

Outbreak Detection in Networks

CS 322: (Social and Information) Network Analysis
Jure Leskovec
Stanford University



Announcements

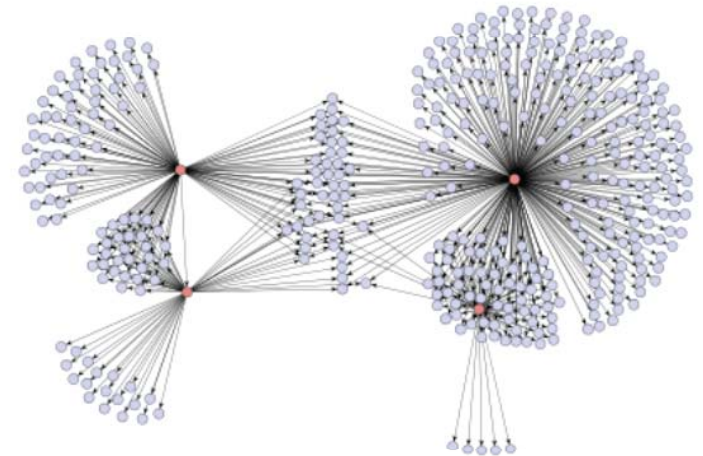
Thursday: guest lecture

Lars Backstrom (Facebook/Cornell)

on networks & geography

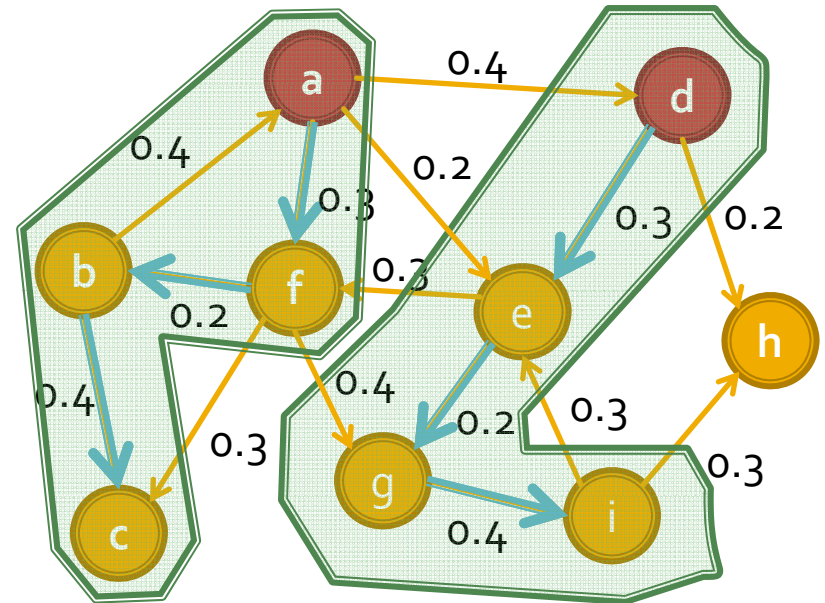
Finding influencers

- Blogs – information epidemics
 - Which are the influential/infectious blogs?
- Viral marketing
 - Who are the trendsetters?
 - Influential people?
- Disease spreading
 - Where to place monitoring stations to detect epidemics?



Most Influential Subset of Nodes

- Most influential set of size k : set S of k nodes producing largest expected cascade size $f(S)$ if activated
[Domingos-Richardson '01]



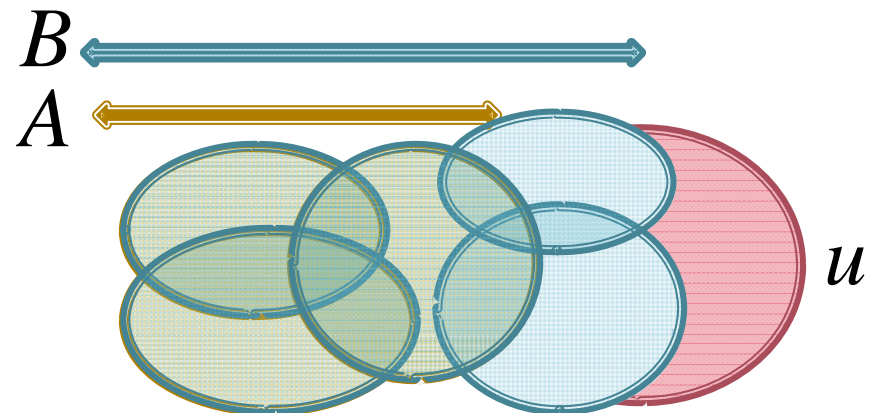
- Optimization problem: $\max_{S \text{ of size } k} f(S)$

Problem structure: Submodularity

- f is **submodular**: $S \subseteq T$

$$\underbrace{f(S \cup \{u\}) - f(S)}_{\text{Gain of adding a node to a small set}} \geq \underbrace{f(T \cup \{u\}) - f(T)}_{\text{Gain of adding a node to a large set}}$$

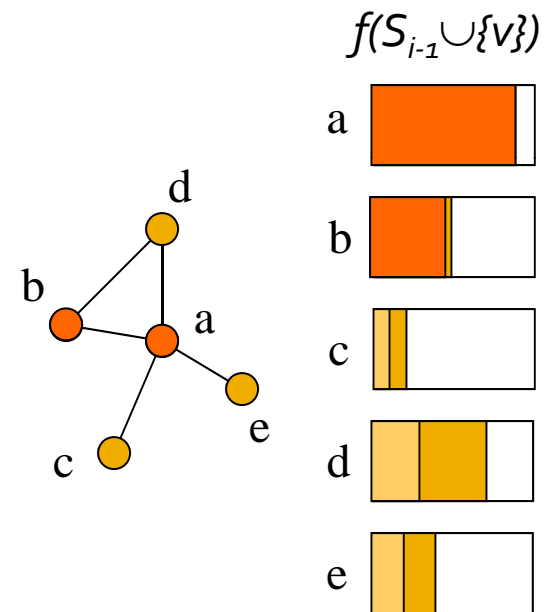
- Natural example:
 - Sets A_1, A_2, \dots, A_n
 - $f(A)$ = size of union of A_i
(size of covered area)



- If f_1, \dots, f_K are submodular, then $\sum p_i f_i$ is submodular

Hill climbing

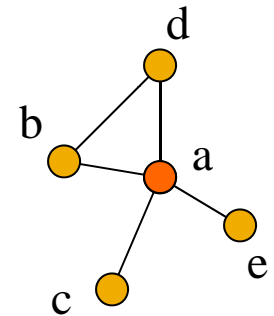
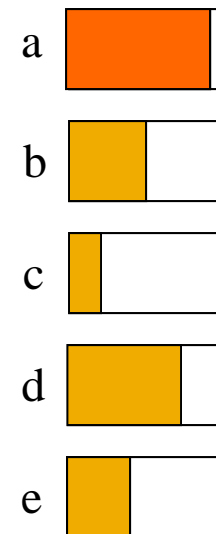
- Start with $S_0 = \{\}$
- For $i=1 \dots k$
 - Choose node v that $\max f(S_{i-1} \cup \{v\})$
 - Let $S_i = S_{i-1} \cup \{v\}$
- Hill climbing produces a solution S where $f(S) \geq (1 - 1/e)$ of optimal value (~63%) when f is monotone and submodular [Hemhauser, Fisher, Wolsey '78]



Lazy evaluation

- **Lazy hill-climbing:**
 - Keep an ordered list of marginal benefits b_i from previous iteration
 - Re-evaluate b_i **only** for top node
 - Re-sort and prune

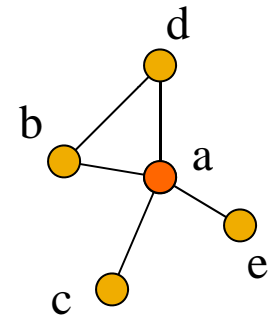
Marginal gain



Lazy evaluation

- **Lazy hill-climbing:**
 - Keep an ordered list of marginal benefits b_i from previous iteration
 - Re-evaluate b_i **only** for top node
 - Re-sort and prune

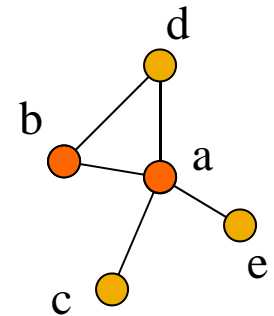
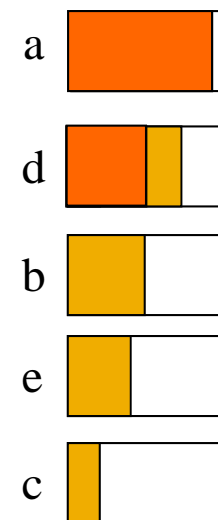
Marginal gain



Lazy evaluation

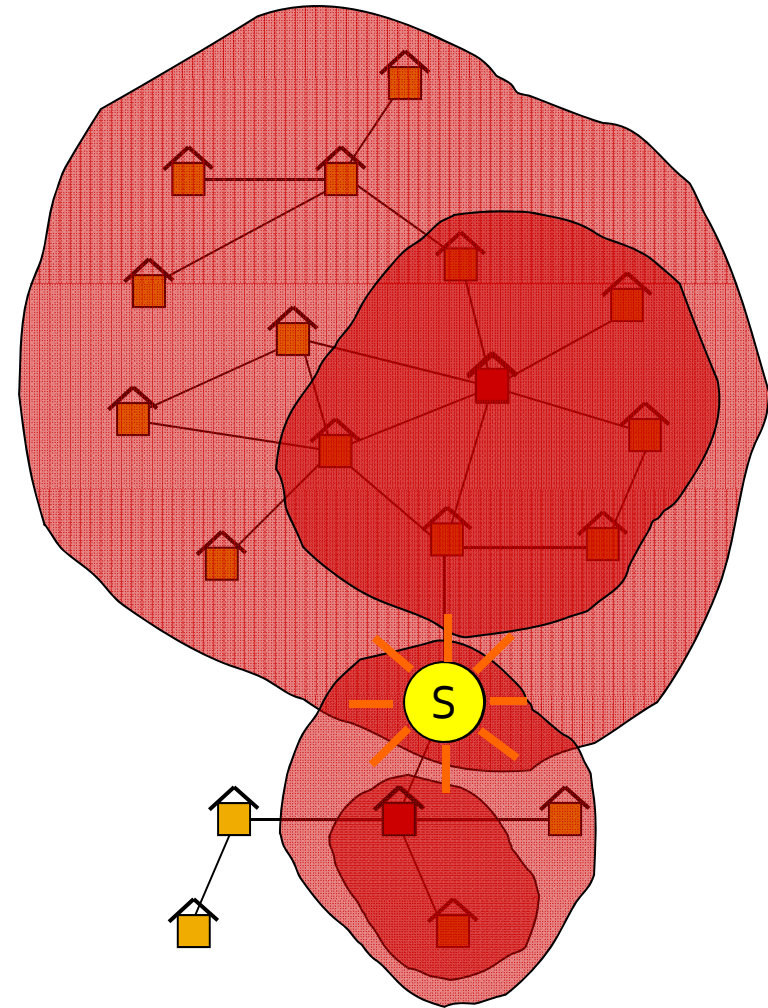
- **Lazy hill-climbing:**
 - Keep an ordered list of marginal benefits b_i from previous iteration
 - Re-evaluate b_i **only** for top node
 - Re-sort and prune

Marginal gain



Problem: Water Network

- Given a real city water distribution network
- And data on how contaminants spread in the network
- Problem posed by *US Environmental Protection Agency*



Problem Setting

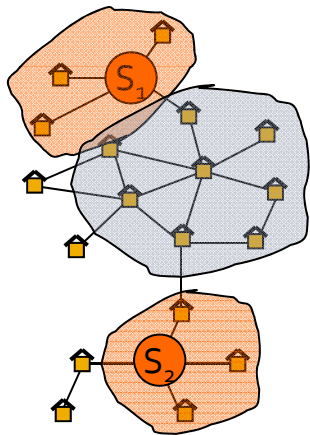
- Given a graph $G(V, E)$
- and a budget B for sensors
- and data on how contaminations spread over the network:
 - for each contamination i we know the time $T(i, u)$ when it contaminated node u
- Select a **subset** of nodes A that **maximize** the **expected reward**

$$\max_{A \subseteq \mathcal{V}} R(A) \equiv \sum_i P(i) \underbrace{R_i(T(i, A))}_{\text{Reward for detecting contamination } i}$$

subject to $cost(A) < B$

Structure of the Problem

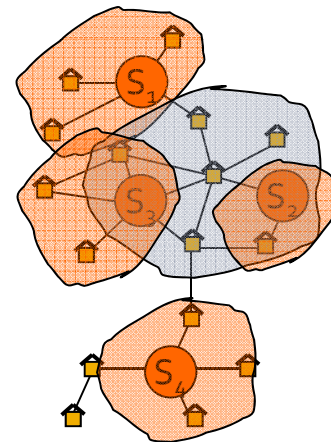
- Observation: **Diminishing returns**



Placement A = $\{S_1, S_2\}$

Adding S' helps a lot

New sensor:

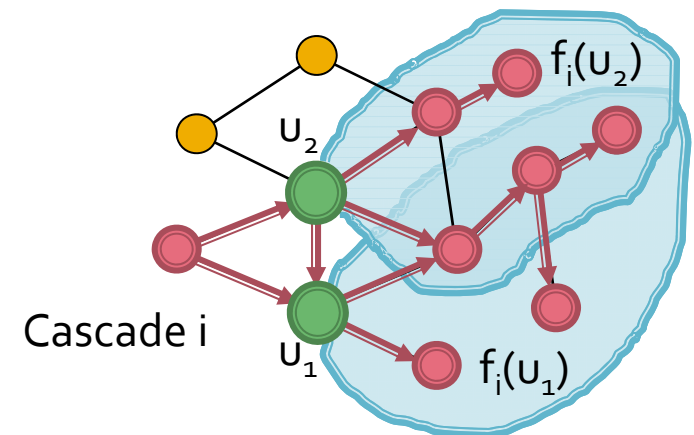


Placement B = $\{S_1, S_2, S_3, S_4\}$

Adding S' helps very little

Reward function is submodular

- Claim:
 - Reward function is submodular
- Consider cascade i :
 - $f_i(u_k)$ = set of nodes saved from u_k
 - $f_i(A)$ = size of union $f_i(u_k)$, $u_k \in A$
 - ⇒ f_i is **submodular**
- Global optimization:
 - $f(A) = \sum \text{Prob}(i) f_i(A)$
 - ⇒ f is **submodular**



Towards a New Algorithm

- Consider: **hill-climbing ignoring the cost**
 - Ignore sensor cost. Repeatedly select sensor with highest marginal gain
- It always prefers more expensive sensor with reward r to a cheaper sensor with reward $r-\epsilon$
 - For variable cost it can **fail arbitrarily badly**

- **Idea:** What if we optimize **benefit-cost ratio**?

$$s_k = \operatorname{argmax}_{s \in \mathcal{V} \setminus \mathcal{A}_{k-1}} \frac{R(\mathcal{A}_{k-1} \cup \{s\}) - R(\mathcal{A}_{k-1})}{c(s)}$$

Benefit-Cost: More Problems

- Benefit-cost ratio can **fail arbitrarily badly**
- Consider: budget B :
 - 2 locations s_1 and s_2 :
 - Costs: $c(s_1)=\varepsilon$, $c(s_2)=B$
 - Only 1 cascade: $R(s_1)=2\varepsilon$, $R(s_2)=B$
 - Then benefit-cost ratio is
 - $bc(s_1)=2$ and $bc(s_2)=1$
 - So, we first select s_1 and then can not afford s_2
- We get reward 2ε instead of B
- Now send ε to 0 and we get **arbitrarily bad**

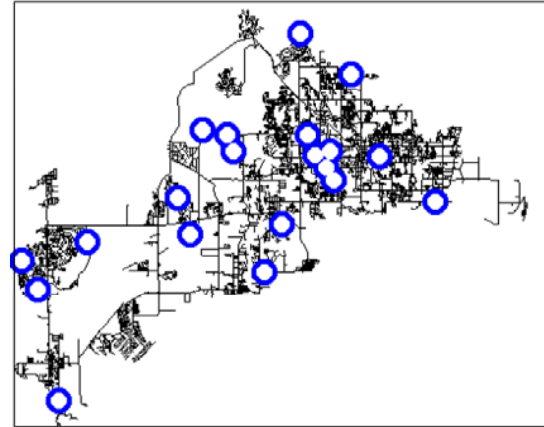
Solution: CELF Algorithm

- **CELF (cost-effective lazy forward-selection):**
A two pass greedy algorithm:
 - Set (solution) A: use benefit-cost greedy
 - Set (solution) B: use unit cost greedy
 - Final solution: $\text{argmax}(R(A), R(B))$

- Theorem: **CELF is near optimal**
 - CELF achieves $1/2(1-1/e)$ factor approximation

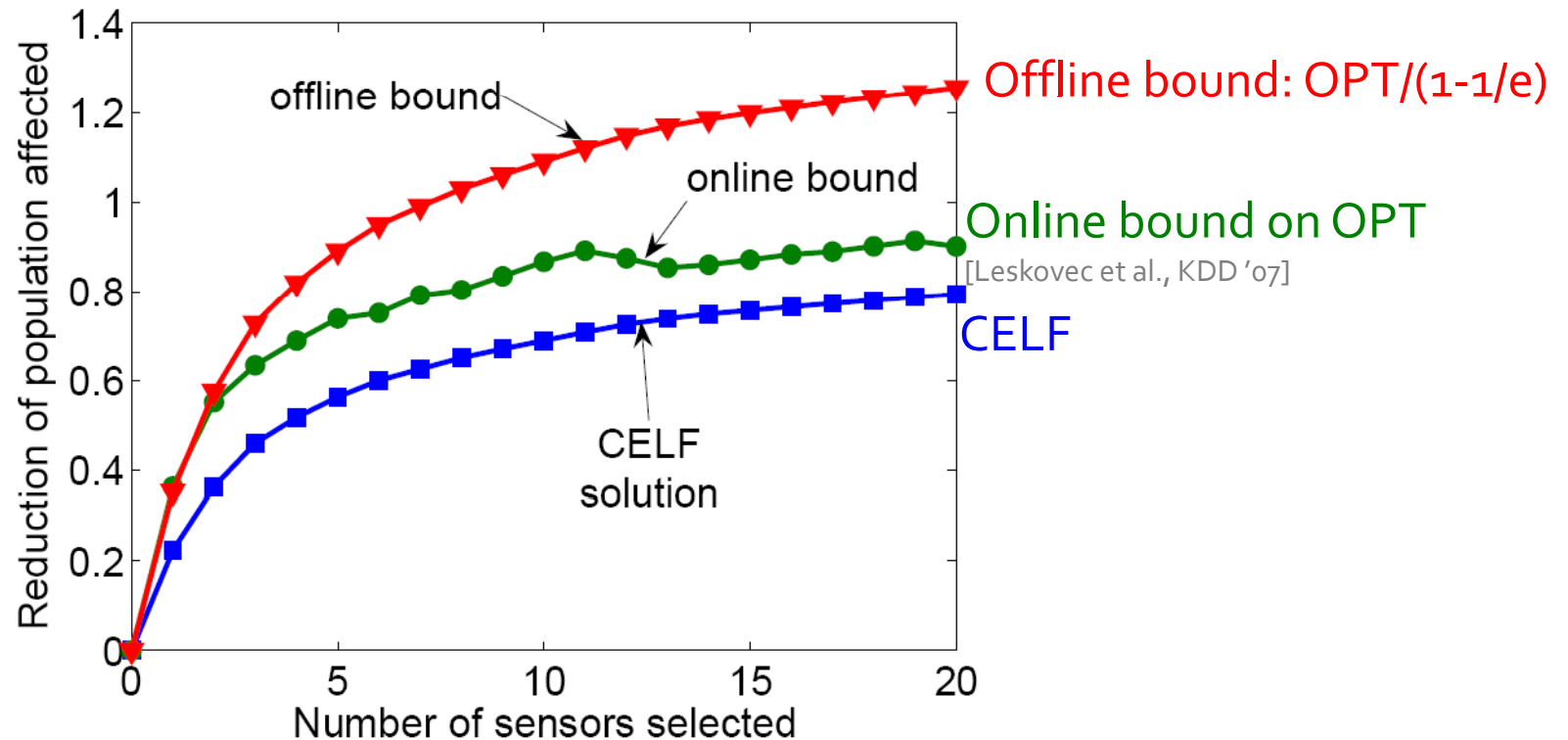
Case study: Water Network

- Real metropolitan area water network
 - $V = 21,000$ nodes
 - $E = 25,000$ pipes



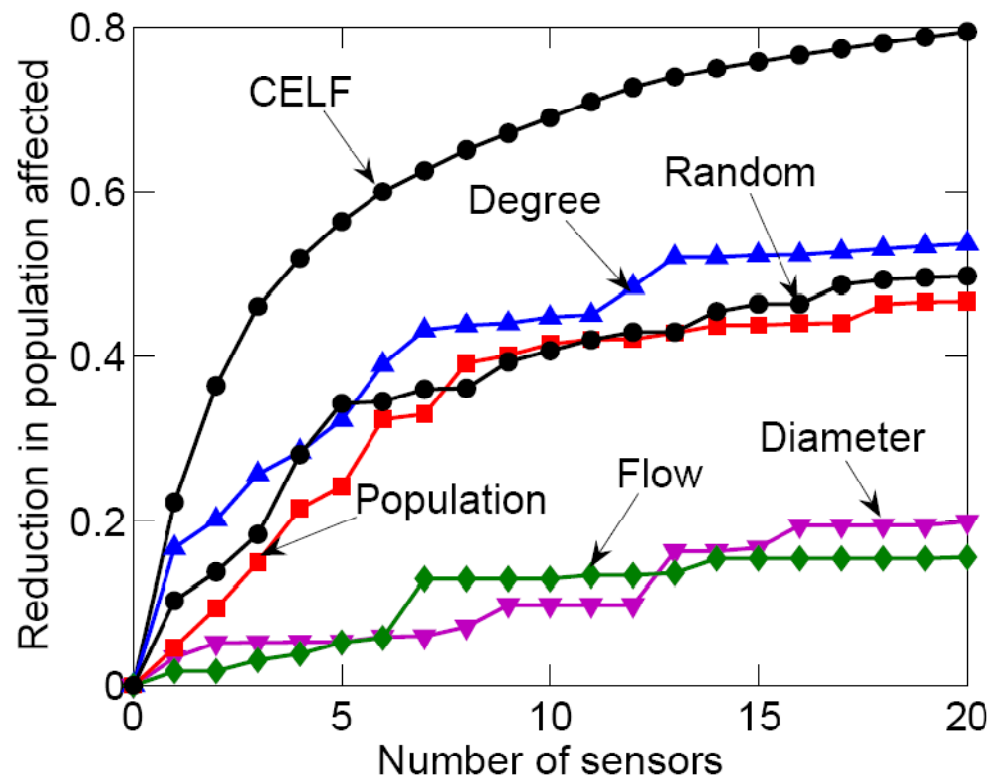
- Use a cluster of 50 machines for a month
- Simulate 3.6 million epidemic scenarios (152 GB of epidemic data)
- By exploiting sparsity we fit it into main memory (16GB)

Water: Solution Quality



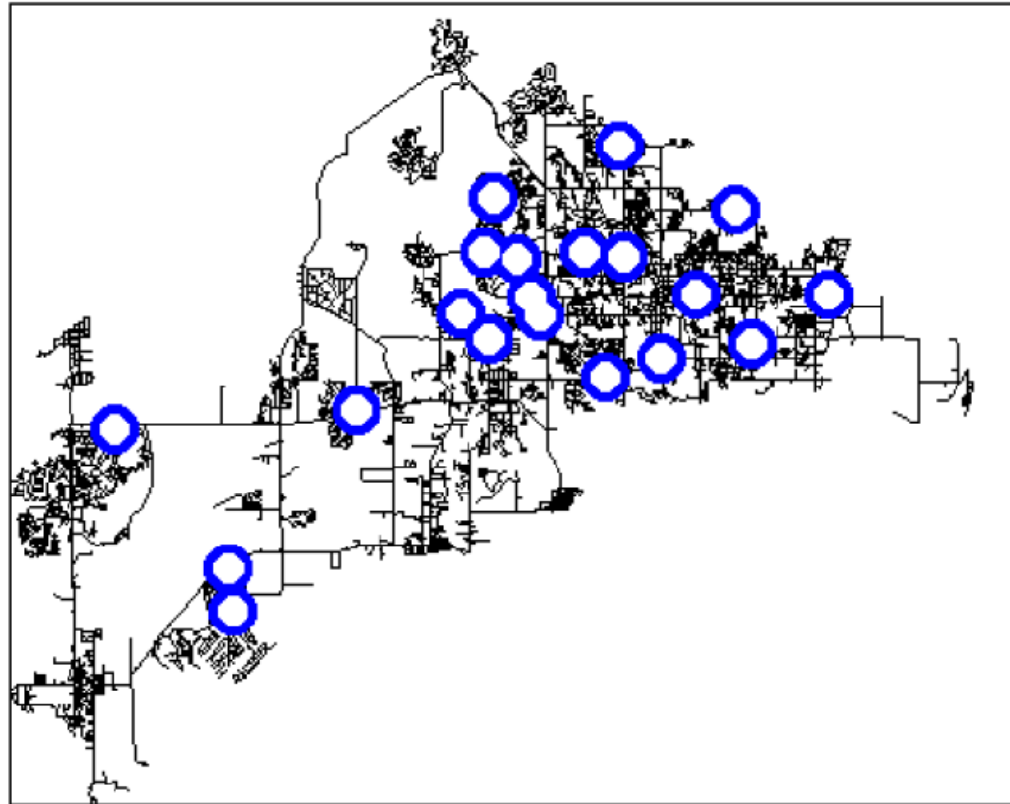
- The online (data dependent) bound gives much better estimate of how far from unknown optimal solution is the CELF solution

Water: Heuristic Placement

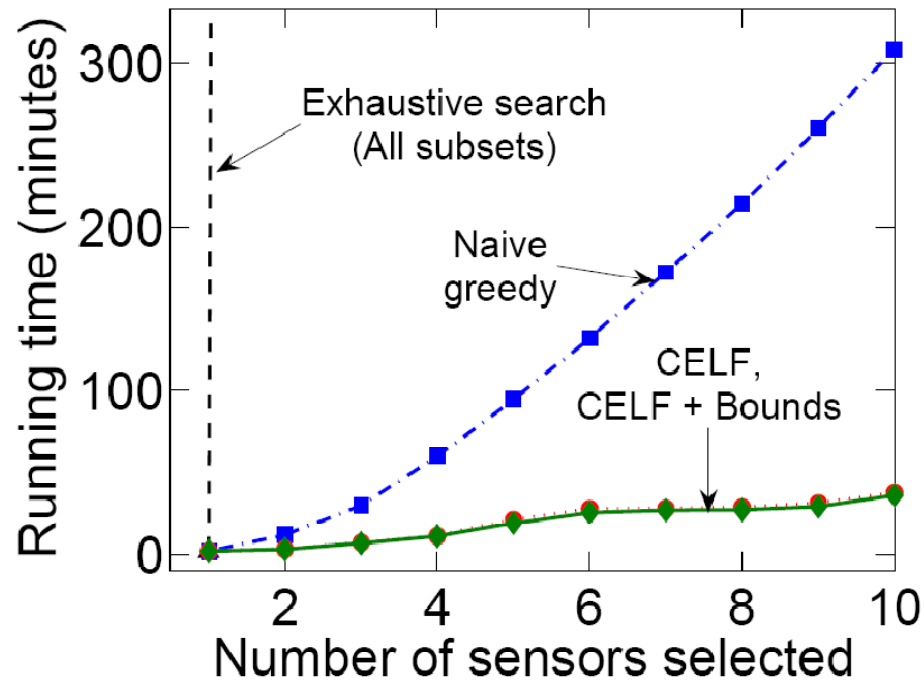


- Heuristics placements perform much worse
- One really needs to consider the spread of epidemics

Water: Placement Visualization



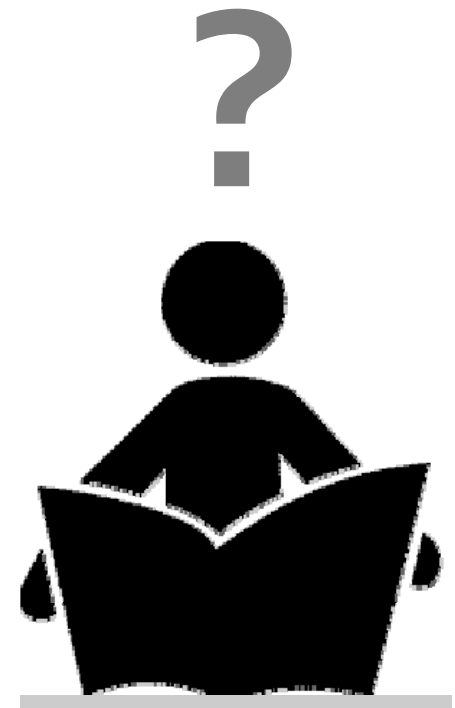
Water: Algorithm Scalability



- CELF is an order of magnitude faster than hill-climbing

Question...

- = I have 10 minutes. Which blogs should I read to be most up to date?
- = Who are the most influential bloggers?



Detecting information outbreaks

Want to read things before others do.

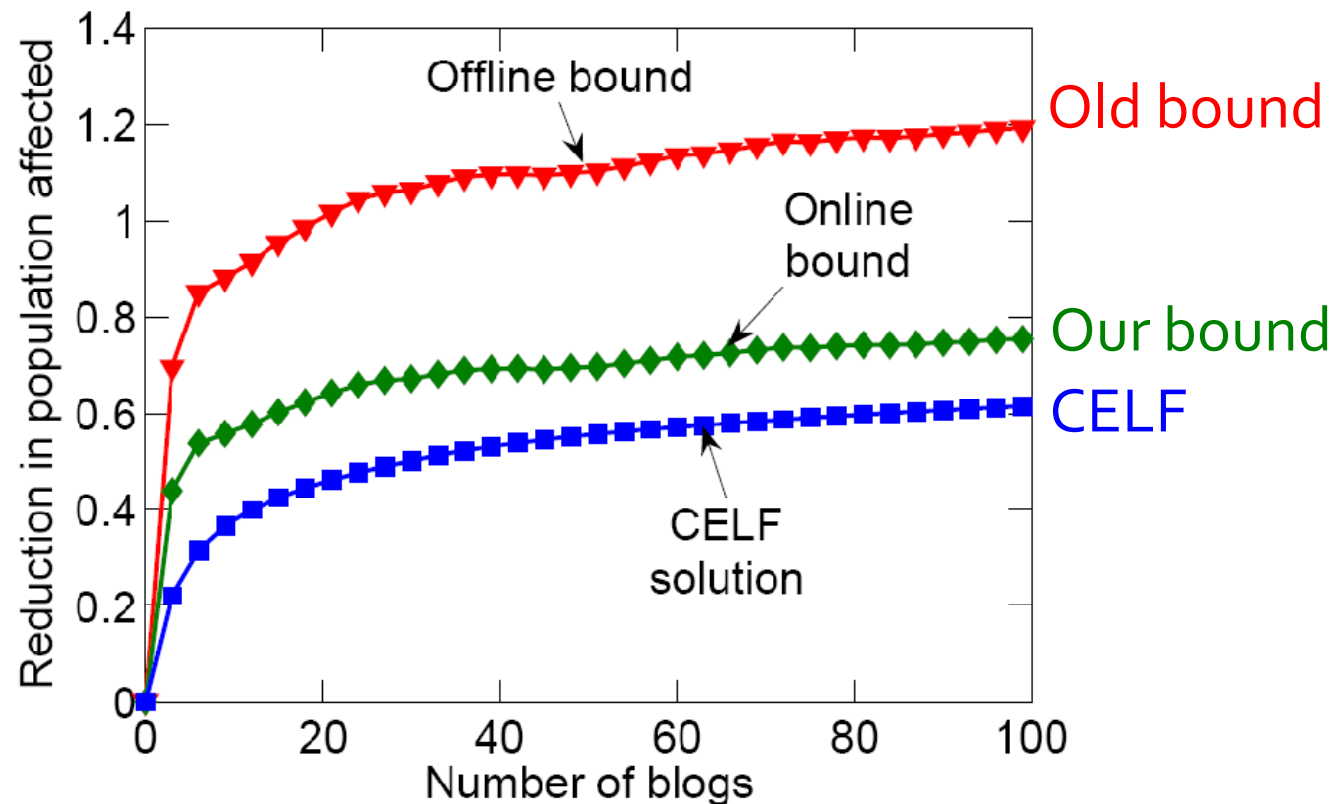
Detect **blue** & **yellow** soon but miss **red**.

Detect **all** stories but **late**.

Blogs: Solution Quality

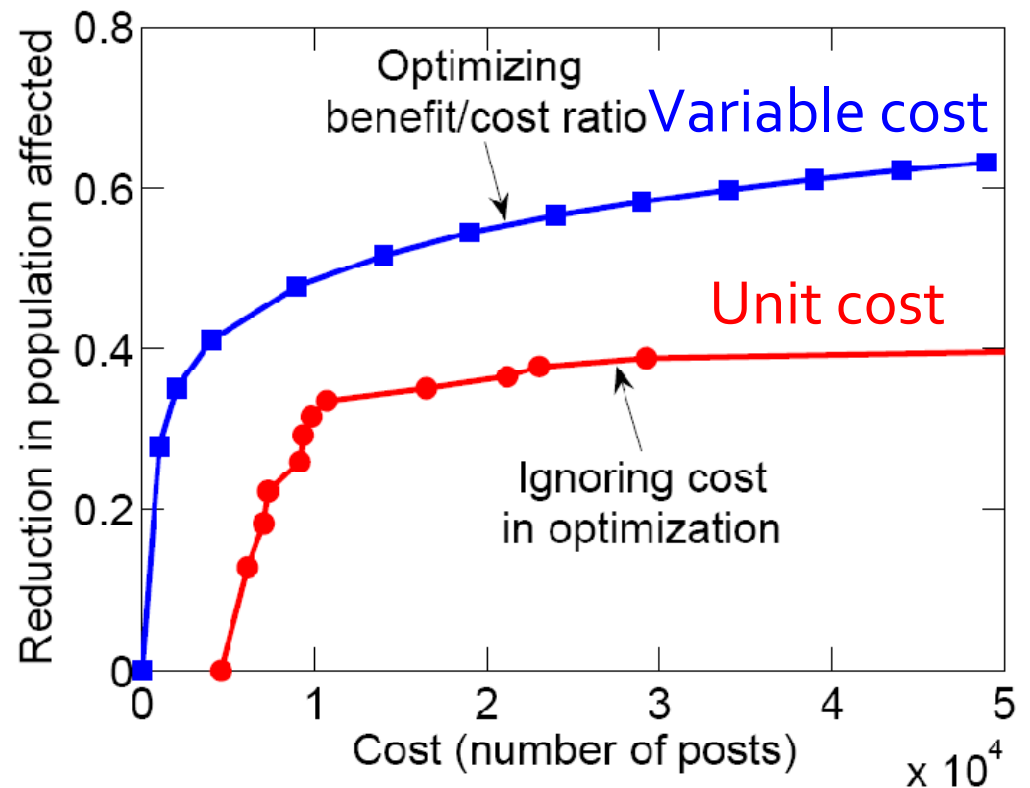
- Online bound [Leskovec et al., KDD '07] is much tighter

- 13%



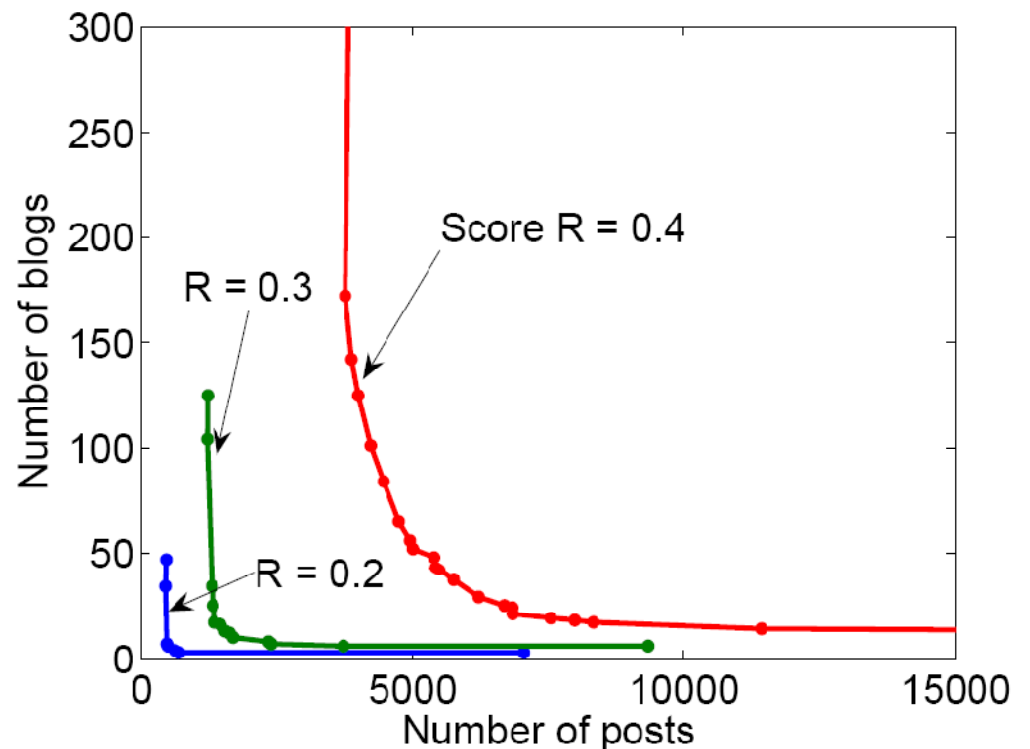
Blogs: Cost of a Blog

- Unit cost:
 - algorithm picks **large popular blogs**:
instapundit.com,
michellemalkin.com
- Variable cost:
 - proportional to the **number of posts**
- We can do much better when considering costs



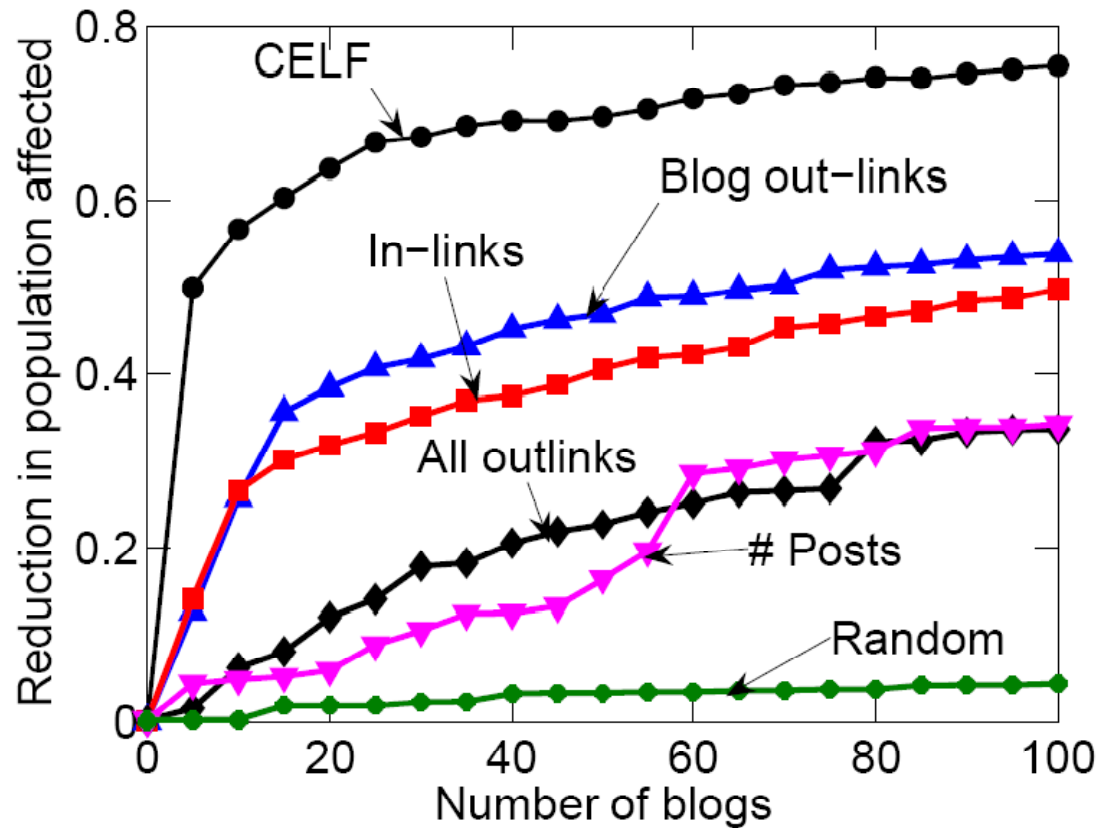
Blogs: Cost of a Blog

- But then algorithm picks **lots of small blogs** that participate in few cascades
- We pick best solution that interpolates between the costs
- We can get good solutions with **few blogs and few posts**



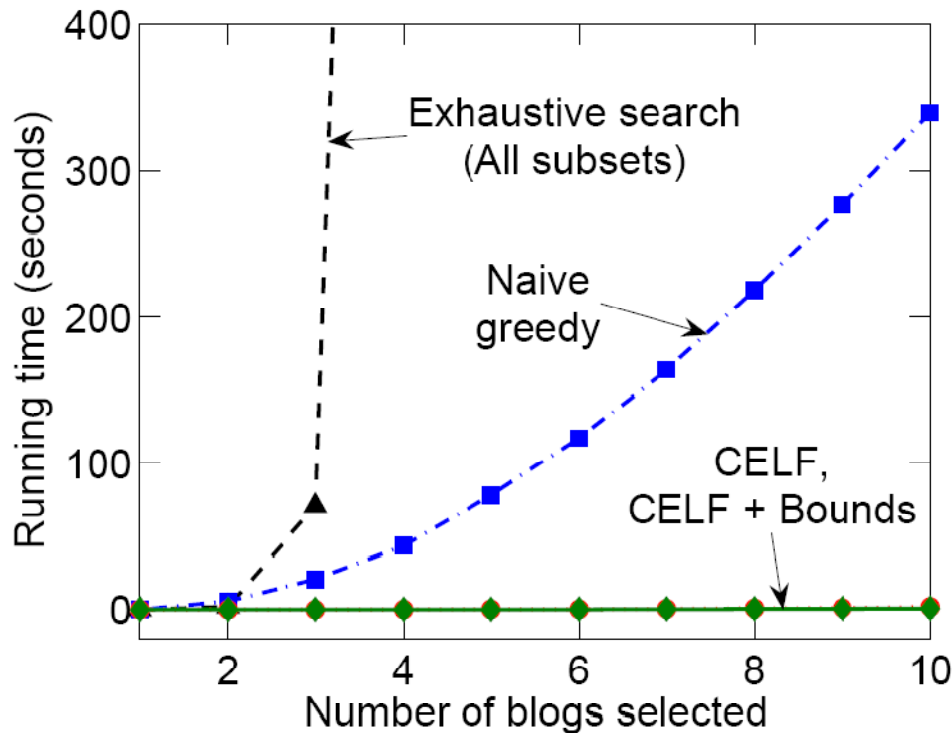
Each curve represents solutions with same final reward

Blogs: Heuristic Selection



- Heuristics perform much worse
- One really needs to perform optimization

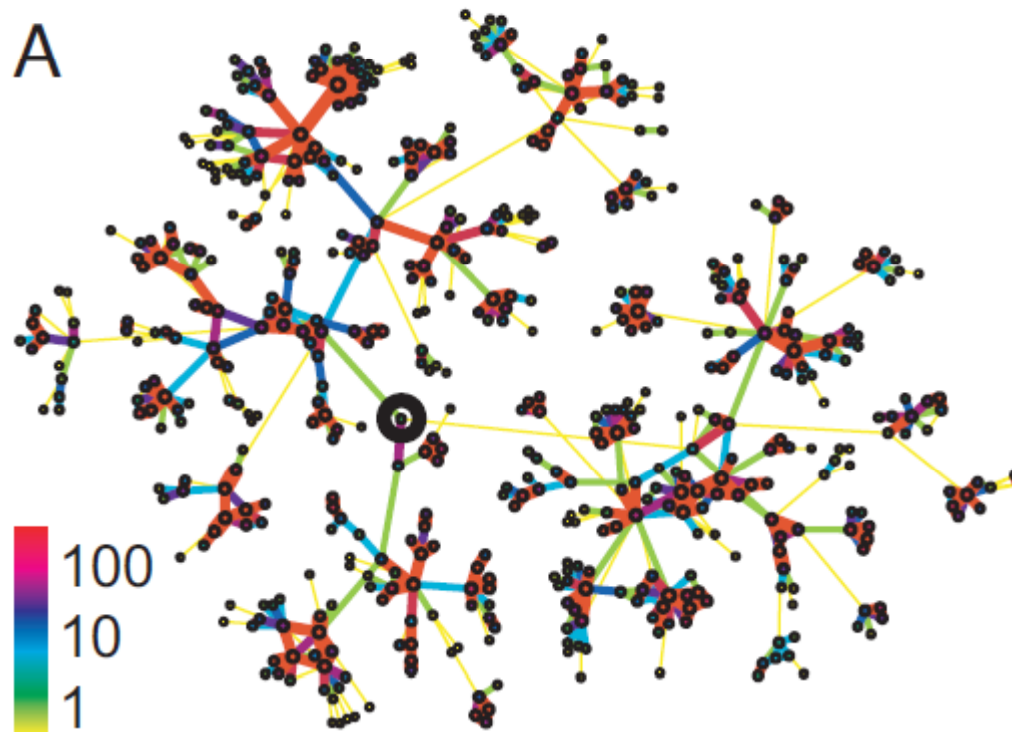
Blogs: Scalability



- **CELF** runs **700** times faster than simple hill-climbing algorithm

Finding communities/clusters in networks

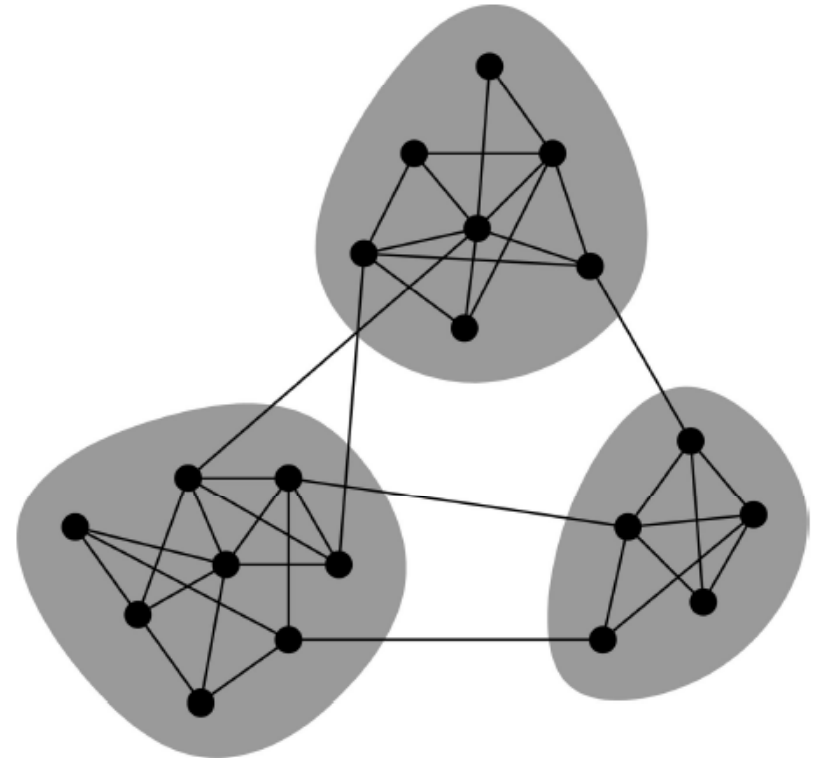
Strength of weak ties



- Real edge strengths in mobile call graph

Network communities

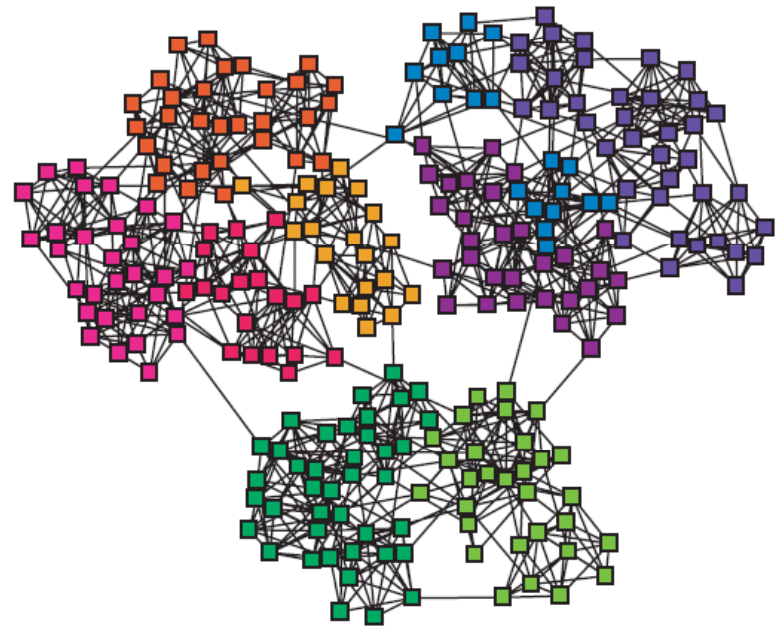
- Findings so far suggest that **network groups are tightly connected**
- **Network communities:**
 - Sets of nodes with **lots** of connections **inside** and **few** to **outside** (the rest of the network)



Communities, clusters,
groups, modules

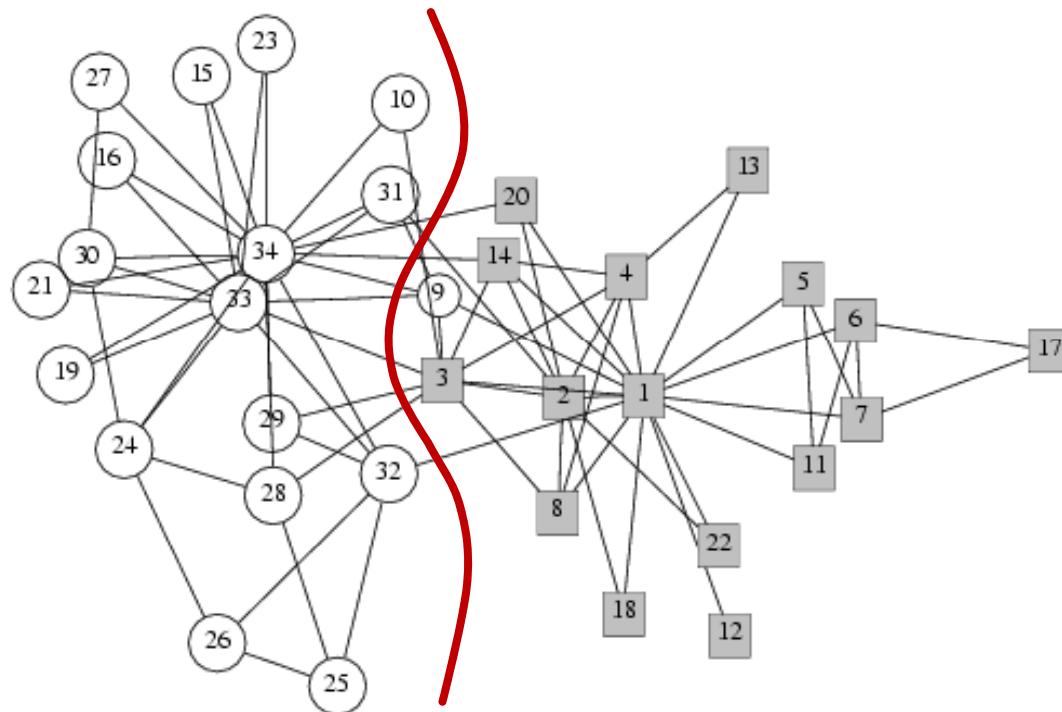
Finding network communities

- How to automatically find such densely connected groups of nodes?
- Ideally such automatically detected clusters would then correspond to real groups
- For example:



Communities, clusters,
groups, modules

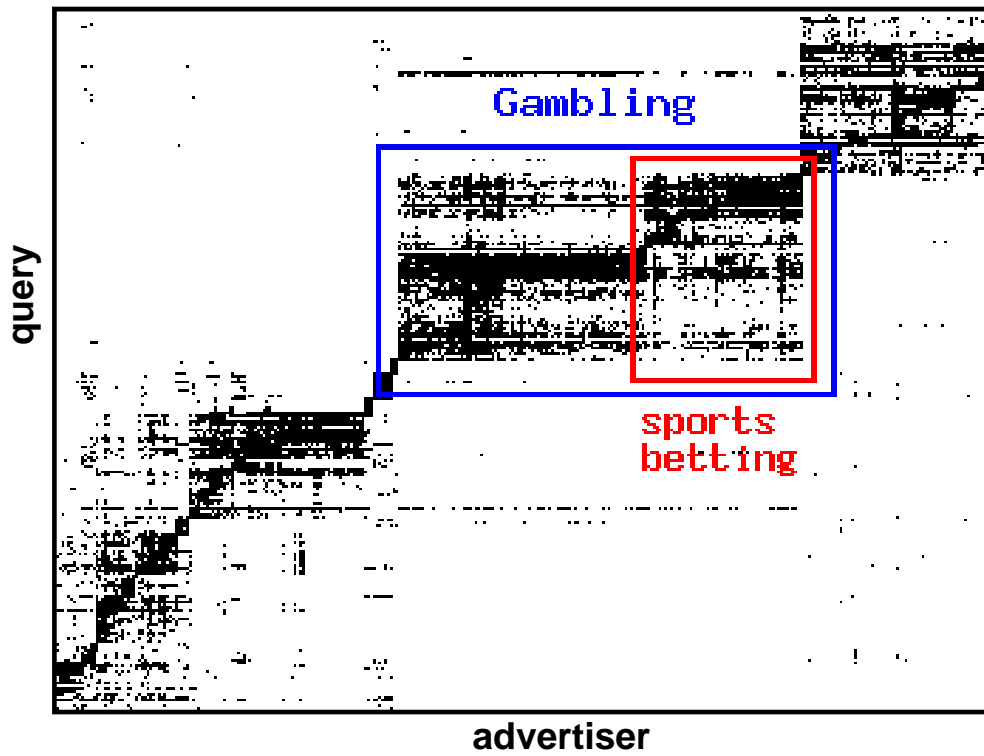
Social Network Data



- Zachary's Karate club network:
 - Observe social ties and rivalries in a university karate club
 - During his observation, conflicts led the group to split
 - Split could be explained by a minimum cut in the network

Micro-markets in sponsored search

Find micro-markets by partitioning the “query x advertiser” graph:



Clustering and Community Finding

Many methods:

- Linear (low-rank) methods:

- If Gaussian, then low-rank space is good

- Kernel (non-linear) methods:

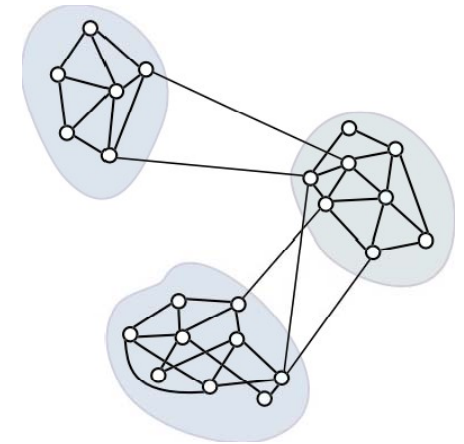
- If low-dimensional manifold, then kernels are good

- Hierarchical methods:

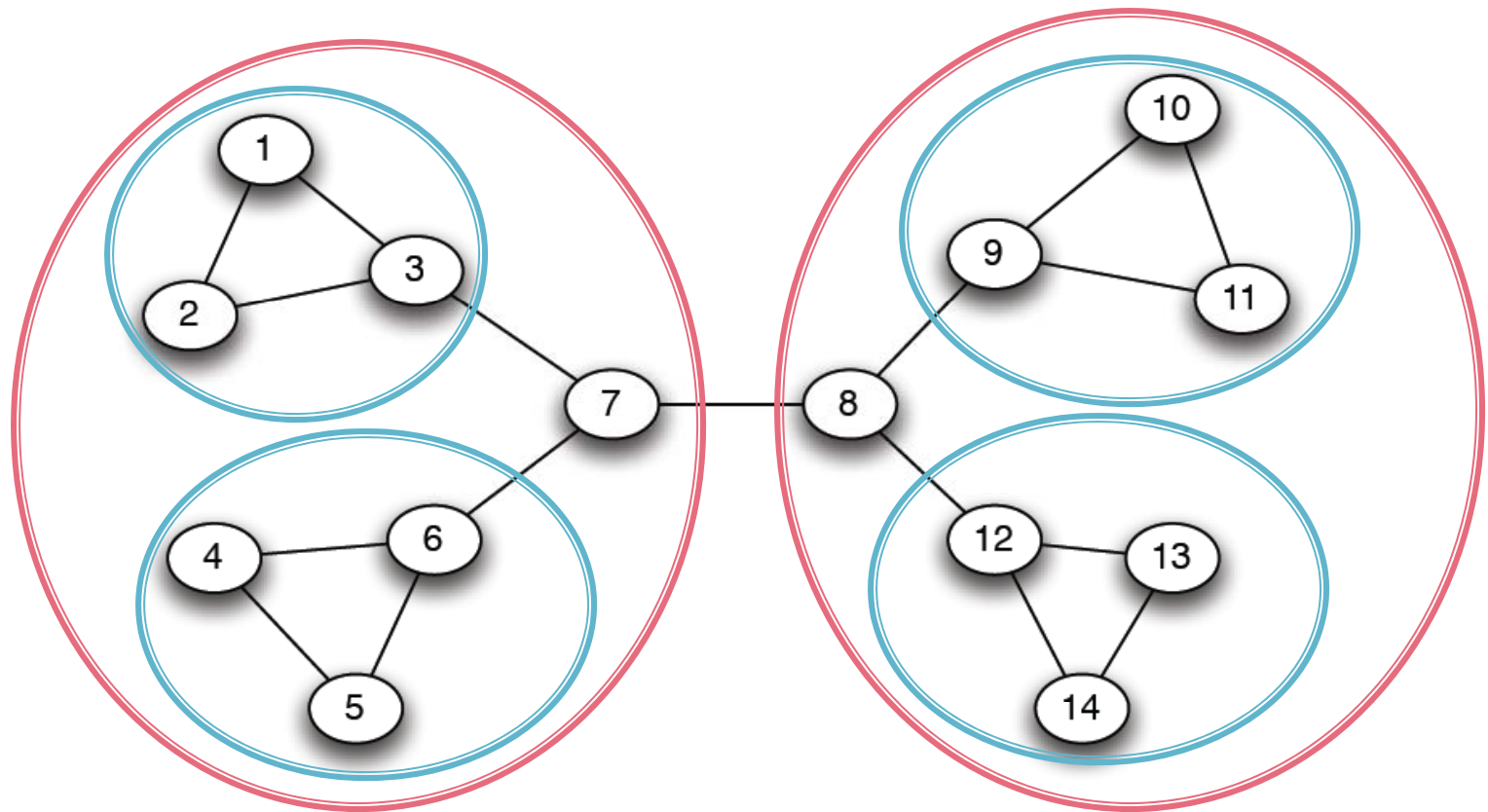
- Top-down and bottom-up – common in social sciences

- Graph partitioning methods:

- Define “edge counting” metric – conductance, expansion, modularity, etc. – and optimize!



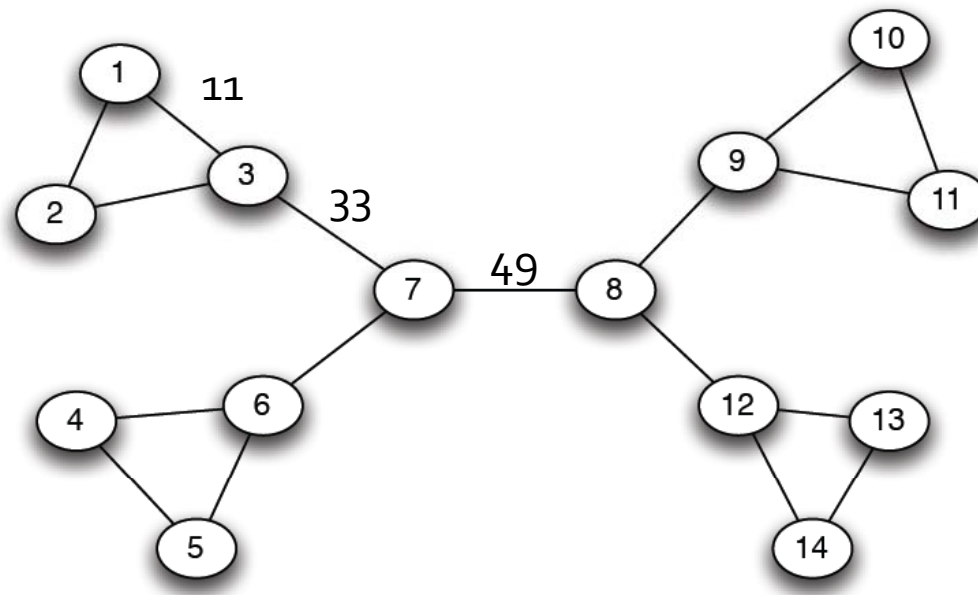
Hierarchically nested communities



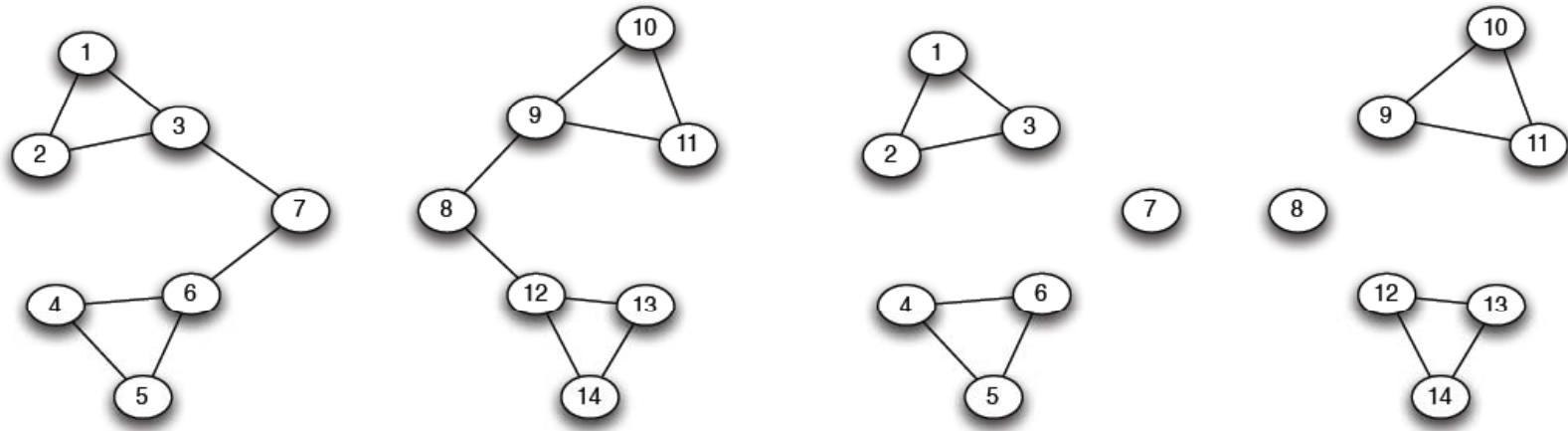
What is a good notion that would extract such clusters?

Algorithm of Girvan-Newman

- Divisive hierarchical clustering based on the notion of edge **betweenness**:
 - Number of shortest paths passing through the edge
- Remove edges in decreasing betweenness

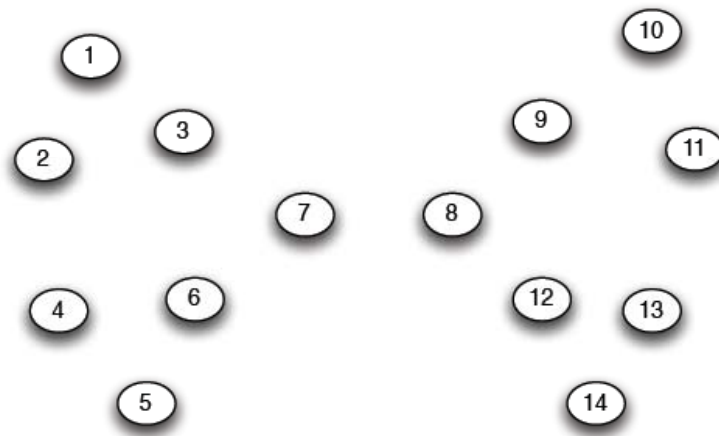


Algorithm of Girvan-Newman



(a) Step 1

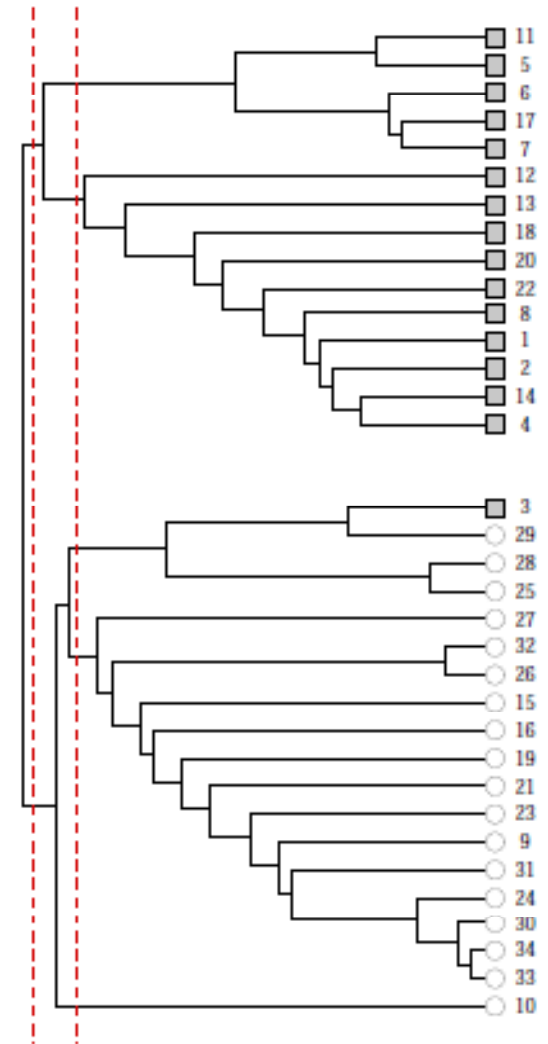
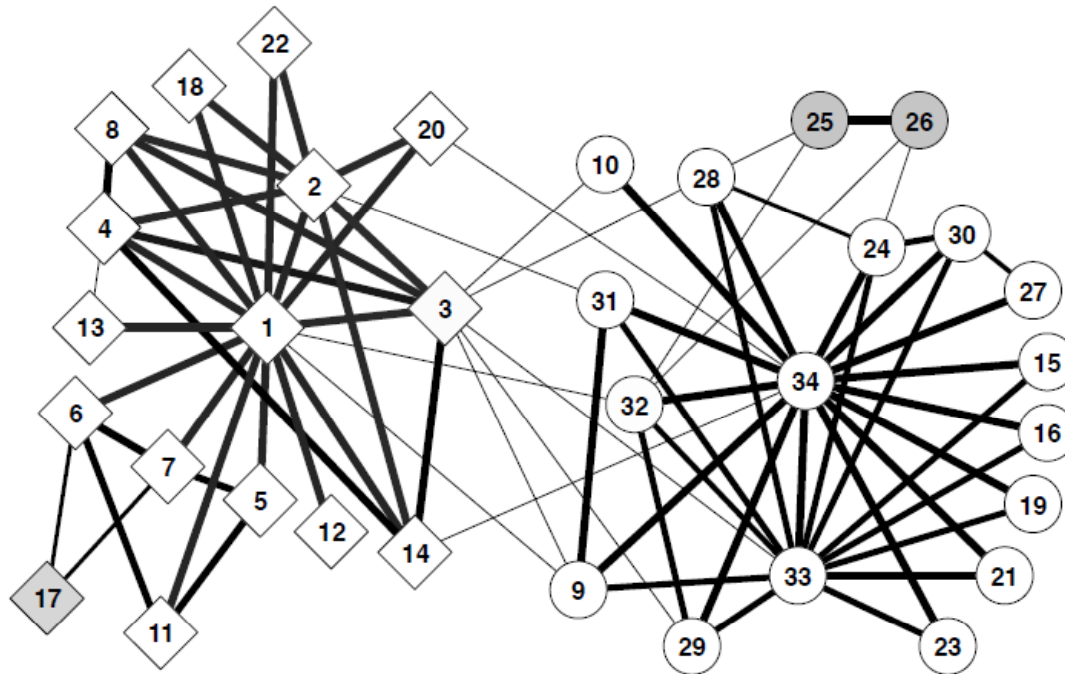
(b) Step 2



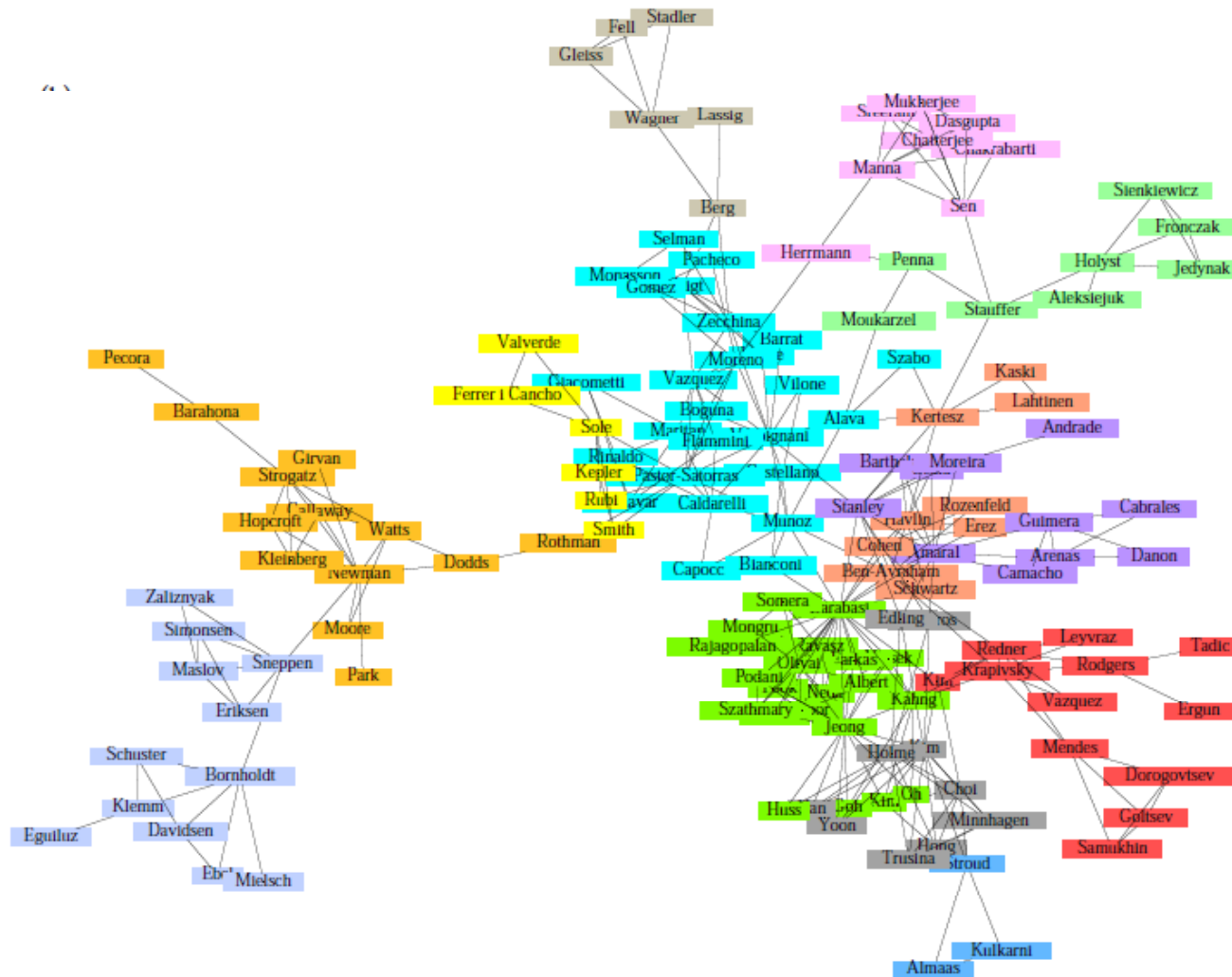
(c) Step 3

Girvan-Newman: Results

- Zachary's Karate club: hierarchic decomposition

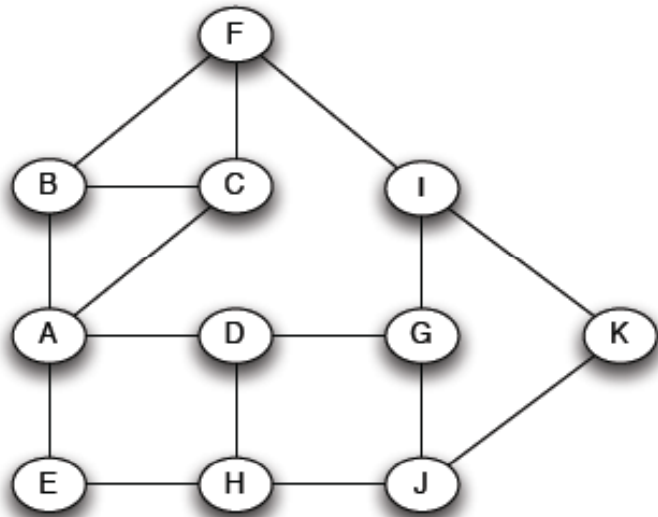


Girvan-Newman: Results



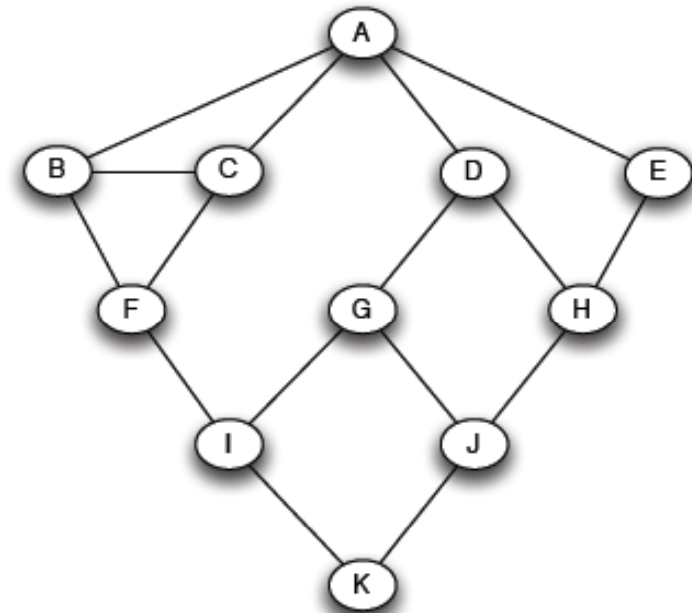
Communities in physics collaborations

How to compute betweenness (1)



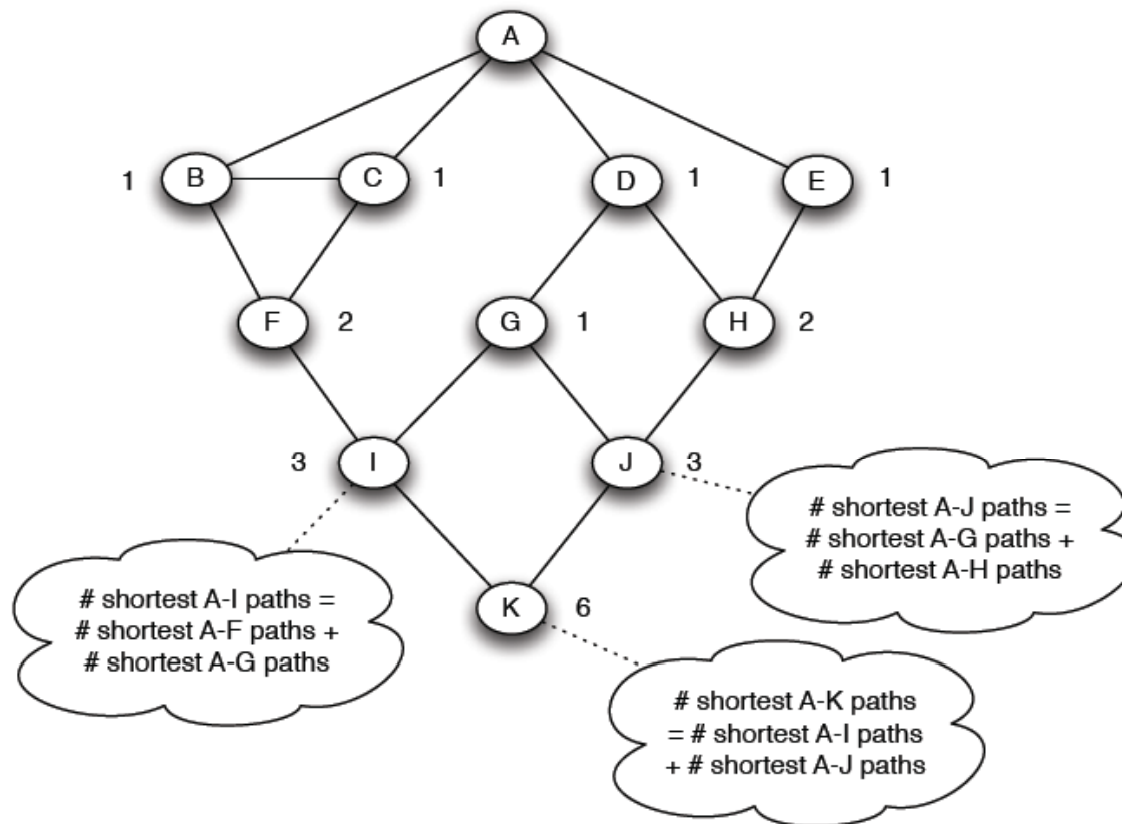
- Want to compute betweenness of paths starting at node A

- Breath first search starting from A:



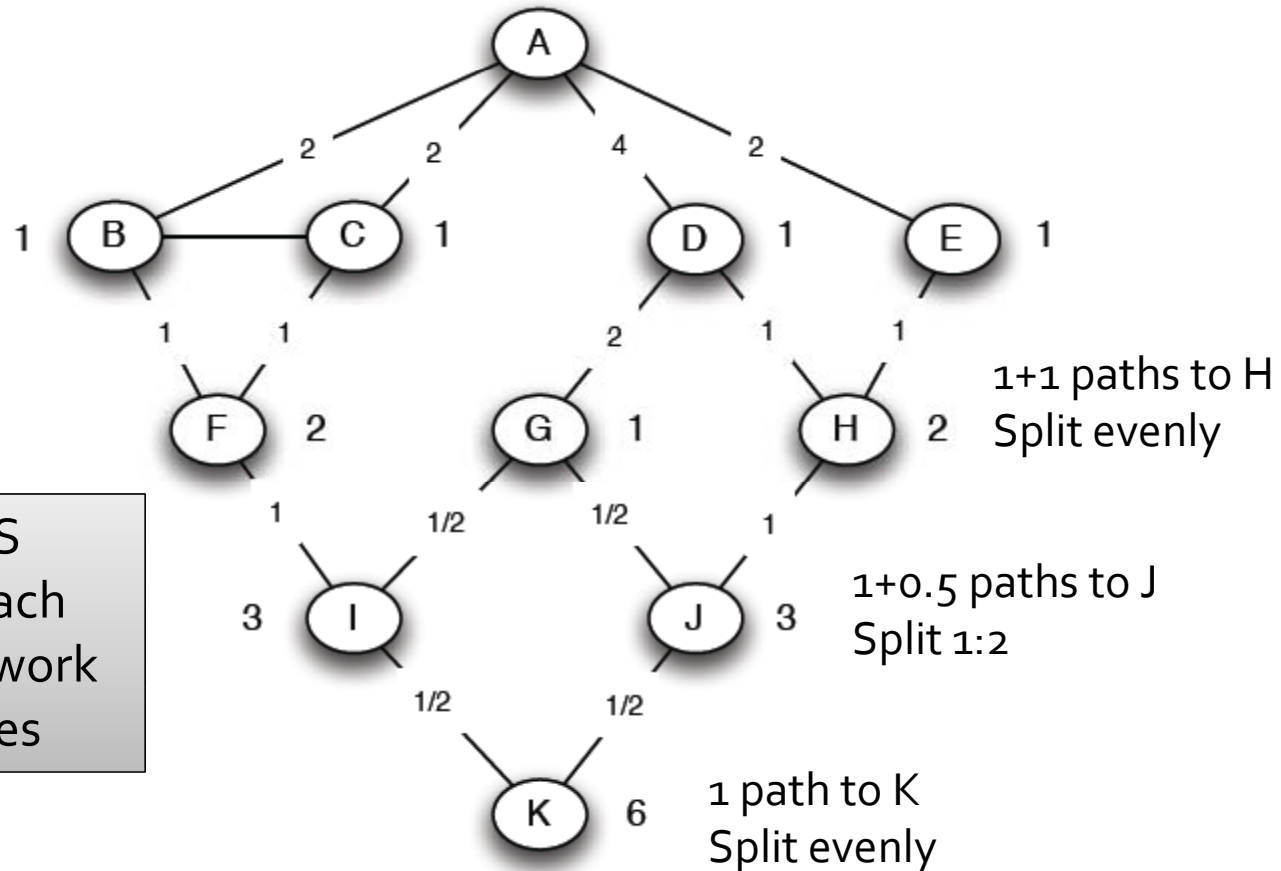
How to compute betweenness (2)

- Count the number of shortest paths from A to all other nodes of the network:



How to compute betweenness (3)

- Compute betweenness by working up the tree: If there are multiple paths count them fractionally



- Repeat the BFS procedure for each node of the network
- Add edge scores