

# The Trouble with Community Detection

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



## The Colorado Index of Complex Networks (ICON)

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.

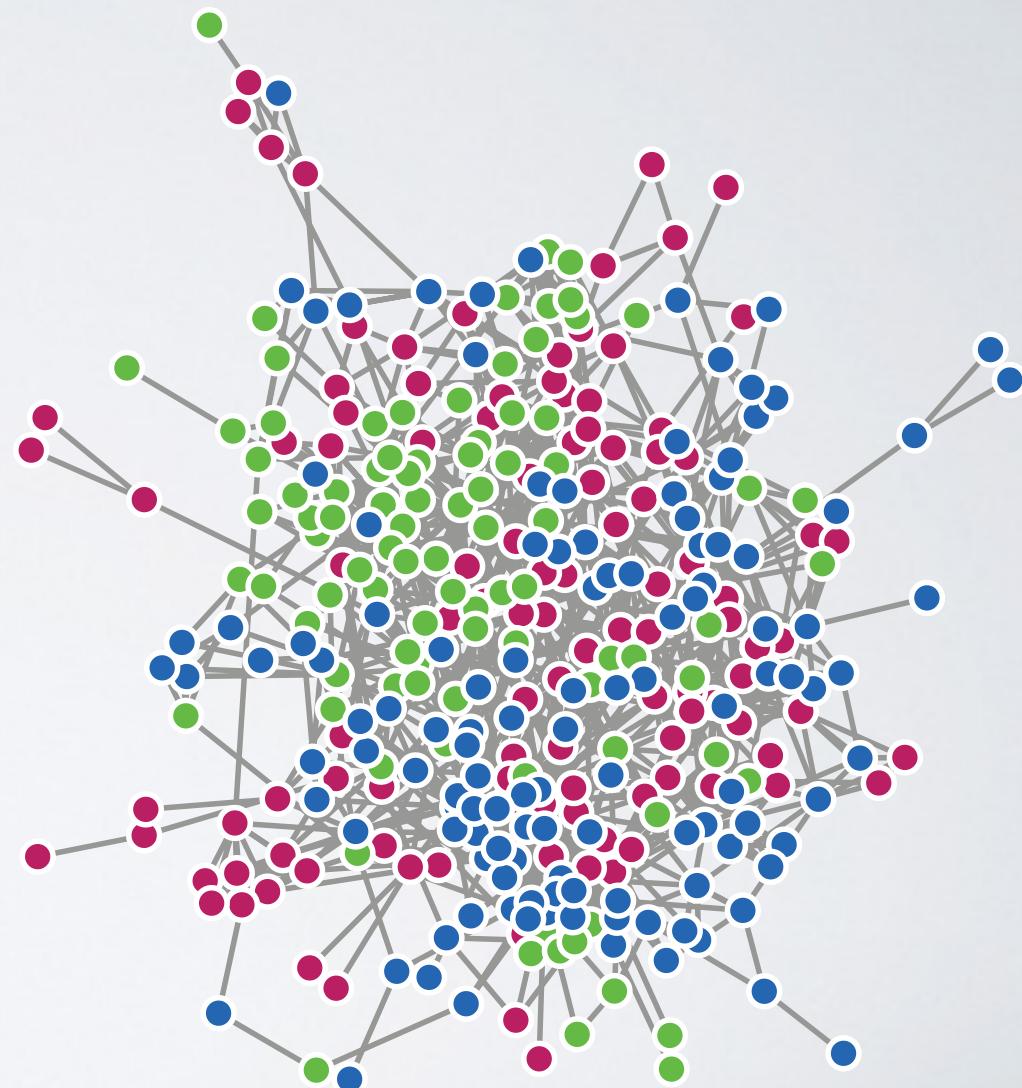


Entries found: 459 Networks found: 3528



# what is community structure?

real networks are complicated objects



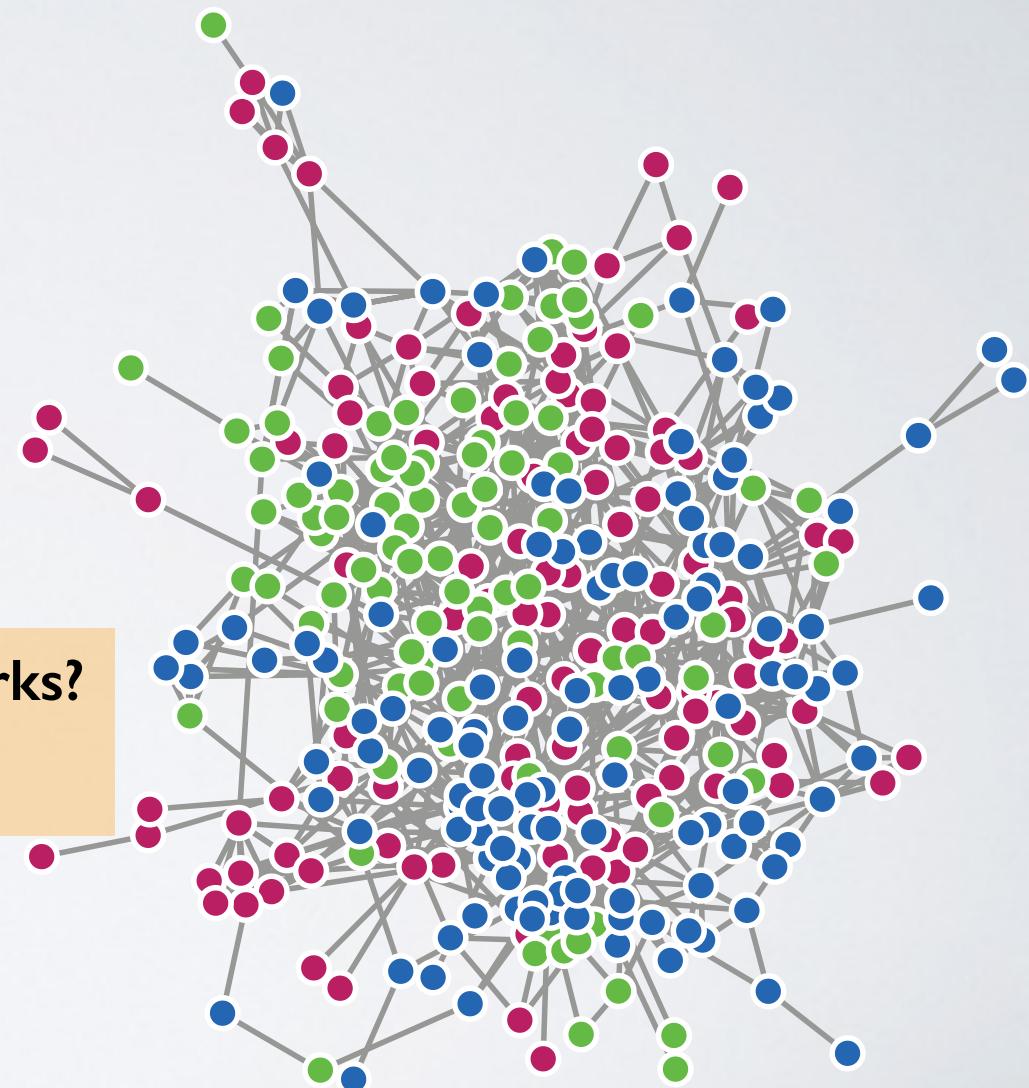
# what is community structure?

real networks are complicated objects

- **how are the edges organized?**
- **how do vertices differ?**
- **does network location matter?**
- **are there underlying patterns?**

what we want to know

- **what processes shape these networks?**
- **how can we tell?**



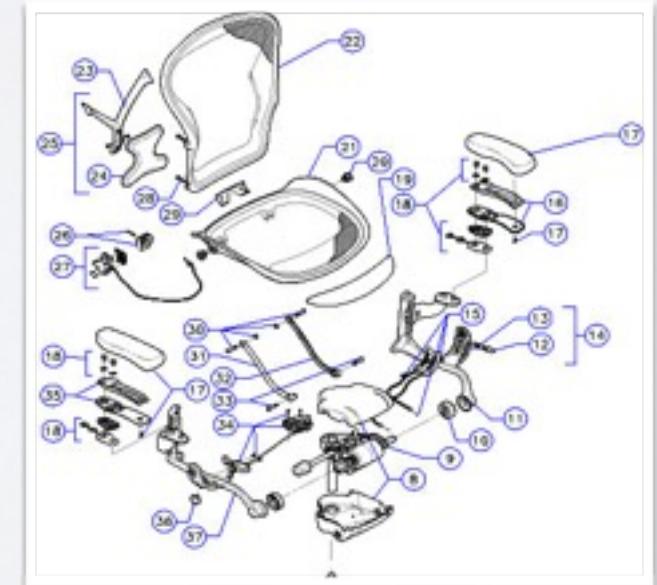
# what is community structure?

what we want : **coarse-grain** the network structure

$$f : G \rightarrow \{\theta_1, \dots, \theta_k\}$$

- what are its building blocks?
  - how are these organized?
  - where are the degrees of freedom  $\vec{\theta}$ ?
  - structure — dynamics — function?

these questions are about ***large-scale structure***



# **what is community structure?**

***large-scale structure = community structure***

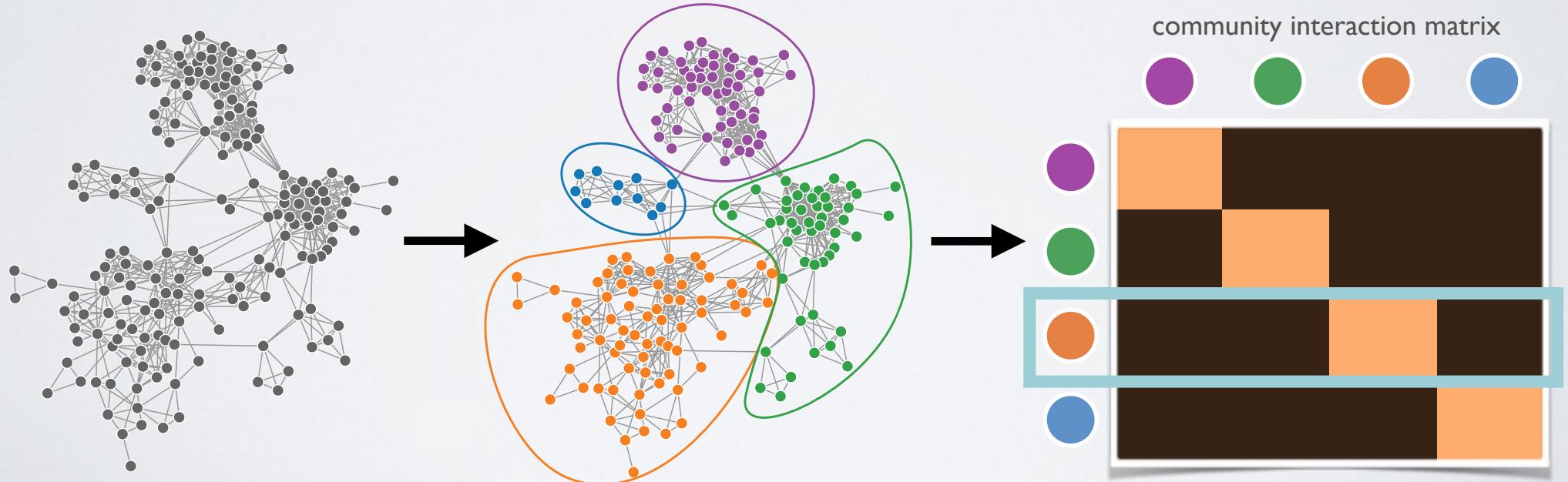
vertices with same pattern of inter-community connections

# what is community structure?

***large-scale structure = community structure***

vertices with same pattern of inter-community connections

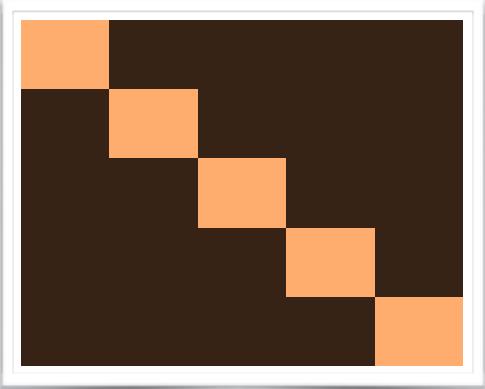
community detection:



# what is community structure?

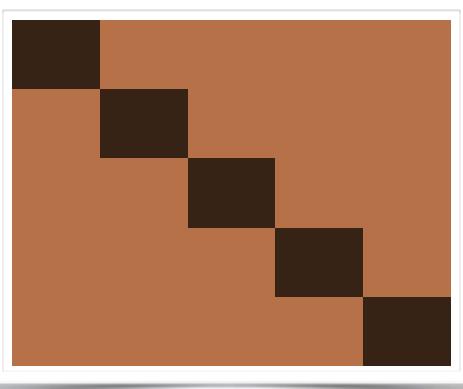
**assortative**

edges within groups



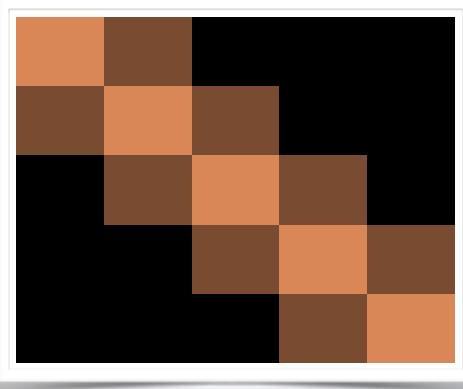
**disassortative**

edges between groups



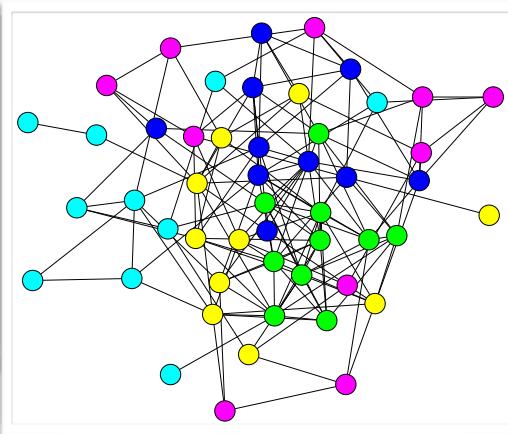
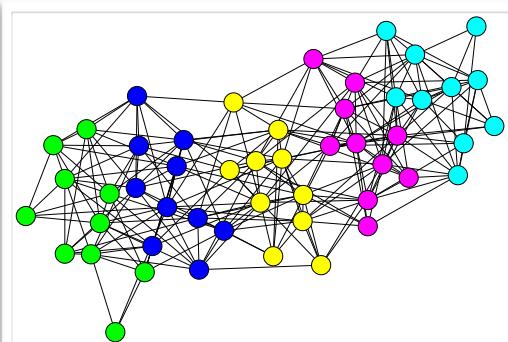
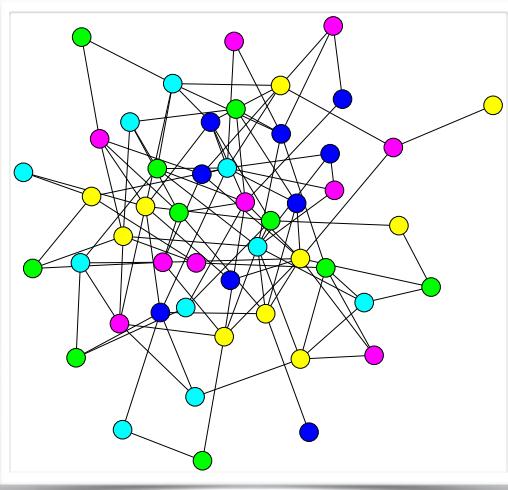
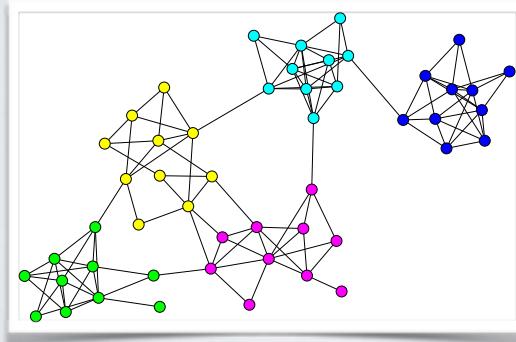
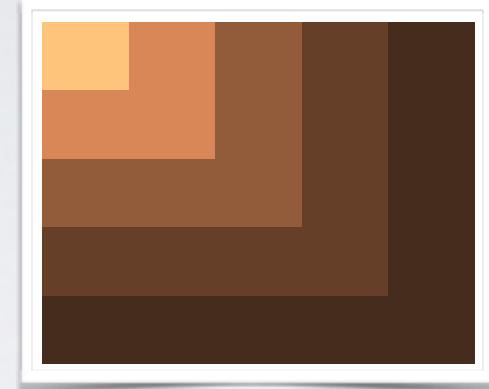
**ordered**

linear group hierarchy



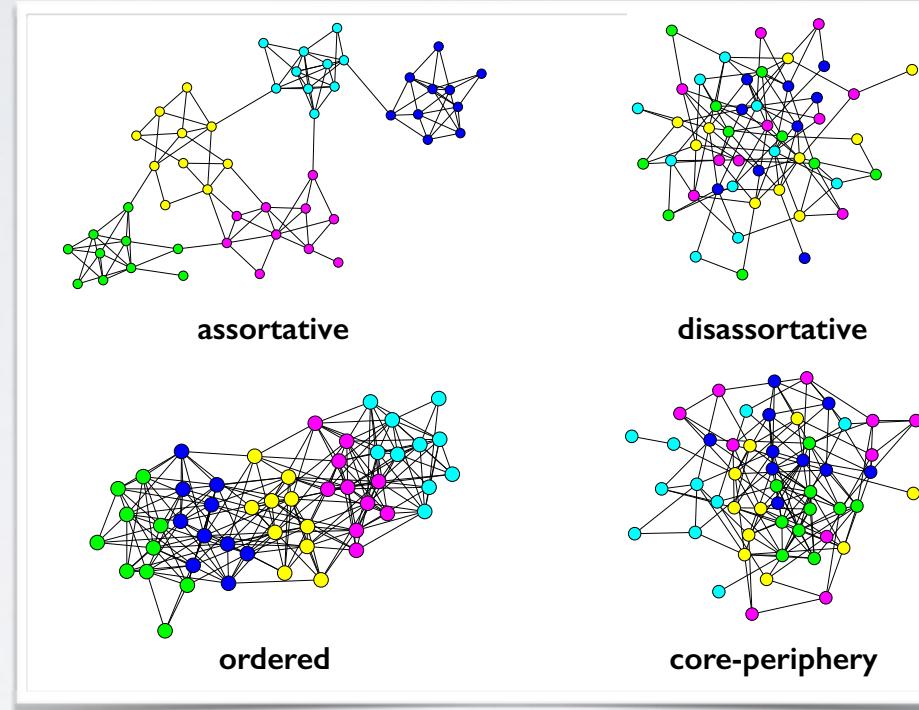
**core-periphery**

dense core, sparse periphery



# what is community structure?

- enormous interest, especially since 2000
- dozens of algorithms for extracting various large-scale patterns
- hundreds of papers published
- spanning Physics, Computer Science, Statistics, Biology, Sociology, and more
- this was one of the first:



## Community structure in social and biological networks

M. Girvan<sup>\*†‡</sup> and M. E. J. Newman<sup>\*§</sup>

PNAS 2002

9500+ citations on Google Scholar

# **the trouble with community detection (part I)**

# the trouble with community detection (part I)

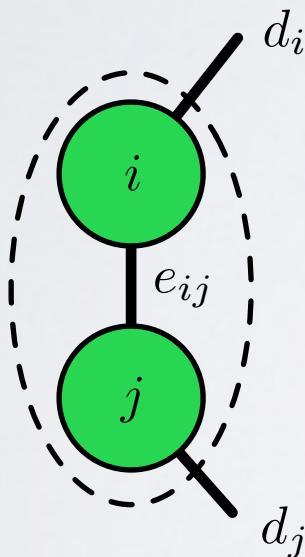
best understood for modularity maximization 
$$Q = \sum_{i=1}^k \left[ \frac{e_i}{m} - \left( \frac{d_i}{2m} \right)^2 \right]$$

- I. ***resolution limit***: in large, unweighted networks, a pair of "true" groups  $i, j$  will be incorrectly merged when maximizing  $Q$

# the trouble with community detection (part I)

best understood for modularity maximization  $Q = \sum_{i=1}^k \left[ \frac{e_i}{m} - \left( \frac{d_i}{2m} \right)^2 \right]$

- I. **resolution limit**: in large, unweighted networks, a pair of "true" groups  $i, j$  will be incorrectly merged when maximizing  $Q$



$$\Delta Q_{ij} = \frac{e_{ij}}{m} - \frac{d_i d_j}{2m^2}$$

merging is favored when

$$e_{ij} > (d_i d_j) / 2m = E[e_{ij}]$$

# the trouble with community detection (part I)

best understood for modularity maximization 
$$Q = \sum_{i=1}^k \left[ \frac{e_i}{m} - \left( \frac{d_i}{2m} \right)^2 \right]$$

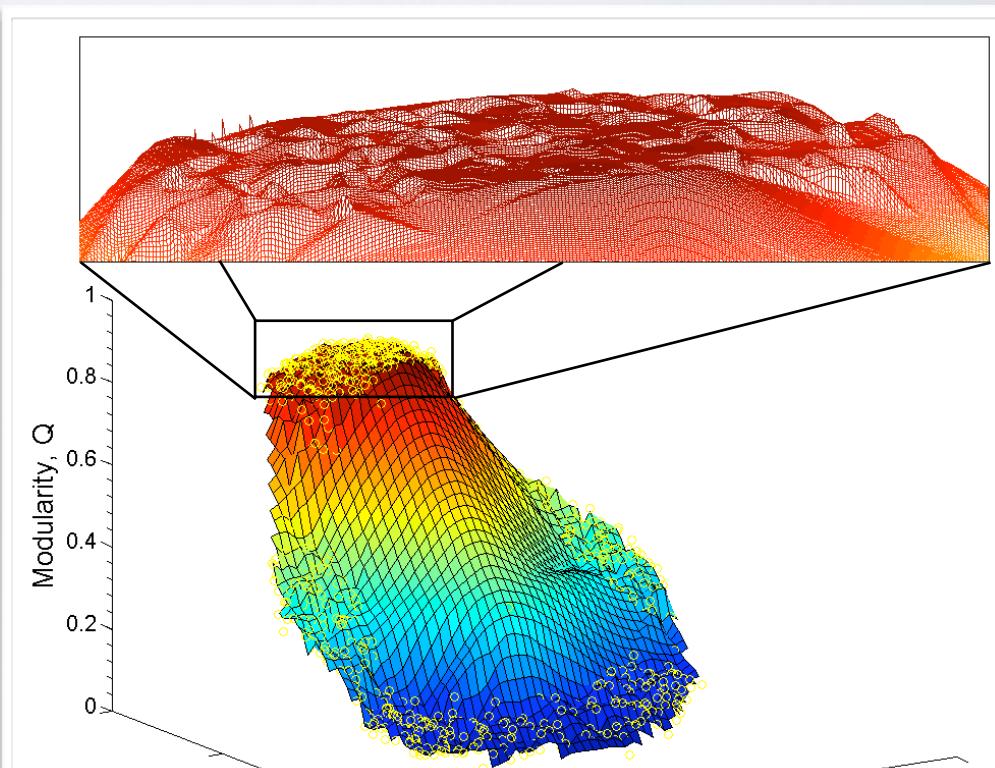
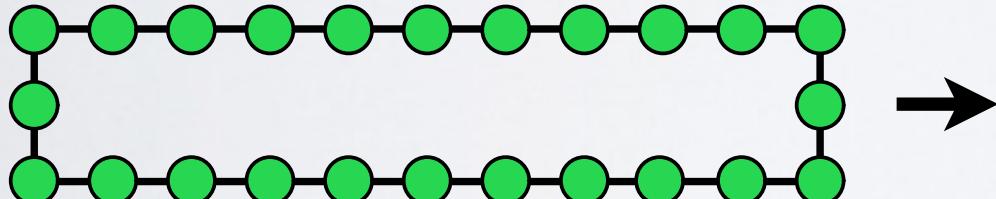
1. ***resolution limit***: in large, unweighted networks, a pair of "true" groups  $i, j$  will be incorrectly merged when maximizing  $Q$
2. ***extreme degeneracies***: there are an exponential number of local optima with scores close to the maximum  $Q$

# the trouble with community detection (part I)

best understood for modularity maximization

$$Q = \sum_{i=1}^k \left[ \frac{e_i}{m} - \left( \frac{d_i}{2m} \right)^2 \right]$$

1. ***resolution limit***: in large, unweighted networks, a pair of "true" groups  $i, j$  will be incorrectly merged when maximizing  $Q$
2. ***extreme degeneracies***: there are an exponential number of local optima with scores close to the maximum  $Q$



# the trouble with community detection (part I)

best understood for modularity maximization 
$$Q = \sum_{i=1}^k \left[ \frac{e_i}{m} - \left( \frac{d_i}{2m} \right)^2 \right]$$

1. ***resolution limit***: in large, unweighted networks, a pair of "true" groups  $i, j$  will be incorrectly merged when maximizing  $Q$
2. ***extreme degeneracies***: there are an exponential number of local optima with scores close to the maximum  $Q$

***the results of community detection should be interpreted with caution:***

the true groups are likely merged in complicated ways

# **the trouble with community detection (part 2)**

# the trouble with community detection (part 2)

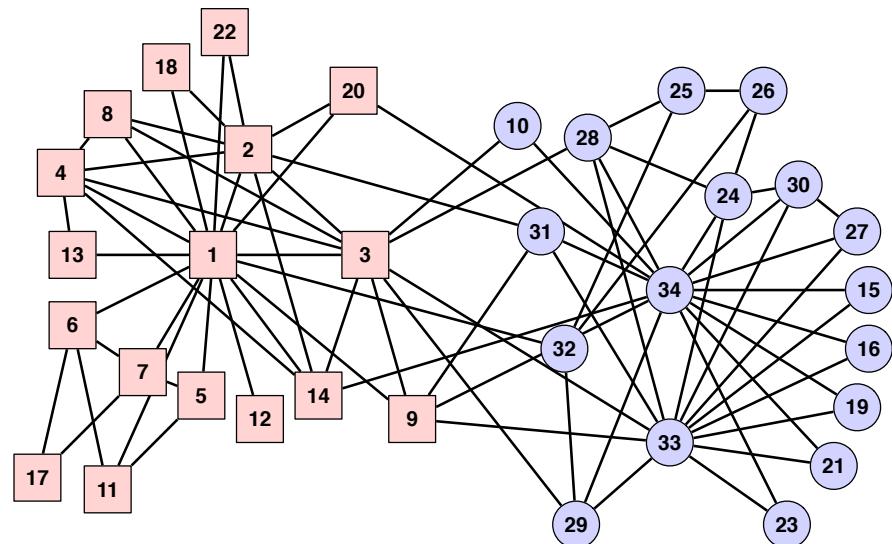
many networks include **metadata** on their nodes:

social networks	age, sex, ethnicity or race, etc.
food webs	feeding mode, species body mass, etc.
Internet	data capacity, physical location, etc.
protein interactions	molecular weight, association with cancer, etc.

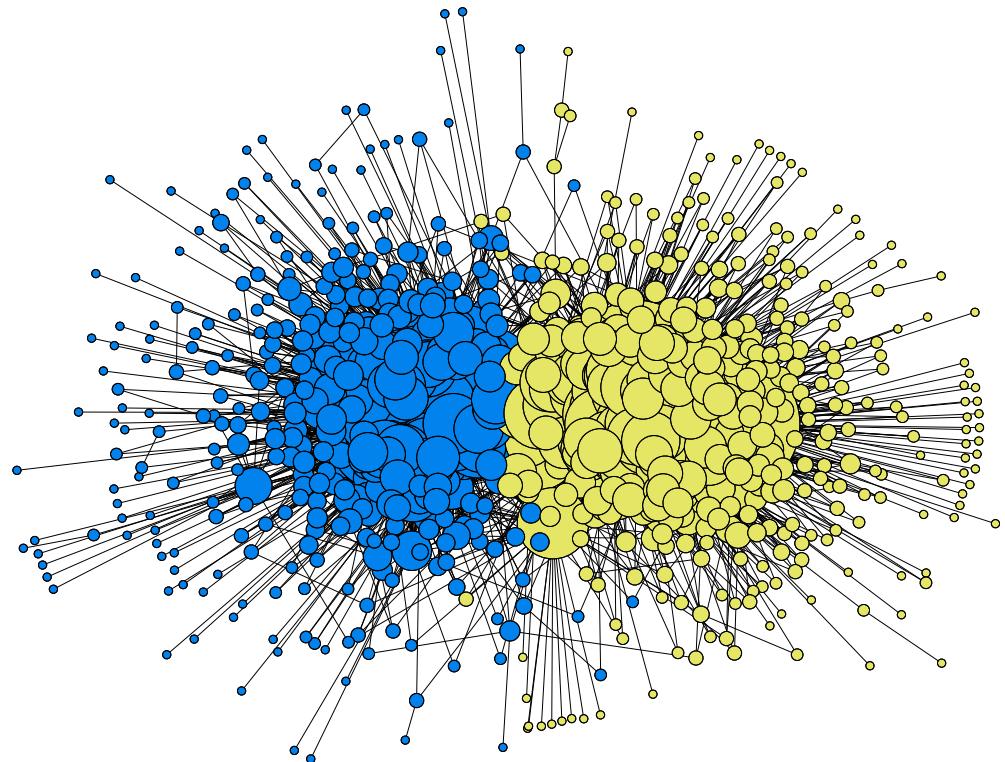
metadata  $\mathbf{x}$  is often used to evaluate the accuracy of community detection algs.

if community detection method  $\mathcal{A}$  finds a partition  $\mathcal{P}$  that correlates with  $\mathbf{x}$   
then we say that  $\mathcal{A}$  is good

# the trouble with community detection

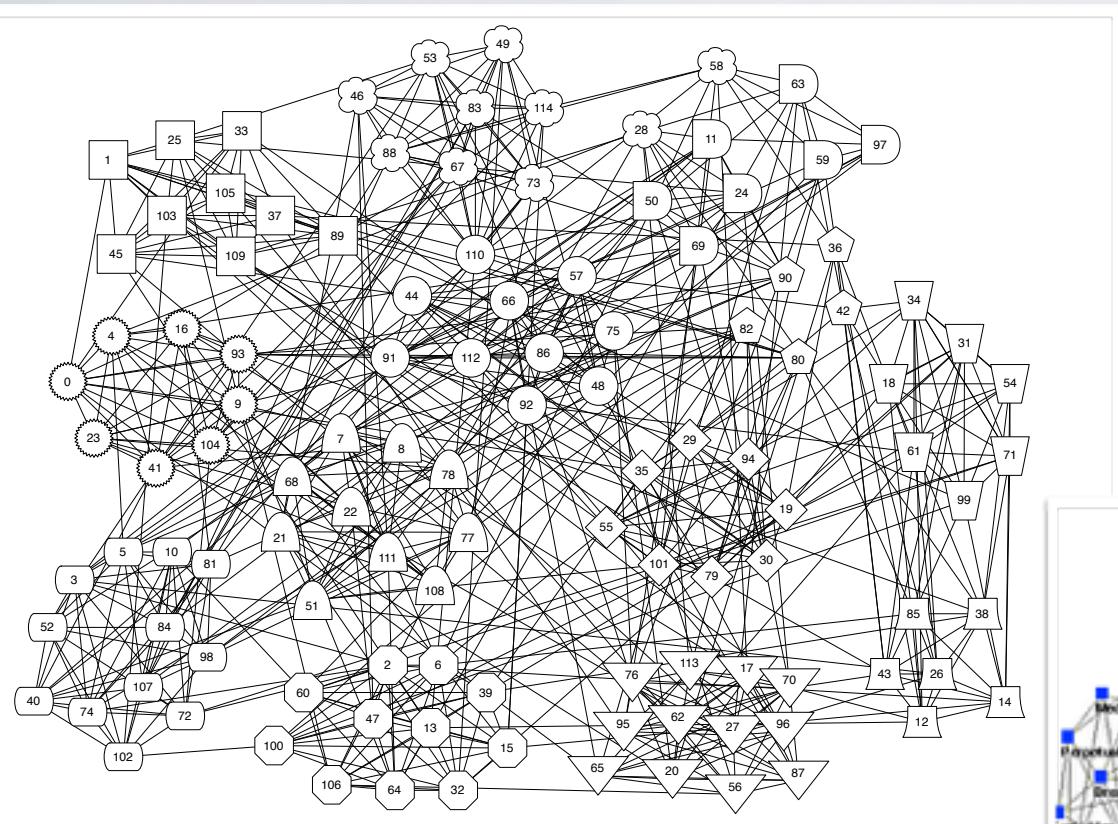


Zachary karate club

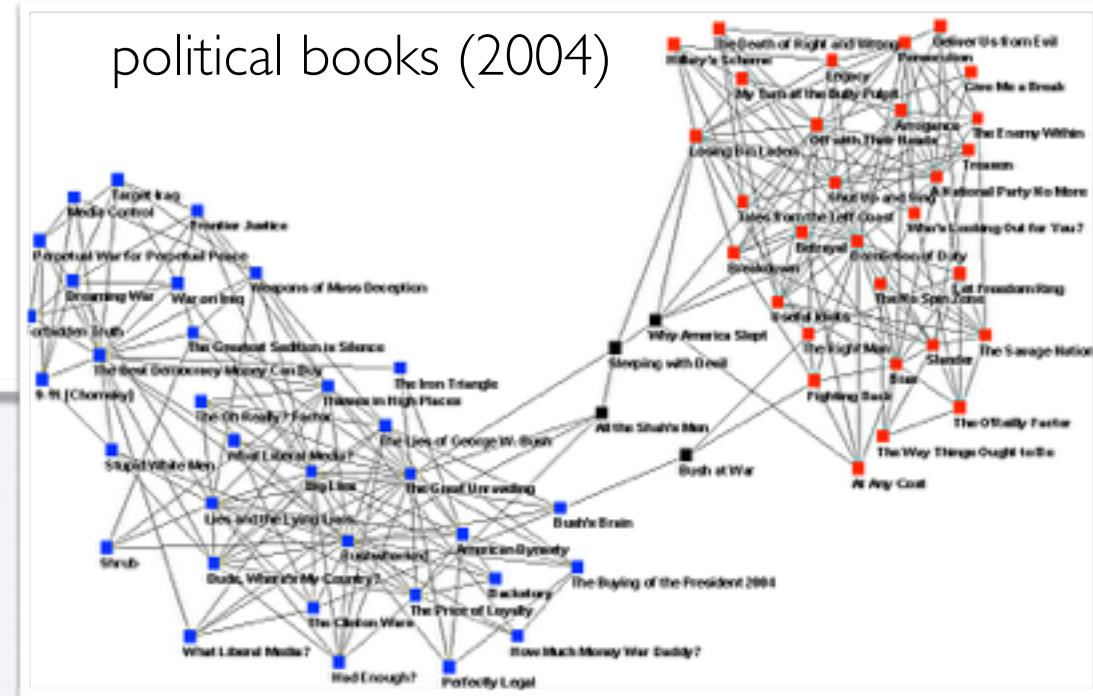


political blogs network

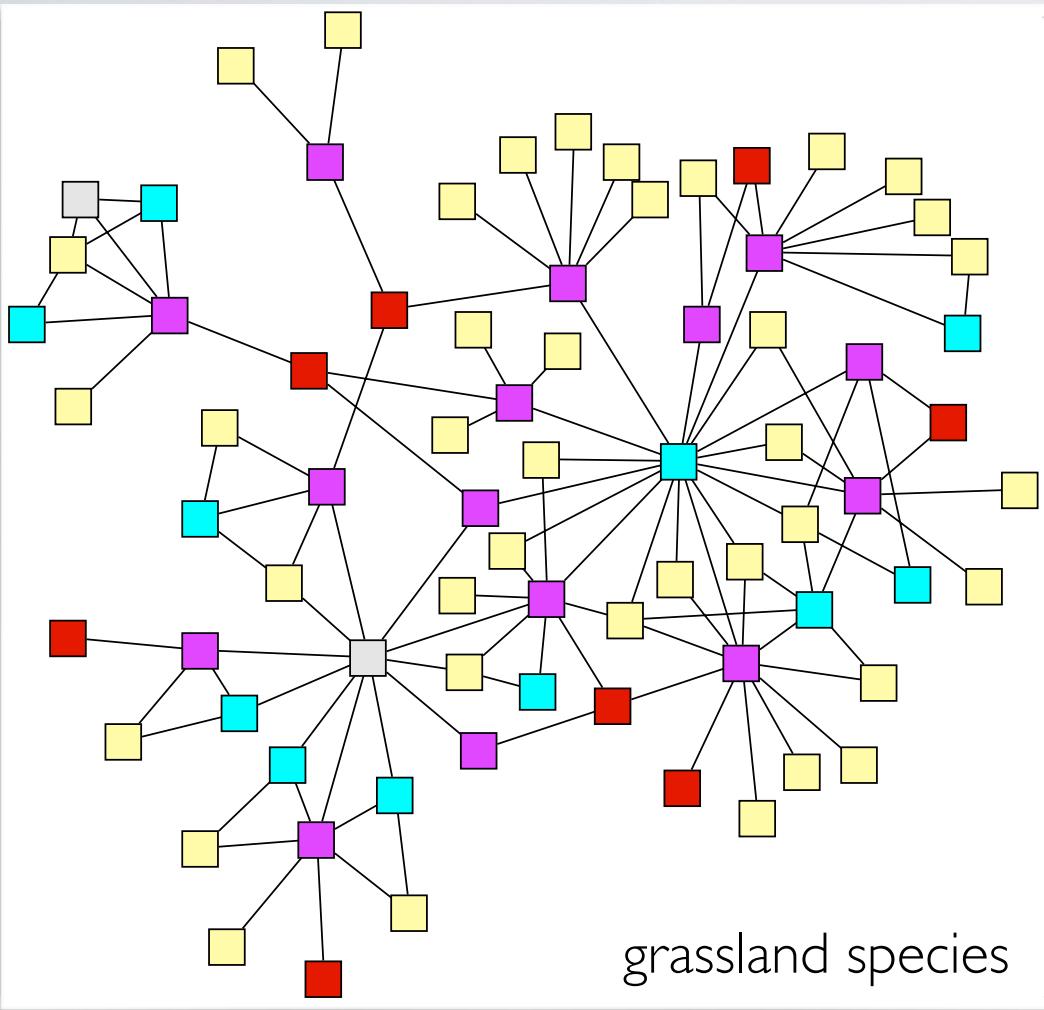
# the trouble with community detection



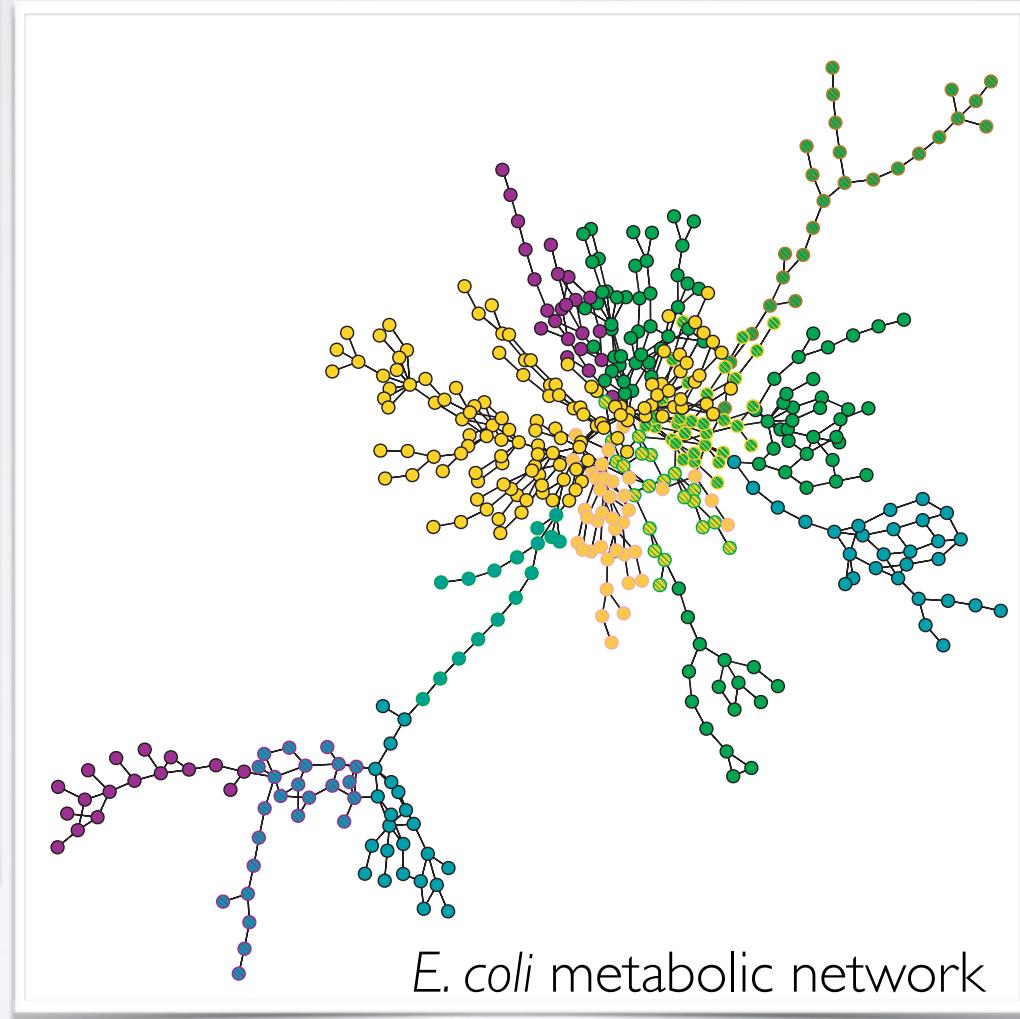
# NCAA 2000 Schedule



# the trouble with community detection



grassland species



*E. coli* metabolic network

# **the trouble with community detection**

**often, groups found by community detection are meaningful**

- allegiances or personal interests in social networks [1]
- biological function in metabolic networks [2]

**but**

[1] see Fortunato (2010), and Adamic & Glance (2005)

[2] see Holme, Huss & Jeong (2003), and Guimera & Amaral (2005)

# the trouble with community detection

often, groups found by community detection are meaningful

- allegiances or personal interests in social networks [1]
- biological function in metabolic networks [2]

but some recent studies claim these are the exception

- real networks **either** do not contain structural communities **or** communities exist but they do not correlate with metadata groups [3]

[1] see Fortunato (2010), and Adamic & Glance (2005)

[2] see Holme, Huss & Jeong (2003), and Guimera & Amaral (2005)

[3] see Leskovec et al. (2009), and Yang & Leskovec (2012), and Hric, Darst & Fortunato (2014)

# the trouble with community detection

Hric, Darst & Fortunato (2014)

- 115 networks with metadata & 12 community detection methods
- compare extracted  $\mathcal{P}$  with observed  $\mathbf{x}$  for each  $\mathcal{A}$

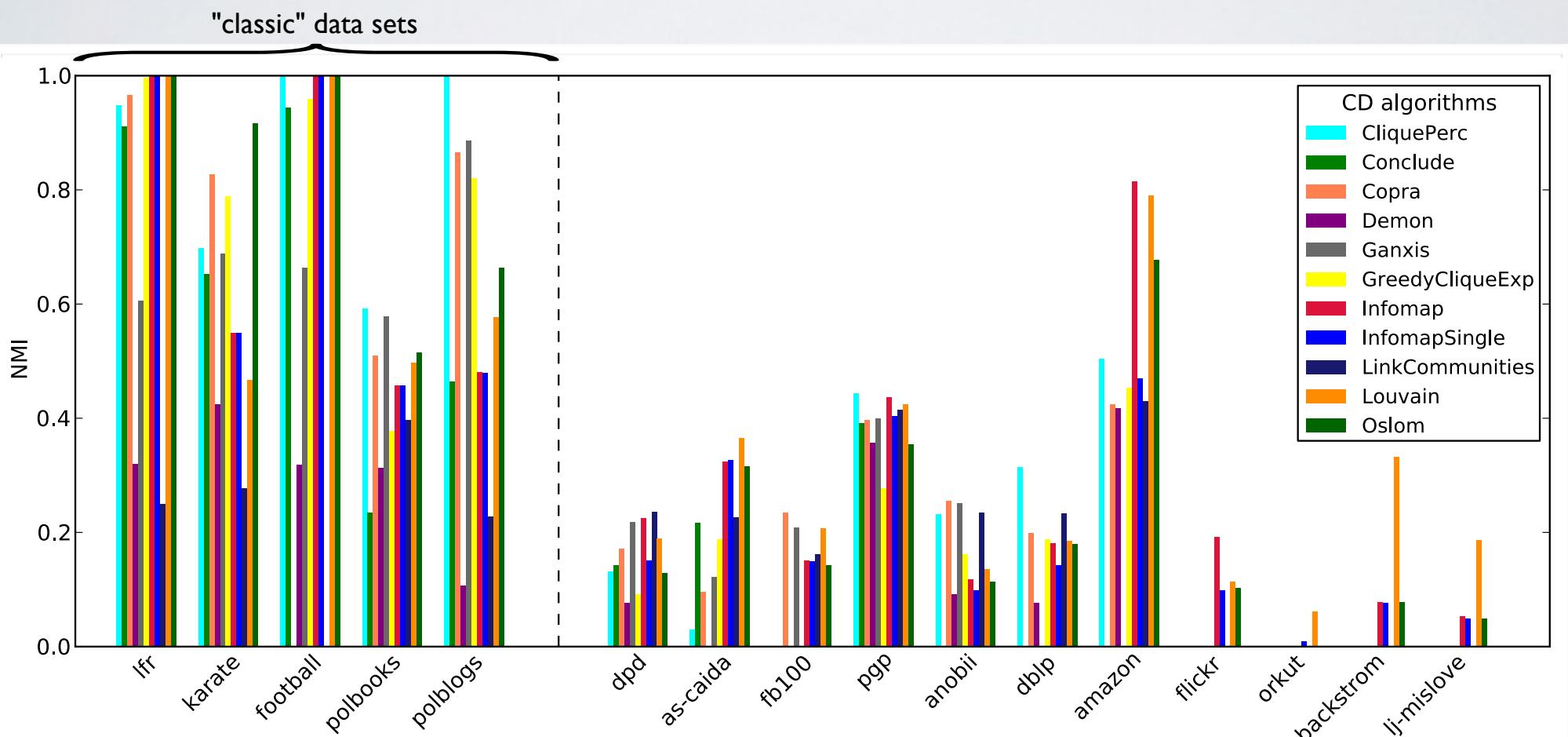
Name	No. Nodes	No. Edges	No. Groups	Description of group nature
lfr	1000	9839	40	artificial network (lfr, 1000S, $\mu = 0.5$ )
karate	34	78	2	membership after the split
football	115	615	12	team scheduling groups
polbooks	105	441	2	political alignment
polblogs	1222	16782	3	political alignment
dpd	35029	161313	580	software package categories
as-caida	46676	262953	225	countries
fb100	762–41536	16651–1465654	2–2597	common students' traits
pgp	81036	190143	17824	email domains
anobii	136547	892377	25992	declared group membership
dblp	317080	1049866	13472	publication venues
amazon	366997	1231439	14–29432	product categories
flickr	1715255	22613981	101192	declared group membership
orkut	3072441	117185083	8730807	declared group membership
lj-backstrom	4843953	43362750	292222	declared group membership
lj-mislove	5189809	49151786	2183754	declared group membership

[!] fb100 is 100 networks

# the trouble with community detection

Hric, Darst & Fortunato (2014)

- evaluate by normalized mutual information  $NMI(\mathcal{P}, \mathbf{x})$



[!] maximum NMI between any partition layer of the metadata partitions and any layer returned by the community detection method

# the trouble with community detection

lies, damned lies, and community detection

*true groups can be merged in complicated ways*

*best partition typically lost in exponential number of local optima*

*the groups we do find often don't correlate with what we think we  
are trying to recover*

**but wait!**



# a solution

**idea:**

use metadata  $\mathbf{x}$  to help select a partition  $\mathcal{P}^* \in \{\mathcal{P}\}$  that correlates with  $\mathbf{x}$ , from among the exponential number of *plausible* partitions



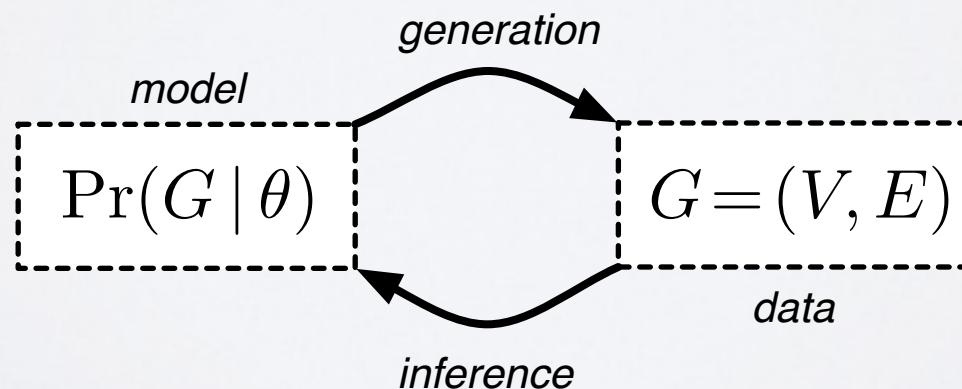
# a solution

**idea:**

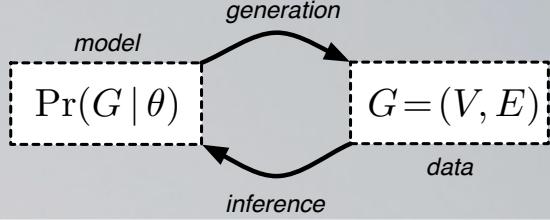
use metadata  $\mathbf{x}$  to help select a partition  $\mathcal{P}^* \in \{\mathcal{P}\}$  that correlates with  $\mathbf{x}$ , from among the exponential number of *plausible* partitions

use a generative model to guide the selection:

- define a parametric probability distribution over networks  $\Pr(G | \theta)$
- *generation* : given  $\theta$  , draw  $G$  from this distribution
- *inference* : given  $G$ , choose  $\theta$  that makes  $G$  likely



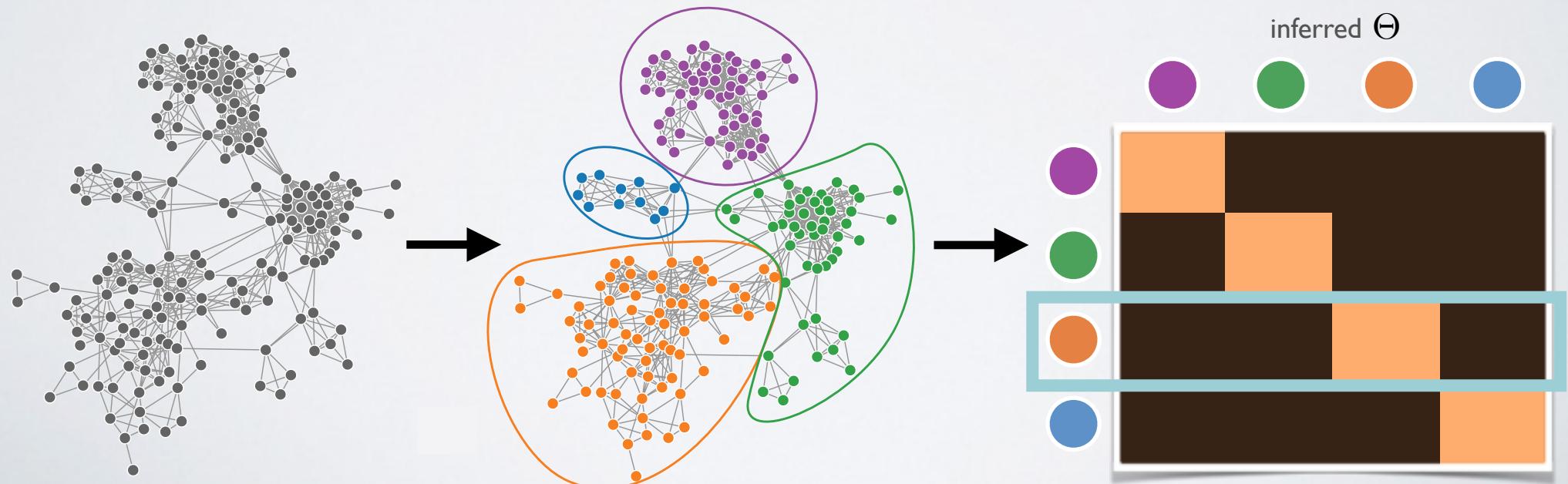
# a solution



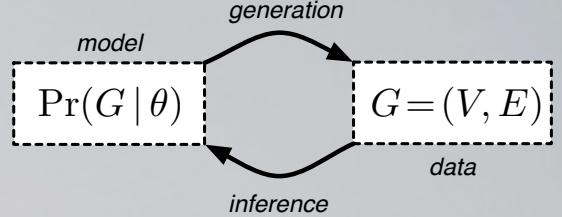
## the stochastic block model

- each vertex  $u$  has type  $s_u \in \{1, \dots, k\}$
- stochastic block matrix  $\Theta$  of group-level connection probabilities
- probability that  $P(A_{uv} = 1) = \theta_{s_u, s_v}$

community = vertices with same pattern of inter-community connections



# a metadata-aware stochastic block model



## generation

given metadata  $\mathbf{x} = \{x_u\}$  and degree  $\mathbf{d} = \{d_u\}$  for each node  $u$

- each node  $u$  is assigned a community  $s$  with probability  $\gamma_{sx}$
- thus, prior on community assignments is  $P(s | \Gamma, \mathbf{x}) = \prod_i \gamma_{s_i, x_i}$
- given assignments, place edges independently, each with probability:

$$p_{uv} = d_u d_v \theta_{s_u, s_v}$$

- where the  $\theta_{st}$  are the stochastic block matrix parameters

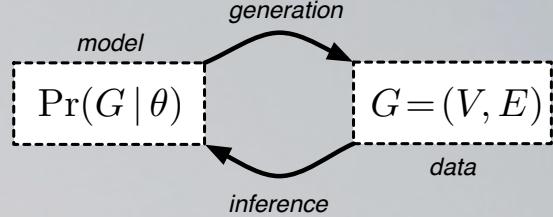
this is a degree-corrected stochastic block model (DC-SBM)

with a metadata-based prior on community labels

[1]  $\Gamma$  is the  $k \times K$  matrix of parameters  $\gamma_{sx}$

[2] Karrer & Newman (2011)

# a metadata-aware stochastic block model



## inference

given observed network  $\mathbf{A}$  (adjacency matrix)

- the model likelihood is

$$\begin{aligned} P(\mathbf{A} | \Theta, \Gamma, \mathbf{x}) &= \sum_{\mathbf{s}} P(\mathbf{A} | \Theta, \mathbf{s}) P(\mathbf{s} | \Gamma, \mathbf{x}) \\ &= \sum_{\mathbf{s}} \prod_{u < v} p_{uv}^{A_{uv}} (1 - p_{uv})^{1 - A_{uv}} \prod_u \gamma_{s_u, x_u} \end{aligned}$$

- where  $\Theta$  is a  $k \times k$  matrix of community interaction parameters  $\theta_{st}$ , and the sum is over all possible assignments  $\mathbf{s}$
- we fit this model to data using expectation-maximization (EM) to maximize  $P(\mathbf{A} | \Theta, \Gamma, \mathbf{x})$  w.r.t.  $\Theta$  and  $\Gamma$

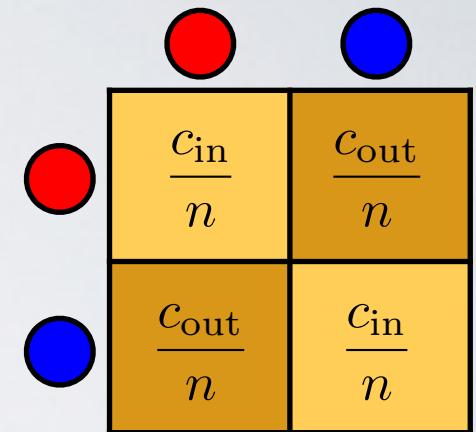
# networks with planted structure

*does this method recover known structure in synthetic data?*

# networks with planted structure

**does this method recover known structure in synthetic data?**

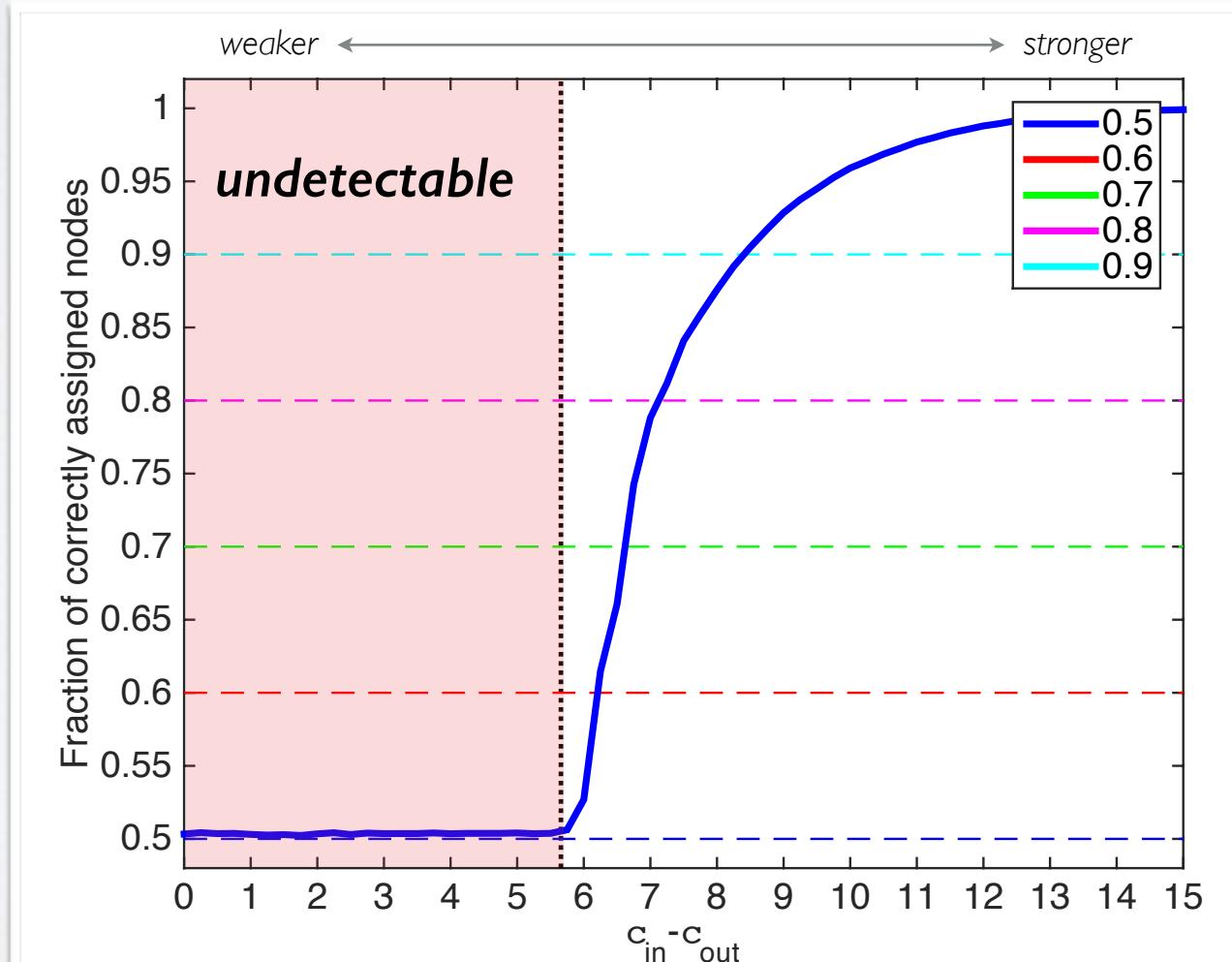
- use SBM to generate *planted partition* networks, with  $k = 2$  equal-sized groups and mean degree  $c = (c_{\text{in}} + c_{\text{out}})/2$
- assign metadata with variable correlation  $\rho \in [0.5, 0.9]$  to true group labels
- vary strength of partition  $c_{\text{in}} - c_{\text{out}}$
- when  $c_{\text{in}} - c_{\text{out}} \leq \sqrt{2(c_{\text{in}} + c_{\text{out}})}$ , no structure-only algorithm can recover the planted communities better than chance (the *detectability threshold*, which is a phase transition)



# networks with planted structure

let mean degree  $c = 8$

- when  $\rho = 0.5$ , metadata isn't useful and we recover regular SBM behavior



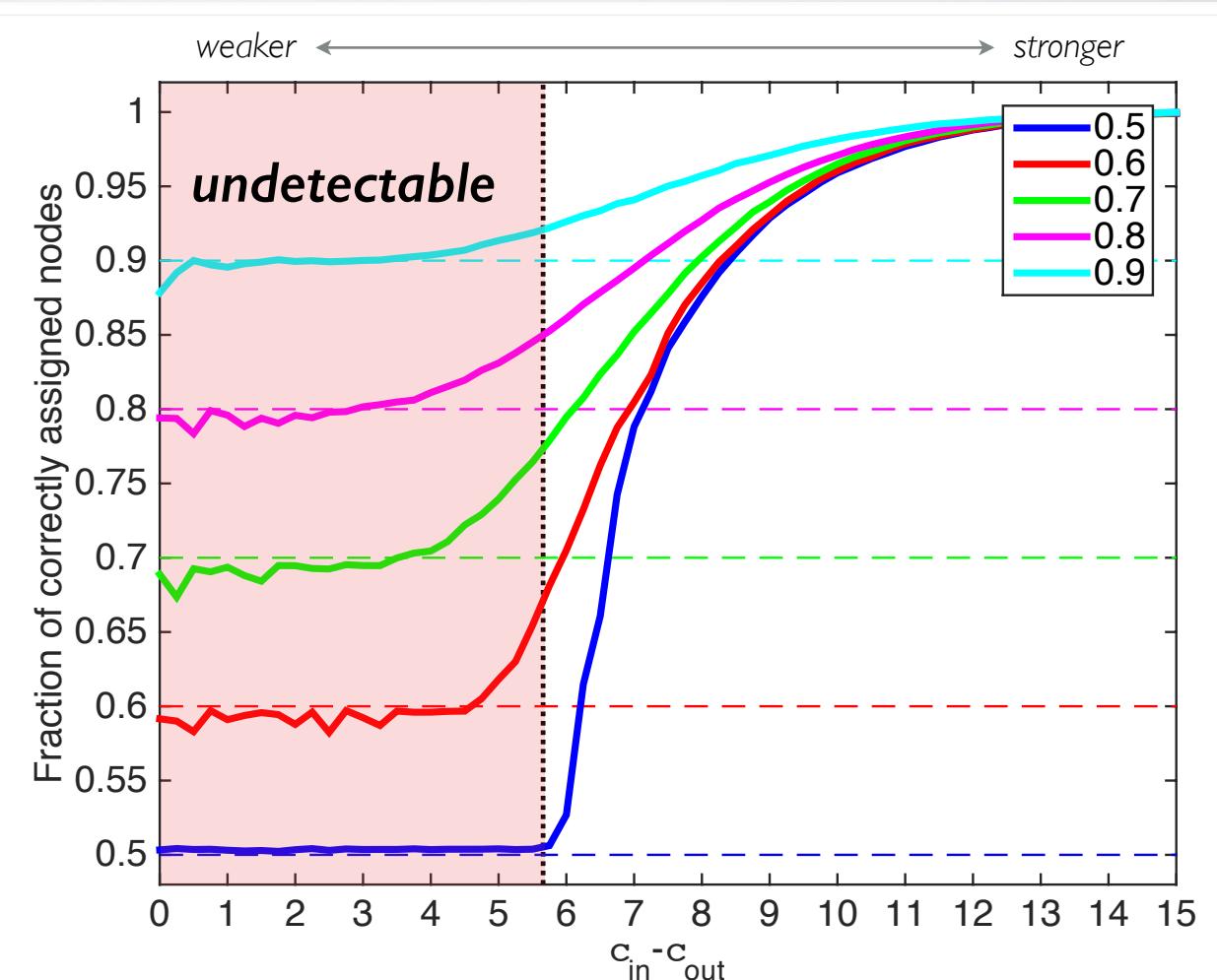
# networks with planted structure

let mean degree  $c = 8$

- when  $\rho = 0.5$ , metadata isn't useful and we recover regular SBM behavior
- when metadata correlates with true groups,  $\rho > 0.5$  accuracy is better than either metadata or SBM alone

**metadata + SBM performs better than either**

- **any algorithm without metadata, or**
- **metadata alone.**



# **real-world networks**

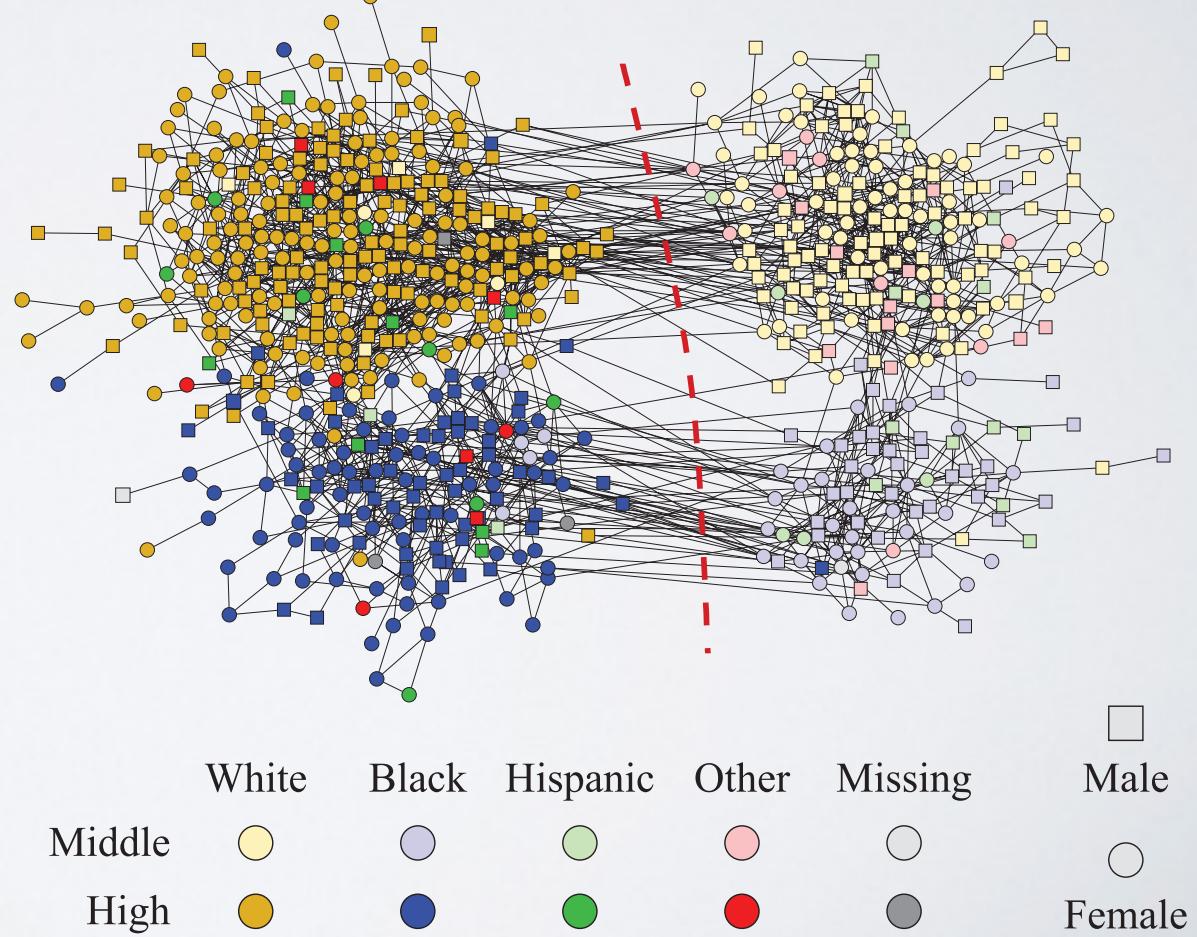
# real-world networks

1. **high school social network:** 795 students in a medium-sized American high school and its feeder middle school
2. **marine food web:** predator-prey interactions among 488 species in Weddell Sea in Antarctica
3. **Malaria gene recombinations:** recombination events among 297 var genes
4. **Facebook friendships:** online friendships among 15,126 Harvard students and alumni
5. **Internet graph:** peering relations among 46,676 Autonomous Systems

# real-world networks

I. **high school social network:** 795 students in a medium-sized American high school and its feeder middle school

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



# real-world networks

I. **high school social network:** 795 students in a medium-sized American high school and its feeder middle school

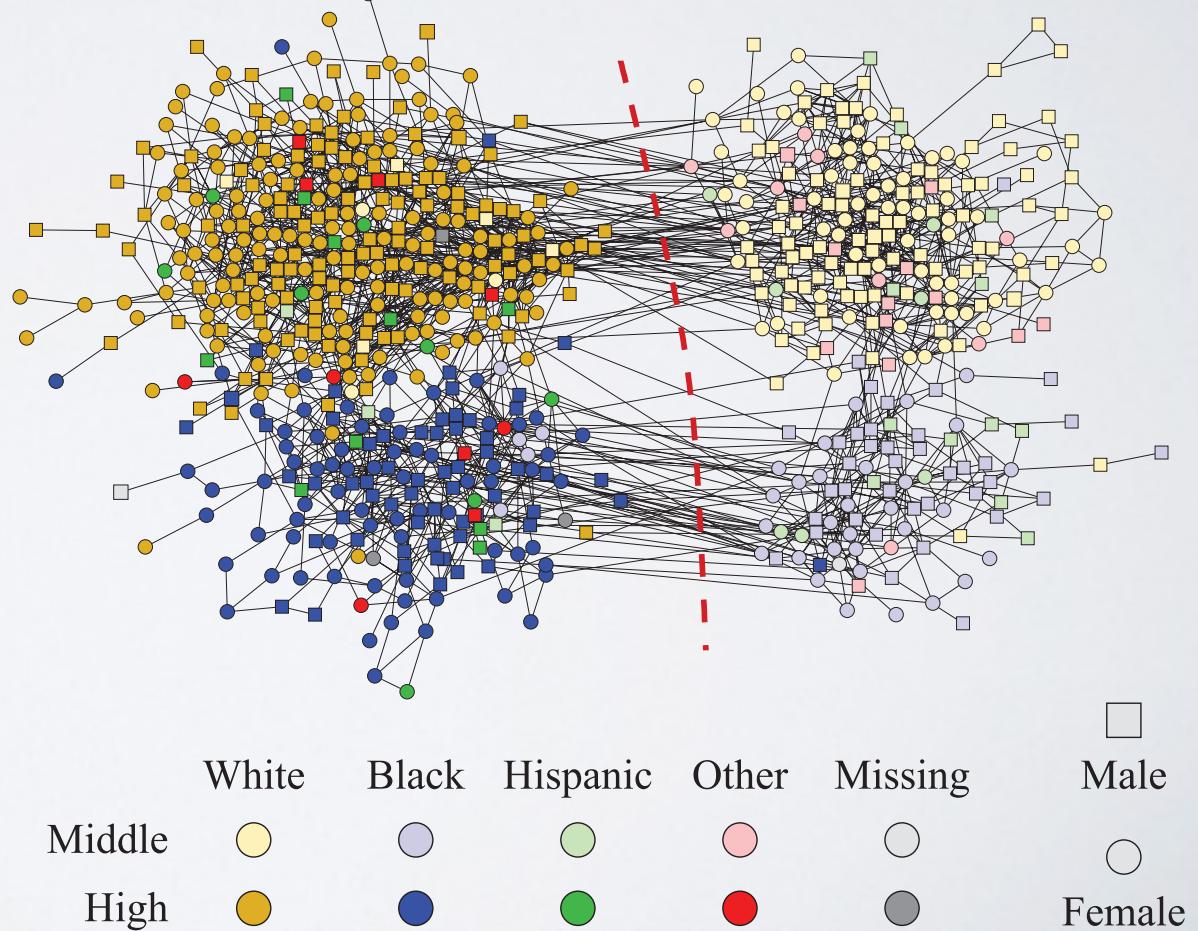
- $x = \{\text{grade 7-12, ethnicity, gender}\}$

- method finds a good partition between high-school and middle-school

$$\text{NMI} = 0.881$$

- without metadata:

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



# real-world networks

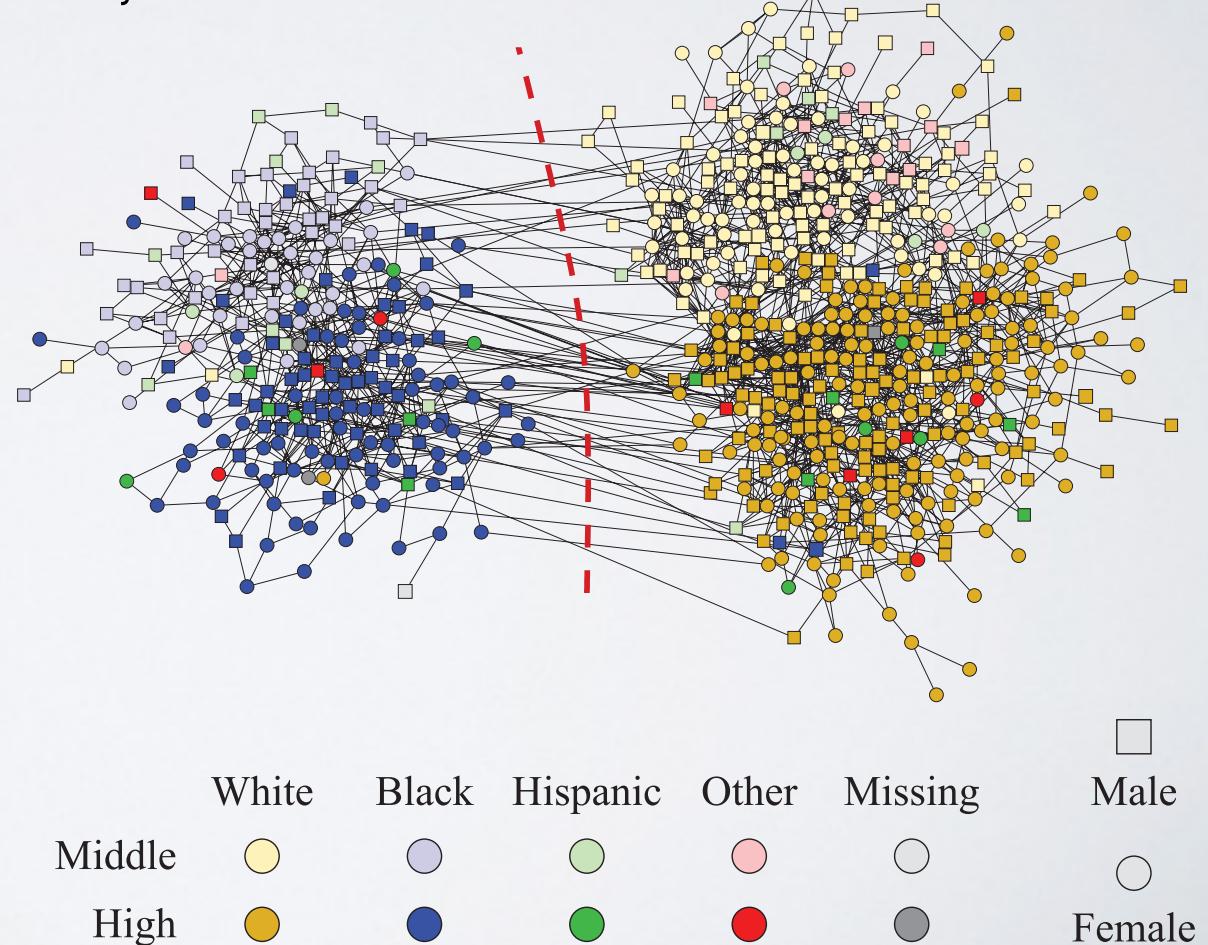
I. **high school social network:** 795 students in a medium-sized American high school and its feeder middle school

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

$\text{NMI} = 0.820$

- without metadata:

$\text{NMI} \in [0.120, 0.239]$



# real-world networks

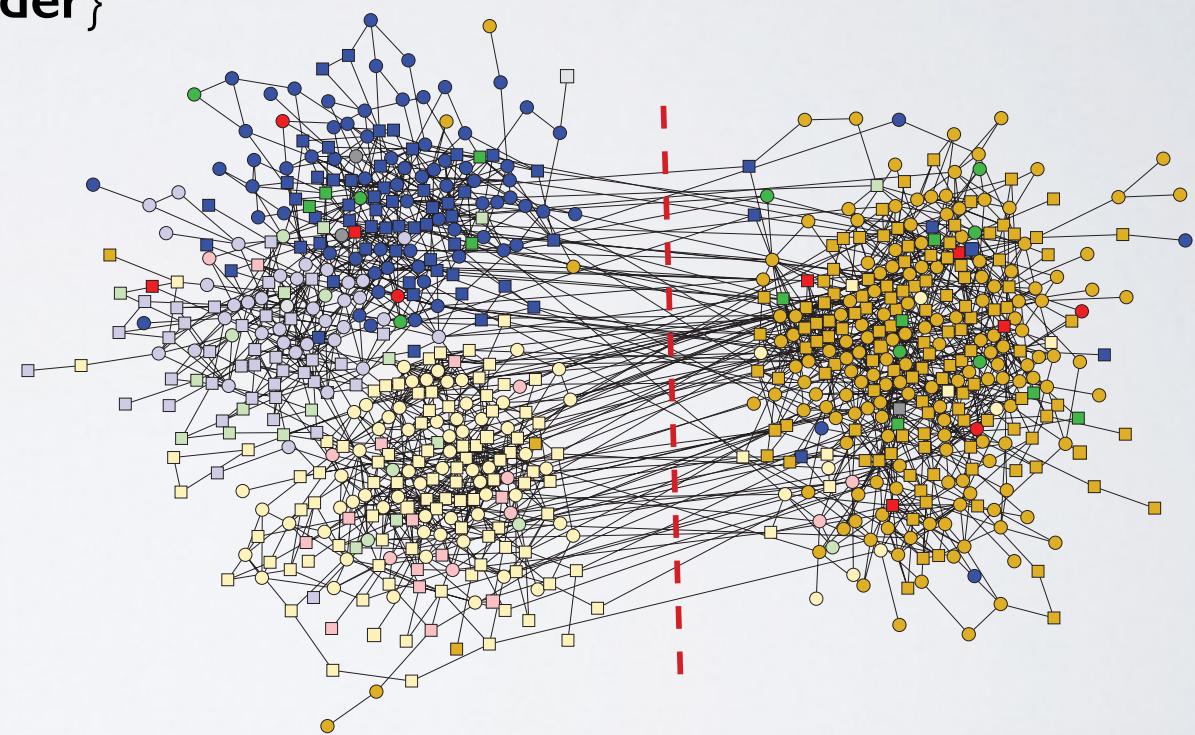
I. **high school social network:** 795 students in a medium-sized American high school and its feeder middle school

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

$$\text{NMI} = 0.003$$

- without metadata:

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

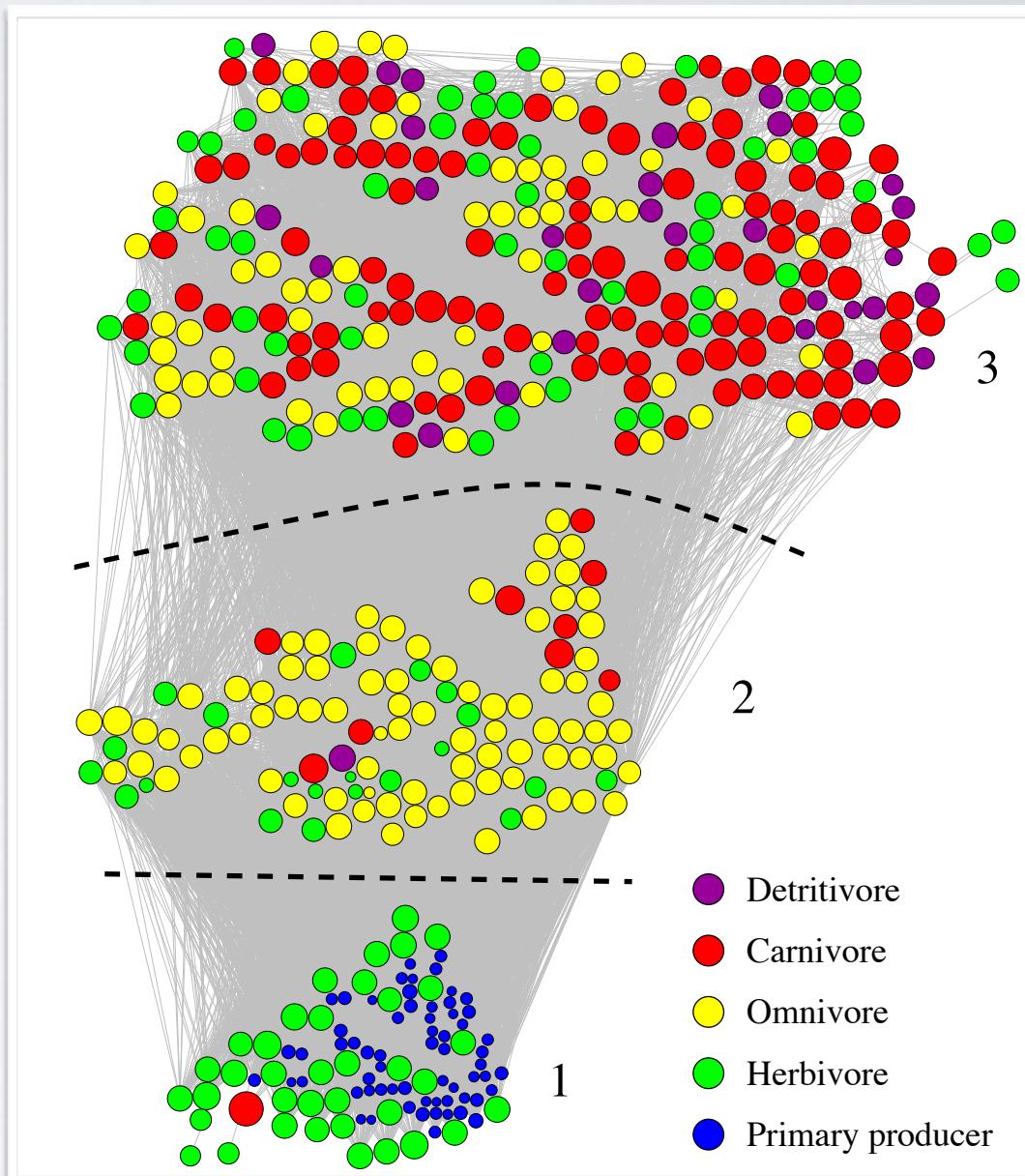
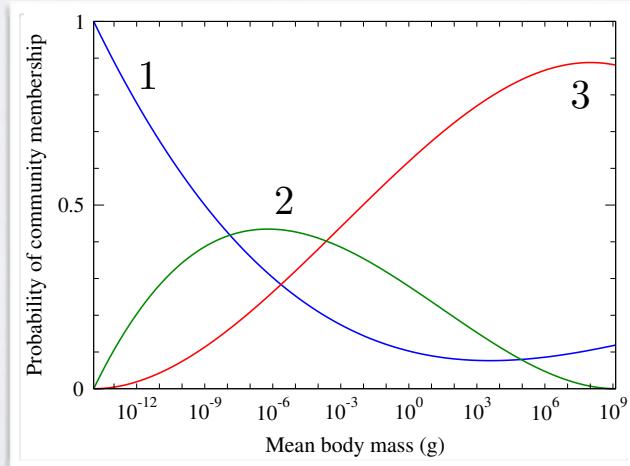


	White	Black	Hispanic	Other	Missing	Male
Middle	○	○	○	○	○	○
High	●	●	●	●	●	●

# 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:



[1] here, we're using a continuous metadata model

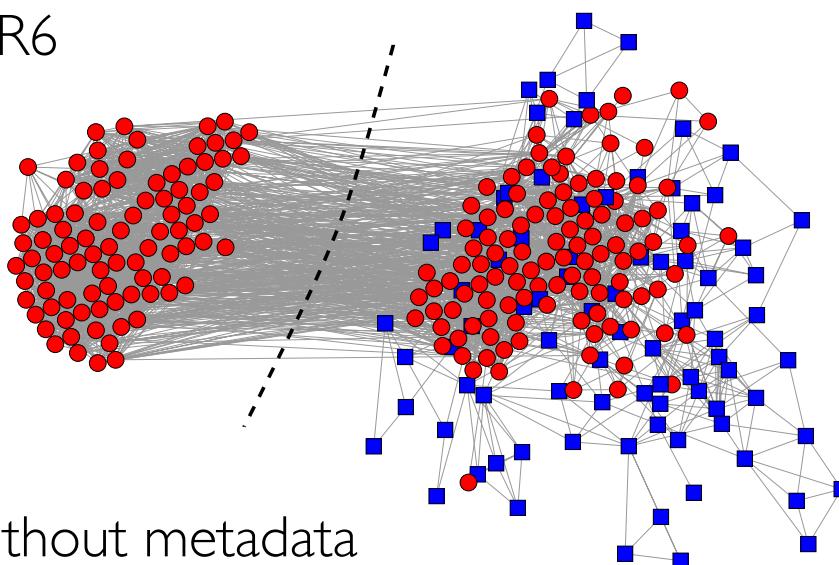
[2] Brose et al. (2005)

# real-world networks

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

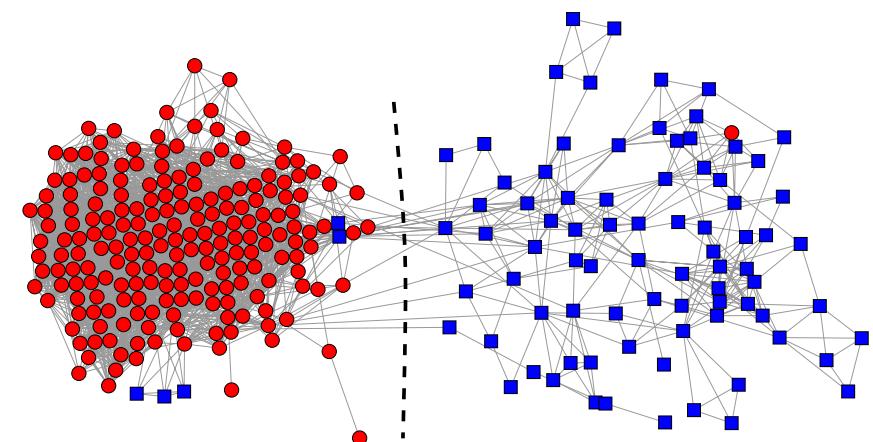
- **x = {Cys-PoLV labels for HVR6 region}**
- with metadata, partition discovers correlation with Cys labels (which are associated with severe disease)

HVR6



without metadata

$$\text{NMI} \in [0.077, 0.675]$$



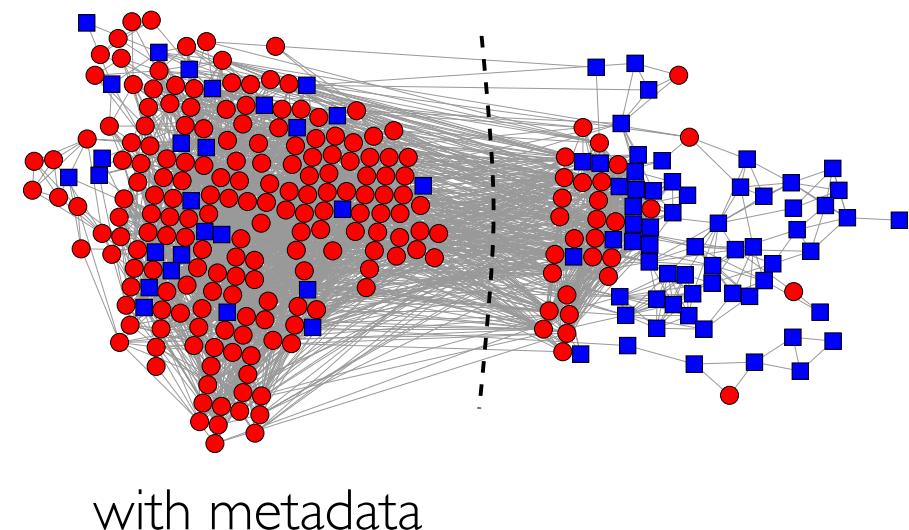
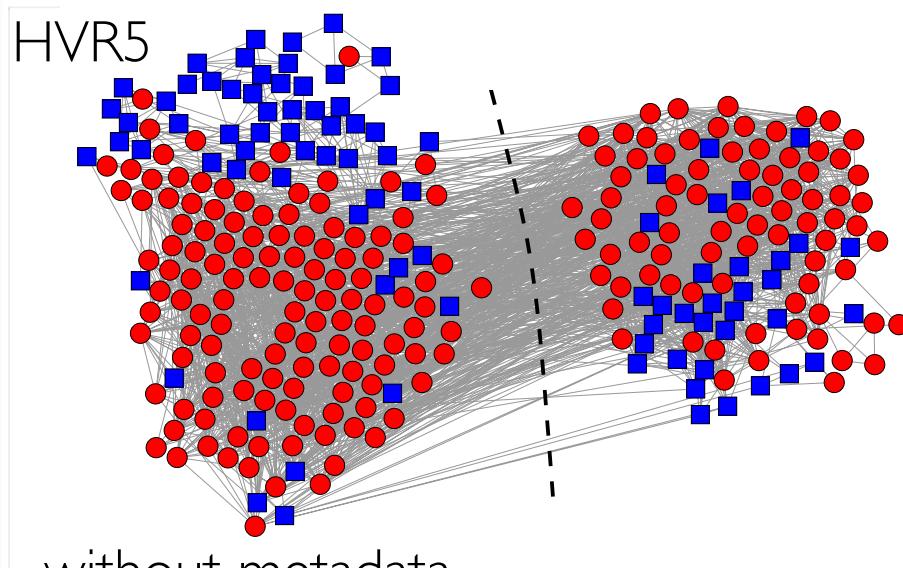
with metadata

$$\text{NMI} = 0.596$$

# real-world networks

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

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



# the ground truth about metadata

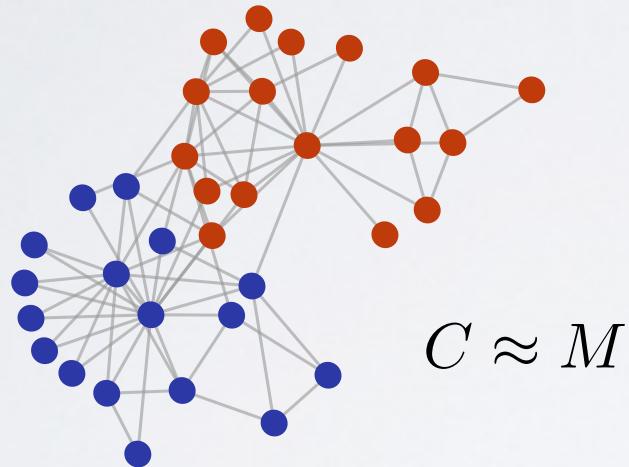
what is the goal of community detection?

network  $G$  + method  $f \rightarrow$  communities  $C = f(G)$  vs.  $M$  metadata

# the ground truth about metadata

what is the goal of community detection?

network  $G$  + method  $f \rightarrow$  communities  $C = f(G)$  vs.  $M$  metadata

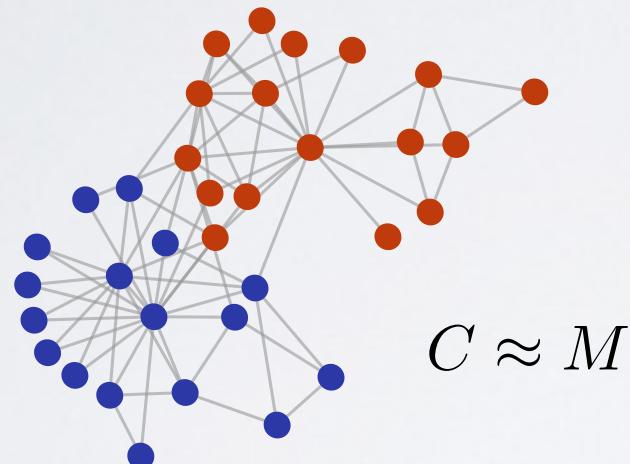


"this method works!"

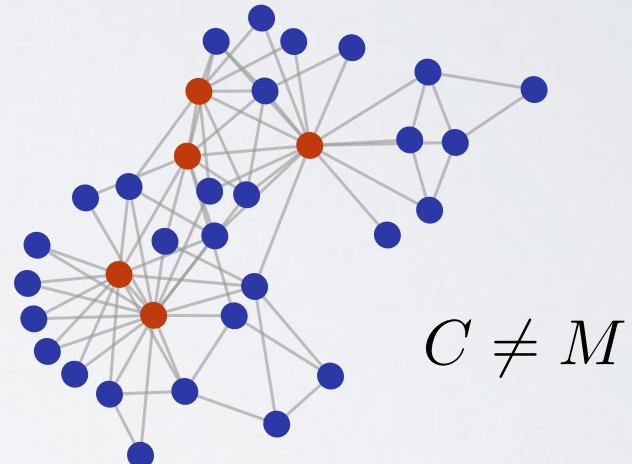
# the ground truth about metadata

what is the goal of community detection?

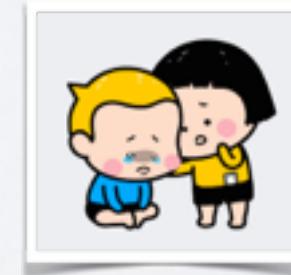
network  $G$  + method  $f \rightarrow$  communities  $C = f(G)$  vs.  $M$  metadata



"this method works!"



$C \neq M$



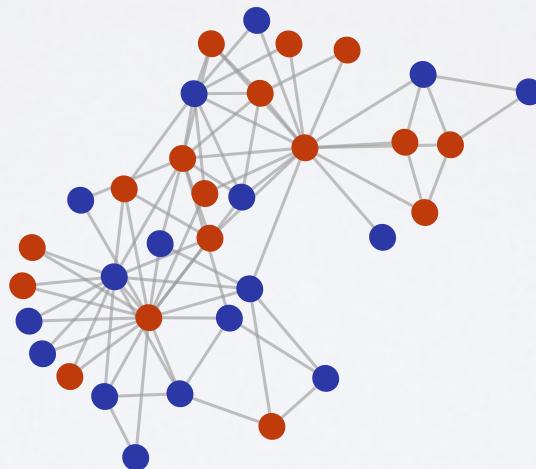
"this method stinks!"

# the ground truth about metadata

what is the goal of community detection?

there are 4 indistinguishable reasons why we might find  $f(G) = C \neq M$  :

- I. metadata  $M$  are unrelated to network structure  $G$

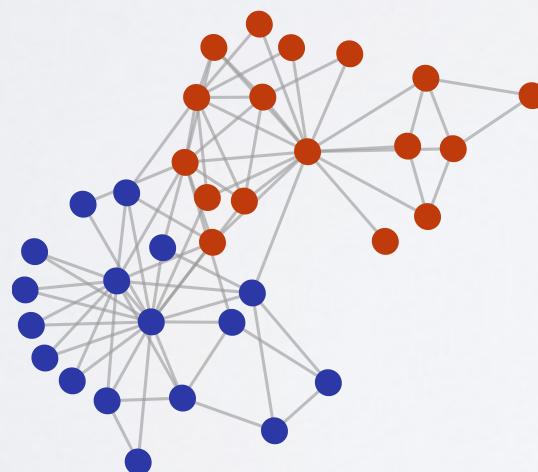


# the ground truth about metadata

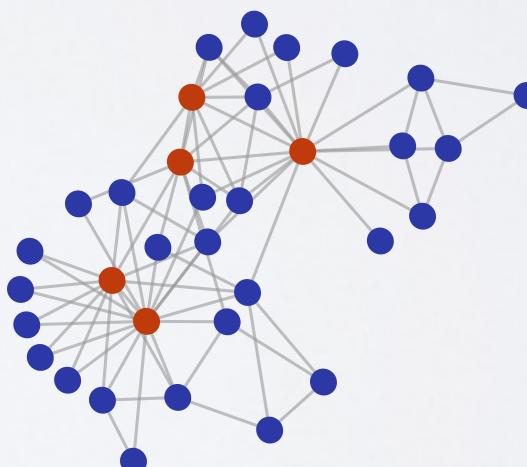
what is the goal of community detection?

there are 4 indistinguishable reasons why we might find  $f(G) = C \neq M$  :

1. metadata  $M$  are unrelated to network structure  $G$
2. metadata  $M$  and communities  $C$  capture different aspects of structure



social groups



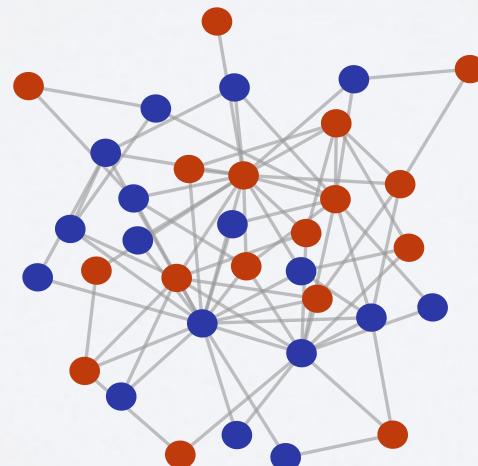
leaders and followers

# the ground truth about metadata

what is the goal of community detection?

there are 4 indistinguishable reasons why we might find  $f(G) = C \neq M$  :

1. metadata  $M$  are unrelated to network structure  $G$
2. metadata  $M$  and communities  $C$  capture different aspects of structure
3. network  $G$  has no community structure

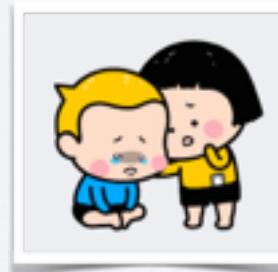


# the ground truth about metadata

## what is the goal of community detection?

there are 4 indistinguishable reasons why we might find  $f(G) = C \neq M$  :

1. metadata  $M$  are unrelated to network structure  $G$
2. metadata  $M$  and communities  $C$  capture different aspects of structure
3. network  $G$  has no community structure
4. algorithm  $f$  is bad



"this method stinks!"

# theorems for community detection

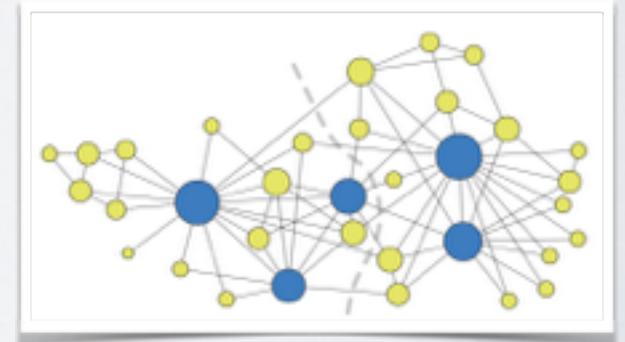
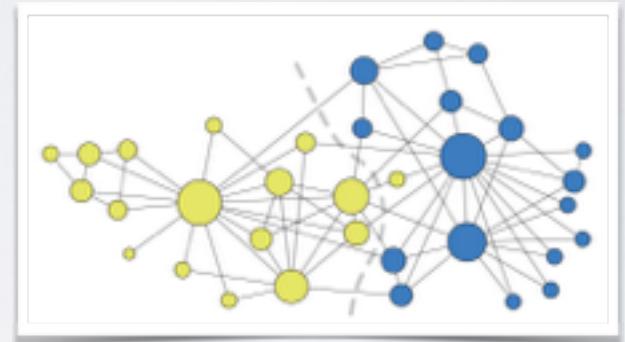
## I. Theorem: no bijection between ground truth and communities

$g(T) \rightarrow G \leftarrow g'(T')$     2 different processes, on 2 different ground truths, can create the same observed network

## 2. Theorem: No Free Lunch in community detection

no algorithm  $f$  has better performance than any other algorithm  $f'$ , when averaged over all possible inputs  $\{G\}$

→ good performance comes from matching algorithm  $f$  to its preferred subclass of networks  $\{G'\} \subset \{G\}$



[1] performance defined as adjusted mutual information (AMI), which is like the normalized mutual information, but adjusted for expected values

[2] original NFL theorem: Wolpert, *Neural Computation* (1996)

[3] proofs of these theorems is in Peel, Larremore, Clauset (2016)

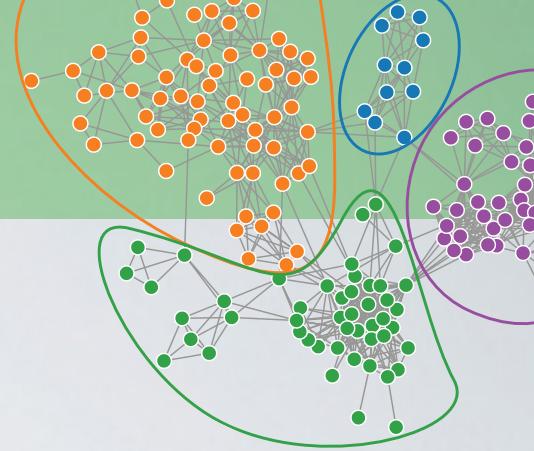
# theorems for community detection



DON'T TRY TO FIND THE GROUND TRUTH

INSTEAD ... TRY TO REALIZE THERE IS NO GROUND TRUTH

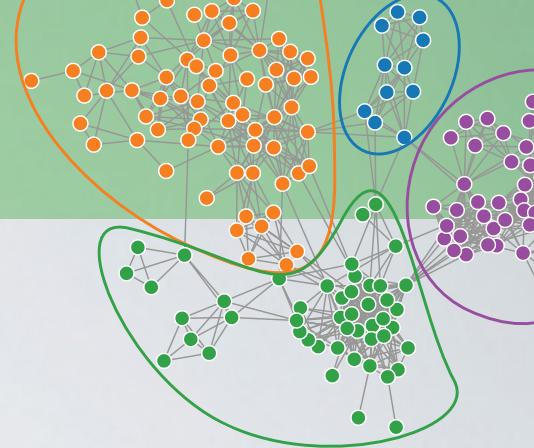
# conclusions



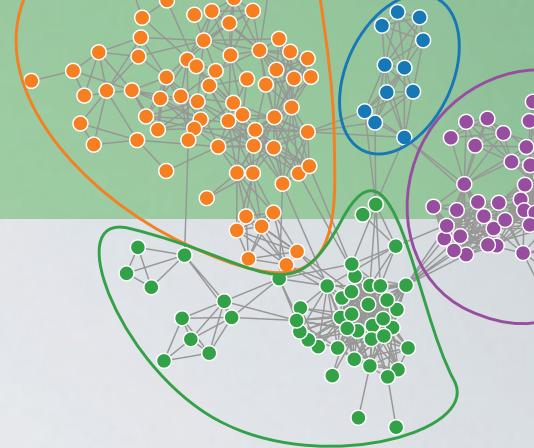
# conclusions

**the trouble with community detection:**

1. there is no ground truth in real networks
2. good performance depends on what we use the communities for



# conclusions



**the trouble with community detection:**

1. there is no ground truth in real networks
2. good performance depends on what we use the communities for

**metadata is just more data**

*use it to select structural communities that correlate with attributes  
(if they exist)*

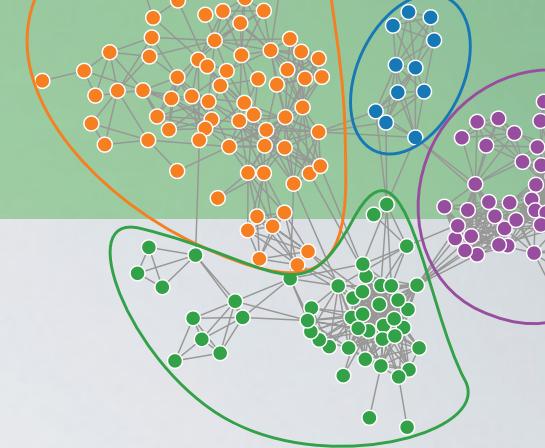
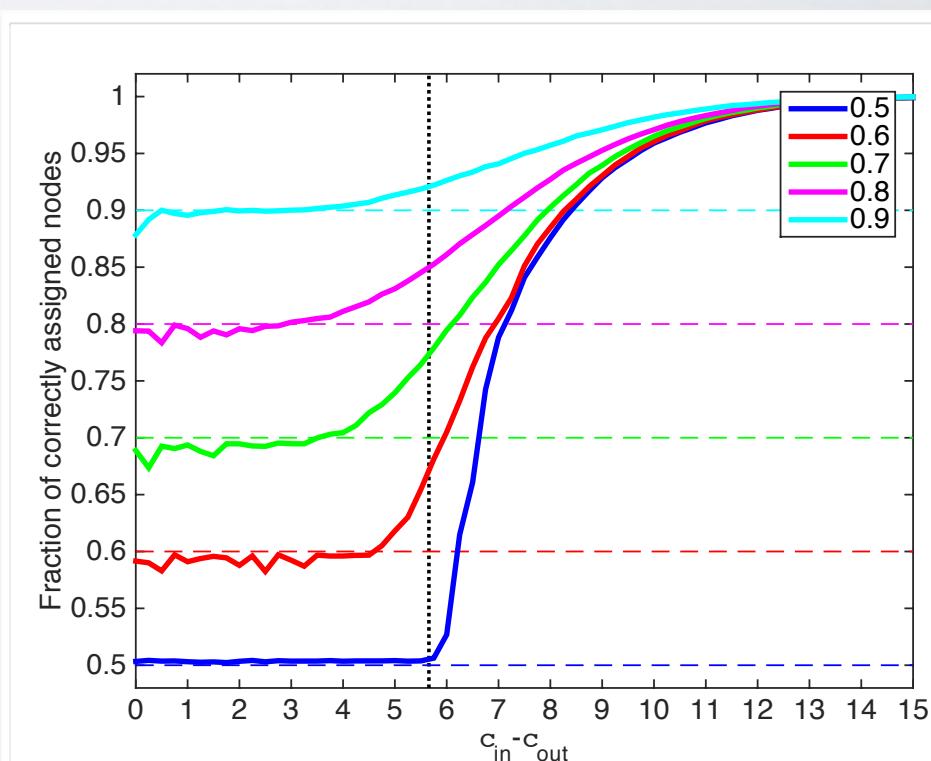
# conclusions

## a metadata-aware stochastic block model

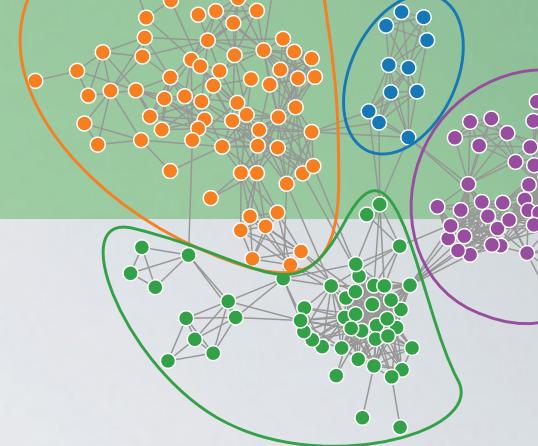
- probabilistic model of community structure and node metadata

$$P(\mathbf{A} | \Theta, \Gamma, \mathbf{x})$$

- **yields**: posterior probabilities of community labels  $q$ , the mixing matrix  $\Theta$ , and learned metadata-community association  $\Gamma$
- metaSBM performs better than any either *structure-only* or *metadata-only* algorithm
- highly scalable, via EM + belief propagation
- works well in practice



# future directions



- **a phase transition**  
does metadata eliminate or simply defer the detectability transition for sparse communities?
- **applications**  
how many networks with metadata have we overlooked because our algorithms didn't appear to work well on them?  
apply to all the data, learn all the science
- **extensions**  
to time-evolving networks, to multiplex networks, to other types of metadata (tags, vectors, etc.), etc.
- **more?**

# acknowledgements

Funding support:



Mark Newman  
(Michigan)



Leto Peel  
(Louvain)



Daniel B. Larremore  
(Santa Fe)



Benjamin H. Good  
(Harvard)



Yva de Montjoye  
(UCL)

## Structure and inference in annotated networks

M. E. J. Newman<sup>1,2</sup> and Aaron Clauset<sup>2,3</sup>

<sup>1</sup>*Department of Physics and Center for the Study of Complex Systems,  
University of Michigan, Ann Arbor, MI 48109, USA*

<sup>2</sup>*Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501*

<sup>3</sup>*Department of Computer Science and BioFrontiers Institute,  
University of Colorado, Boulder, CO 80309, USA*

## The ground truth about metadata and community detection in networks

Leto Peel,<sup>1,2,\*</sup> Daniel B. Larremore,<sup>3,†</sup> and Aaron Clauset<sup>4,5,3,‡</sup>

<sup>1</sup>*ICTEAM, Université Catholique de Louvain, Louvain-la-Neuve, Belgium*

<sup>2</sup>*naXys, Université de Namur, Namur, Belgium*

<sup>3</sup>*Santa Fe Institute, Santa Fe, NM 87501, USA*

<sup>4</sup>*Department of Computer Science, University of Colorado, Boulder, CO 80309, USA*

<sup>5</sup>*BioFrontiers Institute, University of Colorado, Boulder, CO 80309, USA*

PHYSICAL REVIEW E 81, 046106 (2010)

## Performance of modularity maximization in practical contexts

Benjamin H. Good,<sup>1,2,\*</sup> Yves-Alexandre de Montjoye,<sup>3,2,†</sup> and Aaron Clauset<sup>2,‡</sup>

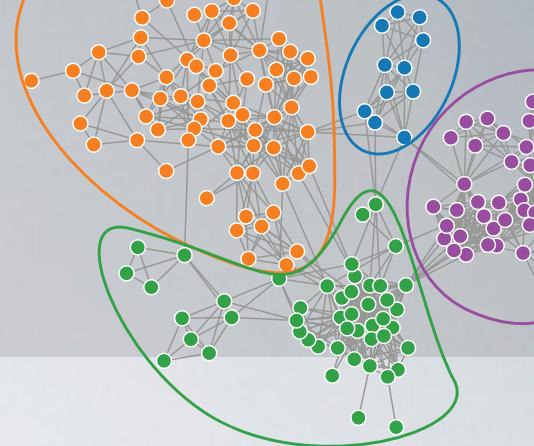
<sup>1</sup>*Department of Physics, Swarthmore College, Swarthmore, Pennsylvania 19081, USA*

<sup>2</sup>*Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, New Mexico 87501, USA*

<sup>3</sup>*Department of Applied Mathematics, Université Catholique de Louvain, 4 Avenue Georges Lemaitre, B-1348 Louvain-la-Neuve, Belgium*

(Received 1 October 2009; published 15 April 2010)

**fin**



**community detection**

~~EVERYTHING IS~~

**AWESOME!**



# real-world networks

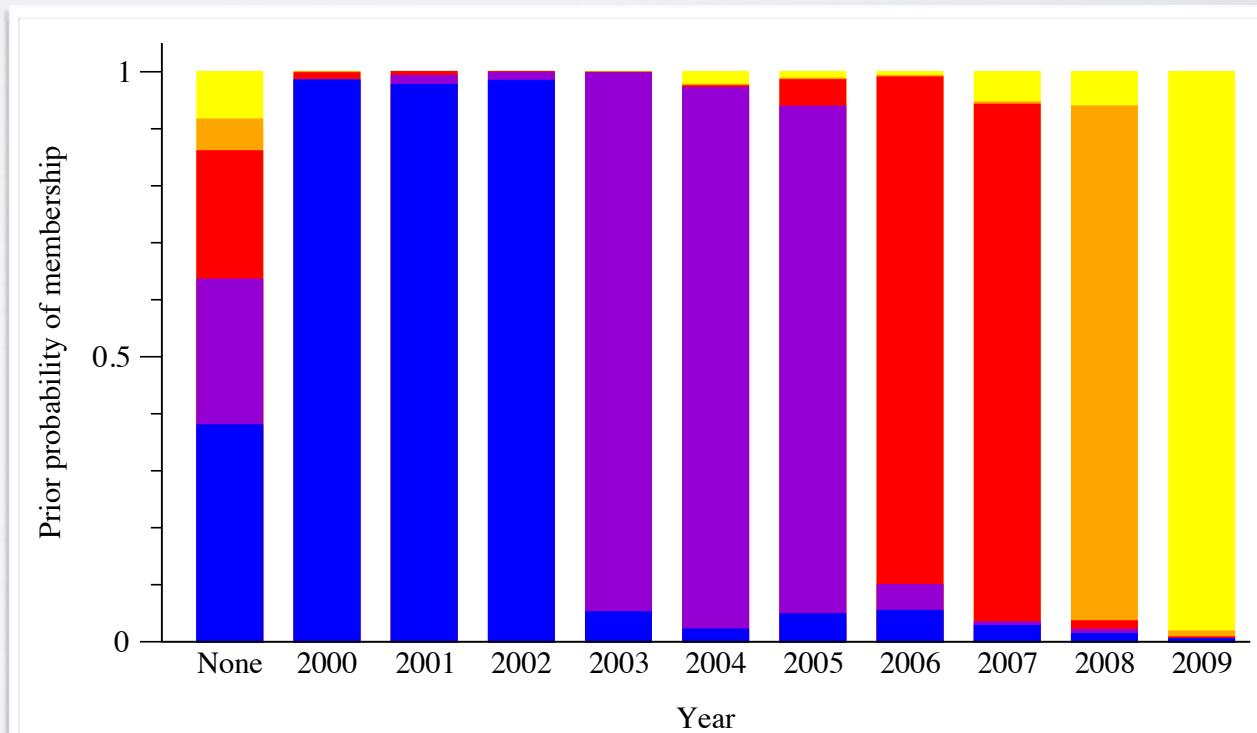
4. **Facebook friendships:** online friendships among 15,126 Harvard students and alumni (in Sept. 2005)

- $x = \{\text{graduation year}, \text{dormitory}\}$
- method finds a good partition between alumni, recent graduates, upperclassmen, sophomores, and freshmen

$$\text{NMI} = 0.668$$

- without metadata:

$$\text{NMI} \in [0.573, 0.641]$$



# real-world networks

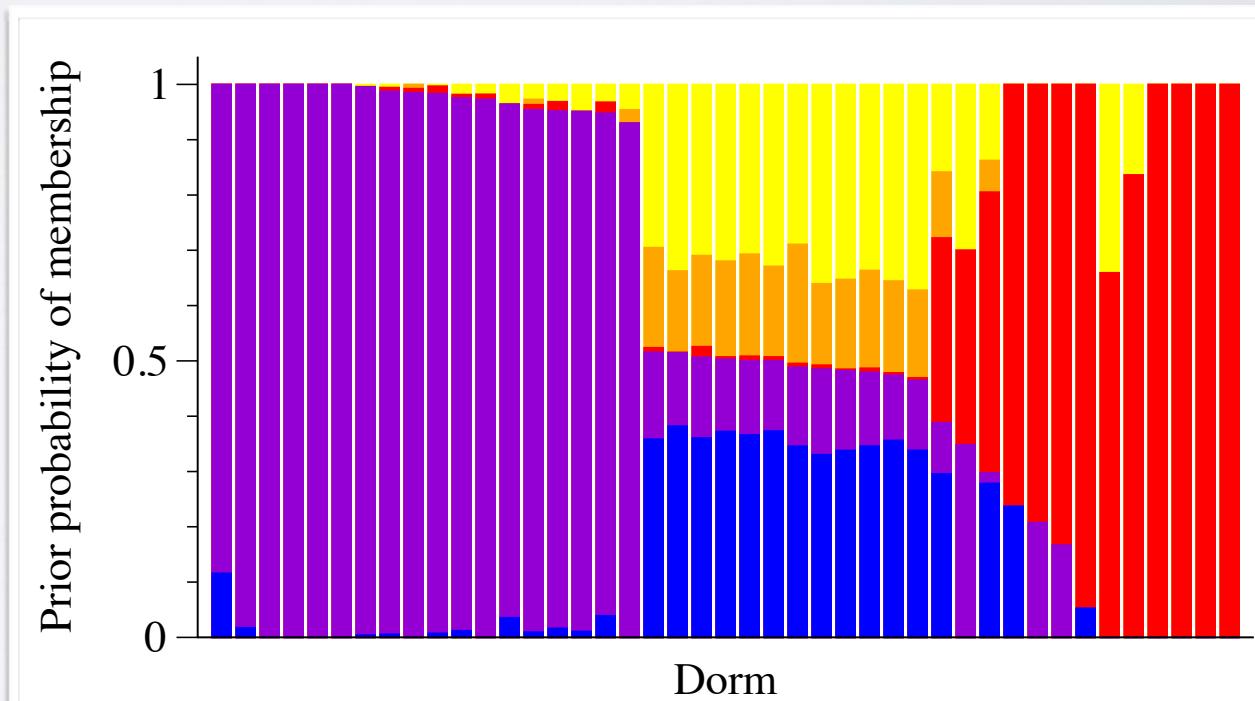
4. **Facebook friendships:** online friendships among 15,126 Harvard students and alumni (in Sept. 2005)

- $\mathbf{x} = \{\text{graduation year, dormitory}\}$
- method finds a good partition among the dorms

$$\text{NMI} = 0.255$$

- without metadata:

$$\text{NMI} \in [0.074, 0.224]$$



# real-world networks

5. **Internet graph:** 262,953 peering relations among 46,676 Autonomous Systems

- $x = \{\text{country location of AS}\}$
- method finds a good partition along the lines of the 173 countries

$$\text{NMI} = 0.870$$

- without metadata:

$$\text{NMI} \in [0.398, 0.626]$$