

# Networks and Community Structure

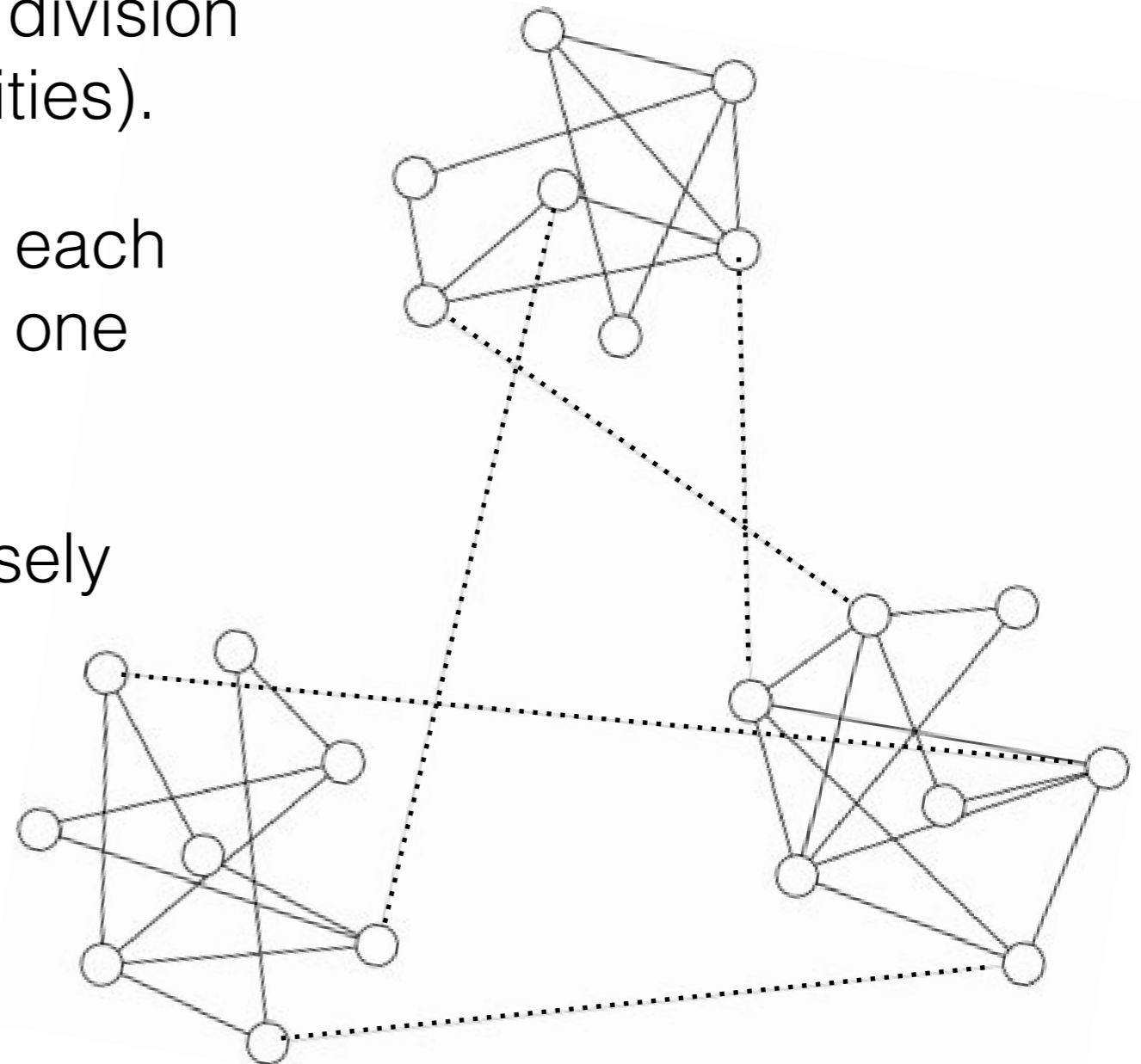
# Community Structure

Community structure is a “natural” division of a network into groups (communities).

Community structure is a *partition*: each node is a member of one and only one group

Within a group, the nodes are densely connected, with only sparse connections between groups

This community structure is sometimes known to the people in the network. Sometimes not.



NB: community structure is (confusingly) sometimes referred to as “clustering”

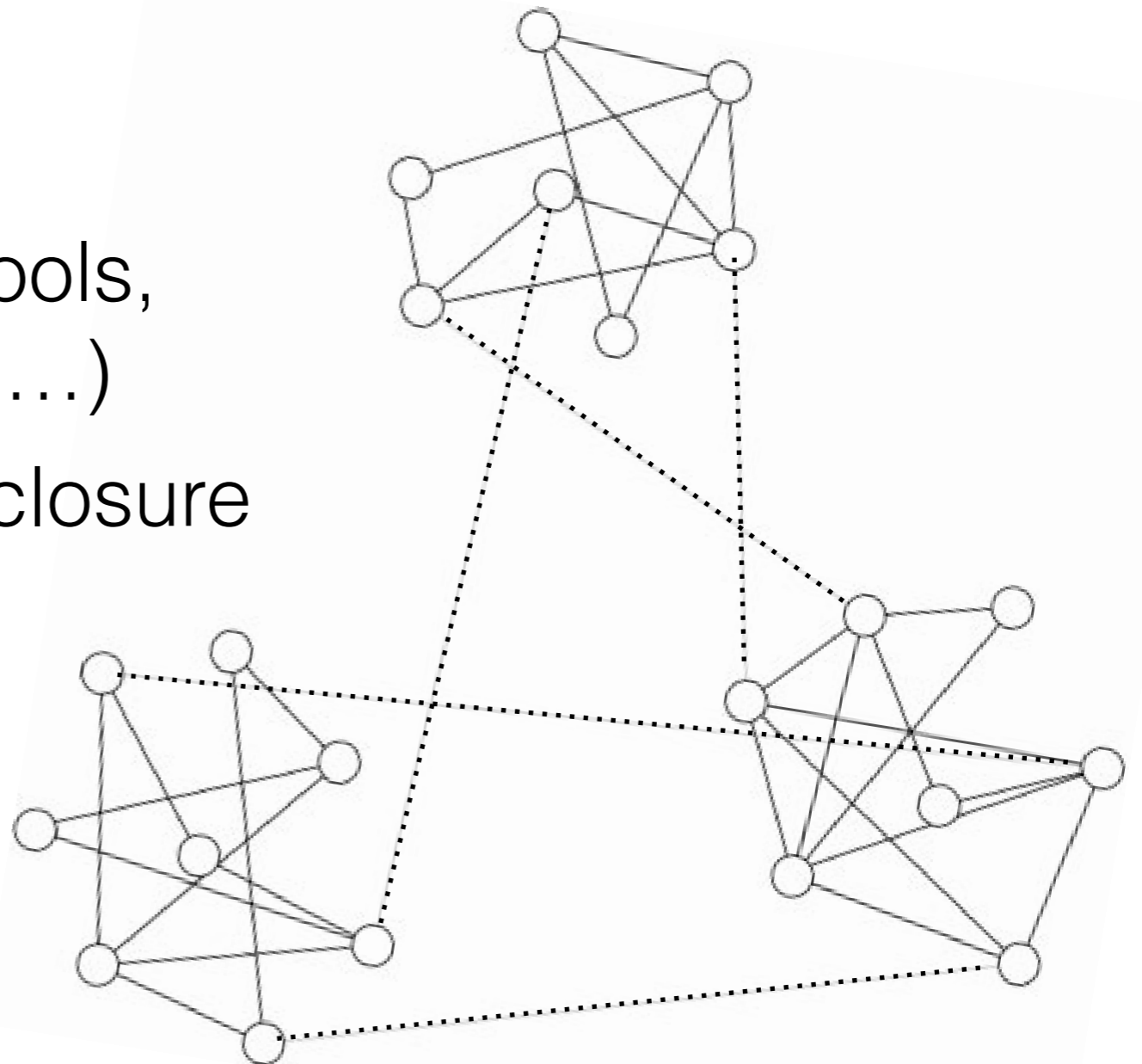
What the heck?



# Community Structure

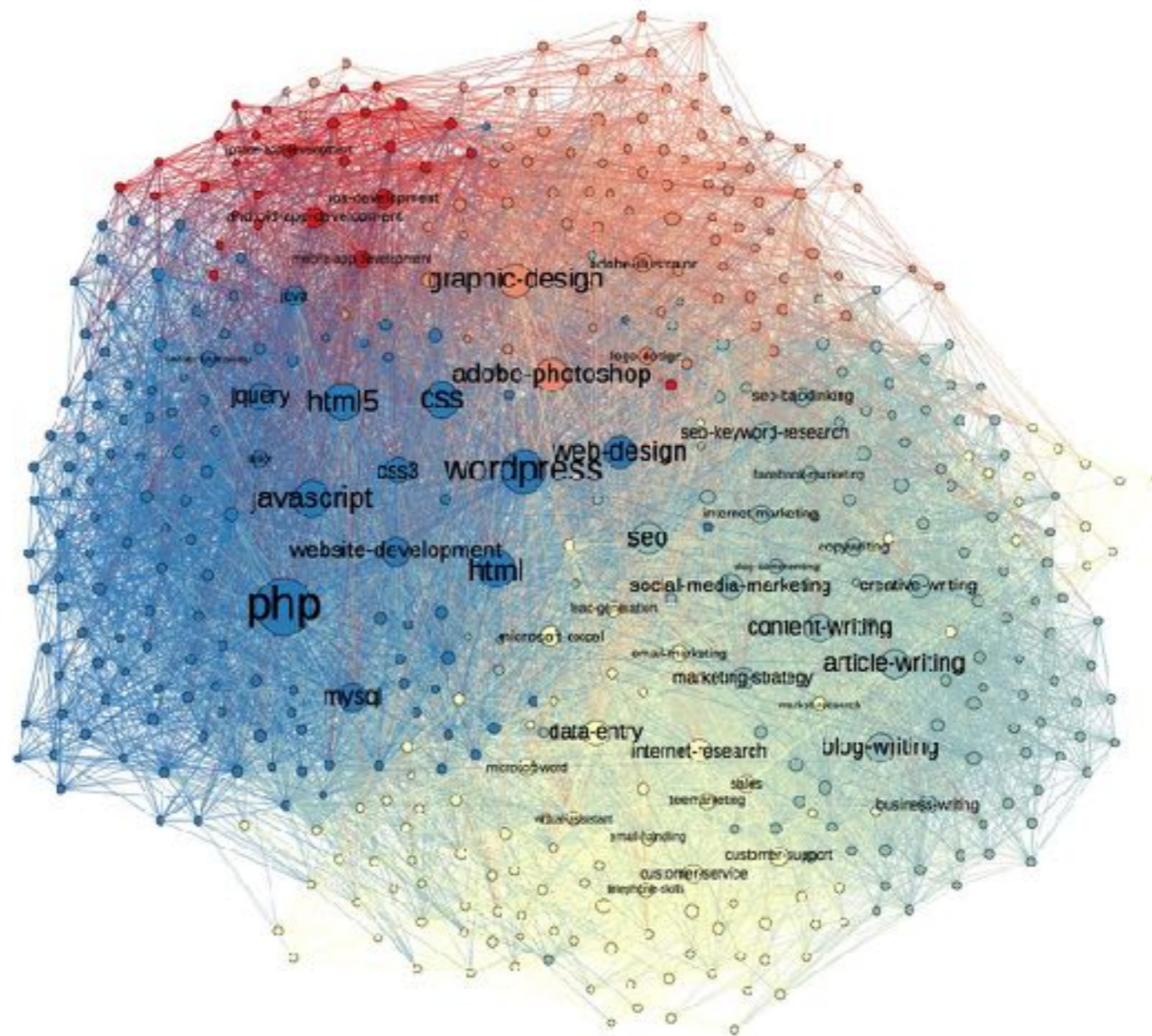
What are some sources of community structure in social groups?

- geography
- family groups
- organizations (e.g. schools, clubs and teams, firms...)
- homophily and triadic closure
- institutional structure
- social norms
- specialization



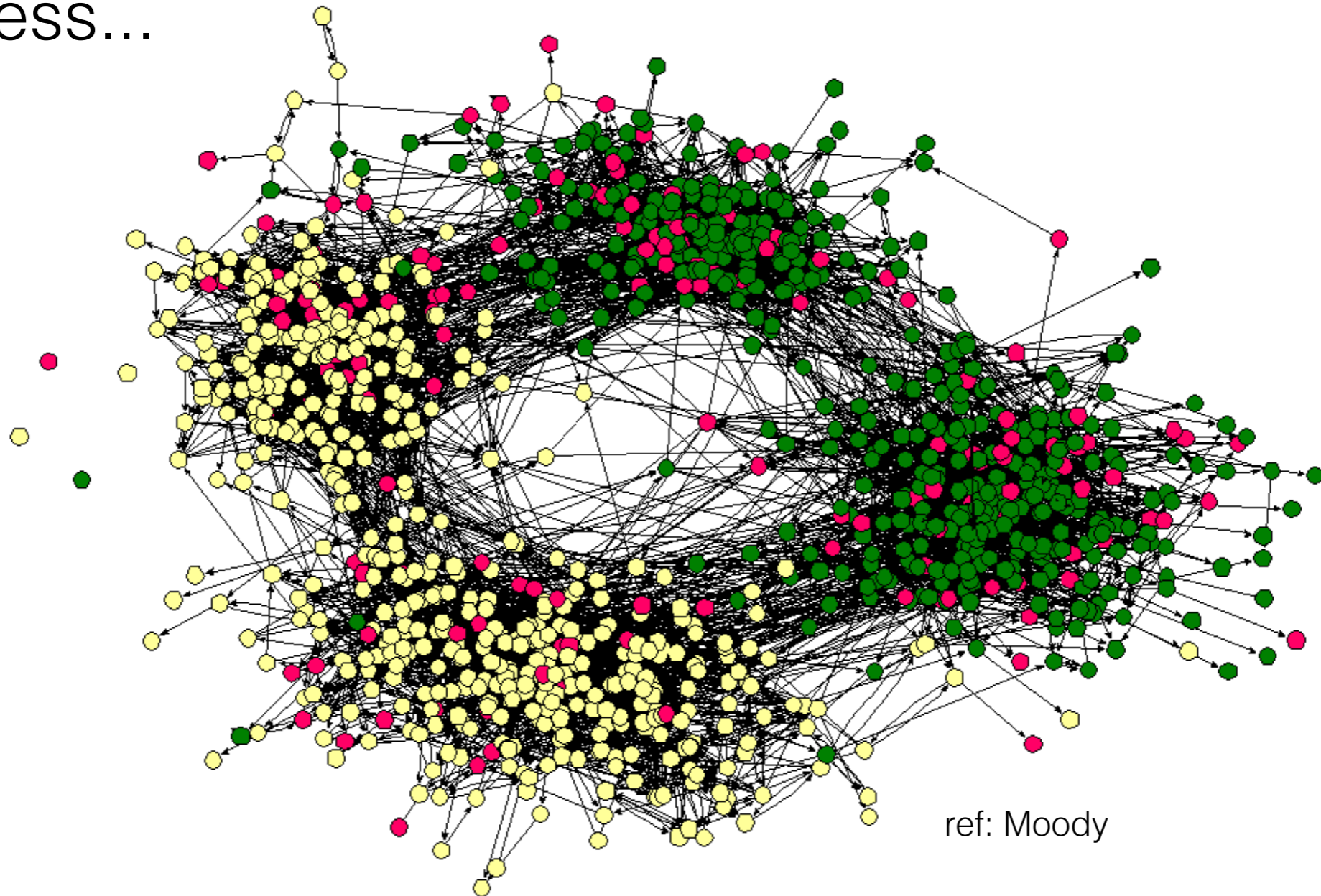
# Community Structure

But why would we care about community structure in a network?



# Community Structure

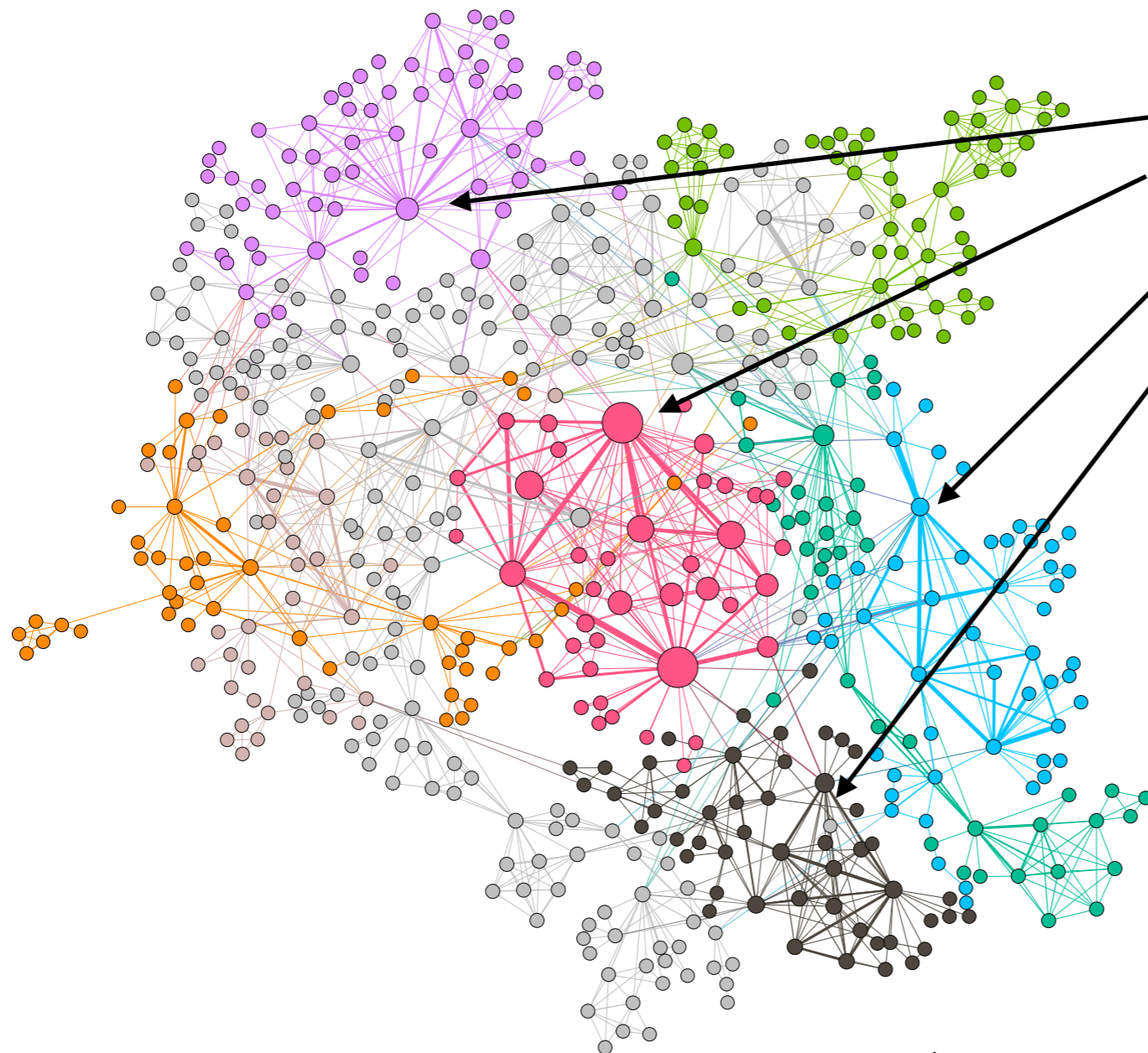
Sometimes it reveals a deeper underlying social process...



...for example, self-segregation

# Community Structure

Sometimes you find something about the function of a group: a coauthorship network for a small scientific field (Physics Education Research)



Communities are centered around the founders of the field

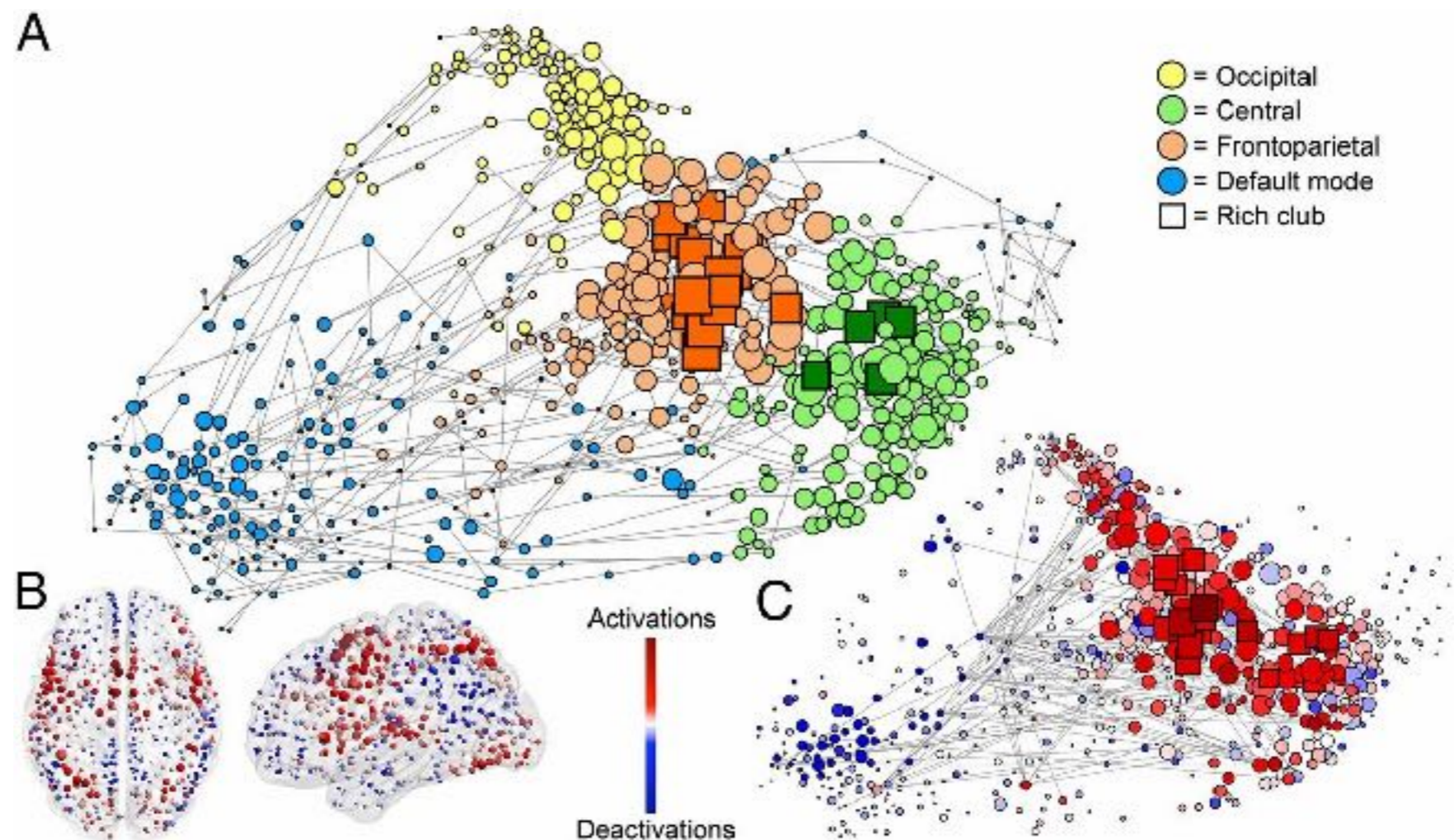
Suggests that scholars are more likely to work with other people who studied with their PhD advisor

Why might that be?

What might change that?

# Community Structure

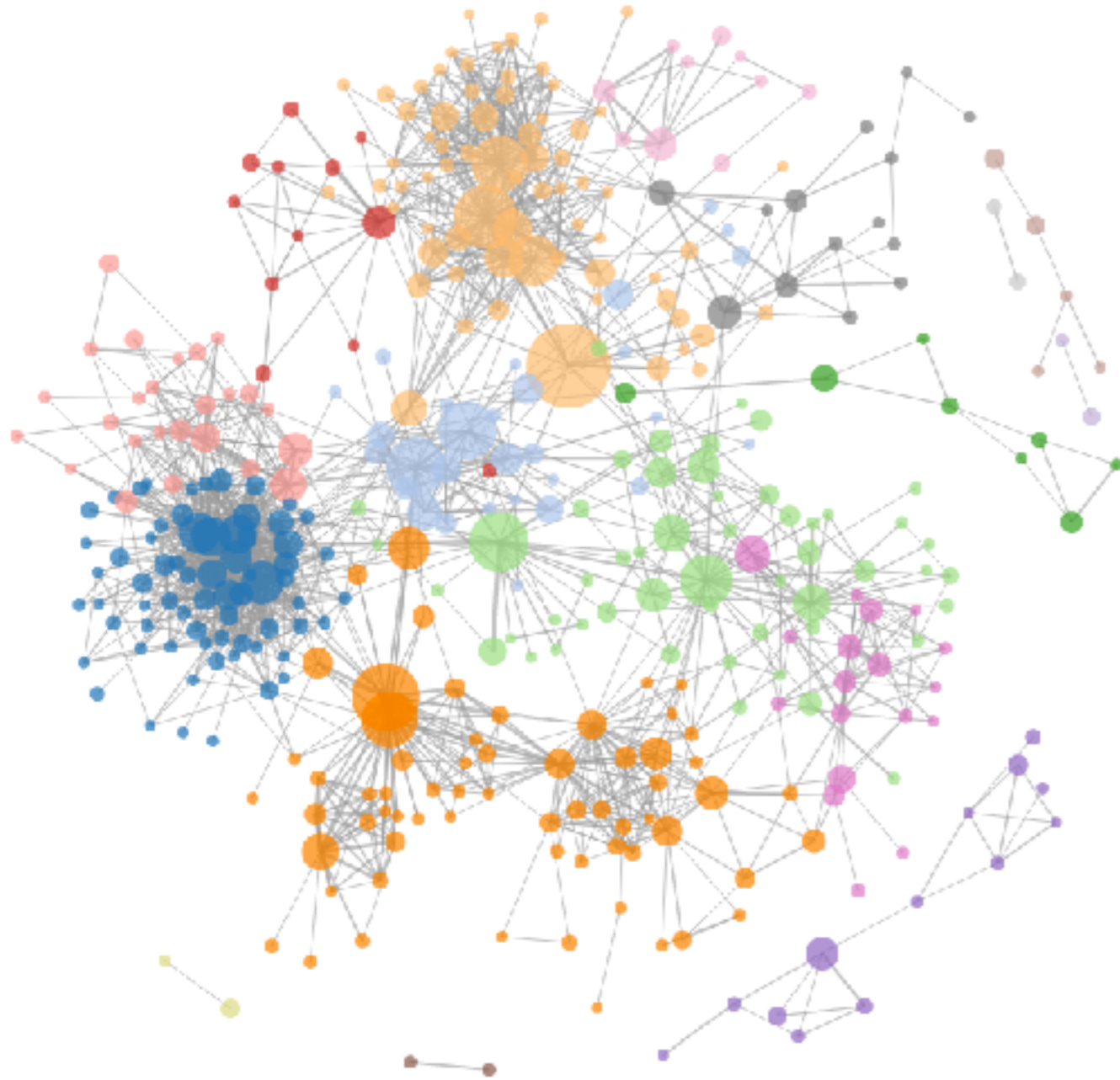
When you are looking at networks that aren't social, there may be different causes for community structure...



ref: Crossley et al. (PNAS 2013)

Example: in neural networks, community structure might reflect the underlying function...

# Community Structure



Citation network (sociology)

- nodes = papers
- $A \rightarrow B$  if paper A cites paper B

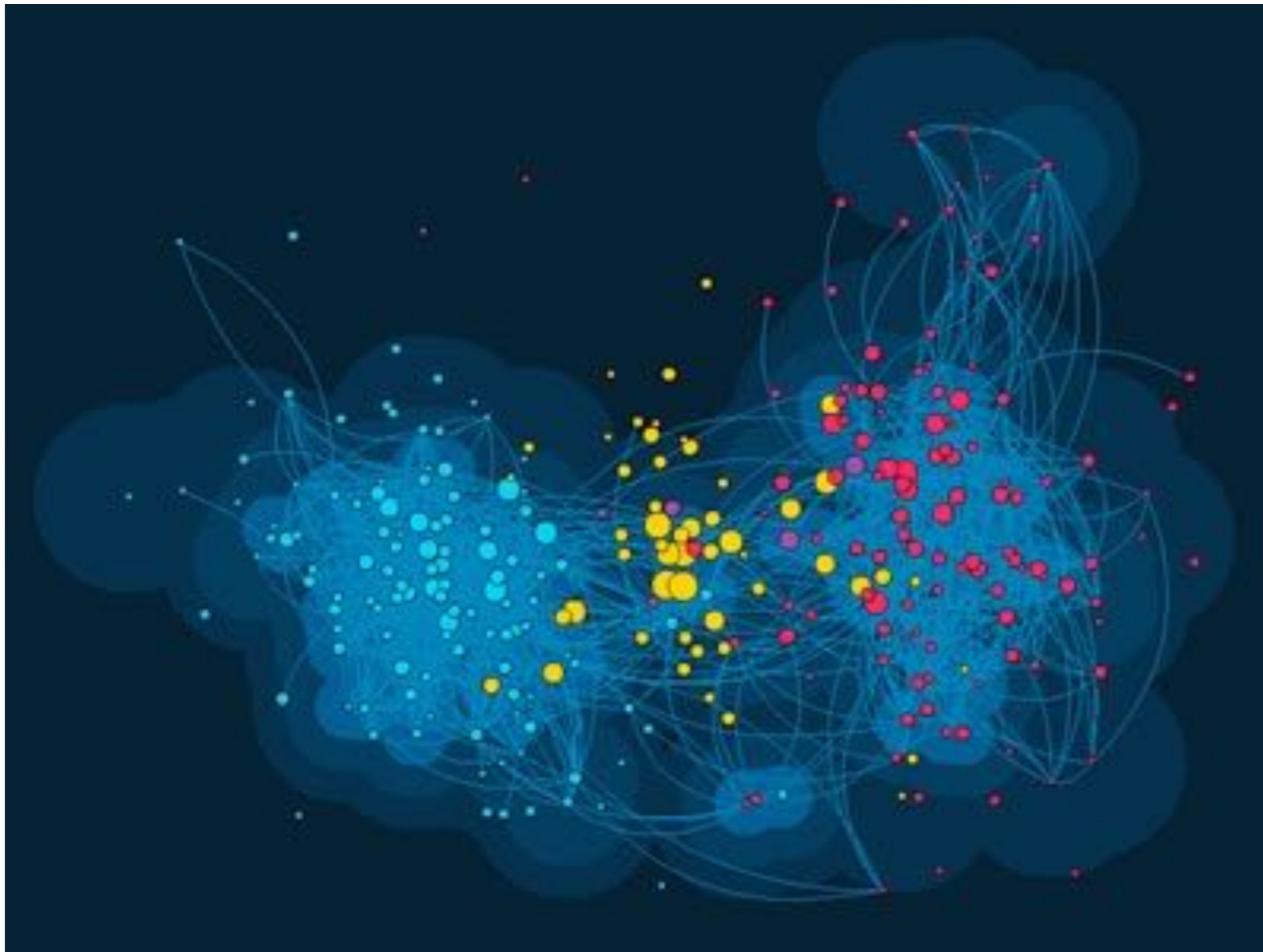
What do the communities represent?

Where might they come from?

ref: <http://goo.gl/L9ars>



# Community Detection

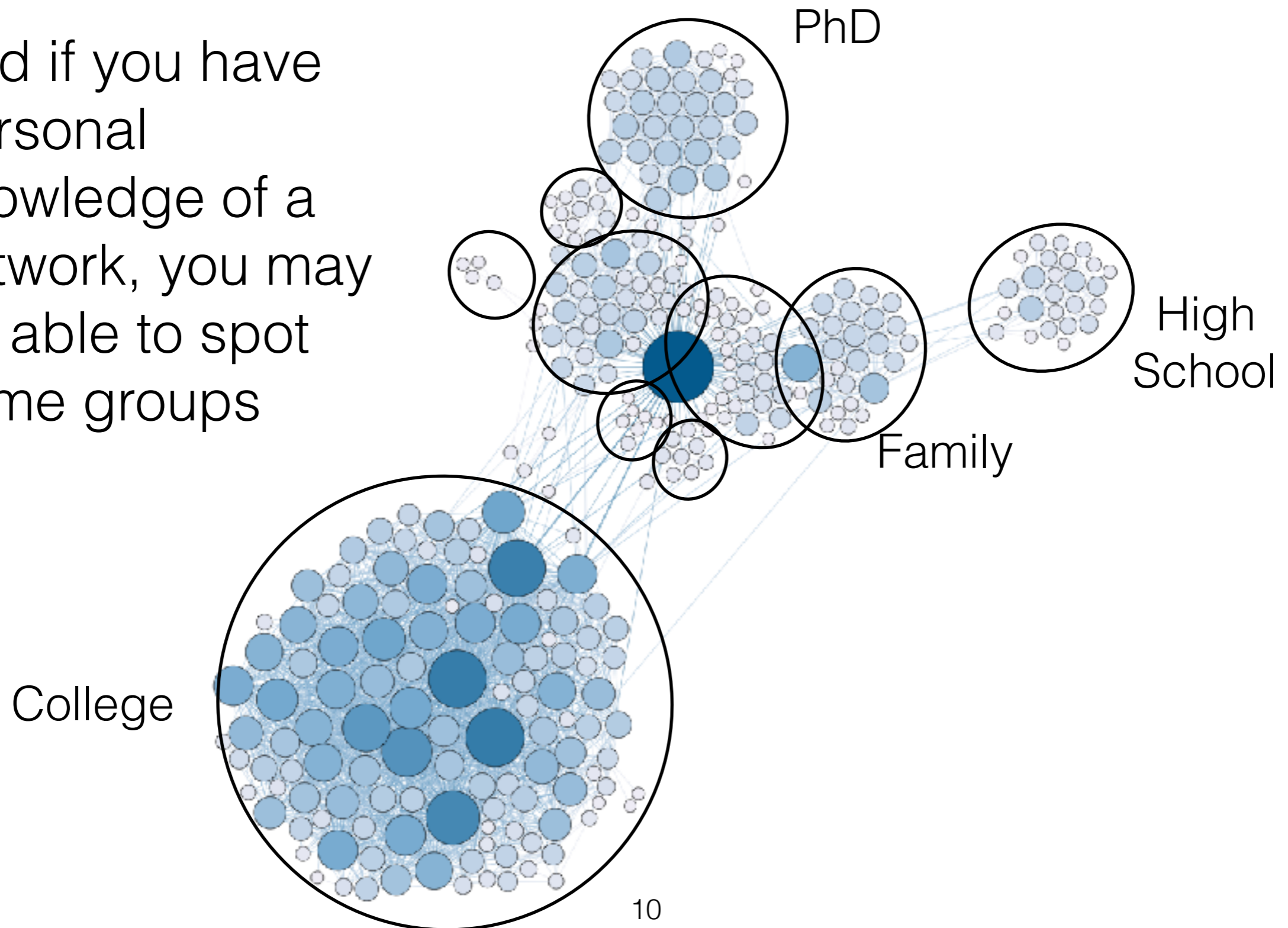


In some cases,  
community  
structure is easy  
to detect by  
eye...

ref: Lada Adamic

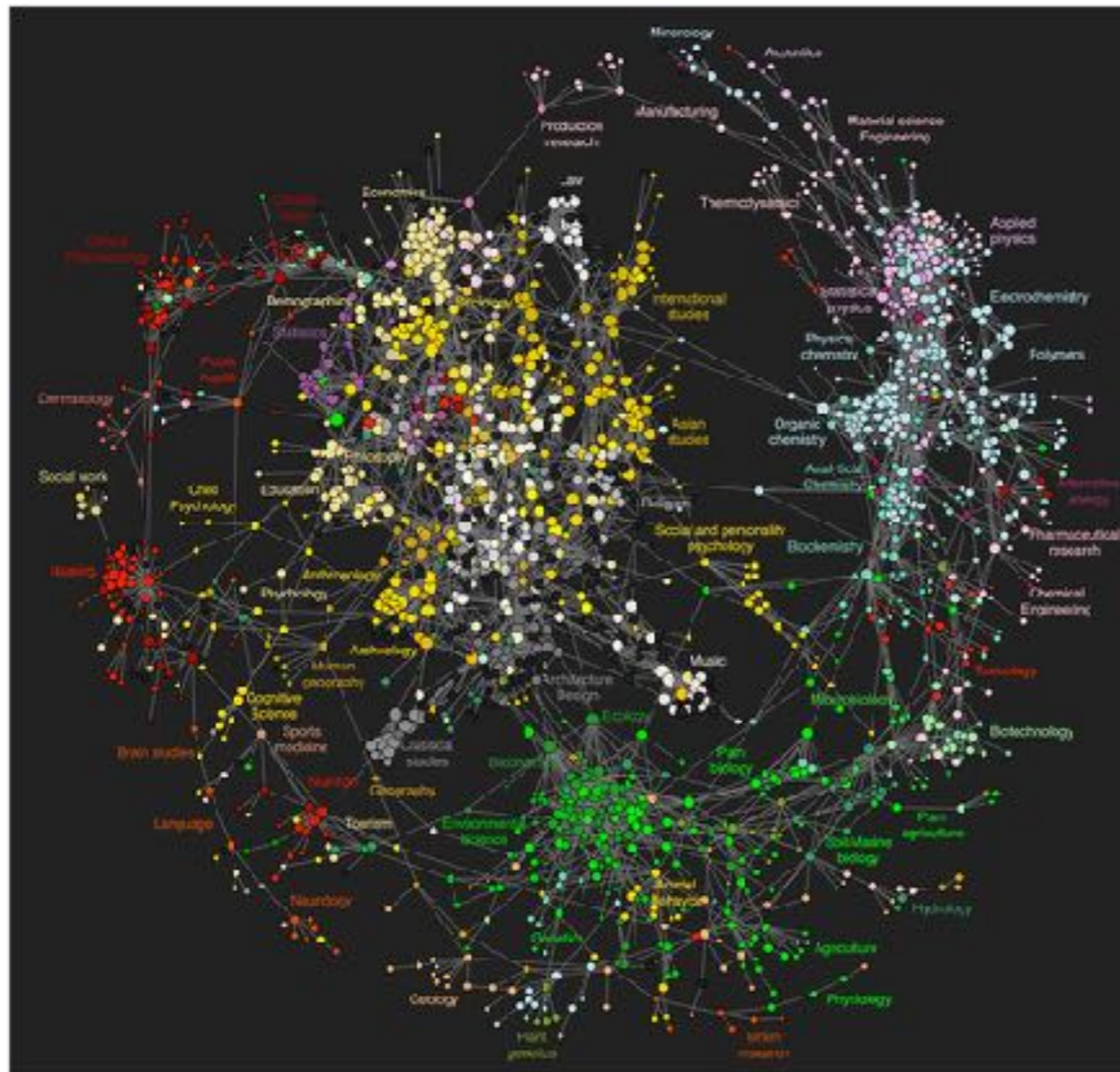
# Community Detection

And if you have personal knowledge of a network, you may be able to spot some groups



# Community Detection

But in many cases, communities are much harder to pick out by eye (or your eyes lie to you)



It can also be difficult to categorize individual nodes.

So we would like to have a more scientific way of dividing the network up...

ref: network of science, Bollen et al (2009)

# Community Detection Algorithms

General idea: create a partition of the nodes, based on where the network “naturally” wants to split

There are lots of ways to do this (we'll look at three):

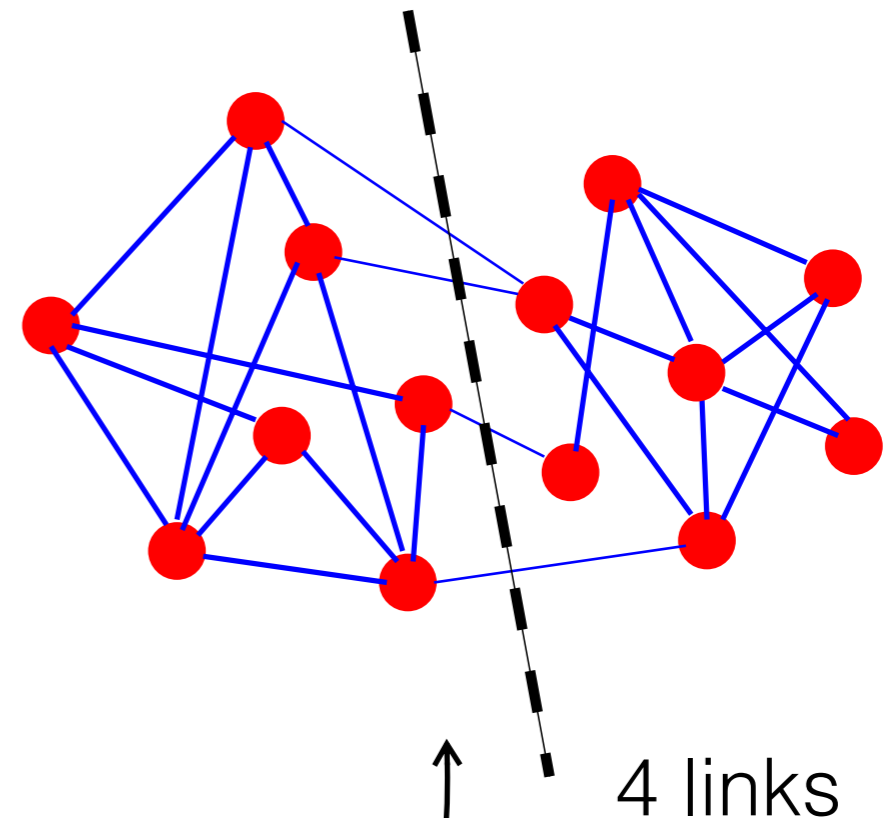
- Graph Partitioning
- Hierarchical Clustering
- Girvan-Newman

Math fact: a *partition* is a division of a set into smaller, non-overlapping sets.

# Method 1: Graph Partitioning

*Graph Partitioning:* divide the network into a pre-defined number of chunks of a pre-defined size

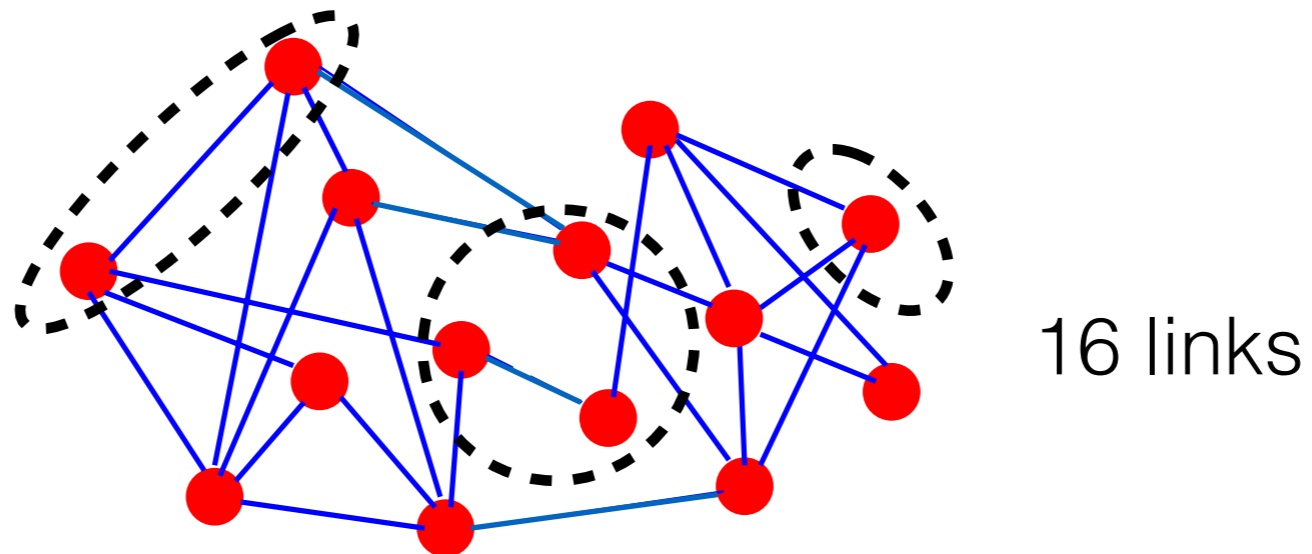
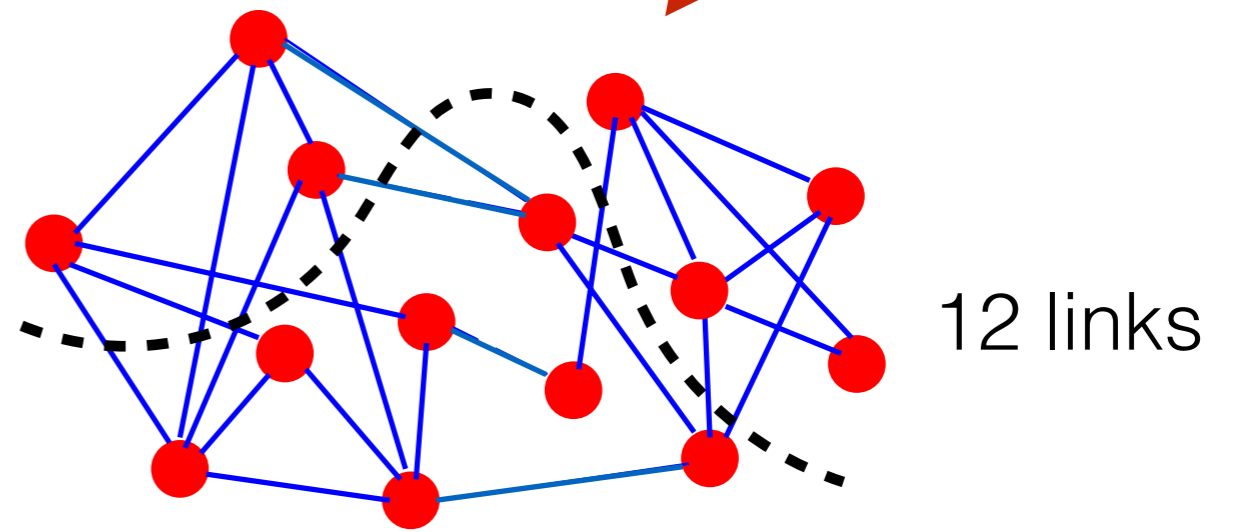
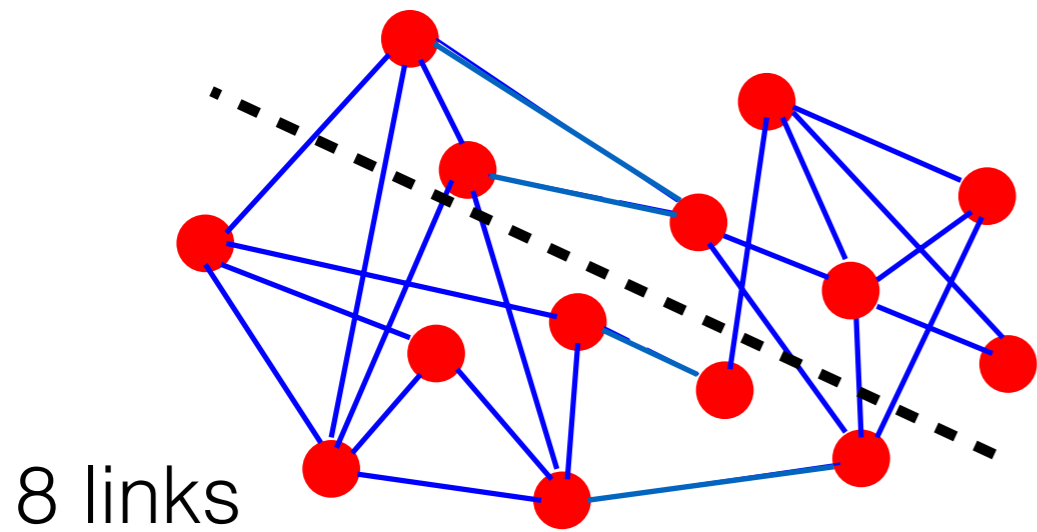
Make the cut in a place that severs the fewest links



a 14 node network:  
divide into two sets of 7

# Method 1: Graph Partitioning

Other divisions require cutting more links

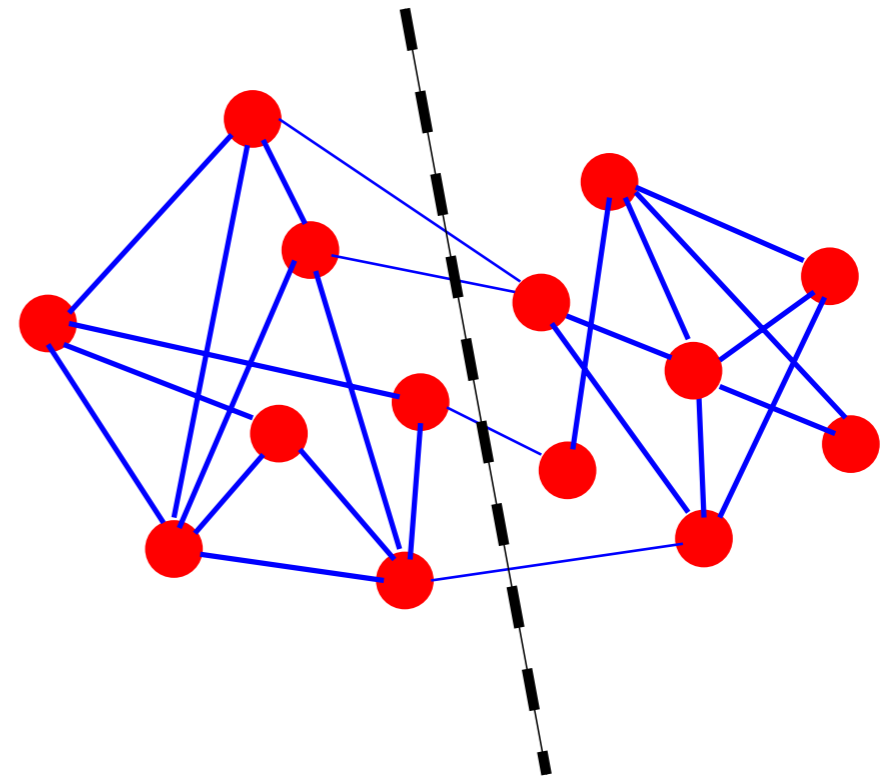


# Community Structure

## Graph Partitioning

Graph partitioning is a very straightforward way to divide the network into communities.

But there is a problem: we need to know how many partitions we want, and how big we want them to be!

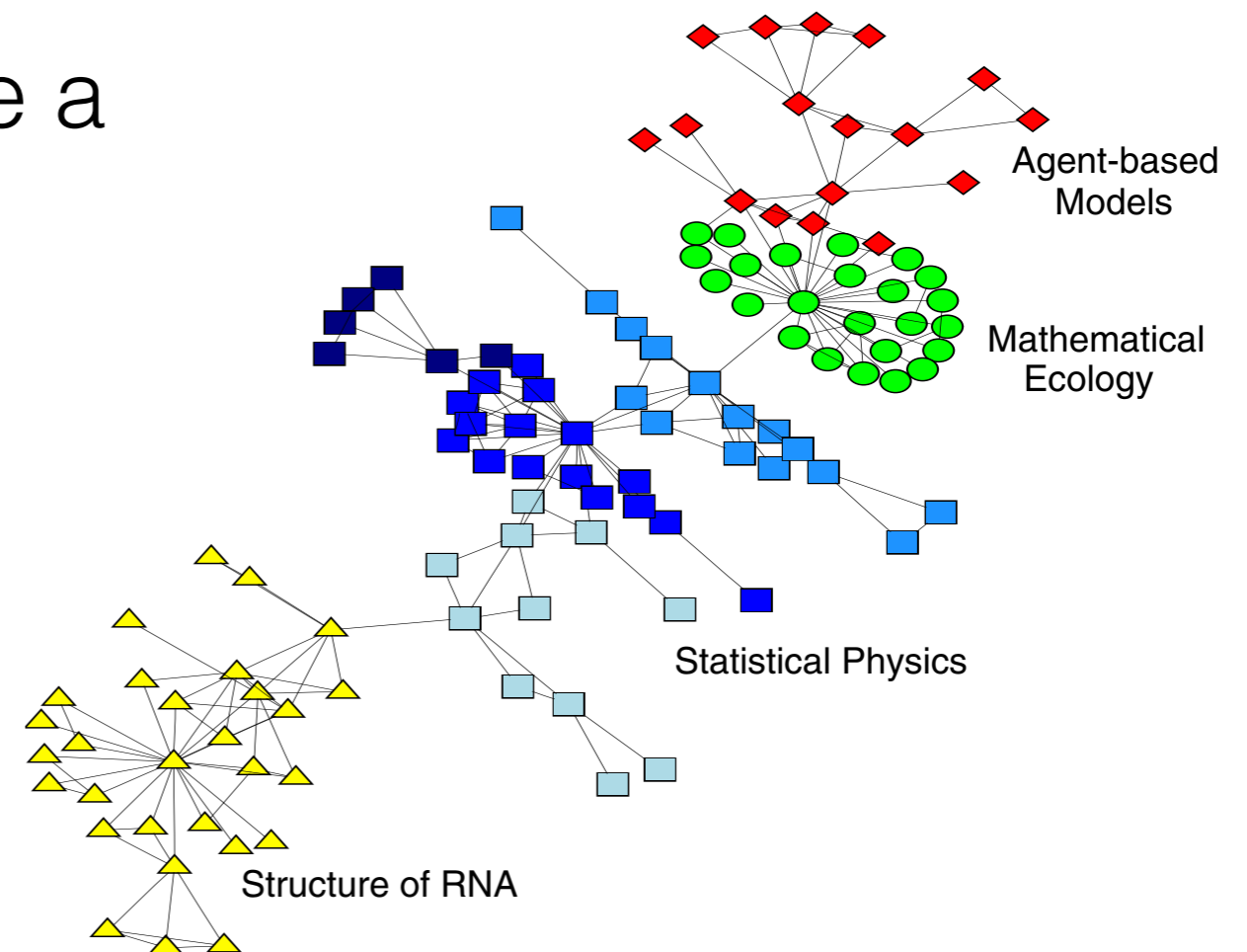


# Community Structure

## Graph Partitioning

In some cases, this may be a reasonable thing to do

But in many cases, the number and size of the communities is exactly what we want to find out...





# Community Structure

## Hierarchical Clustering

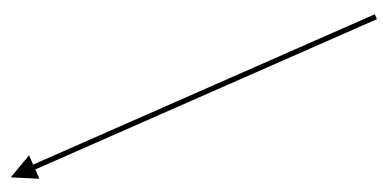
*Hierarchical Clustering* is a method for dividing the network into clusters of sizes determined by the network itself.

# Community Structure

## Hierarchical Clustering

Procedure:

this weight could technically be anything, but probably reflects how closely related the nodes are

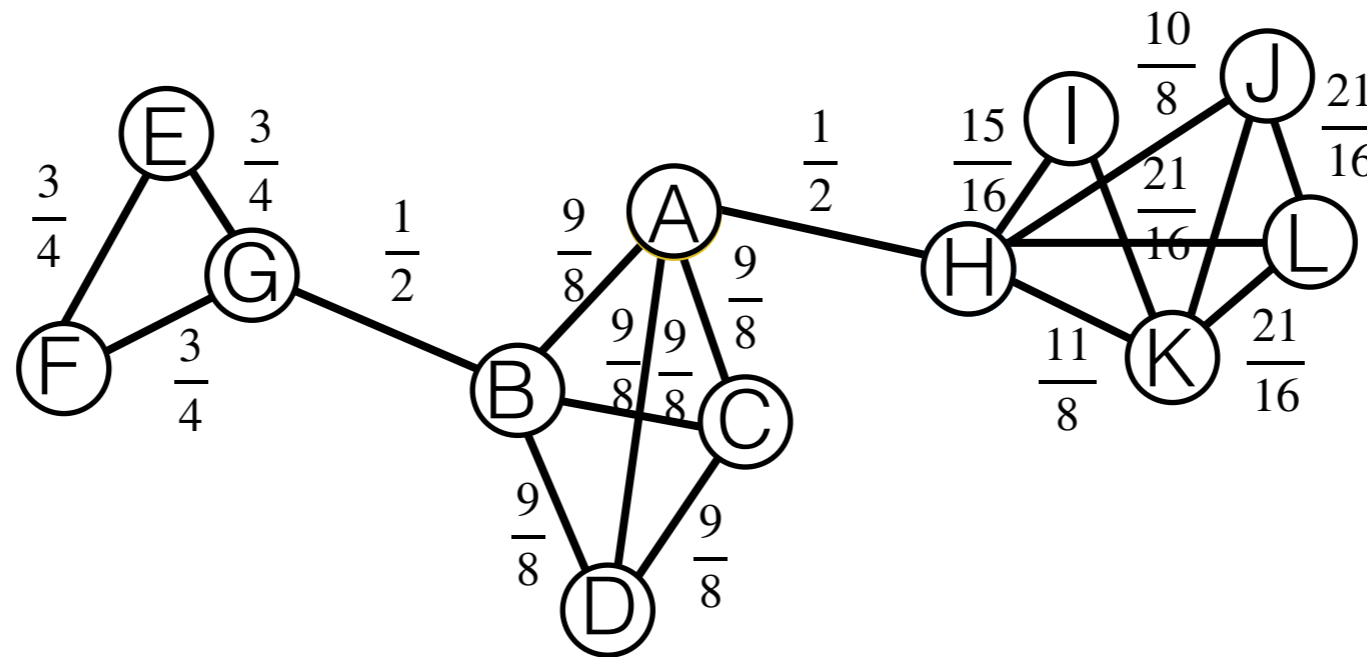


- Assign a weight,  $w_{ij}$ , to each pair of nodes in the network
- Remove all of the edges in the network.
- Reconnect the nodes, starting with the edge that has the highest weight
- As edges are added, the network is connected back together (it may not be the same way it was before, but that's fine)

# Community Structure

## Hierarchical Clustering

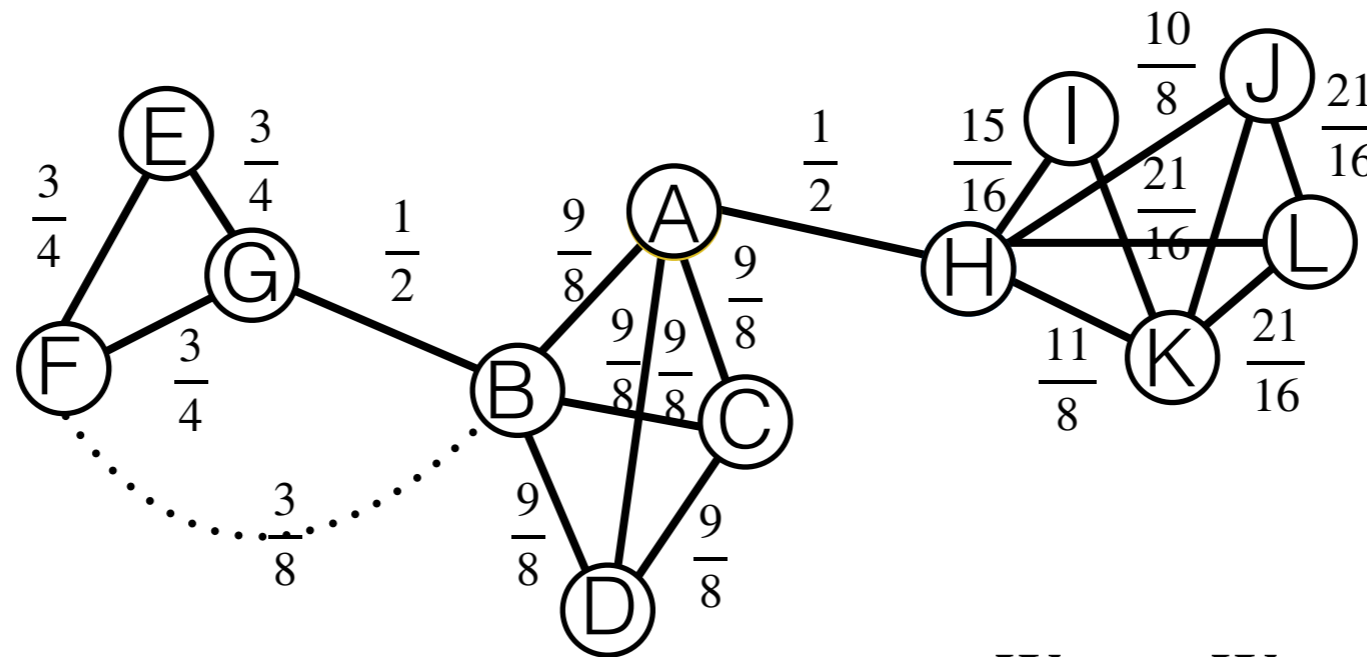
Weights = number of paths between nodes, weighted by length of path



# Community Structure

## Hierarchical Clustering

Weights = number of paths between nodes,  
weighted by length of path



$$W_{AH} = W_{AG} = 1$$

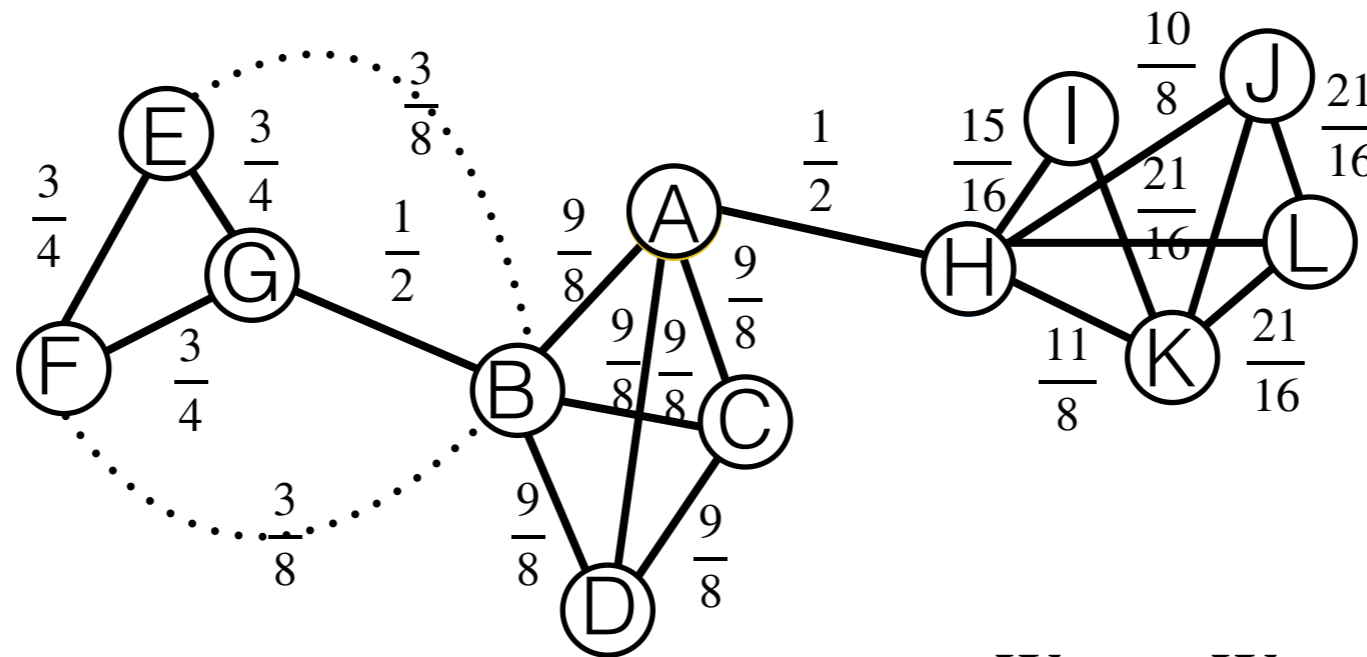
$$W_{AB} = 3$$

$$W_{HK} = 4$$

# Community Structure

## Hierarchical Clustering

Weights = number of paths between nodes, weighted by length of path



$$W_{AH} = W_{AG} = 1$$

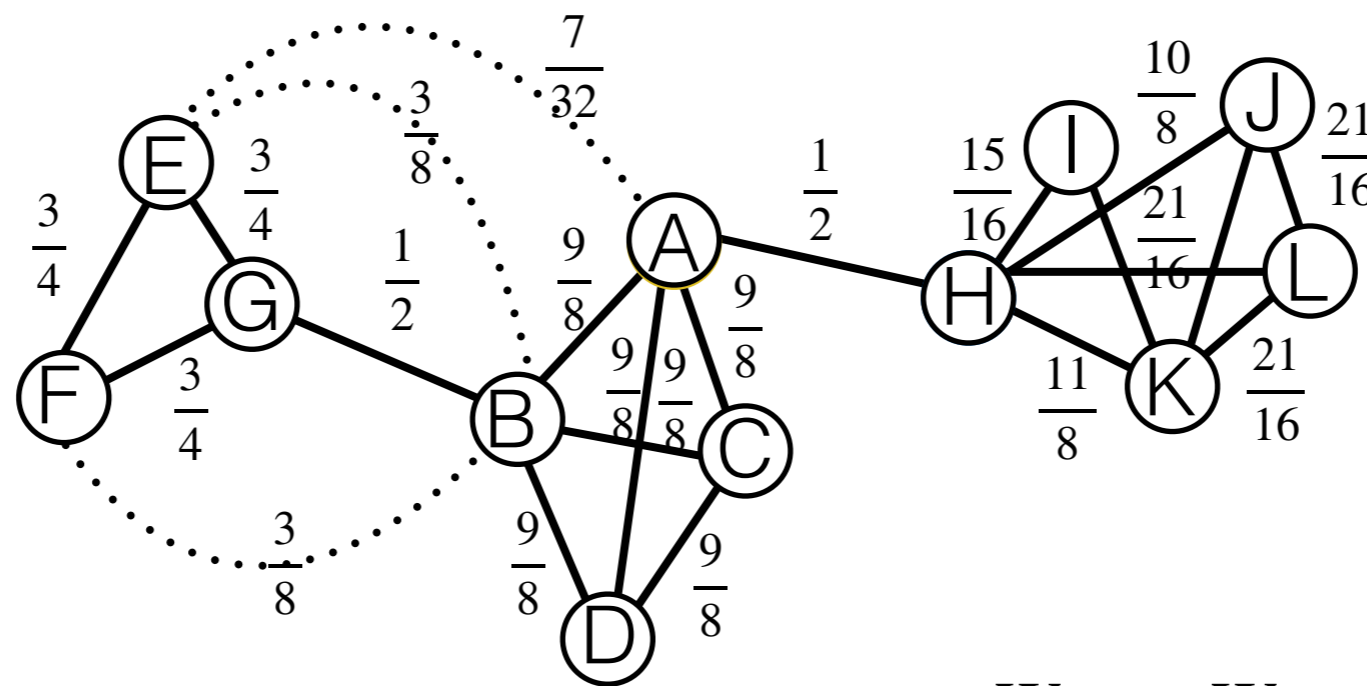
$$W_{AB} = 3$$

$$W_{HK} = 4$$

# Community Structure

## Hierarchical Clustering

Weights = number of paths between nodes,  
weighted by length of path



$$W_{AH} = W_{AG} = 1$$

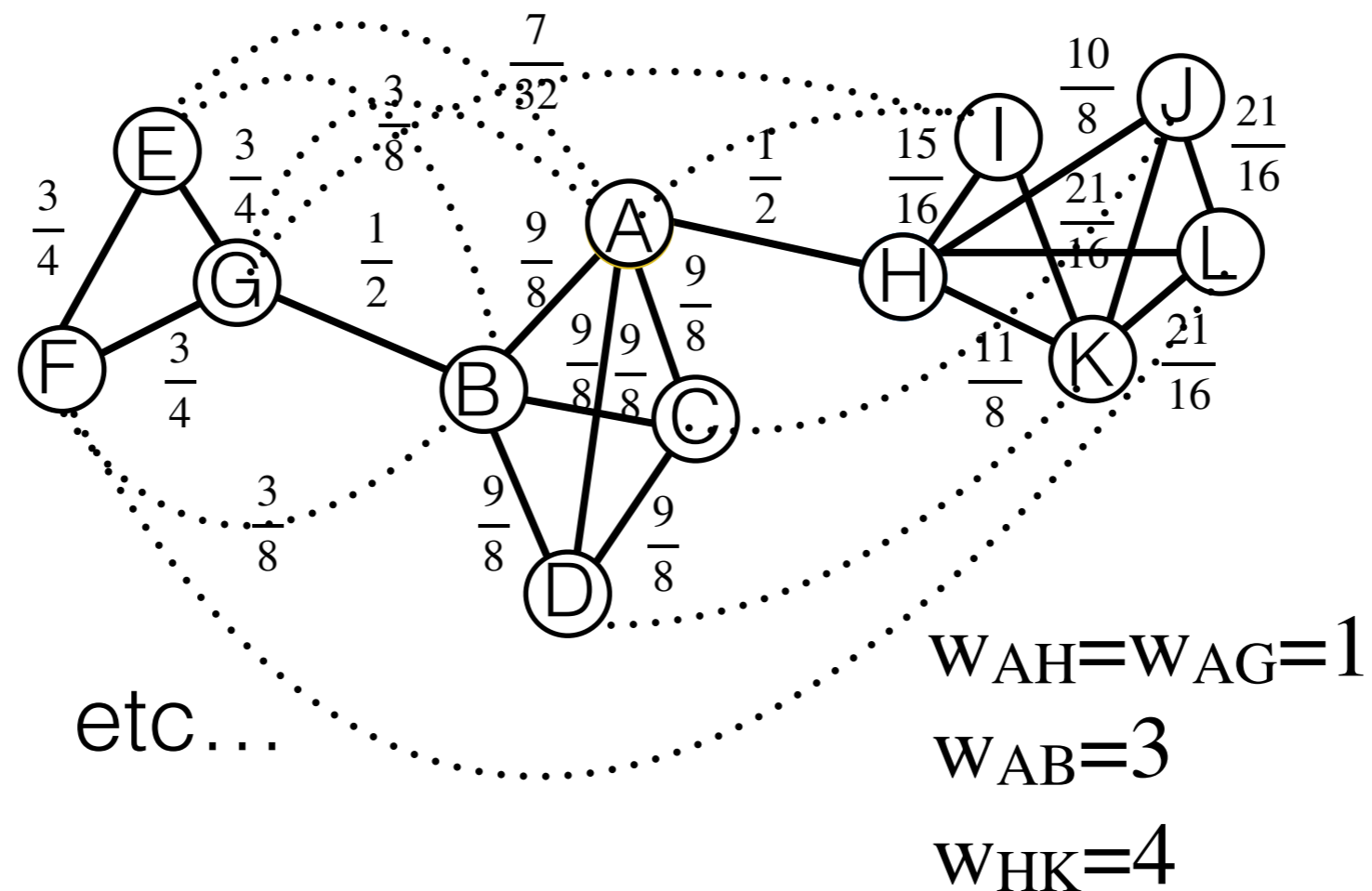
$$W_{AB} = 3$$

$$W_{HK} = 4$$

# Community Structure

## Hierarchical Clustering

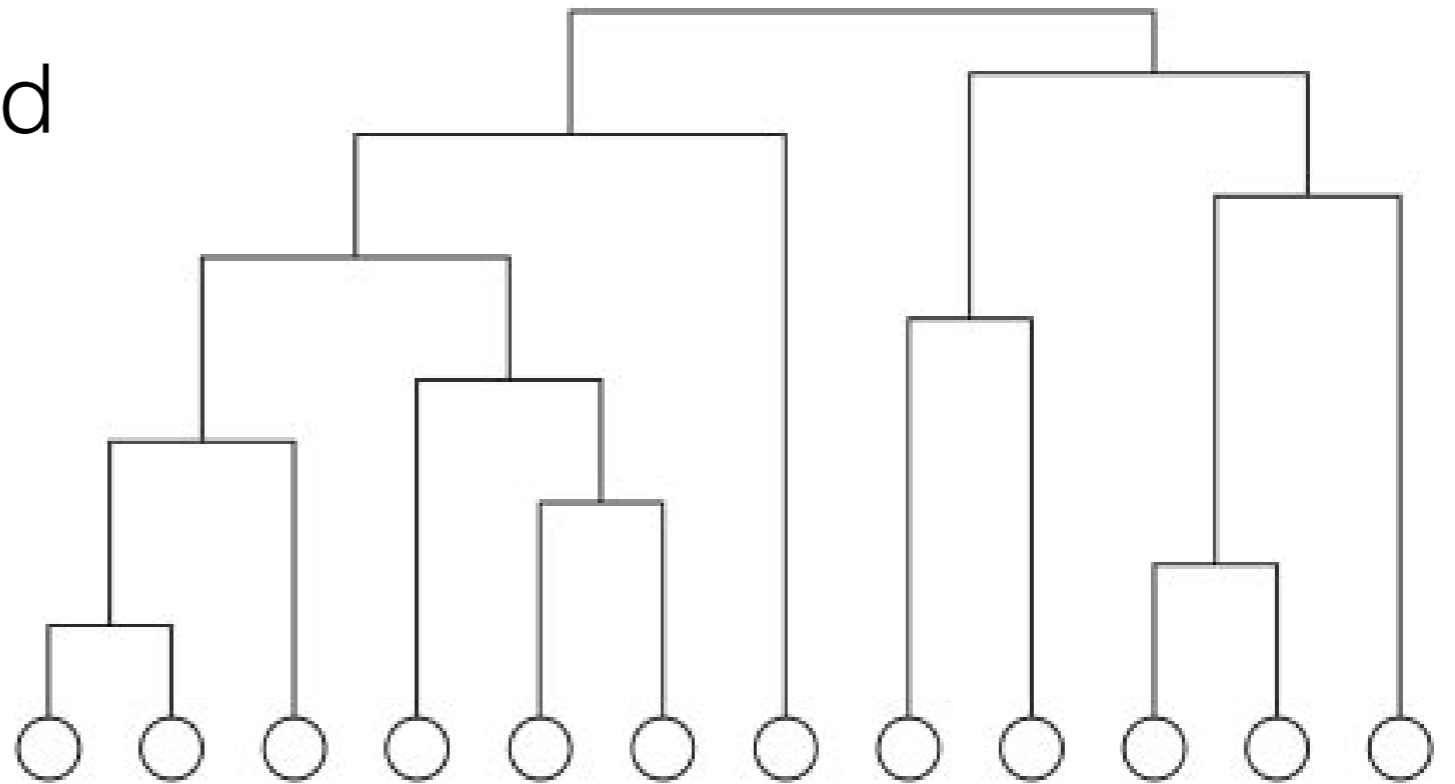
Weights = number of paths between nodes,  
weighted by length of path



# Community Structure

## Hierarchical Clustering

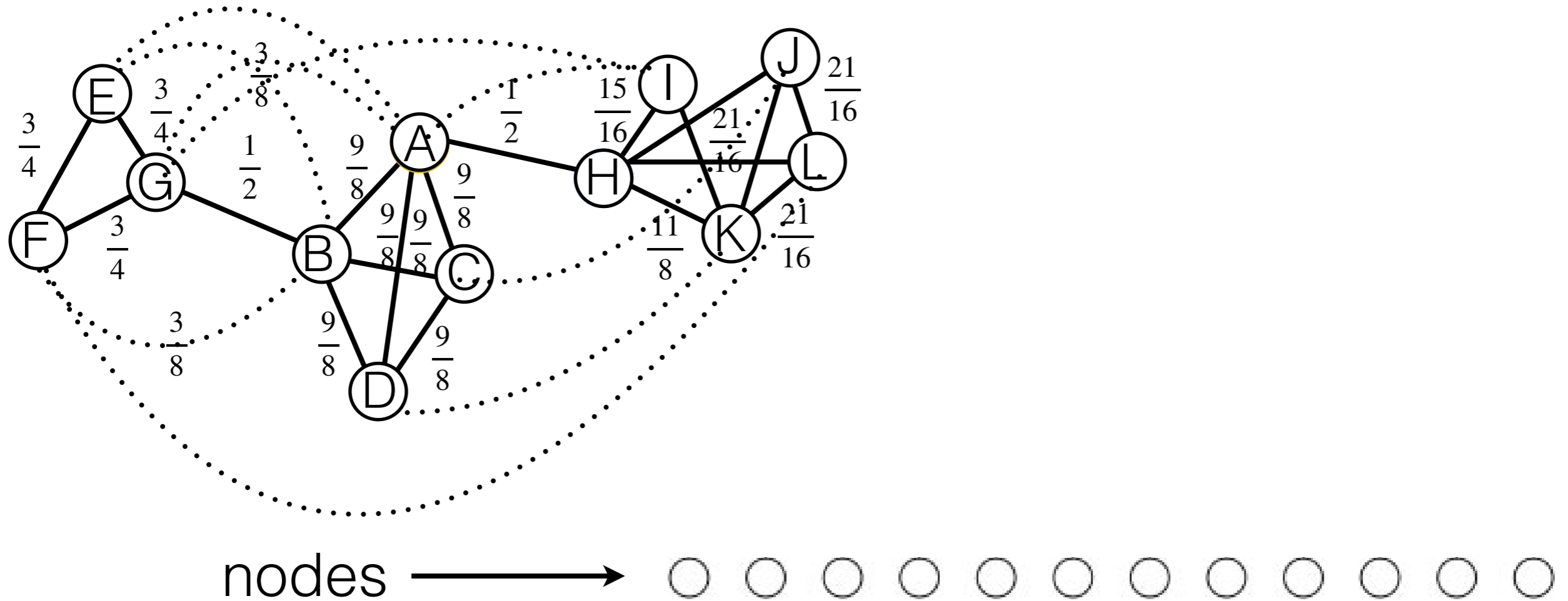
The result of this process is summarized by a *dendrogram*





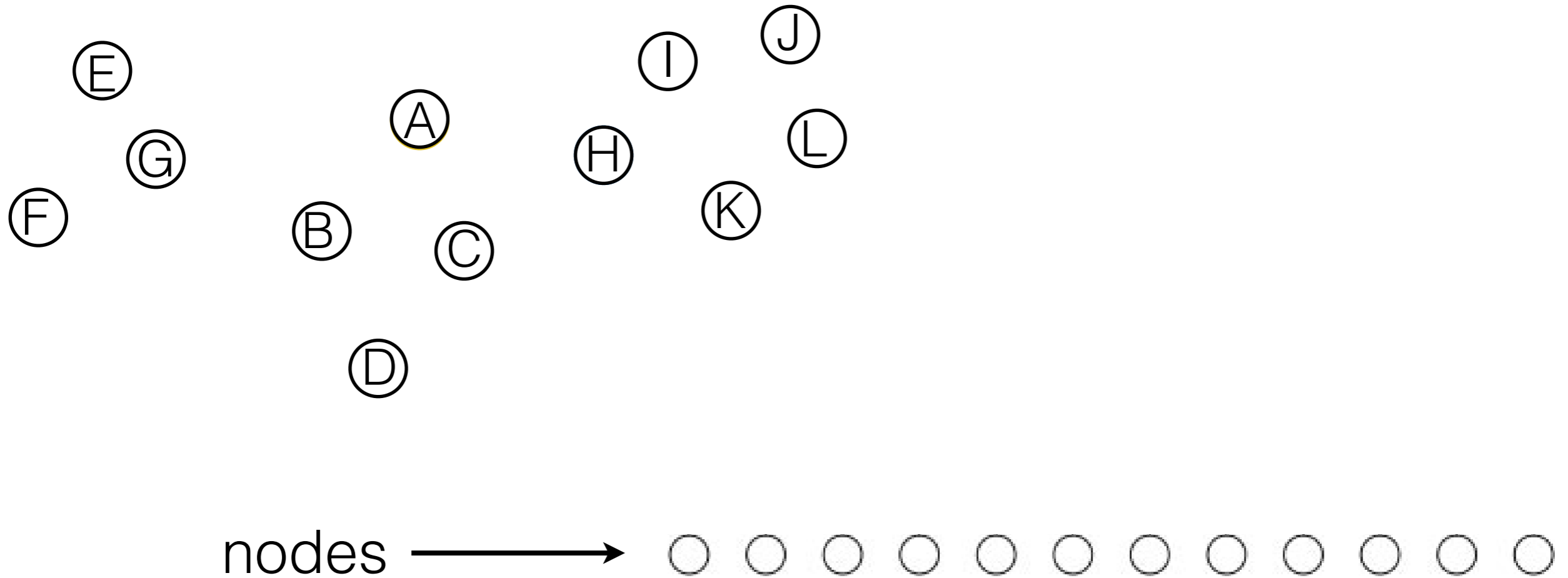
# Community Structure

## Hierarchical Clustering



# Community Structure

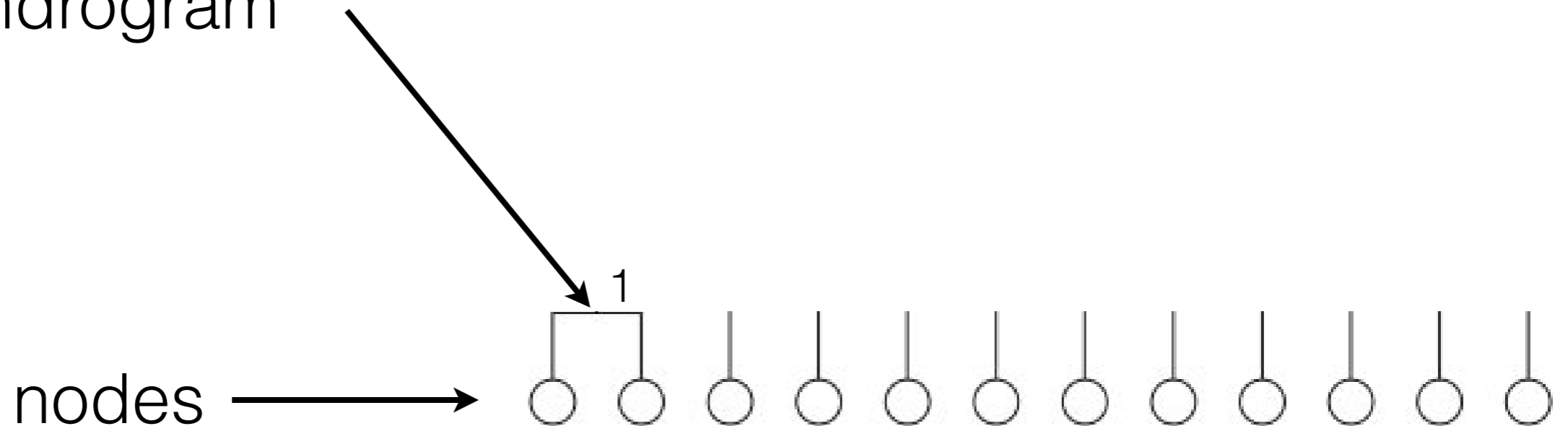
## Hierarchical Clustering



# Community Structure

## Hierarchical Clustering

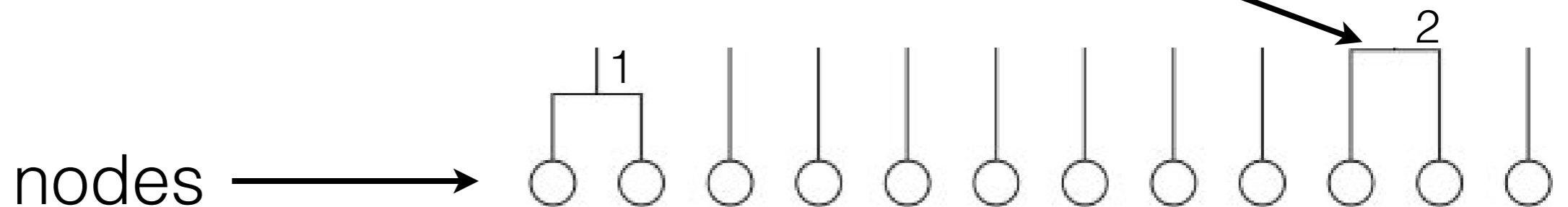
Each edge that is added is represented by a connection in the dendrogram



# Community Structure

## Hierarchical Clustering

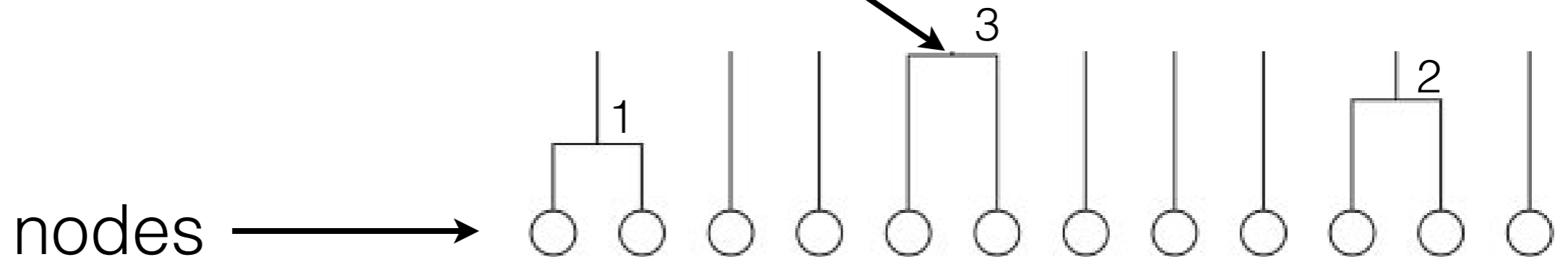
As more edges are added,  
the network becomes more  
connected



# Community Structure

## Hierarchical Clustering

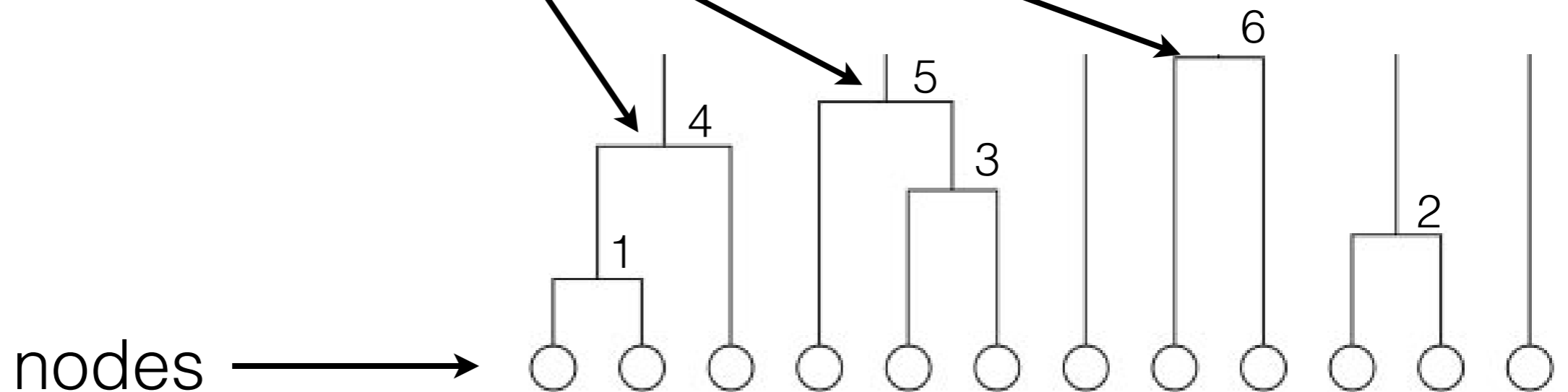
As more edges are added,  
the network becomes more  
connected



# Community Structure

## Hierarchical Clustering

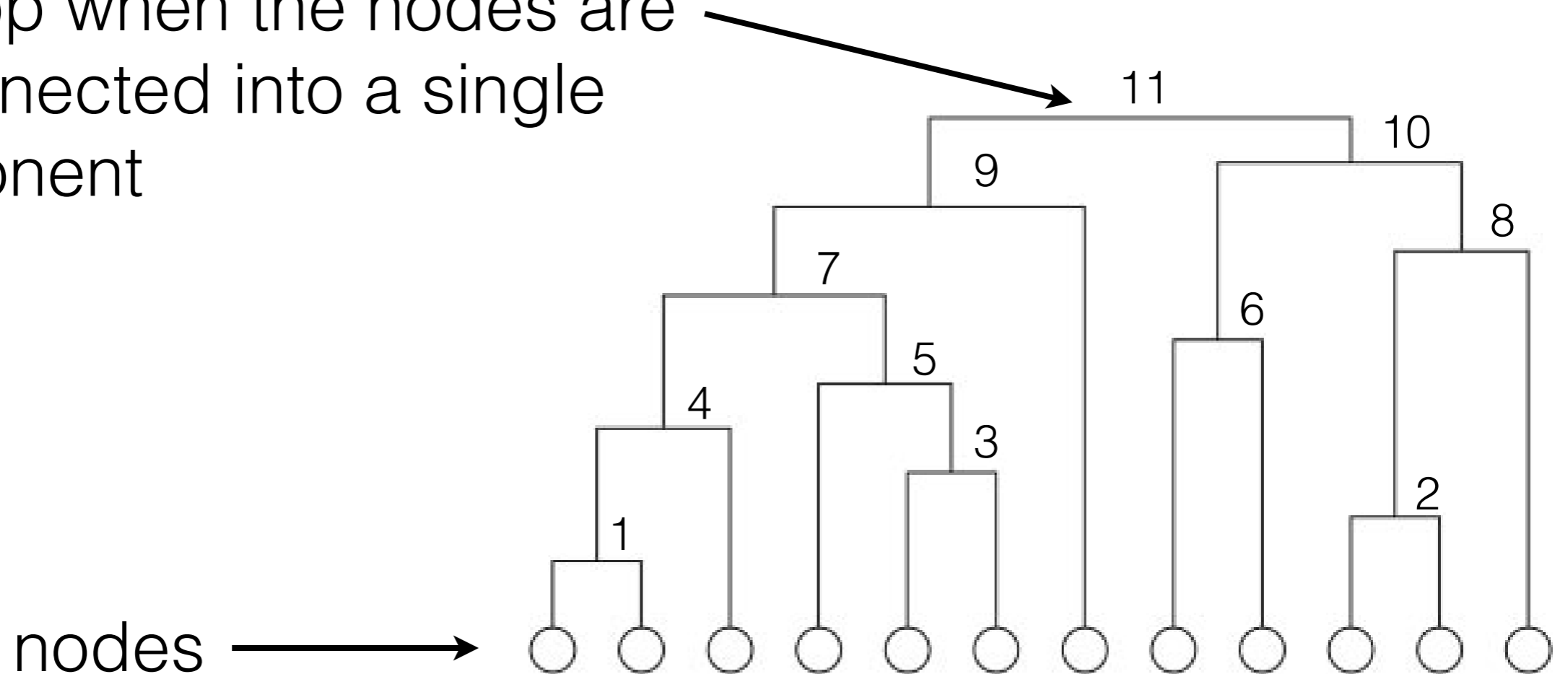
As more edges are added,  
the network becomes more  
connected



# Community Structure

## Hierarchical Clustering

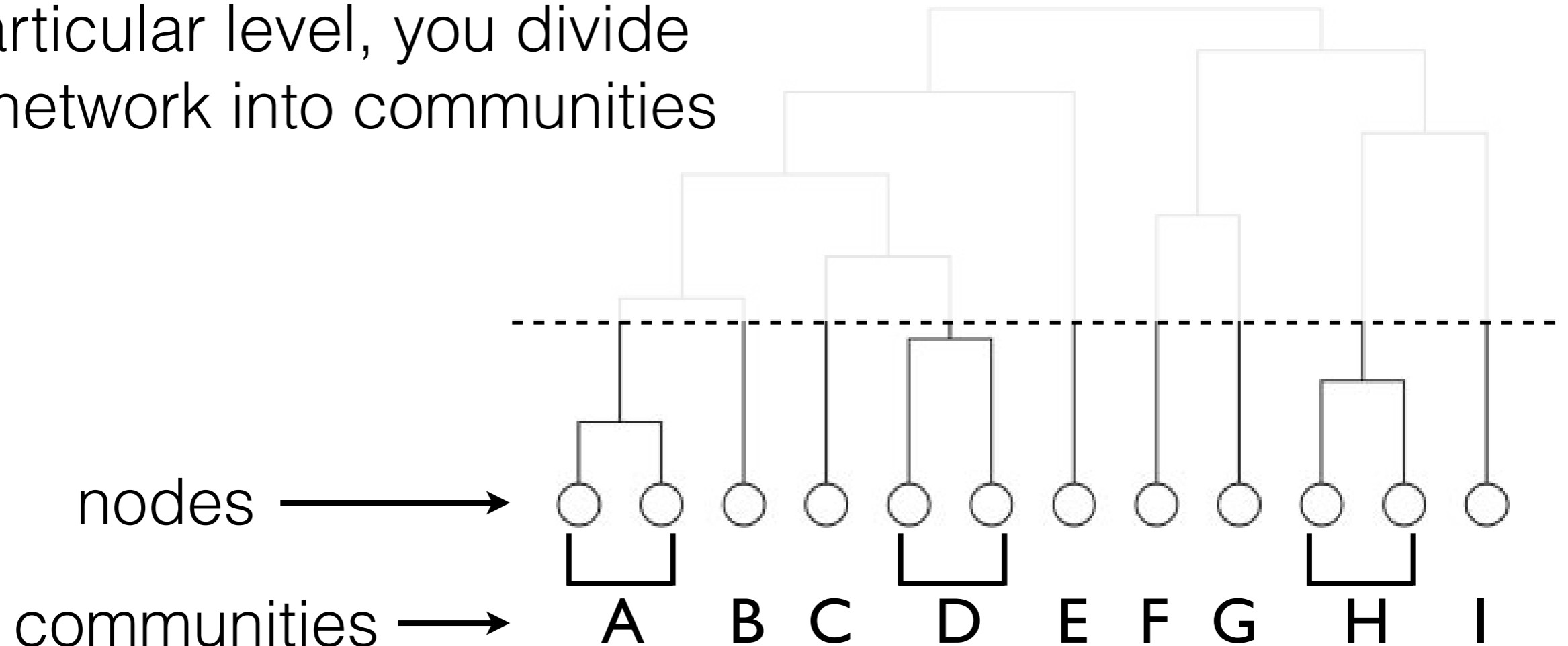
We stop when the nodes are all connected into a single component



# Community Structure

## Hierarchical Clustering

By cutting the dendrogram at a particular level, you divide the network into communities

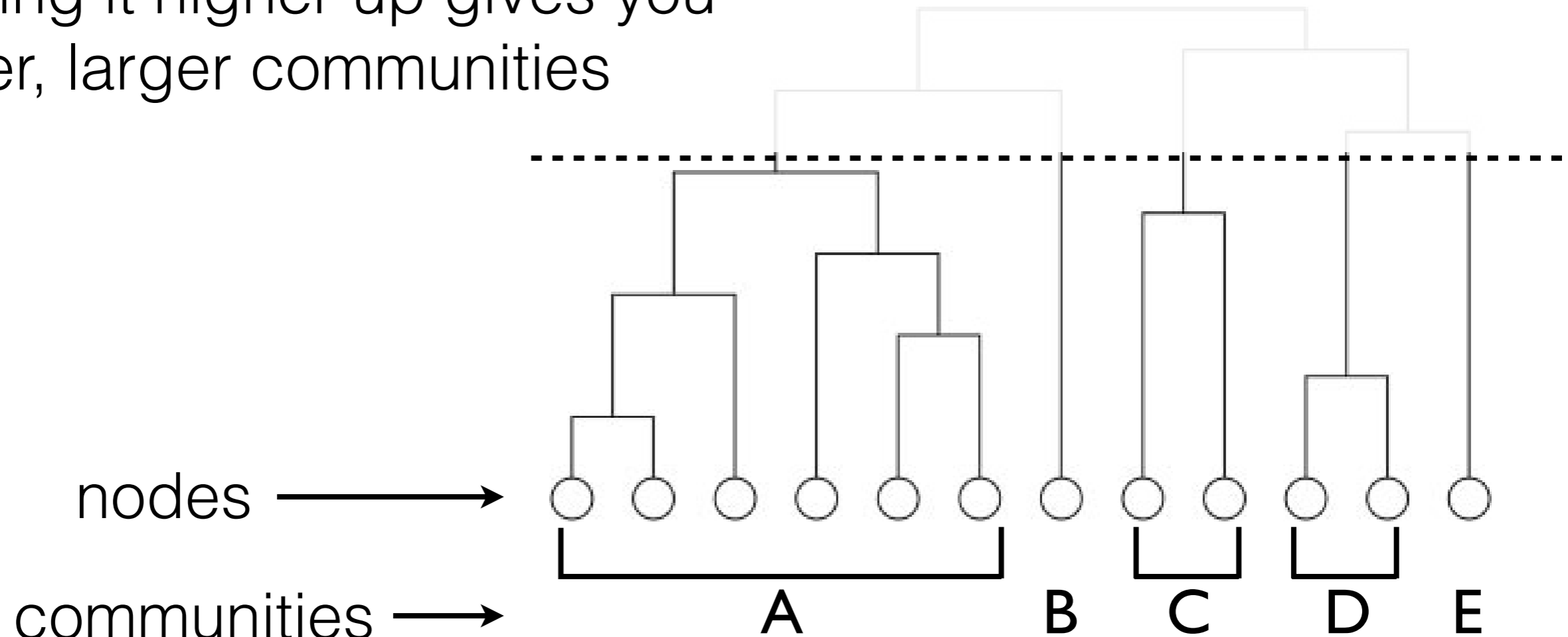




# Community Structure

## Hierarchical Clustering

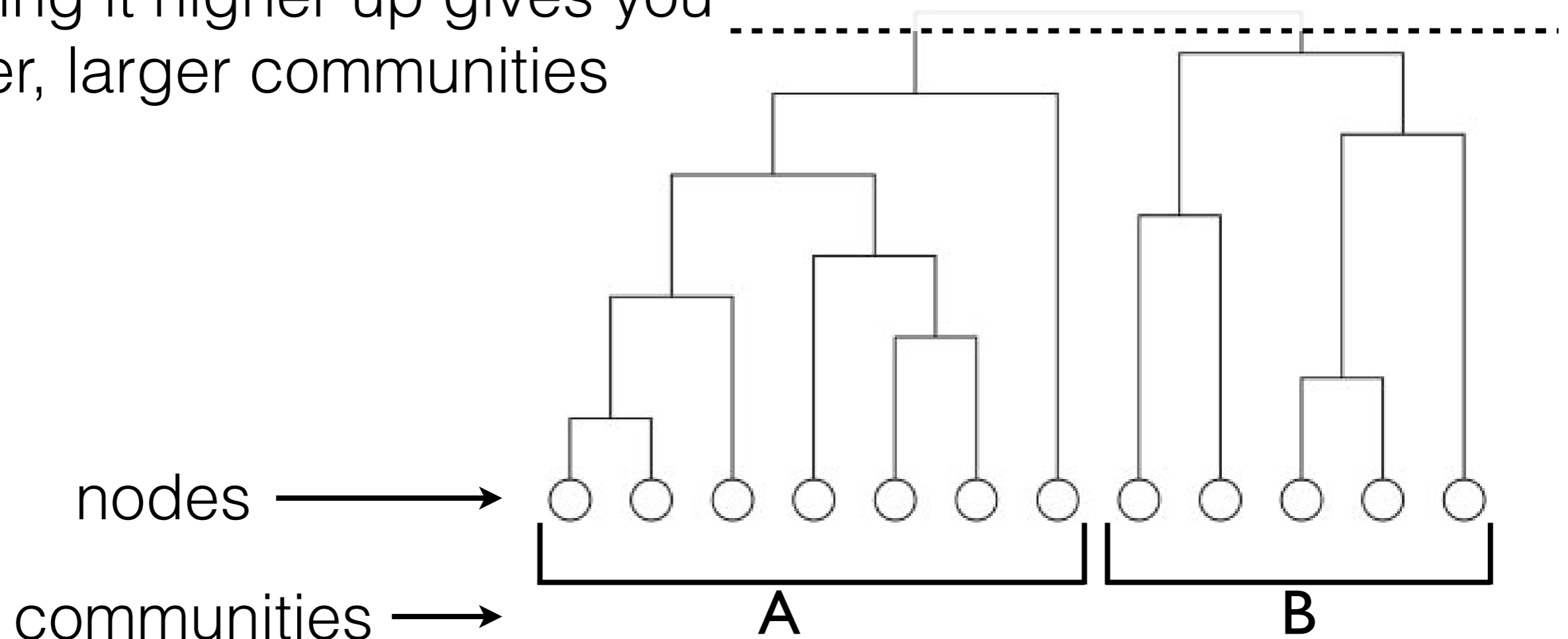
Cutting it higher up gives you fewer, larger communities



# Community Structure

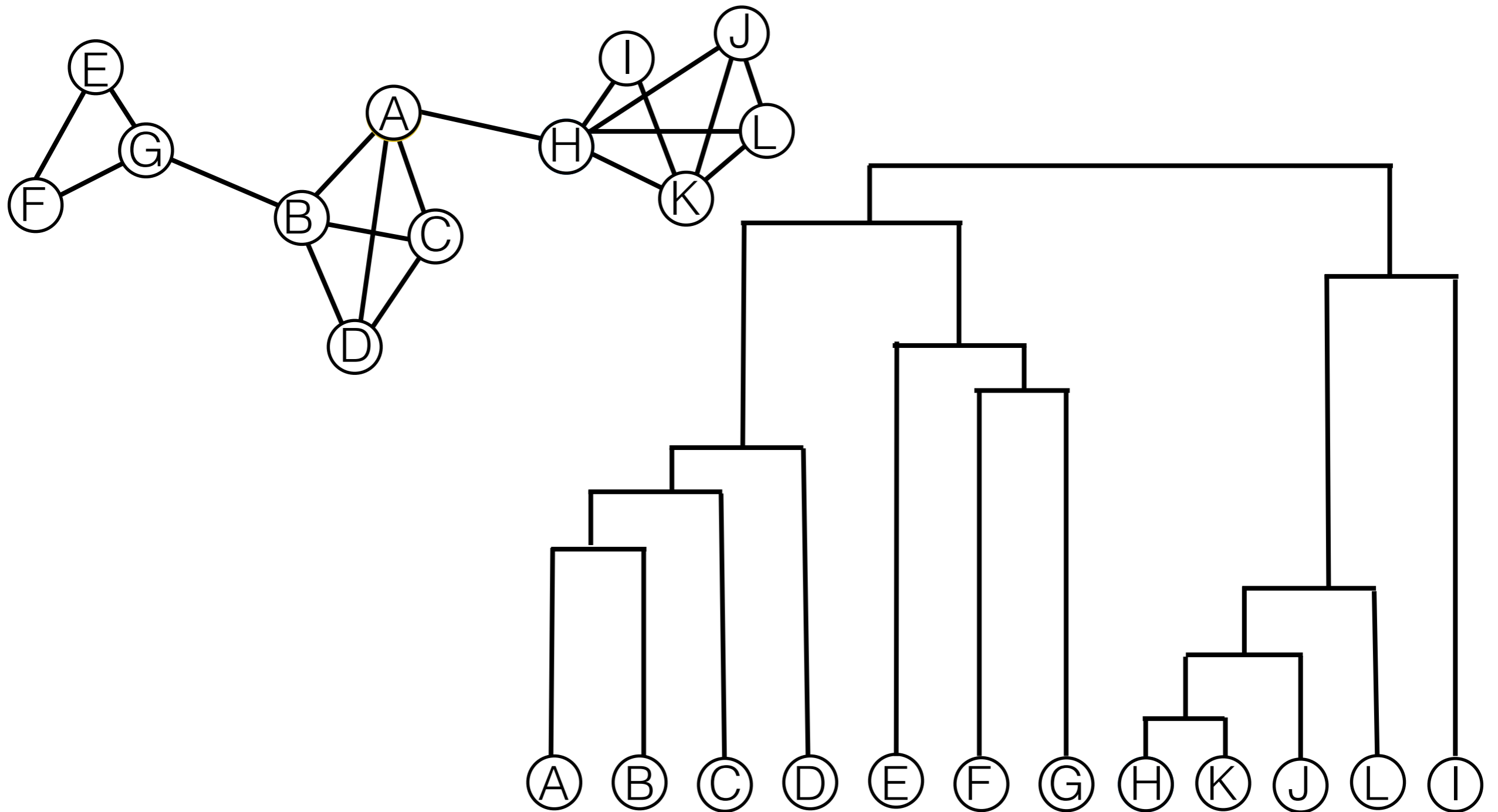
## Hierarchical Clustering

Cutting it higher up gives you fewer, larger communities



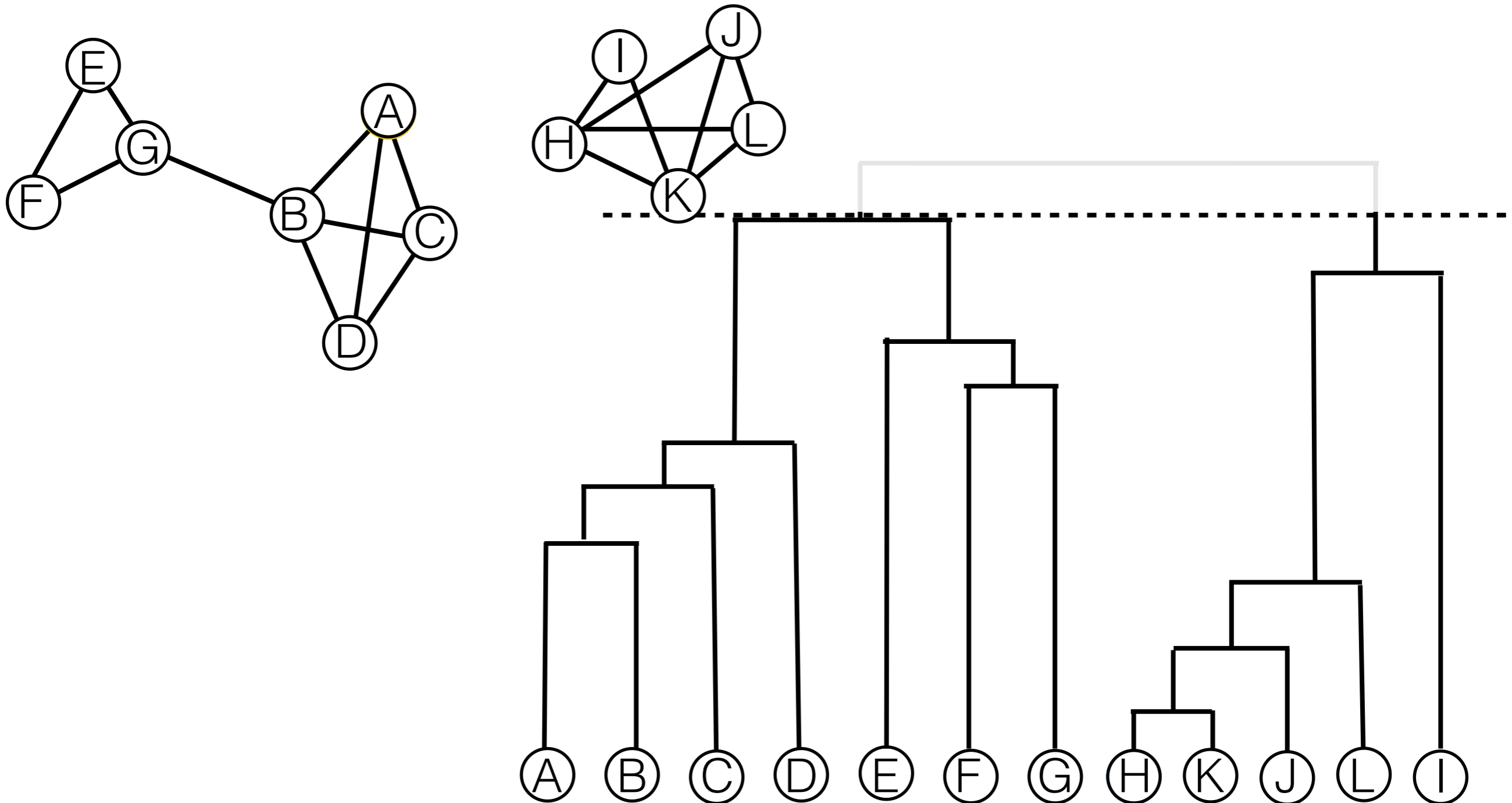
# Community Structure

## Hierarchical Clustering



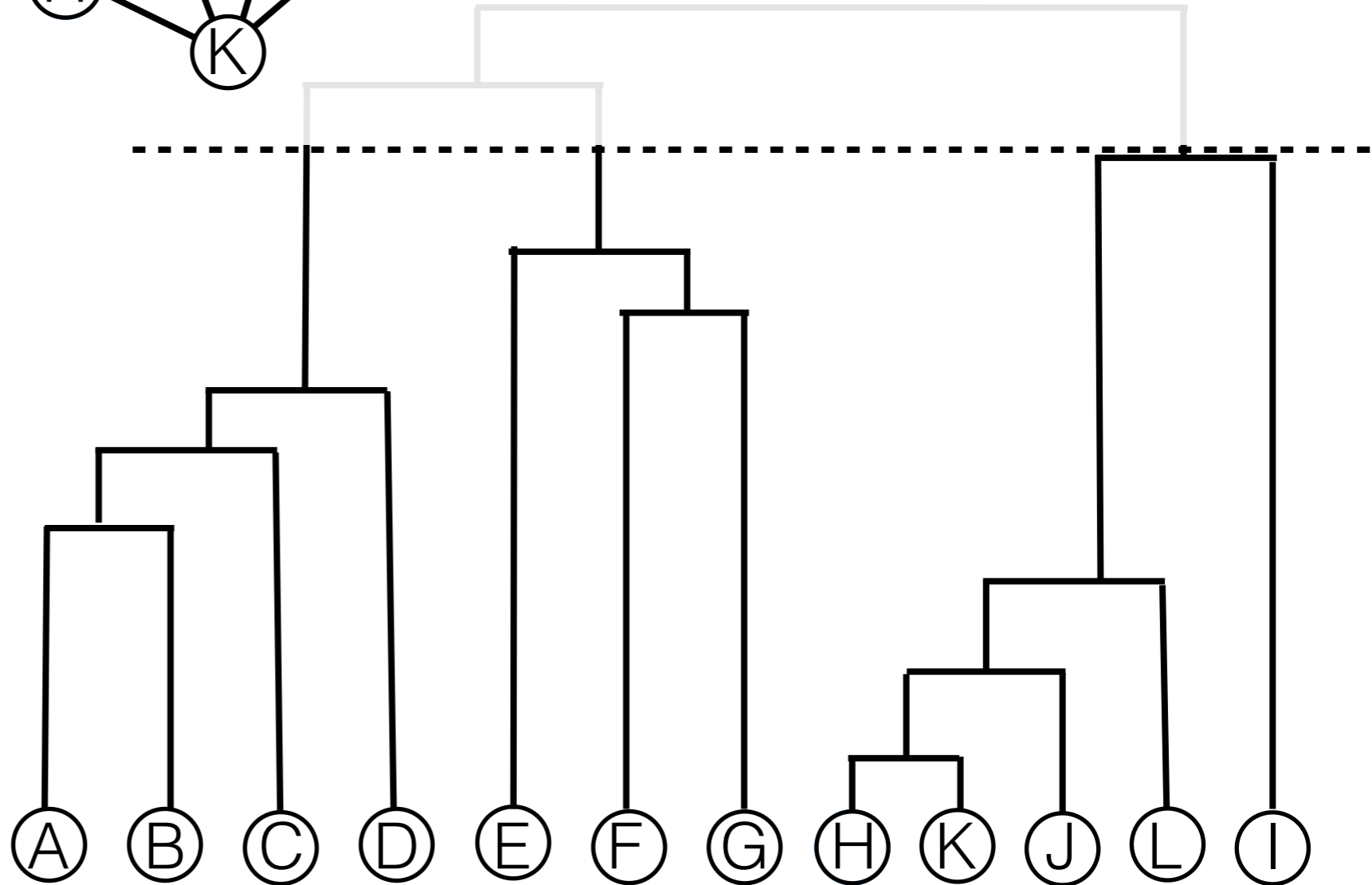
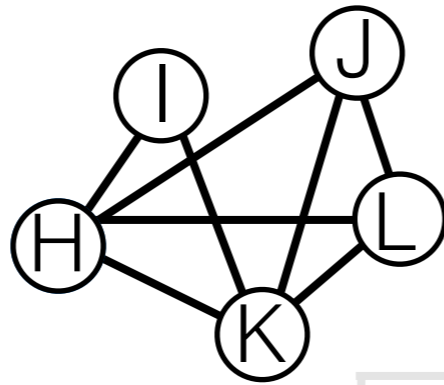
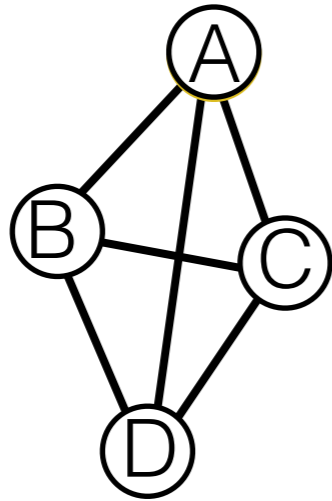
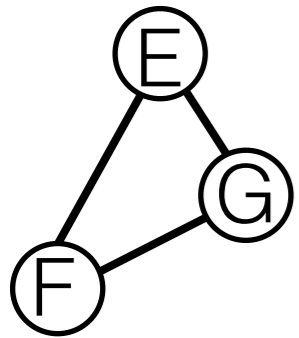
# Community Structure

## Hierarchical Clustering



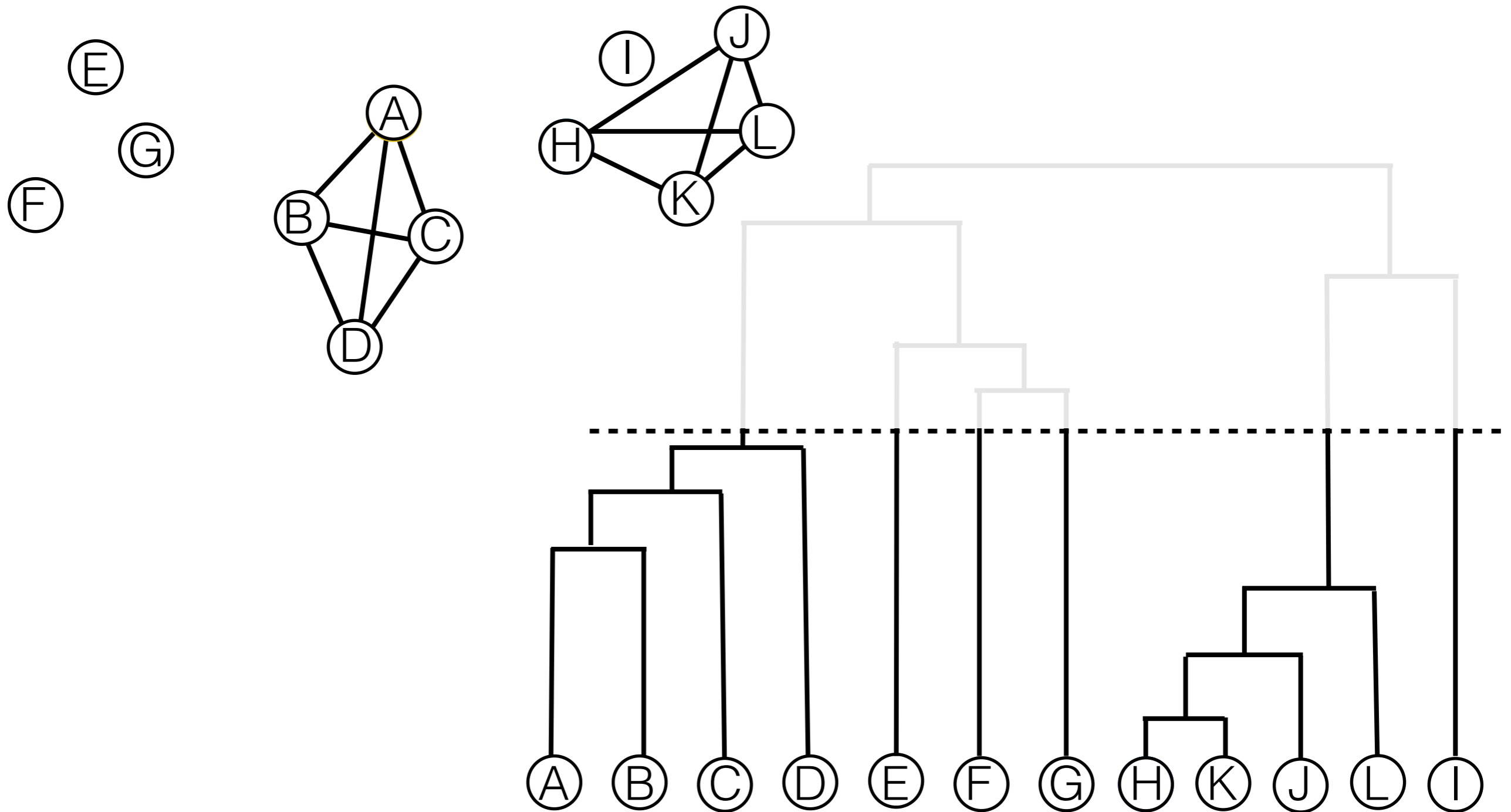
# Community Structure

## Hierarchical Clustering



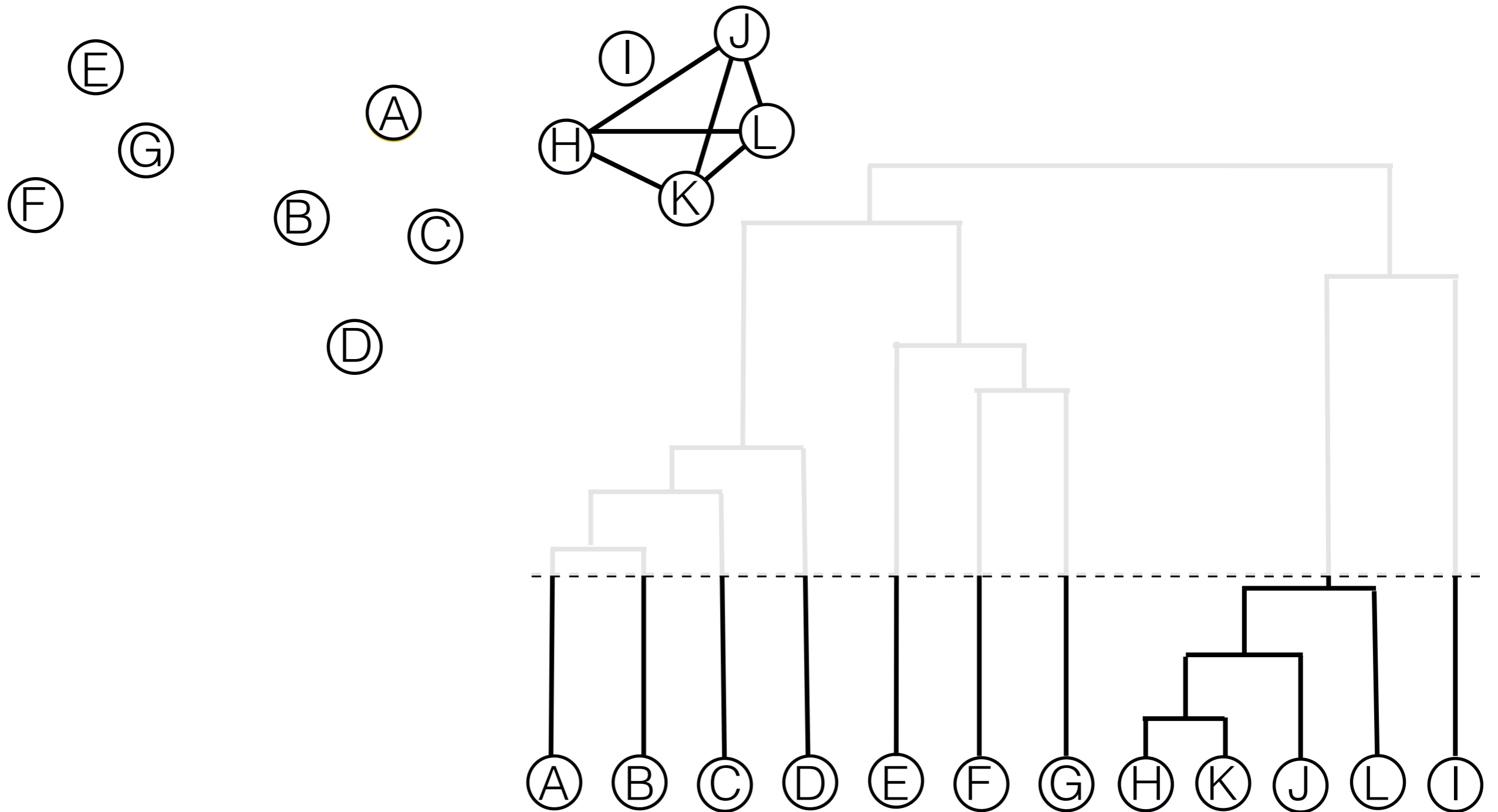
# Community Structure

## Hierarchical Clustering



# Community Structure

## Hierarchical Clustering

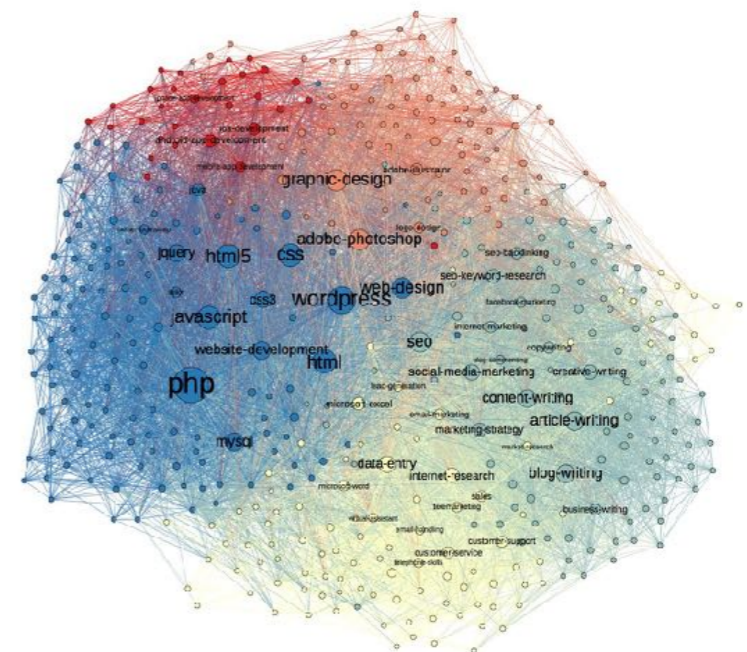
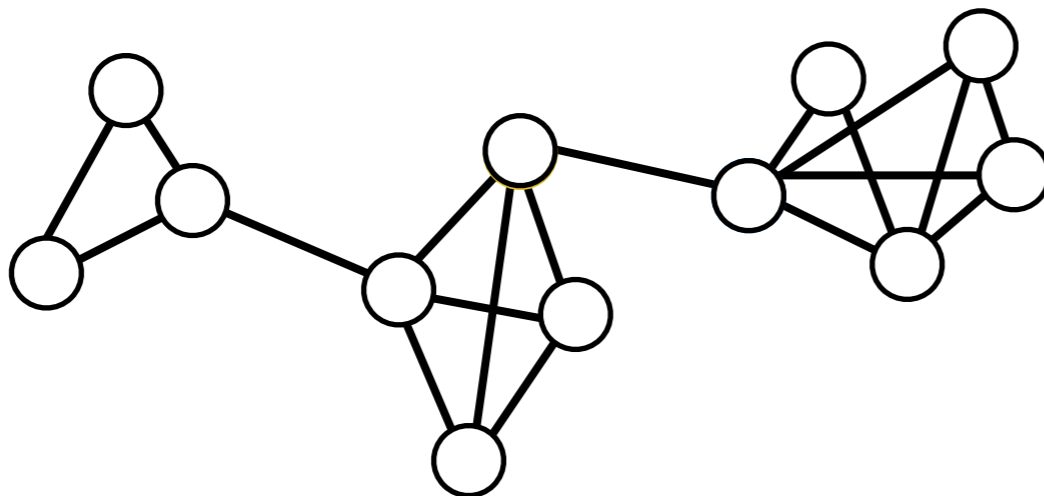


# Community Structure

## Hierarchical Clustering

There are two disadvantages to hierarchical clustering:

1. It tends to chop off “leaf” nodes that are peripheral to a community
2. It works best on networks that have a naturally hierarchical (nested) structure (which is not all networks)



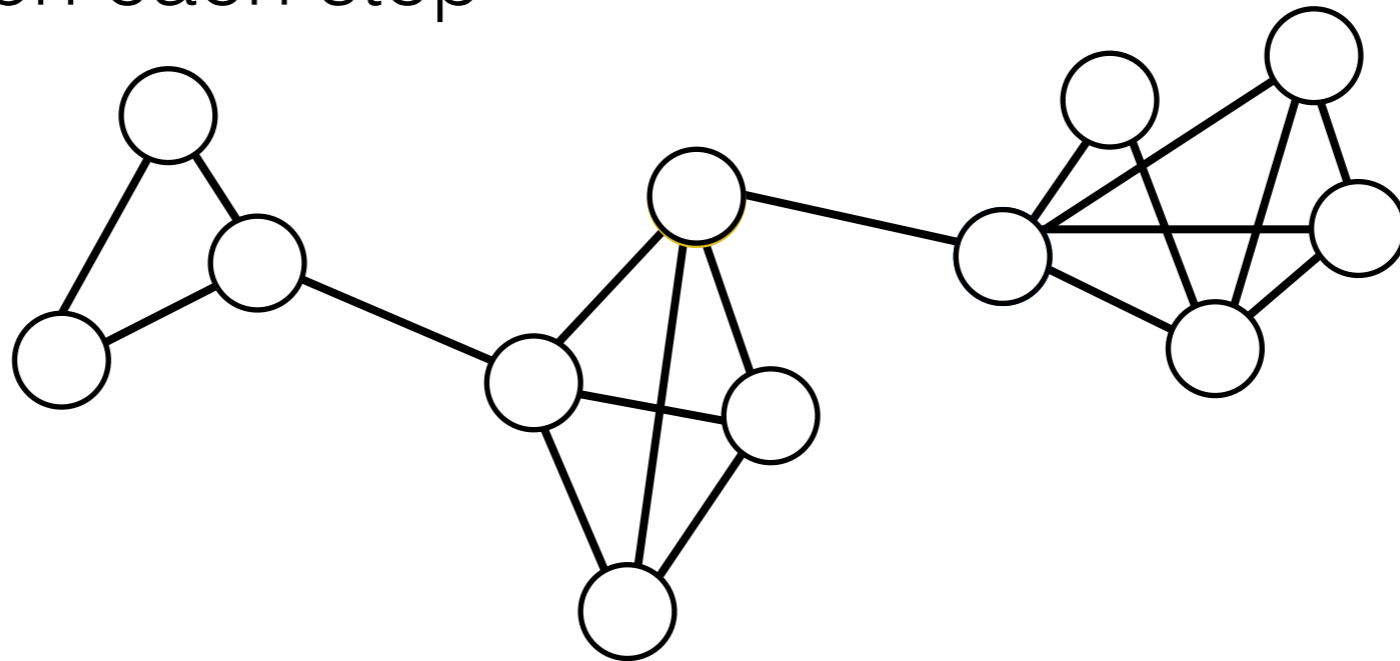


# Community Structure

## Girvan-Newman

The Girvan-Newman Algorithm sequentially removes edges with the highest *edge betweenness*

But unlike hierarchical clustering, it recalculates the edge betweenness on each step



*Edge betweenness*: The number of shortest paths that go along a particular edge

# Community Structure

## Girvan-Newman

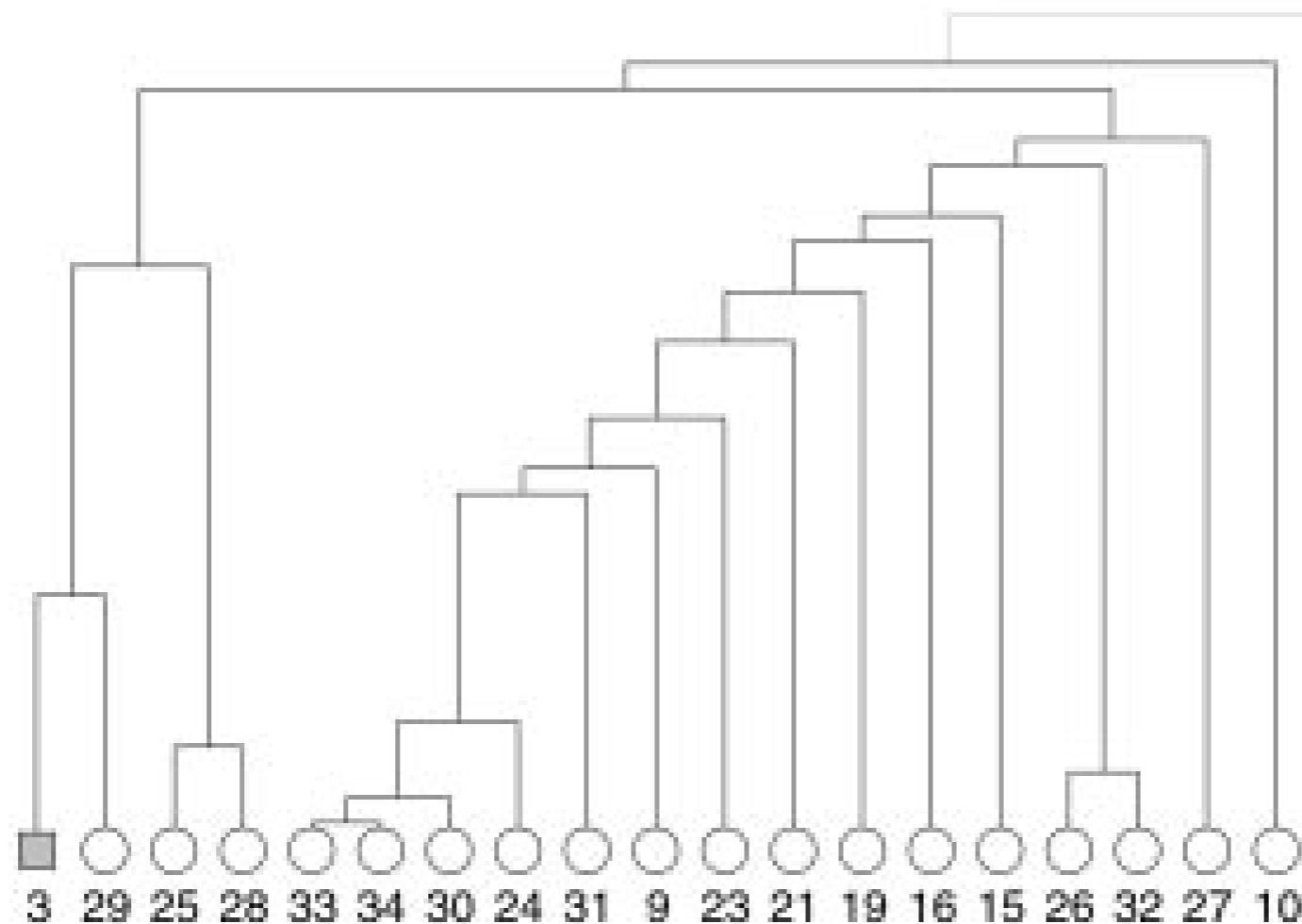
Procedure:

1. Calculate betweenness for all edges
2. Remove the edge with the highest betweenness
3. Recalculate the betweenness of all remaining edges
4. Repeat until no edges remain

# Community Structure

## Girvan-Newman

The result is, again, a dendrogram, which we can cut at different levels to produce different partitions of the network

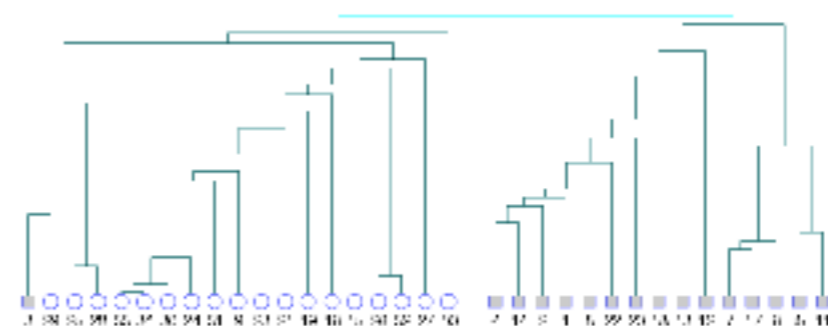
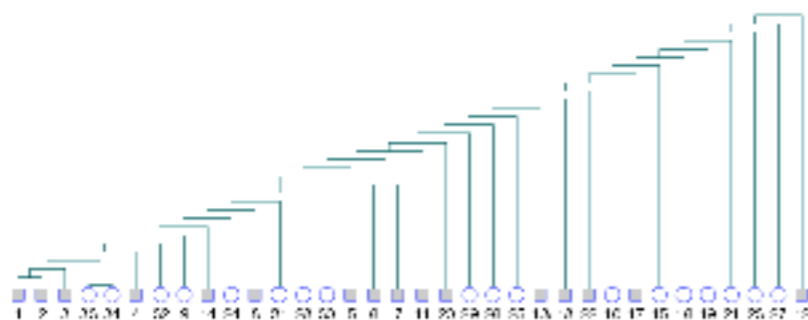


# Community Structure Evaluation

Now we have two algorithms, producing two different community structures. How do we tell which algorithm is best?

→ Answer: there is no definitive answer!

However, there are some tests we can perform that give some insight...



# Evaluating Algorithms

## Random Networks

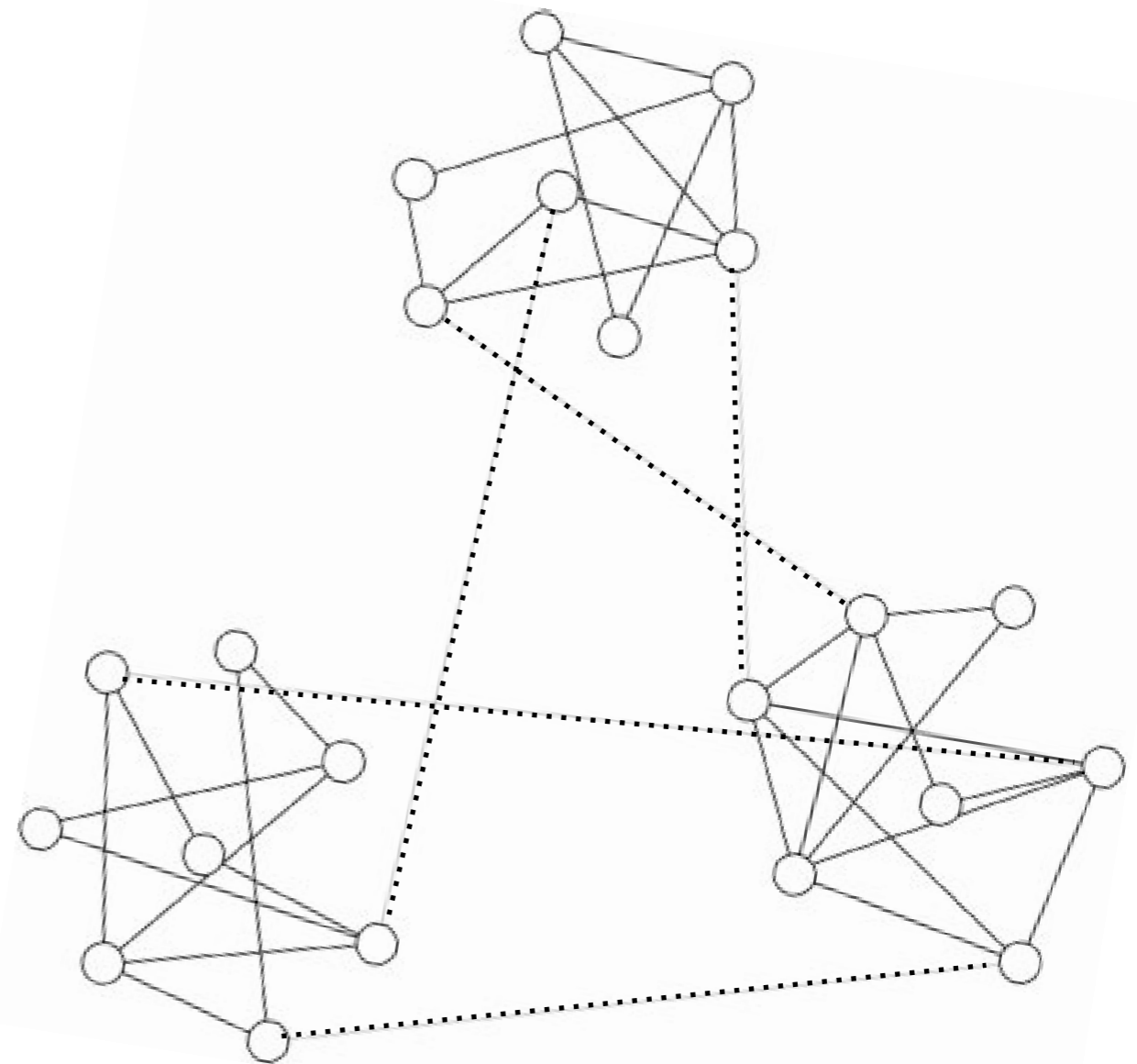
Test 1:

Generate random networks with known communities

- Divide nodes into communities
- Link each node to each other node with a set probability
- Probability of linking within your community greater than outside:

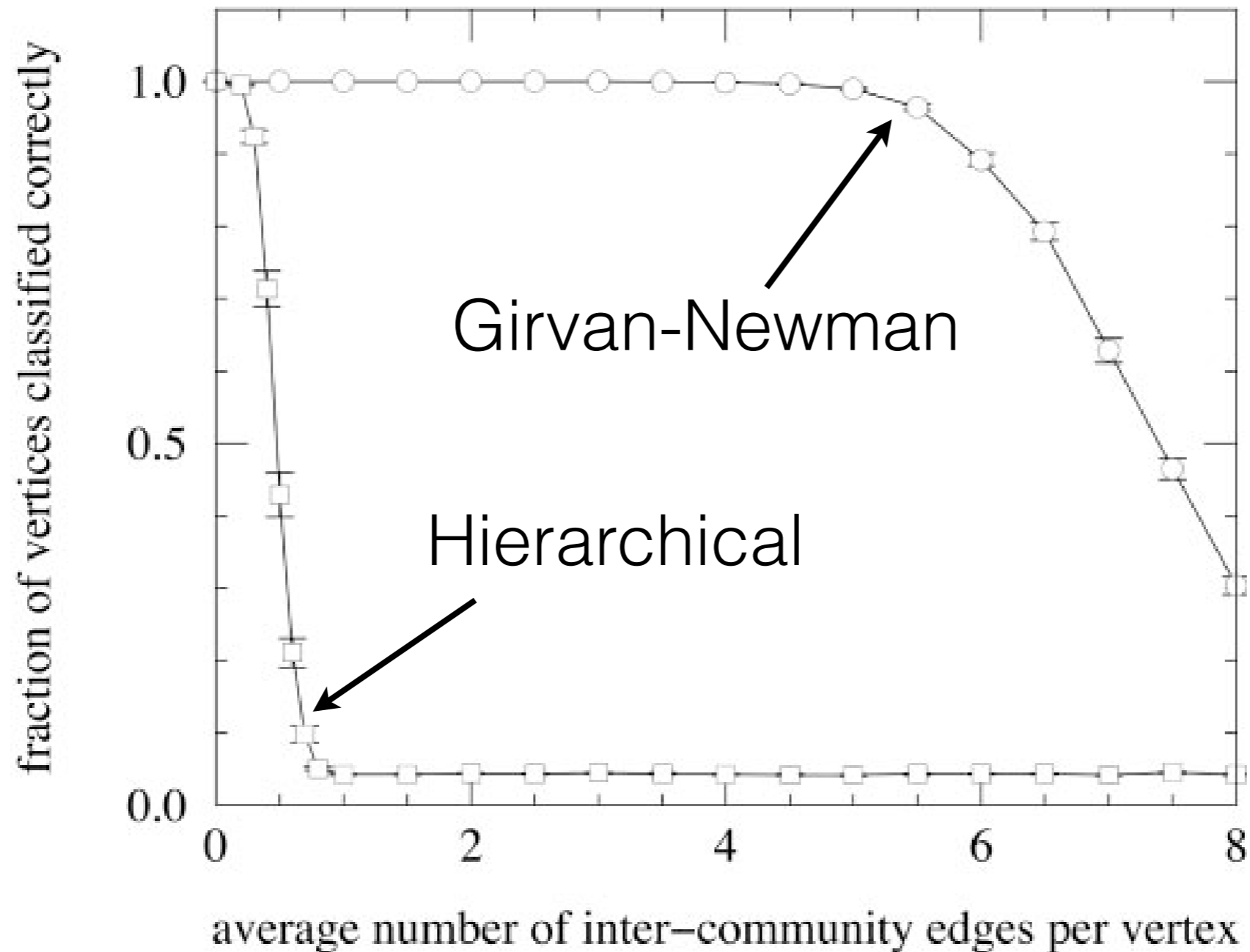
$$p_{\text{out}} < p_{\text{in}}$$

Run your algorithm: do you get out the same communities you put in?



# Evaluating Algorithms

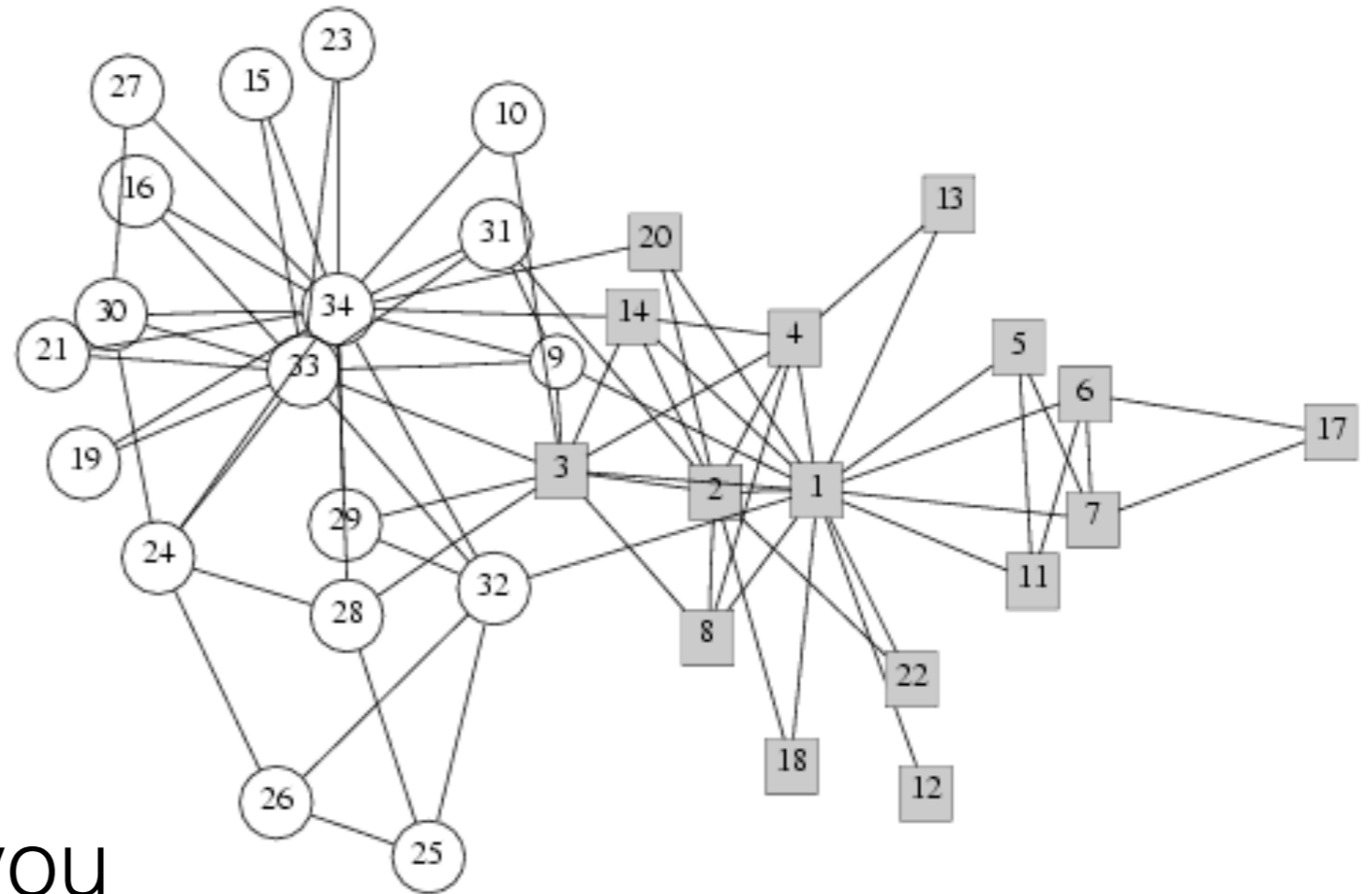
## Random Networks



# Evaluating Algorithms

## Known community structure

Test 2:  
Use a real social network with known community structure



Run your algorithm: do you get out the communities you know exist in the network?

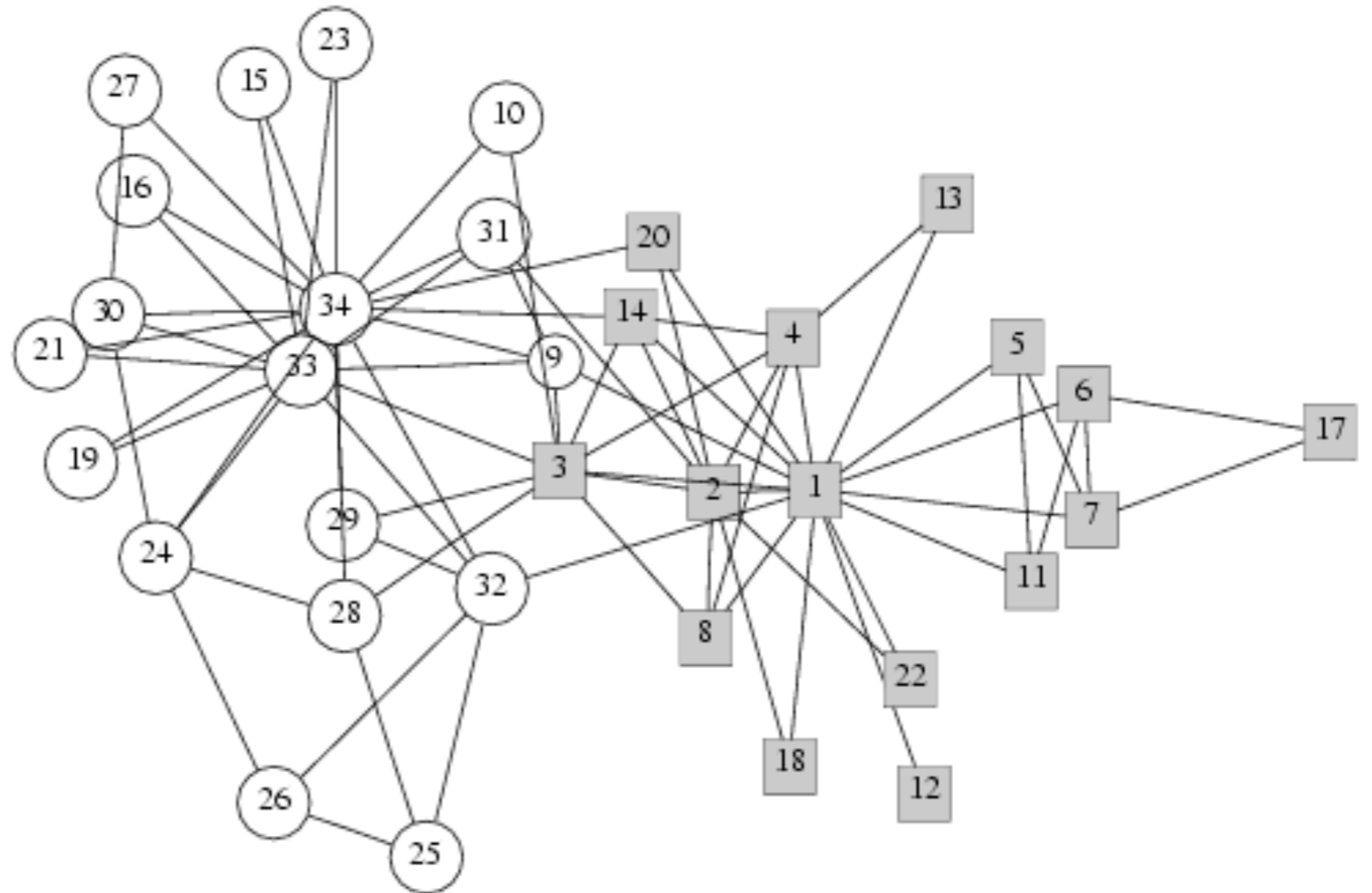
# Evaluating Algorithms

## Zachary's Karate Club

Karate Club with 34 members

During the study, the club split in half due to a disagreement

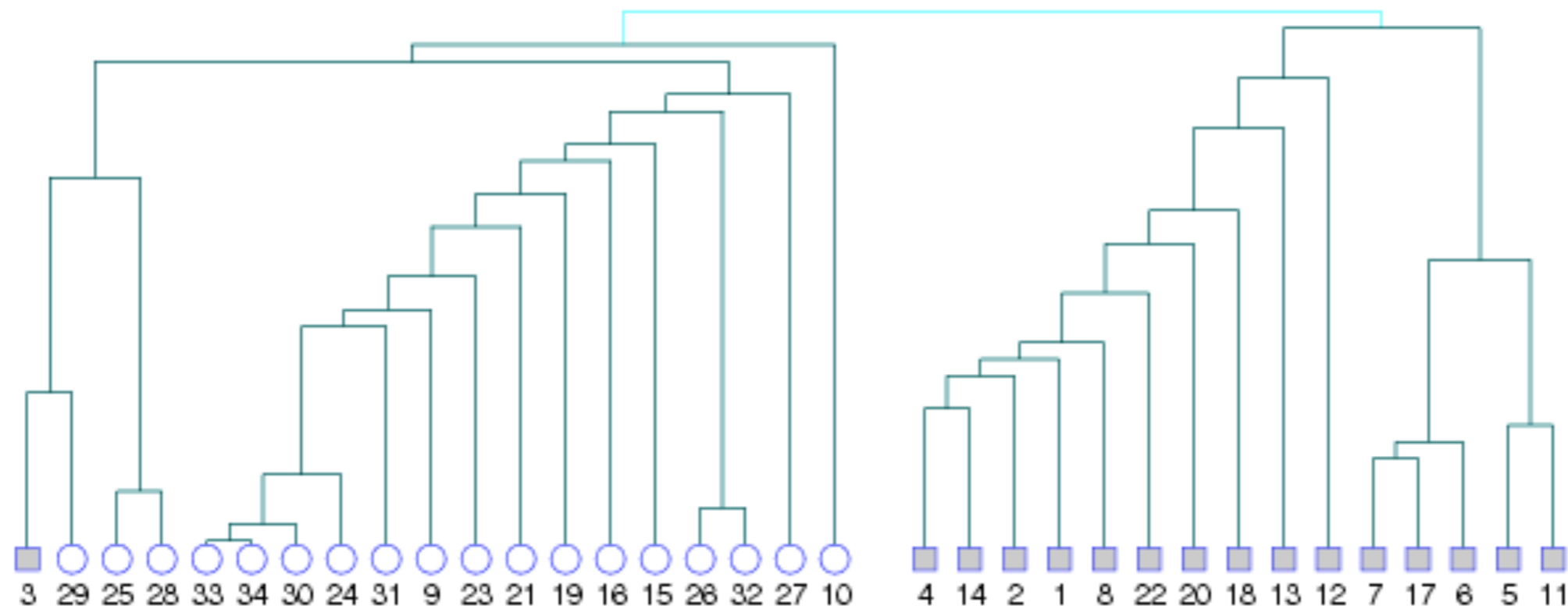
Based on the network, can the algorithms predict the actual split?





# Evaluating Algorithms

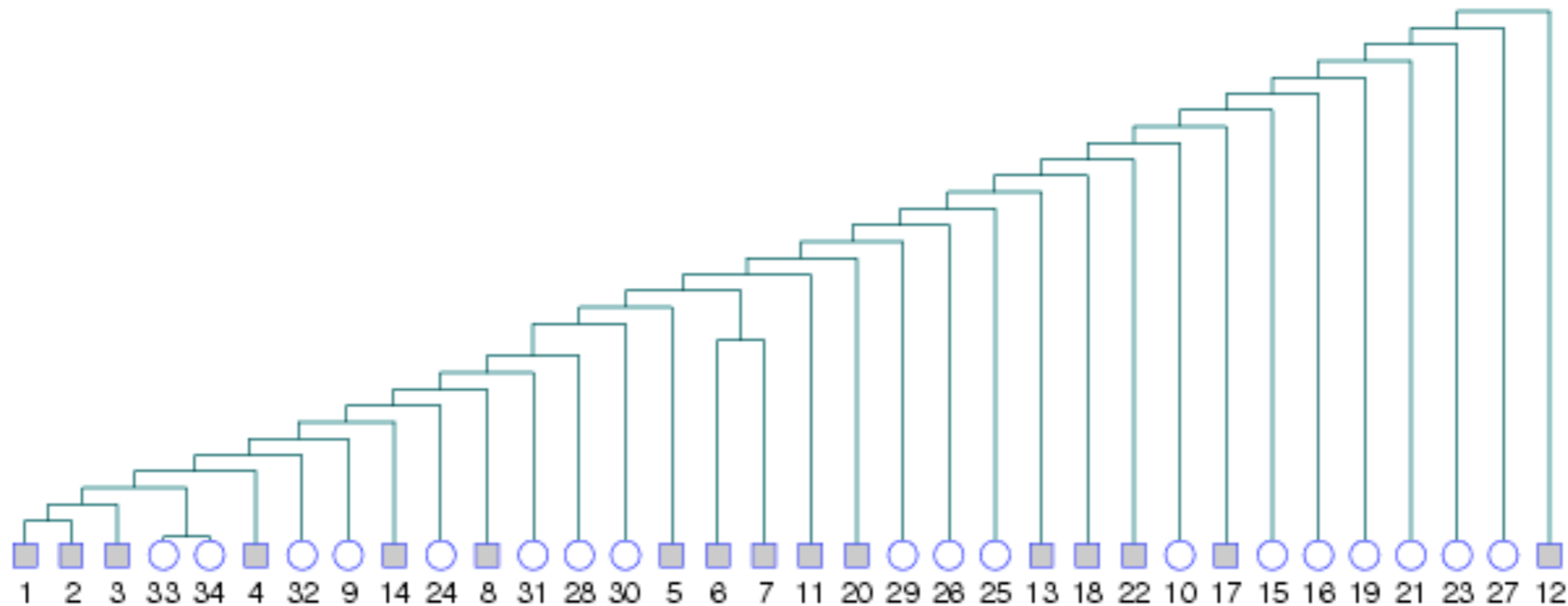
## Zachary's Karate Club



Girvan-Newman does quite well...

# Evaluating Algorithms

## Zachary's Karate Club



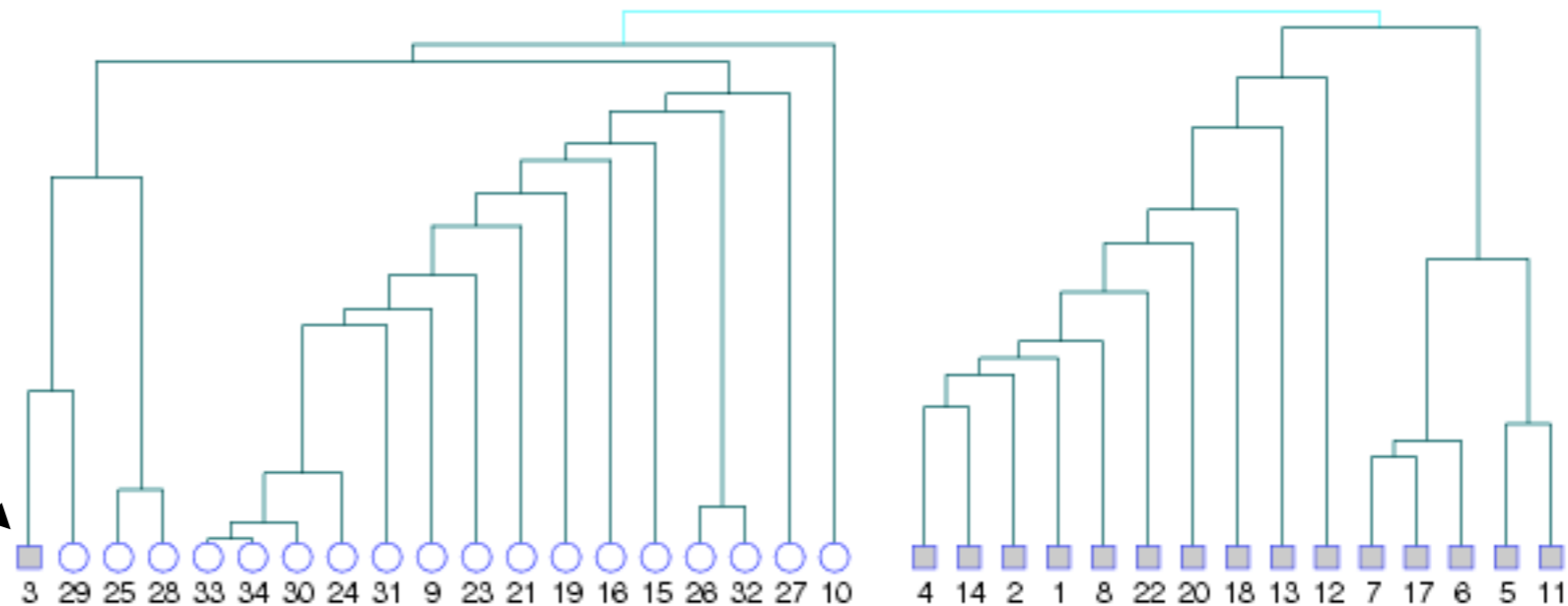
The hierarchical algorithm does quite poorly...

# Evaluating Algorithms

## Zachary's Karate Club

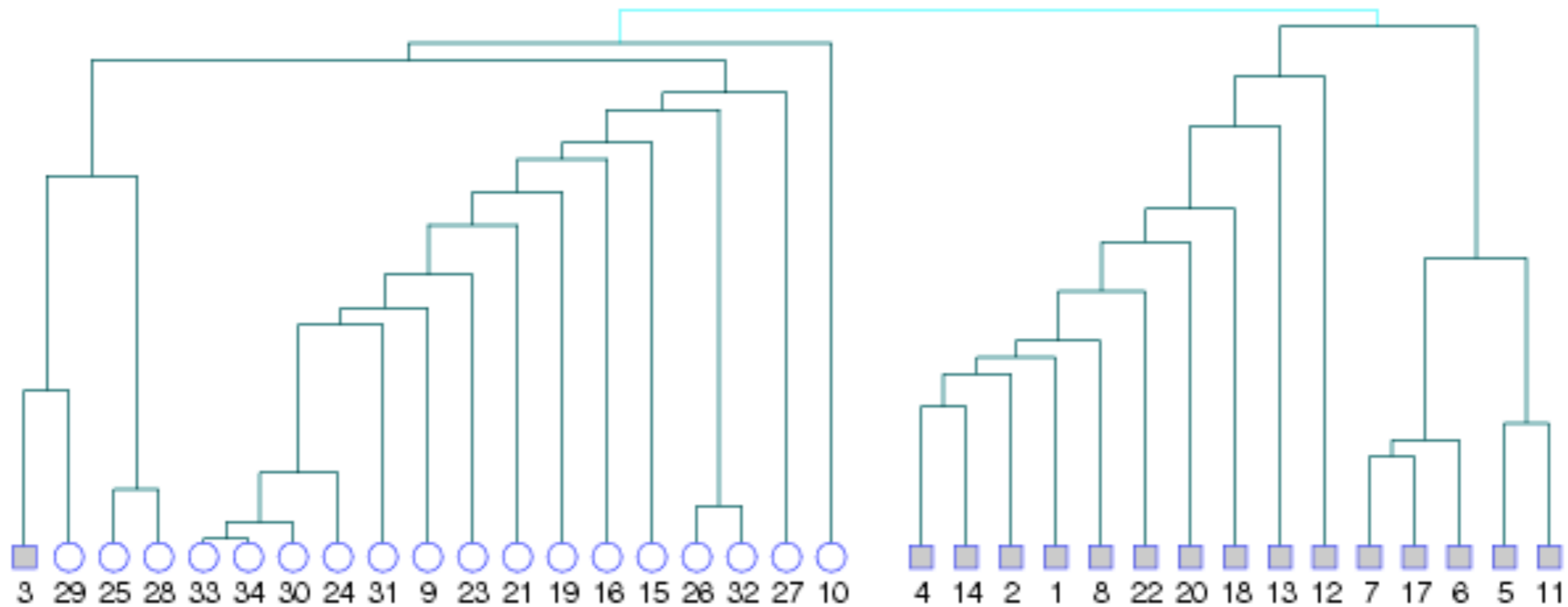
Note: the algorithms tell us about structure, not behavior. They can miss idiosyncrasies...

Node 1's brother

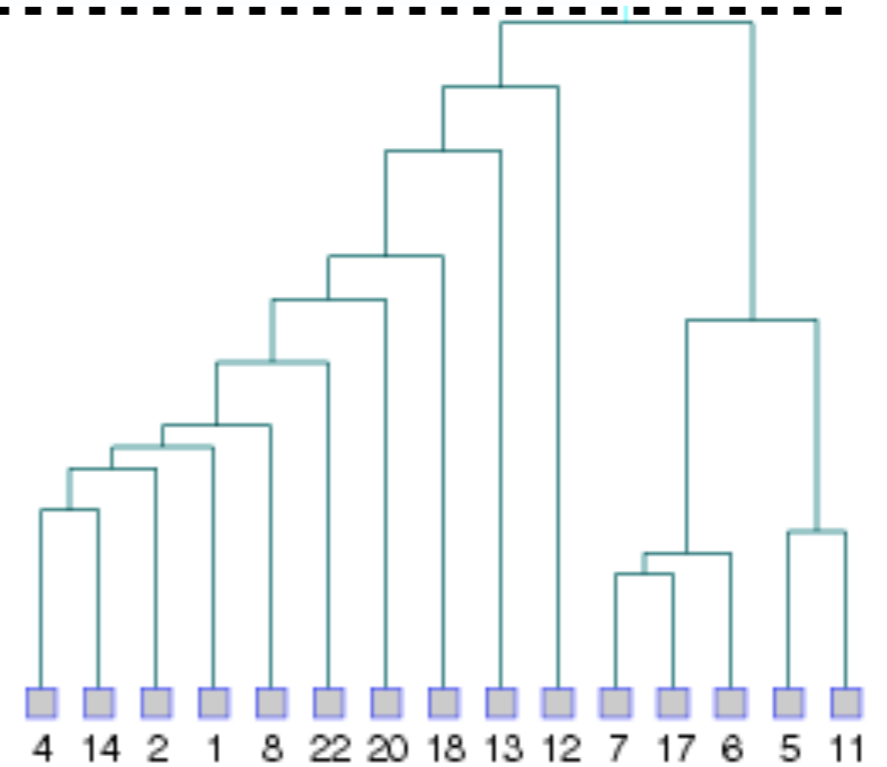
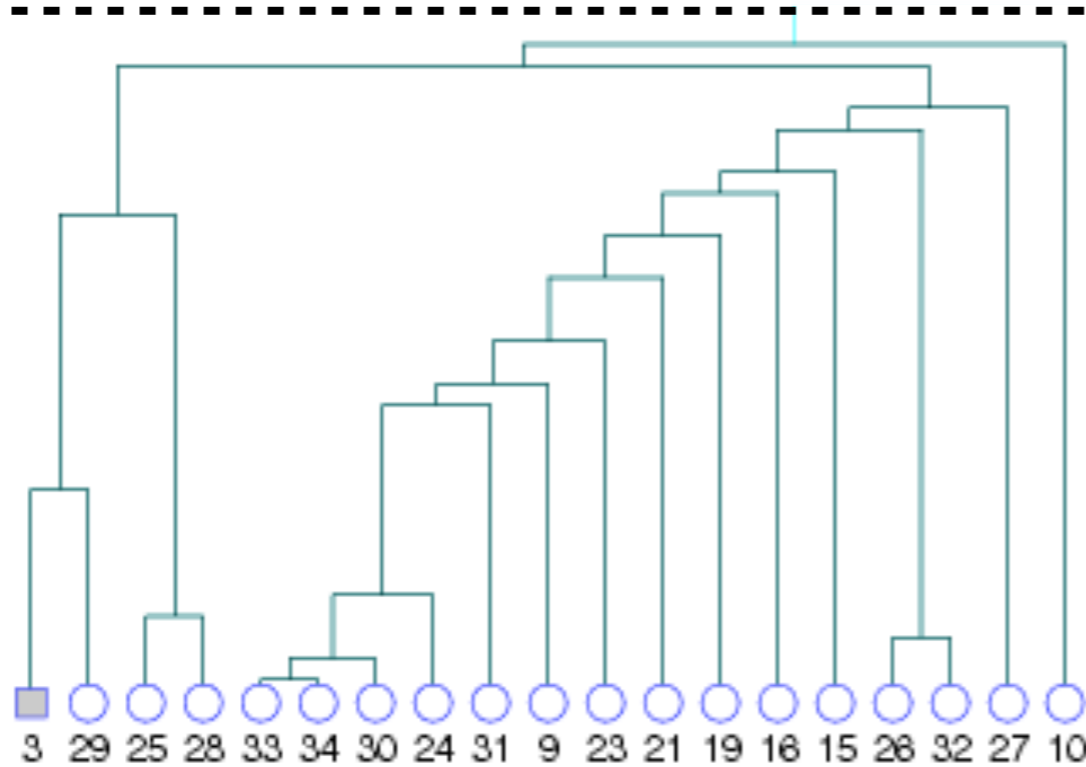
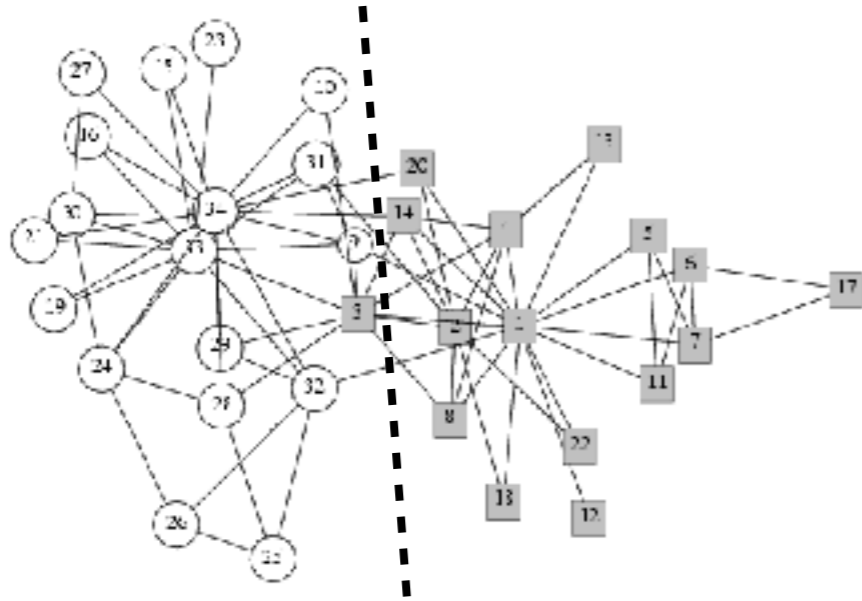


# Community Structure

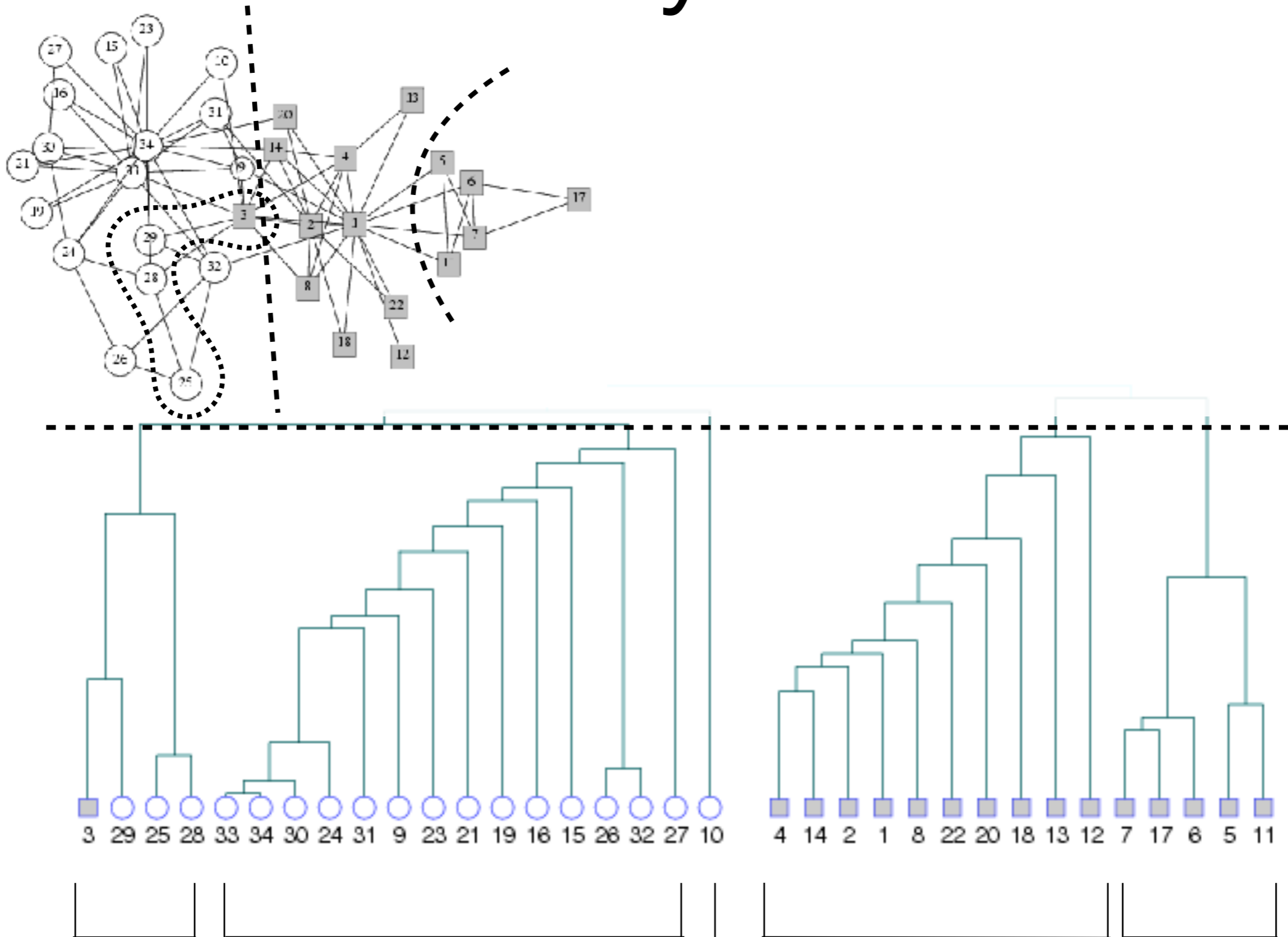
But one issue: these algorithms let us cut the network apart again and again...but when do we stop?



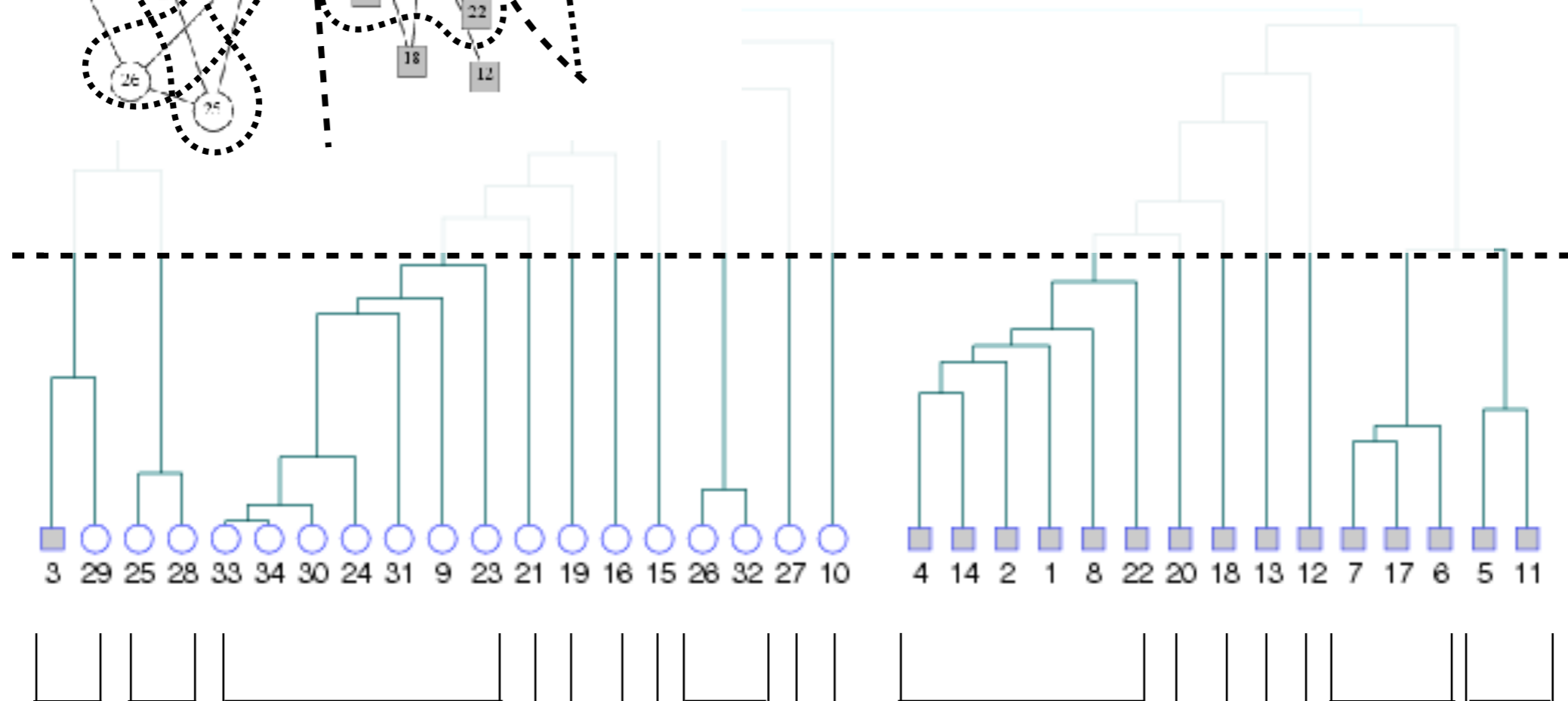
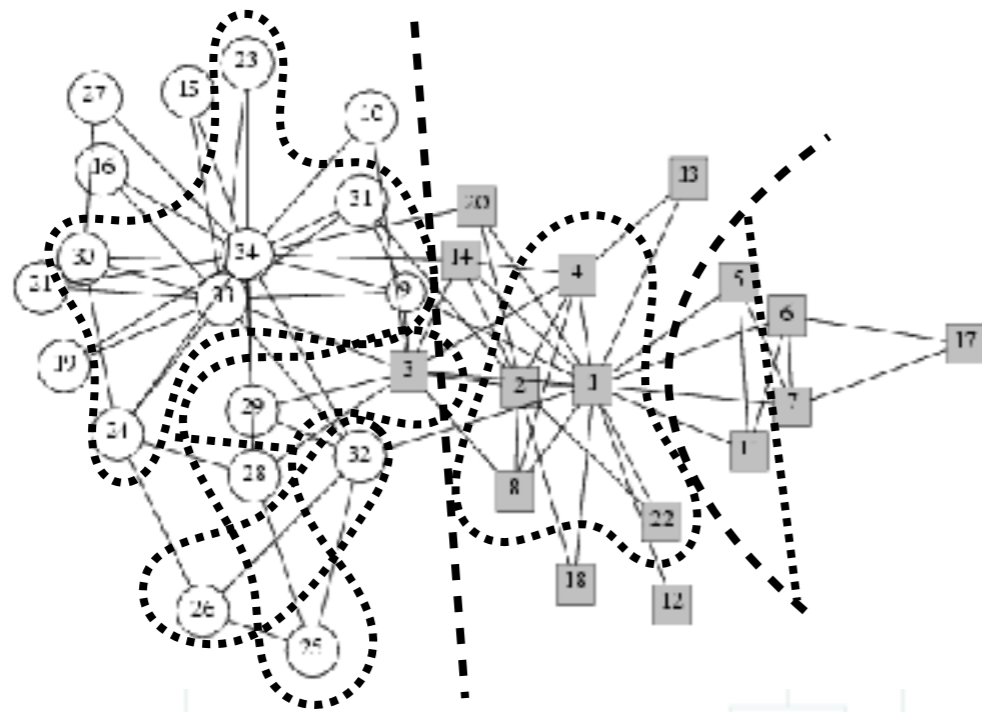
# Community Structure



# Community Structure

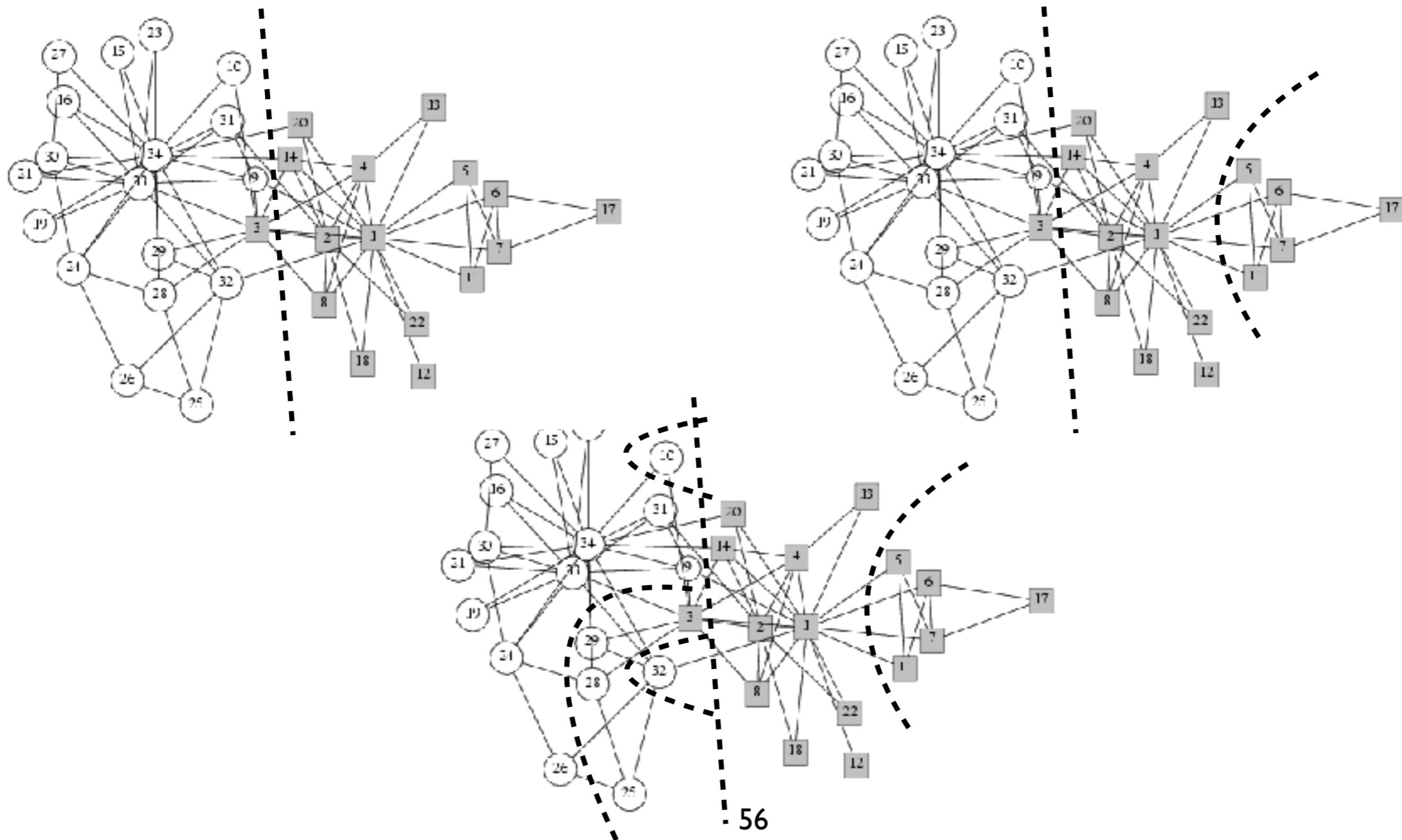


# Community Structure



# Community Structure

That is a problem if you don't have some exogenous information about community structure

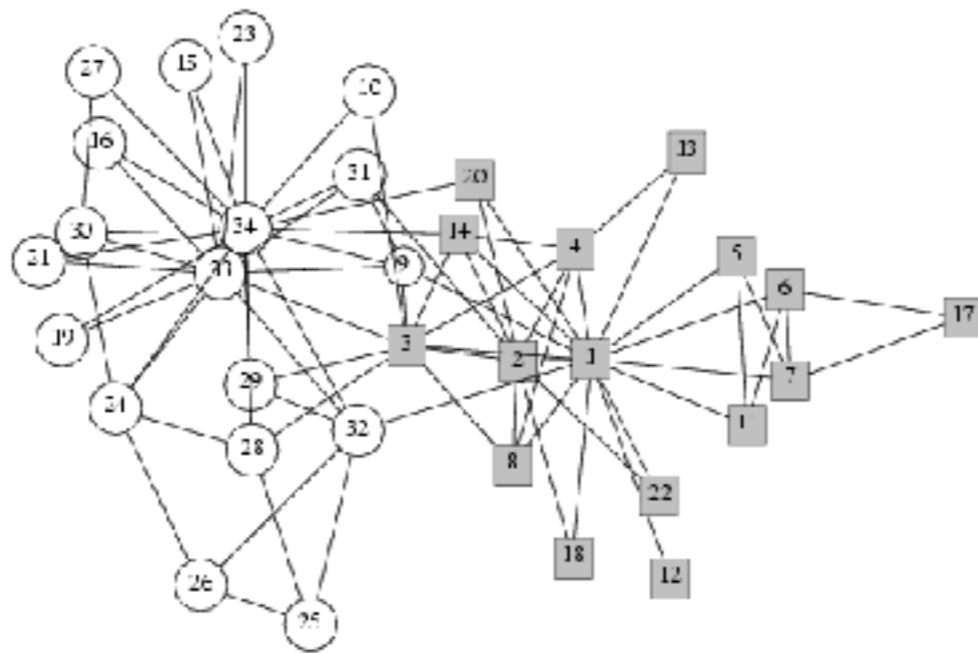




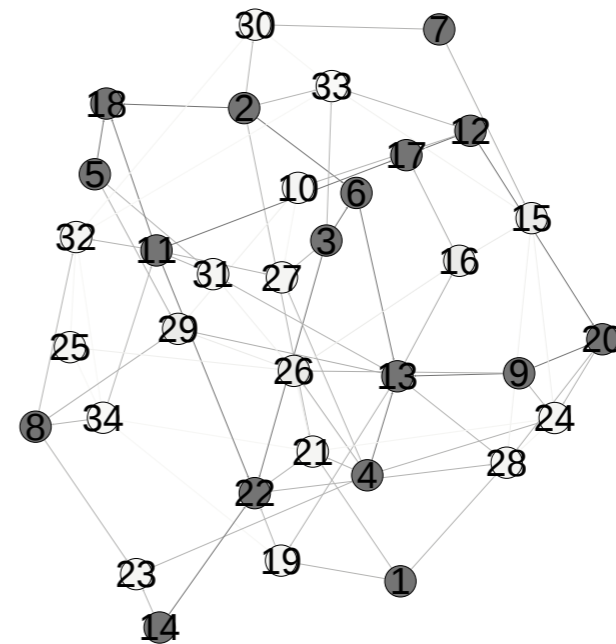
# Community Structure

What we want is to not find communities where they don't exist, but pull them out when they are unusual

One method: compare the partition you make on the actual network with the partition you would get on a similar random network:



Karate Club

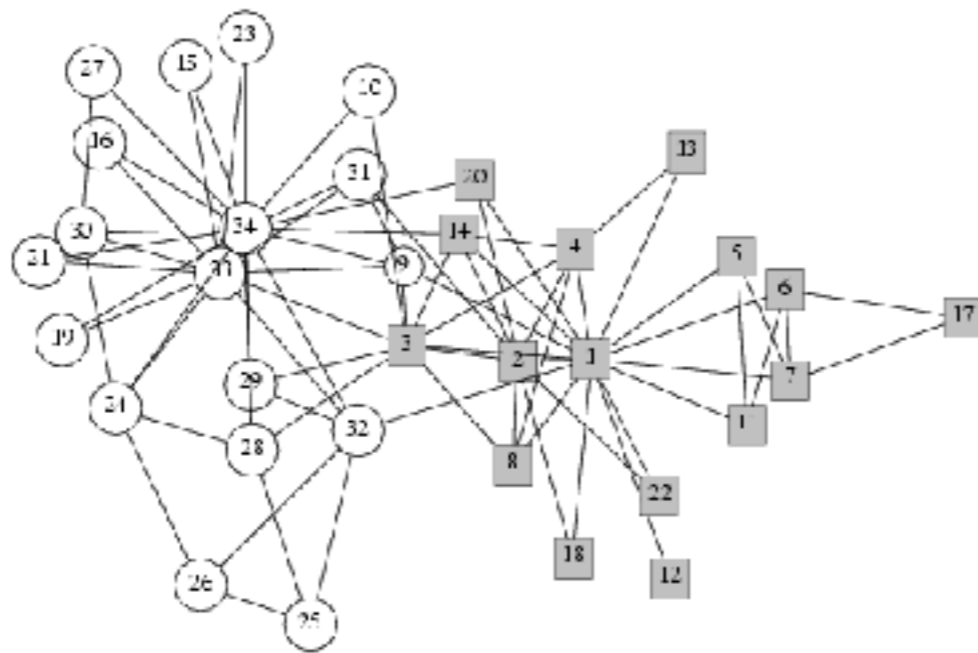


random network with the same number of nodes and same number of links

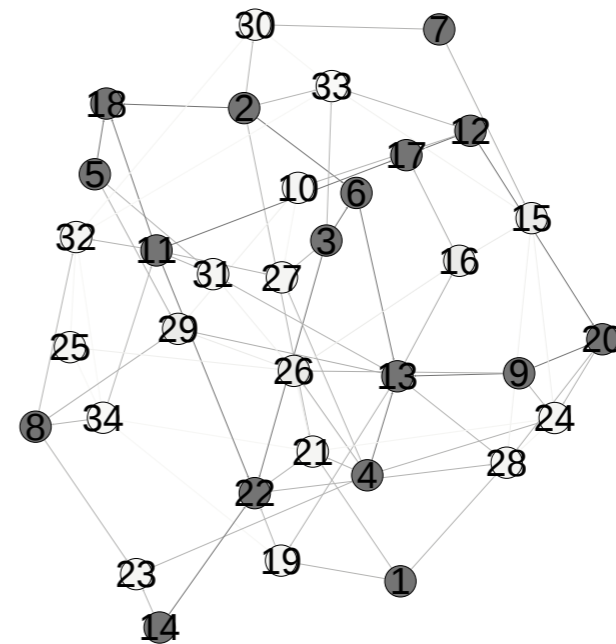
# Community Structure

When you partition your network, what fraction of the links are between communities?

When you use the same partition on a random network, what fraction are between communities?



Karate Club



random network with the same number of nodes and same number of links

# Modularity

Given a partition of the network into groups, *modularity* is a measure of how cohesive those groups are, relative to a random network

$$Q = \sum_i (e_{ii} - \bar{e}_i)$$

Sum over all groups

fraction of edges in the network that fall between nodes in group  $i$

fraction we would expect in a random network

In a random network:  $Q = \sum_i (\bar{e}_i - \bar{e}_i) = 0$

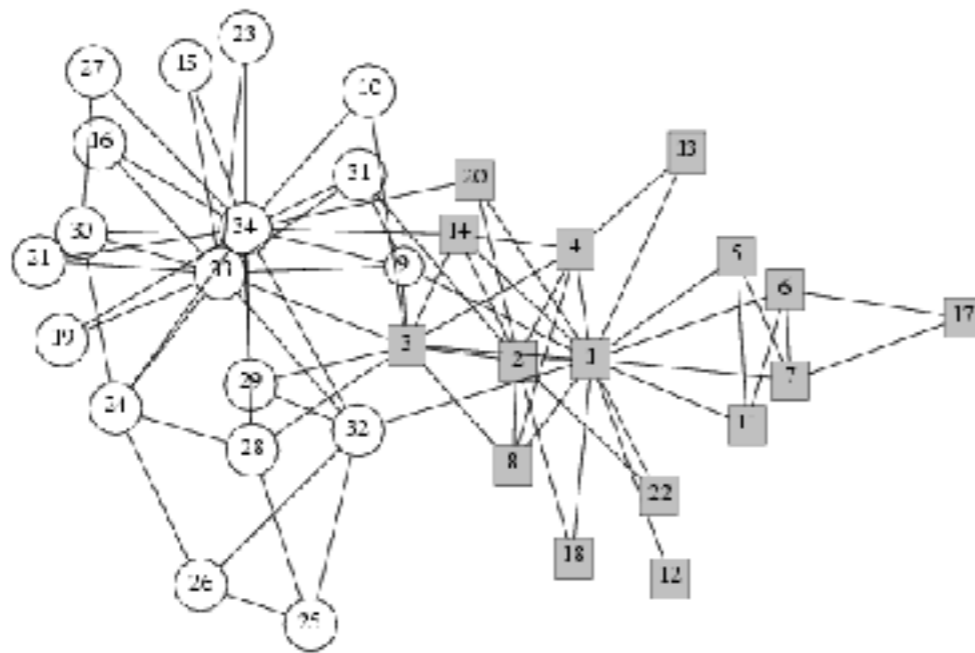
$0.3 < Q < 0.7$  indicates significant community structure

# Community Structure

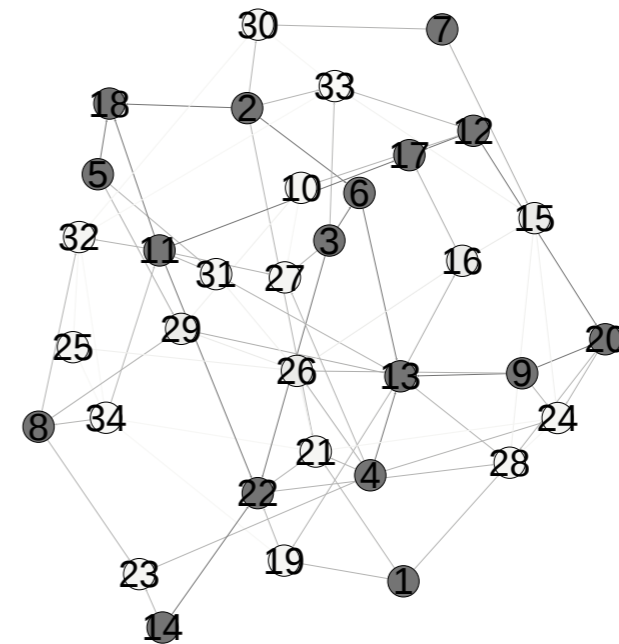
For this partition of the Karate Club Graph, there are 9 links between communities: ~12% of the links

In the random version, on average ~50% are

Modularity:  $(0.5 - 0.12) = .38$



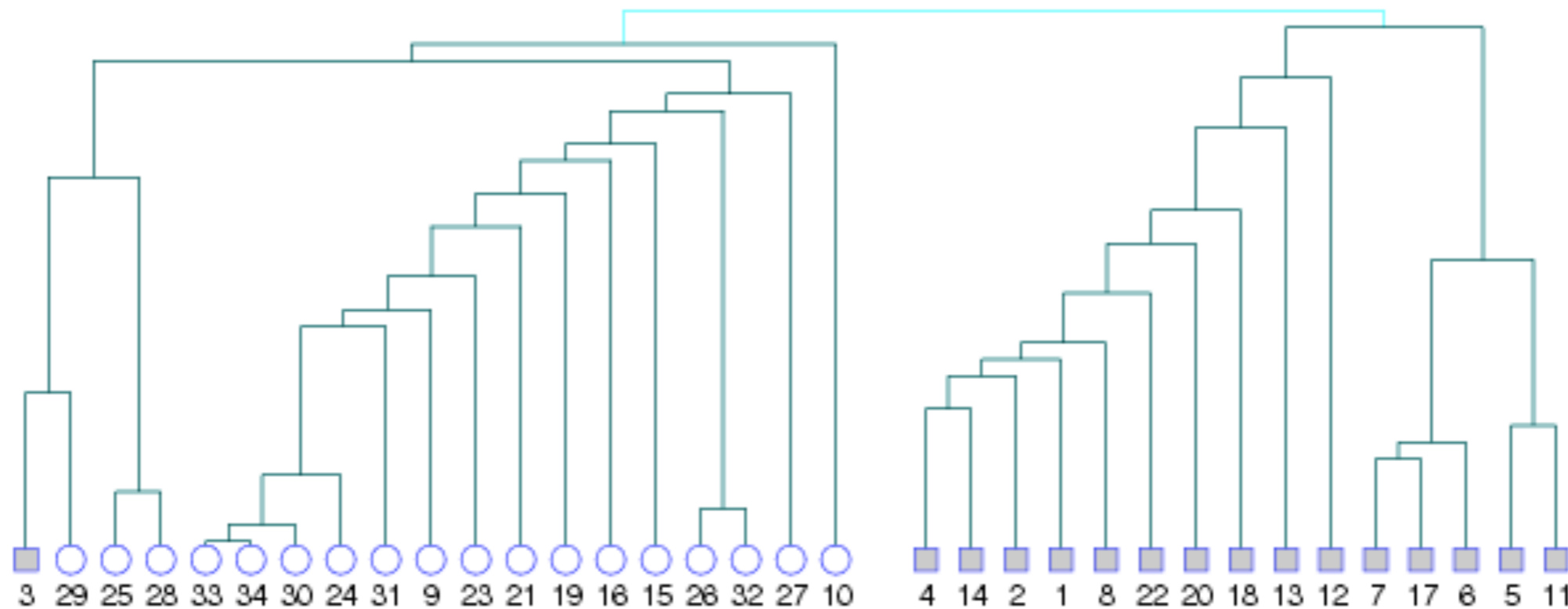
Karate Club



random network with the same number of nodes and same number of links

# Modularity

One idea for choosing when to stop dividing the network: we could just choose a division into communities that maximizes modularity ( $Q$ )



# Community Finding

## Big Picture

- Community structure is an interesting global property of networks
- There are many algorithms that one can use to distinguish communities
- The algorithms play off of different elements in the network, and produce different results
- When you stop dividing is important, and not obvious