

Distributed Estimation of Graph 4-Profiles

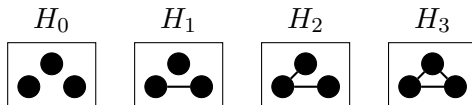
Ethan R. Elenberg, Karthikeyan Shanmugam,
Michael Borokhovich, Alexandros G. Dimakis

University of Texas, Austin, USA

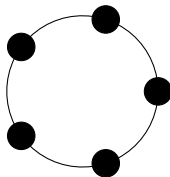
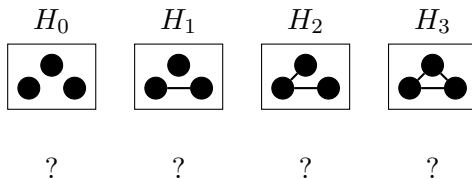
April 14, 2016

- Perform analytics on large graphs
 - How does information spread across the Web?
 - Is this social network user a spam bot?
 - Classify a protein as helpful or harmful
- Generalize existing subgraph analysis
 - Triangle counts, clustering coefficient, graphlet frequencies
- Scalable, distributed algorithms

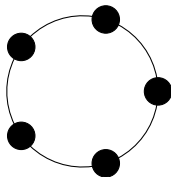
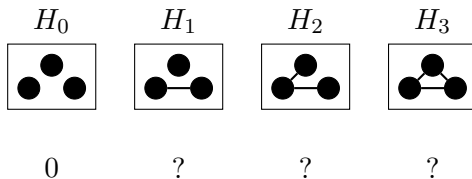
- Count the induced subgraphs formed by selecting all triples of vertices



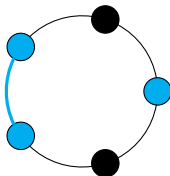
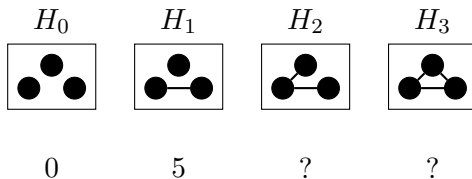
- Count the induced subgraphs formed by selecting all triples of vertices



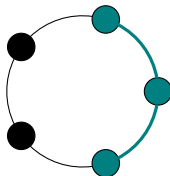
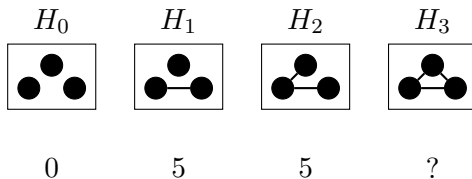
- Count the induced subgraphs formed by selecting all triples of vertices



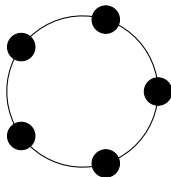
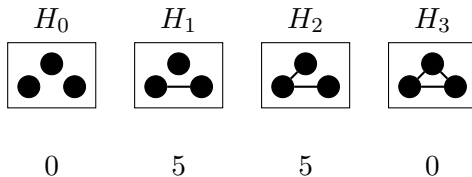
- Count the induced subgraphs formed by selecting all triples of vertices



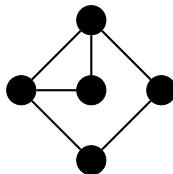
- Count the induced subgraphs formed by selecting all triples of vertices



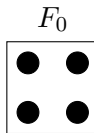
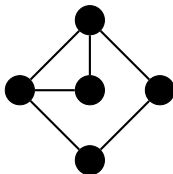
- Count the induced subgraphs formed by selecting all triples of vertices



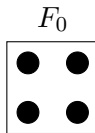
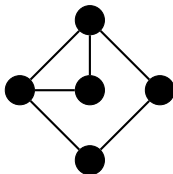
- Similarly, count the induced subgraphs formed by selecting all sets of 4 vertices
- $\mathbf{n}(G) = [?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?]$



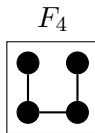
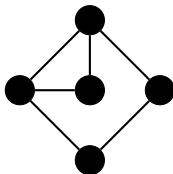
- Similarly, count the induced subgraphs formed by selecting all sets of 4 vertices
- $\mathbf{n}(G) = [?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?]$



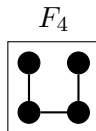
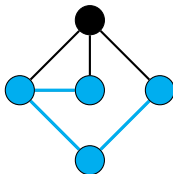
- Similarly, count the induced subgraphs formed by selecting all sets of 4 vertices
- $\mathbf{n}(G) = [0, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?]$



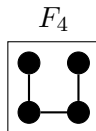
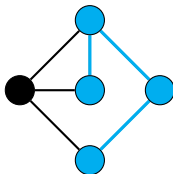
- Similarly, count the induced subgraphs formed by selecting all sets of 4 vertices
- $\mathbf{n}(G) = [0, 0, 0, 0, ?, ?, ?, ?, ?, ?, ?]$



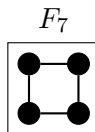
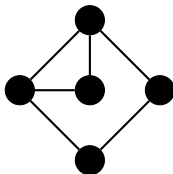
- Similarly, count the induced subgraphs formed by selecting all sets of 4 vertices
- $\mathbf{n}(G) = [0, 0, 0, 0, 1, ?, ?, ?, ?, ?, ?]$



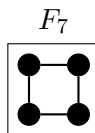
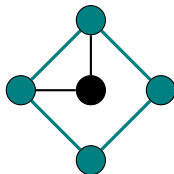
- Similarly, count the induced subgraphs formed by selecting all sets of 4 vertices
- $\mathbf{n}(G) = [0, 0, 0, 0, 2, ?, ?, ?, ?, ?, ?]$



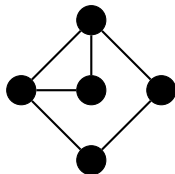
- Similarly, count the induced subgraphs formed by selecting all sets of 4 vertices
- $\mathbf{n}(G) = [0, 0, 0, 0, 2, 0, 0, ?, ?, ?, ?]$



- Similarly, count the induced subgraphs formed by selecting all sets of 4 vertices
- $\mathbf{n}(G) = [0, 0, 0, 0, 2, 0, 0, 1, ?, ?, ?]$



- Similarly, count the induced subgraphs formed by selecting all sets of 4 vertices
- $\mathbf{n}(G) = [0, 0, 0, 0, 2, 0, 0, 1, 2, 0, 0]$



Definition

Let N_i be the number of F_i 's in a graph G . The vector $\mathbf{n}(G) = [N_0, N_1, \dots, N_{10}]$ is called the **global 4-profile** of G .

- Always sums to $\binom{|V|}{4}$, the total number of 4-subgraphs

Definition

Let N_i be the number of F_i 's in a graph G . The vector $\mathbf{n}(G) = [N_0, N_1, \dots, N_{10}]$ is called the **global 4-profile** of G .

- Always sums to $\binom{|V|}{4}$, the total number of 4-subgraphs

Definition

For each $v \in V$, the **local 4-profile** counts how many times v participates in each F_i with 3 other vertices.

- Local 4-profiles embed each vertex into an 11 dimensional feature space
 - Spam detection
 - Generative models
- Global 4-profile concisely describes local connectivity
 - Molecule classification

- **Problem:** Compute (or approximate) 4-profile quantities for a large graph

- **Problem:** Compute (or approximate) 4-profile quantities for a large graph
- Previous approaches have 2 drawbacks
 - Require global communication
 - Many vertices redundantly repeat the same calculations

- ① Design novel, **distributed** algorithm to calculate local 4-profiles
- ② Derive improved concentration bounds for 4-profile sparsifiers
- ③ Evaluate performance on real-world datasets

Well studied across several communities (graphlets, motifs, subgraph frequencies):

- Graph sub-sampling
[Kim, Vu '00] [Tsourakakis, et al. '08 -'11] [Jha, et al. '15]
- Large-scale triangle counting
[Shank '07] [Satish, et al. '14] [Eden, et al. '15]
- Subgraph/graphlet counting equations
[Kloks, et al. '00] [Kowaluk, et al. '13]
ORCA [Hočevár, Demšar '14] [E. '15] [Ahmed, et al. '15]

- ① Introduction
- ② 4-PROF-DIST Algorithm
- ③ 4-profile Sparsifier
 - Edge Sub-sampling Process
 - Concentration Bound
- ④ Experiments
- ⑤ Conclusions

- Message passing algorithm in the Gather-Apply-Scatter framework
 - GraphLab, Pregel, Spark GraphX, etc.
 - Communication only allowed between adjacent vertices
 - Intermediate results stored as edge data

- 1 Each vertex computes its local 3-profile and triangle list

- ① Each vertex computes its local 3-profile and triangle list
- ② Each vertex solves a 17×17 system of equations relating its local 4-profile to its neighbors' local 3-profiles
 - Edge Pivots [Kloks, et al. '00] [Kowaluk, et al. '13] [E. '15]
 - 4-clique Counting
 - 2-hop Histogram

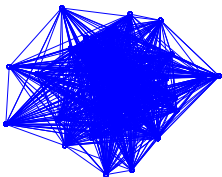
- ① Introduction
- ② 4-PROF-DIST Algorithm
- ③ 4-profile Sparsifier
 - Edge Sub-sampling Process
 - Concentration Bound
- ④ Experiments
- ⑤ Conclusions

Edge Sub-sampling Process

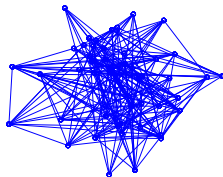
- Sub-sample each edge in the graph independently with probability p
- Relate the original and sub-sampled graphs via a 1-step Markov chain

Edge Sub-sampling Process: 4-clique

Original

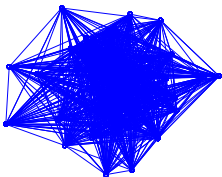


Sub-sampled



Edge Sub-sampling Process: 4-clique

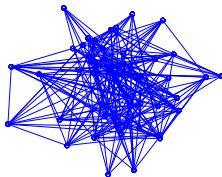
Original



p^6

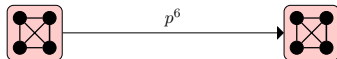


Sub-sampled

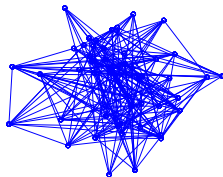


Edge Sub-sampling Process: 4-clique

Original



Sub-sampled

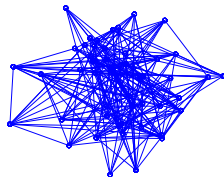
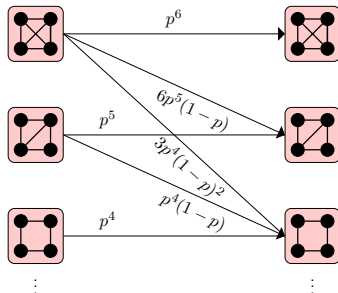
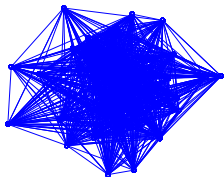


$$\left[\text{Estimator} \right] = \frac{1}{p^6} \left[\text{Sub-sampled} \right]$$

Edge Sub-sampling Process: 4-clique

Original

Sub-sampled



$$\left[\text{Estimator} \right] = \frac{1}{p^6} \left[\text{Sub-sampled} \right]$$

In general, construct an invertible transition matrix \mathbf{H}

$$H_{ij} = \mathbb{P}(\text{sub-sampled} = F_i \mid \text{original} = F_j)$$

$$\begin{bmatrix} \text{Unbiased} \\ \text{Estimators} \end{bmatrix} = \mathbf{H}^{-1} \begin{bmatrix} \text{Sub-sampled} \\ \text{4-profile} \end{bmatrix}$$

Main Concentration Result

- N_{10} - # 4-cliques
- k_{10} - Maximum # 4-cliques sharing a common edge

Theorem (*4-clique sparsifier*)

If the sampling probability

$$p \geq \left(\frac{\log(2/\delta)k_{10}}{2\epsilon^2 N_{10}} \right)^{1/12},$$

then the relative error is bounded by ϵ with probability at least $1 - \delta$.

Main Concentration Result

- N_{10} - # 4-cliques
- k_{10} - Maximum # 4-cliques sharing a common edge

Theorem (4-clique sparsifier)

If the sampling probability

$$p \geq \left(\frac{\log(2/\delta)k_{10}}{2\epsilon^2 N_{10}} \right)^{1/12},$$

then the relative error is bounded by ϵ with probability at least $1 - \delta$.

Proof Sketch:

- 4-clique estimator is associated with a read- k_{10} function family

$$f(G, p) = e_1 e_2 e_4 e_5 e_7 e_8 + e_4 e_5 e_6 e_8 e_9 e_{10} + \dots$$

- ① Introduction
- ② 4-PROF-DIST Algorithm
- ③ 4-profile Sparsifier
 - Edge Sub-sampling Process
 - Concentration Bound
- ④ Experiments
- ⑤ Conclusions

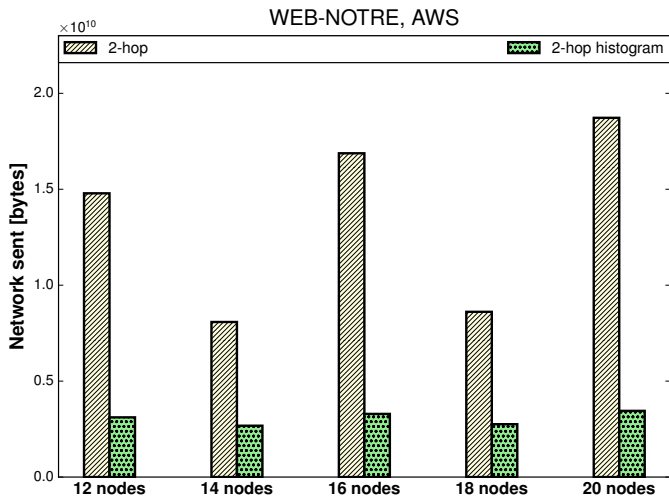
- GraphLab PowerGraph v2.2
- Multicore server
 - 256 GB RAM, 72 logical cores
- EC2 cluster (Amazon Web Services)
 - 20 c3.8xlarge, 60 GB RAM, 32 logical cores each

- GraphLab PowerGraph v2.2
- Multicore server
 - 256 GB RAM, 72 logical cores
- EC2 cluster (Amazon Web Services)
 - 20 c3.8xlarge, 60 GB RAM, 32 logical cores each

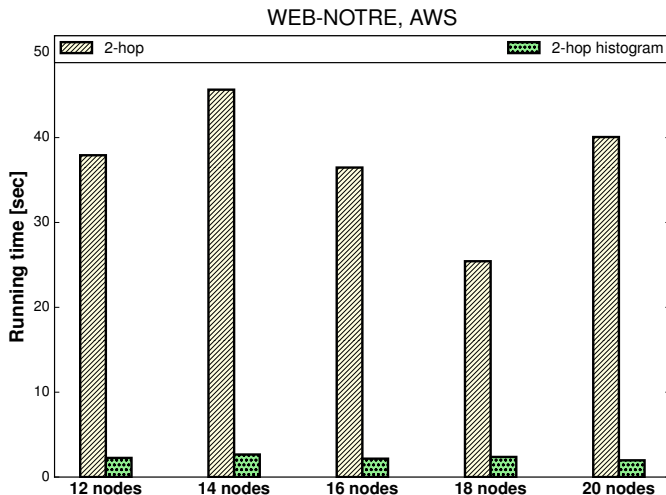
Datasets

Name	Vertices	Edges (undirected)
WEB-NOTRE	325,729	1,090,108
LiveJournal	4,846,609	42,851,237

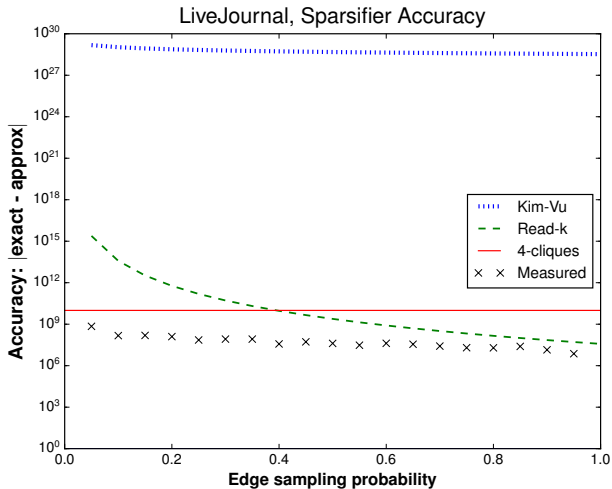
Results: AWS Full Neighborhood vs. Histogram, 10 runs



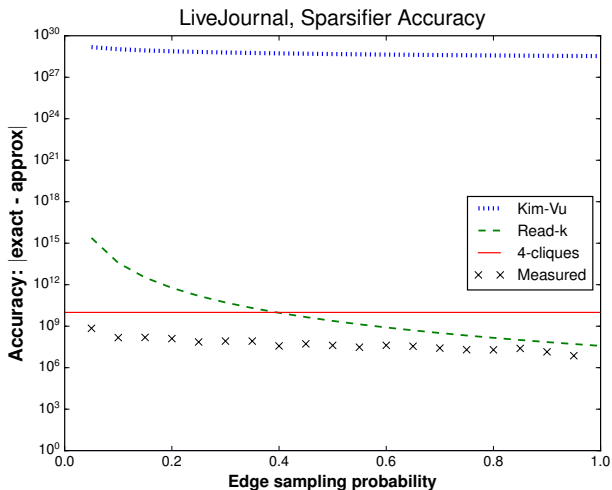
Results: AWS Full Neighborhood vs. Histogram, 10 runs



Results: Concentration Bounds

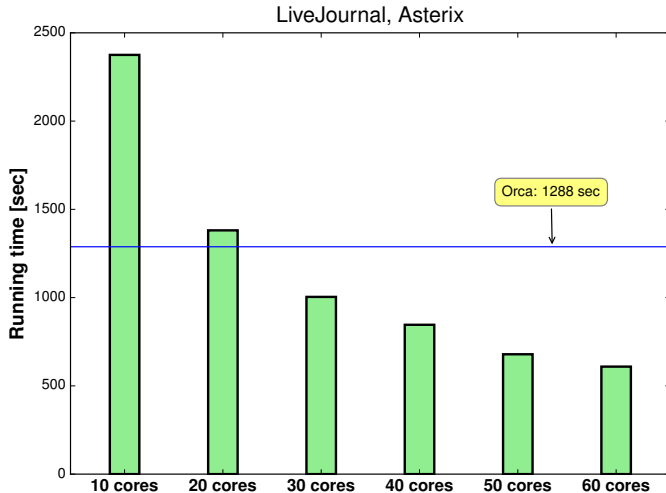


Results: Concentration Bounds



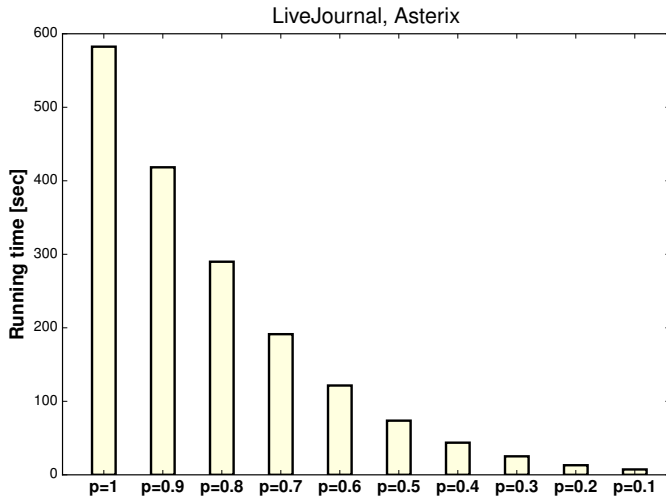
Improves on previous bounds if $p = \Omega(1/\log |E|)$

Results: Multicore Running Time Comparison, 10 runs

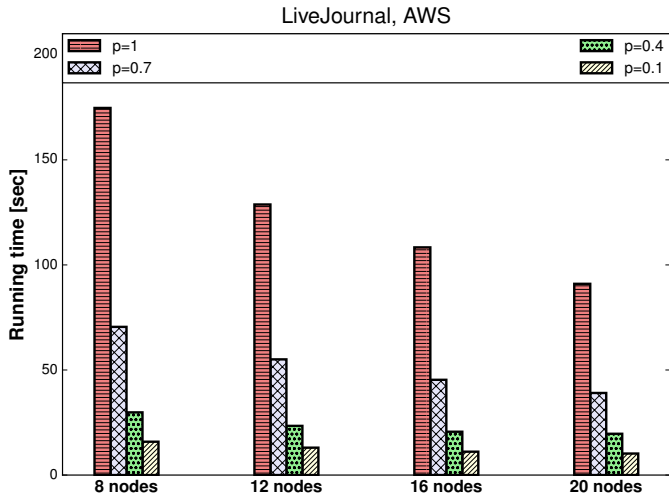


4-PROF-DIST, $p = 1$ vs. ORCA [Hočevár, Demšar '14]

Results: Multicore Running Time Comparison, 10 runs



Results: AWS Running Time Comparison, 10 runs

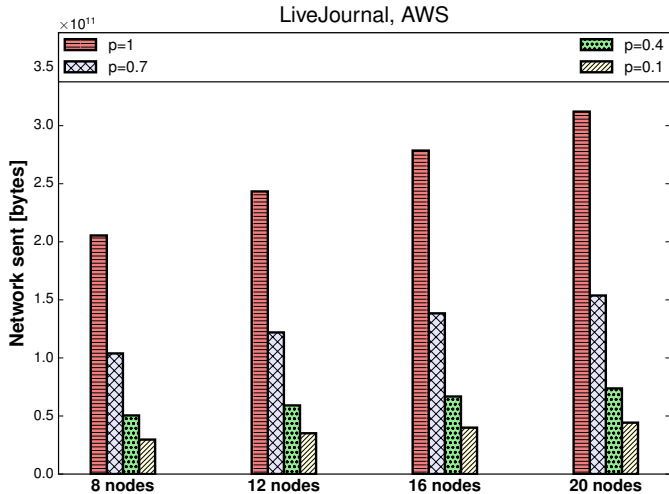


- ① Introduction
- ② 4-PROF-DIST Algorithm
- ③ 4-profile Sparsifier
 - Edge Sub-sampling Process
 - Concentration Bound
- ④ Experiments
- ⑤ Conclusions

- ① 2-hop histogram reduces network communication in a distributed setting
- ② Edge sub-sampling produces fast, accurate 4-profile estimates
 - Bounds for other subgraphs in the full paper
- ③ Distributed/parallel implementation improves performance at scale

`github.com/eelenberg/4-profiles`

(Backup) Results: AWS Network Communication, 10 runs



(Backup) 2-hop Histogram

- At each vertex a , store a $(p, c_a[p])$ pair for each neighbor p
- For any $a \in \Gamma(v)$ and $p \notin \Gamma(v)$,

$$c_a[p] = 1 \quad \Leftrightarrow \quad vap \text{ forms a 2-path}$$

- Gather across neighbors to count the number of distinct 2-paths from v to p

(Backup) 2-hop Histogram

- At each vertex a , store a $(p, c_a[p])$ pair for each neighbor p
- For any $a \in \Gamma(v)$ and $p \notin \Gamma(v)$,

$$c_a[p] = 1 \quad \Leftrightarrow \quad vap \text{ forms a 2-path}$$

- Gather across neighbors to count the number of distinct 2-paths from v to p

The diagram illustrates the decomposition of a 2-hop path from vertex v to vertex p . On the left, a central black vertex is connected to four other vertices: p (top, light orange), v (bottom, light yellow), and two unlabeled black vertices (left and right). This forms a diamond shape. An equals sign follows. To the right of the equals sign are two graphs separated by a plus sign. The first graph, labeled $F_7(v)$, shows a square with vertices v (bottom-left, light yellow), p (top-right, light orange), and two unlabeled black vertices (top-left and bottom-right). The second graph, labeled $F_9(v)$, shows a square with vertices v (bottom-left, light yellow), p (top-right, light orange), and two unlabeled black vertices (top-left and bottom-right), with an additional diagonal edge between the top-left and bottom-right vertices.

$$\sum_{p \notin \Gamma(v)} \left(\bigoplus_{a \in \Gamma(v)} c_a[p] \right)_2 = F_7(v) + F_9(v)$$

(Backup) 2-hop Histogram

- At each vertex a , store a $(p, c_a[p])$ pair for each neighbor p
- For any $a \in \Gamma(v)$ and $p \notin \Gamma(v)$,

$$c_a[p] = 1 \quad \Leftrightarrow \quad vap \text{ forms a 2-path}$$

- Gather across neighbors to count the number of distinct 2-paths from v to p

$$\sum_{p \notin \Gamma(v)} \left(\bigoplus_{a \in \Gamma(v)} c_a[p] \right)_2 = F_7(v) + F_9(v)$$

(Backup) 2-hop Histogram

- At each vertex a , store a $(p, c_a[p])$ pair for each neighbor p
- For any $a \in \Gamma(v)$ and $p \notin \Gamma(v)$,

$$c_a[p] = 1 \quad \Leftrightarrow \quad vap \text{ forms a 2-path}$$

- Gather across neighbors to count the number of distinct 2-paths from v to p

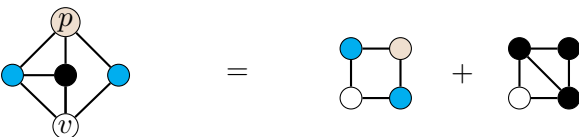
$$\sum_{p \notin \Gamma(v)} \left(\bigoplus_{a \in \Gamma(v)} c_a[p] \right)_2 = F_7(v) + F_9(v)$$

(Backup) 2-hop Histogram

- At each vertex a , store a $(p, c_a[p])$ pair for each neighbor p
- For any $a \in \Gamma(v)$ and $p \notin \Gamma(v)$,

$$c_a[p] = 1 \quad \Leftrightarrow \quad vap \text{ forms a 2-path}$$

- Gather across neighbors to count the number of distinct 2-paths from v to p


$$\sum_{p \notin \Gamma(v)} \left(\bigoplus_{a \in \Gamma(v)} c_a[p] \right)_2 = F_7(v) + F_9(v)$$

Vertex program in the Gather-Apply-Scatter framework [E. '15]

(Backup) Local 3-profile

Vertex program in the Gather-Apply-Scatter framework [E. '15]

- 1 For each vertex v : **Gather** and **Apply** vertex IDs to store $\Gamma(v)$

Vertex program in the Gather-Apply-Scatter framework [E. '15]

- 1 For each vertex v : **Gather** and **Apply** vertex IDs to store $\Gamma(v)$
- 2 For each edge va : **Scatter**

$$n_{3,va} = |\Gamma(v) \cap \Gamma(a)|,$$



$$n_{2,va}^c = |\Gamma(v)| - |\Gamma(v) \cap \Gamma(a)| - 1, \dots$$



Vertex program in the Gather-Apply-Scatter framework [E. '15]

- 1 For each vertex v : **Gather** and **Apply** vertex IDs to store $\Gamma(v)$
- 2 For each edge va : **Scatter**

$$n_{3,va} = |\Gamma(v) \cap \Gamma(a)|,$$



$$n_{2,va}^c = |\Gamma(v)| - |\Gamma(v) \cap \Gamma(a)| - 1, \dots$$



- 3 For each vertex v : **Gather** and **Apply**

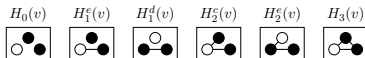
$$n_{3,v} = \frac{1}{2} \sum_{a \in \Gamma(v)} n_{3,va}$$



$$n_{2,v}^c = \frac{1}{2} \sum_{a \in \Gamma(v)} n_{2,va}^c, \dots$$



(Backup) Edge Pivot Equations



$$\sum_{a \in \Gamma(v)} \binom{n_{1,va}^e}{2} = F_1(v) + F_2(v)$$

$$\sum_{a \in \Gamma(v)} n_{1,va}^e n_{3,va} = 2F_5(v) + F_8''(v)$$

$$\sum_{a \in \Gamma(v)} \binom{n_{2,va}^c}{2} = 3F_6'(v) + F_8'(v)$$

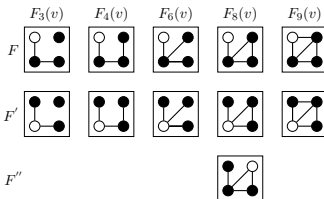
$$\sum_{a \in \Gamma(v)} n_{2,va}^c n_{2,va}^e = F_4'(v) + 2F_7(v)$$

$$\sum_{a \in \Gamma(v)} \binom{n_{3,va}}{2} = F_9'(v) + 3F_{10}(v)$$

$$\sum_{a \in \Gamma(v)} n_{2,va}^c n_{3,va} = 2F_8'(v) + 2F_9'(v)$$

$$\sum_{a \in \Gamma(v)} n_{1,va}^e n_{2,va}^c = 2F_3'(v) + F_4'(v)$$

$$n_{1,v}^d |\Gamma(v)| = F_2(v) + F_4(v) + F_8(v)$$

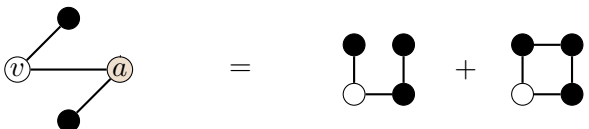


(Backup) Edge Pivot Equations

- Relate local 4-profile at v to neighboring local 3-profiles
- Accumulate pairs of subgraph counts involving edge va , summed over all neighbors

(Backup) Edge Pivot Equations

- Relate local 4-profile at v to neighboring local 3-profiles
- Accumulate pairs of subgraph counts involving edge va , summed over all neighbors

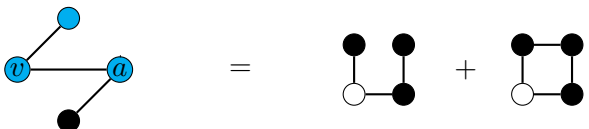


The diagrammatic equation shows a local 4-profile at node v (a white node connected to three black nodes, one of which is a) equal to the sum of two subgraph counts. The first subgraph is a path of length 2 (white node v connected to two black nodes, one of which is a), and the second is a square (white node v connected to two black nodes, one of which is a , and the other two black nodes are connected to each other).

$$\sum_{a \in \Gamma(v)} n_{2,va}^c n_{2,va}^e = F'_4(v) + 2F_7(v)$$

(Backup) Edge Pivot Equations

- Relate local 4-profile at v to neighboring local 3-profiles
- Accumulate pairs of subgraph counts involving edge va , summed over all neighbors

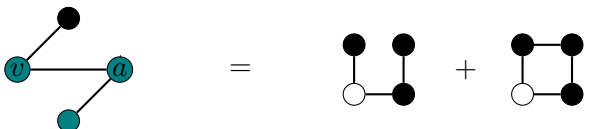


The diagrammatic equation shows a local 4-profile at node v (a blue node connected to three other nodes: one blue, one blue labeled a , and one black) is equal to the sum of two subgraph counts. The first subgraph is a path of length 2 (white node - black node - black node), labeled $F'_4(v)$. The second subgraph is a square (white node - black node - black node - black node), labeled $2F_7(v)$.

$$\sum_{a \in \Gamma(v)} n_{2,va}^c n_{2,va}^e = F'_4(v) + 2F_7(v)$$

(Backup) Edge Pivot Equations

- Relate local 4-profile at v to neighboring local 3-profiles
- Accumulate pairs of subgraph counts involving edge va , summed over all neighbors

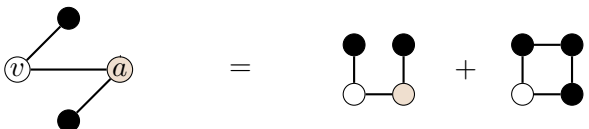


The diagrammatic equation shows a local 4-profile at node v (a teal node connected to a black node, a teal node a , and another teal node) is equal to the sum of two subgraph counts. The first subgraph is a path of length 2 (teal node - black node - black node), labeled $F'_4(v)$. The second subgraph is a square (teal node - black node - black node - teal node), labeled $2F_7(v)$.

$$\sum_{a \in \Gamma(v)} n_{2,va}^c n_{2,va}^e = F'_4(v) + 2F_7(v)$$

(Backup) Edge Pivot Equations

- Relate local 4-profile at v to neighboring local 3-profiles
- Accumulate pairs of subgraph counts involving edge va , summed over all neighbors

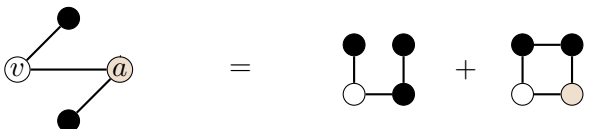


The diagrammatic equation shows a central node v (white) connected to three black nodes. One of these edges is highlighted as the pivot edge va , where a is a pink node. This is equated to the sum of two subgraph counts: $F'_4(v)$, represented by a white node v connected to a pink node a , which is in turn connected to two black nodes; and $2F_7(v)$, represented by a white node v connected to a pink node a , which is part of a square formed by four black nodes.

$$\sum_{a \in \Gamma(v)} n_{2,va}^c n_{2,va}^e = F'_4(v) + 2F_7(v)$$

(Backup) Edge Pivot Equations

- Relate local 4-profile at v to neighboring local 3-profiles
- Accumulate pairs of subgraph counts involving edge va , summed over all neighbors

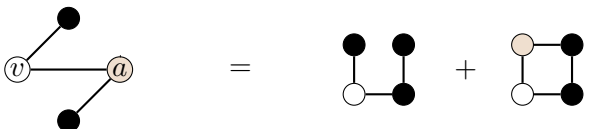


The diagrammatic equation shows a local 4-profile at node v (a white node connected to three black nodes, one of which is node a) equal to the sum of two subgraph counts. The first subgraph is a path of three nodes (white, black, black) with the white node at v . The second subgraph is a path of three nodes (black, black, white) with the white node at v . The white node in the second subgraph is shaded orange.

$$\sum_{a \in \Gamma(v)} n_{2,va}^c n_{2,va}^e = F'_4(v) + 2F_7(v)$$

(Backup) Edge Pivot Equations

- Relate local 4-profile at v to neighboring local 3-profiles
- Accumulate pairs of subgraph counts involving edge va , summed over all neighbors



The diagrammatic equation shows a central node v (white) connected to three black nodes. One of these edges is highlighted in orange, connecting v to a node a (orange). This is equated to the sum of two subgraph counts: $F'_4(v)$ and $2F_7(v)$. $F'_4(v)$ is represented by a square with nodes v (white), a (black), and two other black nodes. $2F_7(v)$ is represented by a square with nodes v (white), a (orange), and two other black nodes.

$$\sum_{a \in \Gamma(v)} n_{2,va}^c n_{2,va}^e = F'_4(v) + 2F_7(v)$$

(Backup) Clique Counting

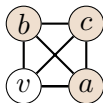
- Directly count the number of 4-cliques, F_{10}
- ① Initially, each vertex stores its triangle list $\Delta(v)$

(Backup) Clique Counting

- Directly count the number of 4-cliques, F_{10}
- ① Initially, each vertex stores its triangle list $\Delta(v)$
- ② For each edge v :
 Gather triangles of a that contain 3 neighbors of v

(Backup) Clique Counting

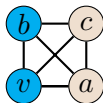
- Directly count the number of 4-cliques, F_{10}
- ① Initially, each vertex stores its triangle list $\Delta(v)$
- ② For each edge v :
 Gather triangles of a that contain 3 neighbors of v



$$\sum_{a \in \Gamma(v)} |(b, c) \in \Delta(a) : b \in \Gamma(v), c \in \Gamma(v)|$$

(Backup) Clique Counting

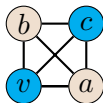
- Directly count the number of 4-cliques, F_{10}
- ① Initially, each vertex stores its triangle list $\Delta(v)$
- ② For each edge v :
 Gather triangles of a that contain 3 neighbors of v



$$\sum_{a \in \Gamma(v)} |(b, c) \in \Delta(a) : b \in \Gamma(v), c \in \Gamma(v)|$$

(Backup) Clique Counting

- Directly count the number of 4-cliques, F_{10}
- ① Initially, each vertex stores its triangle list $\Delta(v)$
- ② For each edge v :
 Gather triangles of a that contain 3 neighbors of v



$$\sum_{a \in \Gamma(v)} |(b, c) \in \Delta(a) : b \in \Gamma(v), c \in \Gamma(v)|$$

- Because there are no large constants, our concentration bounds improve on earlier work in practical settings

Corollary (*Comparison with [Kim, Vu]*)

Let G be a graph with m edges. If $p = \Omega(1/\log m)$ and $\delta = \Omega(1/m)$, then read- k provides better triangle sparsifier accuracy than [Kim, Vu]. If additionally $k_{10} \leq N_{10}^{5/6}$, then read- k provides better 4-clique sparsifier accuracy than [Kim, Vu].

(Backup) Results: 4-profile Sparsifier Accuracy, 10 runs

