

**Algorithmic Approaches
for Analyzing Large Graphs:
Sampling, Sketching, Streaming**

Andrew McGregor

University of Massachusetts, Amherst

- Big Graphs

Social networks, web, call graphs, biological networks... Graphs are a natural way to encode structural information when we have data about both *basic entities* and their *relationships*.

Can't use classic algorithms on such graphs. Often can't store graph in memory, may be distributed across multiple machines, data may change over time...



What are new
“textbook” algorithms for
modern massive graphs?

- **Tutorial Goals and Caveats**

Present some new algorithmic primitives for large graphs.
Widely applicable technique; try to be platform agnostic.
Won't be comprehensive; will cherry pick illustrative results.
Focus on arbitrary graphs rather than specific applications.
Won't focus on proofs but will give basic outline when it helps convey why certain approaches are effective.

- **Resources**



Survey: SIGMOD Record

<http://people.cs.umass.edu/~mcgregor/papers/graphsurvey.pdf>


Tutorial: Slides and Bibliography

<http://people.cs.umass.edu/~mcgregor/graphs>

Lectures: Ten Lectures on Graph Streams

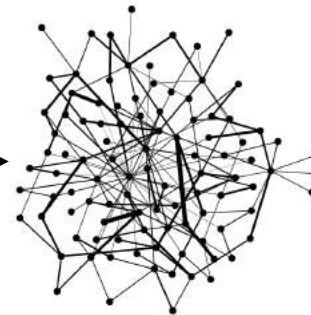
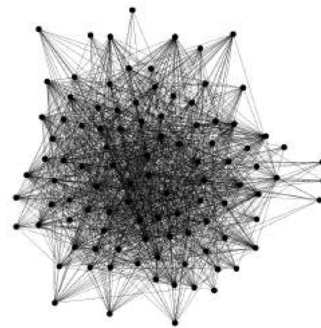
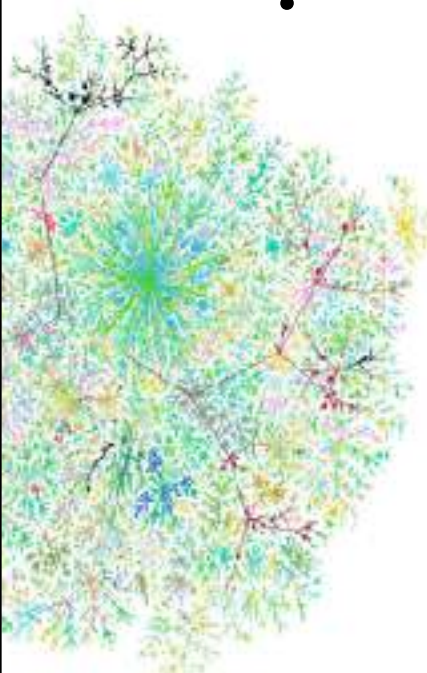
<https://people.cs.umass.edu/~mcgregor/courses/CS711SI8/>

Overview

- 
- Part I: Sampling Sampling for finding densest subgraphs, small matchings, triangles, spectral properties.
“Different sampling techniques for different problems”
 - Part II: Sketching Dimensionality reduction for graph data. Examples include connectivity and cut sparsification.
“Homomorphic compression: sketch first and then run algorithms on the sketched data”
 - Part III: Streaming What can you compute in limited memory with only a few passes over the edges.
“A little inspiration yields a lot less iteration”

Recurring Theme

? What's appropriate notion of lossy compression for graphs?



- If compression is easy, we get faster and more-space efficient algorithms by using existing algorithms on compressed graphs.



Part I

Sampling

Uniform Sampling + Densest Subgraph

Snape Sampling + Matching

Monochromatic Sampling + Clustering Coefficient

Edge-Weighted Sampling + Cuts and Sparsification



Part I

Sampling

Uniform Sampling + Densest Subgraph

Snape Sampling + Matching

Monochromatic Sampling + Clustering Coefficient

Edge-Weighted Sampling + Cuts and Sparsification

- Given a graph G , the *density* of a set of nodes $S \subset V$ is:

$$D_S = \frac{\# \text{ of edges with both endpoints in } S}{\# \text{ of nodes in } S}$$

- ***Problem*** Estimating $D^* = \max_S D_S$ is a basic graph problem with numerous applications. Studied in a variety of models.

See tutorial *Gionis, Tsourakakis* [KDD 15]

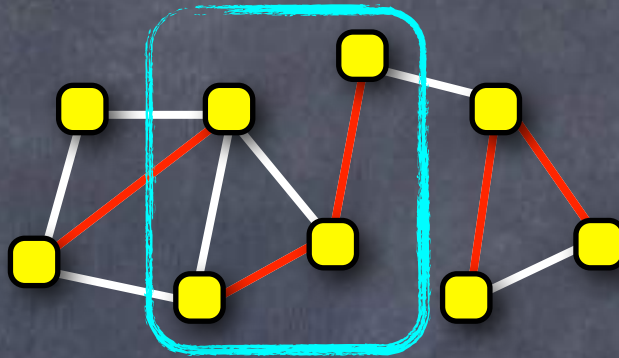
- ***Thm*** Sample $\tilde{O}(\varepsilon^{-2} n)$ edges uniformly and find the densest subgraph in sampled graph. Gives a $(1+\varepsilon)$ -approx whp.

McGregor et al. [MFCS 15], *Esfandiari et al.* [SPAA 16]

Mitzenmacher et al. [KDD 15]

Why **Uniform Sampling** Works...

- Essentially sampling each edge w/p $p \approx \varepsilon^{-2}n/m$.
- Let D'_S be density of S in sampled graph.



$$D_S = 1.0$$
$$D'_S = 0.5$$

- Chernoff:** For each S , $D_S = D'_S/p \pm \varepsilon D^*$ w/p $1 - n^{-3|S|}$
- Union Bound:** Bound applies for all S w/p $1 - n^{-1}$
 - There are $\leq n^k$ subsets of k nodes. So bound fails for some subset of size k w/p $\leq n^k n^{-3k} = n^{-2k}$
 - Bound fails for some subset w/p $\leq n^{-2} + n^{-4} + \dots + n^{-2n} \leq n^{-1}$
- So max density of sampled graph gives $1 + \varepsilon$ approx.



Part I

Sampling

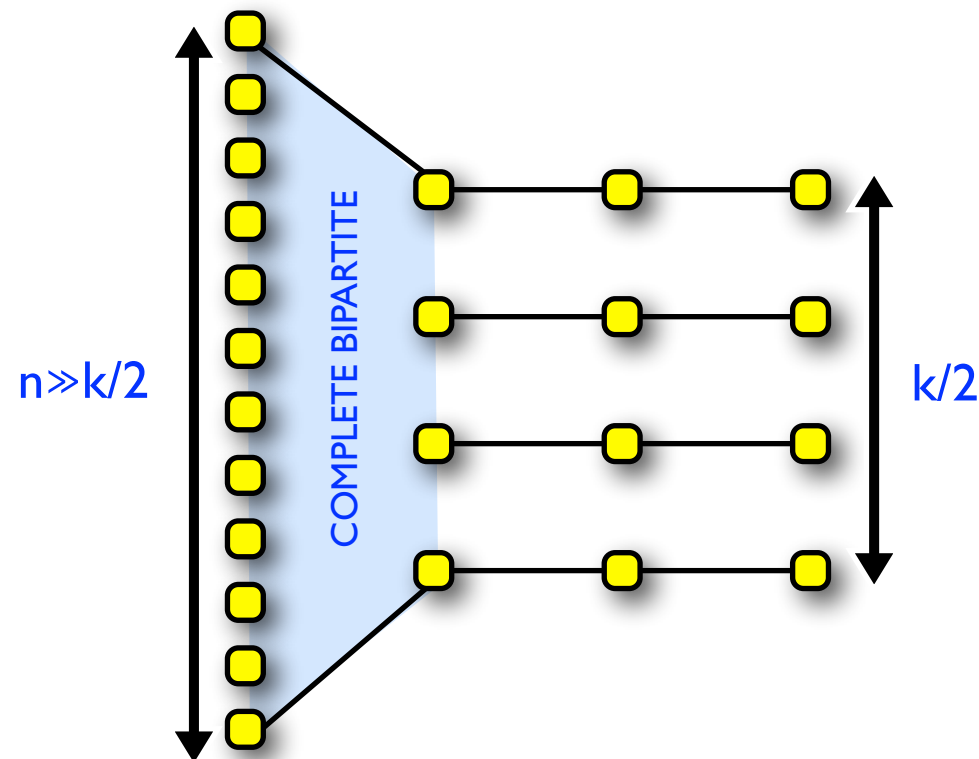
Uniform Sampling + Densest Subgraph

Snape Sampling + Matching

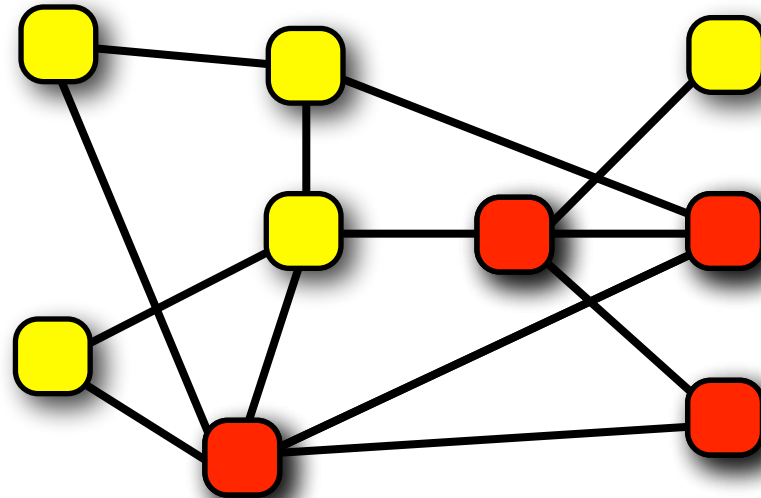
Monochromatic Sampling + Clustering Coefficient

Edge-Weighted Sampling + Cuts and Sparsification

- Matching Problem Find large set of edges such that no two edges share an endpoint.
- How many “samples” are needed to find a matching of size k ?
- Sampling uniformly can be very inefficient...



- SNAPE “Sample Nodes And Pick Edge” Sampling:
- **Sample** each node with probability $1/k$ and **delete** rest
- **Pick** a random edge amongst those that remain.



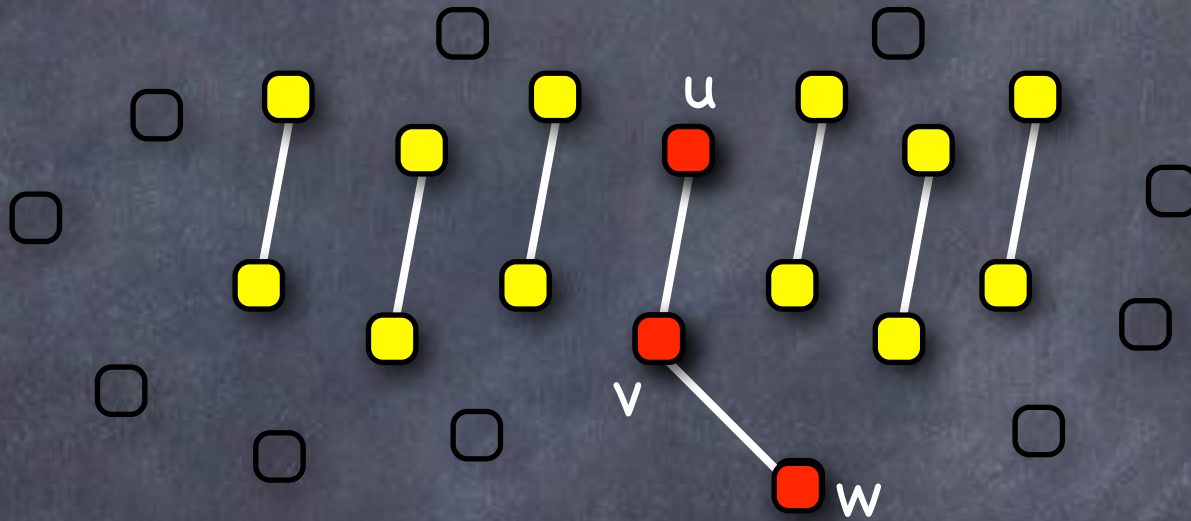
- Theorem If G has max matching size k , then $O(k^2 \log k)$ SNAPE samples will include a max matching from G .

Chitnis et al. [SODA 16], Bury et al. [Algorithmica 18]



Why **SNAPE** Sampling Works...

- Consider a maximum matching M of size k and focus on arbitrary edge $\{u,v\}$ in this matching.



- W/p $\Omega(k^{-2})$ u and v only endpoints of M sampled.
- Hence, when we pick one of the remaining edges it's either $\{u,v\}$ or another edge that's equally useful.
- Take $O(k^2 \log k)$ samples; apply analysis to all edges.



Part I

Sampling

Uniform Sampling + Densest Subgraph

Snape Sampling + Matching

Monochromatic Sampling + Clustering Coefficient

Edge-Weighted Sampling + Cuts and Sparsification

- Given a graph G , the *global clustering coefficient* is

$$\kappa = \frac{3 \times \text{number of triangles}}{\text{number of length 2 paths}}$$

A measure of how much nodes tend to cluster together.

- *Monochromatic Sampling* Randomly color each node from a set of colors. Store all edges with monochromatic endpoints. If length-2 path $\{u,v\}, \{v,w\}$ is stored, $\{u,w\}$ also stored if it exists.
- *Thm* Can additively estimate κ from $\tilde{O}(\sqrt{n})$ samples.

Pagh, Tsourakakis [IPL 12], Jha, Seshadhri, Pinar [KDD 15]

Proof of Lemma

- Let $k = \#$ colors and $T_S = \#$ triangles sampled.
- Edges Sampled: mp
- Expectation: $E[T_S] = Tp^2$ where $p = 1/k$ and T is $\#$ triangles
- Variance: Let $W = \#$ length 2-paths

$$\text{Var}[T_S] \leq Tp^2 + 10 \cdot W^{3/2} p^3$$

- Chebyshev Bound:

$$\Pr[|T_S/p^2 - T| > \epsilon W] \leq \frac{\text{Var}[T_S]}{\epsilon^2 p^4 W^2} \leq \frac{O(1)}{p \epsilon^2 \sqrt{W}}$$

- $k \approx \epsilon^2 \sqrt{(m^2/n)}$ for 0.01 error prob. and $O(\epsilon^{-2} \sqrt{n})$ samples



Part I

Sampling

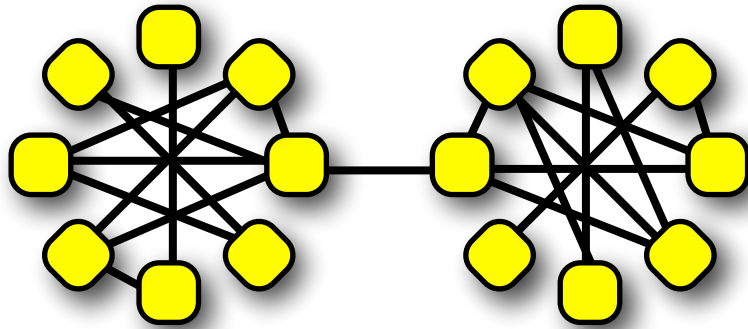
Uniform Sampling + Densest Subgraph

Snape Sampling + Matching

Monochromatic Sampling + Clustering Coefficient

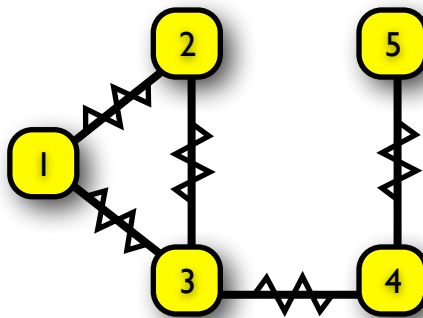
Edge-Weighted Sampling + Cuts and Sparsification

- **Defn** A *sparsifier* of graph G is a weighted subgraph H with:
 - \forall cuts: “size of cut in G ” = $(1 \pm \epsilon)$ “size of cut in H ”
- **Basic Approach** Sample each edge uv with probability p_{uv} and reweight by $1/p_{uv}$. Probabilities depend on edge properties...



- **Thm** If $p_{uv} \approx \epsilon^{-2}/\lambda_{uv}$ or $p_{uv} \approx \epsilon^{-2}r_{uv}$ then result is sparsifier with $\tilde{O}(\epsilon^{-2} n)$ edges. *Fung et al. [STOC 11], Spielman, Srivastava [STOC 08]*

λ_{uv} is the min number of edges whose removal disconnects u and v



r_{uv} is potential difference when unit of flow injected at u and extracted at v

- **Simpler Thm** If min-cut is $\gg \epsilon^{-2} \log n$ then $p_e = 1/2$ works.

Proof Idea of **Simpler Theorem** ...

- **Lemma (Chernoff)** Let k' be number of edges that were sampled across some cut of size k . Then

$$\Pr[k' \neq (1 \pm \epsilon) k/2] \leq \exp(-\epsilon^2 k/6)$$

- **Lemma (Karger)** The number of cuts with k edges is $< \exp(2k \log n / \lambda)$ where λ is size of min-cut.
- Result then follows by substituting bound for λ and applying union bound over all cuts.



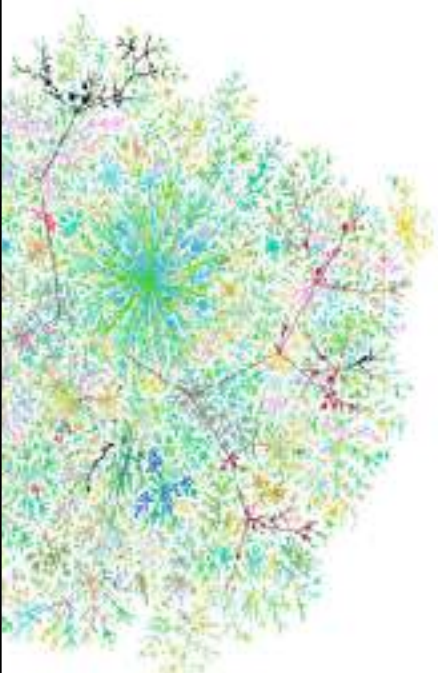
Part II

Sketching

What is sketching?

Surprising connectivity example

Revisiting graph cuts and sparsification



Part II

Sketching

What is sketching?

Surprising connectivity example
Revisiting graph cuts and sparsification

Basic Idea for l_0 sampling...

- Entry in i^{th} row of M is 1 w/p 2^{-i+1} . Some entry of Ma probably corresponds to single entry of a .

$$\begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} x \\ 0 \\ 0 \\ y \\ z \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} x + y + z \\ x + z \\ y \\ 0 \end{pmatrix}$$

← Too many
← Too many
← Just right
← Too few

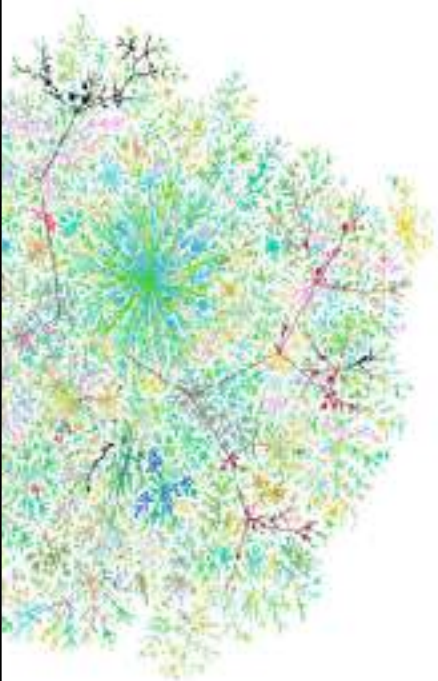
- How do we know when this happens? How do we know which entry of the vector was isolated?

More Details for l_0 sampling...

- Take 3 copies of M : Replace 1's in 2nd copy by #col. Replace 1's in 3rd copy by $r^{\text{\#col}}$ where r is random.

$$\begin{pmatrix}
 \boxed{1} & \boxed{1} & \boxed{1} & \boxed{1} & \boxed{1} & \boxed{1} & \boxed{1} & \boxed{1} \\
 \boxed{1} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{1} & \boxed{1} & \boxed{1} & \boxed{0} \\
 \boxed{0} & \boxed{0} & \boxed{0} & \boxed{1} & \boxed{0} & \boxed{0} & \boxed{1} & \boxed{0} \\
 \boxed{0} & \boxed{1} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0} \\
 \hline
 \boxed{1} & \boxed{2} & \boxed{3} & \boxed{4} & \boxed{5} & \boxed{6} & \boxed{7} & \boxed{8} \\
 \boxed{1} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{5} & \boxed{6} & \boxed{7} & \boxed{0} \\
 \boxed{0} & \boxed{0} & \boxed{0} & \boxed{4} & \boxed{0} & \boxed{0} & \boxed{7} & \boxed{0} \\
 \boxed{0} & \boxed{2} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0} \\
 \hline
 \boxed{r^1} & \boxed{r^2} & \boxed{r^3} & \boxed{r^4} & \boxed{r^5} & \boxed{r^6} & \boxed{r^7} & \boxed{r^8} \\
 \boxed{r^1} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{r^5} & \boxed{r^6} & \boxed{r^7} & \boxed{0} \\
 \boxed{0} & \boxed{0} & \boxed{0} & \boxed{r^4} & \boxed{0} & \boxed{0} & \boxed{r^7} & \boxed{0} \\
 \boxed{0} & \boxed{r^2} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0} & \boxed{0}
 \end{pmatrix}
 \begin{pmatrix}
 x \\
 0 \\
 0 \\
 y \\
 z \\
 0 \\
 0 \\
 0
 \end{pmatrix}
 =
 \begin{pmatrix}
 x + y + z \\
 x + z \\
 \boxed{y} \\
 0 \\
 x + 4y + 5z \\
 x + 5z \\
 \boxed{4y} \\
 0 \\
 rx + r^4y + r^5z \\
 rx + r^5z \\
 \boxed{r^4y} \\
 0
 \end{pmatrix}
 \begin{matrix}
 = u \\
 \\
 = v \\
 \\
 = w
 \end{matrix}$$

- Lemma:** With high probability, if $w = r^{v/u} u$ then we've isolated non-zero element in position v/u .



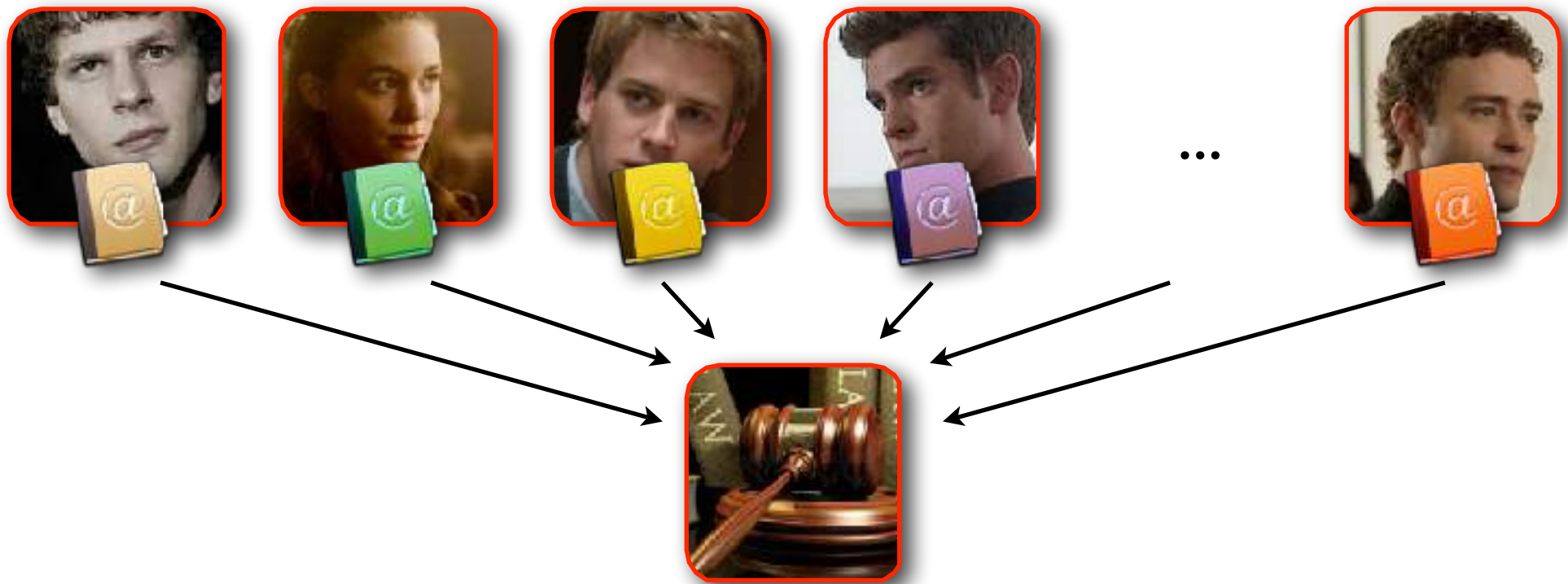
Part II

Sketching

What is sketching?

Surprising connectivity example

Revisiting graph cuts and sparsification

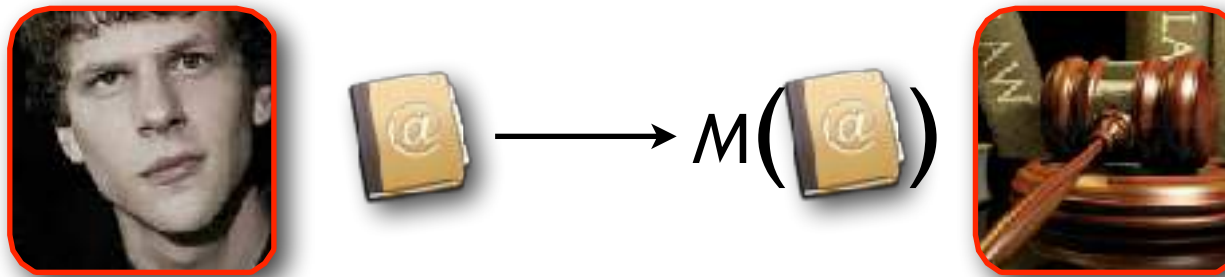


- Communication Problem n players each have a list their friends. Simultaneously, they each send a message to a central player who deduces if underlying graph is connected.
- Thm $O(\text{polylog } n)$ bit message from each player suffices.

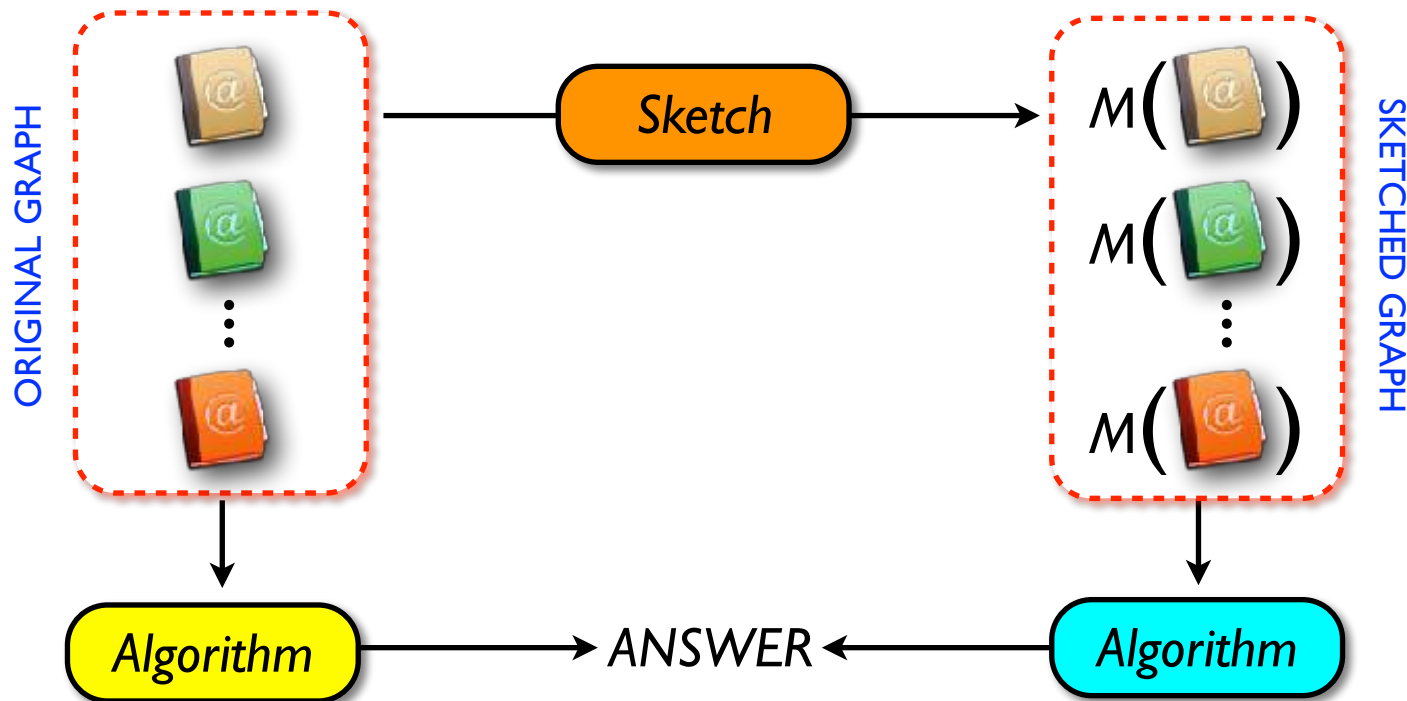
Ahn, Guha, McGregor [SODA 12]



- *Can't be possible!* What if there's a *bridge* (u,v) in the graph, i.e., a friendship that is critical to ensuring the graph is connected.
- It *appears* that at least one player needs to send $\Omega(n)$ bits.
 - Central player needs to know about the special friendship.
 - Participant doesn't know which friendships are special.
 - Participants may have $\Omega(n)$ friends.

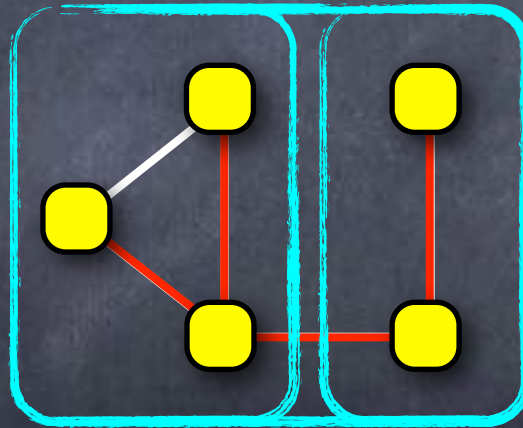


- Players send carefully-designed sketches of address books.
- Homomorphic Compression Instead of running algorithm on original data, run algorithm on sketched data.



Ingredient 1: Basic Algorithm

- Algorithm (Spanning Forest):
 - For each node: pick incident edge
 - For each connected component: pick incident edge
 - Repeat until no edges between connected comp.

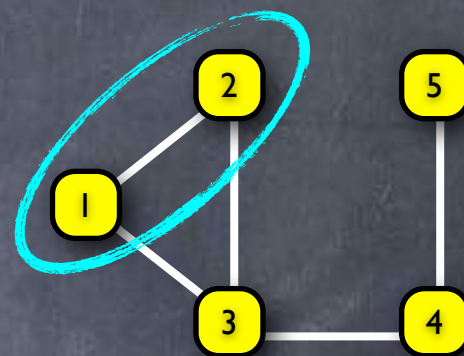


- Lemma** After $O(\log n)$ rounds selected edges include spanning forest.

Ingredient 2: Sketching Neighborhoods

- For node i , let \mathbf{a}_i be vector indexed by node pairs. Non-zero entries: $\mathbf{a}_i[i,j]=1$ if $j>i$ and $\mathbf{a}_i[i,j]=-1$ if $j<i$.

$$\begin{array}{l} \mathbf{a}_1 = \begin{pmatrix} \{1,2\} & \{1,3\} & \{1,4\} & \{1,5\} & \{2,3\} & \{2,4\} & \{2,5\} & \{3,4\} & \{3,5\} & \{4,5\} \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \\ \mathbf{a}_2 = \begin{pmatrix} -1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \\ \mathbf{a}_1 + \mathbf{a}_2 = \begin{pmatrix} 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \end{array}$$



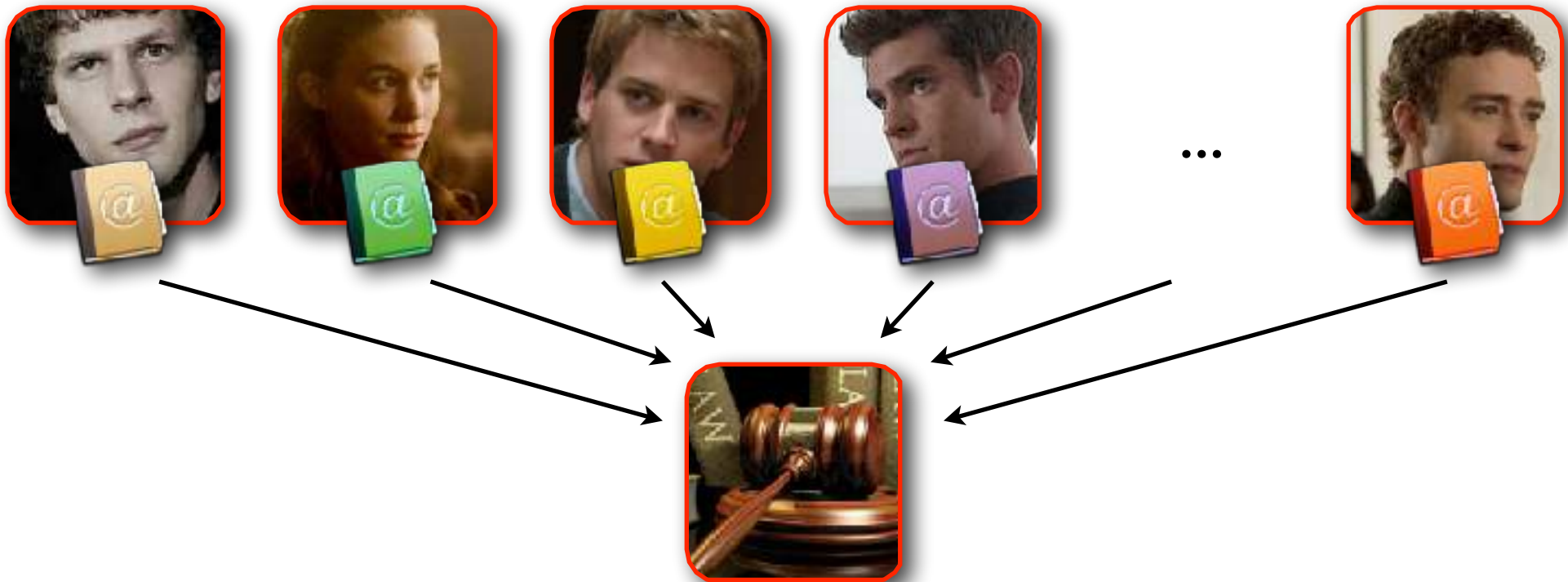
- Lemma** For any subset of nodes $S \subset V$, non-zero entries of $\sum_{j \in S} \mathbf{a}_j$ are edges across cut $(S, V \setminus S)$
- Player j sends $M(\mathbf{a}_j)$ where M is " l_0 sampling" sketch.

Recipe: Sketch & Compute on Sketches

- **Player with Address Books:** Player j sends Ma_j
- **Central Player: "Runs Algorithm in Sketch Space"**
 - Use Ma_j to get incident edge on each node j
 - For $i=2$ to $\log n$:
 - To get incident edge on component $S \subset V$ use:

$$\sum_{j \in S} Ma_j = M\left(\sum_{j \in S} a_j\right) \longrightarrow \text{non-zero elt. of } \sum_{j \in S} a_j = \text{edge across cut}$$

Detail: Actually each player sends $\log n$ independent sketches M_1a_j, M_2a_j, \dots and central player uses M_ia_j when emulating i^{th} iteration of the algorithm.



- ***Thm*** $O(\text{polylog } n)$ bit message from each player suffices.
- ***Extensions*** $\tilde{O}(k)$ bit messages for k -edge connectivity and $\tilde{O}(k^2)$ bit messages for k -node connectivity.

Extending to k -Edge-Connectivity

Basic Algorithm

- **Algorithm:** For $i=1$ to k :
 Let F_i be spanning forest of $G(V, E - F_1 - \dots - F_{i-1})$
 - **Lemma:** $F_1 + \dots + F_k$ is k -edge-connected iff G is.
-

Emulation in Sketch Space

- **Sketch:** Simultaneously construct k independent connectivity sketches $M_1(G), M_2(G), \dots, M_k(G)$.
- **Run Algorithm in Sketch Space:**
 - Use $M_1(G)$ to find a spanning forest F_1 of G
 - Use $M_2(G) - M_2(F_1) = M_2(G - F_1)$ to find F_2
 - Use $M_3(G) - M_3(F_1) - M_3(F_2) = M_3(G - F_1 - F_2)$ to find $F_3 \dots$
- **Extension:** Can recover a set of "weak" edges whose removal leaves connected components with min-cut $> k$.

k-Node-Connectivity

Basic Algorithm

- Algorithm:

Let H_1, H_2, \dots, H_r where $r = k^3 \log n$ and H_i is induced subgraph on random set of n/k nodes.

Let F_i be a spanning forest of H_i .

- Lemma:** Whp, $F_1 + F_2 + \dots + F_r$ is k -node connected iff G is.

Emulation in Sketch Space

- Sketch:** Construct connectivity sketch for each H_i , and use this to find F_i .



Part II

Sketching

What is sketching?

Surprising connectivity example

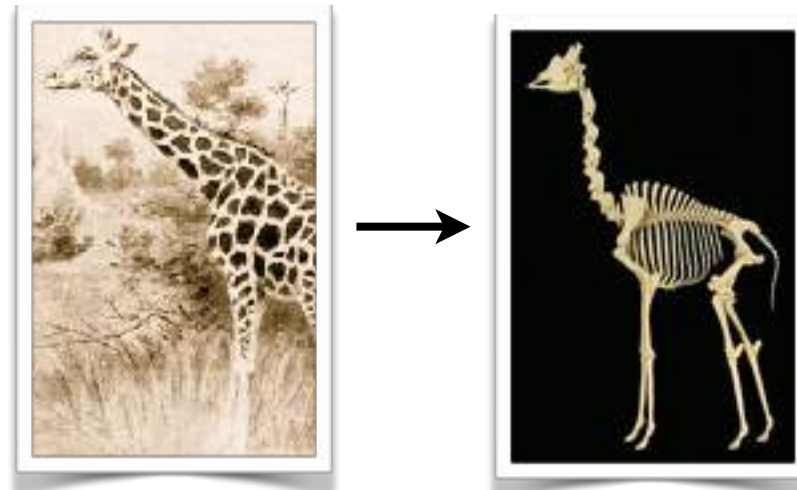
Revisiting graph cuts and sparsification

- Thm $O(\epsilon^{-2} \text{polylog } n)$ bit messages suffice for central player to construct sparsifier and approx all graph cuts.

Guha, McGregor, Tench [PODS 15], Kapralov et al. [STOC 14]

- Main Ideas

1. For a graph G , can find all edges in small cuts.
2. For large cuts, suffices to sample edges with prob. $1/2$.
3. So, sparsifying G reduces to sparsifying sampled graph G' .
4. To sparsify G' recurse... Can do recursion in parallel.





Part III

Streaming

Revisiting Matching

Correlation Clustering

Coloring Graphs

Coverage and Submodular Maximization

- Two Main Graph Stream Models
- *Insert-Only Model*: Input is a stream of edges.
- *Insert-Delete Model*: Edge insertions and edge deletions.



[Mark](#) and [Erica](#) are now friends.



Like · Add Friend

- Goal Using small memory, compute properties of the graph.
- All the earlier algorithms apply in insert-delete model:
 - Maintain sketch Mx where x is characteristic vector of edges.
 - When e inserted, update sketch $Mx \leftarrow Mx + (e^{\text{th}} \text{ column of } M)$
 - When e deleted, update sketch $Mx \leftarrow Mx - (e^{\text{th}} \text{ column of } M)$

- Results so far: One-pass dynamic data stream algorithms for:
 - a) $1+\epsilon$ approx of **densest subgraph** in $\tilde{O}(\epsilon^{-2}n)$ space.
 - b) $1+\epsilon$ approx of **all cuts** in $\tilde{O}(\epsilon^{-2}n)$ space via cut sparsification.
 - c) Additive ϵ approx of **clustering coefficient** in $\tilde{O}(\epsilon^{-2}\sqrt{n})$ space.
 - d) Finding **matching and vertex cover** of size k in $\tilde{O}(k^2)$ space.
 - e) **k -edge connectivity** in $\tilde{O}(nk)$ space.
 - f) **k -node connectivity** in $\tilde{O}(nk^2)$ space.
- Many of above space bounds can be shown to be optimal via reductions from communication complexity.
- Many other algorithms... some only for edge insertions or use multiple passes or only work for certain types of graphs.



Part III

Streaming

Revisiting Matching

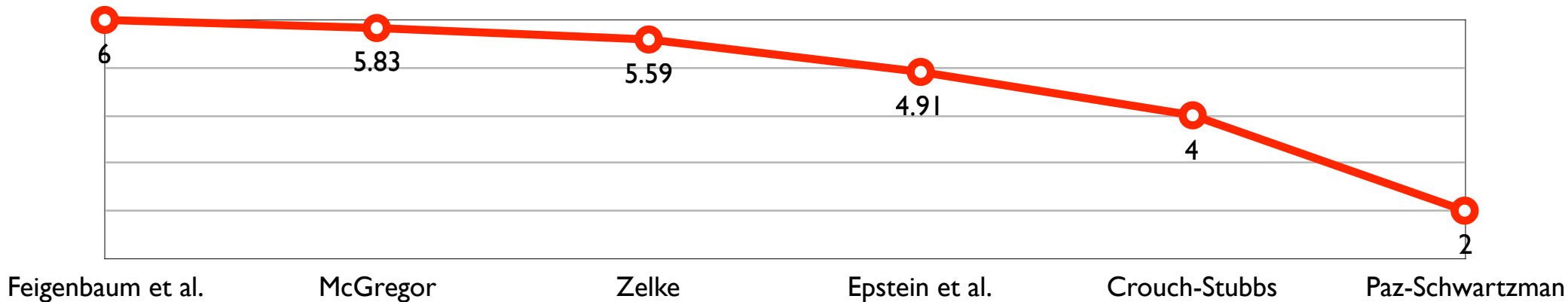
Correlation Clustering

Coloring Graphs

Coverage and Submodular Maximization

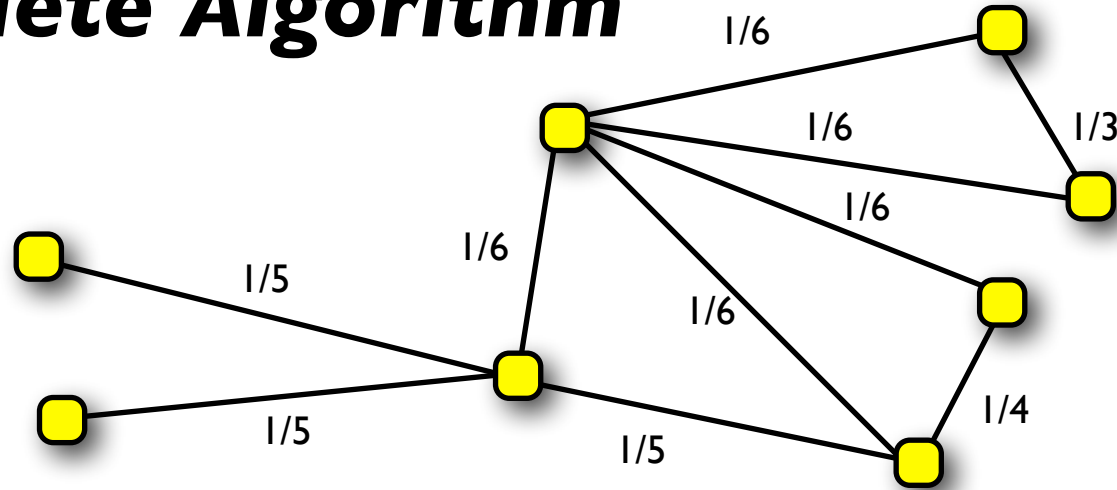
- Unweighted Matching Greedy algorithm returns 2-approx using $\tilde{O}(n)$ space. Embarrassingly, this is best known one-pass result!

Approximation Ratios for Weighted Matching



- Weighted Matching $2+\epsilon$ approx in $\tilde{O}(n/\epsilon)$ space.
Paz, Schwartzman [SODA 17], Ghaffari, Wajc [SOSA 19]
- ? Improve result for sparse graphs? Graph has **arboricity** α if all subgraphs have average degree $< \alpha$. Planar graph has $\alpha=3$.

Insert-Delete Algorithm

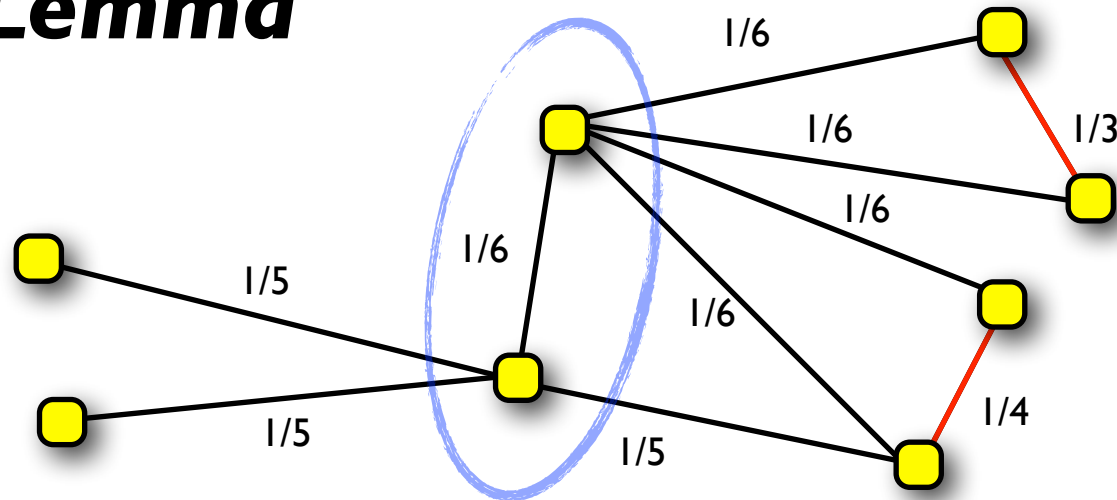


- **Lemma:** $\text{match}(G)/(2+\alpha) \leq A \leq \text{match}(G)$ where A is total edge weight if each edge uv gets weight

$$x_{uv} = \min \left(\frac{1}{\deg(u) + 1}, \frac{1}{\deg(v) + 1} \right)$$

- **Thm:** Can $2+\alpha+\epsilon$ approximate $\text{match}(G)$ using $\tilde{O}(n^{4/5})$ space:
 If $\text{match}(G) \leq n^{2/5}$, can find exactly using earlier algorithm.
 Otherwise, evaluate A on random set of $\approx n^{4/5}$ nodes.

Proof of Lemma



- The edge weights are a **fractional matching**, i.e., for any node u :

$$\sum_{v \in \Gamma(u)} x_{uv} \leq \sum_{v \in \Gamma(u)} \frac{1}{\deg(u) + 1} < 1$$

- **To prove total weight $\leq \text{match}(G)$:** Use Edmond's matching polytope thm since weight on subgraph of r nodes is $\leq (r-1)/2$.
- **To prove total weight $\geq \text{match}(G)/(2+\alpha)$:**

Total weight of edges incident to “high degree” vertices H at least $|H|/(2+\alpha)$ and all other weights are at least $1/(2+\alpha)$.

Matching size is at most $|H| + \text{“edges not incident to } H\text{”}$

Insert-Only Algorithm

- Thm $\alpha+2+\epsilon$ approx of matching size in $O(\text{polylog } n)$ space.
Cormode et al. [ESA 17], McGregor, Vorotnikova [SOSA 18]
- Define Edge $\{u,v\}$ is **special** if $\leq \alpha$ edges incident to u and $\leq \alpha$ edges incident to v later than $\{u,v\}$. Let s be # special edges.
- Lemma $\text{match}(G) \leq s \leq (2+\alpha)\text{match}(G)$.
- Proof Ingredients Graph of special edges has degrees $\leq \alpha+1$. Low arboricity bounds number of almost special edges.
- Algorithm Estimate s up to a factor $1+\epsilon$
 - a) Suppose we have guess g that is 2-approximates s
 - b) Sample each edge w/p $\approx \epsilon^{-2} (\log n)/g$. If you subsequently see $>\alpha$ edges incident to either endpoint, drop the edge.
- Can show a) the current sample size is always small and b) size of final sample and g yields good approx for s .

Final Optimized Algorithm...

1. Initialize $S \leftarrow \emptyset$, $p = 1$, estimate = 0
2. For each edge $e = uv$ in the stream:
 - a. With probability p add e to S and initialize counters $c_e^u \leftarrow 0$ and $c_e^v \leftarrow 0$
 - b. For each edge $e' \in S$, if e' shares endpoint w with e :
 - Increment $c_{e'}^w$
 - If $c_{e'}^w > \alpha$, remove e' and corresponding counters from S
 - c. If $|S| > 40\epsilon^{-2} \log n$:
 - $p \leftarrow p/2$
 - Remove each edge in S and corresponding counters with probability $1/2$
 - d. estimate $\leftarrow \max(\text{estimate}, |S|/p)$
3. Return estimate

- **Best Case Scenario** You put a lot of work/maths into designing and analyzing an algorithm but the final algorithm is simple.



Part III

Streaming

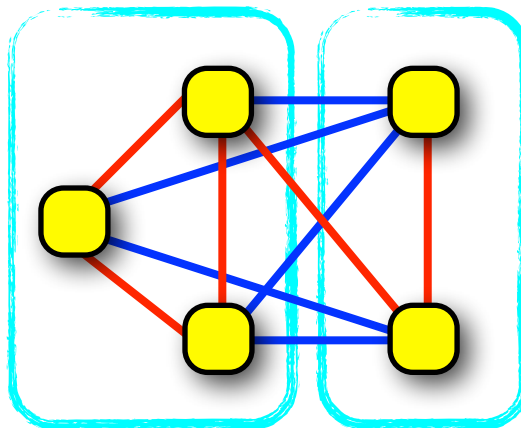
Revisiting Matching

Correlation Clustering

Coloring Graphs

Coverage and Submodular Maximization

- Consider a complete graph where edges are labelled **attractive** or **repulsive**. Given a node partition, an attractive edge is sad if it is cut and a repulsive edge is sad if it is not cut.



- Correlation Clustering Find partition minimizing # sad edges.
See tutorial *Bonchi, Garcia-Soriano, Liberty [KDD 14]*
- 3-Approx Algorithm a) Pick random node. b) Form cluster with it and its attracted neighbors. c) Remove cluster from graph and repeat until nodes remain. *Ailon, Charikar, Newman [J.ACM 08]*

- Emulating algorithm in two passes:
 - Preprocess Randomly order nodes, v_1, v_2, \dots, v_n .
 - First Pass Store all attractive edges incident to $\{v_1, \dots, v_{\sqrt{n}}\}$. Now can emulate first \sqrt{n} iterations of the algorithm.
 - Second Pass Store all remaining attractive edges. Now can emulate remaining steps of the algorithm.
- Thm Algorithm uses $\tilde{O}(n^{1.5})$ space. *Ahn et al. [ICML 16]*
- Proof Idea At most $n^{1.5}$ edges stored in first pass. In second, pass, can show remaining node have at most $n^{0.5}$ neighbors.
- With more work, can get $\tilde{O}(n)$ space with $O(\log \log n)$ passes. Can also find maximal independent sets.



Part III

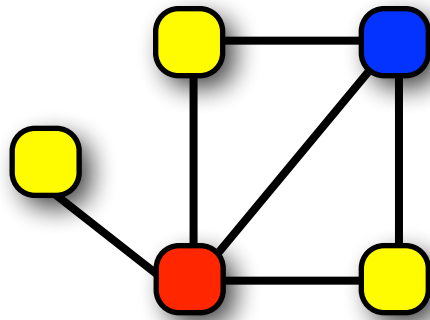
Streaming

Revisiting Matching
Correlation Clustering

Coloring Graphs

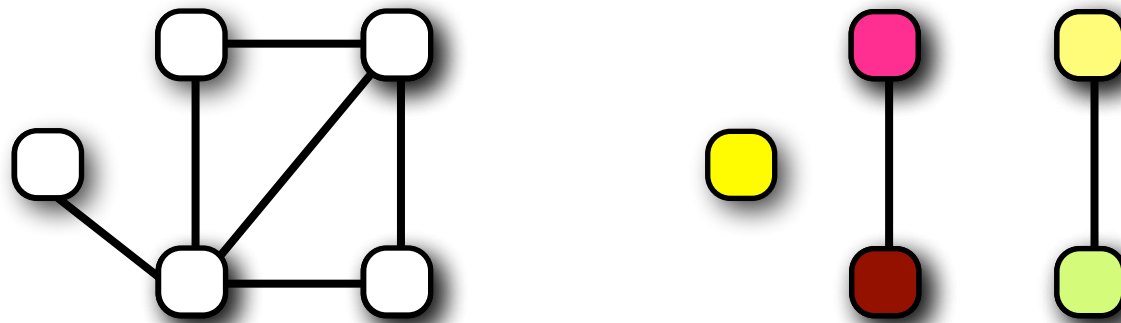
Coverage and Submodular Maximization

- Coloring With min number of colors, assign a color to every node such that no edge has monochromatic endpoints.



- Thm Can color a graph in $\Delta+1$ colors where Δ is max degree.
- ? How can we do this in a few passes with $\tilde{O}(n)$ space?
- $O(\Delta \log \log n)$ passes via independent sets. Let's do better!

- $(1+\epsilon)\Delta$ Coloring a) Randomly color with Δ/r colors. b) Store edges E' with monochromatic endpoints. c) Shade colors such that E' edges no longer monochromatic. Bera, Ghosh [ArXiv 18]



- Space Analysis $|E'| = O(nr)$ since probability edge in E' is r/Δ .
- Colors Analysis If $r \approx \epsilon^{-2} \log n$, max degree in E' is $\Delta_{E'} < (1+\epsilon)r$ and final number of colors is $(1+\Delta_{E'})\Delta/r = (1+\epsilon)\Delta$.
- $\Delta+1$ Coloring Idea For node v , pick $S_v \subset \{1, \dots, \Delta+1\}$ of $O(\log n)$ random colors. May assume v 's color in S_v . Assadi et al. [SODA 19]



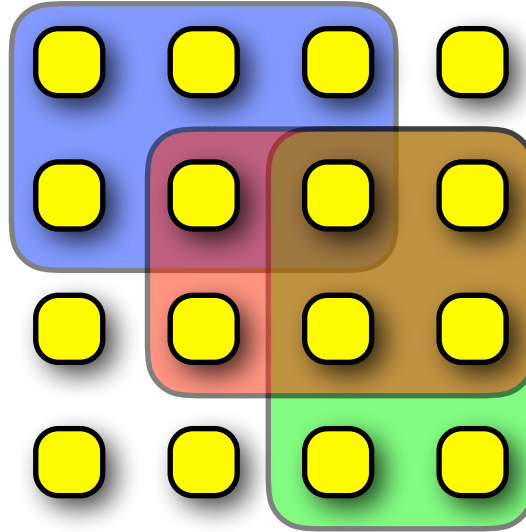
Part III

Streaming

Revisiting Matching
Correlation Clustering
Coloring Graphs

Coverage and Submodular Maximization

- Max-k-Coverage Given a stream of subsets S_1, \dots, S_m of $[n]$, find C that maximizes $f(C) = |\bigcup_{i \in C} S_i|$ subject to $|C| \leq k$.



- Submodular Functions f is sub-modular if for $A \subset B$ and $x \notin B$,

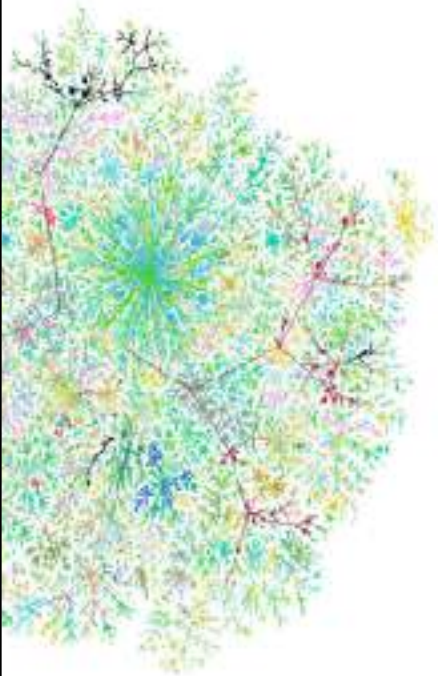
$$f(A \cup \{x\}) - f(A) \geq f(B \cup \{x\}) - f(B)$$

- Thm $(1-\epsilon)/2$ approx. of max-coverage in $\tilde{O}(\epsilon^{-3}k)$ space.

McGregor, Vu [ICDT 17]

- Algorithm Guess g such that $\text{OPT} \leq g \leq (1 + \epsilon)\text{OPT}$. Add first $\leq k$ sets that each cover at least $g/(2k)$ new elements.
- Approx Ratio
 - If k sets added, we cover $g/2 \geq \text{OPT}/2$.
 - If less sets added, each set not added covers $< g/(2k)$ new elements and hence we covered $\text{OPT} - g/2 \geq \text{OPT}(1 - \epsilon)/2$.
- Reducing Space Above algorithm requires $\tilde{O}(\epsilon^{-1} \text{OPT})$ space. Can use subsampling to such that $\text{OPT} = \tilde{O}(\epsilon^{-2} k)$.
- Generalizations Constant passes for $\approx 1 - 1/e$ approx. Extends to other monotone submodular function. Other work on non-monotone functions, beyond cardinality constraints, etc.

McGregor, Vu [ICDT 17], Bateni et al. [SPAA 17], Assadi [PODS 17]



Obrigado!