

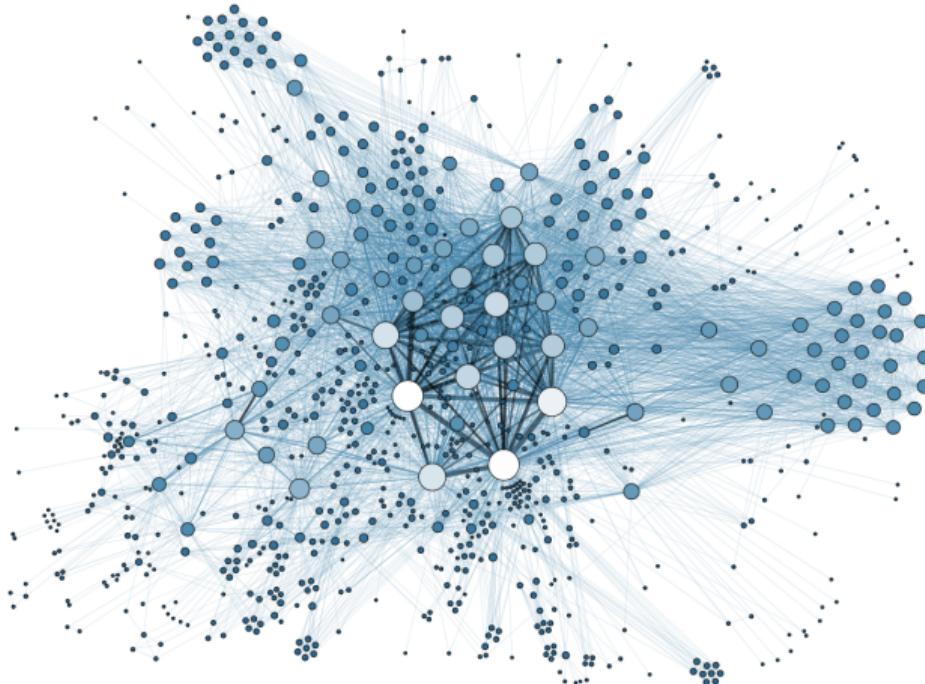
Disentangling homophily, community structure and triadic closure in networks

Tiago P. Peixoto

*Department of Network and Data Science
Central European University, Vienna*

Austin, April 2023

NETWORKS AS RELATIONAL BIG DATA

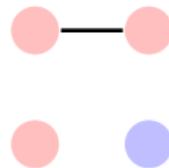


- ▶ Often large
- ▶ High-dimensional
- ▶ Sparse
- ▶ Heterogeneously structured
- ▶ Richly annotated
- ▶ Often dynamic

Often the outcome of multiple generative mechanisms: how do we disentangle them?

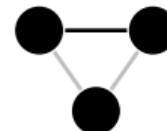
HOMOPHILY VS. TRIADIC CLOSURE

Homophily



Tendency of an edge being placed between nodes of the same kind (e.g. age, race, location, etc.)

Triadic closure



Tendency of an edge being placed between nodes if they share a common neighbor.

Both processes induce the formation of triangles and community structure.

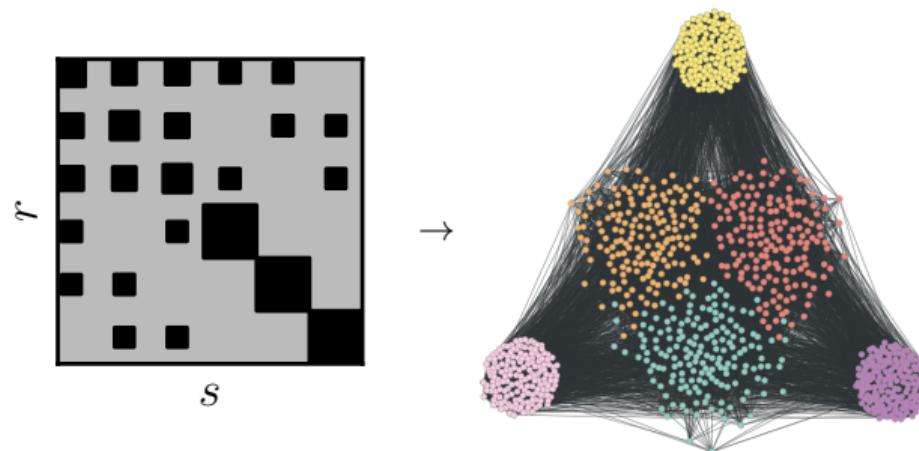
HOW DO WE MODEL HOMOPHILY?

THE STOCHASTIC BLOCK MODEL (SBM)

Planted partition: N nodes divided into B groups.

Parameters: $b_i \rightarrow$ group membership of node i

$\lambda_{rs} \rightarrow$ edge probability from group r to s .

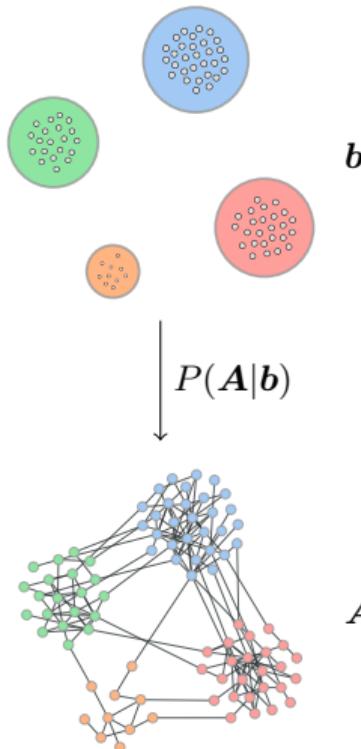


Degree-corrected: Arbitrary degree sequence: $\{k_i\}$

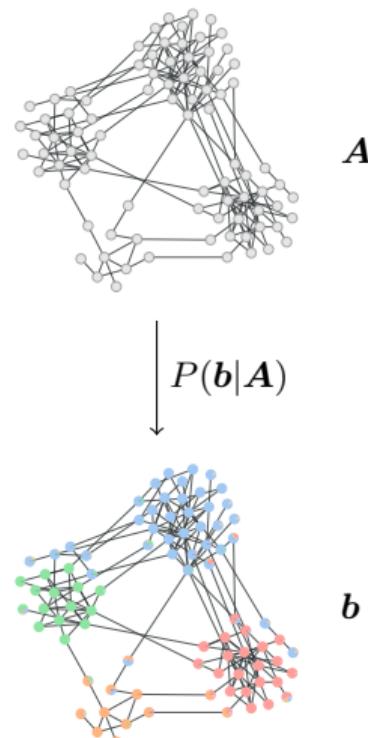
- ▶ Not restricted to homophily; can accommodate arbitrary mixing.

HOW DO WE DETECT HOMOPHILY?

Generative model



Statistical inference



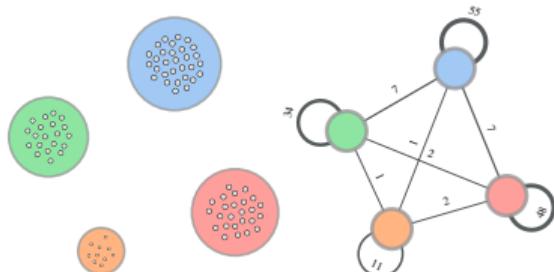
Bayes' rule:

$$P(b|A) = \frac{P(A|b)P(b)}{P(A)}$$

- ✓ **Avoids overfitting**
Incorporates Occam's razor; can be formally mapped into *data compression*. Can also be used to find the number of groups.
- ✓ **Efficient**
Scales to large networks.

MICROCANONICAL DEGREE-CORRECTED STOCHASTIC BLOCK MODEL (DC-SBM)

Microcanonical priors and likelihood

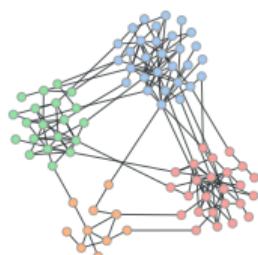


$P(\mathbf{b})$

$P(\mathbf{e}|\mathbf{b})$



$P(\mathbf{k}|e, \mathbf{b})$



$P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})$

$$P(\mathbf{b}) = \frac{\prod_r n_r!}{N!} \times \binom{N-1}{B-1}^{-1} \times \frac{1}{N} \quad P(\mathbf{e}|\mathbf{b}) = \left(\binom{\binom{B}{2}}{E} \right)^{-1}$$

$$P(\mathbf{k}|\mathbf{e}, \mathbf{b}) = \prod_r \binom{n_r}{e_r}^{-1} \quad P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b}) = \frac{\prod_{i < j} A_{ij}! \prod_{ii} A_{ii}!! \prod_r e_r!}{\prod_{r < s} e_{rs}! \prod_r e_{rr}!! \prod_i k_i!}$$

Nonparametric posterior distribution:

$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})P(\mathbf{k}|e, \mathbf{b})P(\mathbf{e}|\mathbf{b})P(\mathbf{b})}{P(\mathbf{A})}$$

Description length:

$$\Sigma = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{Data}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{Model}}$$

Inference \leftrightarrow compression

MINIMUM DESCRIPTION LENGTH (MDL)

$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

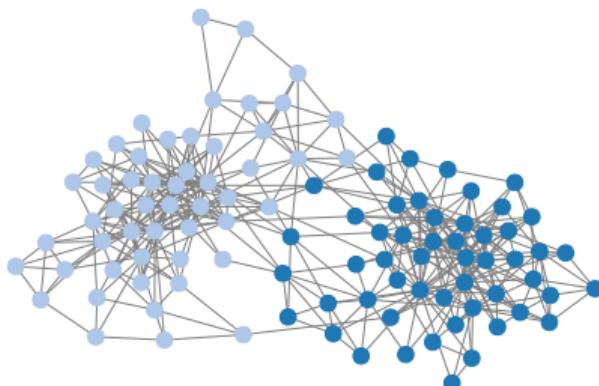
$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$

MINIMUM DESCRIPTION LENGTH (MDL)

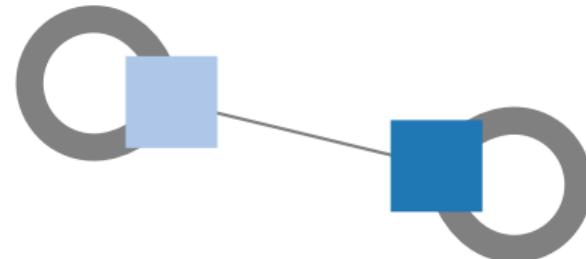
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$$B = 2, \mathcal{S} \approx 1395.8 \text{ bits}$$



$$\text{Model, } \mathcal{L} \approx 581.6 \text{ bits}$$

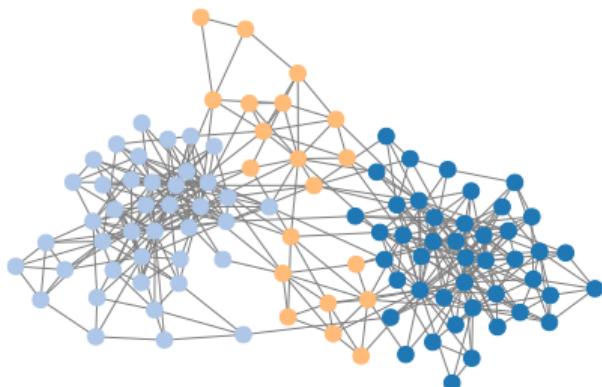
$$\Sigma \approx 1977.5 \text{ bits}$$

MINIMUM DESCRIPTION LENGTH (MDL)

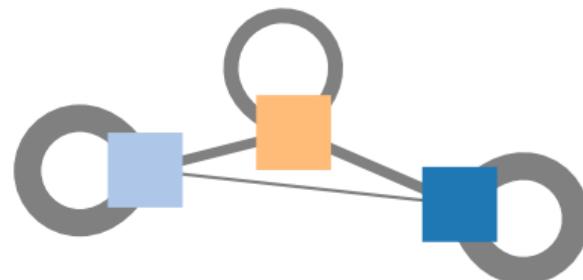
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$$B = 3, \mathcal{S} \approx 1285.7 \text{ bits}$$



$$\text{Model, } \mathcal{L} \approx 654.1 \text{ bits}$$

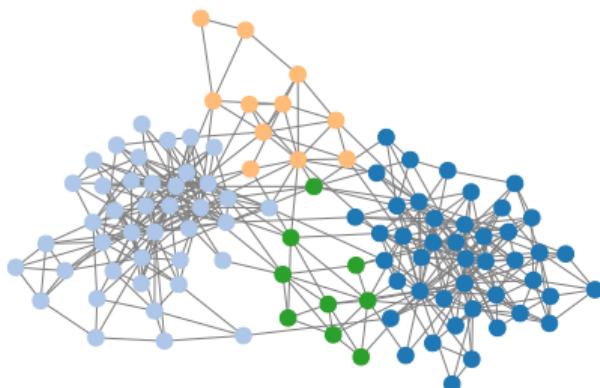
$$\Sigma \approx 1939.7 \text{ bits}$$

MINIMUM DESCRIPTION LENGTH (MDL)

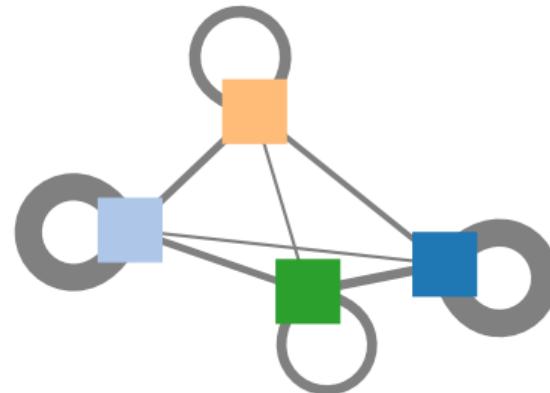
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$$B = 4, \mathcal{S} \approx 1246.9 \text{ bits}$$



$$\text{Model, } \mathcal{L} \approx 693.2 \text{ bits}$$

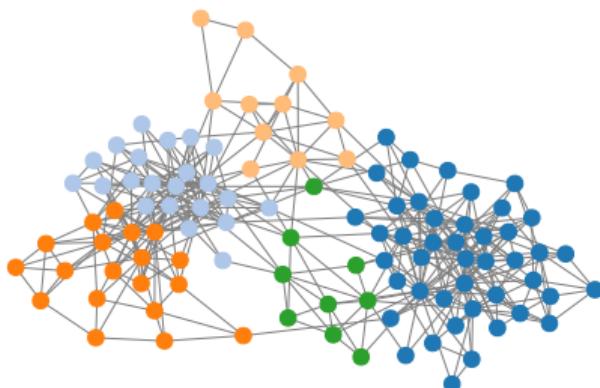
$$\Sigma \approx 1940.2 \text{ bits}$$

MINIMUM DESCRIPTION LENGTH (MDL)

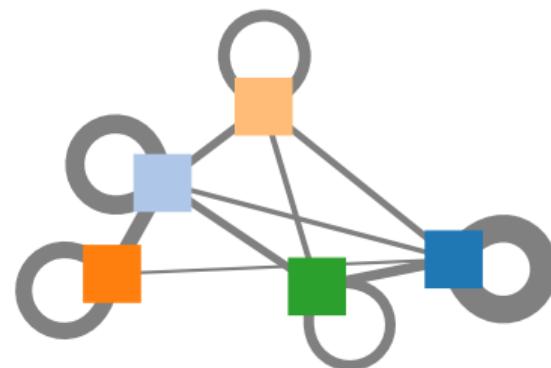
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$B = 5, \mathcal{S} \approx 1195.2$ bits



Model, $\mathcal{L} \approx 747.9$ bits

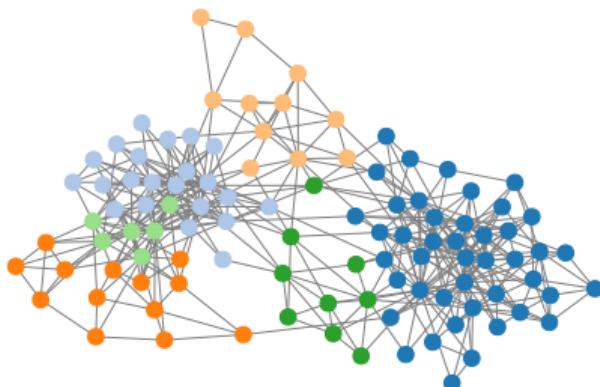
$\Sigma \approx 1943.2$ bits

MINIMUM DESCRIPTION LENGTH (MDL)

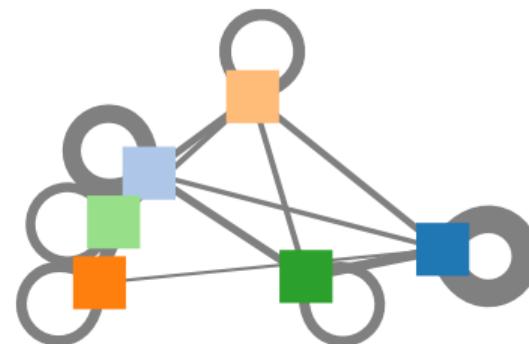
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$B = 6, \mathcal{S} \approx 1165.1$ bits



Model, $\mathcal{L} \approx 780.7$ bits

$\Sigma \approx 1951.9$ bits

MINIMUM DESCRIPTION LENGTH (MDL)

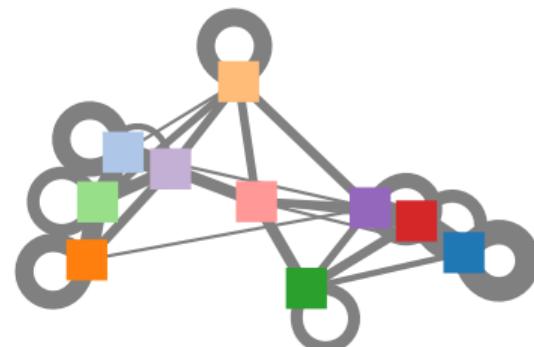
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$$B = 10, \mathcal{S} \approx 1028.4 \text{ bits}$$



$$\text{Model, } \mathcal{L} \approx 942.7 \text{ bits}$$

$$\Sigma \approx 1971.1 \text{ bits}$$

MINIMUM DESCRIPTION LENGTH (MDL)

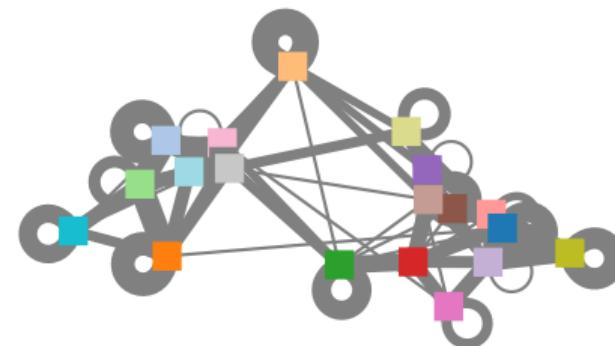
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$$B = 20, \mathcal{S} \approx 808.9 \text{ bits}$$



$$\text{Model, } \mathcal{L} \approx 1280.5 \text{ bits}$$

$$\Sigma \approx 2089.5 \text{ bits}$$

MINIMUM DESCRIPTION LENGTH (MDL)

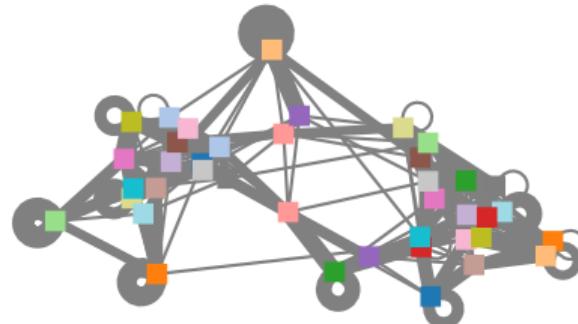
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$B = 40, \mathcal{S} \approx 497.1$ bits



Model, $\mathcal{L} \approx 1748.4$ bits

$\Sigma \approx 2245.6$ bits

MINIMUM DESCRIPTION LENGTH (MDL)

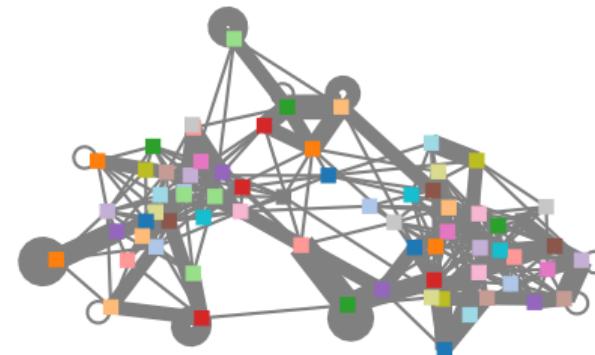
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$$B = 70, \mathcal{S} \approx 178.5 \text{ bits}$$



$$\text{Model, } \mathcal{L} \approx 2272.4 \text{ bits}$$

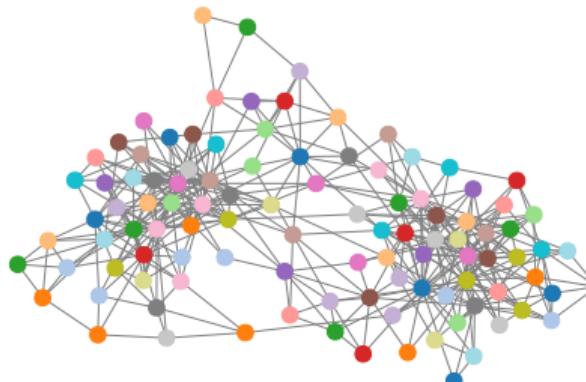
$$\Sigma \approx 2450.9 \text{ bits}$$

MINIMUM DESCRIPTION LENGTH (MDL)

$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



$B = N$, $\mathcal{S} = 0$ bits

$\Sigma \approx 2568.2$ bits



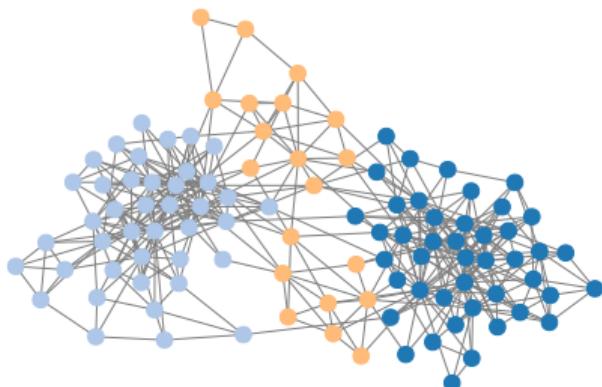
Model, $\mathcal{L} \approx 2568.2$ bits

MINIMUM DESCRIPTION LENGTH (MDL)

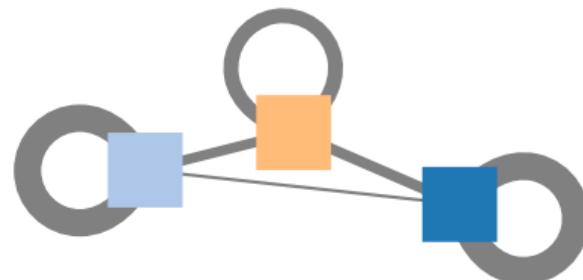
$$P(\mathbf{b}|\mathbf{A}) = \frac{P(\mathbf{A}, \mathbf{b})}{P(\mathbf{A})} = \frac{P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b})}{P(\mathbf{A})} = \frac{2^{-\Sigma}}{P(\mathbf{A})}$$

Description length:

$$\Sigma = -\log_2 P(\mathbf{A}, \mathbf{k}, \mathbf{e}, \mathbf{b}) = \underbrace{-\log_2 P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{data|model, } \mathcal{S}} - \underbrace{\log_2 P(\mathbf{k}, \mathbf{e}, \mathbf{b})}_{\text{model, } \mathcal{L}}$$



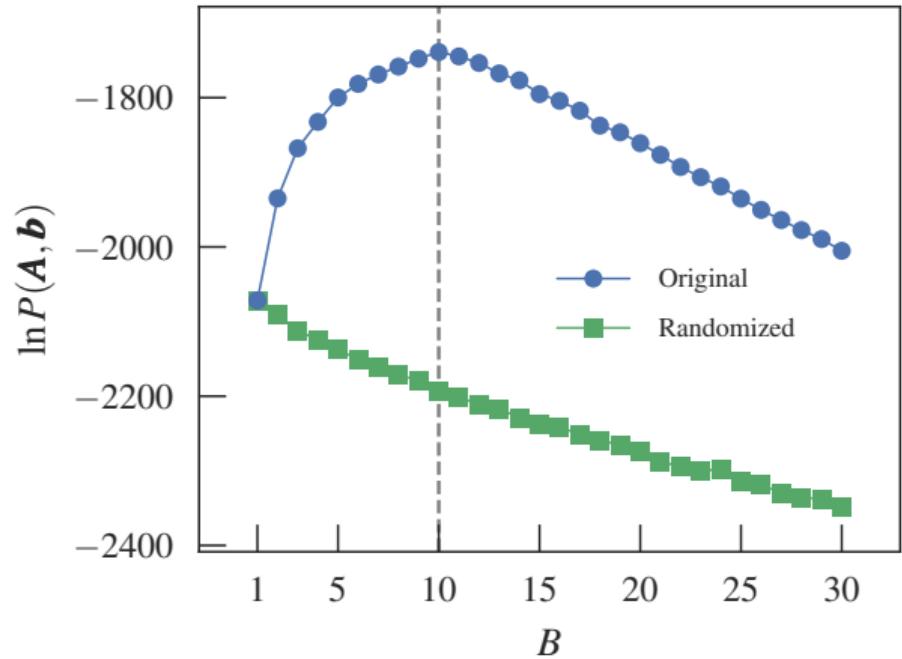
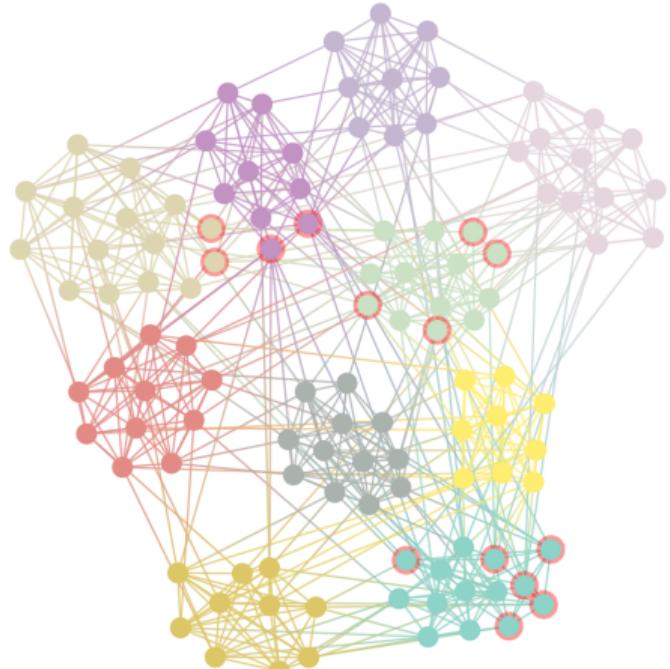
$$B = 3, \mathcal{S} \approx 1285.7 \text{ bits}$$



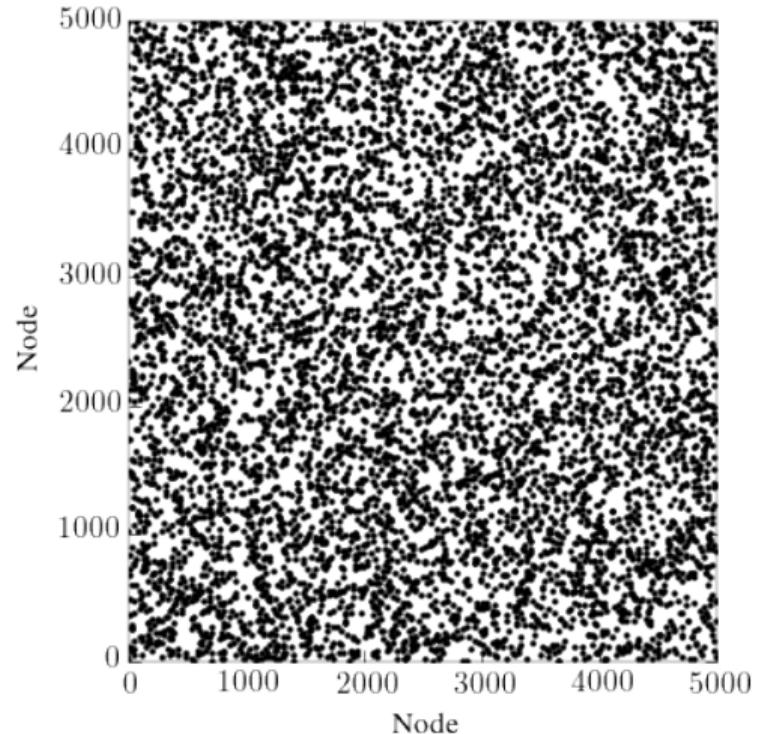
$$\text{Model, } \mathcal{L} \approx 654.1 \text{ bits}$$

$$\Sigma \approx 1939.7 \text{ bits}$$

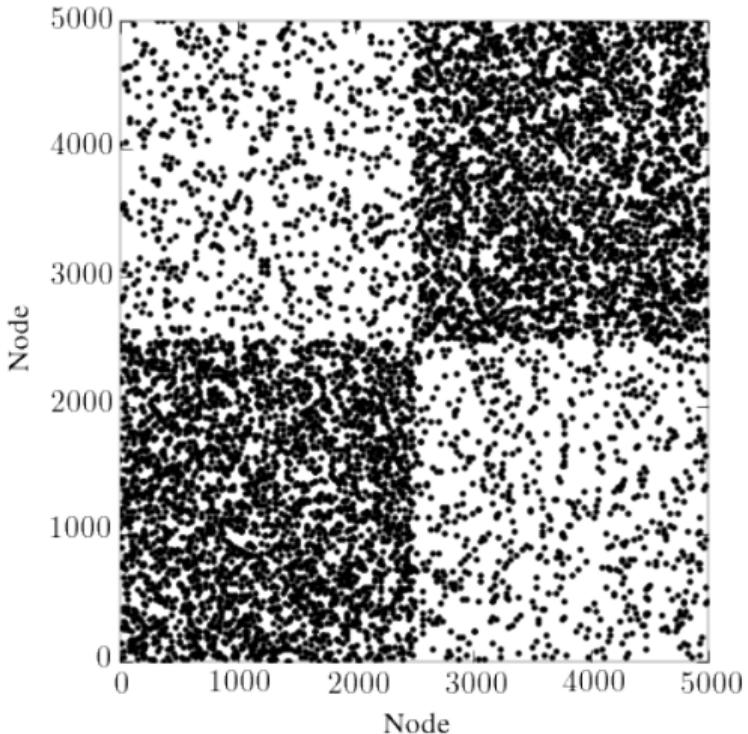
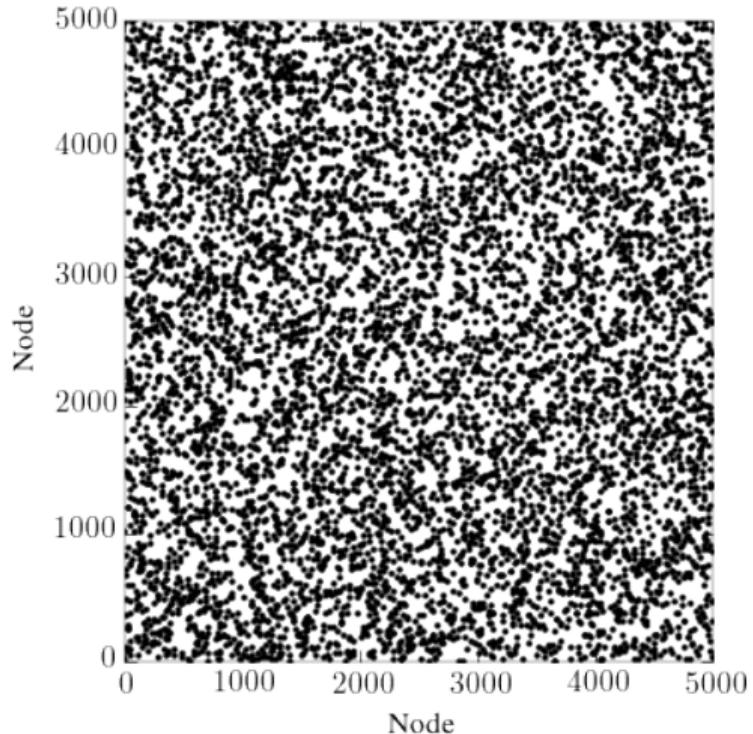
EXAMPLE: AMERICAN COLLEGE FOOTBALL TEAMS



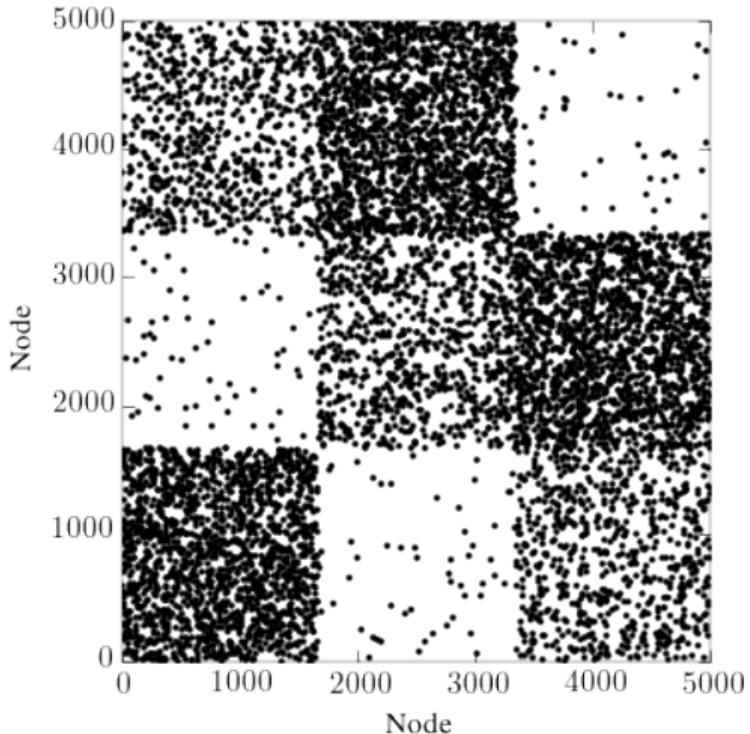
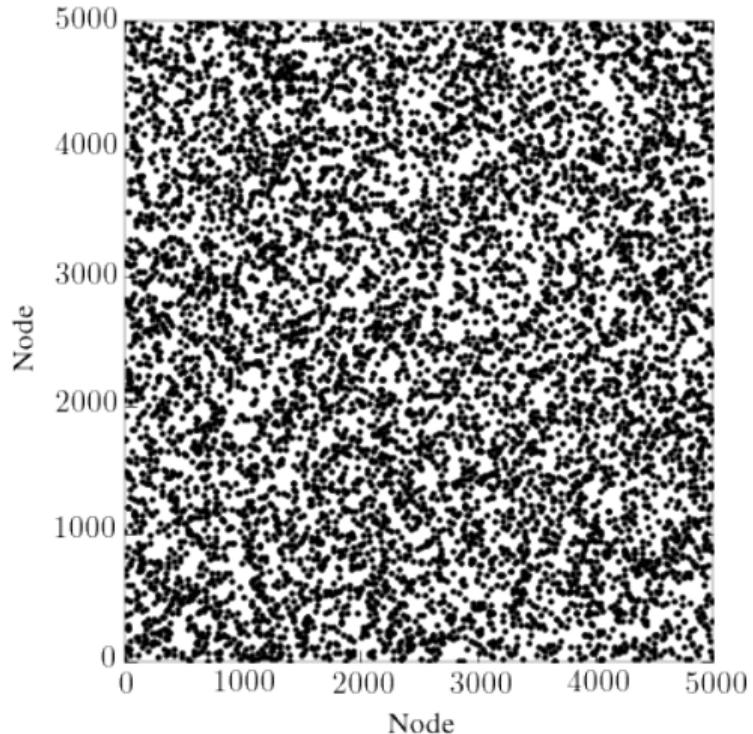
STRUCTURE VS. RANDOMNESS



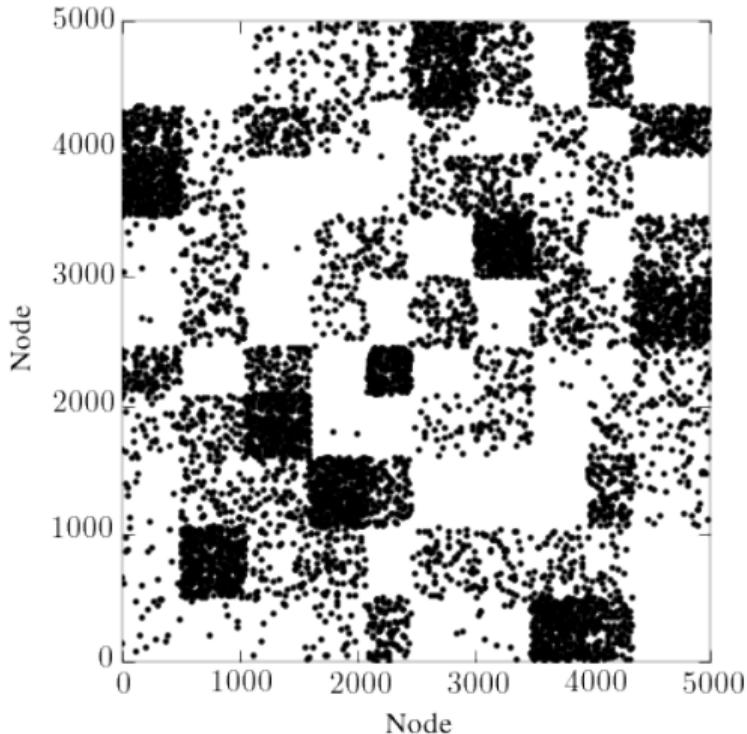
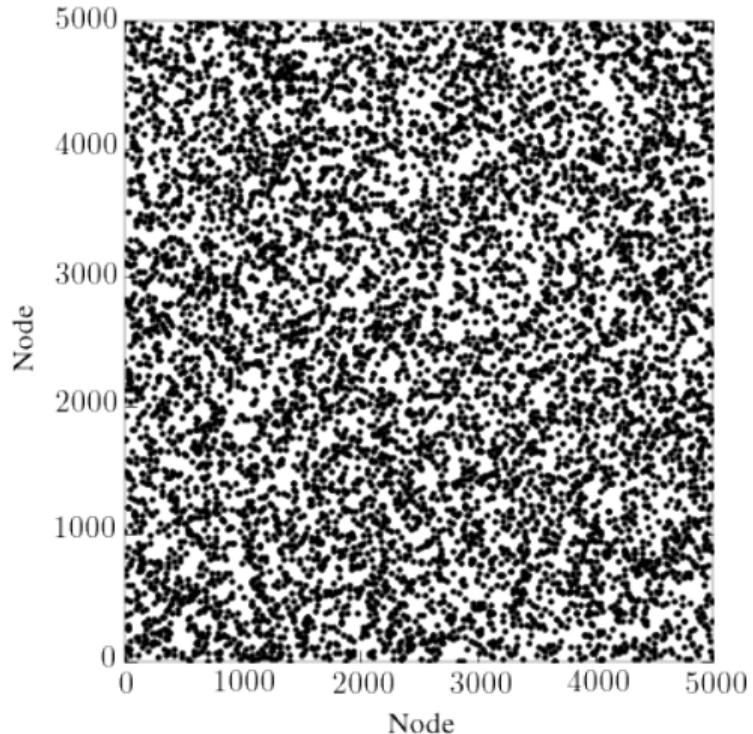
STRUCTURE VS. RANDOMNESS



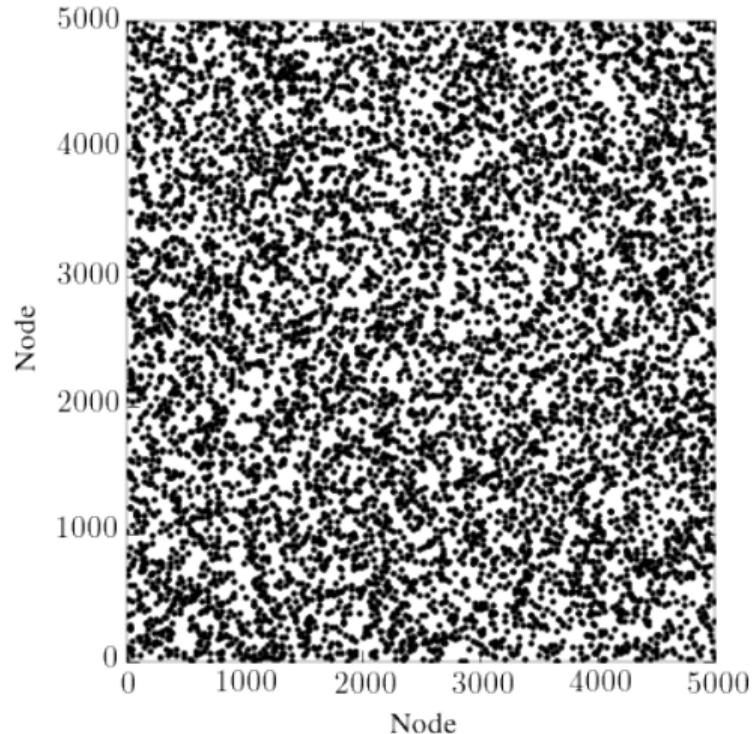
STRUCTURE VS. RANDOMNESS



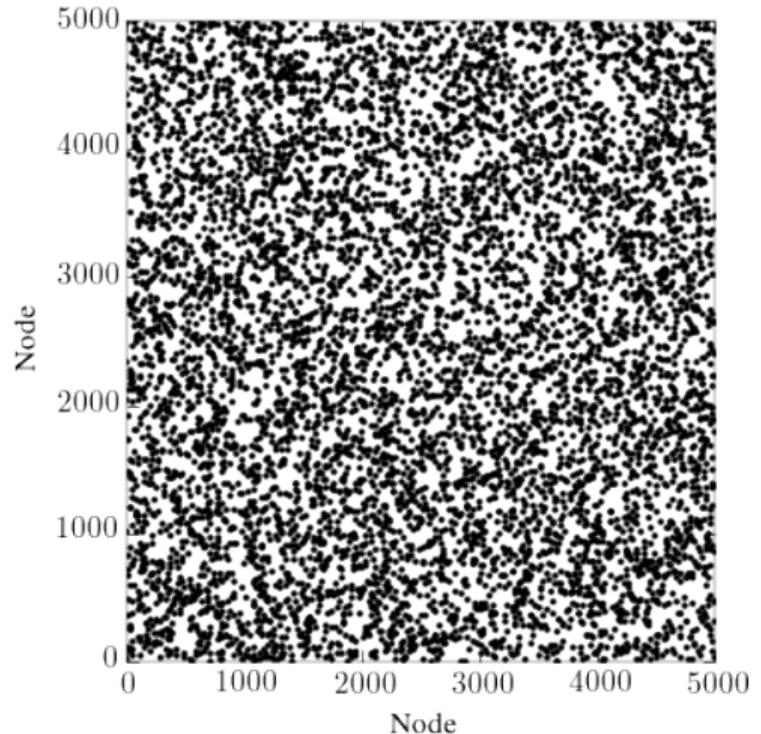
STRUCTURE VS. RANDOMNESS



STRUCTURE VS. RANDOMNESS



STRUCTURE VS. RANDOMNESS



HOMOPHILY CAN INDUCE TRIANGLES



$C \in [0, 1] \rightarrow$ clustering coefficient (fraction of connected triads that form a triangle).

Assuming uniform group size and density:

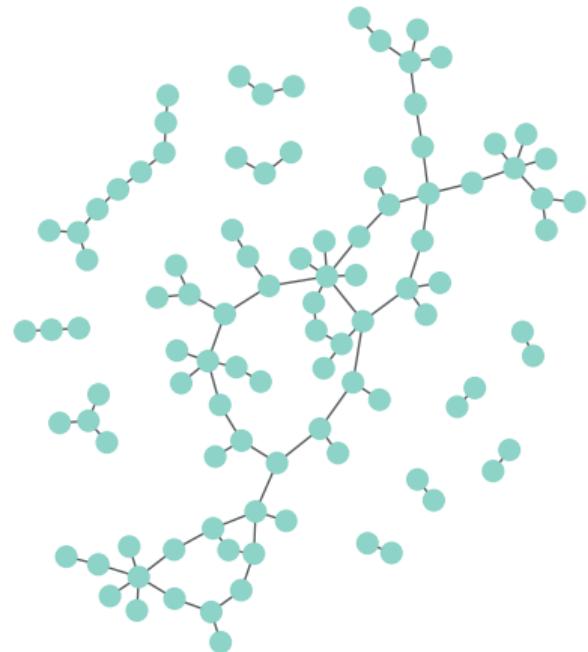
$$C = O\left(\frac{\langle k \rangle B}{N}\right)$$

Large values of C can be obtained for globally sparse networks, with a sufficiently large B .

However, for B and k fixed, and $N \rightarrow \infty$, we have $C \rightarrow 0$.

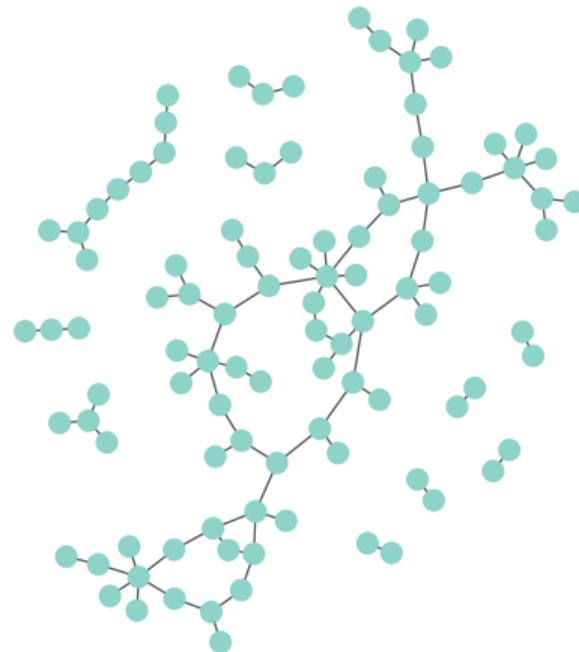
(Typical real networks do not look like this.)

TRIADIC CLOSURE “INDUCES” COMMUNITY STRUCTURE

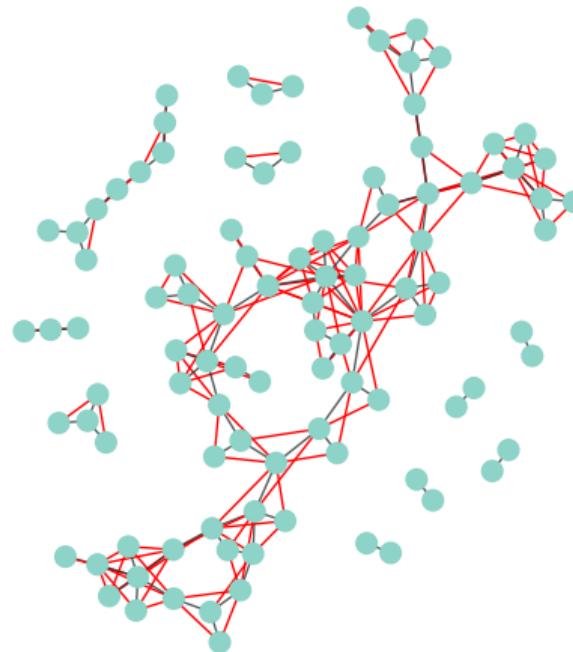


(a) Fully random network (no homophily)

TRIADIC CLOSURE “INDUCES” COMMUNITY STRUCTURE

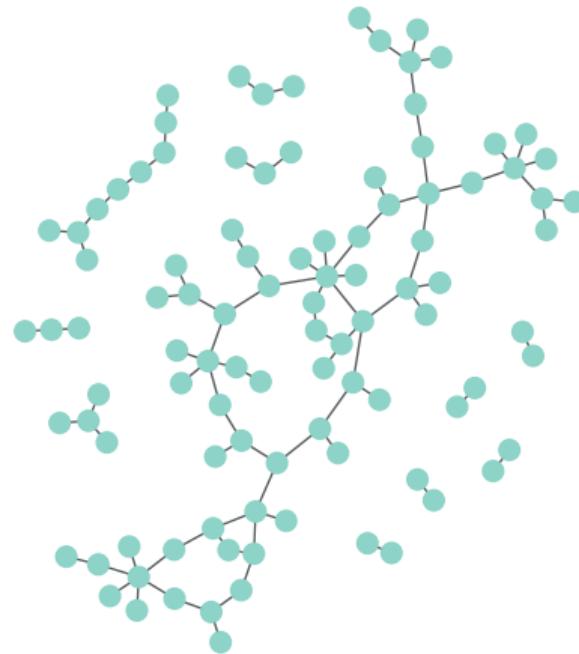


(a) Fully random network (no homophily)

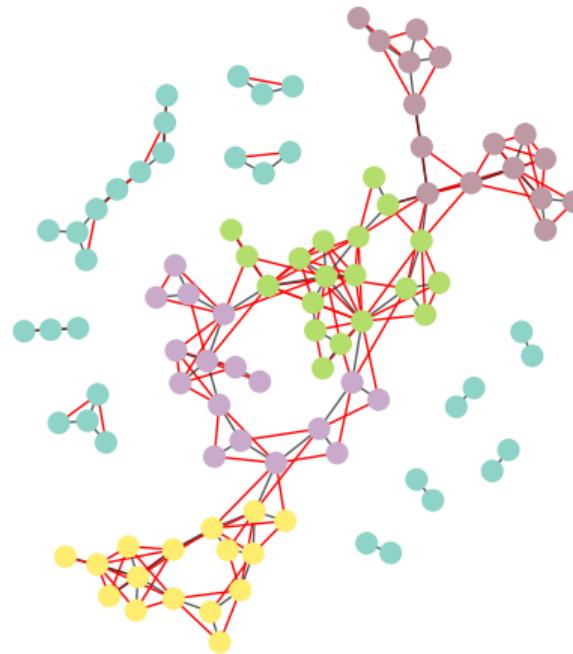


(b) Triadic closure edges

TRIADIC CLOSURE “INDUCES” COMMUNITY STRUCTURE



(a) Fully random network (no homophily)

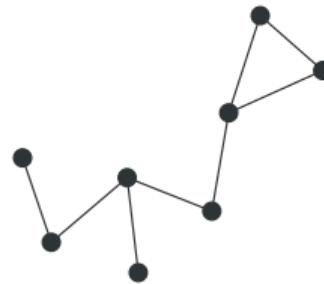


(b) Triadic closure edges

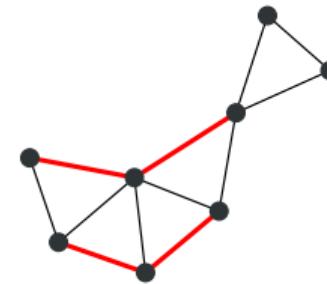
A LATENT TRIADIC CLOSURE MODEL

Generative process

Seminal edges



Triadic closure



Observed network

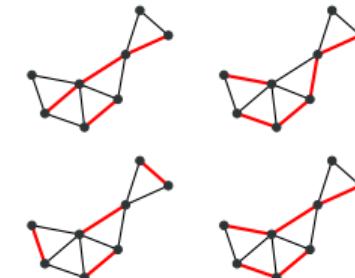


Statistical inference

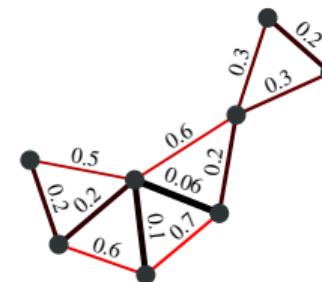
Observed network



Posterior distribution

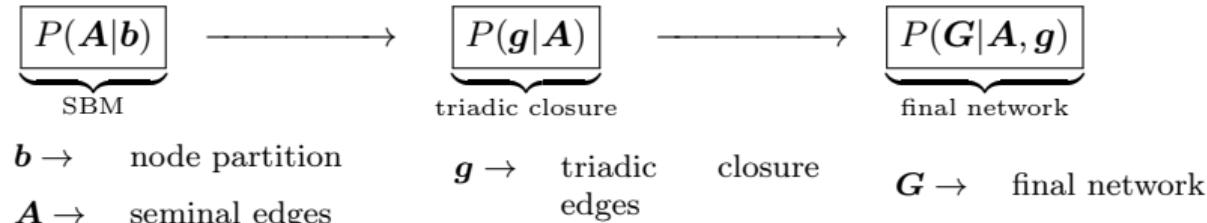


Marginal probabilities



SBM + TRIADIC CLOSURE (SBM/TC)

Generative process



Bayesian inference

$$P(\mathbf{g}, \mathbf{A}, \mathbf{b} | \mathbf{G}) = \frac{P(\mathbf{G} | \mathbf{A}, \mathbf{g}) P(\mathbf{g} | \mathbf{A}) P(\mathbf{A} | \mathbf{b}) P(\mathbf{b})}{P(\mathbf{G})}$$

Model selection & description length

Pure SBM

$$\begin{aligned}\Sigma_{\text{SBM}} &= -\ln P(\mathbf{G}, \mathbf{b}) \\ &= \underbrace{-\ln P(\mathbf{G}|\mathbf{b})}_{\text{Data}} - \underbrace{\ln P(\mathbf{b})}_{\text{Model}}\end{aligned}$$

SBM/TC

$$\begin{aligned}\Sigma_{\text{SBM/TC}} &= -\ln P(\mathbf{G}, \mathbf{A}, \mathbf{g}, \mathbf{b}) \\ &= -\ln P(\mathbf{G}|\mathbf{A}, \mathbf{g}) - \ln P(\mathbf{g}|\mathbf{A}) \\ &\quad - \ln P(\mathbf{A}|\mathbf{b}) - \ln P(\mathbf{b})\end{aligned}$$

HOW DO WE CLOSE TRIANGLES?

Given a network \mathbf{A} sampled from the SBM, with probability $P(\mathbf{A}|\mathbf{b})$.

To each node u we consider its “ego” graph $\mathbf{g}(u)$ with adjacency matrix $g_{ij}(u)$.

An open triad (i, u, j) exists in \mathbf{A} iff $m_{ij}(u) = 1$, with

$$m_{ij}(u) = A_{ui}A_{uj}(1 - A_{ij}).$$

Each open triad is closed with probability p_u , leading to

$$P(\mathbf{g}(u)|\mathbf{A}, p_u) = \prod_{i < j} [p_u m_{ij}(u)]^{g_{ij}(u)} [1 - p_u m_{ij}(u)]^{1-g_{ij}(u)}.$$

We then erase the possible multiple edges from the union of all $\mathbf{g}(u)$, yielding a final graph \mathbf{G} :

$$G_{ij}(\mathbf{A}, \mathbf{g}) = \begin{cases} 1, & \text{if } A_{ij} + \sum_u g_{ij}(u) > 0, \\ 0, & \text{otherwise.} \end{cases}$$

For all ego graphs we can compute the marginal,

$$\begin{aligned} P(\mathbf{g}|\mathbf{A}) &= \prod_u \int_0^1 P(\mathbf{g}(u)|\mathbf{A}, p) P(p) dp \\ &= \prod_u \left[\left(\frac{\sum_{i < j} m_{ij}(u)}{\sum_{i < j} g_{ij}(u)} \right)^{-1} \frac{1}{1 + \sum_{i < j} m_{ij}(u)} \right], \end{aligned}$$

which give us the posterior distribution:

$$P(\mathbf{g}, \mathbf{A}, \mathbf{b}|\mathbf{G}) = \frac{P(\mathbf{G}|\mathbf{A}, \mathbf{g}) P(\mathbf{g}|\mathbf{A}) P(\mathbf{A}|\mathbf{b}) P(\mathbf{b})}{P(\mathbf{g})}$$

Given a final graph \mathbf{G} , we can then recover the seminal network \mathbf{A} , its partition \mathbf{b} , as well as the ego graphs \mathbf{g} which contain the triadic closure edges.

ITERATED TRIADIC CLOSURES

An edge added to close an open triad may **create more open triads.**

We can then iterate triadic closures at successive generations l ,

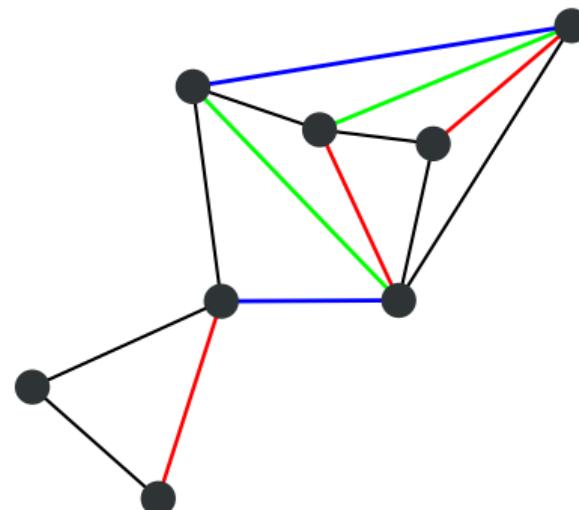
$$P(\mathbf{g}^{(l)} | \mathbf{A}^{(l-1)}, \mathbf{g}^{(l-1)}) = \prod_u \left[\left(\frac{\sum_{i < j} m_{ij}^{(l)}(u)}{\sum_{i < j} g_{ij}^{(l)}(u)} \right)^{-1} \frac{1}{1 + \sum_{i < j} m_{ij}^{(l)}(u)} \right].$$

This gives us a flexible model that can capture large triangle densities.

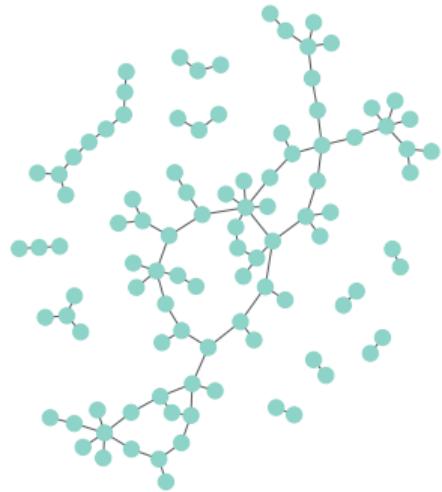
In this way, we have a joint posterior for all triadic closure generations

$$P(\{\mathbf{g}^{(l)}\}, \{\mathbf{A}^{(l)}\}, \mathbf{b} | \mathbf{g}) = \frac{P(\mathbf{g}, \{\mathbf{g}^{(l)}\}, \{\mathbf{A}^{(l)}\} | \mathbf{b}) P(\mathbf{b})}{P(\mathbf{g})}.$$

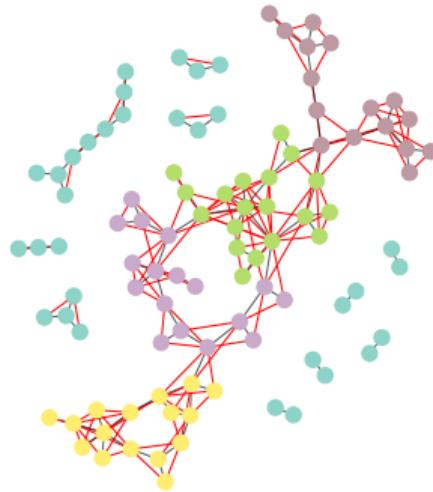
To facilitate identifiability, a new generation can only close a triad that involves at least one edge of the previous generation:



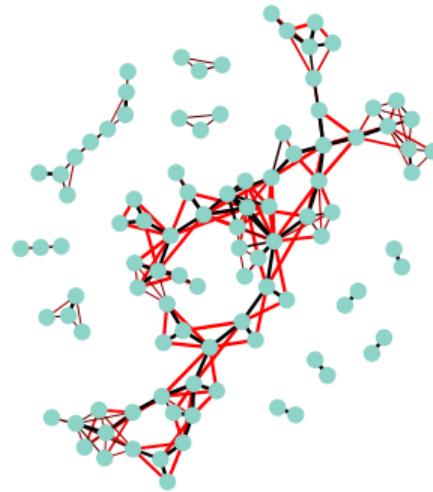
DISTINGUISHING TRIADIC CLOSURE FROM COMMUNITY STRUCTURE



(a) Random seminal edges



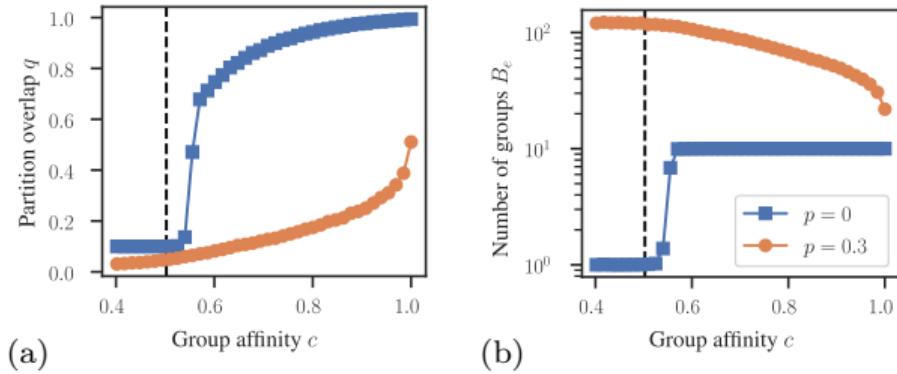
(b) Triadic closure edges and
spurious communities found
with SBM ($\Sigma_{\text{SBM}} = 801.7$ nats)



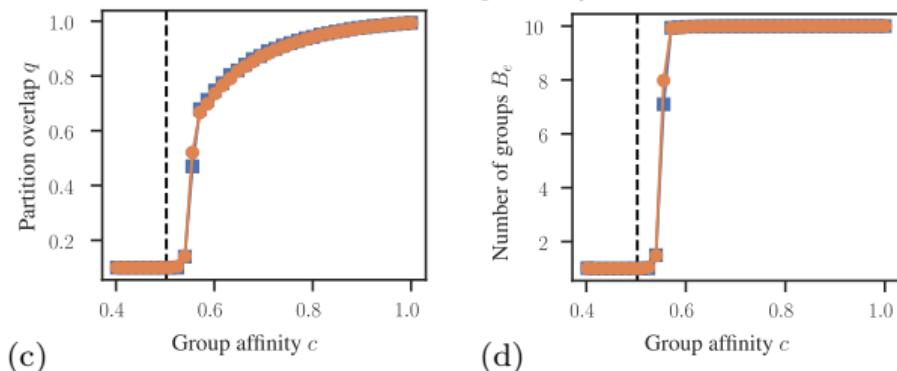
(c) Inference of the SBM/TC
model ($\Sigma_{\text{SBM/TC}} = 590.6$ nats)

ARTIFICIAL NETWORKS

Inference using pure SBM



Inference using SBM/TC



Planted Partition model with triadic closure

The seminal graph \mathbf{A} is generated with edge counts

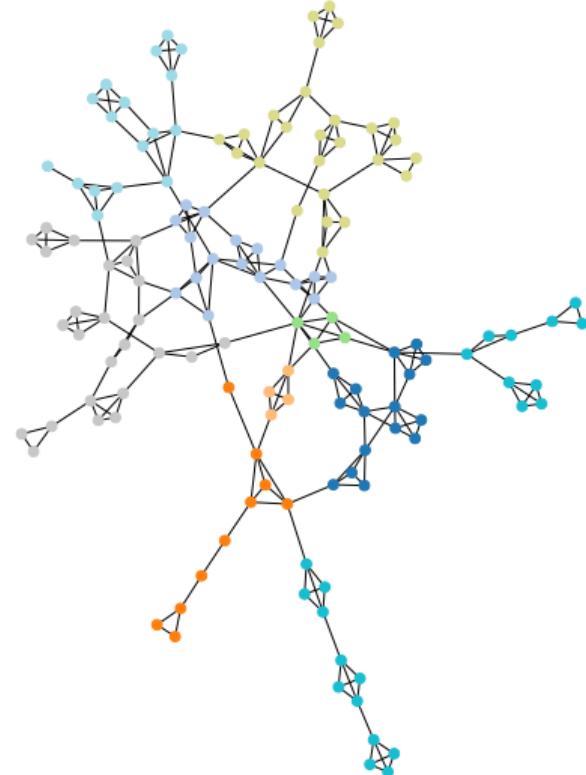
$$e_{rs} = 2E \left[\frac{c}{B} \delta_{rs} + \frac{1-c}{B(B-1)} (1 - \delta_{rs}) \right],$$

where $c \in [0, 1]$ is the group affinity.

The final graph \mathbf{G} is obtained with one round of triadic closures on \mathbf{A} with $p_u = p$.

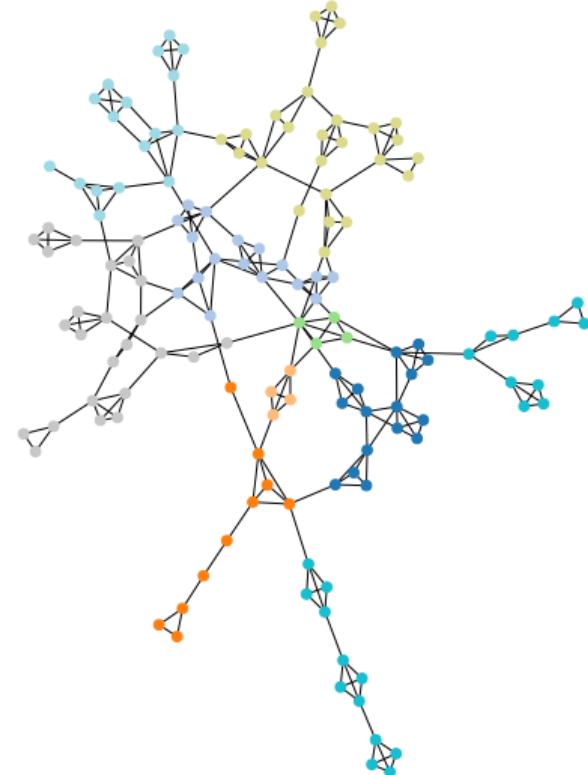
Left: $N = 10^4$, $\langle k \rangle = 5$, $B = 10$

EMPIRICAL NETWORK: STUDENT COOPERATIONS

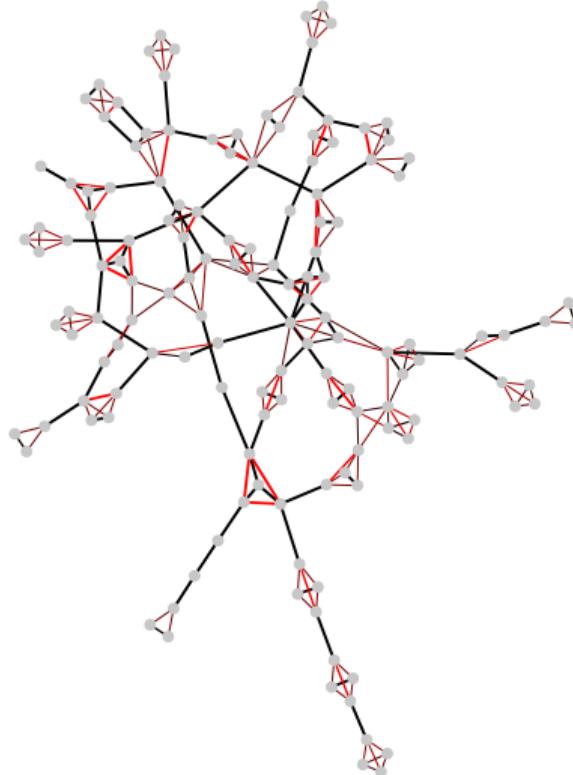


(a) SBM, $B = 9$, $\Sigma_{\text{SBM}} = 1145.6$ nats

EMPIRICAL NETWORK: STUDENT COOPERATIONS

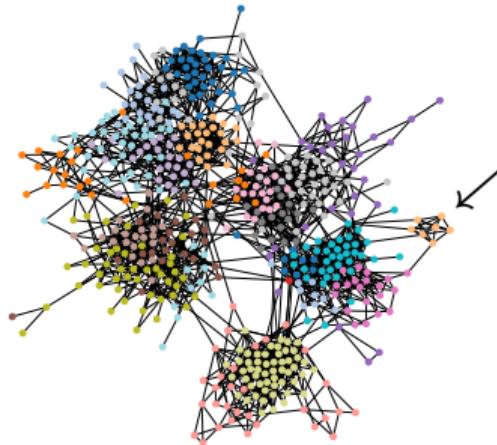


(a) SBM, $B = 9$, $\Sigma_{\text{SBM}} = 1145.6$ nats



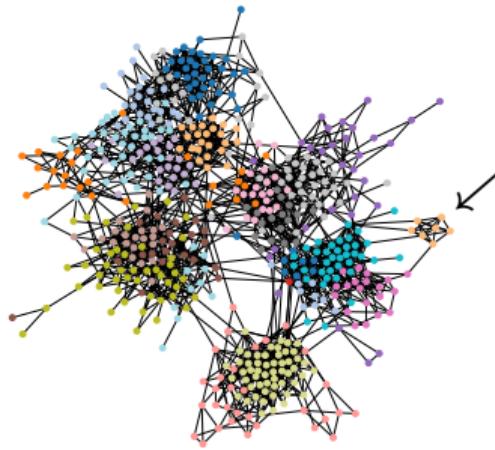
(b) SBM/TC, $\Sigma_{\text{SBM/TC}} = 935.1$ nats

HIGH-SCHOOL STUDENT FRIENDSHIPS

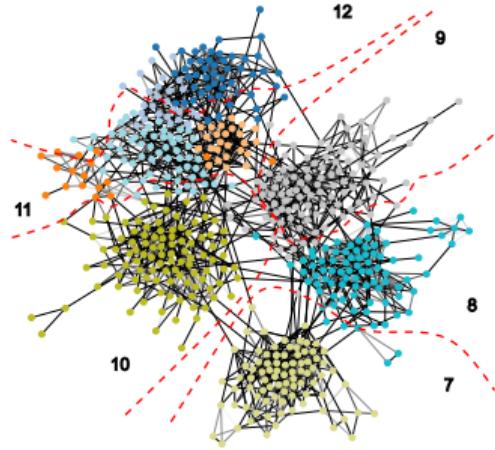


(a) SBM, $B = 26$, $\Sigma_{\text{SBM}} = 8757.0$
nats

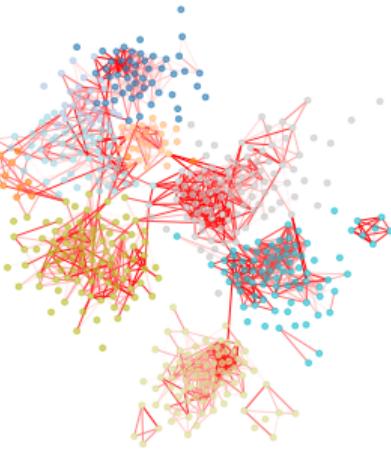
HIGH-SCHOOL STUDENT FRIENDSHIPS



(a) SBM, $B = 26$, $\Sigma_{\text{SBM}} = 8757.0$ nats



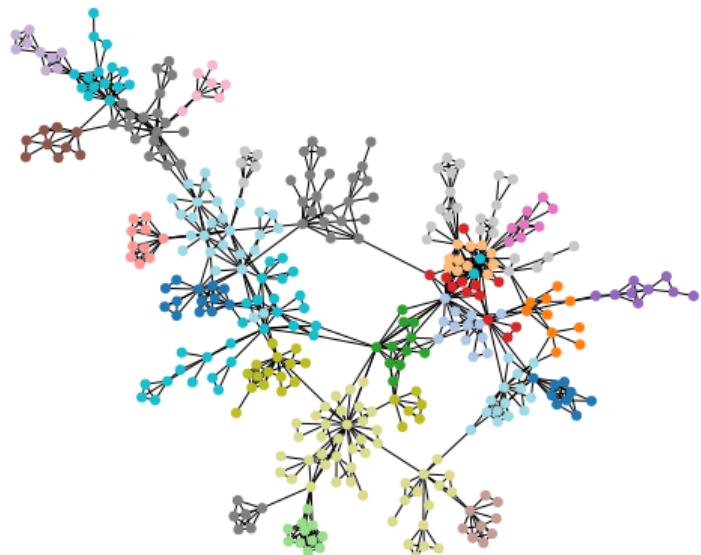
Seminal edges



Triadic closure edges

(b) SBM/TC, $B = 9$, $\Sigma_{\text{SBM/TC}} = 8456.3$ nats

COLLABORATIONS BETWEEN NETWORK SCIENTISTS

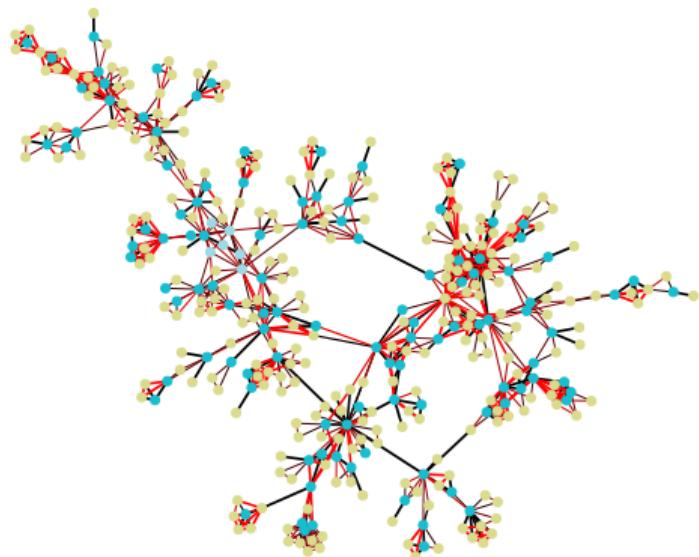


(c) SBM, $B = 27$, $\Sigma_{\text{SBM}} = 3816.3$ nats

COLLABORATIONS BETWEEN NETWORK SCIENTISTS

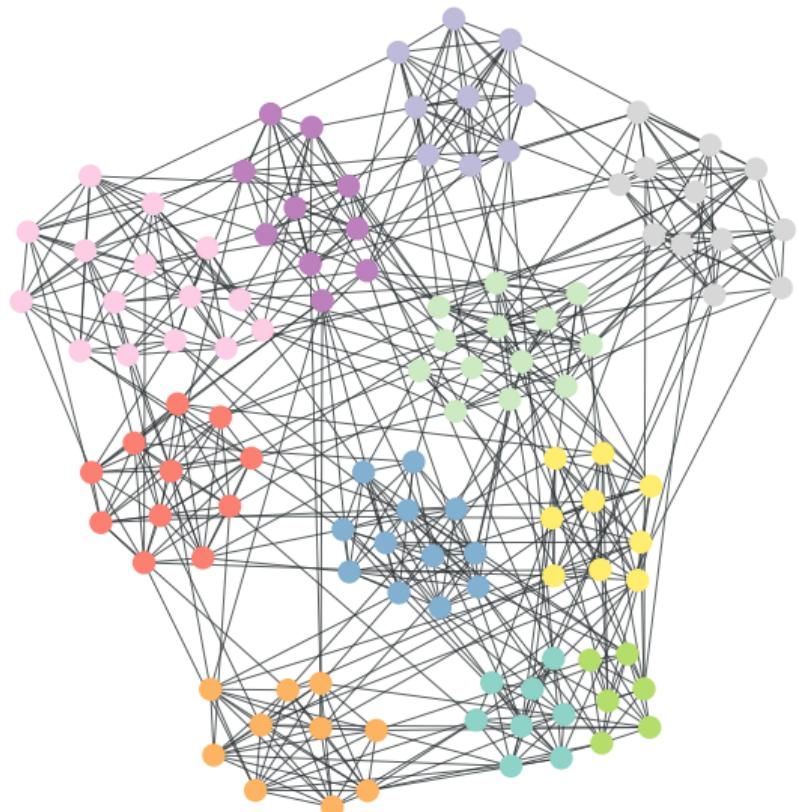


(c) SBM, $B = 27$, $\Sigma_{\text{SBM}} = 3816.3$ nats



(d) SBM/TC, $B = 3$, $\Sigma_{\text{SBM/TC}} = 3009.9$ nats

GAMES BETWEEN AMERICAN FOOTBALL TEAMS



$$B = 11$$

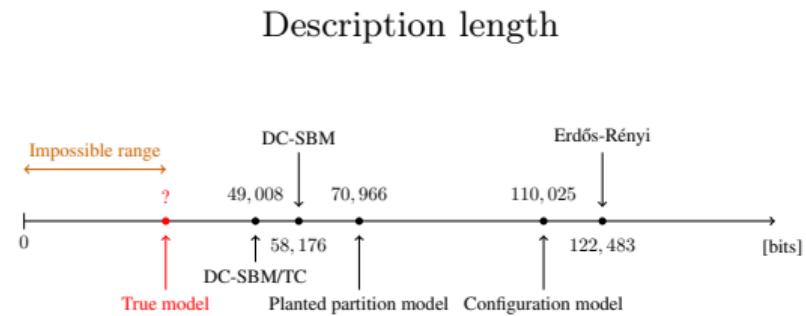
$$\Sigma_{\text{SBM}} = 1761.1 \text{ nats}$$

$$\Sigma_{\text{SBM/TC}} = 1767.6 \text{ nats}$$

$$C = 0.41$$

MODEL SELECTION AND FINDING THE “TRUTH”

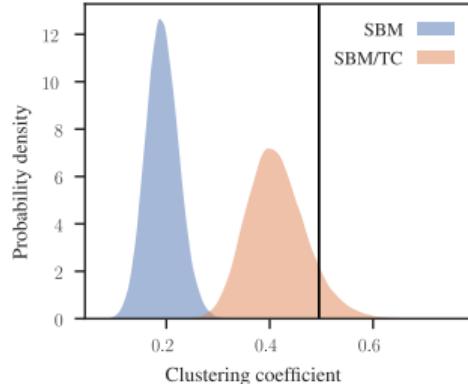
Protein-Protein interaction (wild turkey)



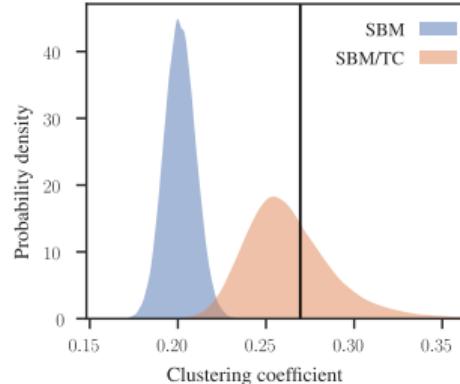
Compression points *toward* what is detectable of the true generative process.

SBM/TC YIELDS PLAUSIBLE CLUSTERING VALUES

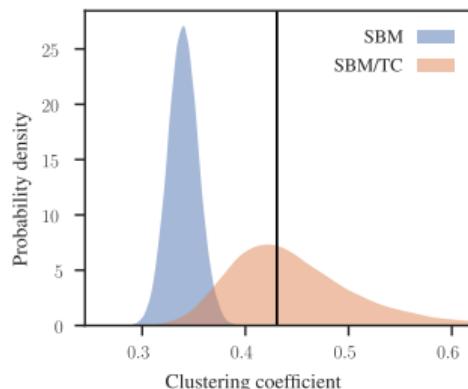
Posterior predictive checks: we sample networks from the fitted model and compare to the data.



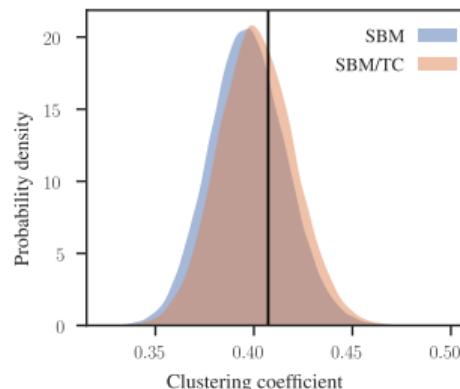
(a) Cooperation between students



(b) Adolescent health (comm26)

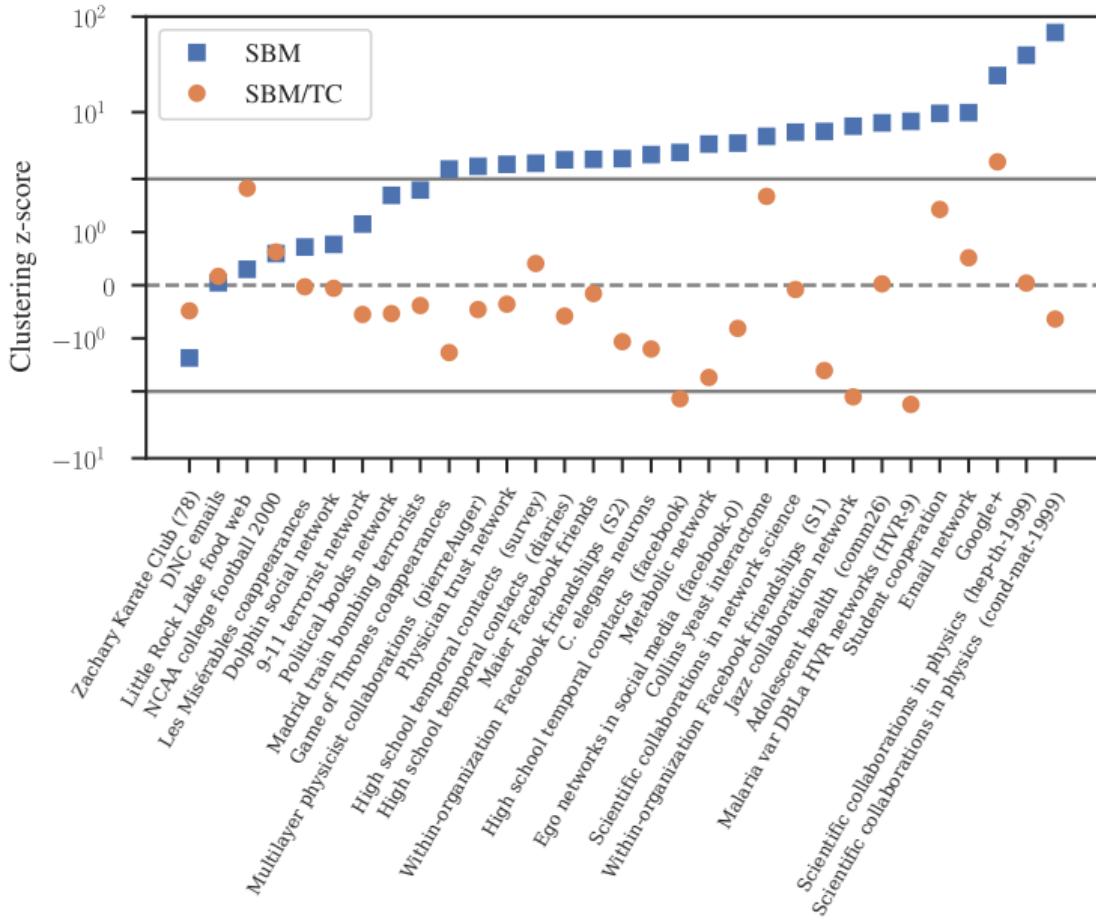


(c) Scientific collaborations in Network Science

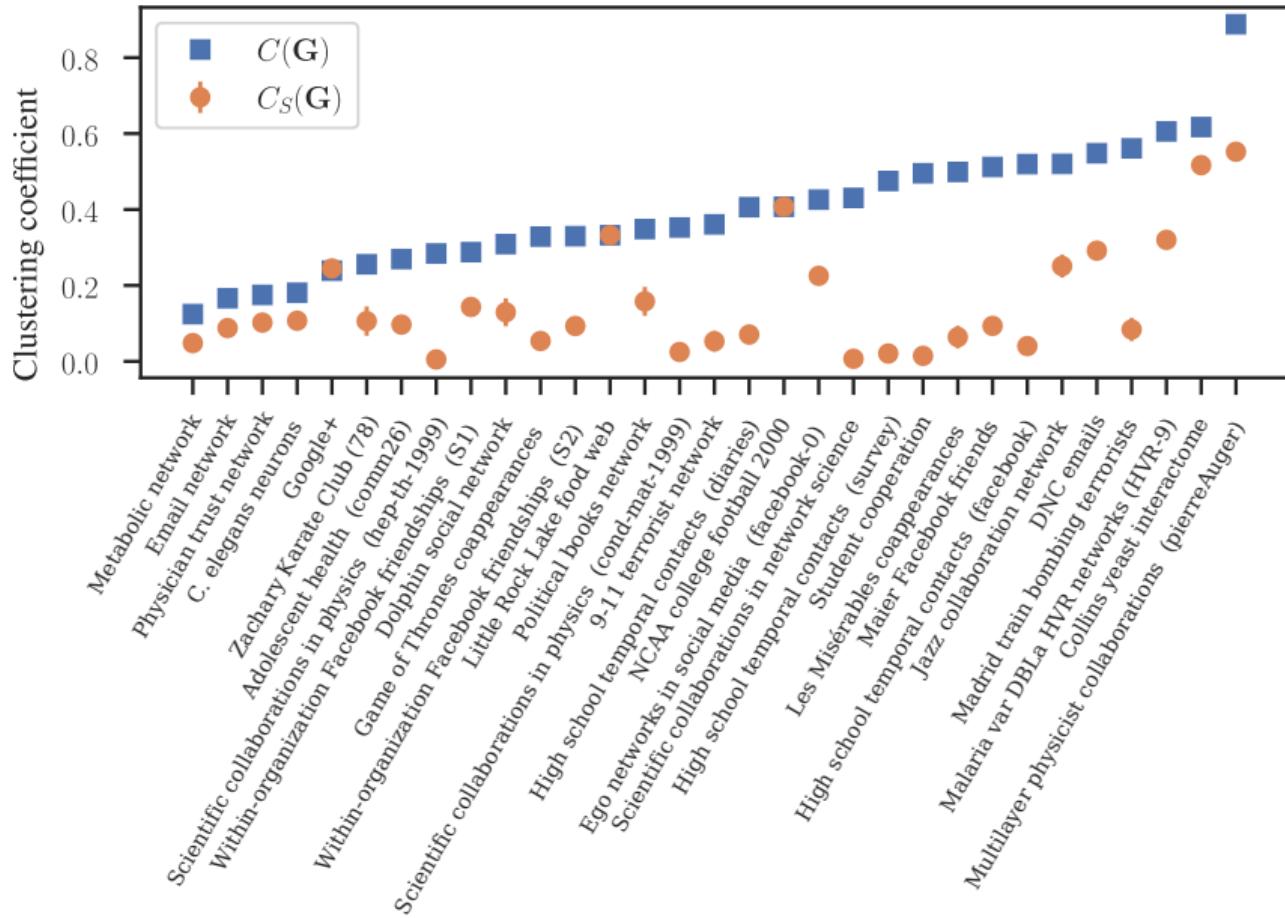


(d) NCAA college football 2000

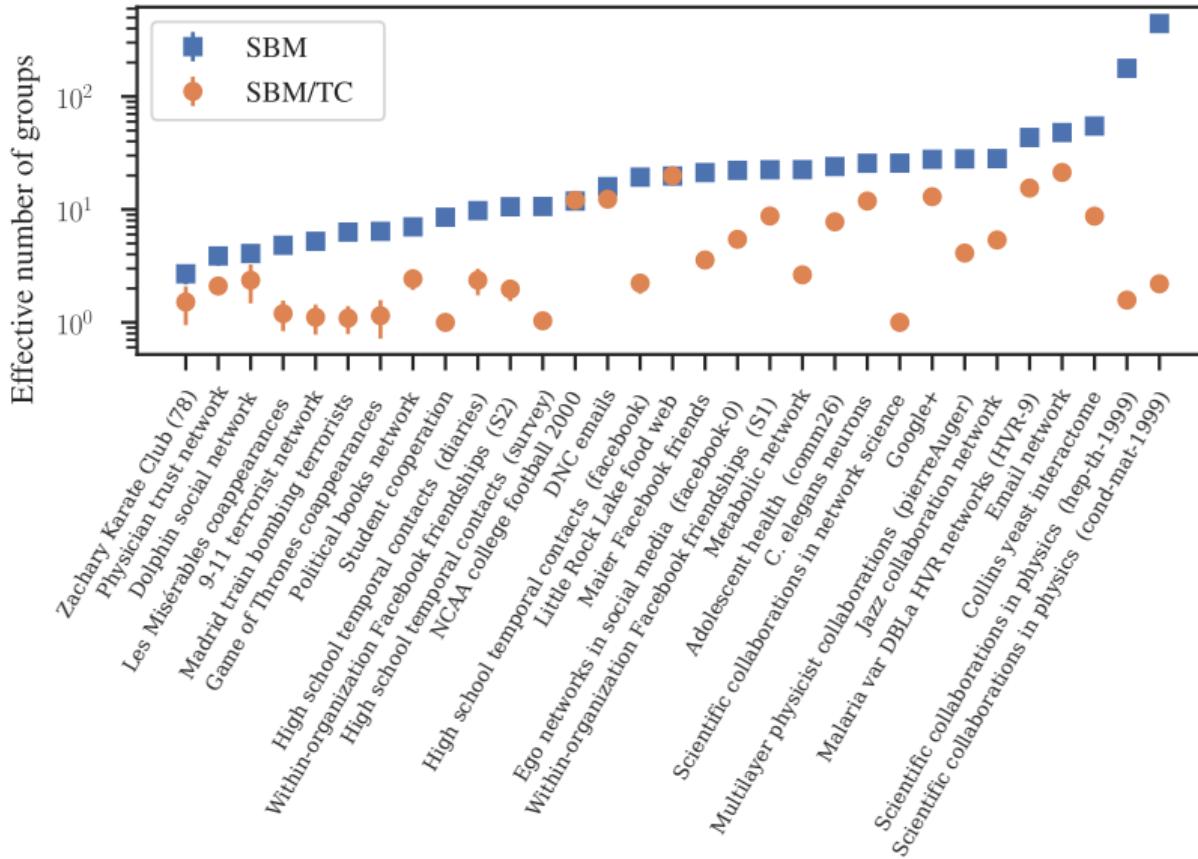
CAPTURING THE CLUSTERING COEFFICIENT



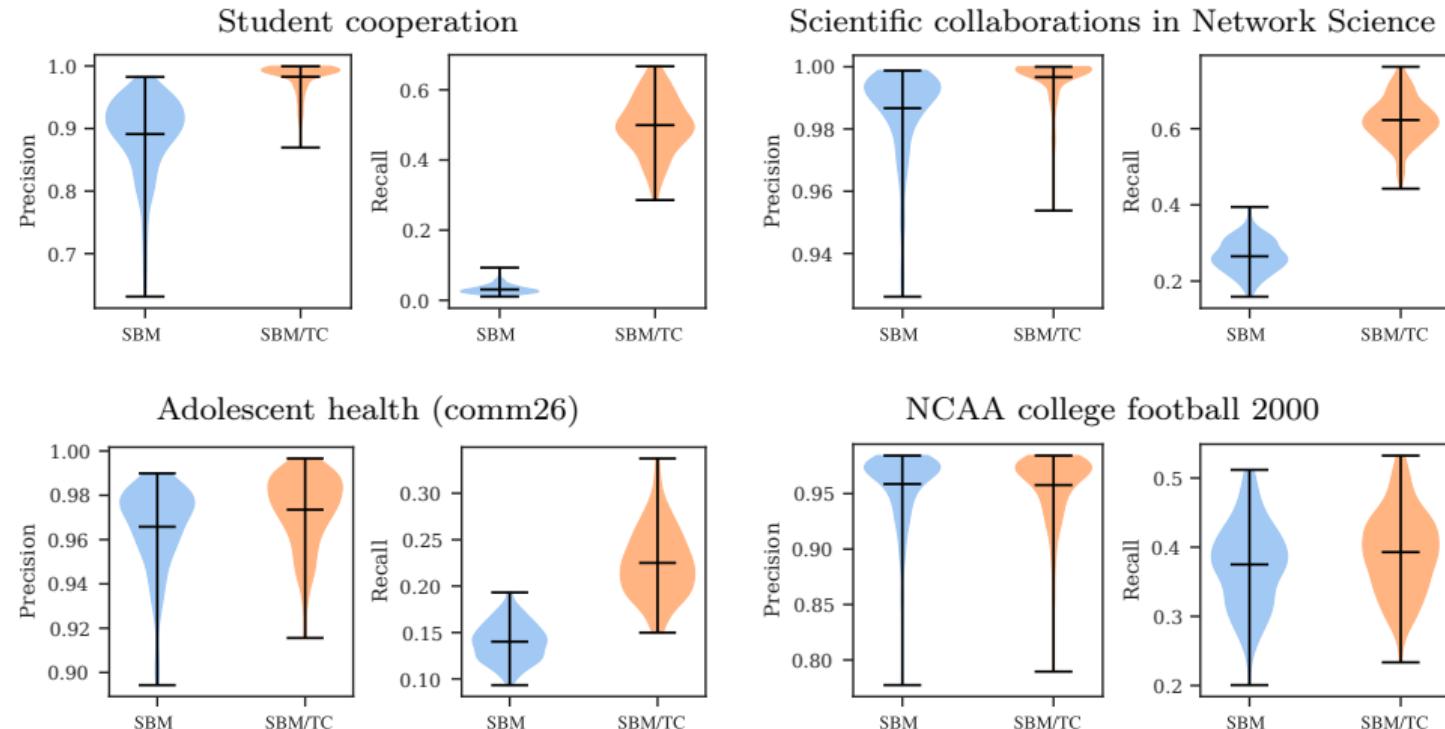
DECOMPOSING THE CLUSTERING COEFFICIENT



COMMUNITY STRUCTURE VS. TRIADIC CLOSURE



SBM/TC IS MORE PREDICTIVE OF MISSING EDGES



SUMMARY

The SBM + triadic closure model allows us to

- ▶ Separate the effect of homophily from triadic closure
- ▶ Identify triadic closure edges
- ▶ Identify communities that *cannot* be explained by triadic closure
- ▶ Improve compression
- ▶ Improve link prediction

Paper: T.P.P., Phys. Rev. X 12, 011004 (2022)

Code: <https://graph-tool.skewed.de>