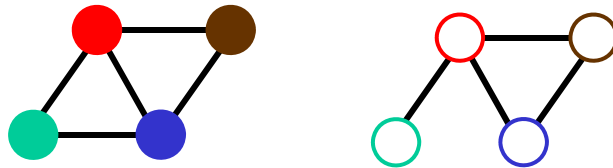


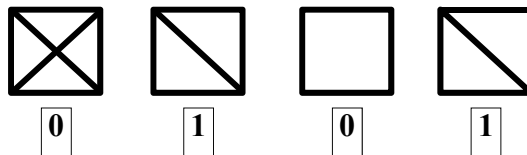
Our Graphs Become Larger

- Simple algorithms do not scale
 - $O(n^k)$ for size k graphlets
- Two approaches:
 - Find clever algorithms for counting small graphlets
 - Approximate count for larger graphlets

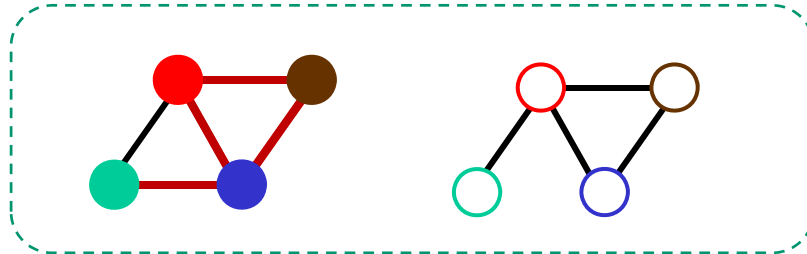
Graph Induced Subgraphs



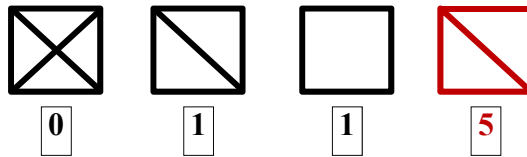
Potential Subgraphs



Graph NON-Induced Subgraphs

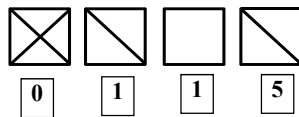
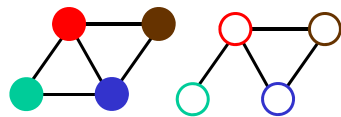


Potential Subgraphs

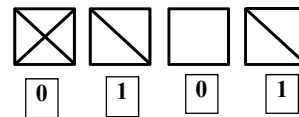
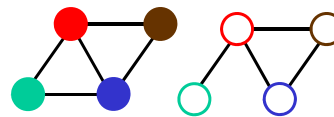


Subgraphs

Non-Induced



Induced



Efficient Counting of Graphlets

Key Observations:

- It is easier to count non-induced graphlets
- Counting all non-induced graphlets can give us all induced graphlets of some size.

Tree Decomposition

A **tree decomposition** of a graph $G = (V, E)$ is a pair (X, T) , where $X = \{X_1, \dots, X_n\}$ is a family of subsets of V , and T is a tree whose nodes are the subsets X_i , satisfying the following properties:

- Each graph vertex is associated with at least one tree node.
- For every edge (v, w) in the graph, there is a subset X_i that contains both v and w .
- If X_i and X_j both contain a vertex v , then all nodes X_k of the tree in the (unique) path between X_i and X_j contain v as well.

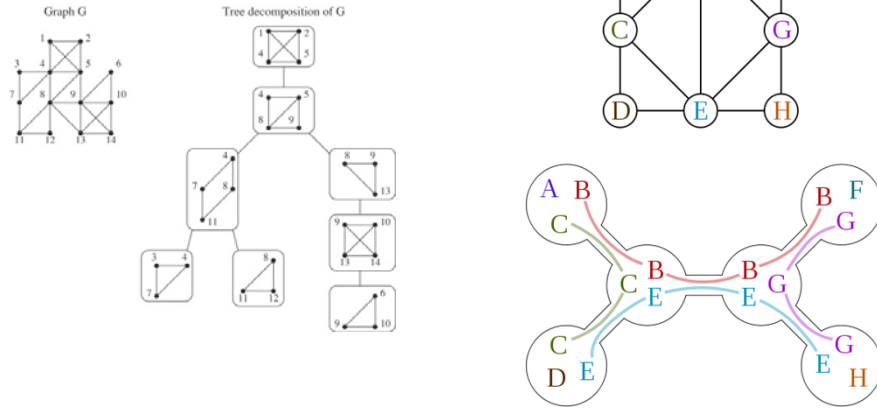
Treewidth

- The *width* of a tree decomposition is the size of its largest set X_i *minus one*.
- The **treewidth** $\text{tw}(G)$ of a graph G is the minimum width among all possible tree decompositions of G .
 - *minus one* in order to make the treewidth of a tree equal to one.

The Treewidth of Graphs

- Every **complete graph** K_n has treewidth $n - 1$.
- A connected graph with at least two vertices has treewidth 1 if and only if it is a **tree**.
- If a graph has a **cycle**, its treewidth is at least two.
- It is NP-complete to determine whether a given graph G has treewidth at most a given variable k .
 - For fixed k , we can check if a graph has a treewidth of k in $O(n^k)$, and find the tree decomposition.

Examples



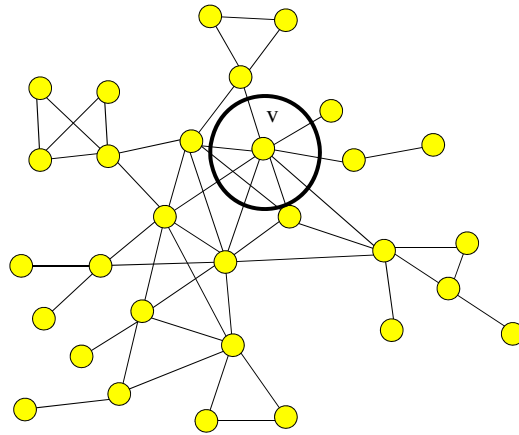
Why Treewidth?

Many algorithmic problems that are NP-complete for arbitrary graphs may be solved efficiently by dynamic programming for graphs of bounded treewidth, using the tree-decompositions of these graphs.

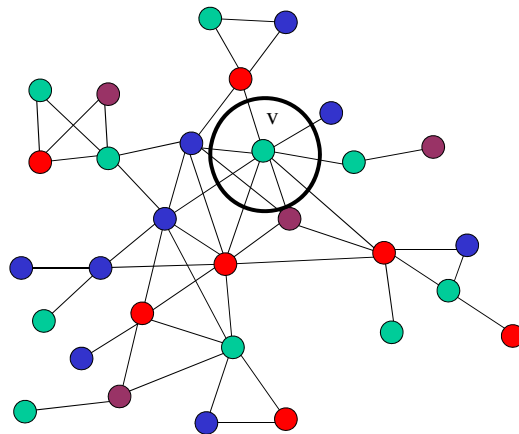
Example: k -coloring a graph

We will approximately count some graphlets based on color-coding alg. [Alon, Yuster, and Zwick, JACM 1995]

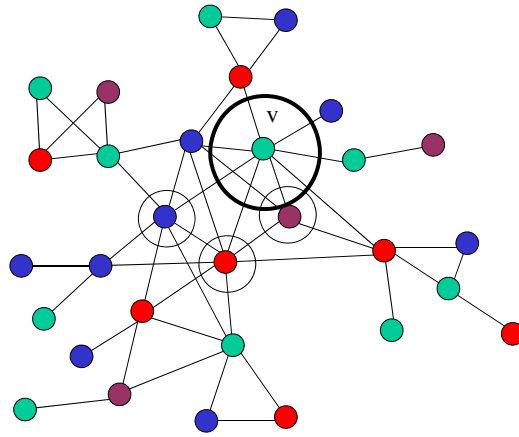
Counting Graphlets Using Color Coding



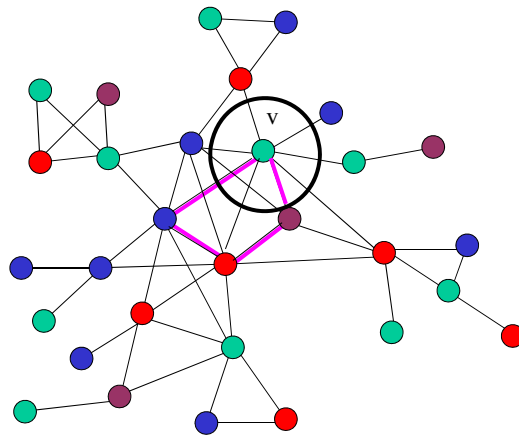
Counting Graphlets Using Color Coding



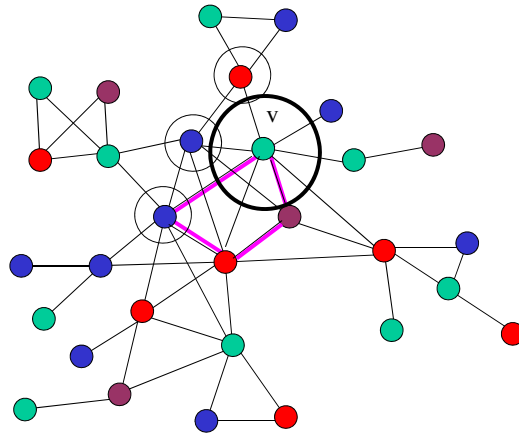
Counting Graphlets Using Color Coding



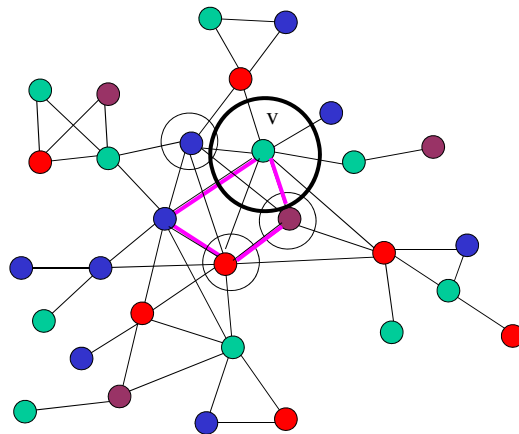
Counting Graphlets Using Color Coding



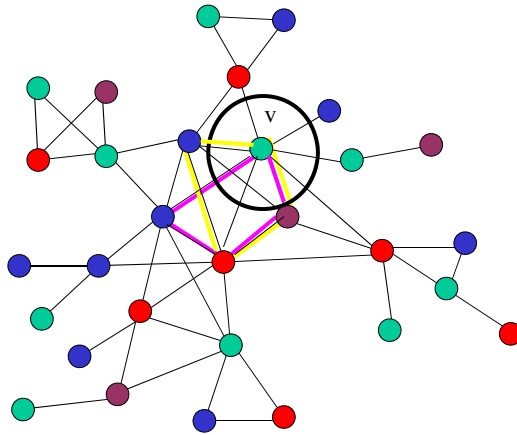
Counting Graphlets Using Color Coding



Counting Graphlets Using Color Coding



Counting Graphlets Using Color Coding



Counting Trees in a Graph (also graphlets with bounded treewidth)

Given a network G with n vertices and a tree T with k vertices, we consider the problem of counting the number of non-induced subtrees of G that are isomorphic to T .

Assume $k \in O(\log n)$

[Alon, Yuster, and Zwick, JACM 1995]

[Alon *et al.*, Bioinformatics 2008]

[Gonen & Shavitt, Internet Mathematics 2009]

The Counting Algorithm

1. **Color coding.** Color each vertex of input graph G independently and uniformly at random with one of the k colors.
2. **Counting.** Apply a dynamic programming routine (explained later) to count the number of non-induced occurrences of T in which each vertex has a unique color.
3. Repeat the above two steps $O(e^k)$ times and add up the number of occurrences of T to get an estimate on the number of its occurrences in G .

Color Coding

r - the total number of copies of T in G .

We assign a color (ind. and uar) to each vertex of G from the color set $[k]=\{1,\dots,k\}$.

For some occurrence of T , the probability that all T 's vertices are assigned unique colors is $p=k!/k^k$, thus the expected number of colorful copies in G is rp .

Let \mathcal{F} denote the family of all copies of T in G .

For each such copy $F \in \mathcal{F}$, let x_F denote the indicator random variable whose value is 1 if and only if the copy is colorful in our random k -coloring of $V(G)$, the vertices of G .

Let $X = \sum_{F \in \mathcal{F}} x_F$ be the random variable counting the total number of colorful copies of T .

By linearity of expectation, the expected value of X is $E(X) = rp$.

Estimating X Variance

For every two distinct copies $F, F' \in \mathcal{F}$, the probability that both F and F' are colorful is at most p .

- in fact strictly smaller unless both copies have exactly the same set of vertices.

\Rightarrow the covariance $\text{Cov}(x_F, x_{F'})$ satisfies:

$$\text{Cov}(x_F, x_{F'}) = E(x_F x_{F'}) - E(x_F)E(x_{F'}) \leq p.$$

Estimating X Variance (cont.)

The variance of X :

$$\begin{aligned}\text{Var}(X) &= \sum_{F \in \mathcal{F}} \text{Var}(x_F) + \sum_{F \neq F' \in \mathcal{F}} \text{Cov}(x_F, x_{F'}) \\ &\leq rp + r(r-1)p = r^2p.\end{aligned}$$

Y = the average of s independent copies of X
(obtained by s independent random colorings)

$$E(Y) = E(X) = rp$$

$$\text{Var}(Y) = \text{Var}(X)/s \leq r^2p/s.$$

Estimating X Variance (cont.)

Applying Chebyshev's inequality

Let X be a random variable with finite expected value μ and finite non-zero variance σ^2 . Then for any real number $k > 0$,

$$\Pr(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}.$$

the probability that Y is smaller than (or bigger than) its expectation by at least ϵrp is at most

$$\frac{r^2p}{\epsilon^2 r^2 p^2 s} = \frac{1}{\epsilon^2 ps}.$$

if $s = 4/\epsilon^2 p$ this probability is at most $1/4$.

Decreasing the Error in Y

Compute Y t times independently.

Let Z be the median.

The probability that the median is less than $(1-\epsilon)rp$ is the probability that at least half of the copies of Y computed will be less than this quantity, which is at most $\binom{t}{t/2}4^{-t} \leq 2^{-t}$.

A similar estimate holds for the probability that Z is bigger than $(1+\epsilon)rp$.

Decreasing the Error in Y (cont.)

If $t = \log(1/\delta)$ then with probability $1-2\delta$ the value of Z will lie in $[(1-\epsilon)rp, (1+\epsilon)rp]$.

Note that the total number of colorings in the process is ts :
$$o\left(\frac{\log(1/\delta)}{\epsilon^2 p}\right) = o\left(\frac{e^k \log(1/\delta)}{\epsilon^2}\right).$$

Our estimate for r is $Z/p = Zk^k/k!$

Counting Paths of length k

Let $C(v, S)$ be the number of colorful paths for which one of the endpoints is v .

S is a subset of the color set $\{1, \dots, k\}$,

$\text{col}(v)$ is the color of vertex v .

Given a color ℓ , for all $v \in V(G)$:

$$C(v, \{\ell\}) = \begin{cases} 1 & \text{if } \text{col}(v) = \ell \\ 0 & \text{otherwise.} \end{cases}$$

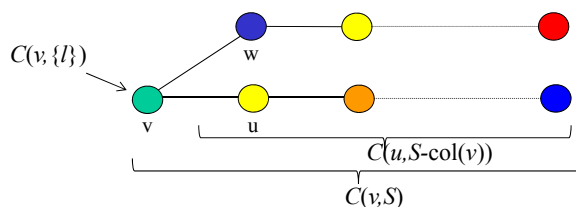
Counting Paths of length k (cont.)

For each vertex v and color set S where $|S| > 1$:

$$C(v, S) = \sum_{u: (u, v) \in E(G)} C(u, S - \text{col}(v)).$$

The number of single colorful paths of length k is

$$\frac{1}{2} \sum_v C(v, \{1, \dots, k\}).$$



Counting Trees

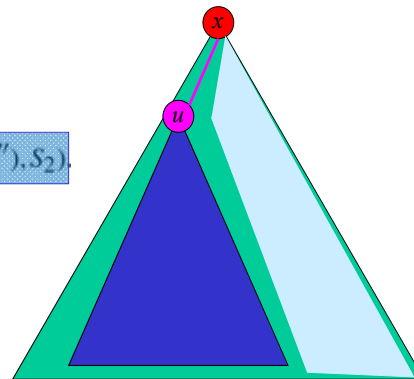
- ρ – an arbitrary vertex of the tree selected to be the root.
- $\tau(\rho)$ – the rooted tree
- For each vertex v of the graph G , we compute $c(v, \tau(\rho), [k])$, the number of $[k]$ -colorful rooted subtrees with root v , which are isomorphic to $\tau(\rho)$.

- The actual number of $[k]$ -colorful occurrences of T in G is $\frac{1}{q} \sum_v c(v, \tau(\rho), [k])$

q is equal to the number of vertices u in T , for which the rooted tree $\tau(u)$ is isomorphic to $\tau(\rho)$.

The Dynamic Program

$$c(x, \tau'(\rho'), S) = \frac{1}{d} \sum_{\forall S_1, S_2: |S_1 \cap S_2| = \emptyset} c(x, \tau_1(\rho'), S_1) \cdot c(u, \tau_2(\rho''), S_2)$$



Complexity of the dynamic program: $O(2^k |E|)$
 Overall complexity:

$$O(2^k |V| |E|) \cdot O\left(|E| \cdot 2^k \cdot e^k \log(1/\delta) \cdot \frac{1}{\epsilon^2}\right)$$

Remarks

- The algorithm can be extended for graphlets with bounded treewidth.
- If the graphlet sized in $O(\log n)$ the $(2e)^k$ term in the complexity is $O(n)$
- We can use this algorithms to count graphlets in a graph, graphlets attached to a node, and orbits (at least for trees).

Counting Cycles

Algorithm 1. (An (ϵ, δ) -approximation algorithm for counting simple cycles of length k adjacent to a vertex v)

1. For $j = 1$ to t
 - (a) For $i = 1$ to s
 - i. Color each vertex of G independently and uniformly at random with one of the k colors.
 - ii. For all $x \in V$, $C_i(v, x, [k]) = \text{count-path}(v, x, k)$. [See Algorithm 2.]
 - iii. Let $CY_i(v, [k]) = \frac{1}{2} \sum_{u \in N(v)} C_i(v, u, [k])$.
 - iv. Let $X_i^v = CY_i(v, [k])$.
 - (b) Let $Y_j^v = \frac{\sum_{i=1}^s X_i^v}{s}$.
 2. Let Z^v be the median of Y_1^v, \dots, Y_t^v .
 3. Return $Z^v \cdot k^k / k!$.
-

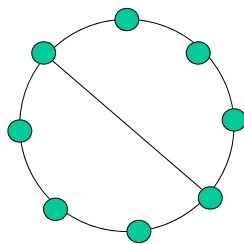
Time Complexity

As before

$O(2^k n |E|)$ – for calculating the paths for all v

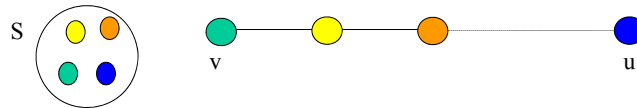
Overall complexity: $O((2e)^k \cdot |E| \cdot |V| \cdot \log(1/\delta) / \epsilon^2)$

Cycle with a Chord

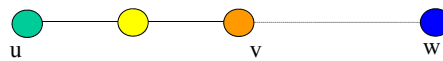


Mathematical Notations

- $C(v,u,S)$ = the number of colorful paths from v to u in a specific coloring, using the colors in S .



- $P(v,u,w,S)$ = the number of colorful paths from u to w that are adjacent to v in a specific coloring, using the colors in S .



The Algorithm

Algorithm's Input:

A graph $G(V,E)$, a vertex v , fault-tolerance ϵ , error probability δ

Notation: let $A_{v,z,b}(S)$ be the set of all pairs (S_1, S_2) such that the following hold:

- $|S_1| = z+1$,
- $|S_2| = b-z+1$,
- $S_1 \cup S_2 = S$,
- $S_1 \setminus \{\text{col}(u) | u \in V'\} \cap S_2 \setminus \{\text{col}(u) | u \in V'\} = \phi$

The Algorithm

Repeat
 $t = \log(1/\delta)$
 times

1. Color each vertex of G independently and uniformly at random with one of the k colors.
2. Compute the number of cycles with a chord in the coloring; there are two cases:

Case 1

k-l

Case 2

k-l

 - 2.1 Compute $X_{1,v}$ = the number of k -length cycles with a chord in case 1.
 - 2.2 Compute $X_{2,v}$ = the number of k -length cycles with a chord in case 2.
3. Let Y_v = the average of all the s $X_{1,v} + X_{2,v}$.
4. Return the median of all the t Y_v multiplied by $k^k/k!$.

Repeat
 $s = 4k^k/\epsilon^2 k!$
 times

Computing the number of paths between v, w , for every color-set S

- For all $S \subseteq [k]$ s.t. $S = \{\ell\}$ $C(v, w, S) = 1$ if $\text{col}(v) = \text{col}(w) = \ell$, and 0 otherwise.
- For $q = 2$ to k , for all $S \subseteq [k]$ s.t. $|S| = q$

$$C(v, w, S) = \sum_{u \in N(v)} C(u, w, S \setminus \{\text{col}(v)\})$$

Case 1

- $P(v,u,w,S) = \sum_{1 \leq z \leq l-1} \sum C(v,w,S_1) \cdot C(v,u,S_2)$

z-length colorful paths between u and w using colors in S_1

$l-z$ -length colorful paths between v and u using colors in S_2

The sum is over all (S_1, S_2) in $A_{v,z,l}(S)$

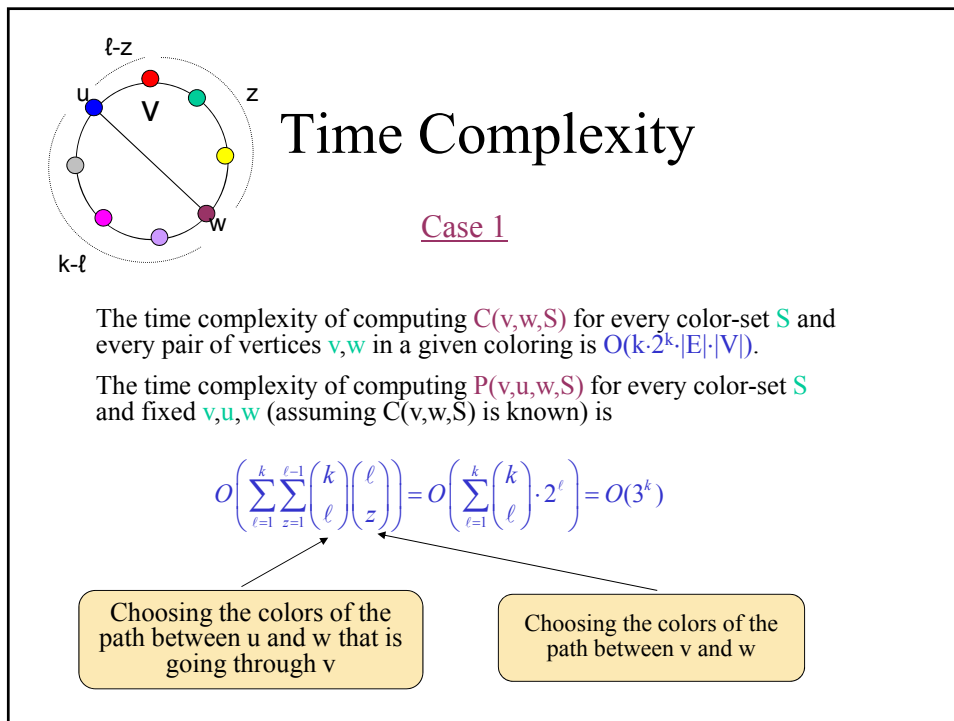
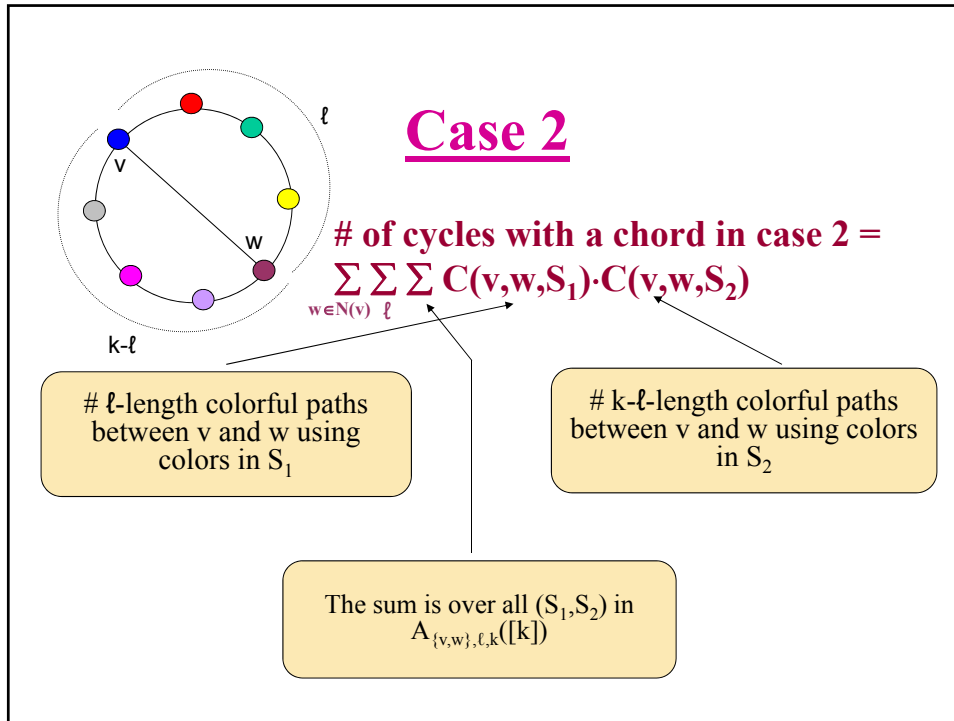
Case 1

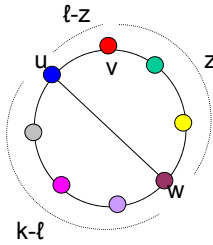
- # of cycles with a chord in case 1 = $\sum_{(u,w) \in E} \sum_{\ell} \sum P(v,u,w,S_3) \cdot C(u,w,S_4)$

l -length colorful paths between u and w that are adjacent to v using colors in S_3

$k-l$ -length colorful paths between u and w using colors in S_4

The sum is over all (S_3, S_4) in $A_{\{u,w\}, \ell, k}([k])$ and over all (S_3, S_4) in $A_{\{u,w\}, k-\ell, k}([k])$



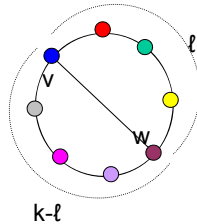


Time Complexity

Case 1

→ The time complexity of the first case, for every edge (u,w) , every vertex v , and every color-set S , in a given coloring is

$$O\left(\sum_{v \in V} \sum_{(u,w) \in E} 3^k\right) = O(3^k |E| |V|)$$



Time Complexity

Case 2

The time complexity of the second case, for every vertex v , every neighbor of v , and every color-set S , in a given coloring is

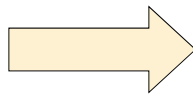
$$O\left(\sum_{v \in V} \sum_{w \in N(v)} \sum_{\ell=1}^k \binom{k}{\ell}\right) = O(2^k \cdot |E|)$$

Choosing the colors of the path between v and w

Time Complexity

→The total time complexity, for every vertex v and every color-set S , in a given coloring, is

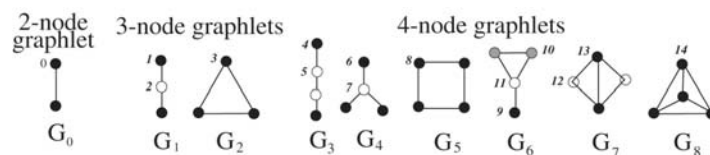
$$O(3^k \cdot |E| \cdot |V|).$$



$$\text{Time complexity} = O((3e)^k \cdot |E| \cdot |V| \cdot \log(1/\delta) / \epsilon^2)$$

Small Graphlets

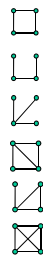
- In many practical cases we are only interested in the graphlets of small size, ≤ 5
- We have already seen how to count triangles (and 3 node paths)
- How efficient can we find 4-node graphlets?



[Marcus & Shavitt, Computer Networks 2012]

Maybe we can use our approx.?

Motif



Time Complexity

$O(|E| \cdot |V| \cdot \log(1/\delta) / \epsilon^2)$

$O(|E| \cdot \log(1/\delta) / \epsilon^2)$

$O(|E| \cdot \log(1/\delta) / \epsilon^2)$

$O(|E| \cdot |V| \cdot \log(1/\delta) / \epsilon^2)$

(ϵ, δ) – Approximation:

$\Pr[(1-\epsilon) \cdot \#f \leq y \leq (1+\epsilon) \cdot \#f] \geq 1-2\delta$

Optimal Algorithms

- Counting all orbits
- Assume nodes labels are $\{1, 2, \dots, n\}$



Triangle Counting Per Node

All algs use **Merge**

In **Merge**:

NodeArray is not useful for counting Δ s, will be needed by other algorithms.

$O(d|E|)$

Algorithm 1. Counting triangles

```

1: procedure TRIANGLECOUNT( $G$ )
2:   for all  $v \in V(G)$  do
3:     for all  $u \in N(v), v < u$  do
4:        $Merged \leftarrow \text{Merge}(G, v, u, V_{arr})$ 
5:        $m_7[v] \leftarrow m_7[v] + |Merged|$ 
6:        $m_7[u] \leftarrow m_7[u] + |Merged|$ 
7:   for all  $v \in V(G)$  do
8:      $m_7[v] \leftarrow m_7[v]/2$ 
9: procedure MERGE( $G, v, u, NodeArray$ )
10:  for all  $w \in N(v)$  do
11:     $NodeArray[w] \leftarrow 0$ 
12:  for all  $w \in N(u)$  s.t.  $w \neq v$  do
13:     $NodeArray[w] \leftarrow 1$ 
14:  for all  $w \in N(v), w \neq u$  do
15:    if  $NodeArray[w] = 1$  then
16:       $NodeArray[w] \leftarrow 3$ 
17:       $list \leftarrow \text{AppendToList } list, w$ 
18:    else
19:       $NodeArray[w] \leftarrow 2$ 
20:  return list
  
```

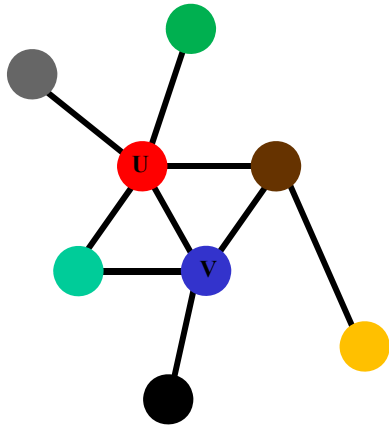
Counting Cycles With Chord

Algorithm 2. Counting non-induced 4-cycles with a chord

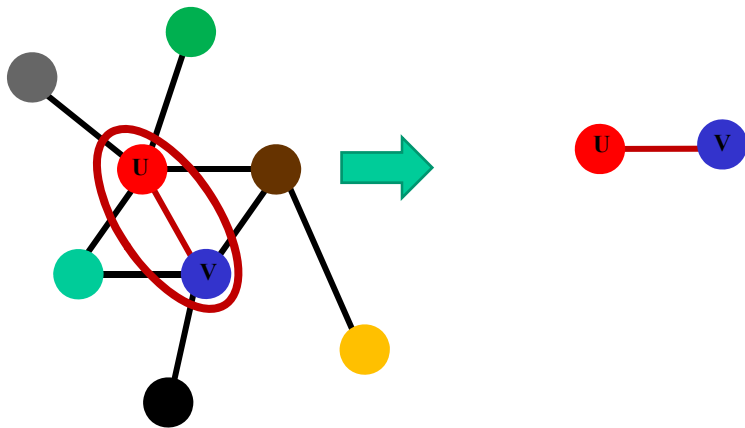
```

1: procedure CHORDCYCLECOUNT( $G$ )
2:   for all  $v \in V(G)$  do
3:     for all  $u \in N(v), v < u$  do
4:        $Merged \leftarrow \text{Merge}(G, v, u, V_{arr})$ 
5:        $m_{2,1}[v] \leftarrow m_{2,1}[v] +$ 
6:          $|Merged| \cdot (|Merged| - 1)/2$ 
7:        $m_{2,1}[u] \leftarrow m_{2,1}[u] +$ 
8:          $|Merged| \cdot (|Merged| - 1)/2$ 
9:       for all  $w \in Merged$  do
10:         $m_{2,2}[w] \leftarrow m_{2,2}[w] + (|Merged| - 1)$ 
  
```

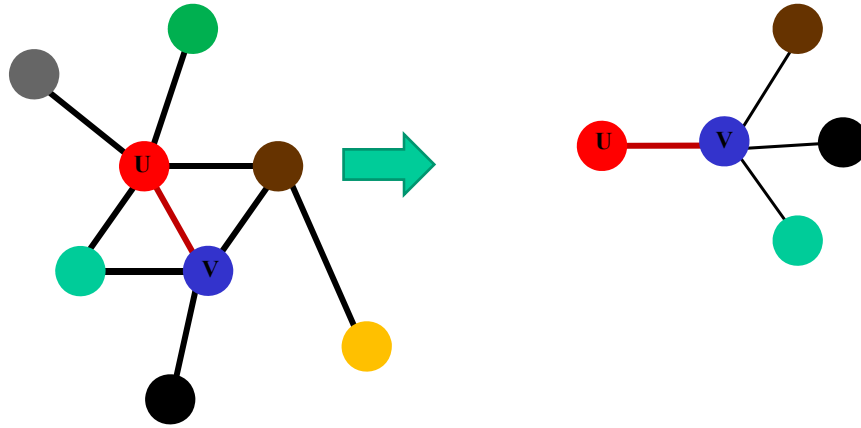
Counting Cycles With Chord



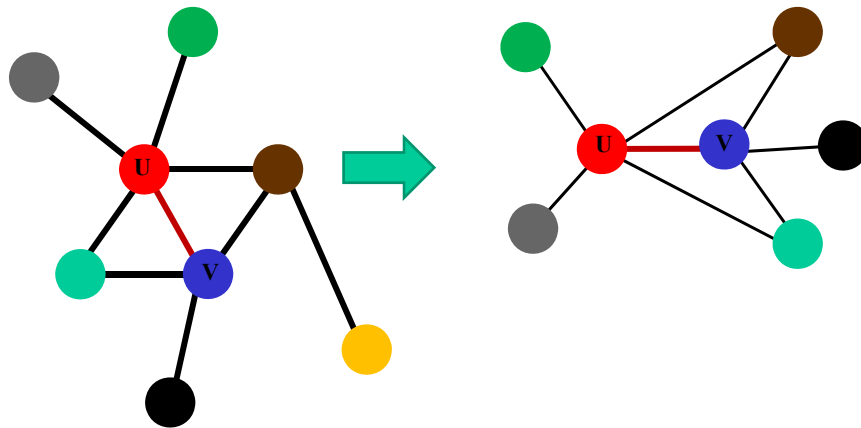
Counting Cycles With Chord



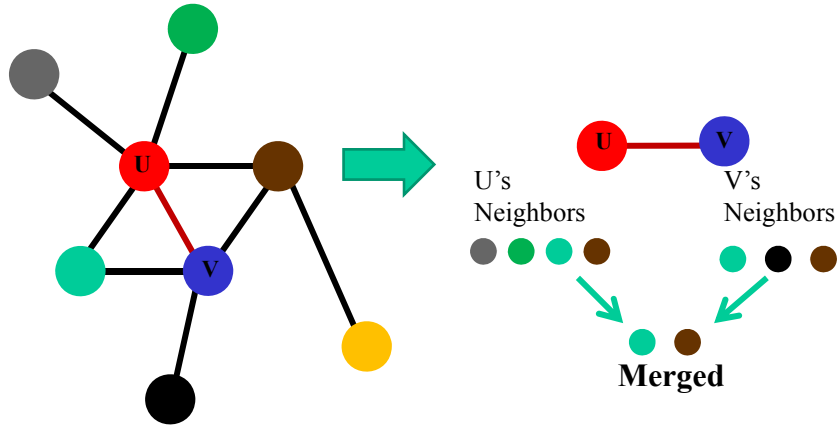
Counting Cycles With Chord



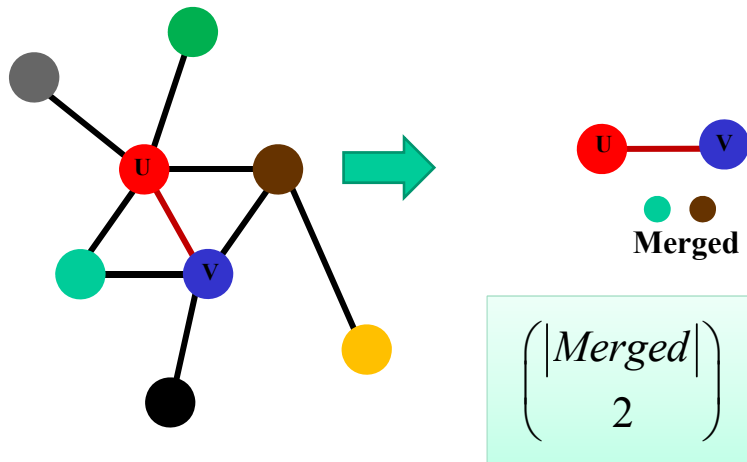
Counting Cycles With Chord Example



Counting Cycles With Chord

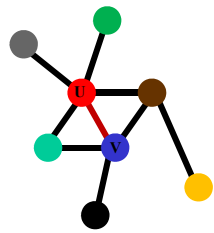


Counting Cycles With Chord



Counting Cycles With Chord

Runtime Analysis



$$O(|E| * (2 * d \log d + d + d)) = O(|E| * d \log d)$$

Iterate
over all
Edges

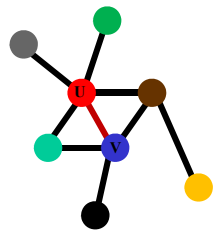
Sort both
Neighbor
lists

Merge
Neighbor
lists

Update count
for relevant
nodes

Counting Cycles With Chord

Runtime Analysis



$$O(|E| * (2 * d \log d + d + d)) = O(|E| * d \log d)$$

*we can get time complexity of **$O(|E| * d)$**
Assuming graph vertices are labeled by the integers $\{1 \dots n\}$

Counting Tail Triangles

Algorithm 3. Counting non-induced tailed triangles

```

1: procedure TAILTRIANGLECOUNT( $G$ )
2:   TriangleCount( $G$ )
3:   for all  $v \in V(G)$  do
4:     for all  $u \in N(v), v < u$  do
5:        $Merged \leftarrow \text{Merge}(G, v, u, V_{arr})$ 
6:        $tails_v \leftarrow \max\{0, (|N(v)| - 2)\}$ 
7:        $tails_u \leftarrow \max\{0, (|N(u)| - 2)\}$ 
8:       for all  $w \in Merged$  do
9:          $m_{4.2}[w] \leftarrow m_{4.2}[w] + tails_v + tails_u$ 
10:    for all  $v \in V(G)$  do
11:       $m_{4.1}[v] \leftarrow \max\{0, (m_7[v] \cdot (|N(v)| - 2))\}$ 
12:      for all  $u \in N(v)$  do
13:         $m_{4.3}[v] \leftarrow m_{4.3}[v] + m_7[u]$ 
14:         $m_{4.3}[v] \leftarrow m_{4.3}[v] - 2 \cdot m_7[v]$ 

```

$O(d|E|)$

Counting 4 Cliques

Algorithm 4. Counting 4-cliques

```

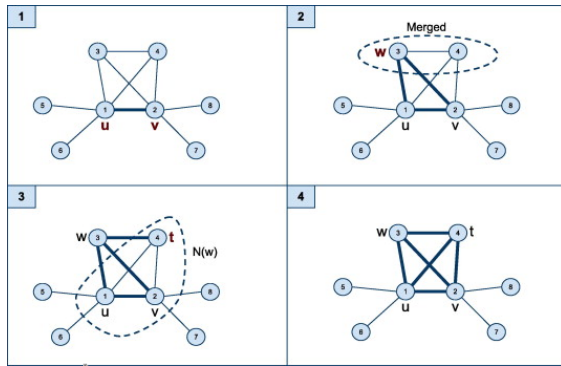
1: procedure CLIQUECOUNT( $G$ )
2:   for all  $v \in V(G)$  do
3:      $V_{arr}[v] \leftarrow 0$ 
4:   for all  $v \in V(G)$  do
5:     for all  $u \in N(v), v < u$  do
6:        $Merged \leftarrow \text{Merge}(G, v, u, V_{arr})$ 
7:       for all  $w \in Merged$  do
8:         for all  $t \in N(w), w < t$  do
9:           if  $V_{arr}[t] = 3$  then  $((*)$ 
10:             $m_1[v] \leftarrow m_1[v] + 1$ 
11:             $m_1[u] \leftarrow m_1[u] + 1$ 
12:           for all  $w \in N(v) \cup N(u)$  do
13:              $V_{arr}[w] \leftarrow 0$ 
14:           for all  $v \in V(G)$  do
15:              $m_1[v] \leftarrow m_1[v] / 3$ 

```

$((*) V_{arr}$ is set in the MERGE call (line 6). $V_{arr}[t] = 3$ iff t is a neighbor of both v and u

Complexity

$$\begin{aligned}
 & O\left(\sum_{v \in V} \sum_{u \in N(v)} \left(d + \sum_{w \in N(v) \cap N(u)} |N(w)|\right)\right) \\
 &= O\left(\sum_{v \in V} \sum_{u \in N(v)} (d) + \sum_{v \in V} \sum_{u \in N(v)} \sum_{w \in N(v) \cap N(u)} |N(w)|\right) \\
 &= O(d|E|) + O(|E|^2). \quad \square
 \end{aligned}$$



In bounded degree graphs $\sum_w N(w) = d^2$
 (instead of $|E|$) so the complexity is $O(d^2|E|)$

Counting 4 Cycles



Algorithm 5. Counting 4-cycles

```

1: procedure CYCLECOUNT(G)
2:   for all v in V(G) do
3:     V_arr[v] ← 0
4:   for all v in V(G) do
5:     for all u in N(v) s.t. v < u do
6:       Merge G, v, u, V_arr
7:       for all w in N(v) \ {u} do
8:         for all t in N(w) do
9:           if V_arr[t] = 1 or V_arr[t] = 3 then
10:            m_3[v] ← m_3[v] + 1
11:            m_3[u] ← m_3[u] + 1
12:       for all w in N(v) ∪ N(u) do
13:         V_arr[w] ← 0
14:   for all v in V(G) do
15:     m_3[v] ← m_3[v] / 2
    
```

$$\left(\sum_{v \in V} \sum_{u \in N(v)} \left(d + \sum_{w \in N(v)} |N(w)|\right)\right) = O(d|E|) + O(|E|^2)$$



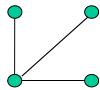
$O(d|E|)$

Algorithm 6. Counting 4-node paths

```

1: procedure PATHCOUNT( $G$ )
2:   for all  $v \in V(G)$  do
3:      $V_{arr}[v] \leftarrow 0$ 
4:   for all  $v \in V(G)$  do
5:     for all  $u \in N(v)$  s.t.  $v < u$  do
6:        $Merged \leftarrow MergeG, u, v, V_{arr}$ 
7:        $m_{6,2}[v]$ 
8:      $\leftarrow m_{6,2}[v] + (|N(v) \setminus \{u\}| \cdot |N(u) \setminus \{v\}|) - |Merged|$ 
9:        $m_{6,2}[u]$ 
10:     $\leftarrow m_{6,2}[u] + (|N(v) \setminus \{u\}| \cdot |N(u) \setminus \{v\}|) - |Merged|$ 
11:    for all  $w \in N(v) \setminus \{u\}$  do
12:      if  $V_{arr}[w] = 3$  then
13:         $m_{6,1}[w] \leftarrow m_{6,1}[w] + |N(u) \setminus \{v\}| - 1$ 
14:      else
15:         $m_{6,1}[w] \leftarrow m_{6,1}[w] + |N(u) \setminus \{v\}|$ 
16:    for all  $w \in N(u) \setminus \{v\}$  do
17:      if  $V_{arr}[w] = 3$  then
18:         $m_{6,1}[w] \leftarrow m_{6,1}[w] + |N(v) \setminus \{u\}| - 1$ 
19:      else
20:         $m_{6,1}[w] \leftarrow m_{6,1}[w] + |N(v) \setminus \{u\}|$ 
21:    for all  $w \in N(v) \cup N(u)$  do
22:       $V_{arr}[w] \leftarrow 0$ 

```



Counting Claws

$$m_{5,1}[v] = \binom{|N(v)|}{3},$$

$$m_{5,2}[v] = \sum_{u \in N(v)} \binom{|N(u)| - 1}{2}$$

Summary for Graphlets

Motif

Time Complexity



$O(|E|^2)$



$O(d \cdot |E|)$



$O(|V|)$



$O(d \cdot |E|)$

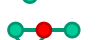
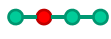


$O(d \cdot |E|)$



$O(|E|^2) \text{ or } O(|E|d^2)$

Runtime Analysis for Orbits



$O(d \cdot |E|)$



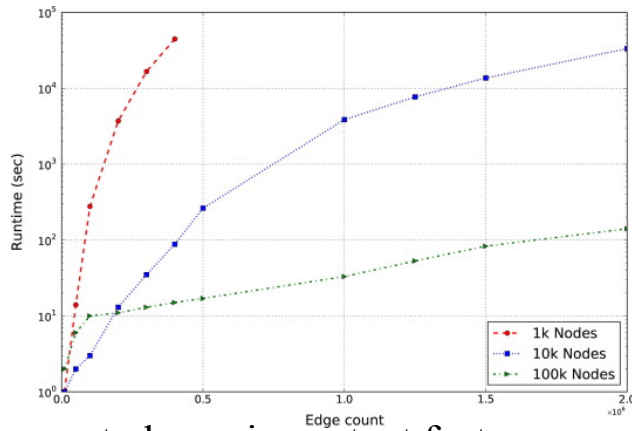
$O(|V|)$

$O(|E|)$

$O(d \cdot |E| + |E|^2)$

For any real-world graph
 $O(|E|) = O(|V|)$

Performance



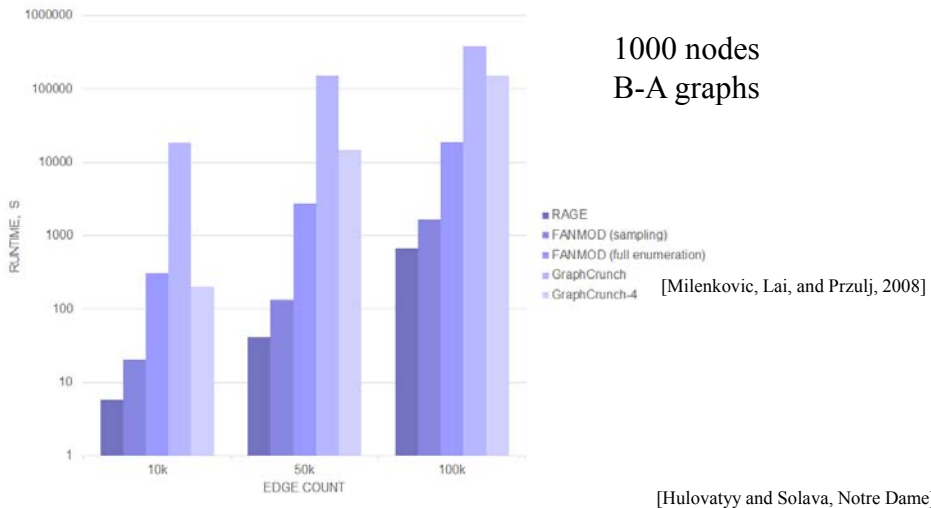
- Density seems to be an important factor

Comparison

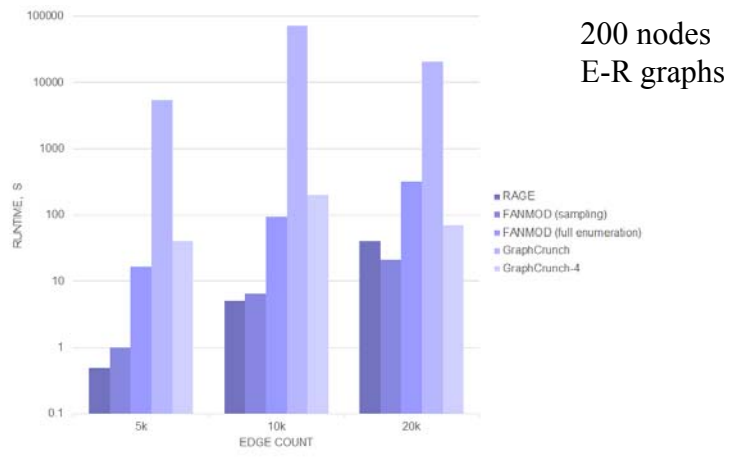
	8k edges	20k edges	52k edges	92.6k edges
Method	5k nodes	10k nodes	20k nodes	AS-Graph (26k nodes)
RAGE	11	64	720	2400
FANMOD (sampling)	57	210	1020	7200
FANMOD (full-enum)	420	2040	25200	172800
IML	140	1203	9850	25000

Rage is faster even than FANMOD sampling

Comparison



Comparison

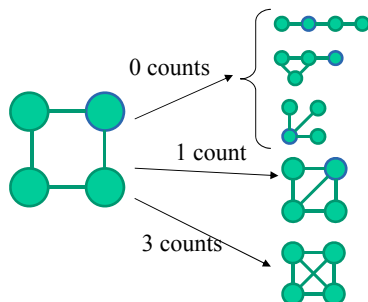


Inferring Induced Count

4 Node Cycles



- Fix the over-count occurred due to subgraphs with access edges



$$\begin{aligned} \text{Induced Count} &= \text{NonInduced}(\text{4-node cycle}) - 3 * \text{Induced}(\text{3-node path}) - 1 * \text{Induced}(\text{2-node path}) \\ &= \text{NonInduced}(\text{4-node cycle}) + 3 * \text{NonInduced}(\text{3-node path}) - 1 * \text{NonInduced}(\text{2-node path}) \end{aligned}$$

Inferring Induced Counts of Orbits

$$\text{Induced Count}(\text{4-node cycle with 1 red node}) = \text{Non Induced Count}(\text{4-node cycle with 1 red node}) - 3 * \text{Non-Induced Count}(\text{3-node path with 1 red node})$$

$$\text{Induced Count}(\text{4-node cycle with 2 red nodes}) = \text{Non Induced Count}(\text{4-node cycle with 2 red nodes}) - 3 * \text{Non-Induced Count}(\text{3-node path with 2 red nodes})$$

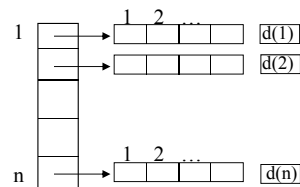
$$\begin{aligned} \text{Induced Count}(\text{4-node cycle with 3 red nodes}) &= \text{Non Induced Count}(\text{4-node cycle with 3 red nodes}) \\ &\quad - 2 * \text{Induced Count}(\text{3-node path with 3 red nodes}) \\ &\quad - 3 * \text{Induced Count}(\text{2-node path with 3 red nodes}) \end{aligned}$$

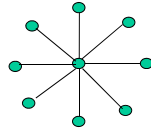
Sublinear Algorithm

- Assume a very large graph
- We can query the graph
 - Minimize the query number
- Sublinear approximation

Counting Stars: Model

- We assume graphs are represented by the **incidence lists** of the vertices, where each list is accompanied by its length.
- **Allowed queries:**
 - what is the degree, $d(v)$, of any vertex v ?
 - who is the i 'th neighbor of v , for any vertex v and index $1 \leq i \leq d(v)$?





Approximating Stars

Upper Bound:

- Given an approximation parameter $0 < \epsilon < 1$ and query access to a graph G , our algorithm outputs an estimate v'_s such that, with high constant probability,

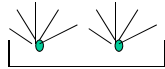
$$(1-\epsilon)v_s(G) \leq v'_s \leq (1+\epsilon)v_s(G),$$

where $v_s(G)$ denotes the number of stars of size $s+1$ in the graph.

- The expected query complexity and running time of the algorithm are

$$O\left(\frac{n}{v_s(G)^{\frac{1}{s+1}}} + \min\left\{n^{\frac{1}{s}}, \frac{n^{\frac{s-1}{s}}}{v_s(G)^{\frac{1}{s}}}\right\} \cdot \text{poly}(\log n, 1/\epsilon)\right)$$

Upper Bound-Main Idea

- Consider a partition of the graph vertices into $O(\log n/\epsilon)$ buckets where in each bucket all vertices have the same degree (with respect to the entire graph) up to a multiplicative factor of $(1 \pm O(\epsilon))$. The degree in bucket B_i is $\sim (1+\beta)^i$, $\beta = O(\epsilon)$.
- If we could get a good estimate of the size of each bucket by sampling, then we would have a good estimate of the number of s -stars (since the vertices in each bucket are the centers of approximately the same number of stars). 
- The difficulty is that some buckets may be very small and we might not even hit them when sampling vertices. However, these buckets can significantly contribute to the number of stars in the graph. 