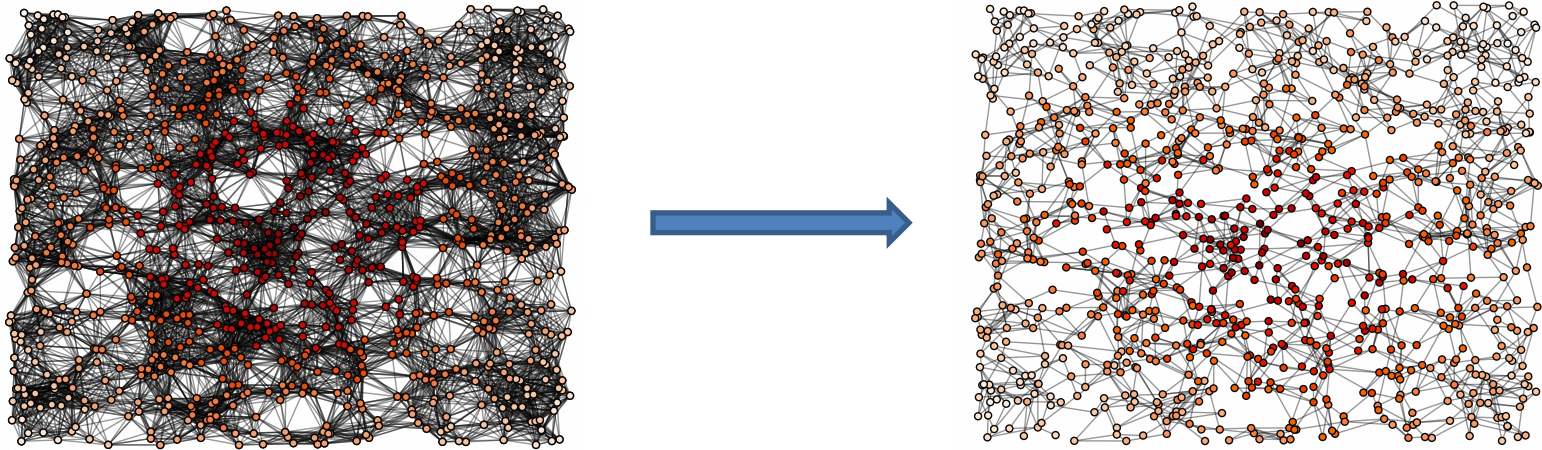


Graph Sparsification

Debmalya Panigrahi

Graph Sparsifiers



“Smaller” graph that (approximately) preserves the values of some set of graph parameters

Graph Sparsifiers

- Spanners
- Emulators
- Small stretch spanning trees
- Vertex sparsifiers
- ...
- Spectral sparsifiers
- Cut sparsifiers

Spectral Sparsification

- Undirected graph $\mathbf{G} = (\mathbf{V}, \mathbf{E})$; error parameter ε
- Goal: $\mathbf{G}_\varepsilon = (\mathbf{V}, \mathbf{E}_\varepsilon)$ with $\tilde{O}(n/\varepsilon^2)$ edges such that for all n -dimensional vectors \mathbf{x} ,

$$(1-\varepsilon) \mathbf{x}^\top \mathbf{L}(\mathbf{G}) \mathbf{x} \leq \mathbf{x}^\top \mathbf{L}(\mathbf{G}_\varepsilon) \mathbf{x} \leq (1+\varepsilon) \mathbf{x}^\top \mathbf{L}(\mathbf{G}) \mathbf{x}$$

- Graph Laplacian $\mathbf{L} = \mathbf{D} - \mathbf{A}$, where
 - \mathbf{D} = Diagonal Degree Matrix of the graph
 - \mathbf{A} = Adjacency Matrix of the graph

Spectral Sparsification: Previous work

Running time of the sparsification algorithm	Number of edges in the sparsifier	
$O(n^3m)$	$O(n/\epsilon^2)$	[Batson-Spielman-Srivastava '09]
$O(n^2m \log^3n + n^4 \log n)$		[Zouzias '12]
$O(m \log^{O(1)} n)$	$O(n \log^{O(1)} n/\epsilon^2)$	[Spielman-Teng '04]
$O(m \log^{O(1)} n)$	$O(n \log n/\epsilon^2)$	[Spielman-Srivastava '08]
$O(m \log^3 n)$		SS + [Koutis-Miller-Peng '10, '11]
$O(m \log^2 n)$		[Koutis-Levin-Peng '12]
$O(m \log n)$	$O(n \log^3 n/\epsilon^2)$	[Koutis-Levin-Peng '12]
$O(m) ???$	$???$	

Spectral to Cut Sparsifiers

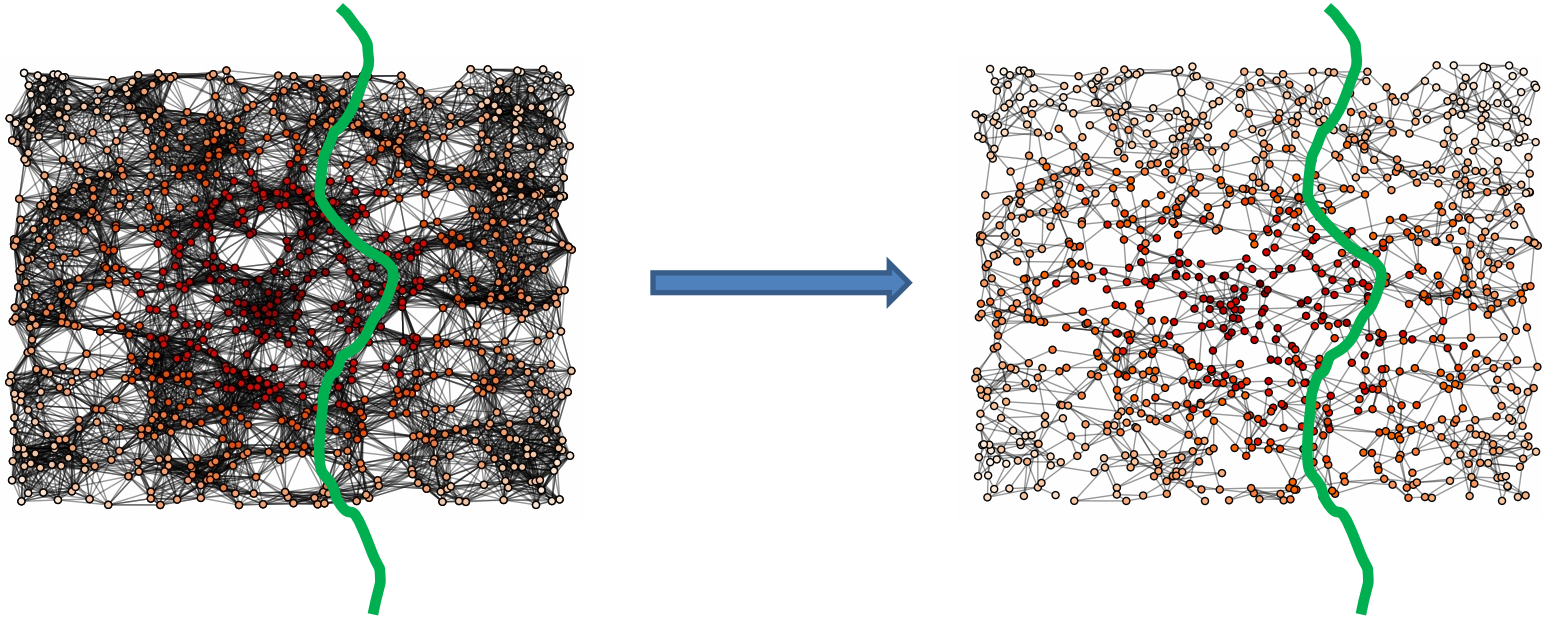
- $\mathbf{G}_\varepsilon = (\mathbf{V}, \mathbf{E}_\varepsilon)$ is a spectral sparsifier of $\mathbf{G} = (\mathbf{V}, \mathbf{E})$ if for all n -dimensional vectors \mathbf{x} ,

$$(1-\varepsilon) \mathbf{x}^\top \mathbf{L}(\mathbf{G}) \mathbf{x} \leq \mathbf{x}^\top \mathbf{L}(\mathbf{G}_\varepsilon) \mathbf{x} \leq (1+\varepsilon) \mathbf{x}^\top \mathbf{L}(\mathbf{G}) \mathbf{x}$$

- $\mathbf{x}^\top \mathbf{L} \mathbf{x} = \sum_{(i,j) \in \mathbf{E}} (\mathbf{x}_i - \mathbf{x}_j)^2$
- Suppose $\mathbf{x} \in \{0, 1\}^n$; $\mathbf{S} = \{i \in \mathbf{V} : \mathbf{x}_i = 1\}$. Then,

$$\mathbf{x}^\top \mathbf{L} \mathbf{x} = \sum_{(i,j) \in \mathbf{E}} (\mathbf{x}_i - \mathbf{x}_j)^2 = \sum_{(i,j) \in (\mathbf{S}, \mathbf{V} - \mathbf{S})} 1 = \mathbf{E}(\mathbf{S})$$

Cut Sparsification



Weight of every cut is preserved up to a multiplicative error of $(1 \pm \epsilon)$

Cut Sparsification

- Undirected (unweighted) graph $\mathbf{G} = (\mathbf{V}, \mathbf{E})$; error parameter ϵ
- Goal: $\mathbf{G}_\epsilon = (\mathbf{V}, \mathbf{E}_\epsilon)$ with $\mathbf{O}(n \log n / \epsilon^2)$ edges such that for all cuts $(\mathbf{S}, \mathbf{V} - \mathbf{S})$,
$$(1 - \epsilon) E(\mathbf{S}) \leq E_\epsilon(\mathbf{S}) \leq (1 + \epsilon) E(\mathbf{S})$$
- Introduced by Benczur-Karger '96
 - $\mathbf{O}(m \log^2 n)$ -time algorithm to find a cut sparsifier (with high probability) containing $\mathbf{O}(n \log n / \epsilon^2)$ edges in expectation

Fung-Hariharan-Harvey-P.:

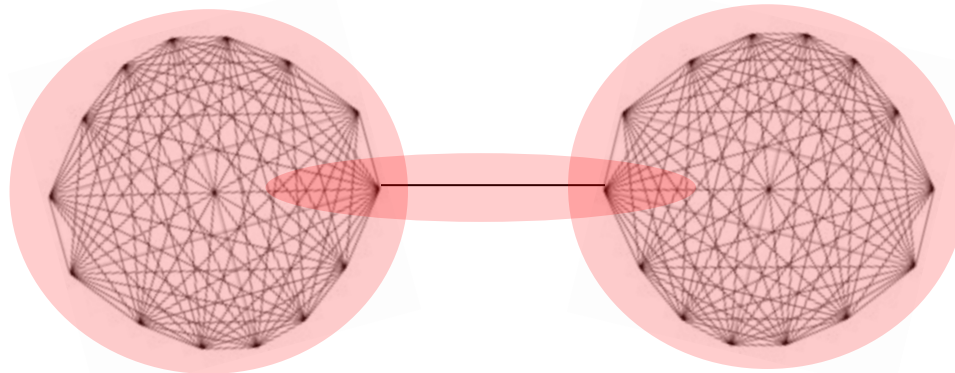
A linear-time, i.e. $O(m)$, algorithm that produces a cut sparsifier (whp) containing $O(n \log n / \epsilon^2)$ edges in expectation

Cut Sparsification by Sampling

Non

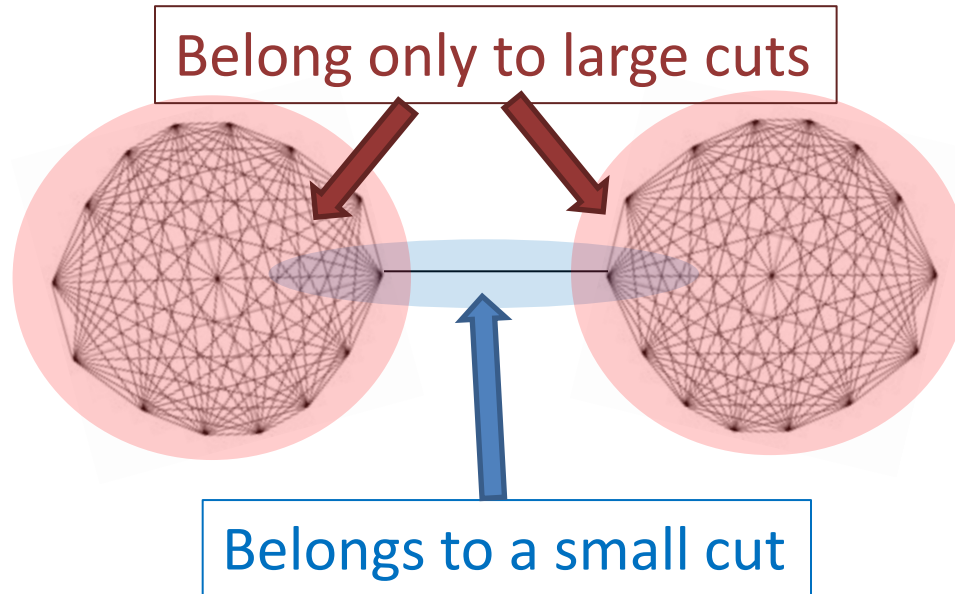
edge e with prob p_e

- **Uniformly** sample ~~all edges with prob $p \approx n/m$~~
 - Selected edge is given weight ~~$1/p$~~ $1/p_e$



$p \approx 1/n$; graph gets disconnected

Sampling Probabilities



Edge Connectivity λ_e = size of smallest cut containing e

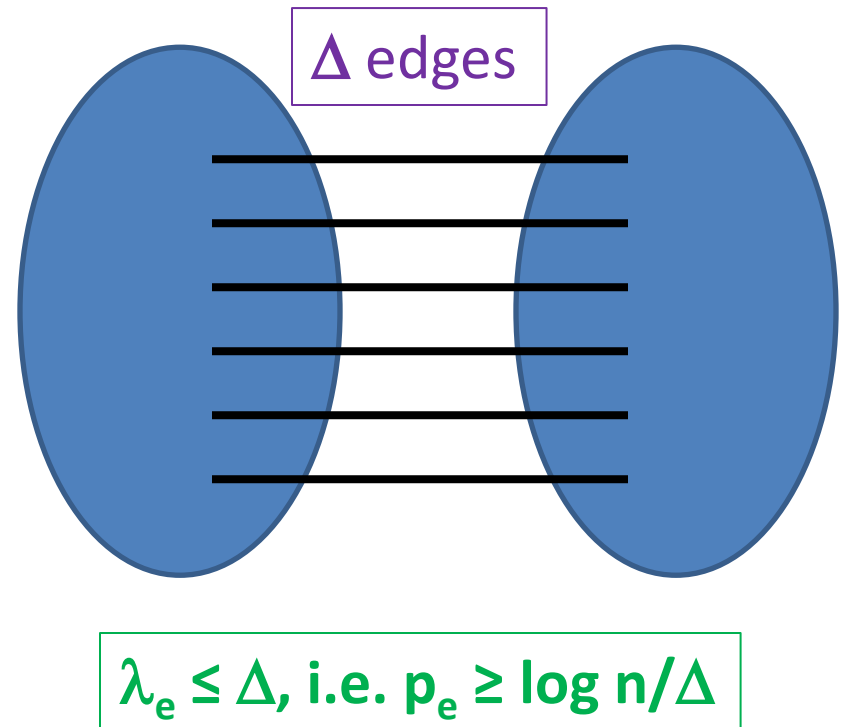
$$p_e = \log n / \lambda_e$$

Sampling by Edge Connectivity

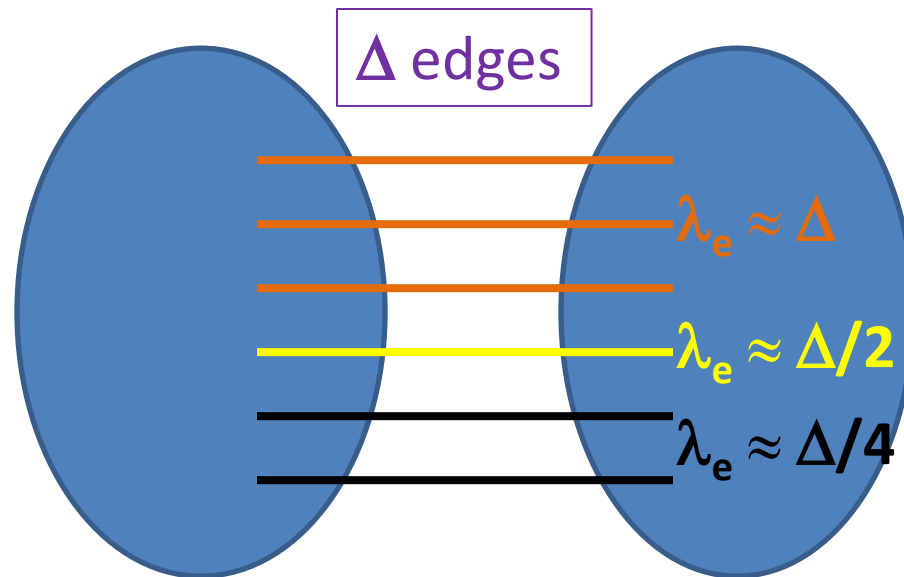
- Sample edge e independently (of other edges) with probability $p_e \approx \log n / \lambda_e$
- If edge e is selected, it is given a weight of $1/p_e$ in the sparsifier
- **Sparsifier has $O(n \log n)$ edges in expectation**
$$\lambda_e \geq 1/r_e \implies \sum_{e \in E} 1/\lambda_e \leq \sum_{e \in E} r_e = n - 1$$
- **$\Pr[E_\varepsilon(S) \in (1 \pm \varepsilon) E(S)$ for all cuts $(S, V - S)$]?**

Bounding Deviation

- Expected number of edges in the cut $\geq \log n$
- Chernoff bounds:
Probability of $\epsilon\Delta$ error $\leq 1/\text{poly}(n)$
- Exponential number of cuts!



Bounding Deviation



- For $\lambda_e \approx \Delta/k$ cut projection, $p_e = k \log n / \Delta$
- Probability of $\epsilon\Delta$ error $\leq \exp(-k \log n) = n^{-\Omega(k)}$

Cut Projections

Lemma: There are $\leq n^{O(k)}$ distinct (Δ, k) cut projections in cuts of size Δ

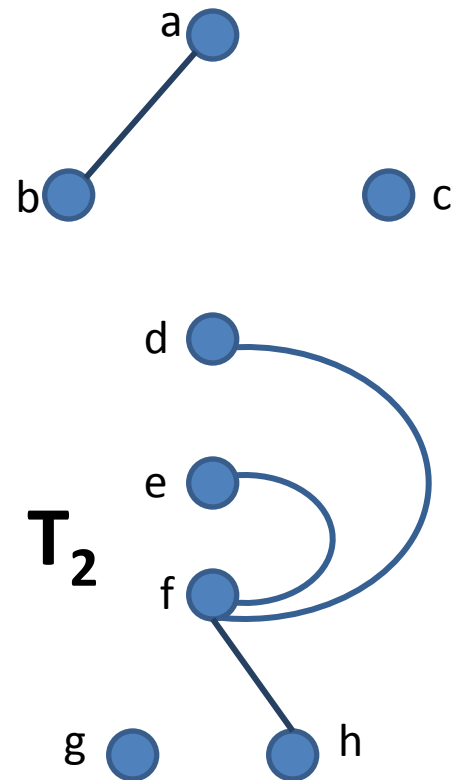
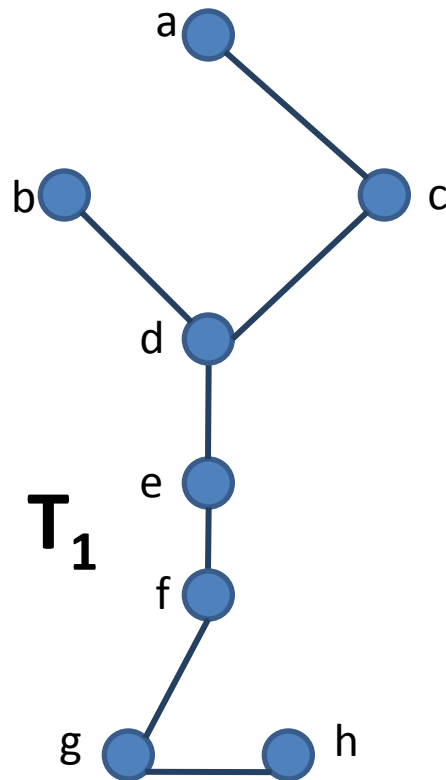
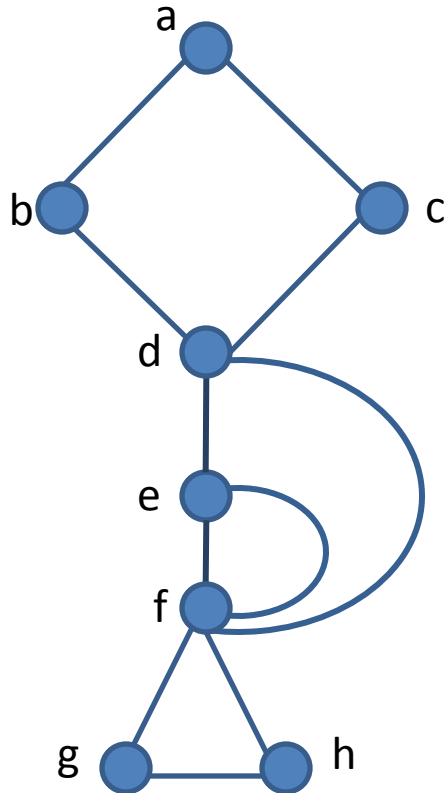
+ union bound on k, Δ



Theorem: Sampling edge e with probability $\log^2 n / \lambda_e$ produces a cut sparsifier

Difficulty: Edge connectivities (λ_e)
are time-consuming to calculate
(Gomory-Hu tree takes $\tilde{O}(mn)$ time
[Bhalgat-Hariharan-Kavitha-P., '07])

Greedy Spanning Forest packing



Sampling by NI Index

Nagamochi-Ibaraki (NI) index of edge e

y_e = index of e in an arbitrary but fixed greedy spanning forest packing

Proposed Cut Sparsification Algorithm

- Sample edge e with probability $p_e \approx \log^2 n / y_e$
- If edge e is selected, it is given a weight of $1/p_e$ in the sparsifier

Sampling by NI Index: Cut preservation

Lemma: The graph $G_\varepsilon = (V, E_\varepsilon)$ produced by sampling using NI indices is a cut sparsifier, i.e., with high probability, for all cuts $(S, V-S)$

$$(1-\varepsilon) E(S) \leq E_\varepsilon(S) \leq (1+\varepsilon) E(S)$$

- For each edge e , $y_e \leq \lambda_e$ (if edge e is in i^{th} forest, then its endpoints are connected by disjoint paths in the previous $i-1$ forests)
- Now piggyback on the proof for sampling using edge connectivities

Sampling by NI Index: Sparsification

Lemma: The sparsifier has $O(n \log^3 n)$ edges in expectation

$$\begin{aligned}\sum_{e \in E} \mathbf{1}/y_e &= \sum_k |T_k|/k = (n-1) \sum_k \mathbf{1}/k \\ &= O(n \log n)\end{aligned}$$

Sampling by NI Index: Running time

Lemma [Nagamochi-Ibaraki '92]: The running time of the sampling algorithm (i.e., time taken to estimate the NI indices of all edges) is $O(m)$

We have shown:

An $O(m)$ -time algorithm that produces a cut sparsifier containing $O(n \log^3 n)$ edges

We will now show:

An $O(m)$ -time algorithm that produces a cut sparsifier containing $O(n \log^2 n)$ edges

We had promised (see the paper):

An $O(m)$ -time algorithm that produces a cut sparsifier containing $O(n \log n)$ edges

Sampling by NI Index: New Algorithm

Previous Algorithm

- Sample edge e with probability $p_e \approx \log^2 n / y_e$
- If edge e is selected, it is given a weight of $1/p_e$ in the sparsifier

New Algorithm

- Sample edge e with probability $p_e \approx \log n / y_e$
- If edge e is selected, it is given a weight of $1/p_e$ in the sparsifier

Sampling by NI Index: New Algorithm

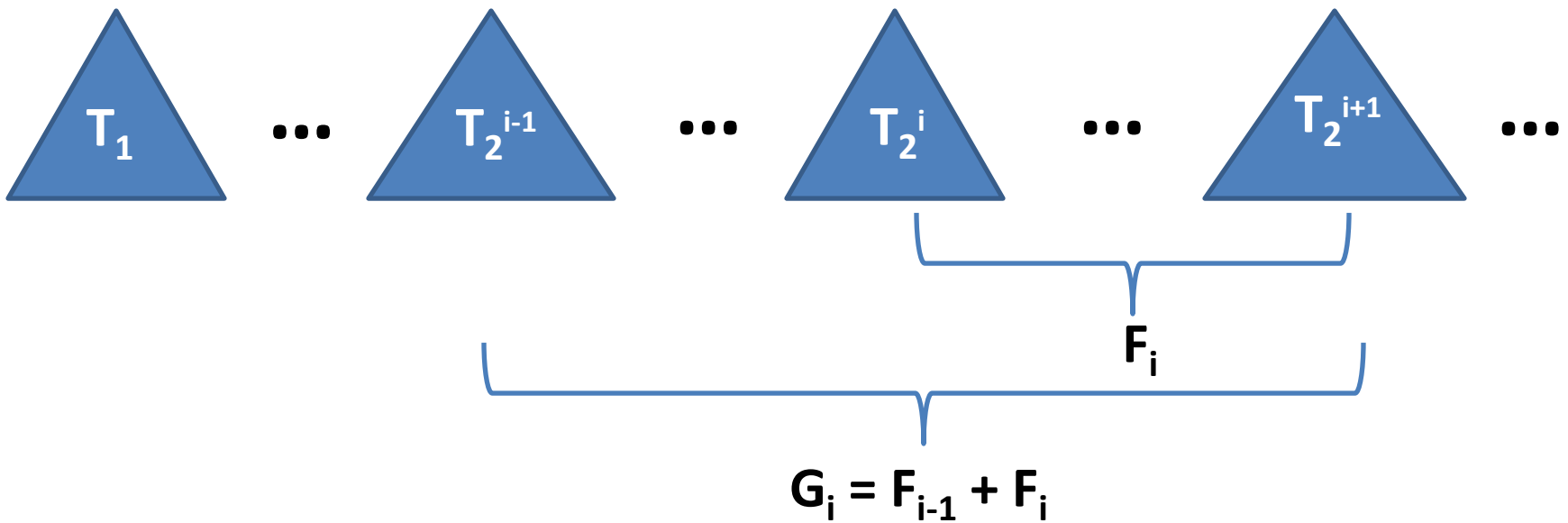
New Algorithm

- Sample edge e with probability $p_e \approx \log n / y_e$
- If edge e is selected, it is given a weight of $1/p_e$ in the sparsifier

- Running time remains $O(m)$
- The expected number of edges is $O(n \log^2 n)$
- **Is the sample a cut sparsifier?**

[Note: We can no longer piggyback on the analysis for sampling with edge connectivity]

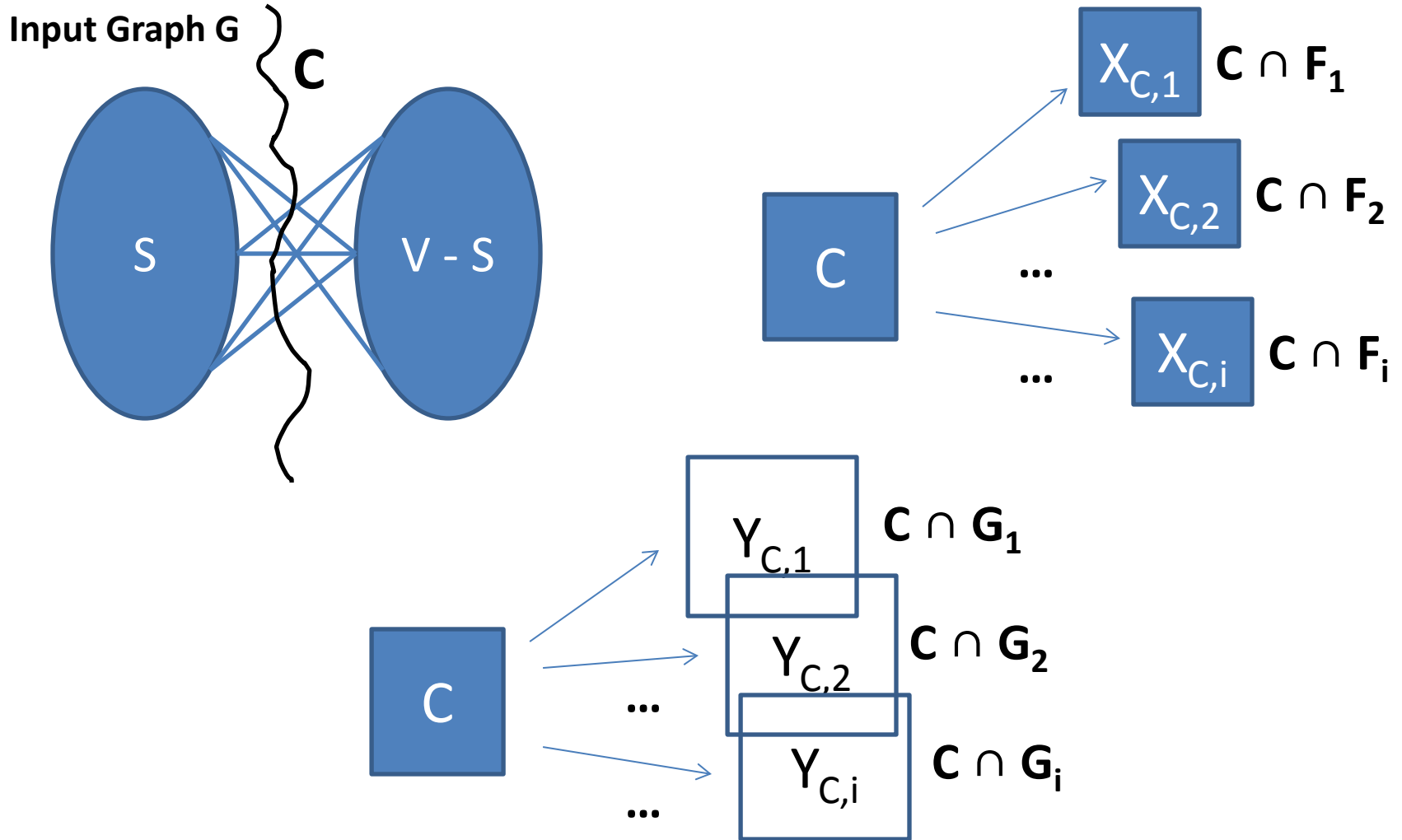
Bucketing the forests



Properties of the bucketing

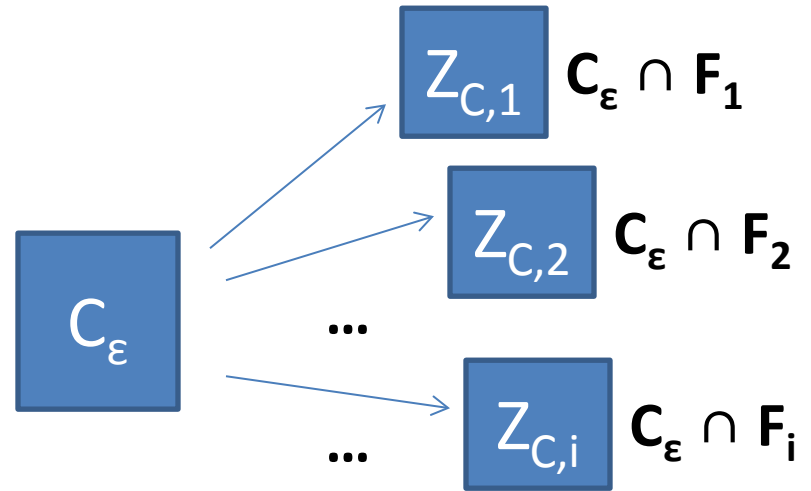
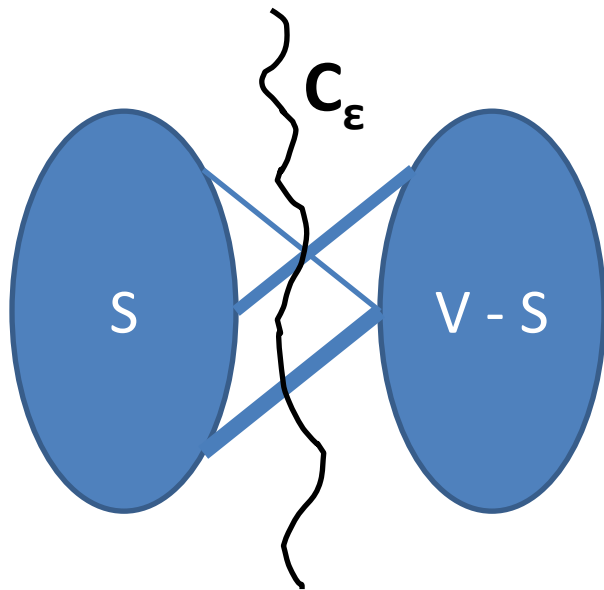
- Similarity property: All edges in F_i have sampling probability $p_e \approx \log n / 2^{i-1}$ (up to a factor of 2)
- Overlap property: Every edge appears in G_i for at most two values of i
- Connectivity property: Every edge in F_i has edge connectivity $\geq 2^{i-1}$ in G_i
 - The endpoints of the edge have 2^{i-1} disjoint paths between them, one in each forest, in G_i

Analysis of a cut



Tail Bounds on Deviation

Sampled graph G_ϵ



Tail Bounds on Deviation

- Need to show: whp, $|\mathbf{C} - \mathbf{C}_\varepsilon| < \varepsilon \mathbf{C}$ for all cuts \mathbf{C}



- whp, $|\mathbf{X}_{\mathbf{C},i} - \mathbf{Z}_{\mathbf{C},i}| < \varepsilon \mathbf{X}_{\mathbf{C},i}$ for all cuts \mathbf{C} and all i




$\sum_i \mathbf{Y}_{\mathbf{C},i} = 2\mathbf{C}$ by the
overlap property

Lemma: whp, $|\mathbf{X}_{\mathbf{C},i} - \mathbf{Z}_{\mathbf{C},i}| < \varepsilon \mathbf{Y}_{\mathbf{C},i}$
for all cuts \mathbf{C} and all i

Tail Bounds on Deviation

- Lemma: whp, $|X_{C,i} - Z_{C,i}| < \epsilon Y_{C,i}$ for all cuts C and all i
- Let C_k be cuts for which $Y_{C,i} = |C \cap G_i| = 2^{i+k}$
- By the connectivity property, every edge in $X_{C,i}$ is 2^{i-1} -connected in $Y_{C,i}$
- By Cut Projection Counting Lemma,
There are at most $n^{2^{i+k}/2^i} = n^{2^k}$ distinct $X_{C,i}$ in C_k

Tail Bounds on Deviation

- Lemma: whp, $|X_{C,i} - Z_{C,i}| < \epsilon Y_{C,i}$ for all cuts C and all i
- There are at most n^{2^k} distinct $X_{C,i}$ in C_k
- By the similarity property + Chernoff bounds,
 $\Pr[|X_{C,i} - Z_{C,i}| > \epsilon Y_{C,i}] < \exp(-2^{i+k} (\log n / 2^i)) = n^{-2^k}$

union bound over distinct $X_{C,i}$ in C_k , all values of k and i

Open Problems

- Linear-time spectral sparsification algorithm
- (Near)-linear time construction of $O(n/\epsilon^2)$ -sized cut/spectral sparsifiers
 - Edge sampling has fundamental limitations (connectivity of Erdos-Renyi random graph has a probability threshold of $\log n/n$)
 - Cut/spectral sparsifiers from spanning trees?
[Goyal-Rademacher-Vempala '09, Fung-Harvey '10]
 - Cut/spectral sparsifiers from spanners?
[Kapralov-Panigrahy '12, Koutis '14]

Thank You

Questions?