

A network diagram with several nodes (circles) connected by lines (edges). One node is highlighted with a blue and white circular logo. The background is dark gray.

HIERARCHICALLY MODULAR STRUCTURE IN COMPLEX NETWORKS

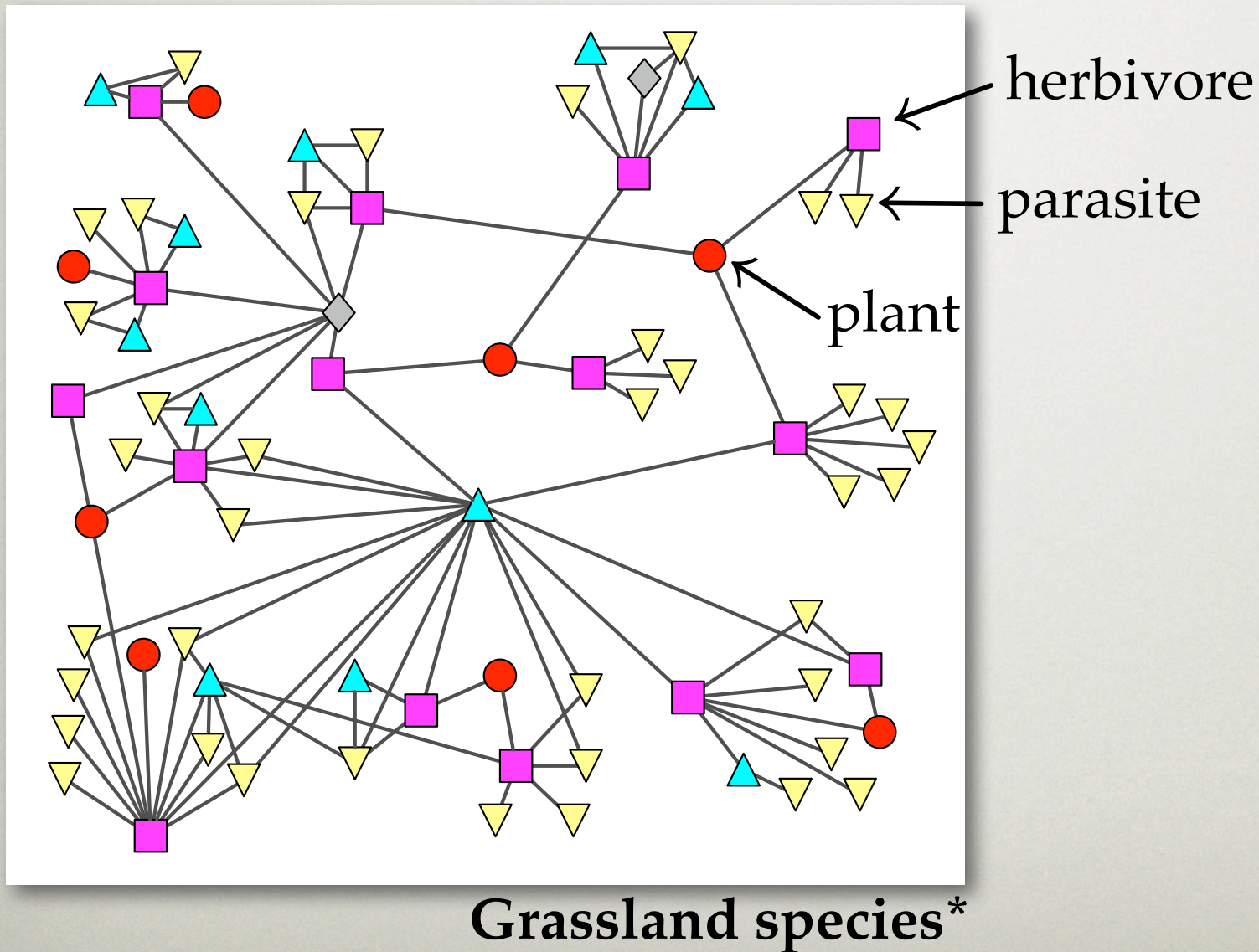
Aaron Clauset
Santa Fe Institute

3 December 2008
SFI Workshop

“Statistical Inference for Complex Networks”

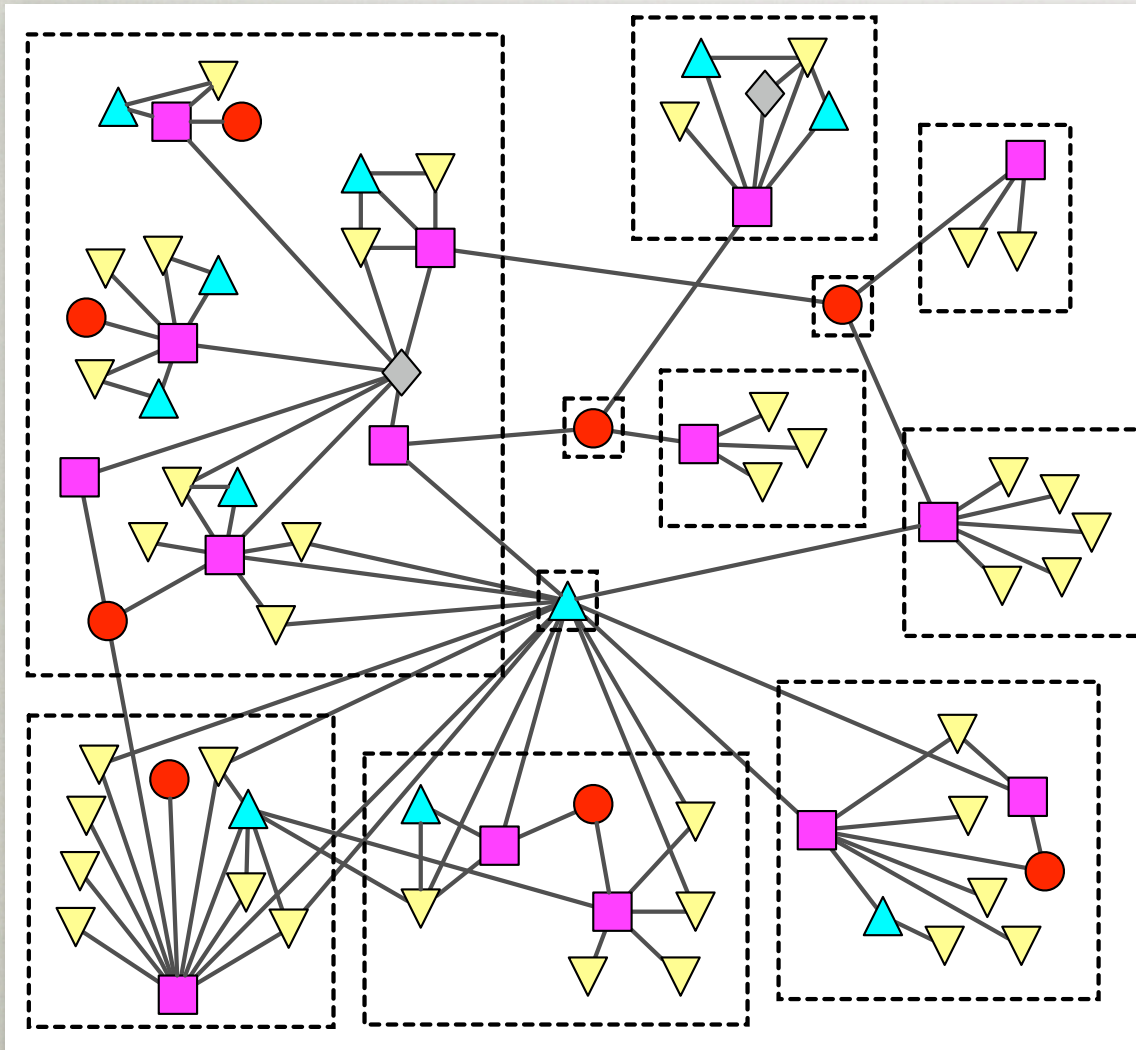
joint with C. Moore, M.E.J. Newman, T.A.S. Pierce

HIERARCHIES OF MODULES



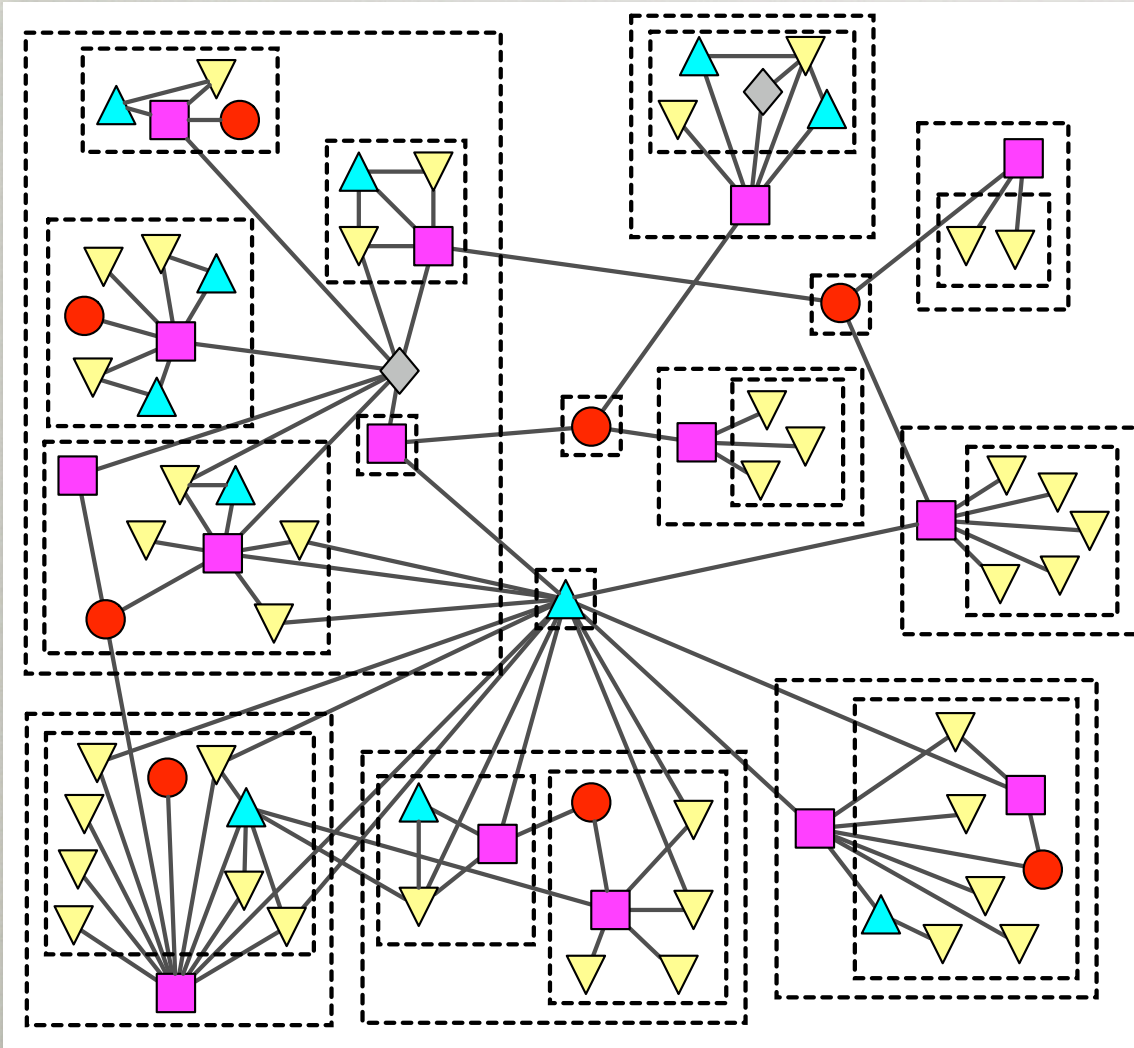
*thank you: Jennifer Dunne

HIERARCHIES OF MODULES



one level

HIERARCHIES OF MODULES



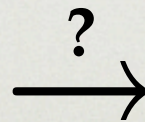
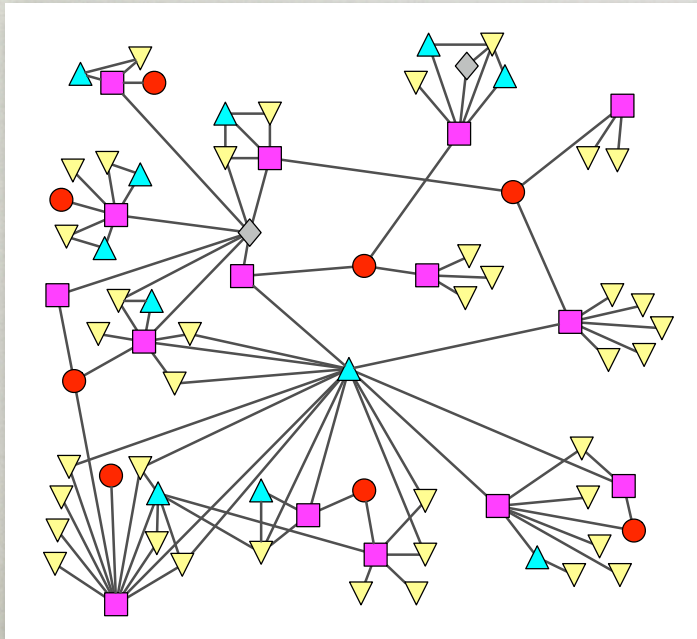
many levels

THE TASK

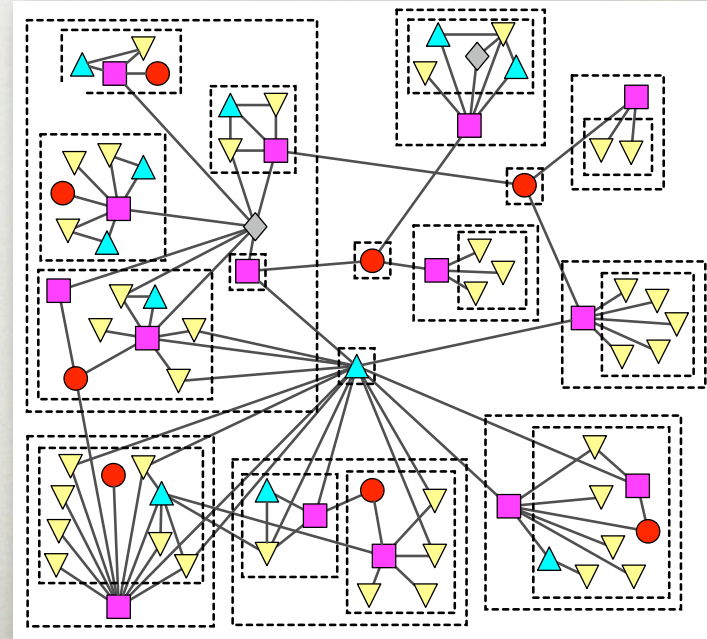
How can we extract

- **this hierarchical (multi-scale) structure** from complex networks?

network data



hierarchical structure



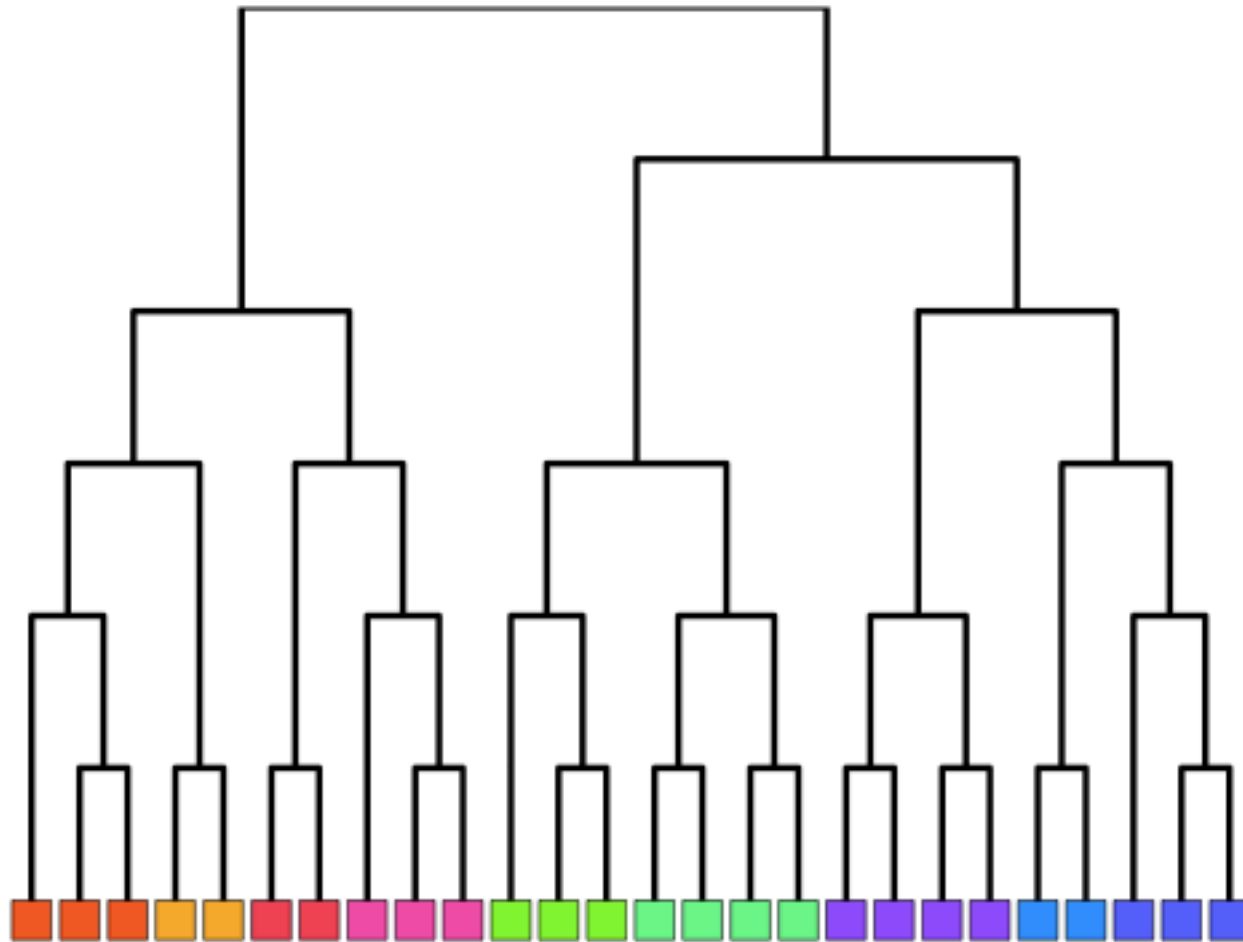
ONE APPROACH

Model-based inference

1. generative model of hierarchies (a model)
2. find “good” instances of this model
3. predict missing information

A MODEL OF HIERARCHY

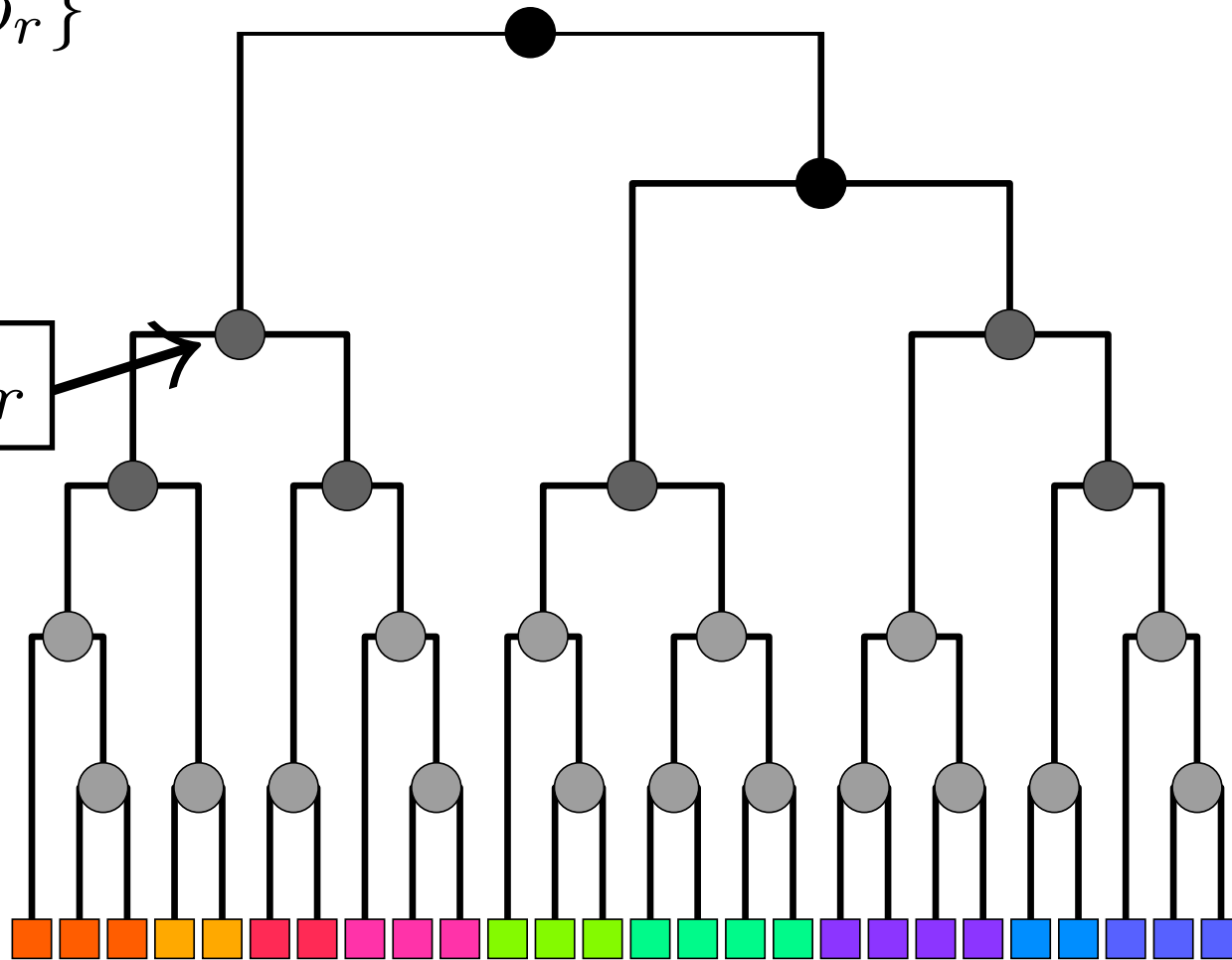
\mathcal{D}



A MODEL OF HIERARCHY

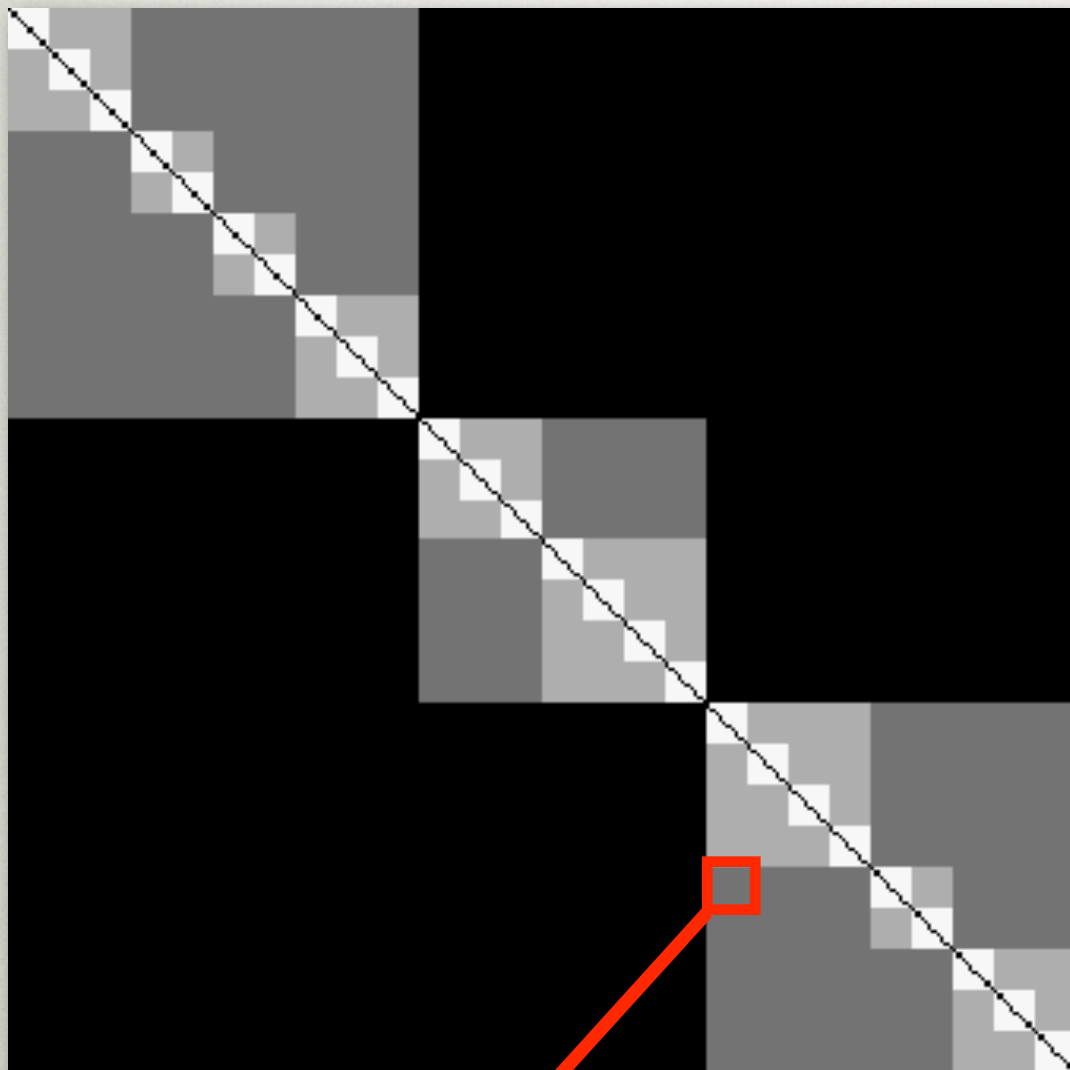
$\mathcal{D}, \{p_r\}$

probability p_r



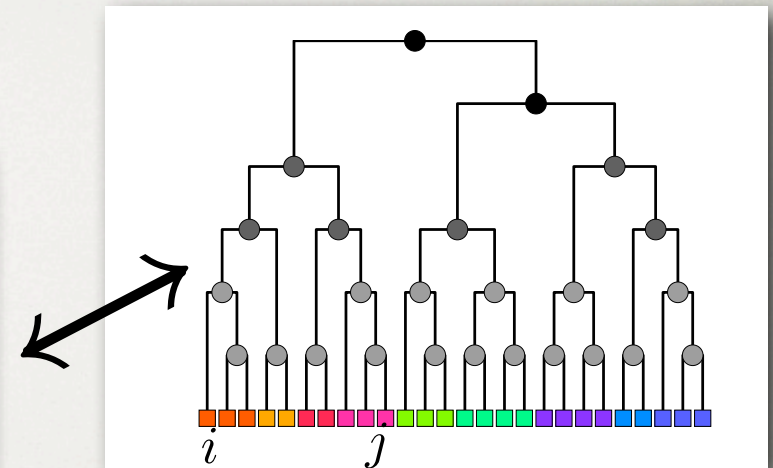
assortative modules

“inhomogeneous” random graph

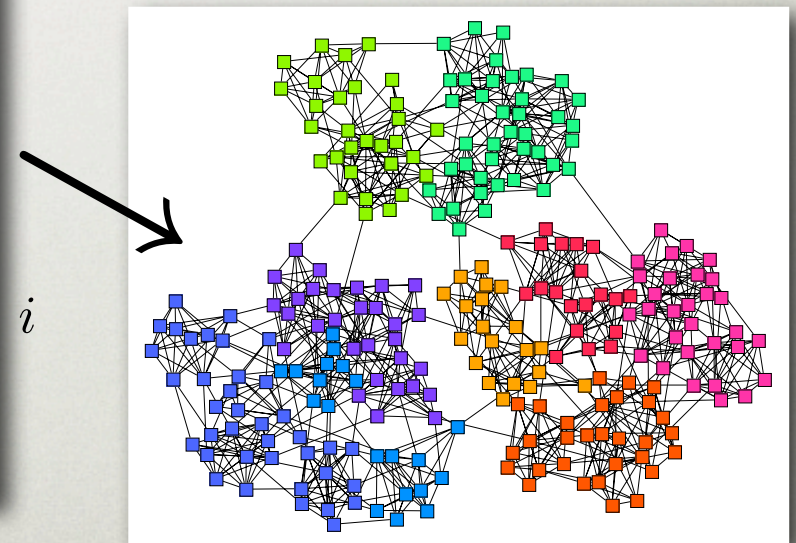


$$\begin{aligned} \Pr(i, j \text{ connected}) &= p_r \\ &= p_{(\text{lowest common ancestor of } i, j)} \end{aligned}$$

model



instance



MODEL FEATURES

- explicit model = explicit assumptions
- very flexible (many parameters)
- captures structure at all scales
- arbitrary mixtures of assortativity, disassortativity
- learnable directly from data

LEARNING FROM DATA

a direct approach

- **likelihood function** $\mathcal{L} = \text{Pr}(\text{ data } | \text{ model })$
(\mathcal{L} scores **quality** of model)
- **sample** the **good** models
via Markov chain Monte Carlo
- technical details in
Nature **453** (2008) and *physics/0610051*

MISSING LINKS

A test: can model predict missing links?

GUESSING IS BAD

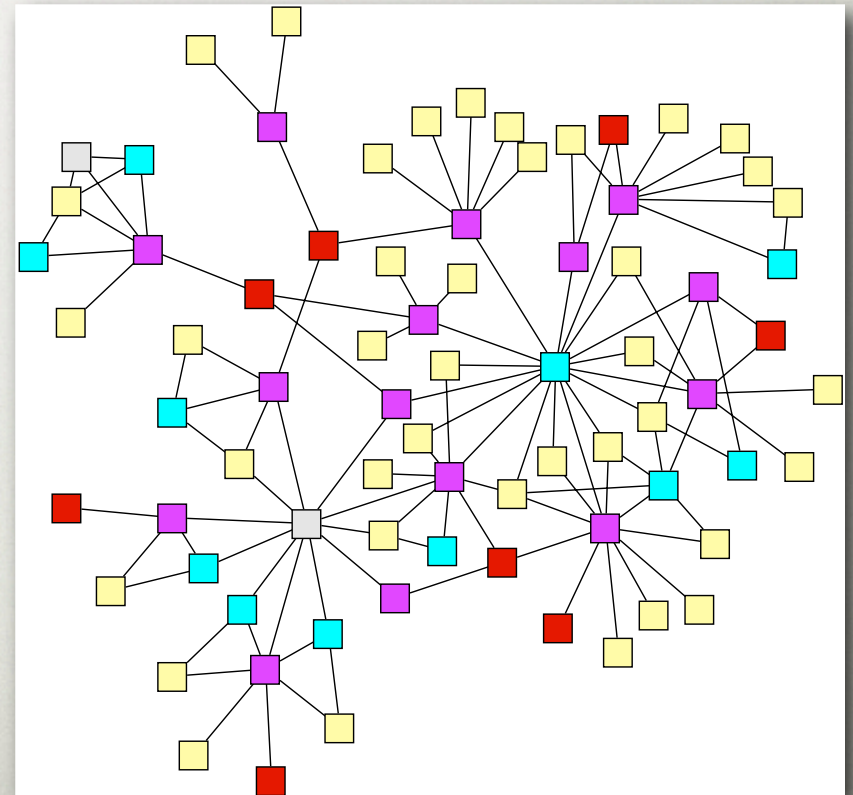
- remove k edges from G
- how easy to guess a missing link?

$$p_{\text{guess}} \approx \frac{k}{n^2 - m + k}$$
$$= O(n^{-2})$$

$$n = 75$$

$$m = 113$$

$$p_{\text{guess}} = k / (2662 + k)$$



PREDICTING MISSING LINKS

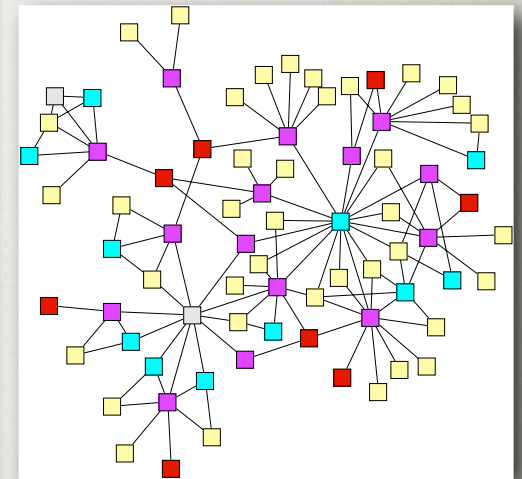
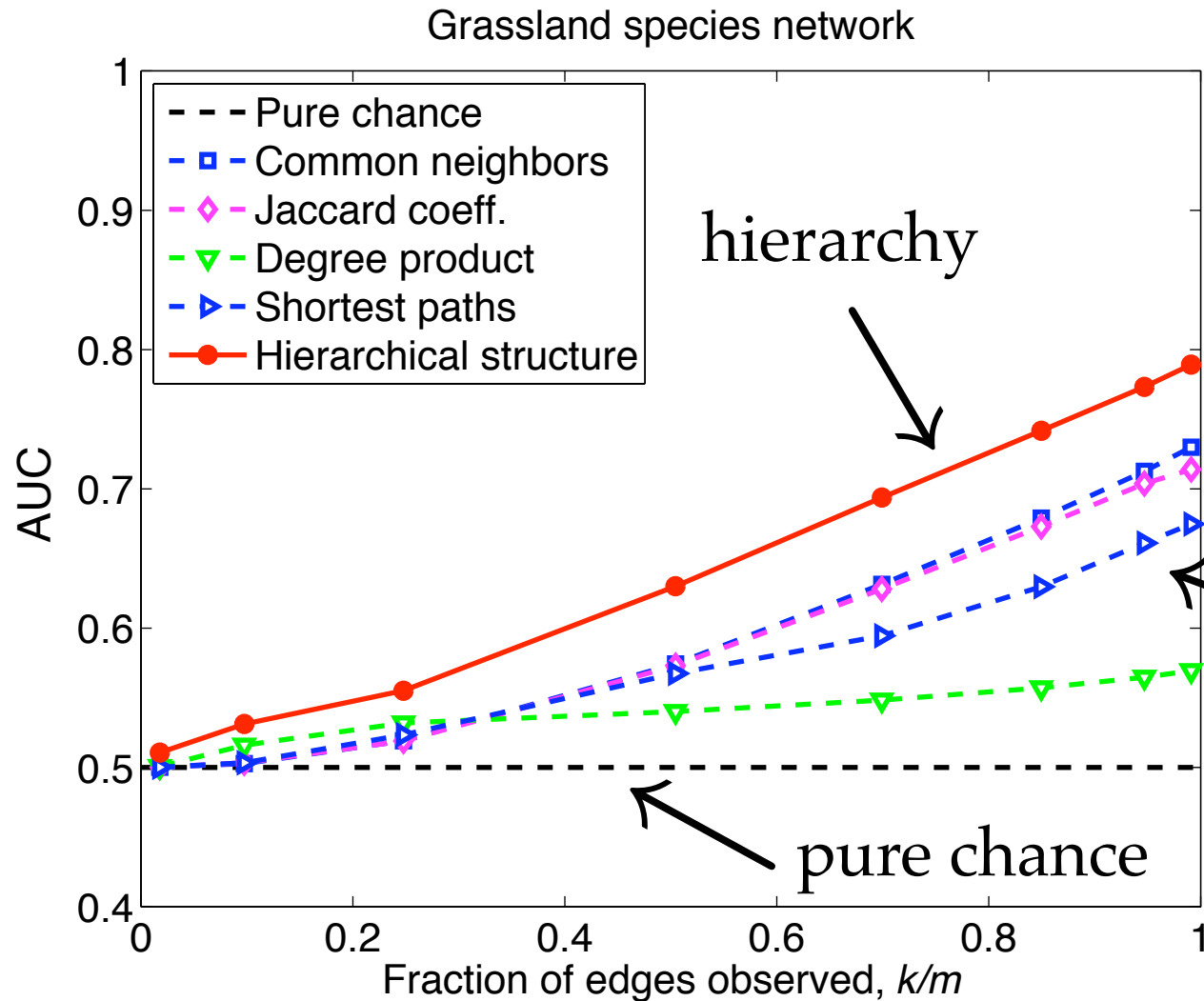
- Given incomplete graph G
- run MCMC to equilibrium
- then, over sampled \mathcal{D} , compute average $\langle p_r \rangle$ for links $(i, j) \notin G$
- predict links with high $\langle p_r \rangle$ values are missing

Test idea via leave- k -out cross-validation

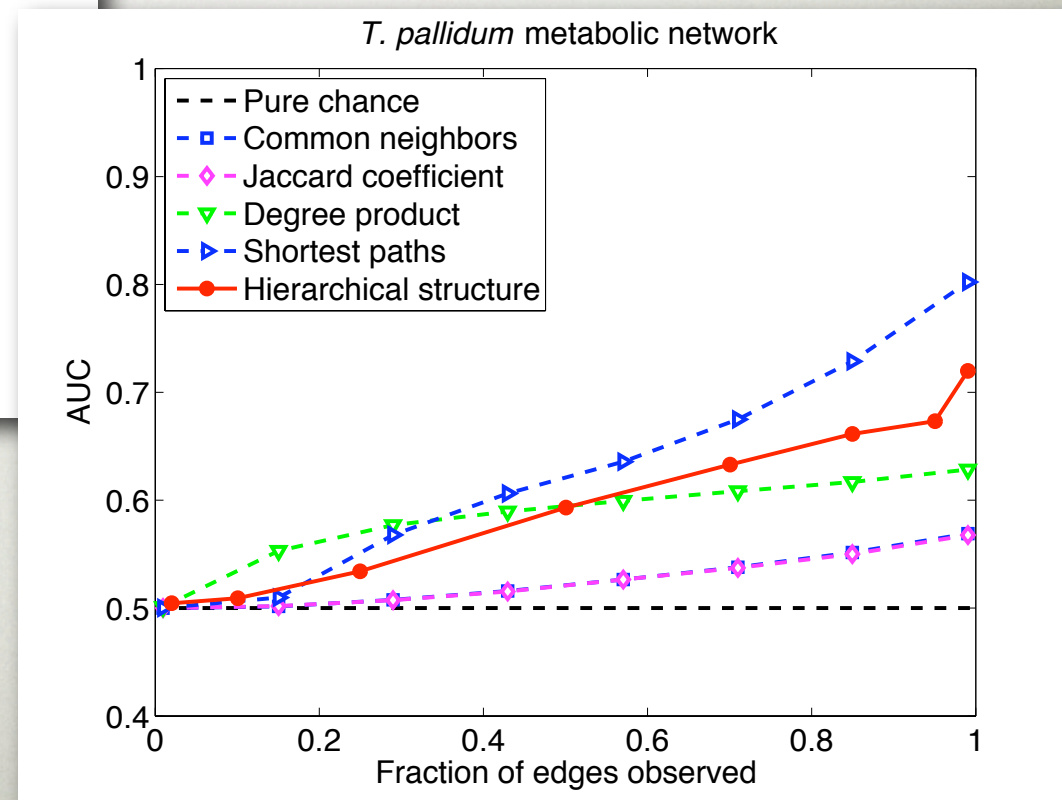
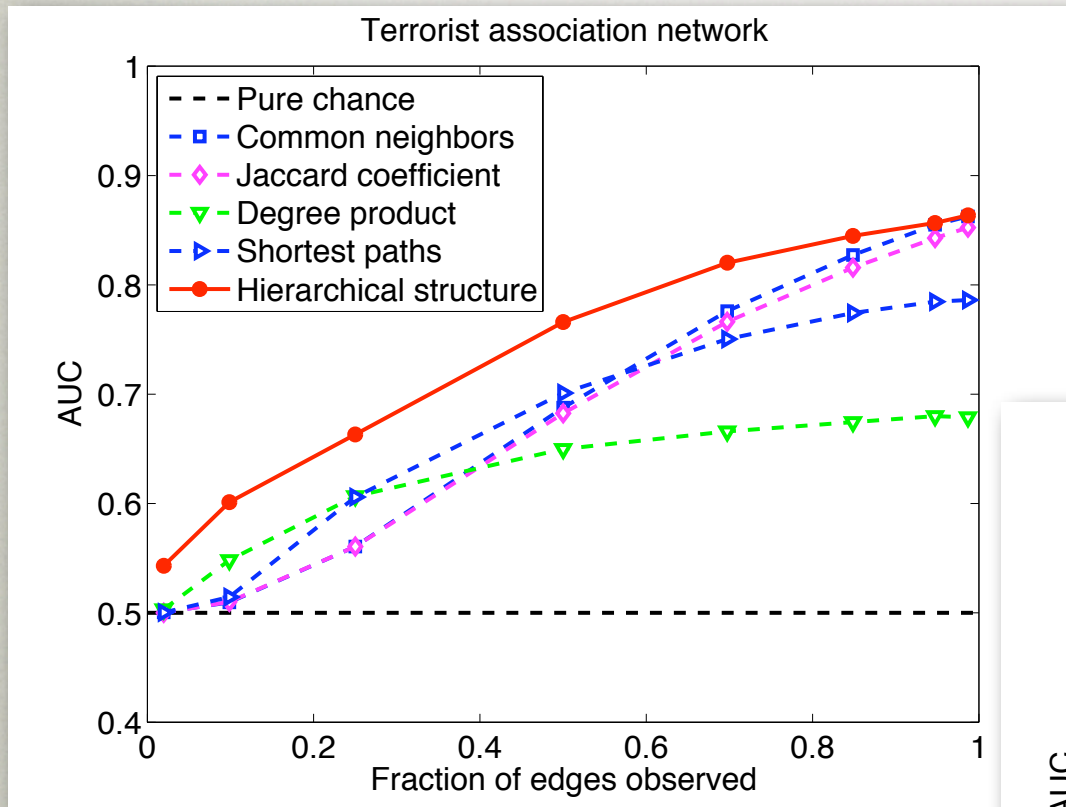
perfect accuracy: $\text{AUC} = 1$

no better than chance: $\text{AUC} = 1/2$

MISSING STRUCTURE



OTHER NETWORKS



OTHER DATA

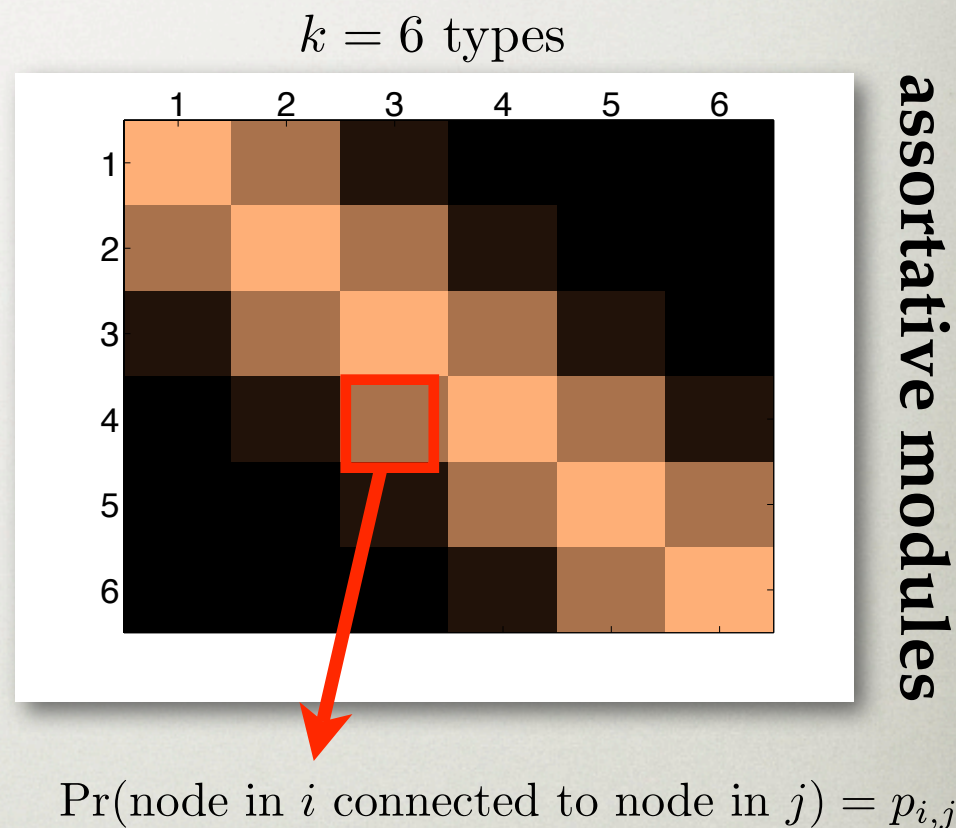
- node labels, attributes, weights
- edge labels, attributes, weights
- geographical structure
- temporal information
- combinations of these

Node labels only:

a simple (“stochastic block”) model can do well

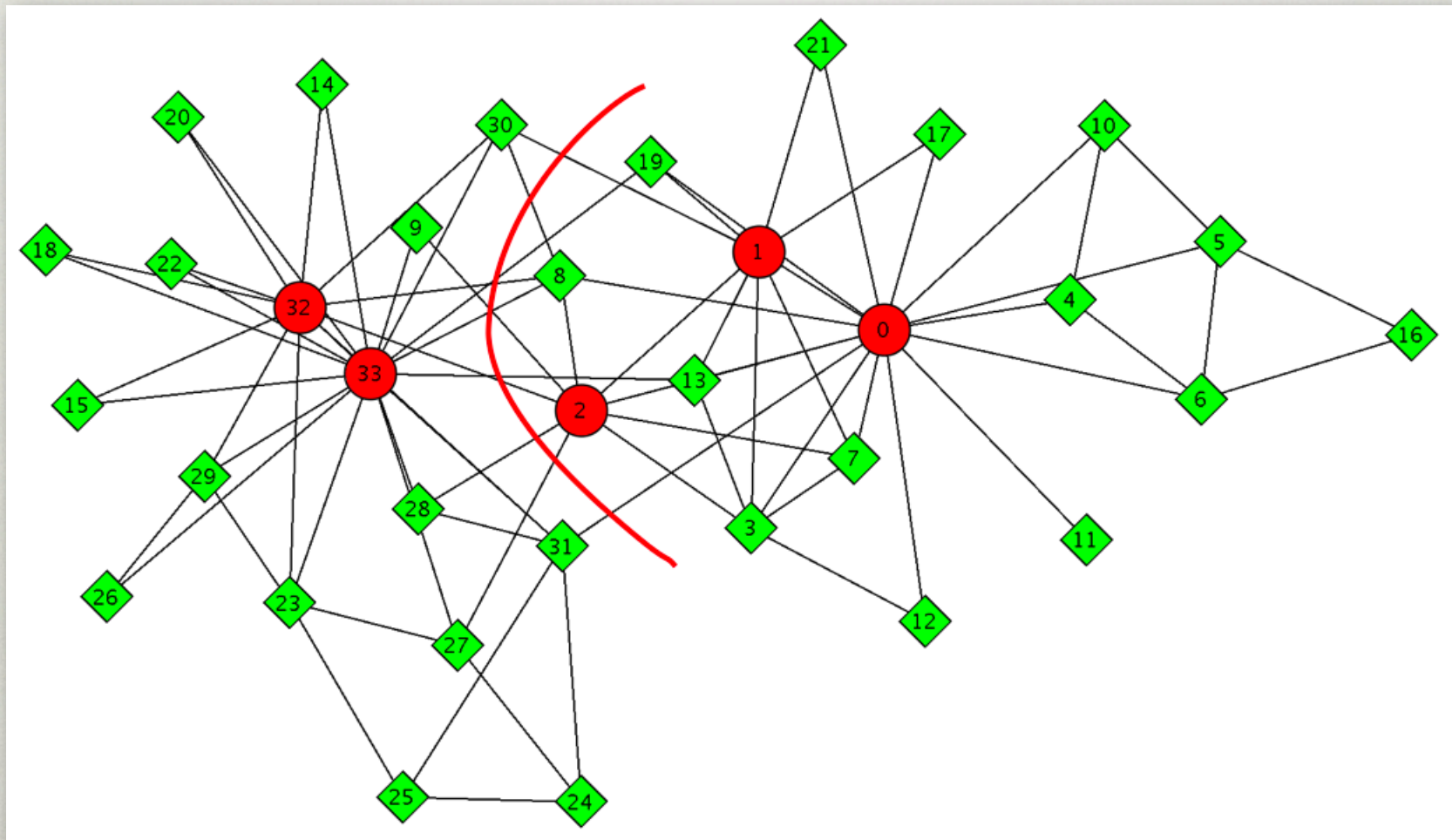
STOCHASTIC BLOCK MODEL

- k types of nodes
- $k \times k$ matrix $p_{i,j}$ of module connectivity
- no assumptions on structure of $p_{i,j}$
- we use MCMC to search over all node labelings
- can we predict *missing* labels?



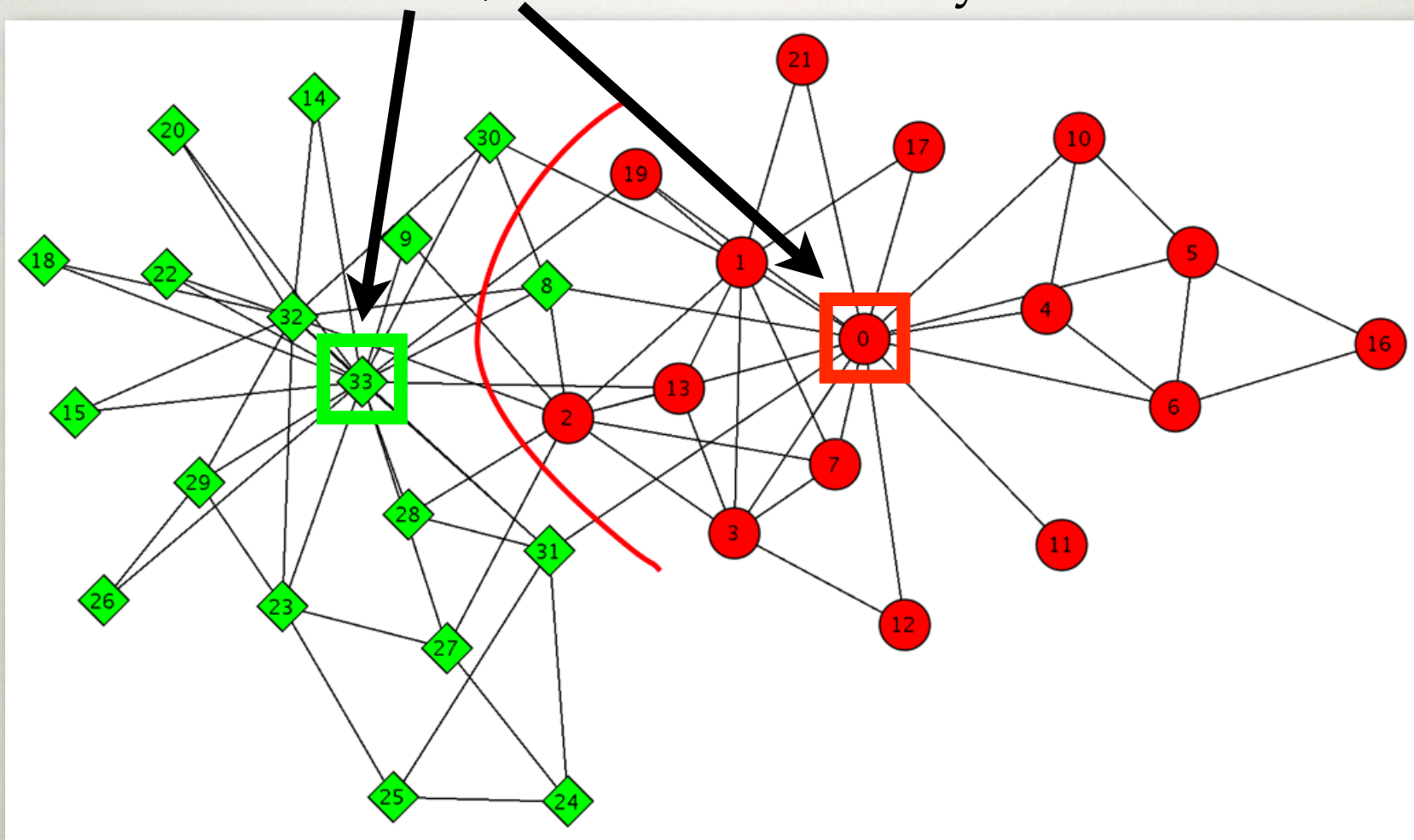
UNCONSTRAINED INFERENCE

let all labels vary



CONSTRAINED INFERENCE

fix these labels; let the others vary



SUMMARY

- generative models of
 - hierarchical modules
 - simple modules
 - potentially many others
- predict missing information (edges, types)
- principled way to fit and test structural theories

Acknowledgments:

C. Moore, M.E.J. Newman, T.A.S. Pierce,
C.H. Wiggins, C.R. Shalizi

FIN
