



HIERARCHICAL STRUCTURE AND PREDICTING MISSING LINKS

Cristopher Moore
University of New Mexico
and the Santa Fe Institute

joint work with
Aaron Clauset (UNM/SFI)
and Mark Newman (Michigan)

THREE GOALS

#1: Inferring network structure from
observed data

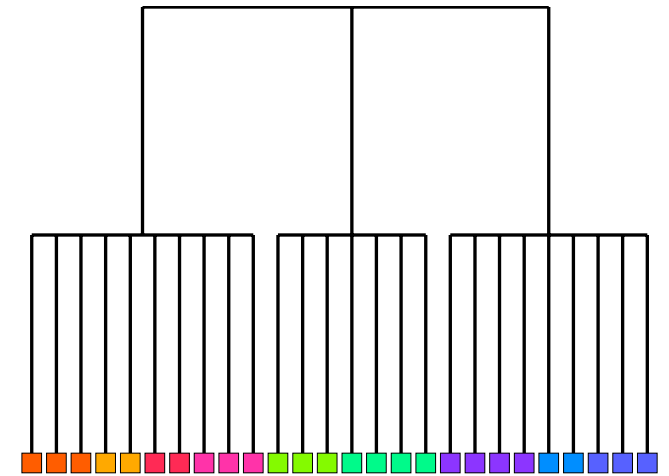
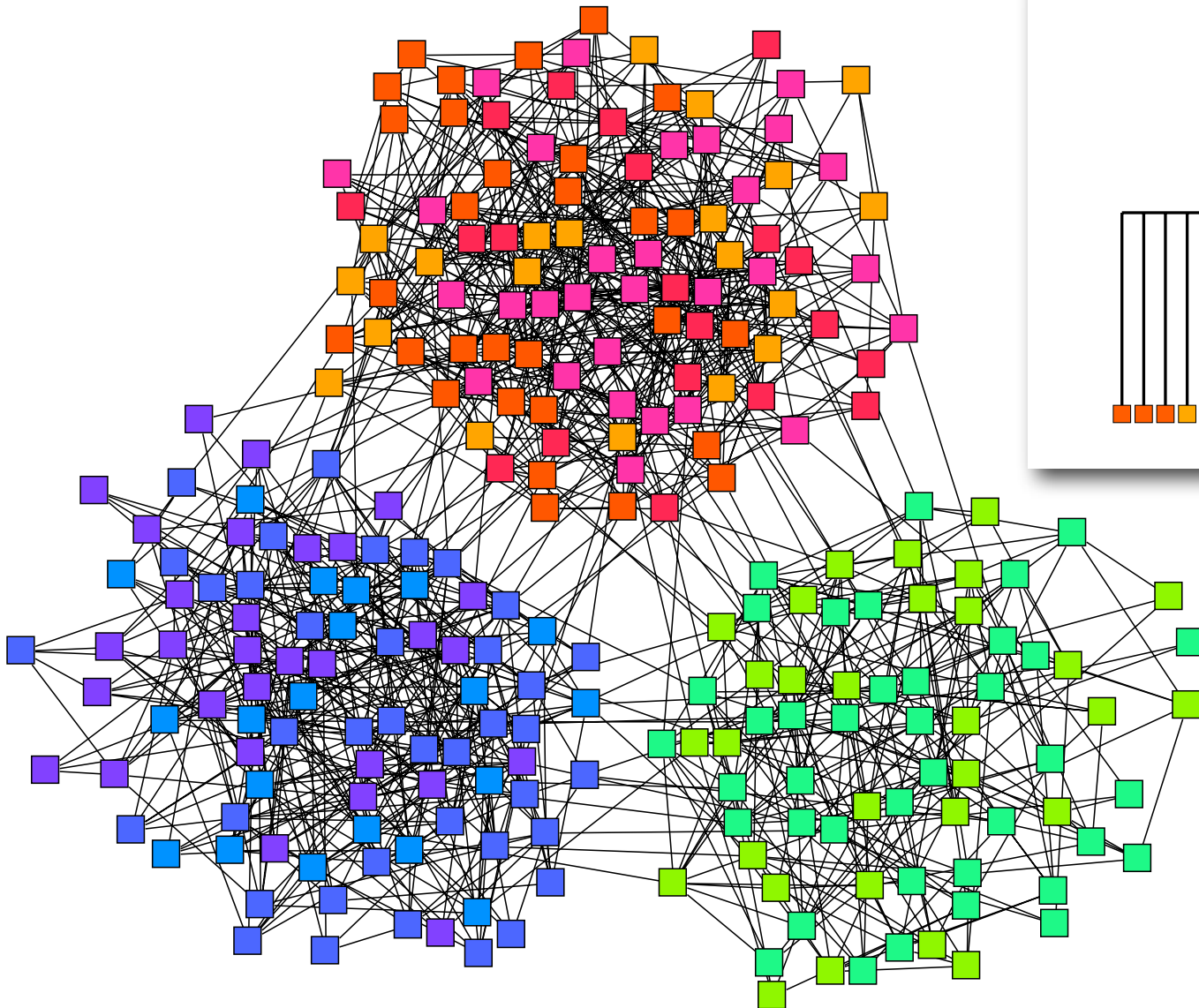
#2: Generating random graphs which are
statistically similar to real ones

#3: Predicting missing links

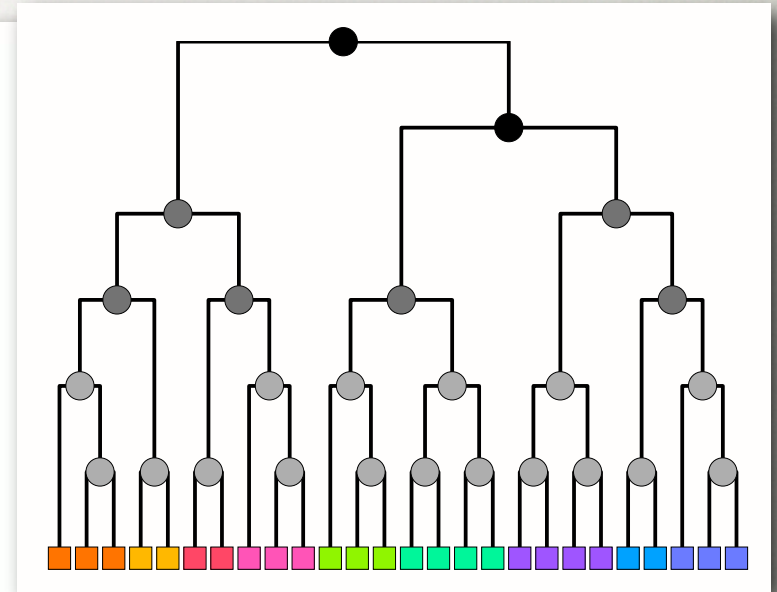
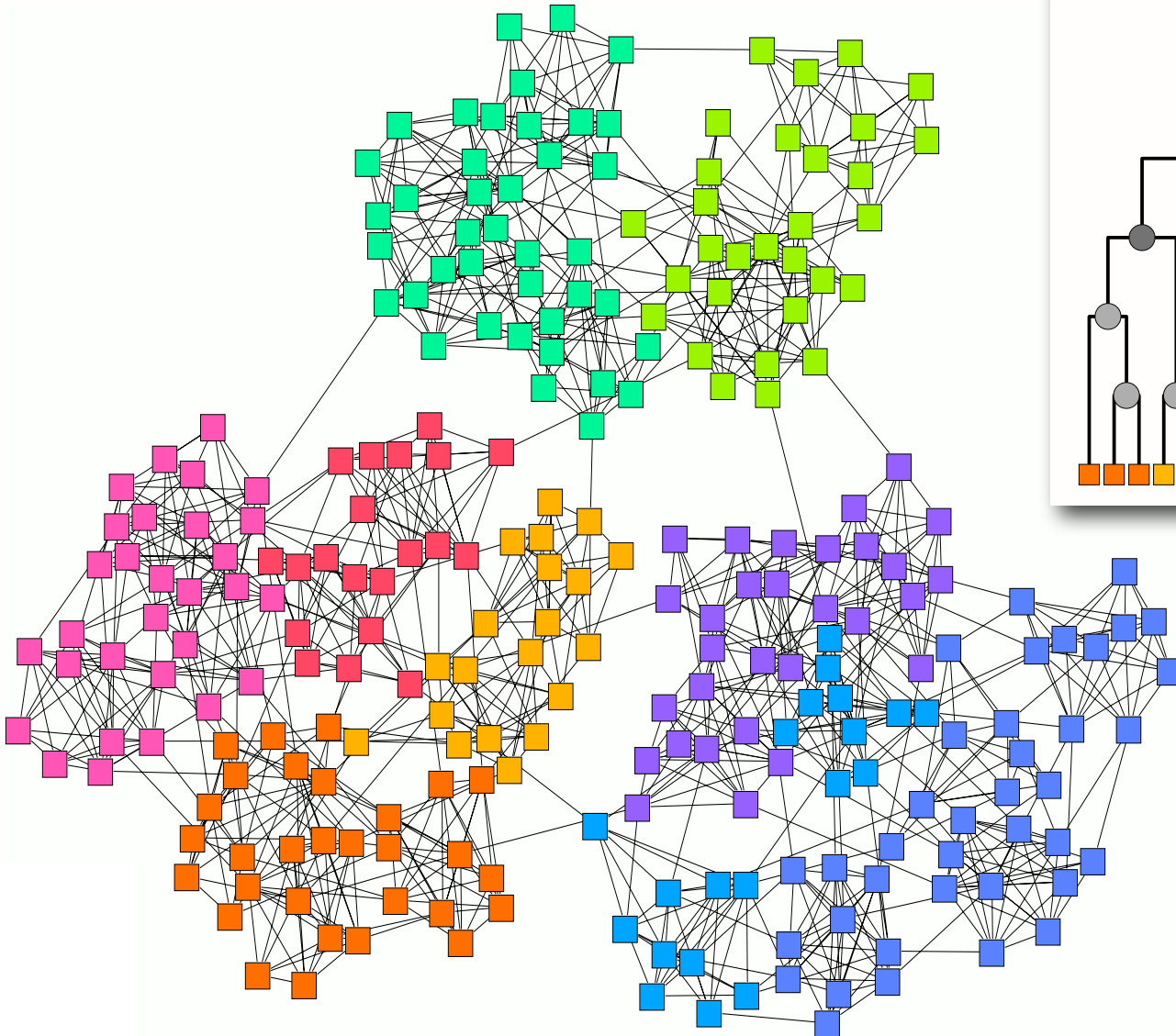
TWO EXTREMES

- Proving theorems about simple models
 - *Captures some aspect of network structure,*
 - *But leaves lots of other structure out.*
- Observing patterns in real networks
 - *Good “natural history”*
 - *But often ad-hoc.*
- One possible bridge: Machine Learning.

CLUSTERING: ONE LEVEL



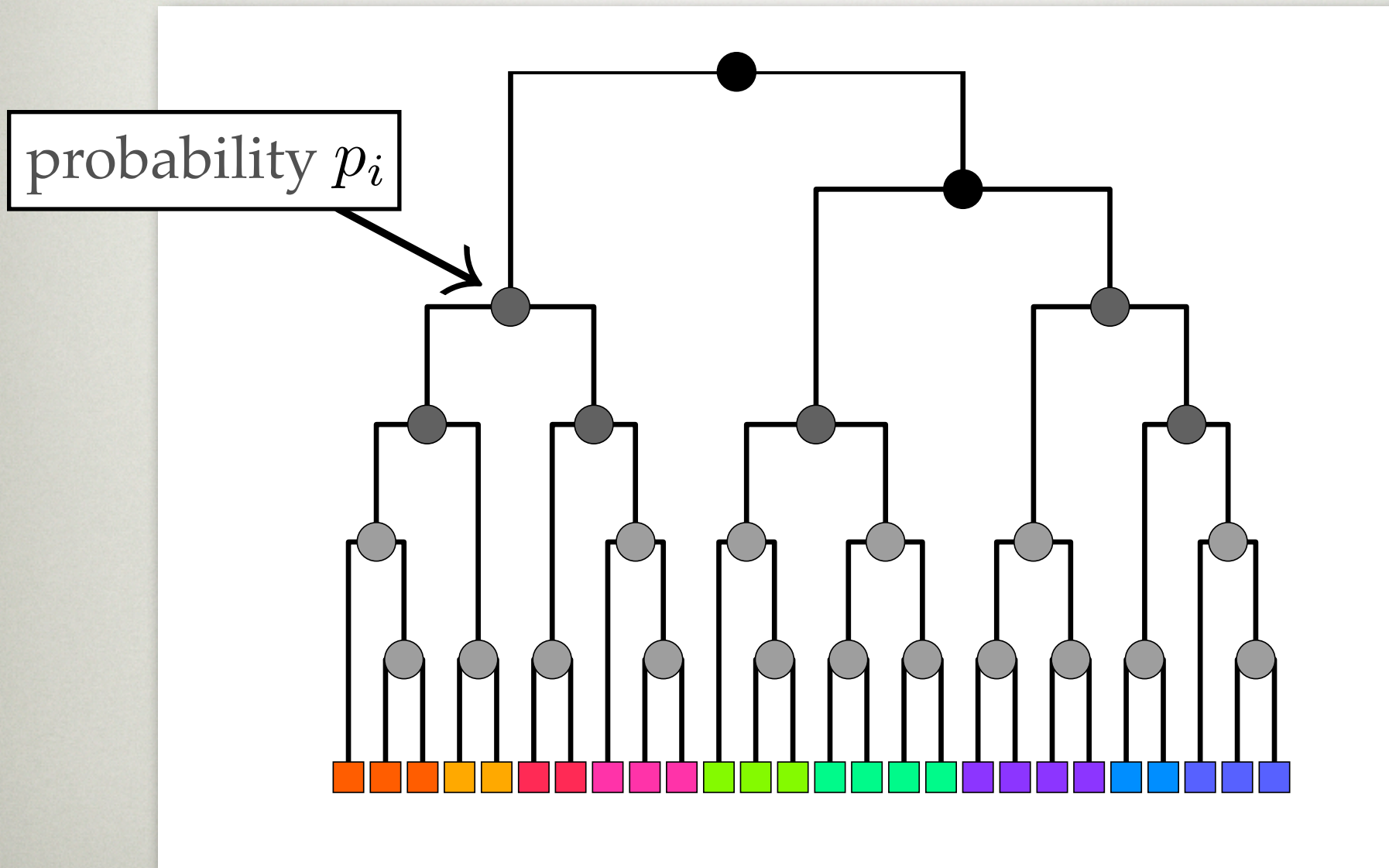
HIERARCHY: MANY LEVELS



A PROBABILISTIC MODEL (WITH LOTS OF PARAMETERS!)

- A binary tree T : leaves are original vertices, internal nodes represent communities
- Each internal node has a probability p_i
- Two vertices are connected with probability p_i where i is their lowest common ancestor
- Allows *assortative* or *disassortative* structure, or any mix across scales and subtrees

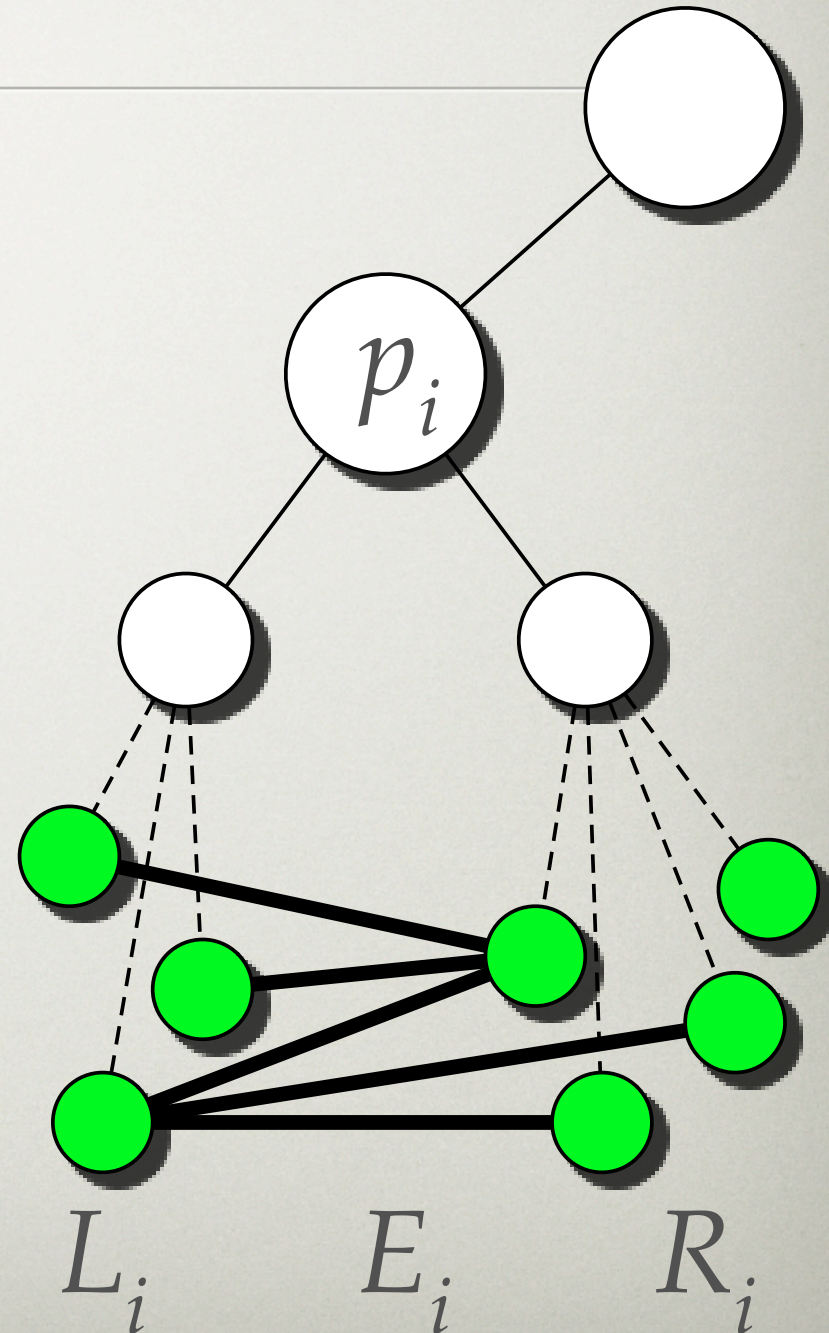
A PROBABILISTIC MODEL (WITH LOTS OF PARAMETERS!)



MAXIMUM LIKELIHOOD

- For each internal node i ,
 - L_i and $R_i = \#$ descendants
 - $E_i = \#$ edges between them
- Likelihood these edges exist, and not others, is

$$\mathcal{L}_i = p_i^{E_i} (1 - p_i)^{L_i R_i - E_i}$$



MAXIMUM LIKELIHOOD

- Each \mathcal{L}_i is maximized by $p_i = E_i / L_i R_i$
- The likelihood of the entire tree is then

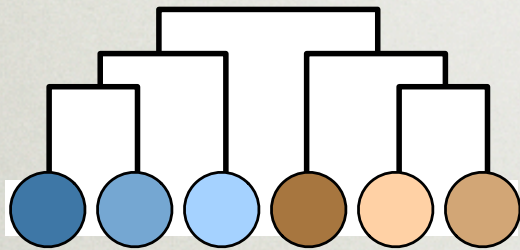
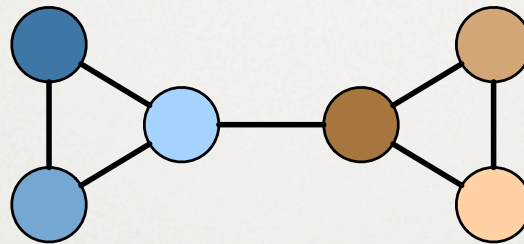
$$\mathcal{L}(T) = \prod_i \mathcal{L}_i$$

- The *log-likelihood* is

$$\ln \mathcal{L}(T) = - \sum_i L_i R_i h(E_i / L_i R_i)$$

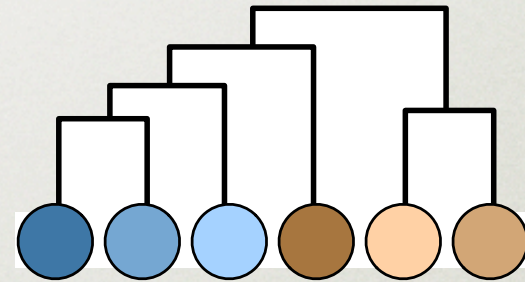
where $h(p) = -p \ln p - (1 - p) \ln(1 - p)$

MAXIMUM LIKELIHOOD



$$\mathcal{L} = \left(\frac{1}{9}\right) \left(\frac{8}{9}\right)^8$$

$$= 0.0433$$

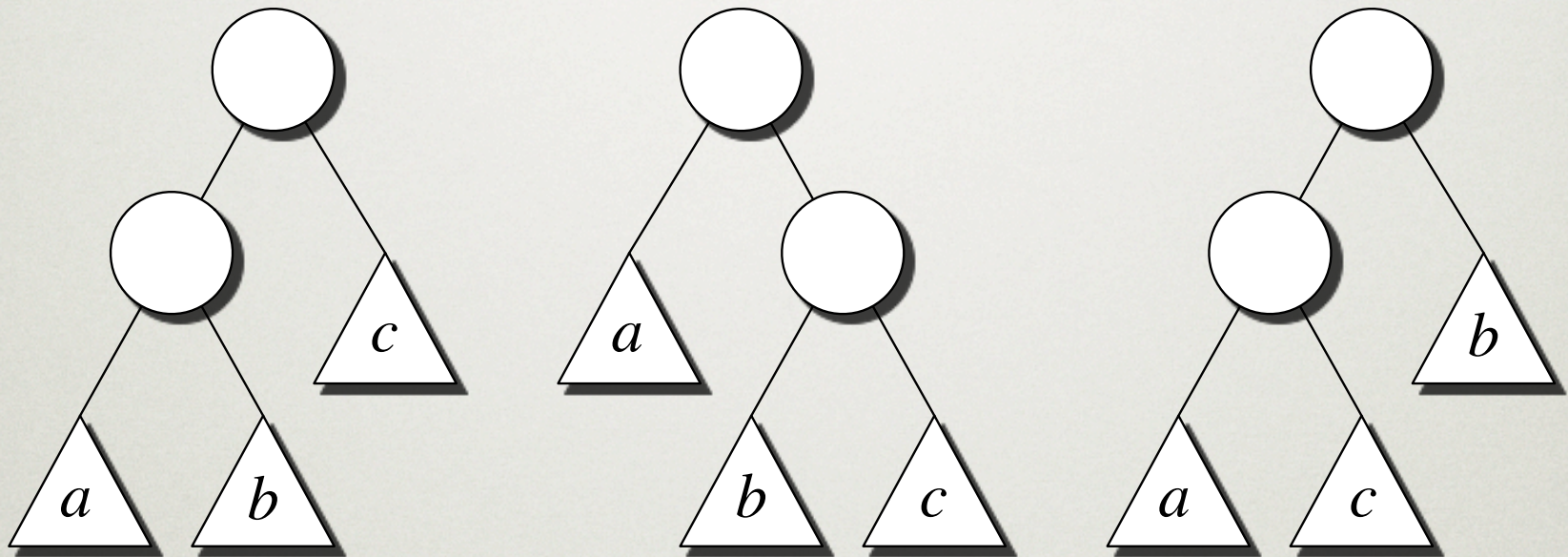


$$\mathcal{L} = \left[\left(\frac{1}{3}\right) \left(\frac{2}{3}\right)^2 \right] \cdot \left[\left(\frac{2}{8}\right)^2 \left(\frac{6}{8}\right)^6 \right]$$

$$= 0.0016$$

A MARKOV CHAIN

- We update the tree T with rotations:

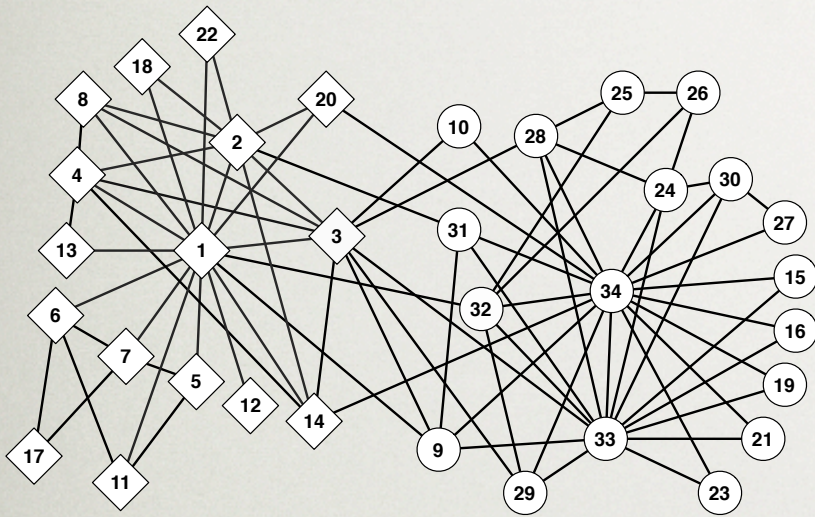


- We move with probability 1 if $\Delta \ln \mathcal{L} \geq 0$ and probability $\exp(\Delta \ln \mathcal{L})$ otherwise

SOME NICE PROPERTIES

- Easy to calculate $\Delta \ln \mathcal{L}$: just “local” terms
- Moves entire chunks of the tree, while keeping their internal structure the same
- Allows us to sample trees with probability proportional to \mathcal{L} , instead of just the one with max likelihood (helps avoid overfitting)

LET'S TRY IT OUT

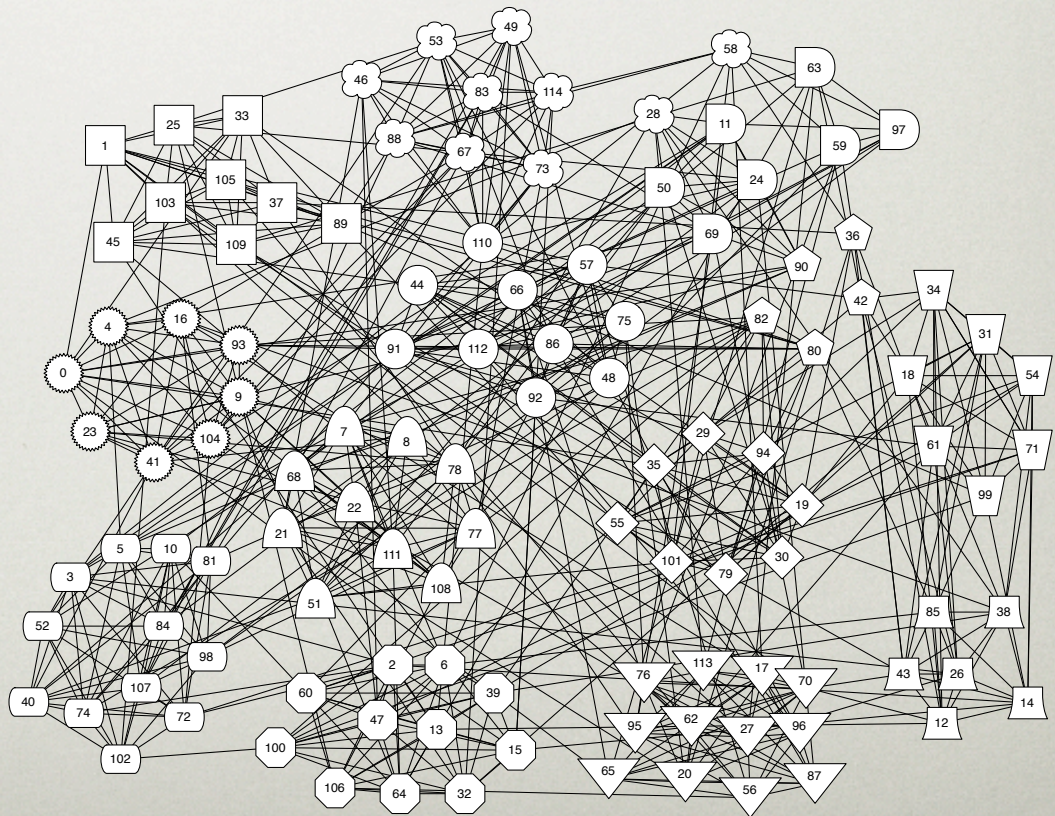


Zachary's Karate Club

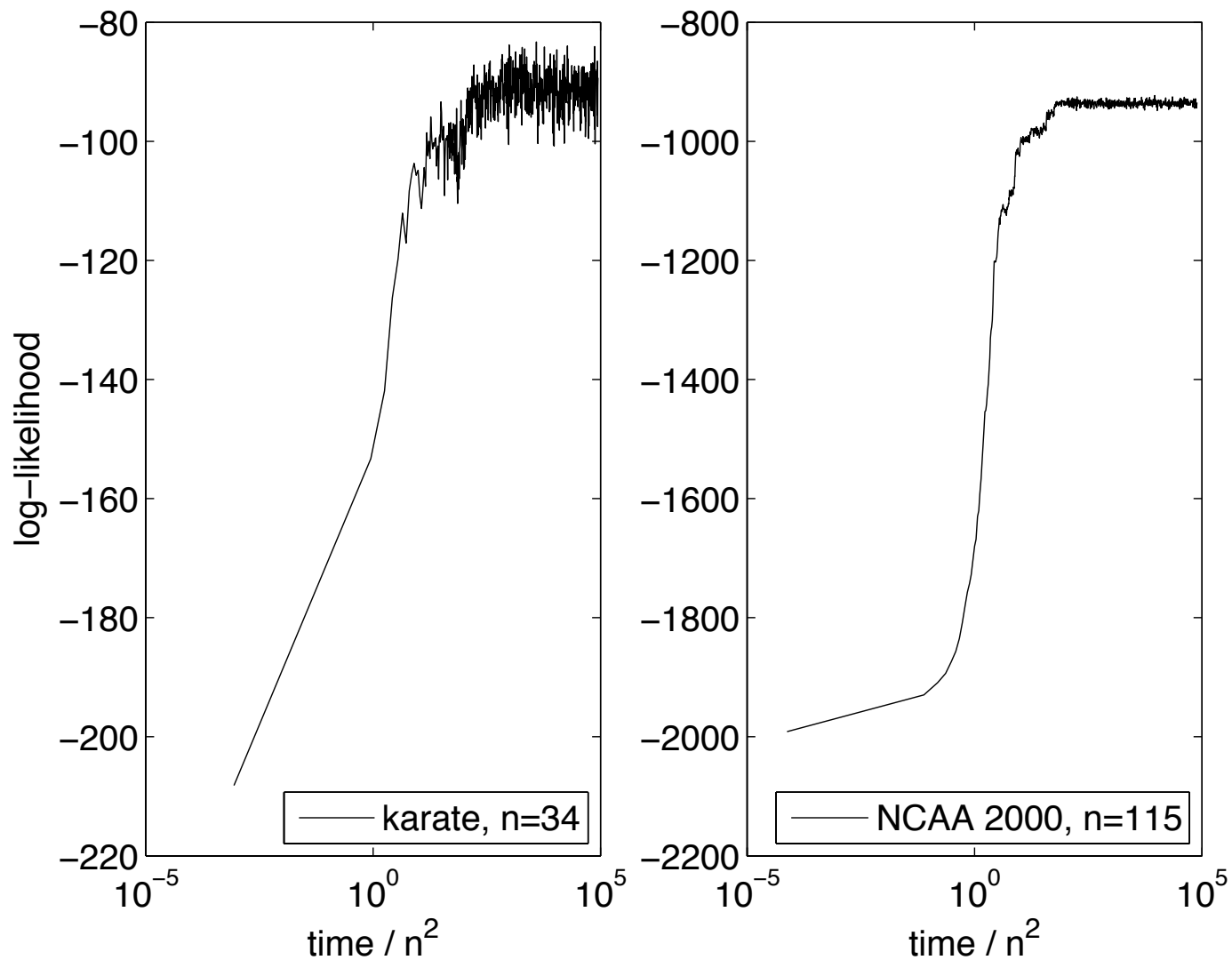
$$n = 34 \quad m = 78$$

NCAA Schedule 2000

$$n = 115 \quad m = 613$$



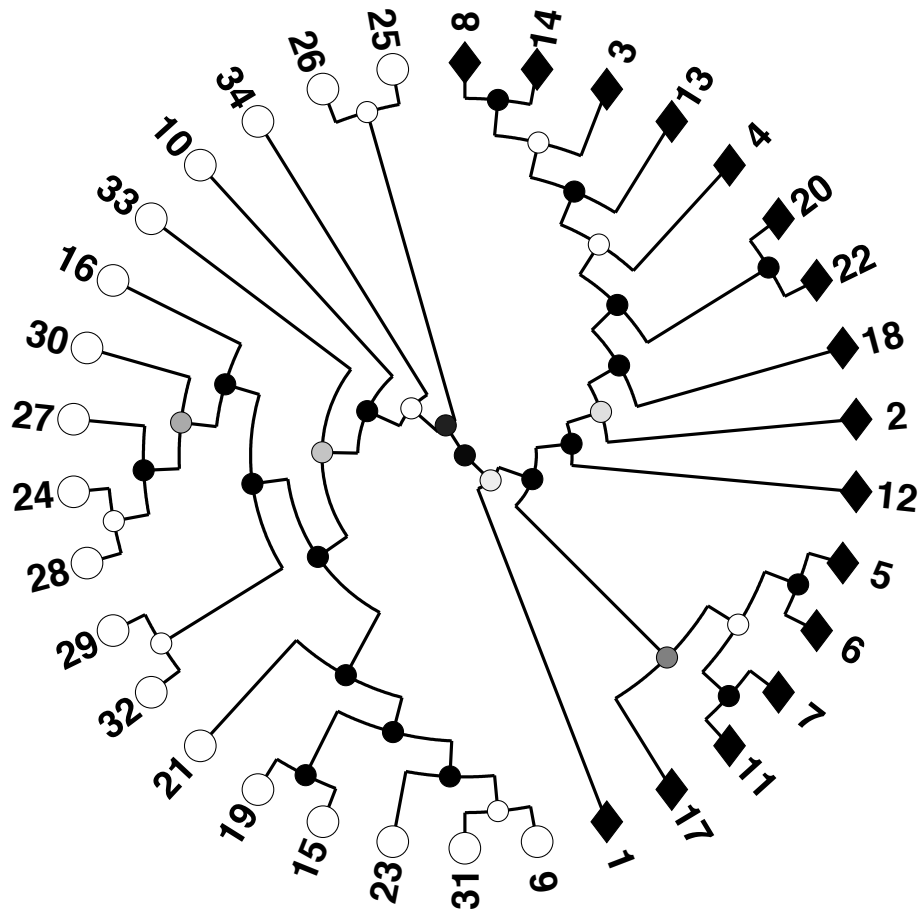
EXPERIMENTS



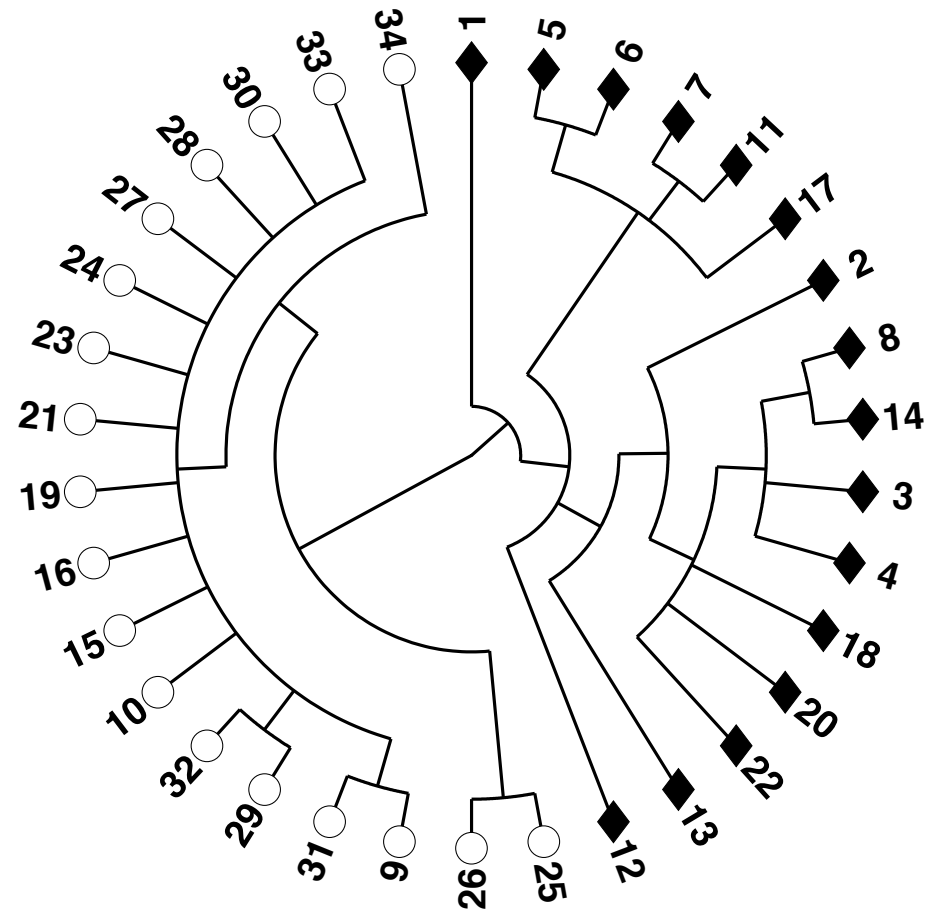
CONSENSUS DENDROGRAMS

- Let's sample many trees instead of one.
- From phylogeny construction: combine these into a *consensus* hierarchy, which includes the (weighted) majority of splits
- More appropriate than any single tree (even the most likely one)

CONSENSUS DENDROGRAM: KARATE CLUB



point estimate

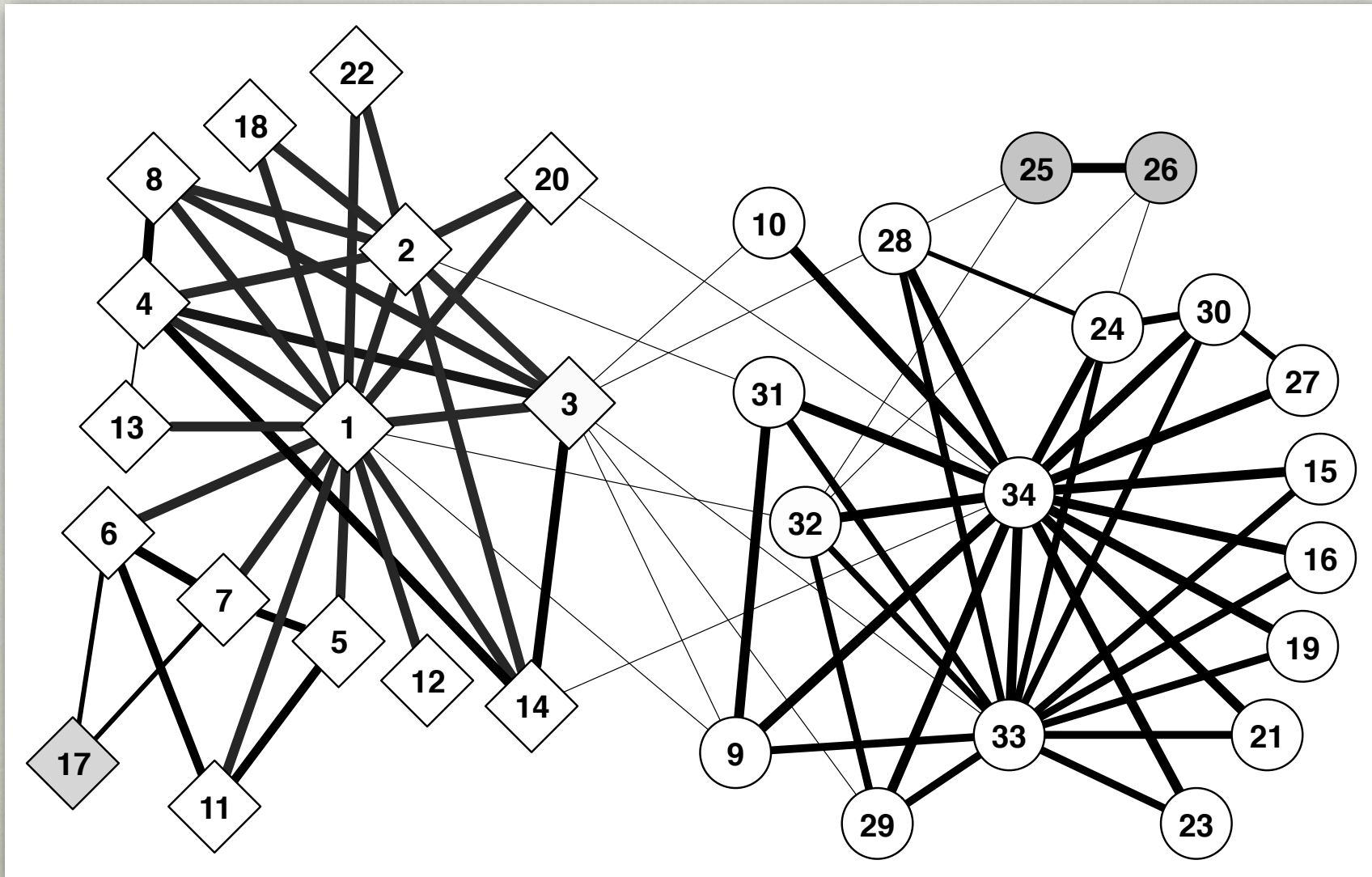


consensus hierarchy

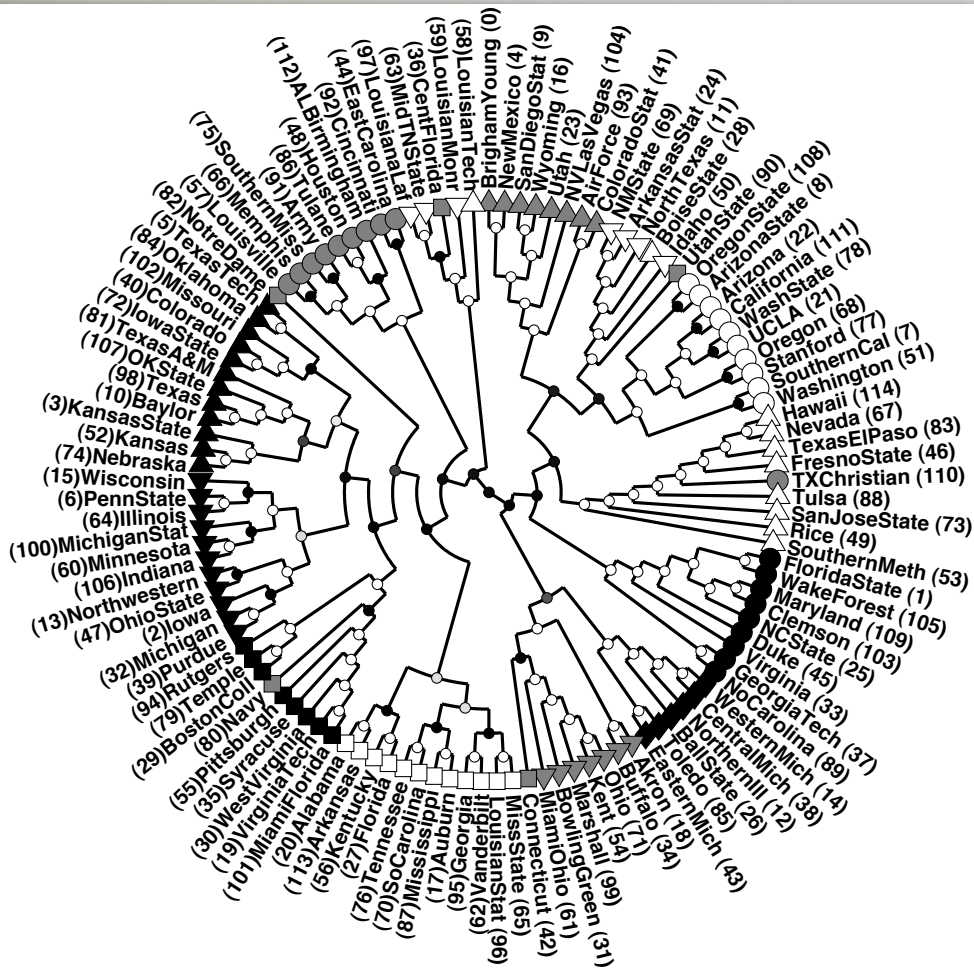
SAMPLING THE DISTRIBUTION

- Again, sample many trees instead of one.
- Calculate the (weighted) average probability that two vertices are connected, or that a vertex is part of a given community.
- Lets us classify how strongly a vertex is part of a group, or how “surprising” an edge is.

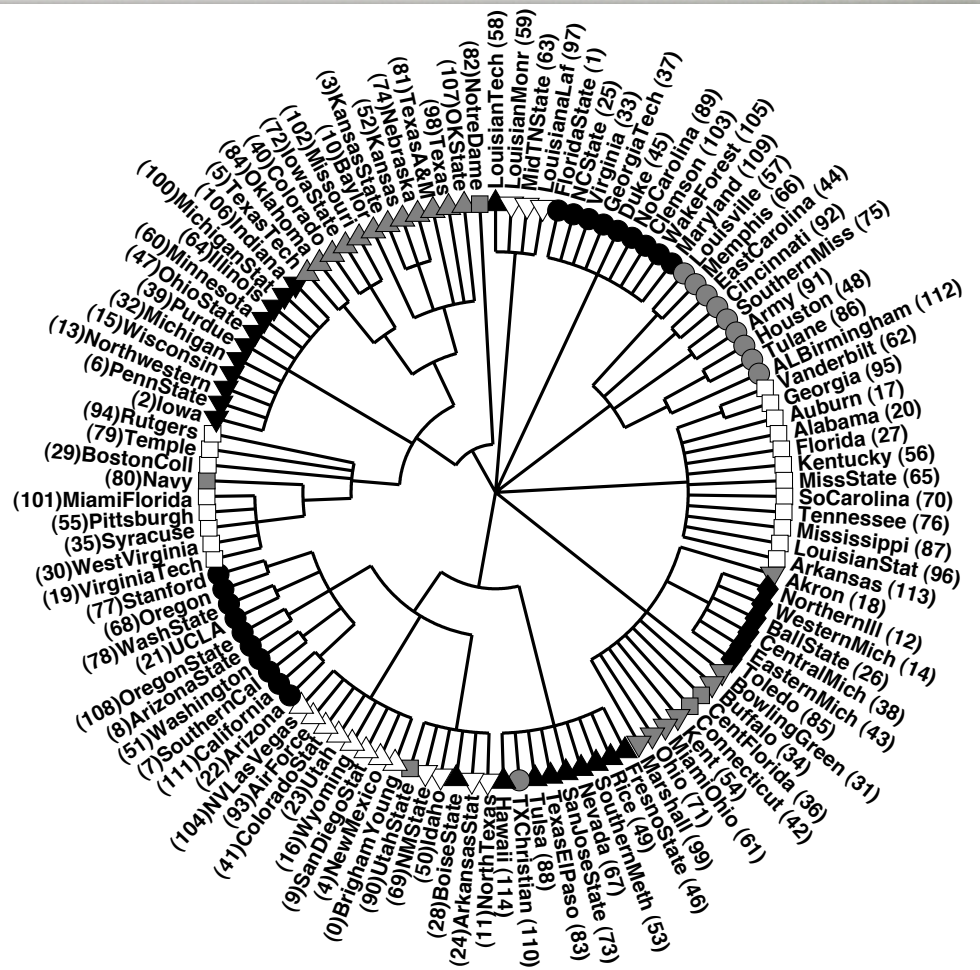
AVERAGE PROBABILITY OF EDGES AND VERTICES



CONSENSUS DENDROGRAM: NCAA

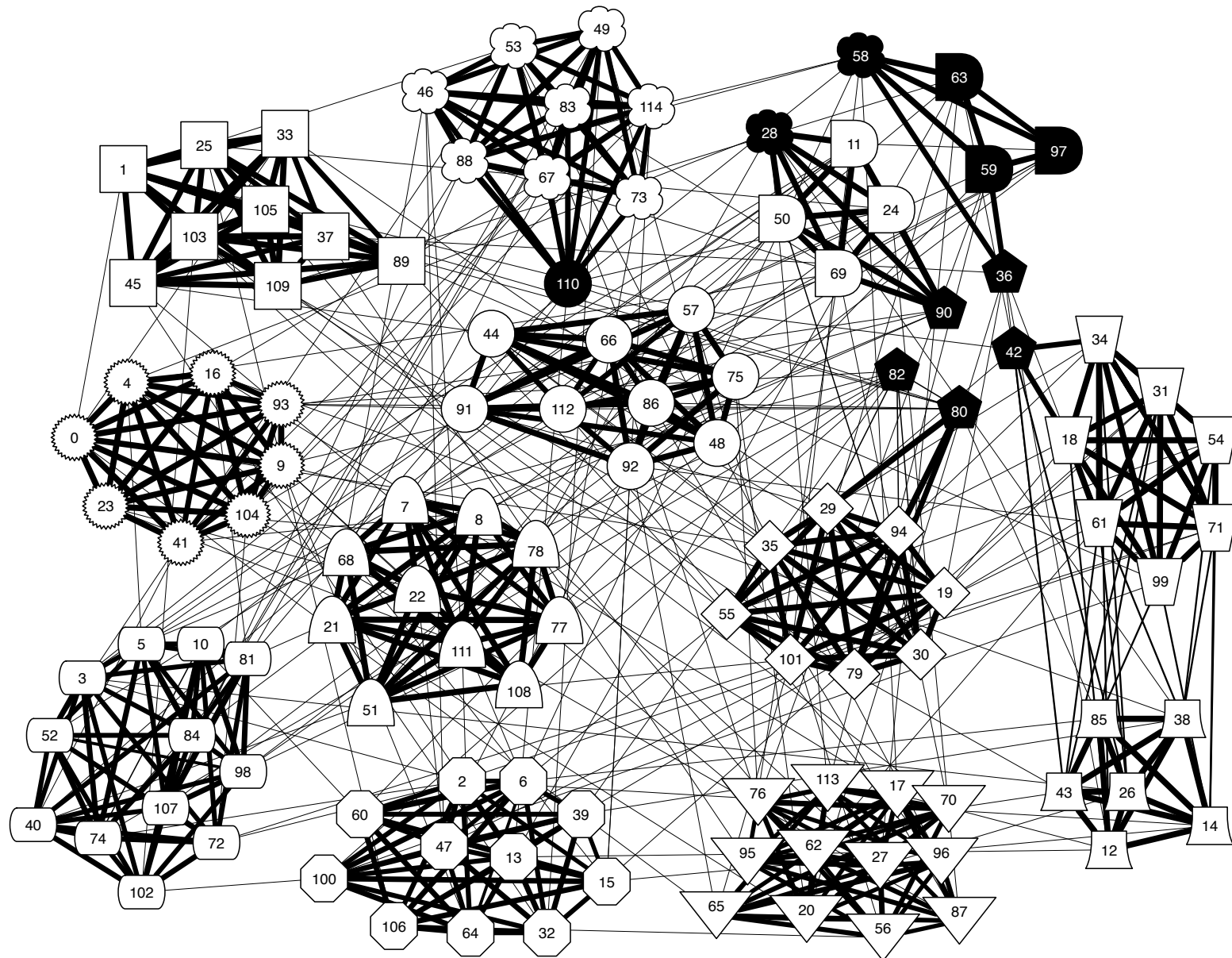


point estimate



consensus hierarchy

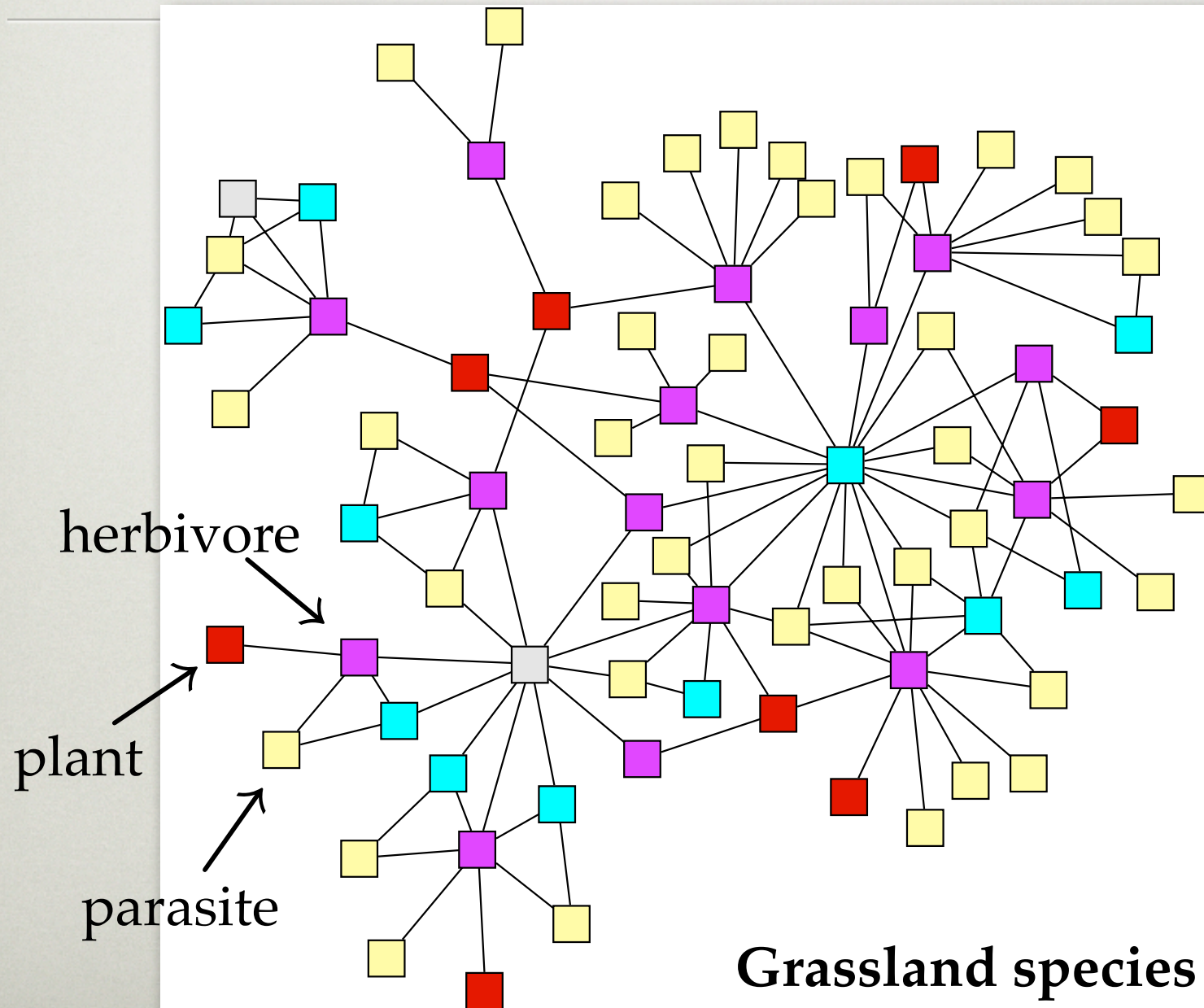
AVERAGE PROBABILITY OF EDGES AND VERTICES



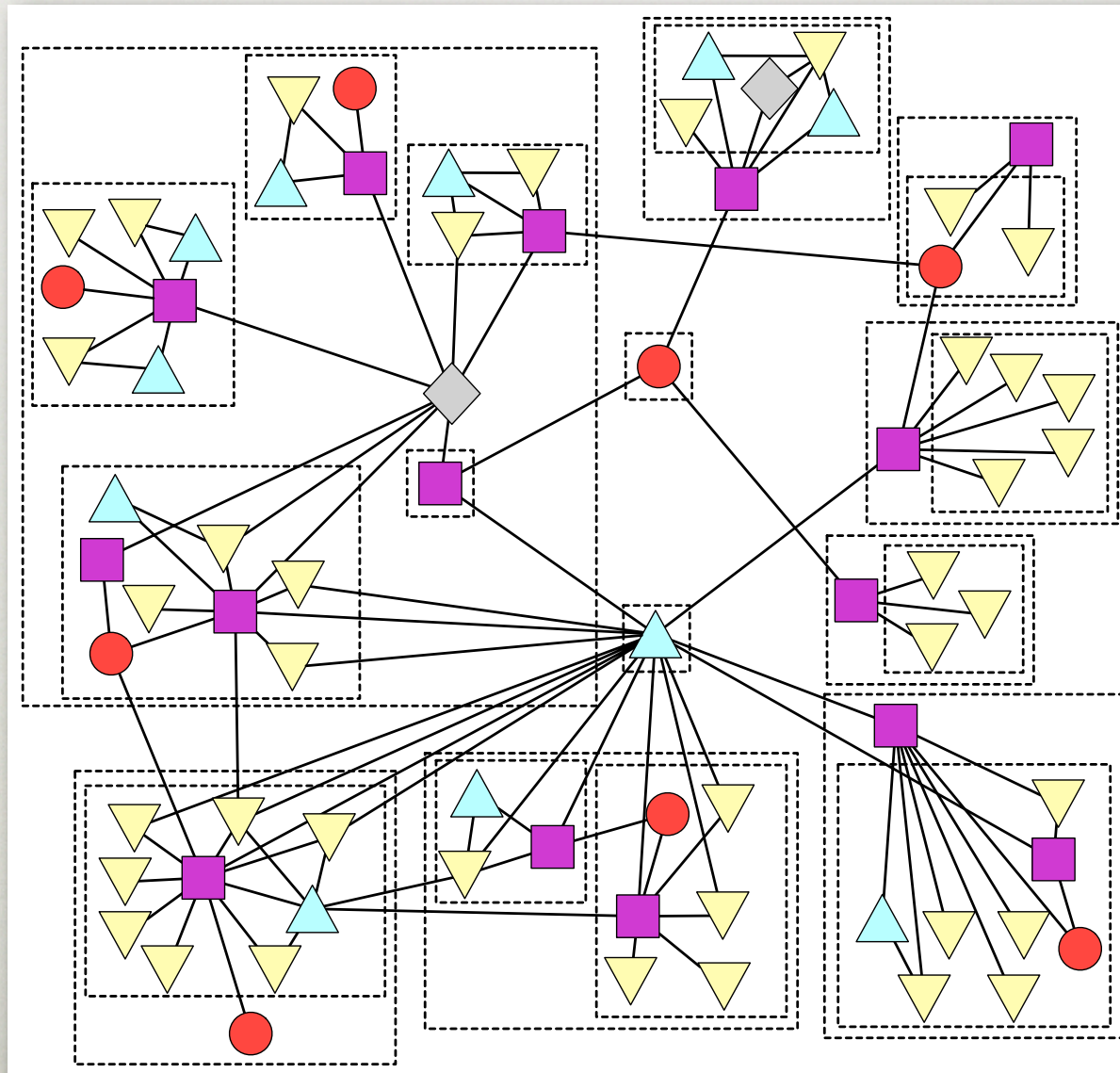
GOAL #2: GENERATING SIMILAR RANDOM GRAPHS

- Idea: once we infer a hierarchical structure (or a distribution of them) we can use it to generate new random graphs.
- If these graphs are “similar” to the original (*n.b.: application-dependent!*) then we may have captured part of its structure.
Generalizing from a single example!
- If they aren’t, we have falsified our model, which is a good thing to be able to do...

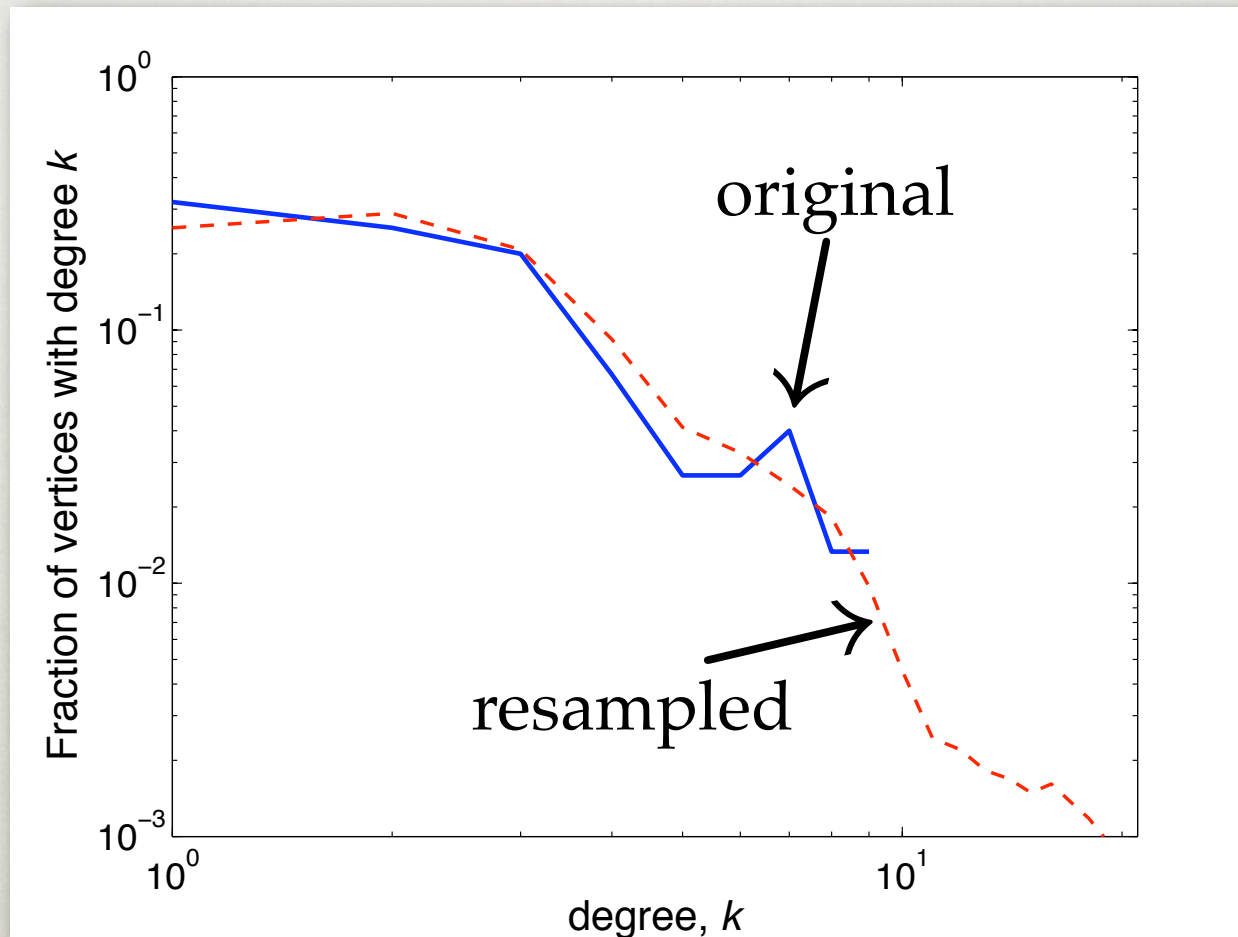
A FOOD WEB



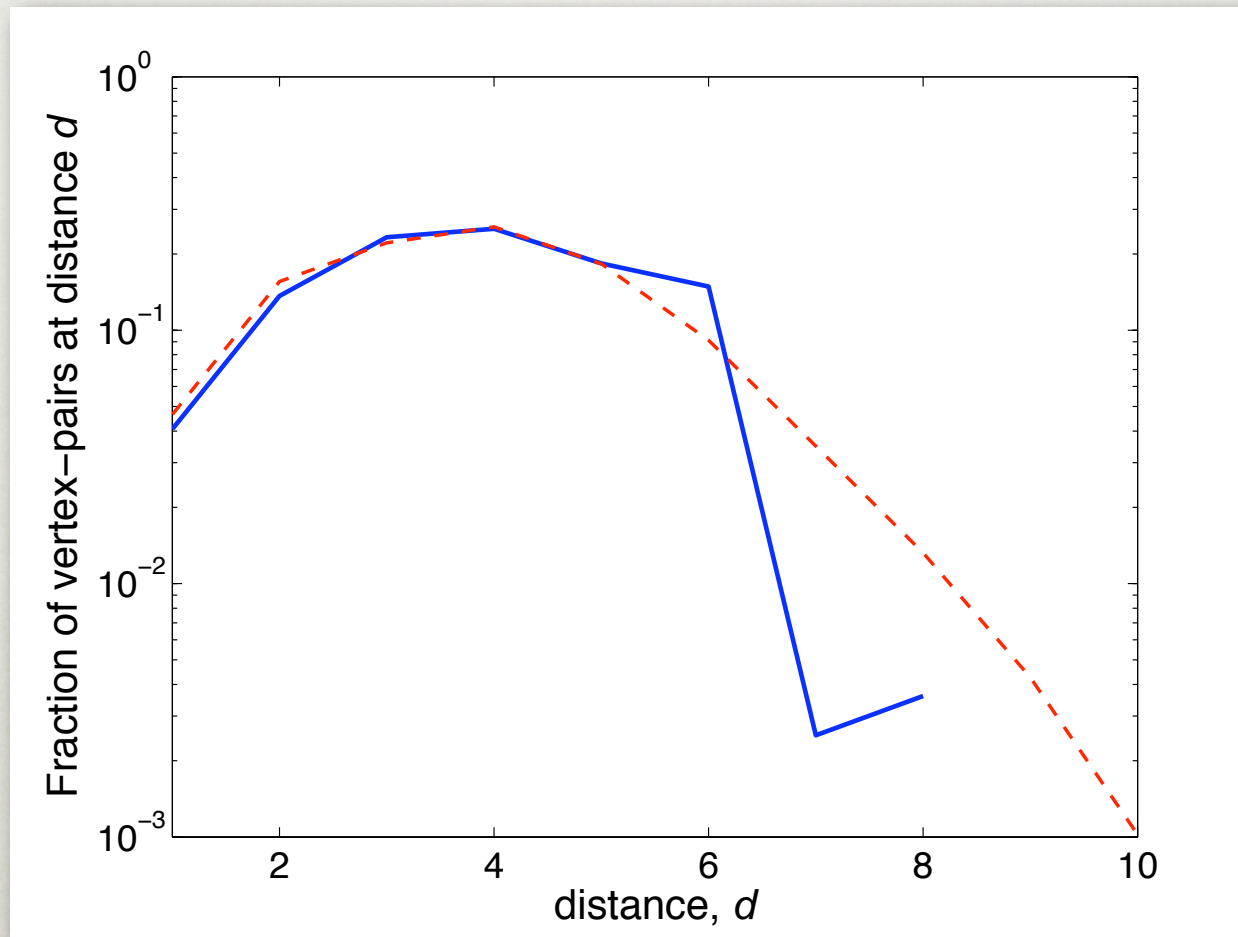
AN INFERRED HIERARCHY



RESAMPLING THE ENSEMBLE: DEGREE DISTRIBUTION



RESAMPLING THE ENSEMBLE: DISTANCE DISTRIBUTION



GOAL #3: LINK PREDICTION

- For many networks, edges are discovered one at a time, using difficult work in the field or laboratory
- Given the edges observed so far, can we predict missing ones better than chance?
- If so, we can focus our attention on pairs of vertices likely to be connected.

OUR APPROACH

- Sample hierarchies using observed edges
- Sort remaining pairs according to the average probability they are connected
- Predict the top few of these
- Cross-validation: remove some fraction of edges randomly, and try to re-predict the ones you removed

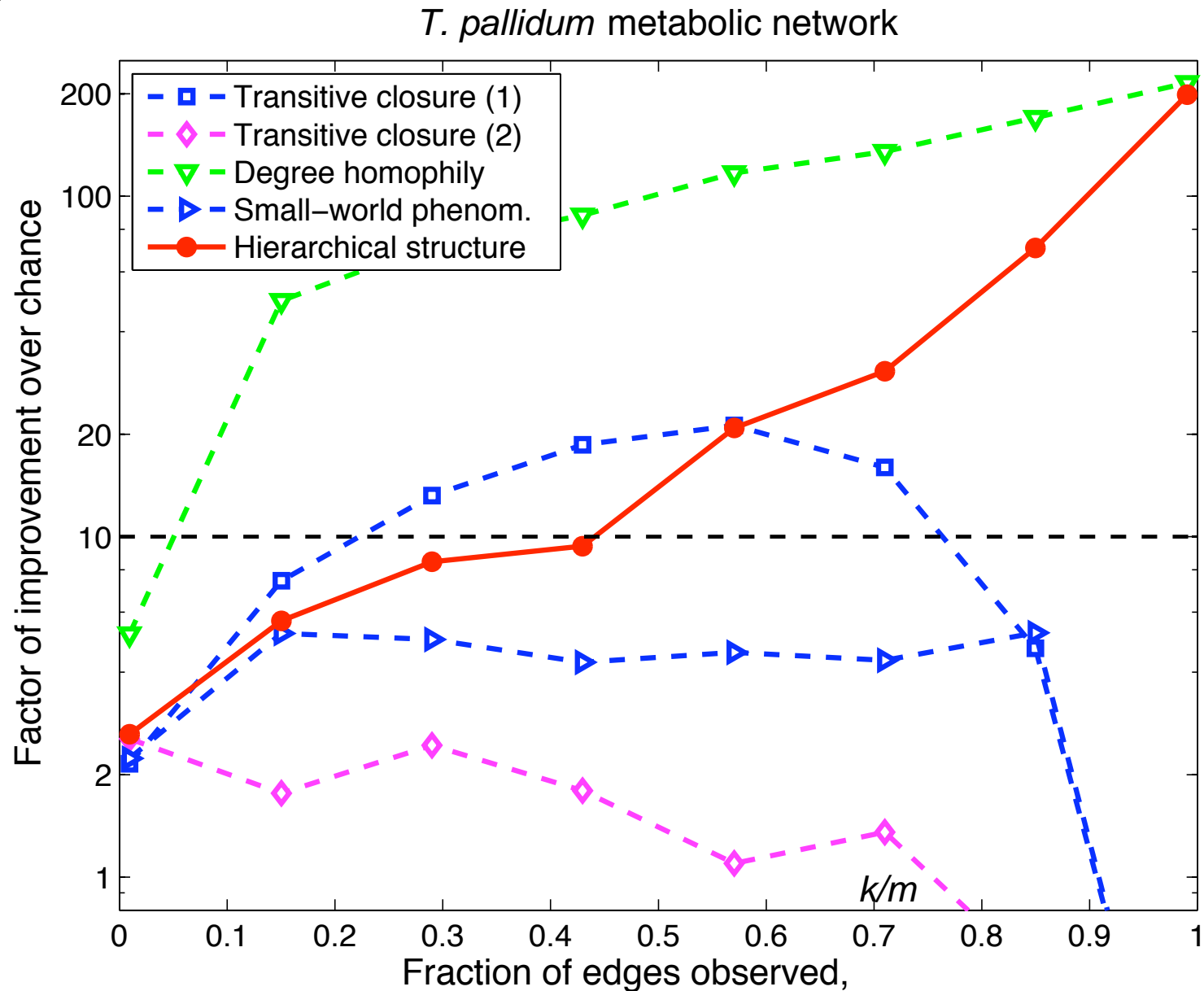
SOME SIMPLER METHODS

“Link Prediction” problem [Liben-Nowell and Kleinberg, 2003: predicting collaborations]

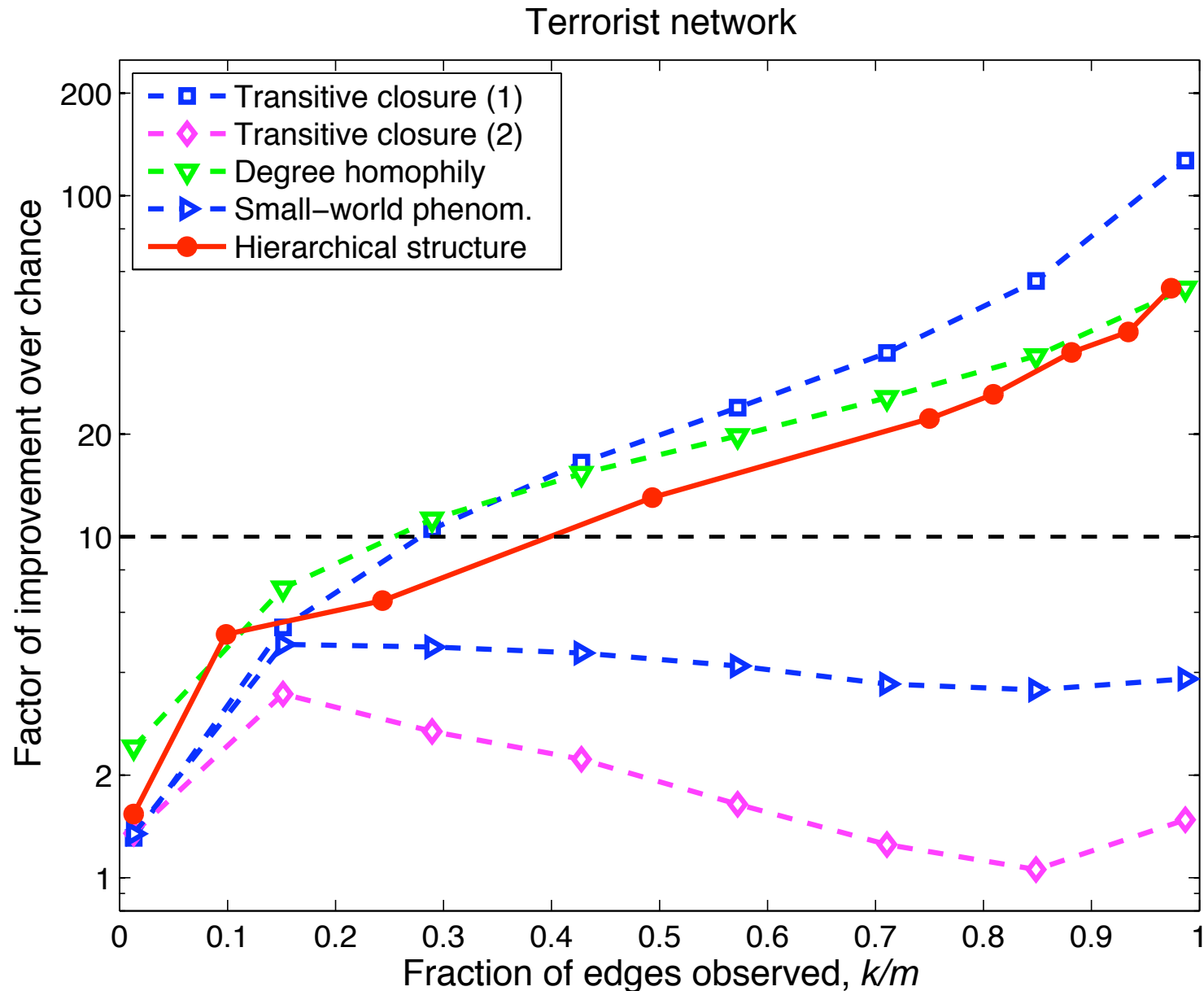
Guess that two vertices are connected if:

- They have many common neighbors
- They share a large fraction of their neighbors (Jaccard coefficient)
- The product of their degrees is large
- They have many short paths between them

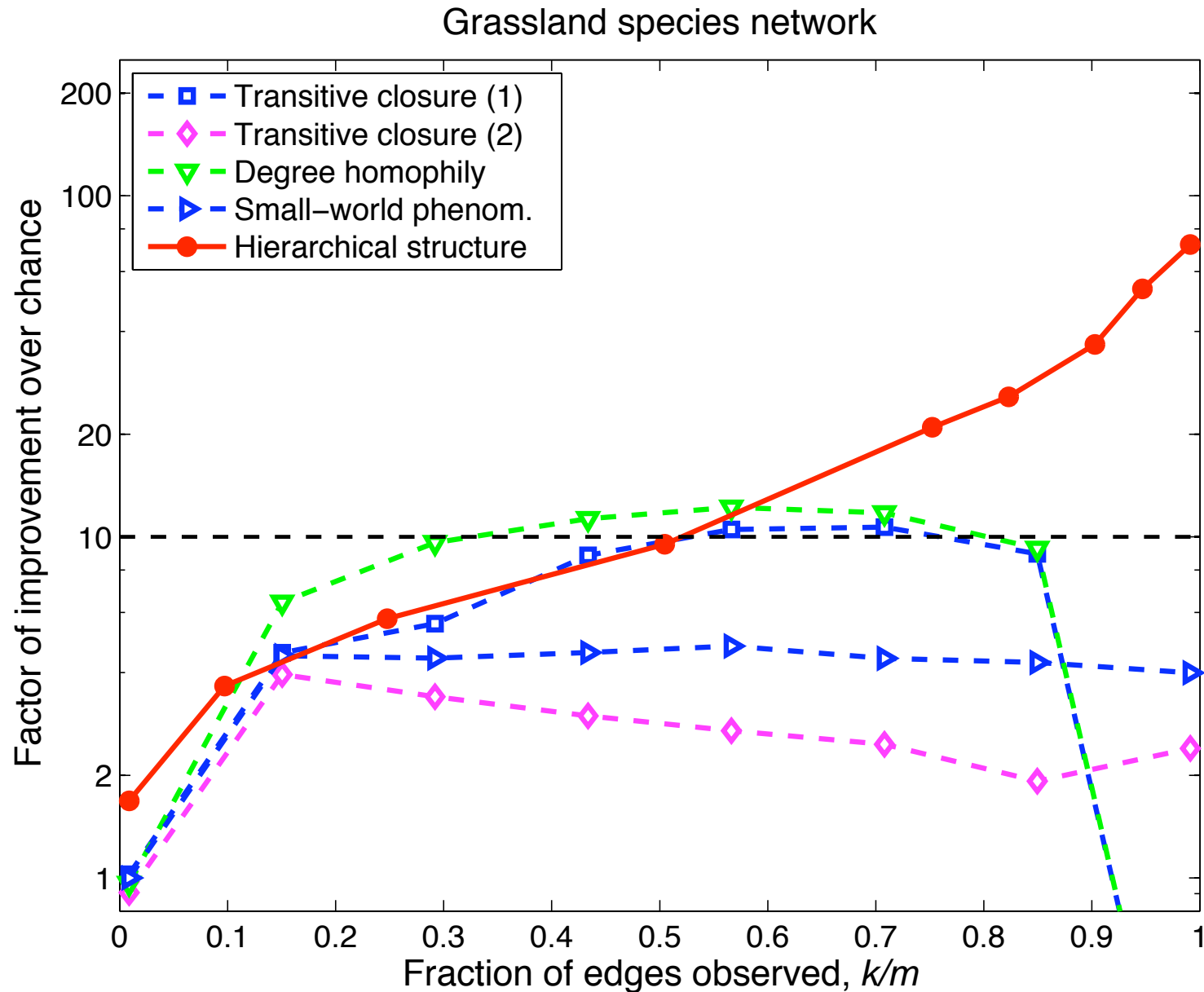
OTHER METHODS: METABOLIC



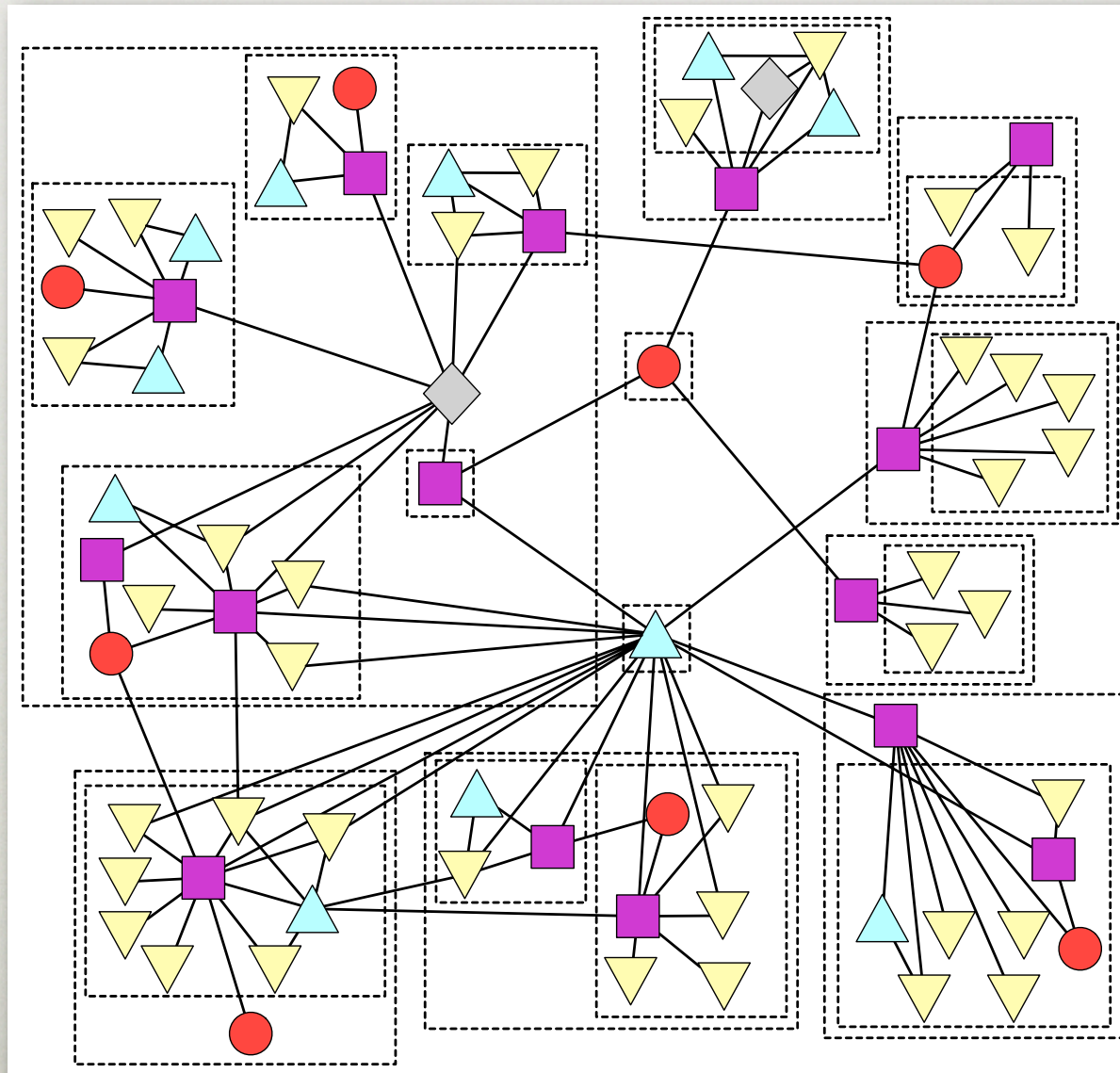
OTHER METHODS: BAD GUYS



OTHER METHODS: FOOD WEB



DISASSORTATIVITY: PREDATORS SHARE PREY



CONCLUSION

- We have a rich space of hierarchical random graph models, not just assortative
- We use a Markov chain to sample the likely models from this space based on data
- This lets us generate random graphs which are statistically similar to an observed one
- It also lets us predict missing links with probability much better than chance.

SHAMELESS PLUG

Computational Complexity and Statistical Physics

Edited by
Allon Percus
Gabriel Istrate
Cristopher Moore



A VOLUME IN THE
SANTA FE INSTITUTE STUDIES IN THE SCIENCES OF COMPLEXITY

SHAMELESS PLUG

Computational Complexity and Statistical Physics

Edited by
Allon Percus
Gabriel Istrate
Cristopher Moore



A VOLUME IN THE
SANTA FE INSTITUTE STUDIES IN THE SCIENCES OF COMPLEXITY

The Nature of Computation



Mertens and Moore

ACKNOWLEDGMENTS

