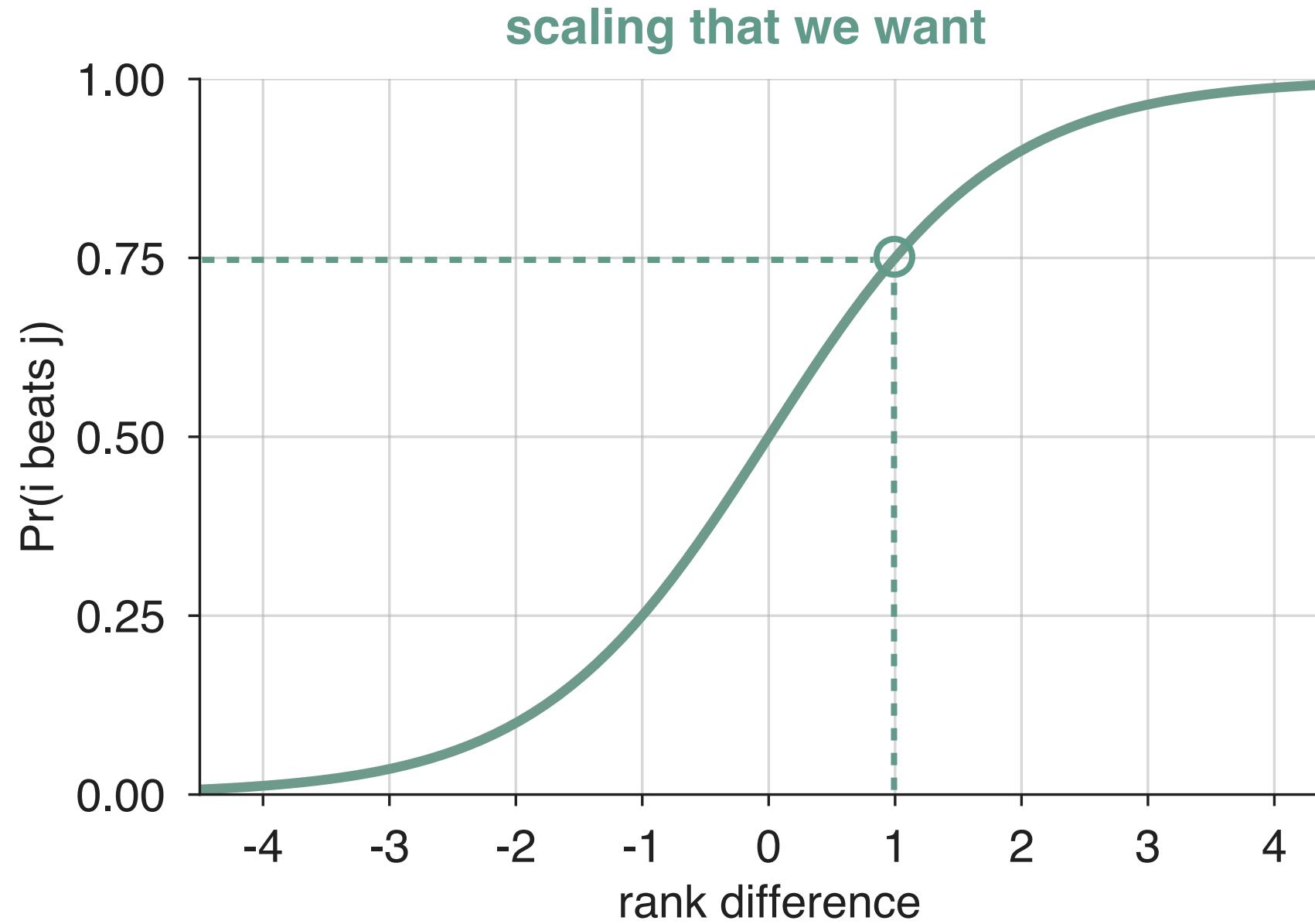
Now consider the following experiment:

- (1) Take the best player in the world (in the case of chess, it's Gary Kasparov). Call him player 1.
- (2) Find someone that the best player beats 70-75% of the time. Call him player 2.
- (3) Call the difference between players 1 and 2 one skill differential.
- (4) Find someone that player 2 can beat 70-75% of the time. Call him player 3. The difference between players 2 and 3 is another skill differential.
- (5) Continue this process until you have taken the chain down to an absolute beginner.
- (6) Count the number of skill differentials involved. This is the **complexity number** of the game.

In the case of chess, this number is about 14.



# Choosing a scale for interpretability



scaling that we have

$$P(i \rightarrow j \mid i \leftrightarrow j) = \frac{1}{1 + e^{-2\hat{\beta}(c_i - c_j)}}$$

enforce the desired scale,  
then solve for a rescaling constant  $k$

$$\bar{c}_i - \bar{c}_j = 1, \quad 0.75 = \frac{1}{1 + e^{-2\hat{\beta}(k\bar{c}_i - k\bar{c}_j)}}$$
$$k = \frac{\log \text{odds } 0.75}{2\hat{\beta}}$$

# Reviewing the approach.

1

$M^{(1)}$ : network layer of first messages

$M^{(R)}$ : network layer of replies

2

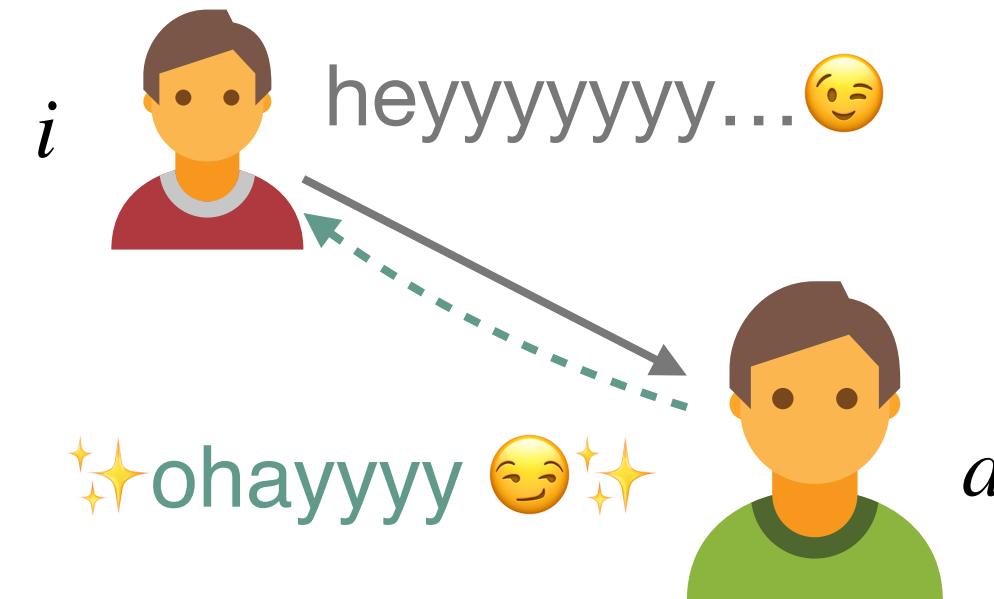
$$C_{ij} = \sum_a M_{ia}^{(1)} M_{ja}^{(1)} M_{ai}^{(R)} \left( 1 - M_{aj}^{(R)} \right)$$

3

$$H(c) = \frac{1}{2} \sum_{ij} C_{ij} \left( c_i - c_j - 1 \right)^2$$

4

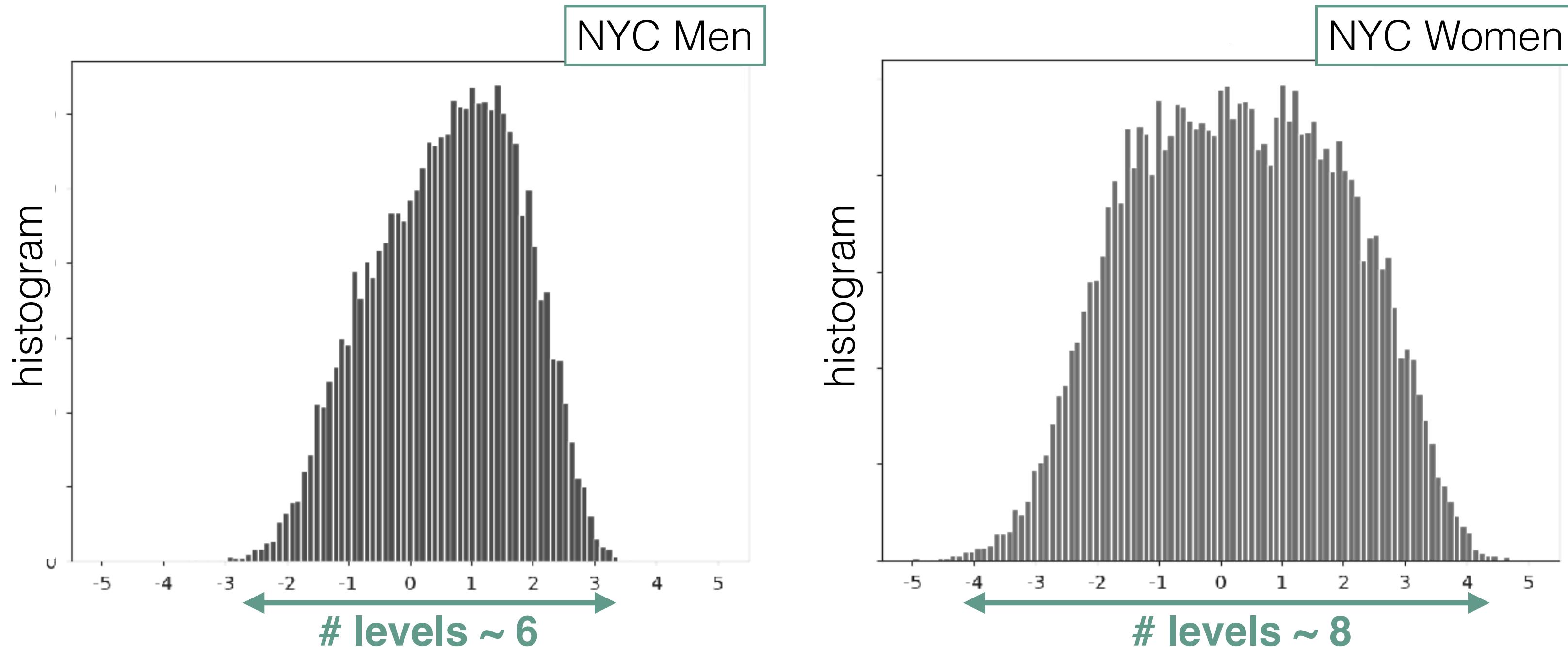
$$\bar{c} = c \frac{\log \text{ odds } 0.75}{2\hat{\beta}}$$



## Result:

An embedding of individuals in a linear hierarchy, such that a one-unit difference predicts a 75% “win rate” in the competition for attention.

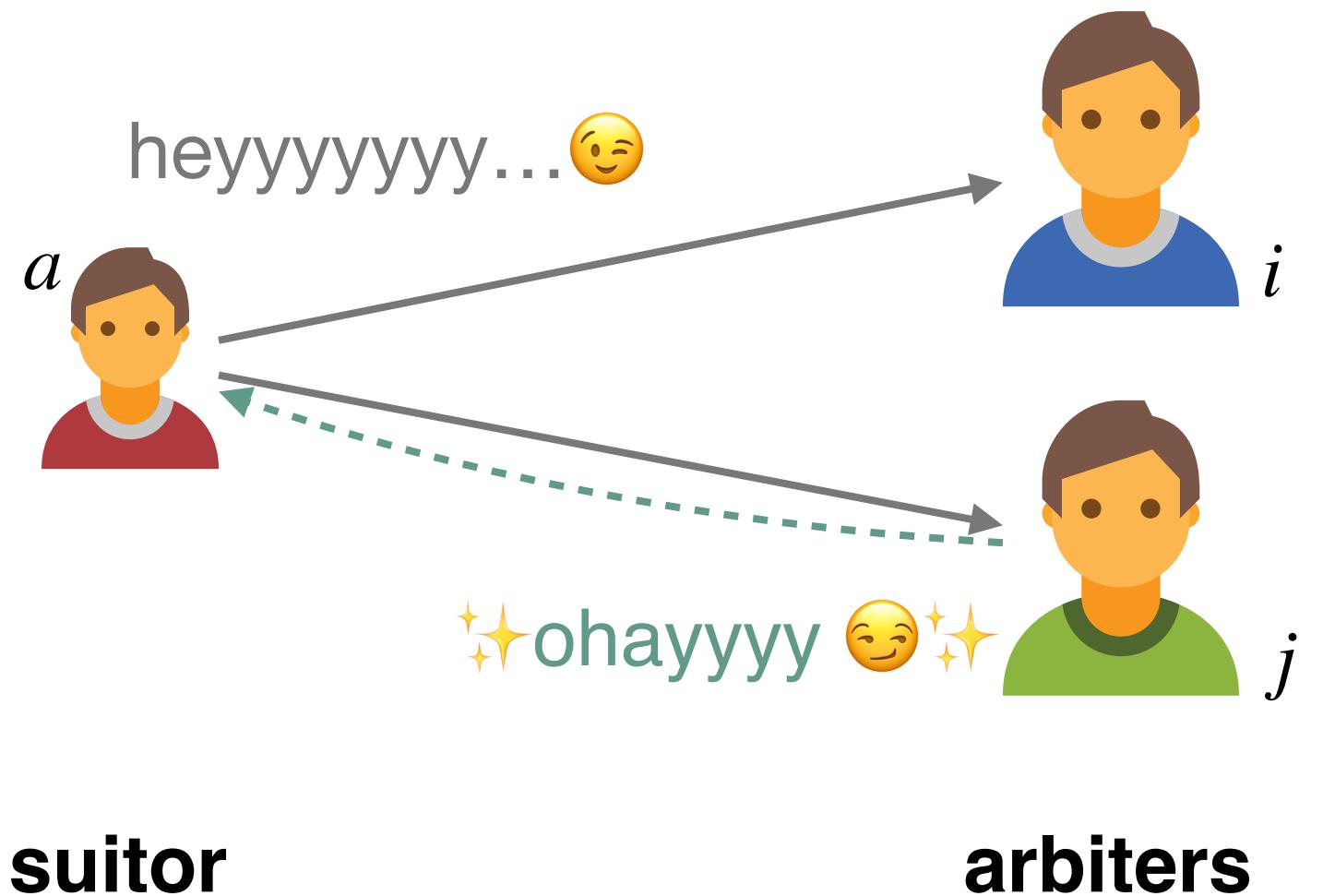
# Real data: New York City rank distributions



- Those who had 100% or 0% reply rates not shown.
- NYC Women's competition space has a deeper "strategic complexity" than NYC Men's.

# Message responses also reveal arbiter preferences!

Which arbiter is more **selective**?



1  $M^{(1)}$ : network layer of first messages

$M^{(R)}$ : network layer of replies

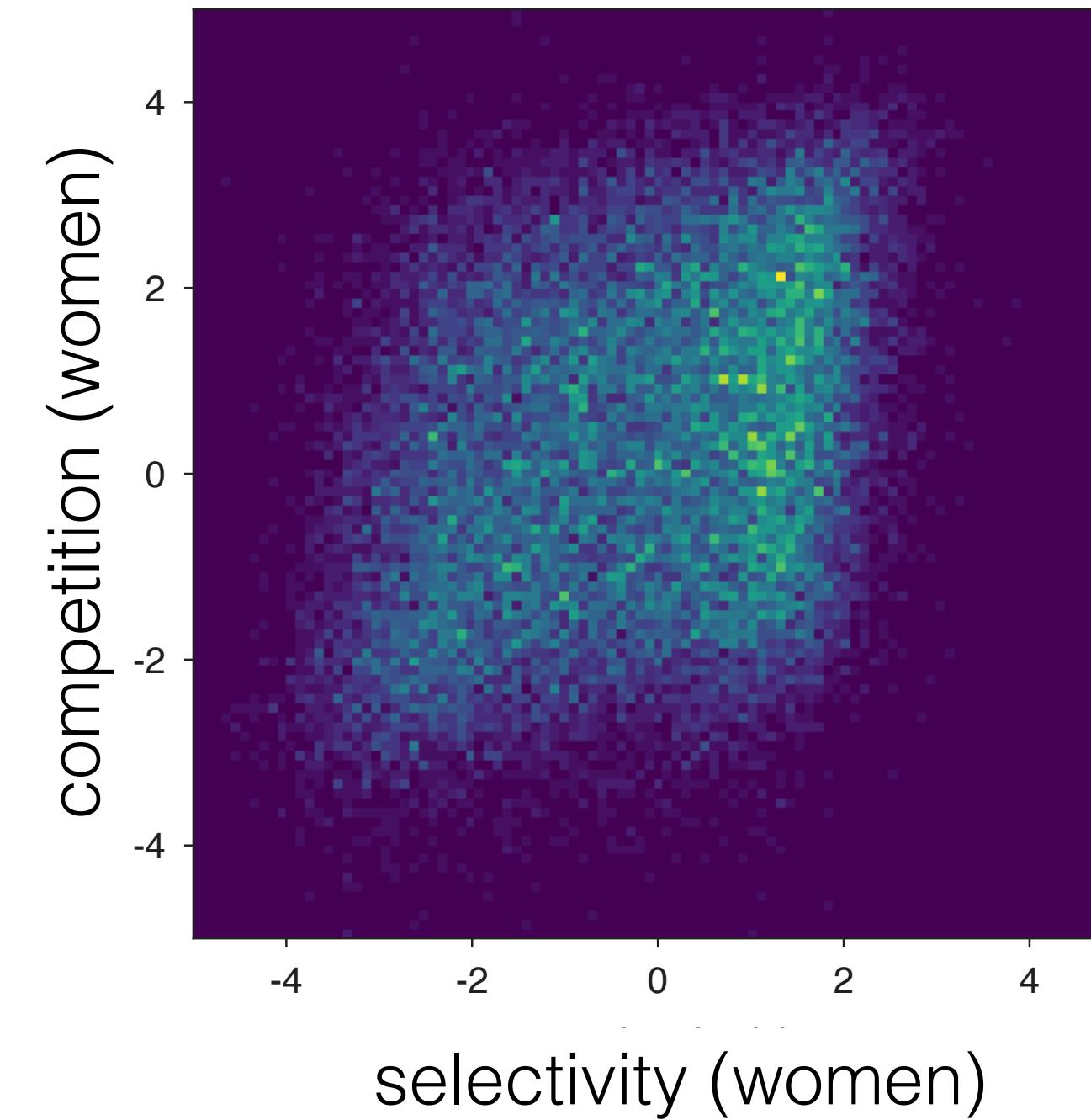
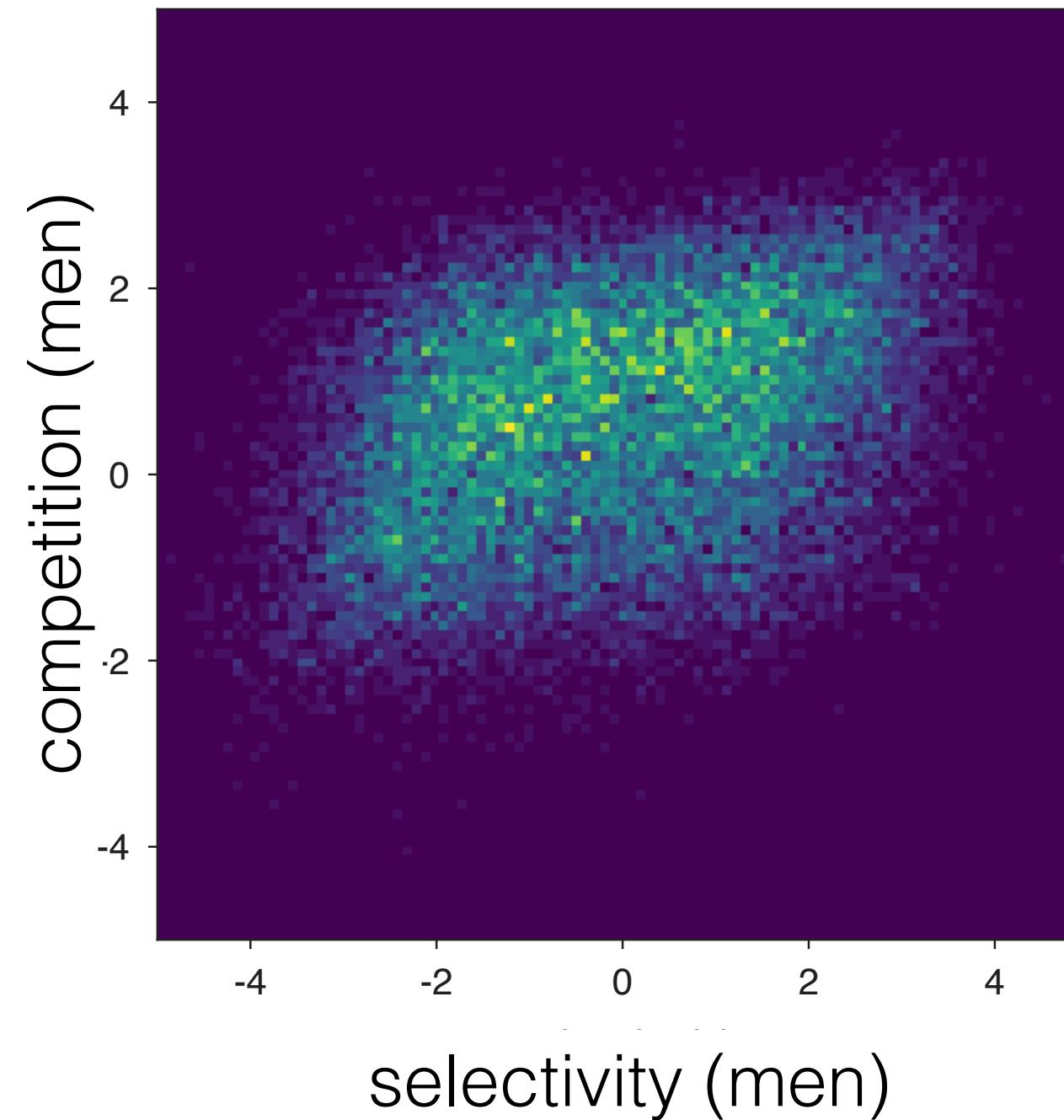
2  $S_{ij} = \sum_a M_{ai}^{(1)} M_{aj}^{(1)} M_{ja}^{(R)} (1 - M_{ia}^{(R)})$

3  $H(s) = \frac{1}{2} \sum_{ij} S_{ij} (s_i - s_j - 1)^2$

4  $\bar{s} = s \frac{\log \text{ odds } 0.75}{2\hat{\beta}}$

**Result:** selectivity  $s$  reflects your perception of yourself. (internal)  
competition  $c$  reflects others' perception of you. (external)

# How well do internal & external perceptions match?



- Spearman and Pearson correlations for both plots are around 0.34.
- The ability to get replies (competition) is only weakly predictive of replying (selectivity).

# Conclusions, Outlook

1. Desirability hierarchies can be inferred via a linear embedding of networks that encode (c) competition and (s) selectivity.
2. Competition & Selectivity (external & internal revealed rank) are correlated, not identical. Why? Will this replicate for homosexual users?
3. This is just one example of a two-sided mutual ranking system.  
cf: high schoolers & colleges; postdocs & faculty jobs.

---

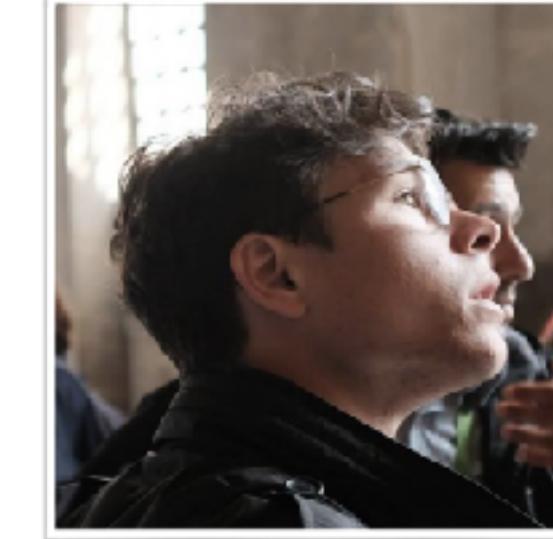
**Currently in the kitchen:** How do preferences over self-declared demographics (age, race/ethnicity, education) interact with desirability?  
And how do such interactions vary across US cities?



Elizabeth Bruch  
[empirical analyses]



Swapnil Gavade



K. Hunter Wapman  
[simulations]



SCIENCE ADVANCES | RESEARCH ARTICLE

SOCIAL SCIENCES

## Aspirational pursuit of mates in online dating markets

Elizabeth E. Bruch<sup>1,2\*</sup> and M. E. J. Newman<sup>2,3</sup>

Recommended reading if you are interested in this subject!



Hierarchy and cognition

What are the mechanisms that create large-scale patterns from many small interactions?

**Confront models with data** to reveal cognitively-accessible social mechanisms in parakeets.





1

→ aggressor (winner)

→ target (loser)

○ gawking, staring, shameless witnesses!



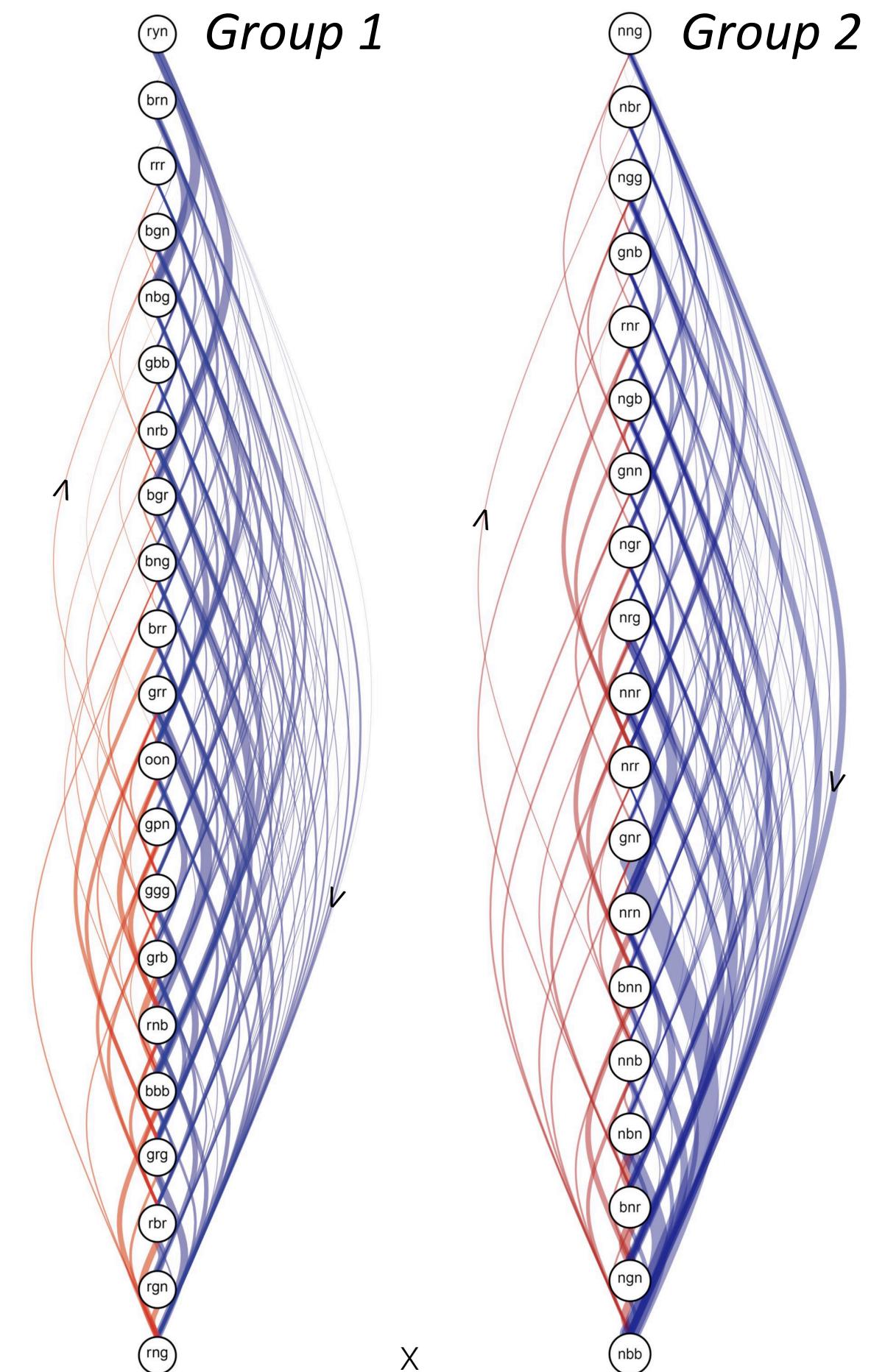
2

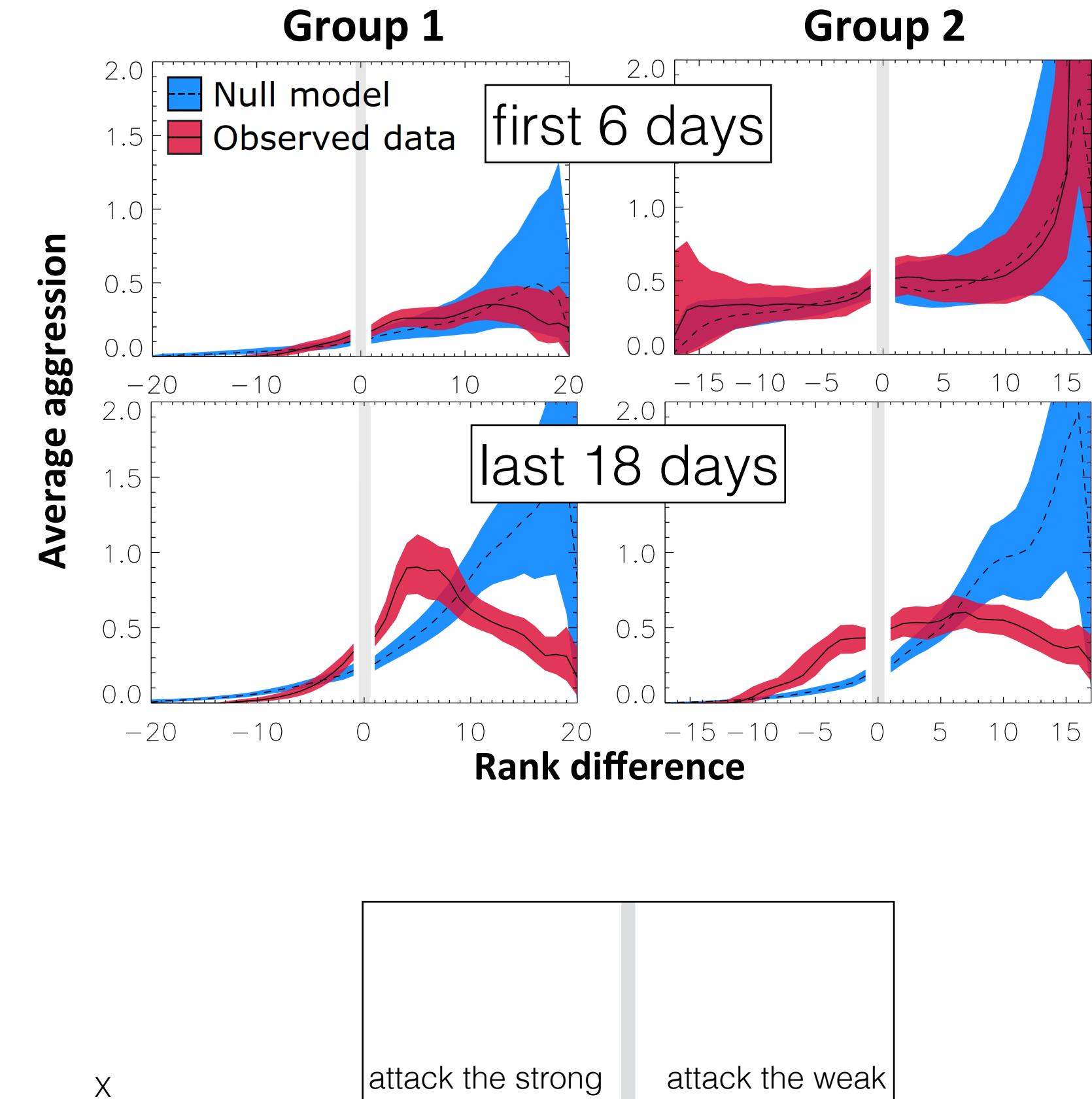
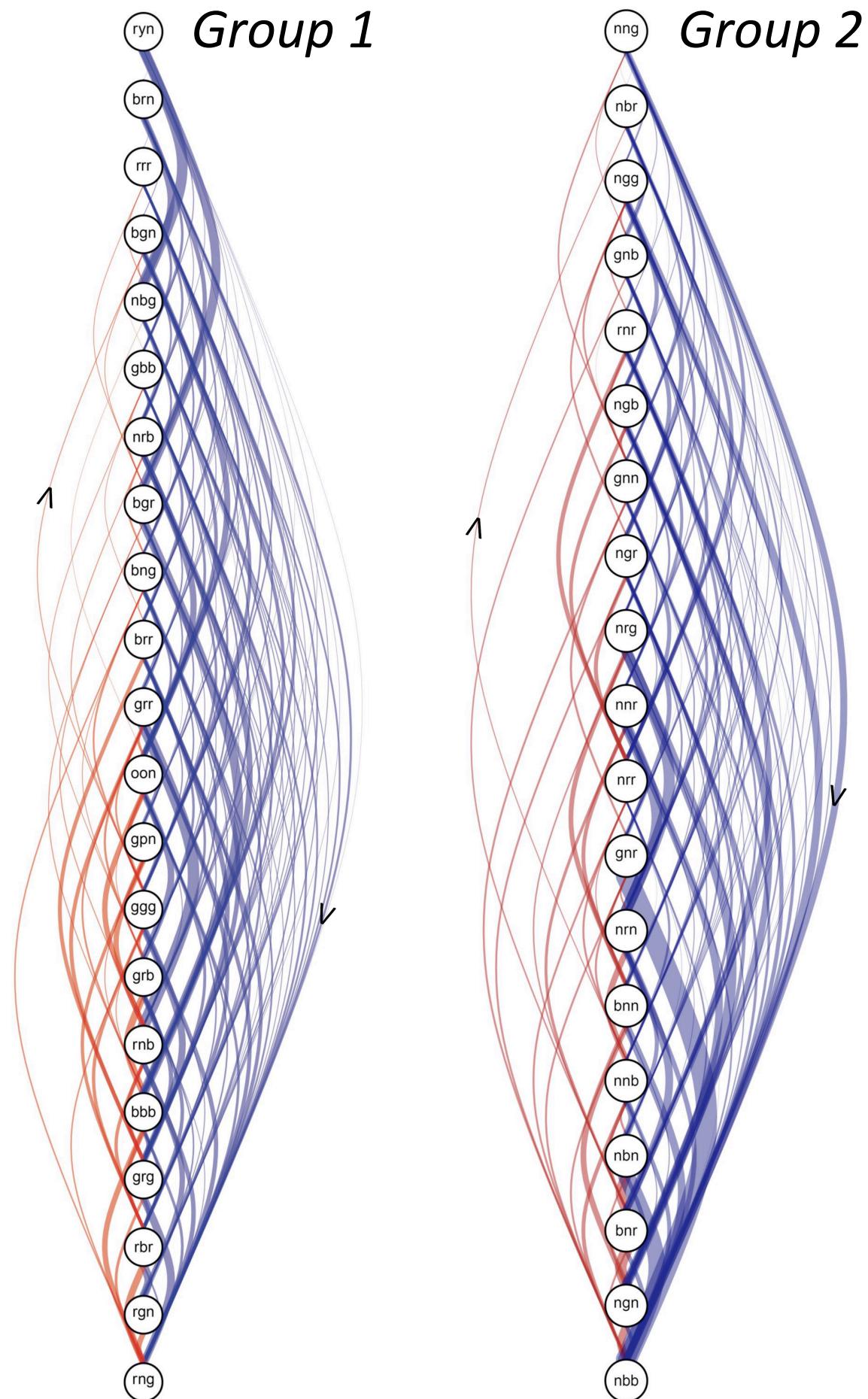


3

x

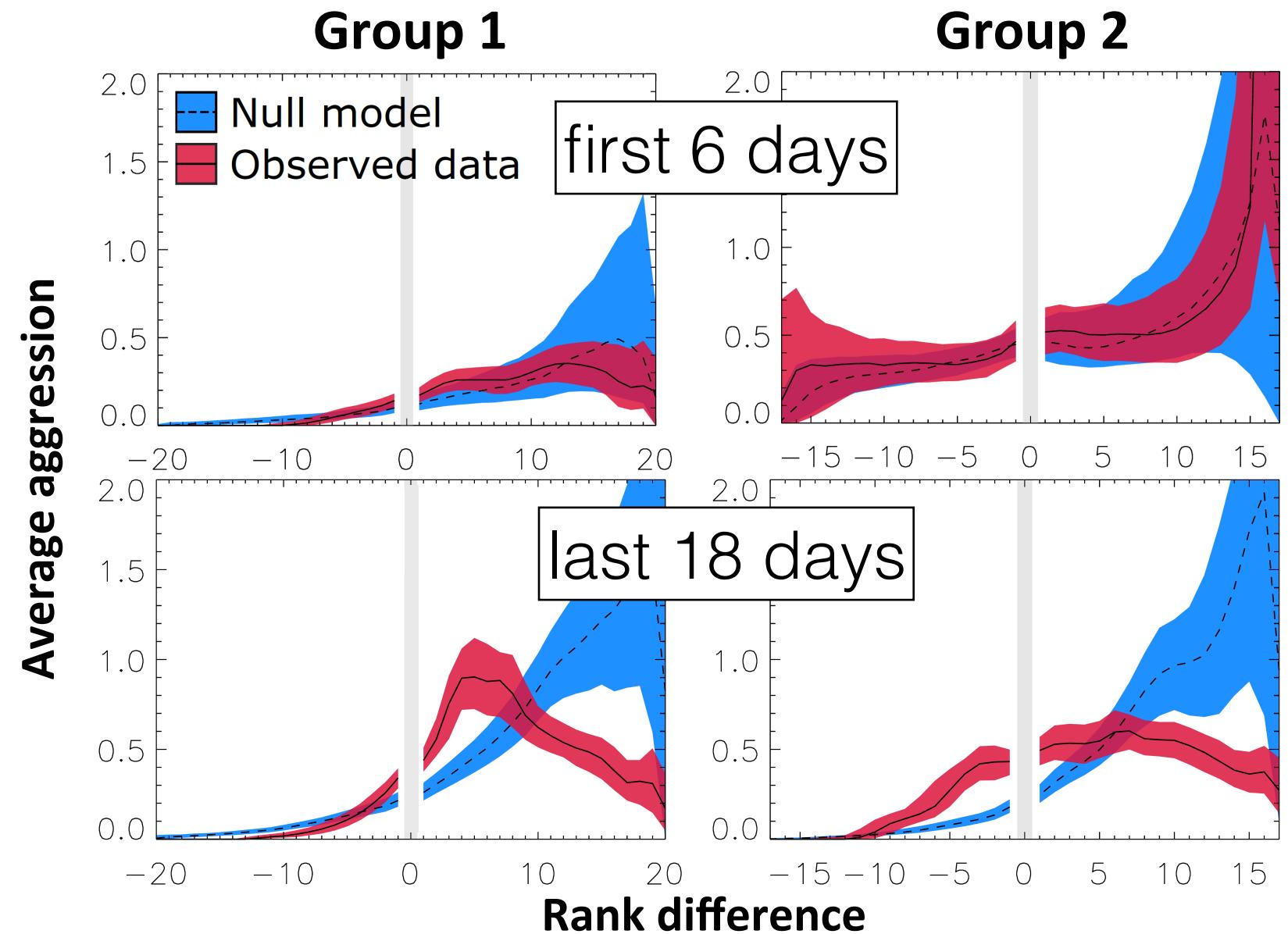
elapsed time < 1 second





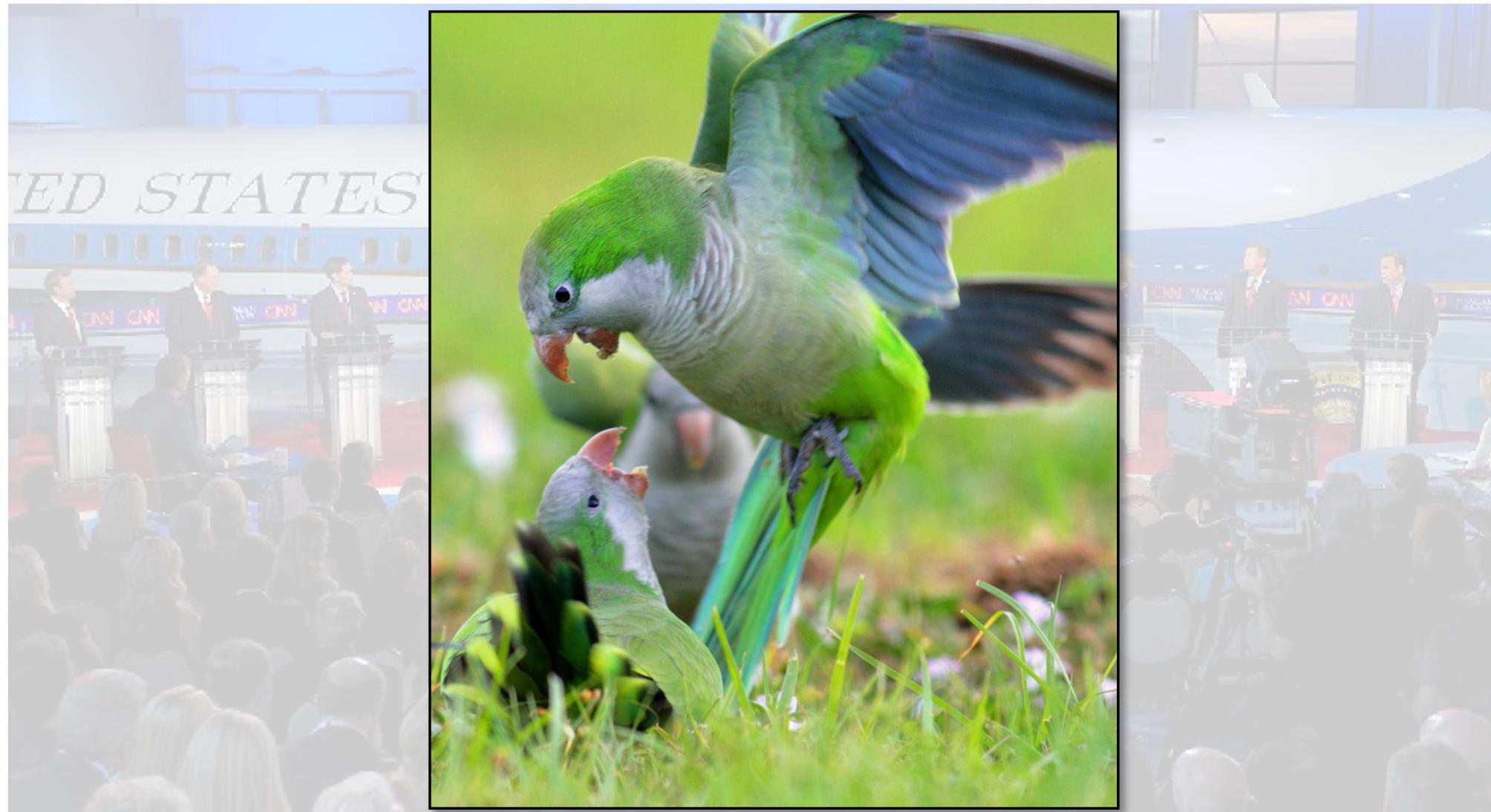
Parakeets know which individuals are ranked above and below themselves.

Parakeets know their own rank and the ranks of others.



Confronting models, which incorporate different complexities of bird-knowledge, with meticulous data, reveals clues about mechanisms of hierarchy formation.

# Complex models reveal complex behaviors



## Pile-on

Target the most recent loser.

[kick 'em while they're down]

## Pass-along

Target lower-ranked after losing.

[hurt people hurt people]

## Opportunism

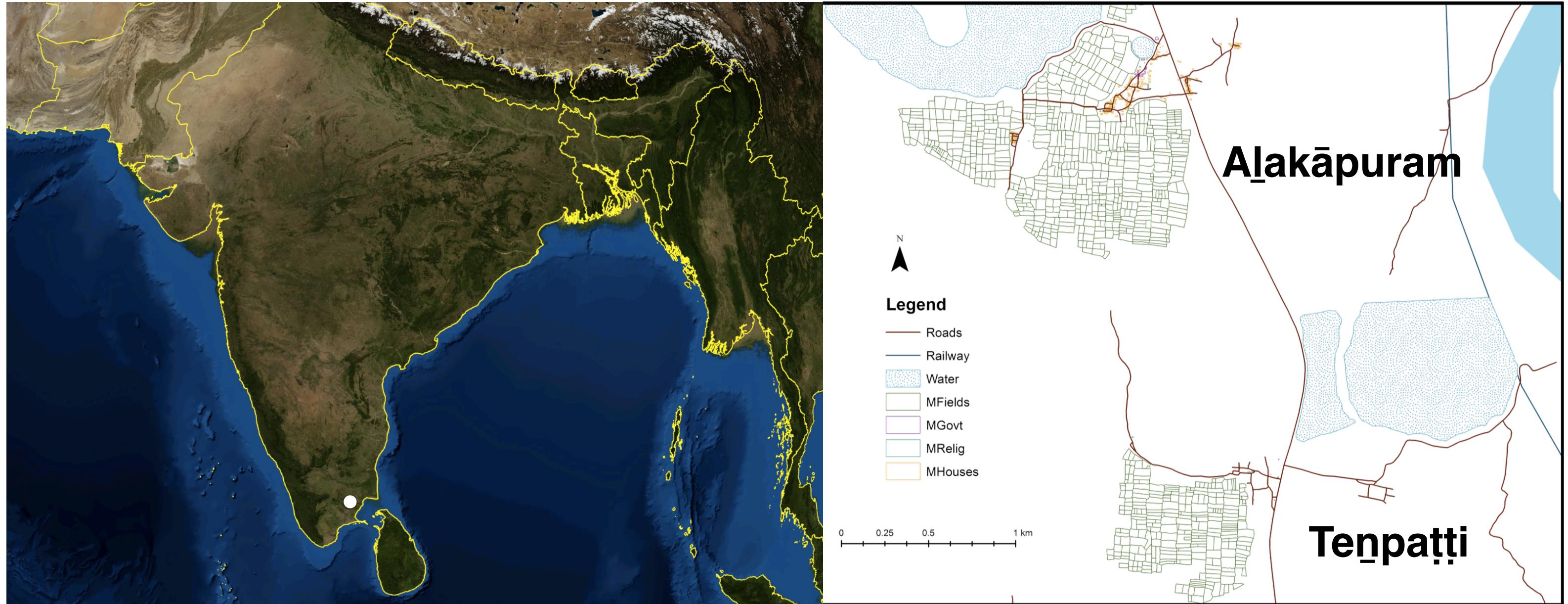
Target a recent loser of higher rank.

[now's my chance!]



# Groups and ordered structures

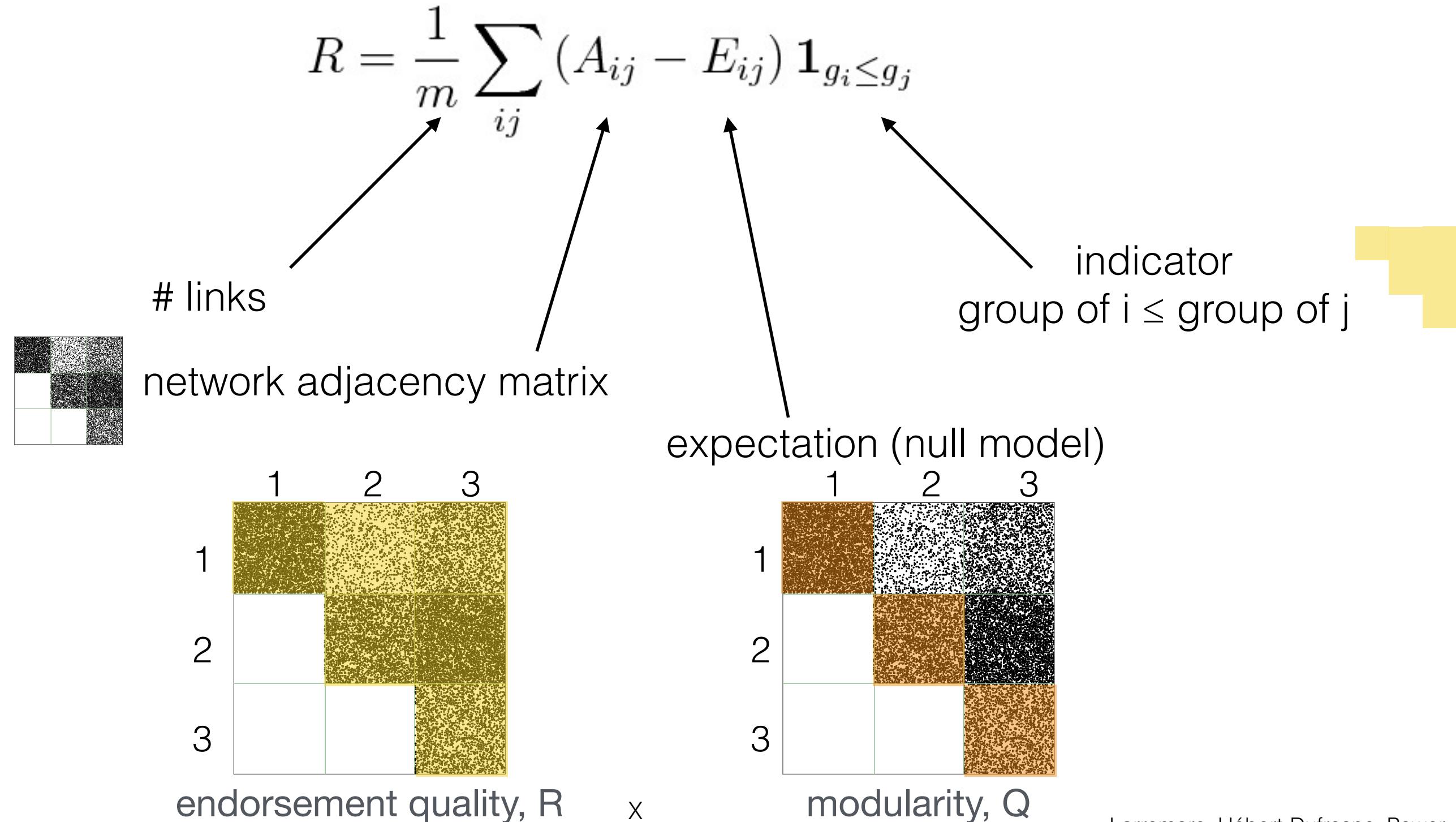
# South Indian networks: Tenpatti and Alakāpuram



1964 question of Srinivas and Béteille: beyond ethnographic investigations?

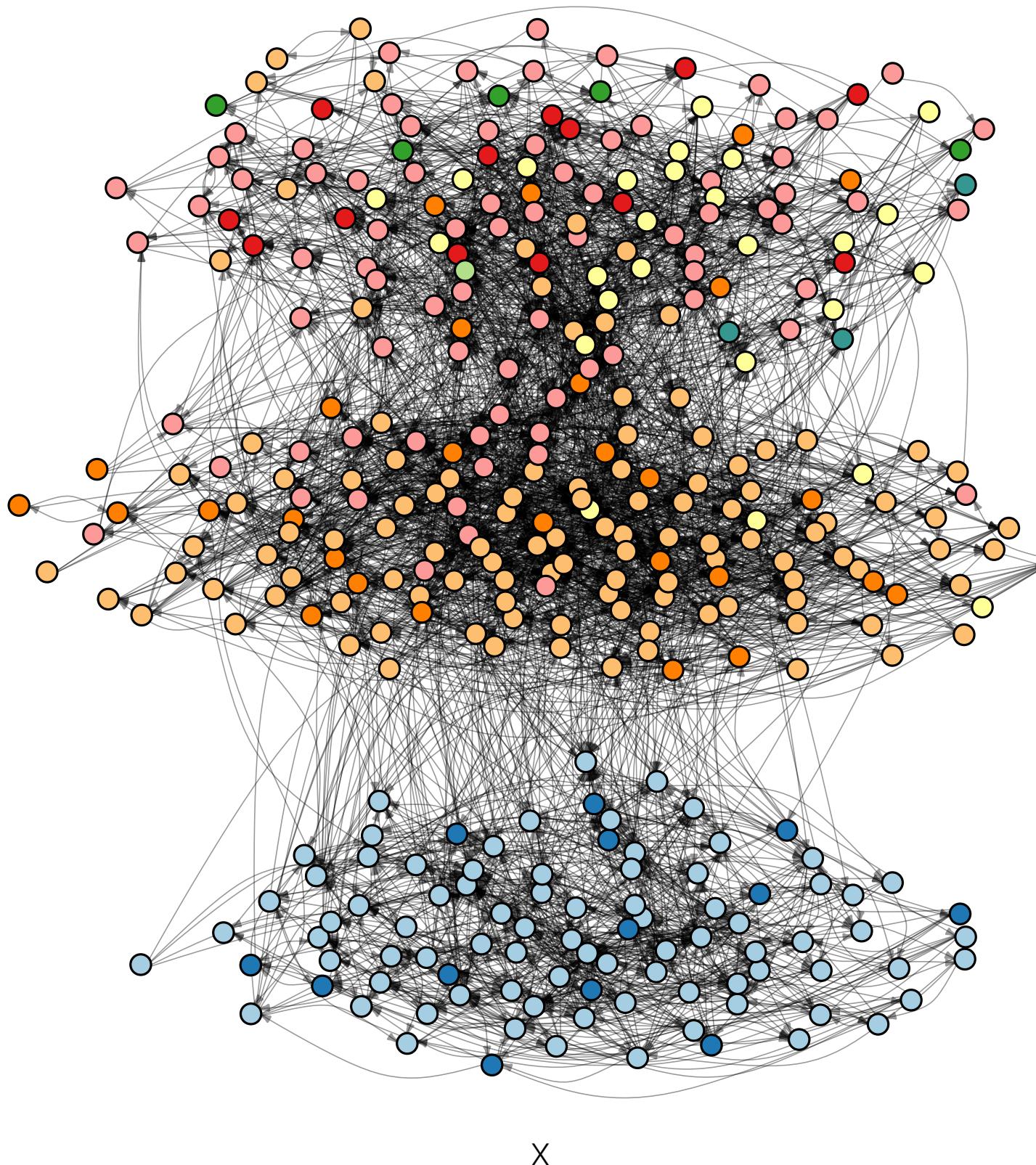
# Ranked order quality, $R$

We propose to measure the **quality** of a ranked ordering by  $R$



# Tenpatti

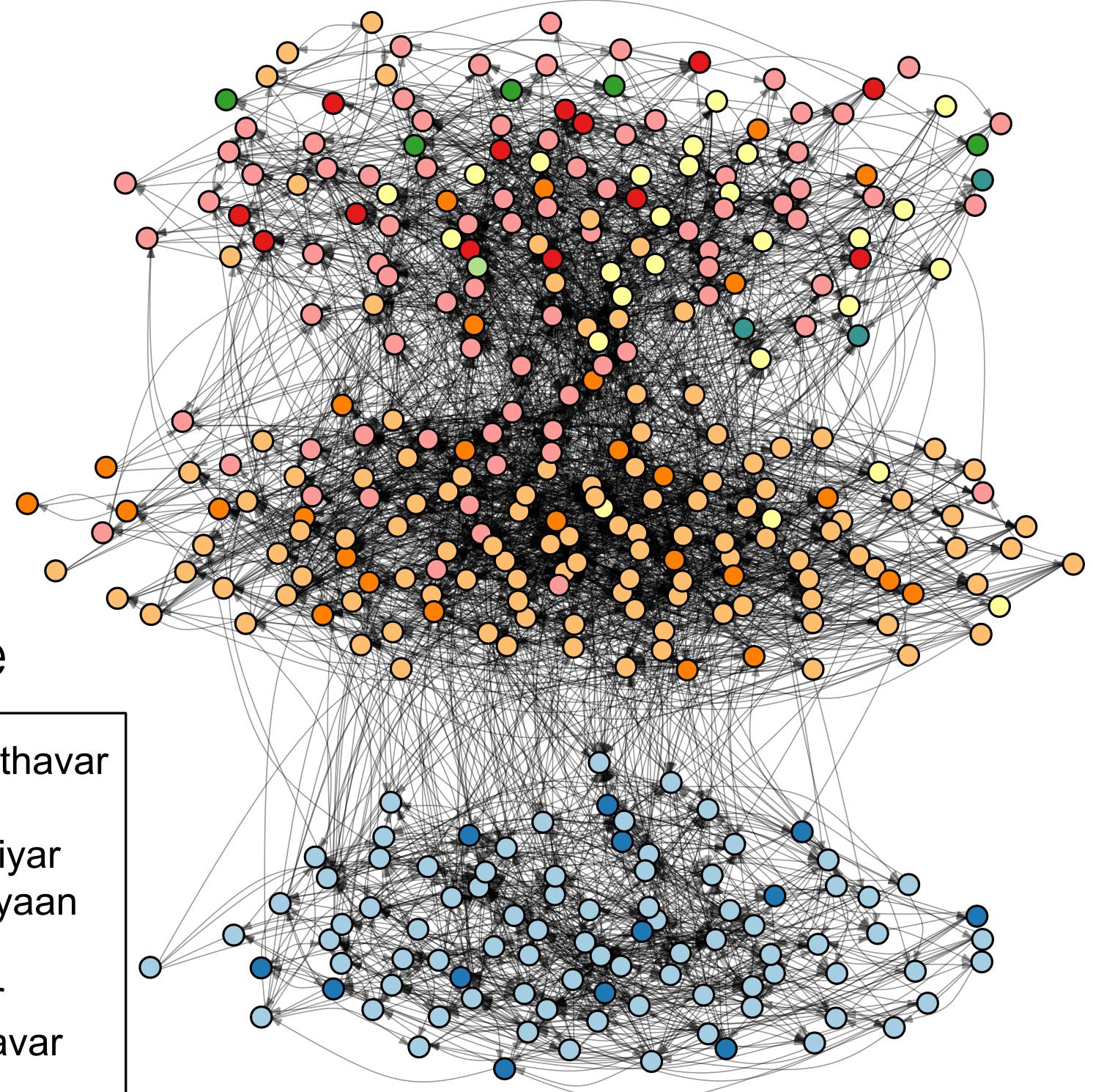
Caste



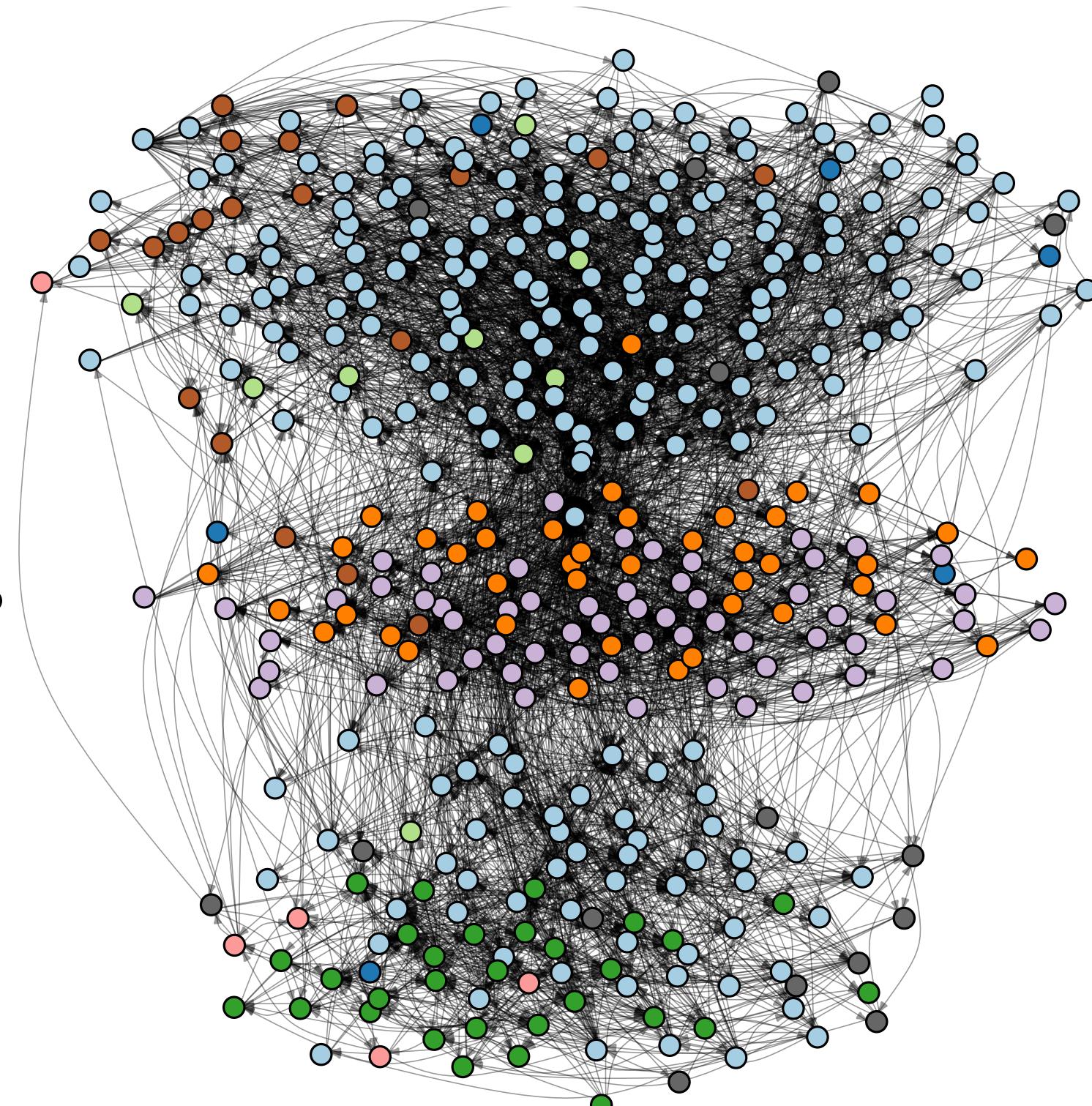
Scheduled castes—*dalit*  
“untouchable”

# Tenpaṭṭi

Caste



# Alakāpuram



X



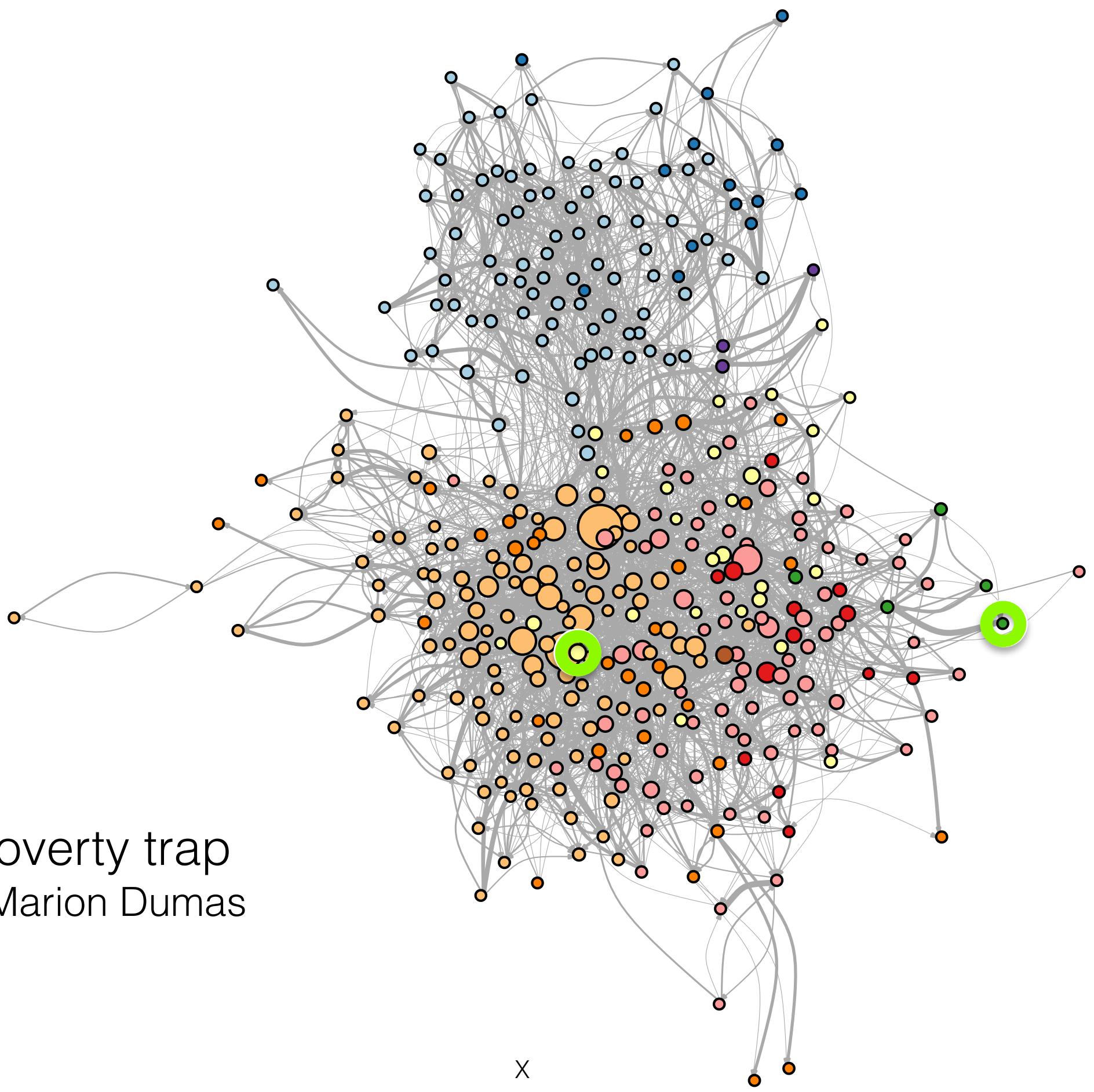
X

Photo: Eleanor A. Power



X

Photo: Eleanor A. Power



The reputational poverty trap  
Eleanor Power & Marion Dumas

X

# We talked about 6 methods:

1. Minimum Violations Rankings & Agony
2. Random Utility Models (economics & marketing)
3. SpringRank (physics)
4. PageRank (random walks & the www)
5. Generative models (social science)
6. Niche models (ecology)

**Beyond pictures: these things matter.**

Inequalities, forecasting, cognition,  
courtship, & social organization.

# And 5 applications:

1. Faculty hiring networks and prestige (computational social science)
2. Sales predictions (e-commerce)
3. Bird hierarchies and cognition (animal behavior)
4. Online dating & desirability (sociology)
5. Group-level social hierarchies (anthropology)

# Many uses for the same techniques. cf regression

## Treat the network like a system:

**Extrapolation.** Make predictions for as-yet unseen nodes (in “space” or time).

**Interpolation.** Identify missing links.

**Generalization.** Nodes of this type are like others of the same type.

## Treat the network like an artifact:

**Mechanisms.** How did this network arise? What rules governed its assembly?

**Explanations.** Coarse-graining or compression.

## Treat the network like a means to an end; an intermediate data structure:

**Useful division.** Need groups so that we can assign treatments in an A/B test.

**Simplification.** Downstream regression model needs ranks or groups.

**PDF** of slides available → <http://LarremoreLab.github.io>

## **Goals** for this talk:

1. **Why** do we look for large-scale structure? 🤔
2. **How** do we find communities and hierarchies? 😊
3. **Where** can we read more details? 📚

# Thank you

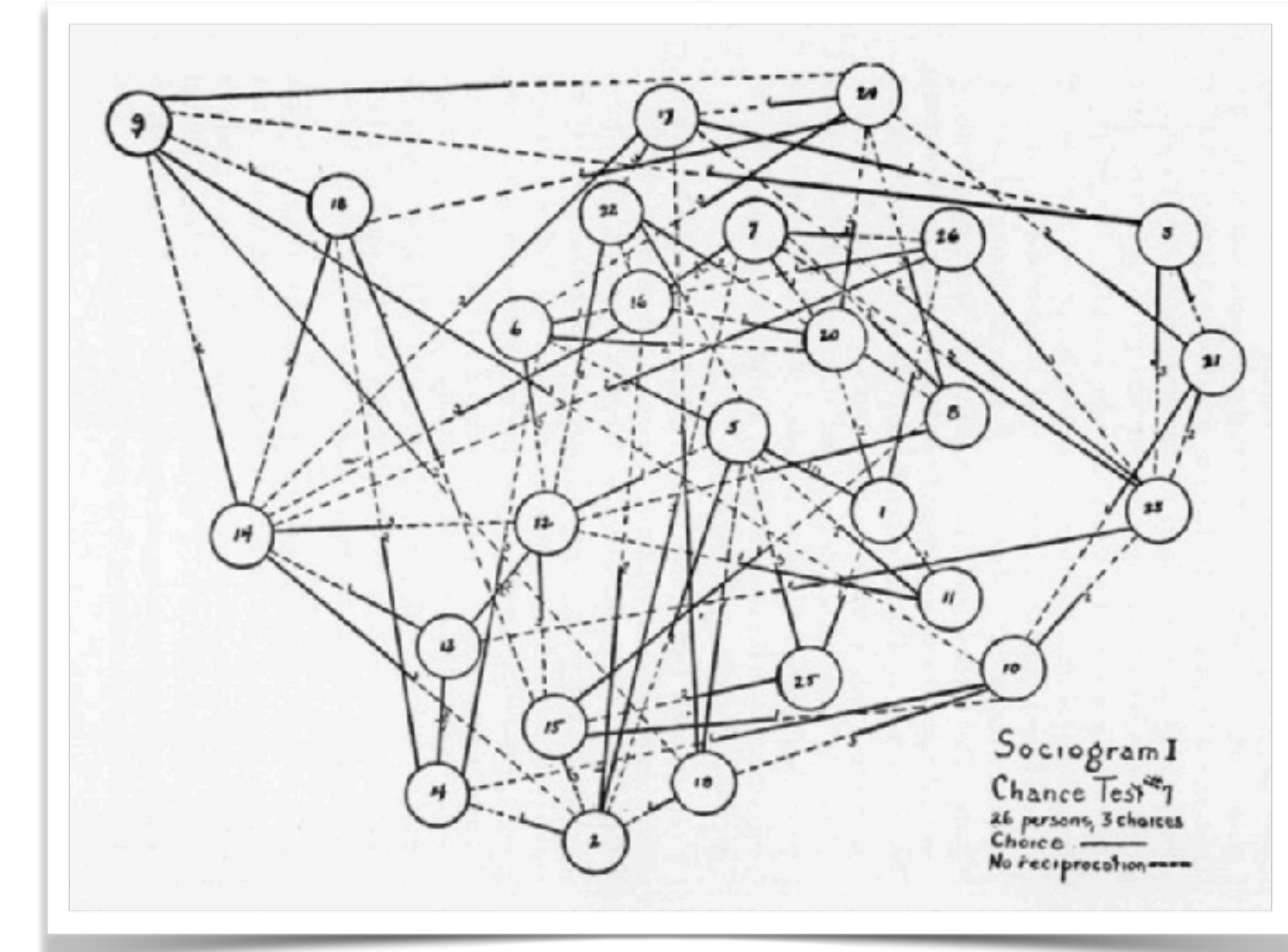
@danlarremore  
daniel.larremore@colorado.edu

# Aside: the birth of null models & chance sociograms

## Who shall survive? Moreno, 1936

Moreno wondered if there were structural explanations for why certain young girls were running away from the school, and thought that sociographic analysis might hold an answer.

chance sociograms



SIAM Review: Configuring random graph models with fixed degree sequences. <http://arxiv.org/abs/1608.00607>

The Book: <http://www.asgpp.org/docs/wss/Book%20VI/index.html>

Johan Ugander's Post: <https://jugander.wordpress.com/2014/08/07/computational-perspectives-on-large-scale-social-networks-a-brief-history/> 93

Aside:

Here is one of my favorite papers of all time:

JMLR: Workshop and Conference Proceedings 27:65–79, 2012      Workshop on Unsupervised and Transfer Learning

## Clustering: Science or Art?

**Ulrike von Luxburg**

*Max Planck Institute for Intelligent Systems, Tübingen, Germany*

ULRIKE.LUXBURG@TUEBINGEN.MPG.DE

**Robert C. Williamson**

*Australian National University and NICTA, Canberra ACT 0200, Australia*

BOB.WILLIAMSON@ANU.EDU.AU

**Isabelle Guyon**

*ClopiNet, 955 Creston Road, Berkeley, CA 94708, USA*

ISABELLE@CLOPINET.COM

<http://proceedings.mlr.press/v27/luxburg12a/luxburg12a.pdf>