# Beyond Triangles: A Distributed Framework for Estimating 3-profiles of Large Graphs
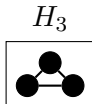
**Ethan R. Elenberg**, Karthikeyan Shanmugam,
Michael Borokhovich, Alexandros G. Dimakis

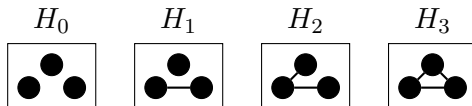University of Texas, Austin, USA

August 12, 2015

- Perform analytics on large graphs
  - World Wide Web, social networks, bioinformatics

- More descriptive than triangle count, clustering coefficient
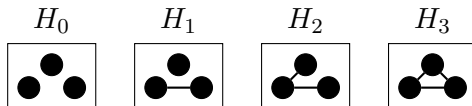
- Scalable, distributed algorithms

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

$H_3$

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



$H_0$     $H_1$     $H_2$     $H_3$

## 3-profile

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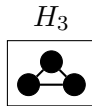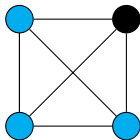


$H_0$     $H_1$     $H_2$     $H_3$

### Definition
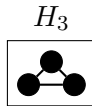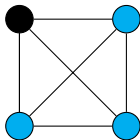
Let $n_i$ be the number of $H_i$'s in a graph $G$. The vector $\mathbf{n}(G) = [n_0, n_1, n_2, n_3]$ is called the 3-profile of $G$.

- Always sums to $\binom{|V|}{3}$, the total number of 3-subgraphs

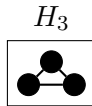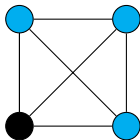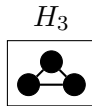- 4-clique: $\mathbf{n}(K_4) = [0, 0, 0, 4]$



$H_3$

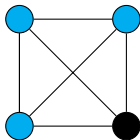- 4-clique: $\mathbf{n}(K_4) = [0, 0, 0, 4]$



$H_3$

- 4-clique: $\mathbf{n}(K_4) = [0, 0, 0, 4]$
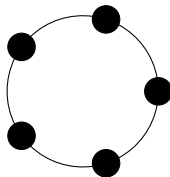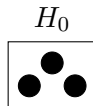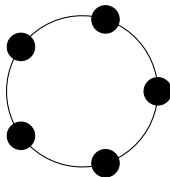


$H_3$

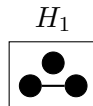- 4-clique: $\mathbf{n}(K_4) = [0, 0, 0, 4]$



$H_3$

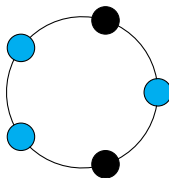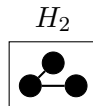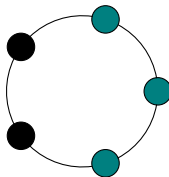- 5-cycle: $\mathbf{n}(C_5) = [?, ?, ?, ?]$
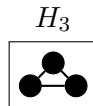
- 5-cycle: $\mathbf{n}(C_5) = [0, ?, ?, ?]$
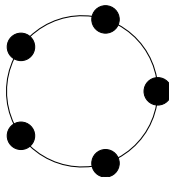


$H_0$

- 5-cycle: $\mathbf{n}(C_5) = [0, 5, ?, ?]$



$H_1$

- 5-cycle: $\mathbf{n}(C_5) = [0, 5, 5, ?]$



$H_2$

- 5-cycle: $\mathbf{n}(C_5) = [0, 5, 5, 0]$



$H_3$

For each $v \in V$:

### Definition

The local 3-profile counts how many times $v$ participates in each $H_i$ with 2 other vertices.

# Related Terms

For each $v \in V$:

## Definition

The local 3-profile counts how many times $v$ participates in each $H_i$ with 2 other vertices.

## Definition

The ego 3-profile is the 3-profile of ego graph $N(v)$.

- Graph induced by set of neighbors $\Gamma(v)$

- Global 3-profile concisely describes local connectivity
  - Molecule classification

- Local and ego 3-profiles are feature vectors for each vertex
  - Spam detection
  - Generative models

- Problem: Compute (or approximate) 3-profile quantities for a large graph

- Problem: Compute (or approximate) 3-profile quantities for a large graph

- Approach: Edge sub-sampling and distributed implementation

1. Derive a 3-profile *sparsifier* with provable guarantees

2. Design distributed, graph engine algorithms to calculate local and ego 3-profiles

3. Evaluate performance on real-world datasets

# Related Work

Well studied across several communities:

- Graph sub-sampling
  [Kim, Vu '00] [Tsourakakis, et al. '08 -'11] [Ahmed, et al. '14]
- Large-scale triangle counting
  [Satish, et al. '14] [Shank '07] [Suri, Vassilvitskii '11]
- Subgraph counting
  [Alon, et al. '97] [Kloks, et al. '00] [Kowaluk, et al. '13]
- Graphlets
  [Pržulj '07] [Shervashidze, et al. '09]

# Outline

- Sub-sample each edge in the graph independently with probability $p$

- Relate the original and sub-sampled graphs via a 1-step Markov chain
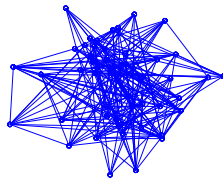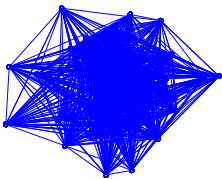
Original

Sub-sampled

$p^3$

Original

Sub-sampled

$p^2$

$p^3$

Original

Sub-sampled

$p^2$

$3p^2(1-p)$

$p^3$

Original
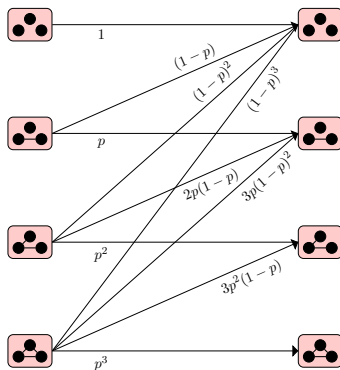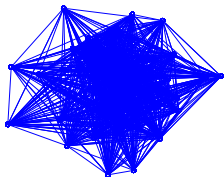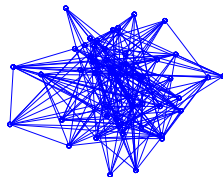
Sub-sampled

Original

Sub-sampled

$$\begin{bmatrix} \text{Estimator} \end{bmatrix} = \begin{bmatrix} 1 & 1-p & (1-p)^2 & (1-p)^3 \\ 0 & p & 2p(1-p) & 3p(1-p)^2 \\ 0 & 0 & p^2 & 3p^2(1-p) \\ 0 & 0 & 0 & p^3 \end{bmatrix}^{-1} \begin{bmatrix} \text{Sub-sampled} \end{bmatrix}$$

# Main Result

### Theorem (*3-profile sparsifiers*)

*For all ($\epsilon$,$p$)-balanced graphs*[*], *the $l_\infty$-norm of the 3-profile sparsifier error is bounded by $\epsilon\binom{|V|}{3}$ with high probability.*

## Main Result

### Theorem (*3-profile sparsifiers*)

*For all $(\epsilon,p)$-balanced graphs*[*], *the $l_\infty$-norm of the 3-profile sparsifier error is bounded by $\epsilon\binom{|V|}{3}$ with high probability.*

### Definition

A graph is $(\epsilon,p)$-balanced if the majority of "triangles," "wedges," or "single-edges" do not depend on one common edge.

## Main Result

### Theorem (*3-profile sparsifiers*)

*For all ($\epsilon$,p)-balanced graphs*, the $l_\infty$-norm of the 3-profile
sparsifier error is bounded by $\epsilon \binom{|V|}{3}$ with high probability.*

### Definition

A graph is ($\epsilon$,p)-balanced if the majority of "triangles," "wedges,"
or "single-edges" do not depend on one common edge.

*Proof Sketch:*

- Apply multivariate polynomial concentration inequalities [Kim,
  Vu '00] to each estimator

$$f(G, p) = e_1 e_2 e_4 + e_4 e_5 e_6 + \ldots$$

# Outline

Vertex program in the Gather-Apply-Scatter framework

Vertex program in the Gather-Apply-Scatter framework

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

Vertex program in the Gather-Apply-Scatter framework

**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

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

**2** For each edge $va$: Scatter

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

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

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

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

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

# Outline

## Implementation

- 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

# Implementation

- 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) |
|------|----------|--------------------|
| Twitter | $41,652,230$ | $1,202,513,046$ |
| PLD | $39,497,204$ | $582,567,291$ |
| LiveJournal | $4,846,609$ | $42,851,237$ |
| Wikipedia | $3,515,067$ | $42,375,912$ |
| DBLP | $317,080$ | $1,049,866$ |

Compare *3-PROF* to GraphLab's default triangle count



Twitter and PLD, Multicore (p=1)

Compare *EGO-PAR* to naive, serial algorithm (*EGO-SER*)



LiveJournal, AWS c3_8xlarge

LiveJournal, AWS c3_8xlarge

Beyond Triangles

# Outline

## Summary

1. Edge sub-sampling produces fast, accurate 3-profile estimates

2. 3-profile counting consumes roughly the same resources as triangle counting

3. Distributed algorithms scale well over large data and large computing clusters

```
github.com/eelenberg/3-profiles
```

$$\sum_{a \in \Gamma(v)} \binom{n_{3,va}}{2} = F_2(v) + 3F_3(v)$$

$$\sum_{a\in\Gamma(v)}\binom{n_{3,va}}{2} = F_2(v) + 3F_3(v)$$

$$\sum_{a \in \Gamma(v)} \binom{n_{3,va}}{2} = F_2(v) + 3F_3(v)$$

Twitter, Accuracy, 3-profiles

LiveJournal Running Time



PLD Running Time

LiveJournal Network Usage



PLD Network Usage

*EGO-SER*



*EGO-PAR*