

Interactive Large-Scale Data and Graph Analytics

Graph Analytics

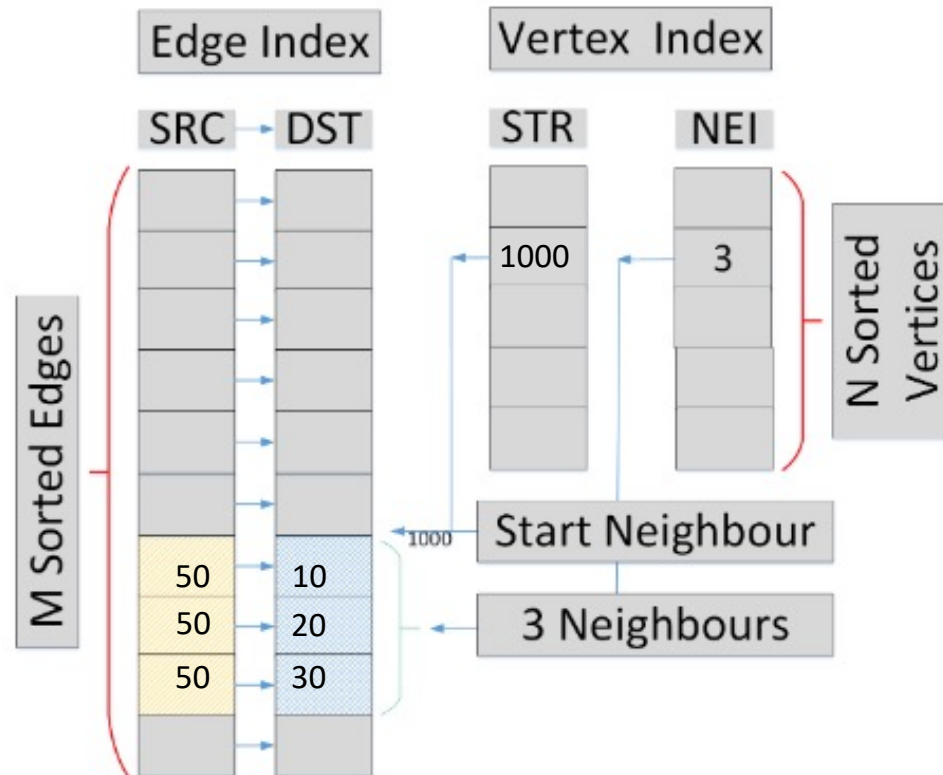
Oliver Alvarado Rodriguez, Naren Khatwani, Zhihui Du, David Bader

Outline

1. Explanation of the data structure.
2. Overview of our algorithmic graph theory.
3. Example of graph analytics from a Jupyter notebook.

Arachne Double-Index (DI) Data Structure

[Alvarado Rodriguez, Du, Patchett, Li, Bader 2022]



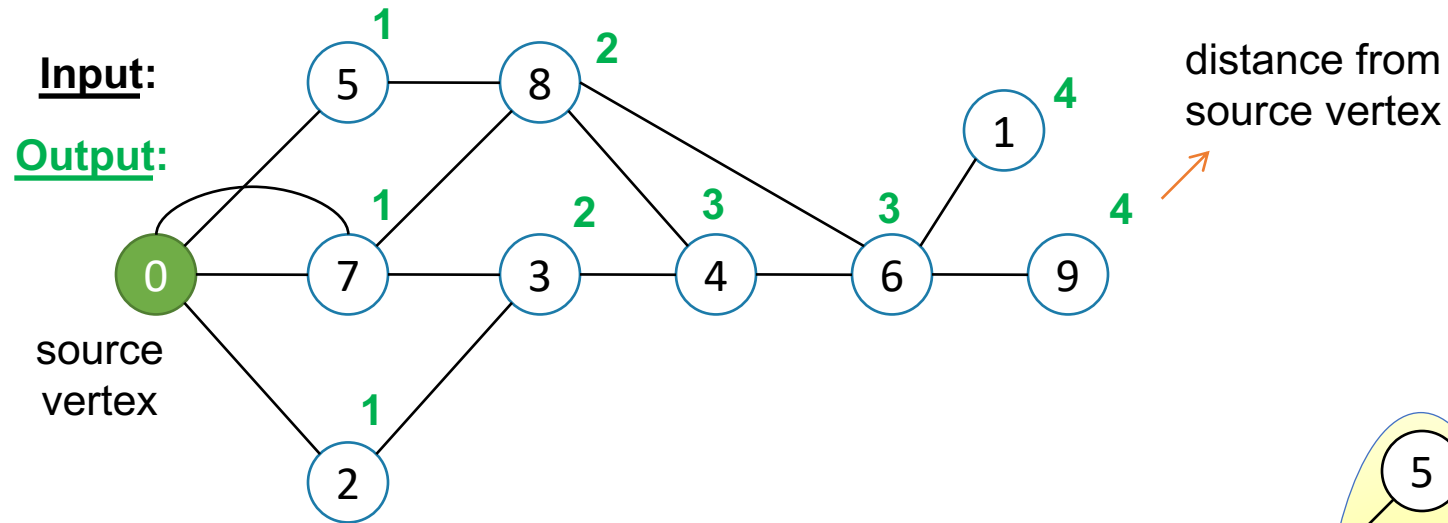
Advantages of DI over CSR:

1. $O(1)$ time complexity:
 - locating a vertex from a given edge ID.
 - locating the adjacency list from a given vertex ID.
2. We can search from edge ID to vertex ID, this is not possible in CSR.
3. DI can support both edge-centric and vertex-centric algorithms whereas CSR can only support the latter.
4. DI can easily achieve load balancing with the edge array being distributed equally amongst many locales.

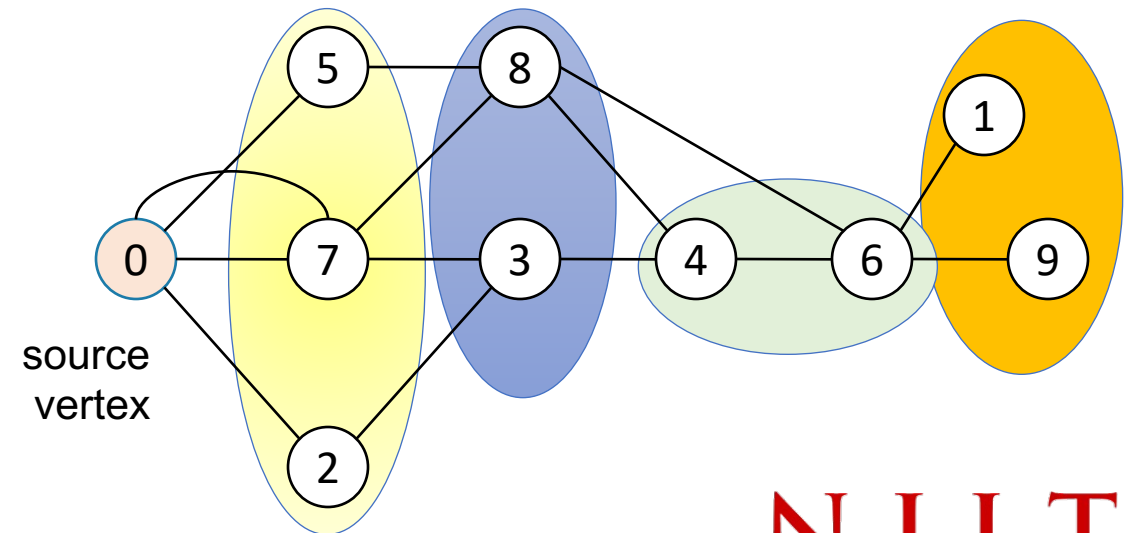
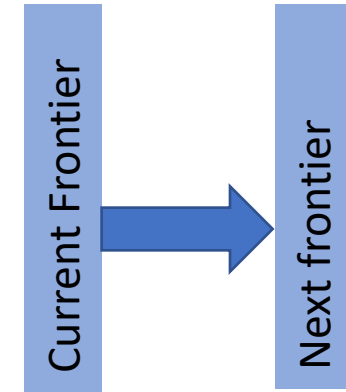
Graph Algorithms in Arachne

- **Breadth-first search** [Du, Alvarado Rodriguez, Bader 2021]
Returns an array of size n with how many hops away some vertex v is from an initial vertex u .
- **Connected components** [Du, Alvarado Rodriguez, Bader 2021]
Returns an array of size n where all vertices who belong to the same component have the same value x . The value of x is the label of the largest vertex in the component.
- **Triangle counting** [Du, Alvarado Rodriguez, Patchett, Bader 2021]
Returns the number of triangles in a graph.
- **Truss Analytics** [Du, Patchett, Bader 2021][Du, Patchett, Alvarado Rodriguez, Li, Bader 2022]
K-truss returns every edge in the truss where each edge must be a part of $k - 2$ triangles that are made up of nodes in that truss. Max truss returns the maximum k . Truss decomposition returns the maximum k for each edge.
- **Triangle centrality** [Patchett, Du, Bader 2022][Patchett, 2022]
Returns an array of size n with the proportion of triangles centered at a vertex v .

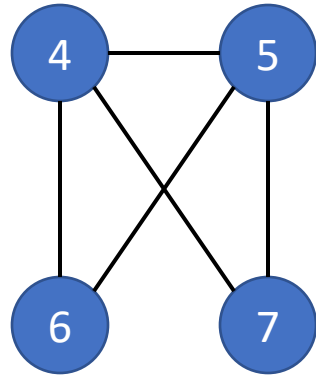
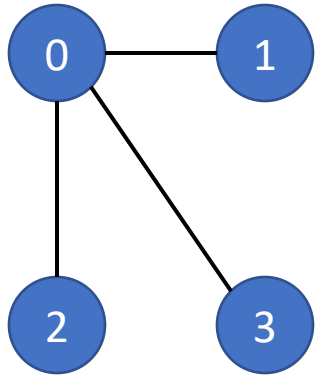
Parallel Breadth-First Search Example



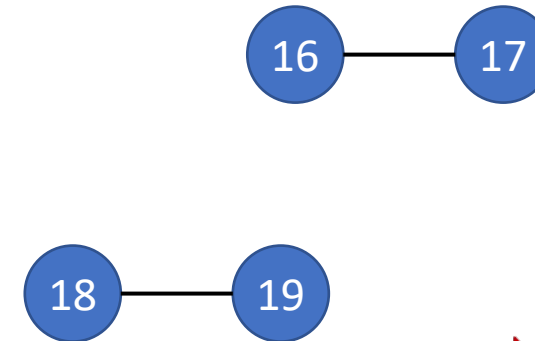
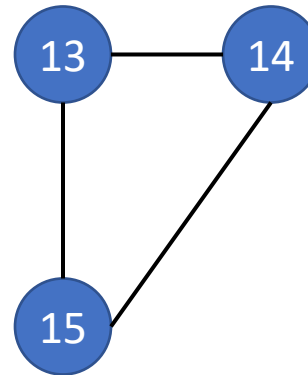
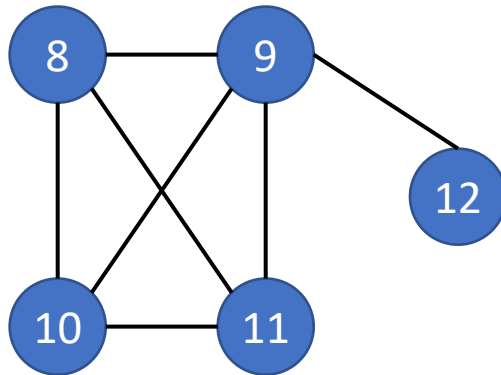
Output: $D = [0, 4, 1, 2, 3, 1, 3, 1, 2, 4]$



Connected Components Example



- Connected subgraphs of a graph that is not part of a larger connected subgraph.
- If u is in connected component 1 and v is in connected component 2, there NOT a possible path $u \rightarrow v$.



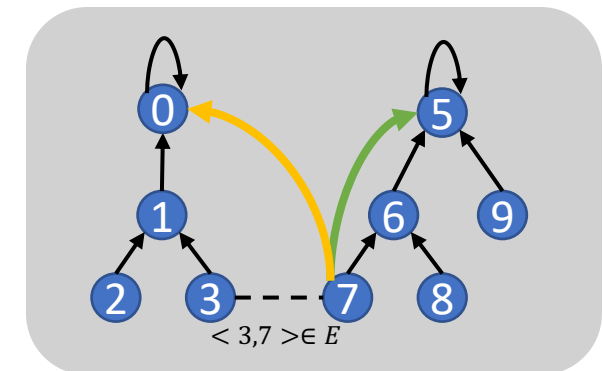
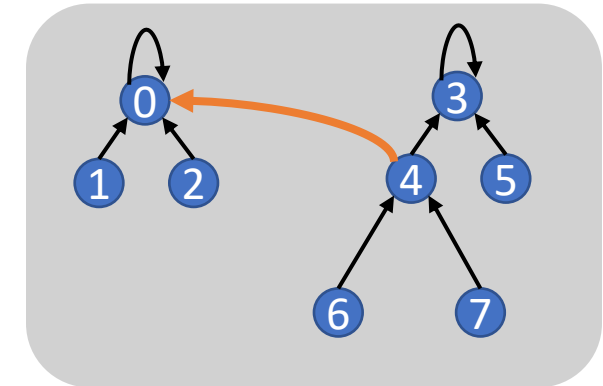
FastSV Connected Components

Algorithm 3.2 Parallel Shiloach-Vishkin connected components.

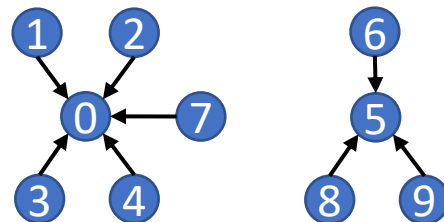
```

1: procedure FASTSV( $G$ )                                      $\triangleright$  A graph  $G$ .
2:   forall  $u \in V$  do
3:      $L[u] = u$ 
4:      $L_n[u] = u$ 
5:   end forall
6:   while  $L$  is changing do
7:     forall  $e = \langle u, v \rangle \in E$  do
8:        $L_n[L[u]] = \min(L_n[L[u]], L[L[v]])$             $\triangleright$  stochastic hooking
9:     end forall
10:    forall  $e = \langle u, v \rangle \in E$  do
11:       $L_n[u] = \min(L_n[u], L[L[v]])$                   $\triangleright$  aggressive hooking
12:    end forall
13:    forall  $u$  in  $V$  do
14:       $L_n[u] = \min(L_n[u], L[L[u]])$                   $\triangleright$  shortcutting
15:    end forall
16:  end while
17:  return  $L$ 
18: end procedure

```



Eventually L_n will make a graph that looks like this:



$L_n =$

0	0	0	0	0	5	5	0	5	5
---	---	---	---	---	---	---	---	---	---

Contour Connected Components

Why have a bunch of hooking procedures if we are just propagating labels?

Algorithm 3.3 Parallel fast-spreading voltage-based connected components.

```

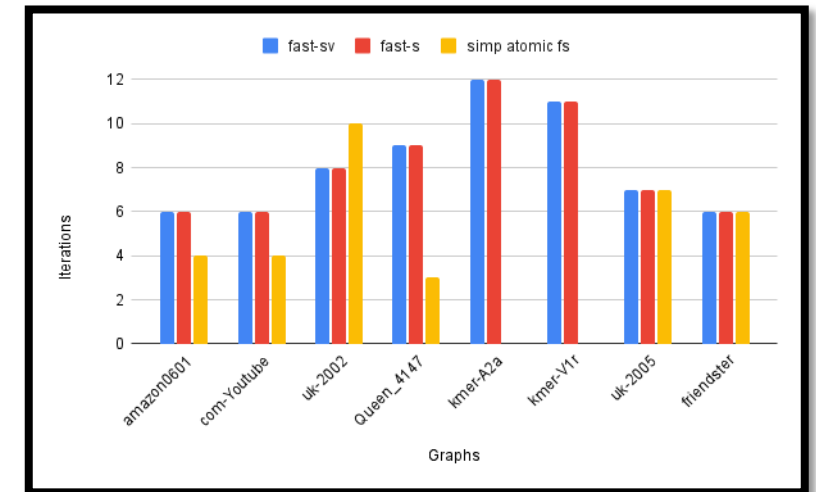
1: procedure FASTSPREADING( $G$ )                                     ▷ A graph  $G$ .
2:   forall  $i$  in  $0..n - 1$  do
3:      $L[i] = i$ 
4:      $L_n[i] = i$ 
5:   end forall
6:   while  $L$  is changing do
7:     forall  $e = \langle u, v \rangle \in E$  do
8:        $VO^2(L_n, L, u, v)$ 
9:     end forall
10:     $L = L_n$ 
11:  end while
12:  return  $L$ 
13: end procedure

```

$$VO^2(L_n, L, u, v) : \begin{bmatrix} L_n[u] \\ L_n[v] \\ L_n[L[u]] \\ L_n[L[v]] \end{bmatrix} \leftarrow z^2.$$

$$z^2 = \min(L[L[u]], L[L[v]])$$

Takeaway: Fast-spreading algorithm completes in $O(\log(d_{max}) + 1)$ iterations where d_{max} is the diameter for the largest connected component.



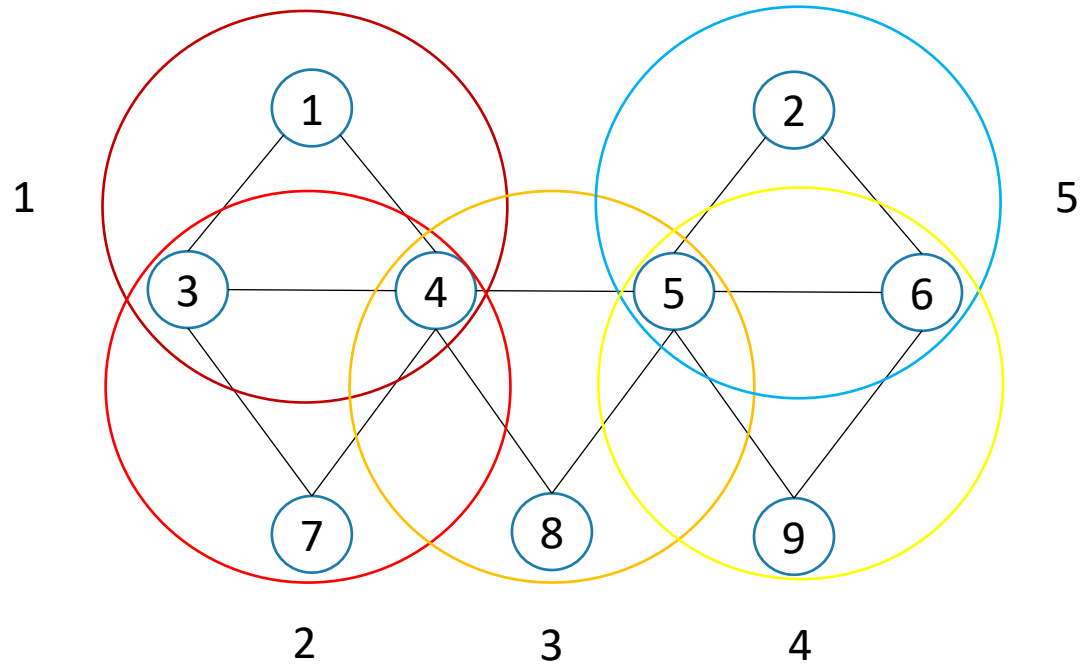
RED=BLUE for iterations

Basically, arriving to same Shiloach-Vishkin result but with less synchronizations.

Preliminary results show same number of iterations as FastSV and faster for some graphs.

Triangle Counting Example

Input:



Triangle_count = 5

Minimum Search Triangle Counting Algorithm

[Du, Patchett, Alvarado Rodriguez, Li, Bader, 2022]

Algorithm 3: *Parallel Minimum Search based Triangle Counting*

Input: A graph $G = (V, E)$.

Output: An integer value of the number of triangles.

```
1 coforall loc in Locales do
2   forall (edge  $e = u, v \in E$ ) && (e is local) do
3     // We assume that  $|Adj(u)| < |Adj(v)|$ 
4     var count:int=0;
5     forall  $w \in Adj(u)$  with (+ reduce Count) do
6       if ( $|Adj(w)| < |Adj(v)|$ ) then
7         if ( $v \in Adj(w)$ ) then
8           count ++;
9         end
10      end
11      else
12        if ( $w \in Adj(v)$ ) then
13          count ++;
14        end
15      end
16    end
17 end
18 return count
```

← get smallest adjacency list

← get closing triangle edge

If we assume $|Adj_u| < |Adj_v|$, and spawn threads for every $w \in |Adj_u|$, then the running time is:

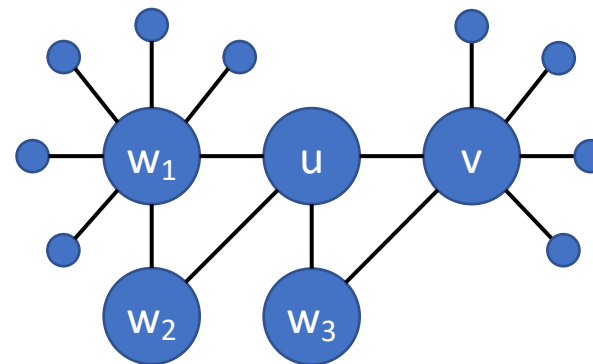
Minimum search: $\max_{w \in Adj_u} \log_2(\min(|Adj_w|, |Adj_v|))$

List Intersection: $\log_2(|Adj_v|)$

Minimum Search Triangle Counting Example

[Du, Patchett, Alvarado Rodriguez, Li, Bader, 2022]

1. Given an edge (u, v) we assume that $|Adj(u)| \leq |Adj(v)|$.
2. Then, for $\forall w \in Adj(u)$ we spawn $|Adj(u)| - 1$ parallel threads to check if we can form a complete triangle with (u, v, w) .
3. If $|Adj(w)| < |Adj(v)|$ we will check if $v \in Adj(w)$, else, we check if $w \in Adj(v)$.



Adj(x)	Value
Adj(u)	4
Adj(v)	6
Adj(w1)	7
Adj(w2)	2
Adj(w3)	2

Thread w_1 : search for w_1 in $Adj(v)$, no match, kill.

Thread w_2 : search for v in $Adj(w_2)$, no match, kill.

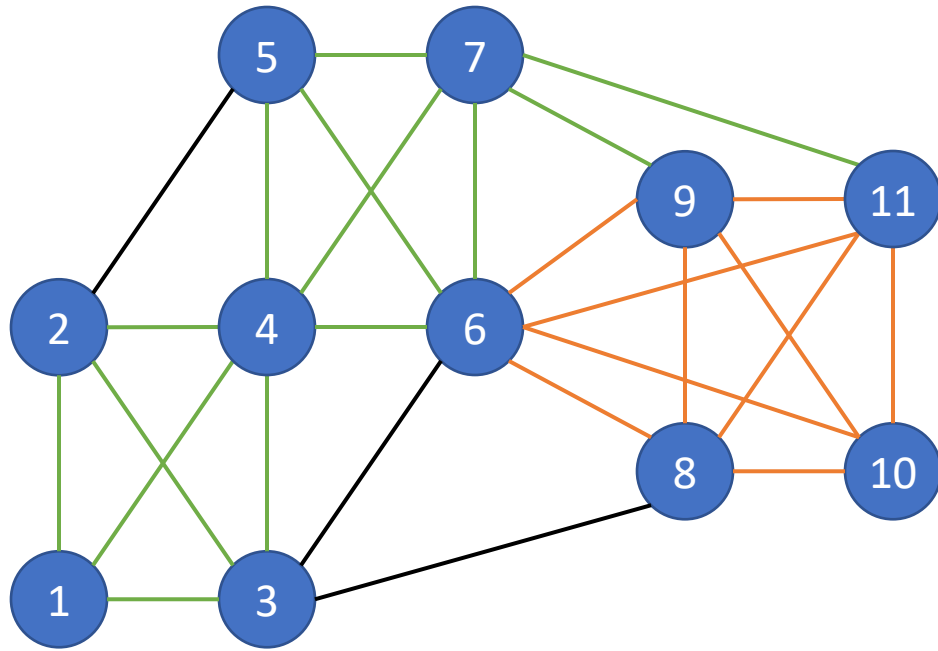
Thread w_3 : search for v in $Adj(w_3)$, match! Increment count.

Minimum Search Triangle Counting Operation Count Comparison

[Du, Patchett, Alvarado Rodriguez, Li, Bader, 2022]

- Assume $|Adj_u| < |Adj_v|$ and we spawn threads for every $w \in |Adj_u|$
 - Minimum search: $\max_{w \in Adj_u} \log_2(\min(|Adj_w|, |Adj_v|))$
 - List Intersection: $\log_2(|Adj_v|)$
 - Say we have the following information for our vertices:
 - $|Adj_u| = 4$ and $|Adj_v| = 1024$
 - For every w in Adj_u , $|Adj_w| \leq 8$
- **List intersection:** 4 threads amounting to $\lceil \log_2 1024 \rceil = 10$ operations each.
 - **Minimum search:** 4 threads amounting to $\lceil \log_2 8 \rceil = 3$ operations each.

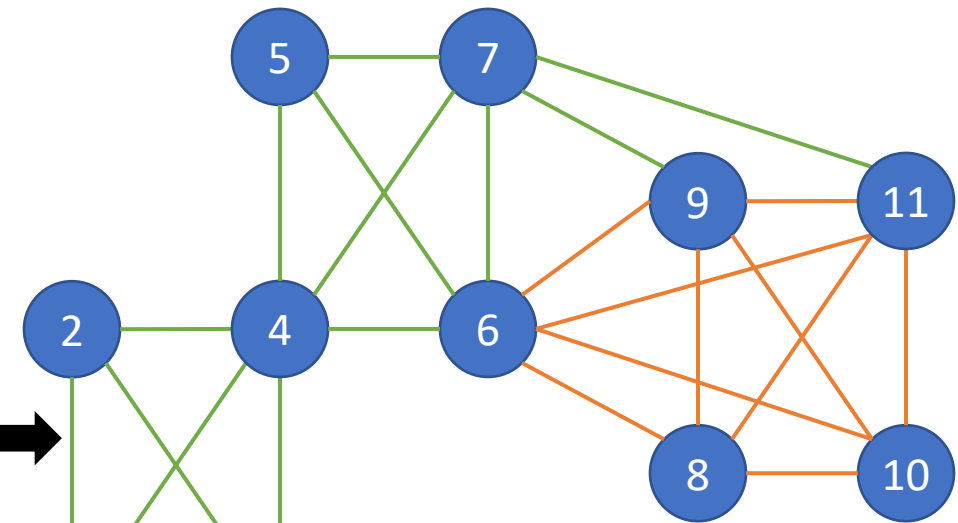
Truss Analytics Example



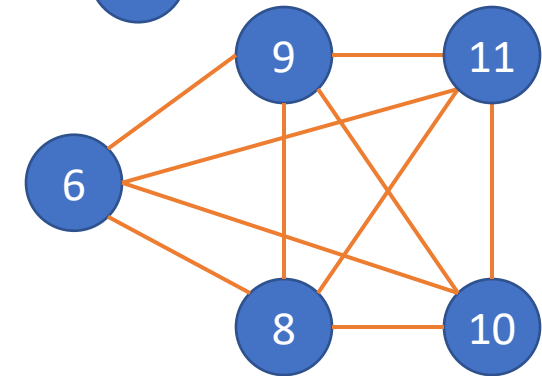
3-truss —————
4-truss —————
5-truss ————— = max-truss

coloring = truss decomposition

k=4 truss
→



k=5 = max-truss
→



truss decomposition
→

[3, 3, 3,, 4, 5, 5, 5, 5, 5, 5, 5, 5, 5]

edge id with trussness

Every edge is part of at least k-2 triangles.

Algorithm 4: Optimized k -Truss Parallel Algorithm

```

1  OptKTruss( $G, k$ )
2  var JustDelEBag = new DistBag(int, Locales);
3  AtoSup  $\leftarrow 0$  and EDel  $\leftarrow -1$ 
4  TriangleCounting( $G, EDel, AtoSup$ )
5  coforall loc in Locales do
6    forall ( $e \in E$ ) && ( $e$  is local) do
7      if (EDel[ $e$ ] == -1) && (AtoSup[ $e$ ].read() <  $k-2$ )
8        then
9          EDel[ $e$ ] =  $1-k$ 
10         JustDelEBag.add( $e$ );
11      end
12    end
13  while (JustDelEBag.getSize() > 0) do
14    SupportUpdate( $G, EDel, JustDelEBag, AtoSup$ )
15    coforall loc in Locales do
16      forall ( $e \in JustDelEBag$ ) && ( $e$  is local) do
17        if (EDel[ $e$ ] ==  $1-k$ ) then
18          EDel[ $e$ ] =  $k-1$ 
19        end
20      end
21      JustDelEBag.clear();
22      forall ( $e \in E$ ) && ( $e$  is local) do
23        if (EDel[ $e$ ] == -1) &&
24          (AtoSup[ $e$ ].read() <  $k-2$ ) then
25          EDel[ $e$ ] =  $1-k$ 
26          JustDelEBag.add( $e$ );
27        end
28      end
29    end
30  return EDel

```

distribution and edge parallelism

generate new support
only with undeleted edges

keep removing edges with
support less than $k-2$ until
all proper supports have
been updated

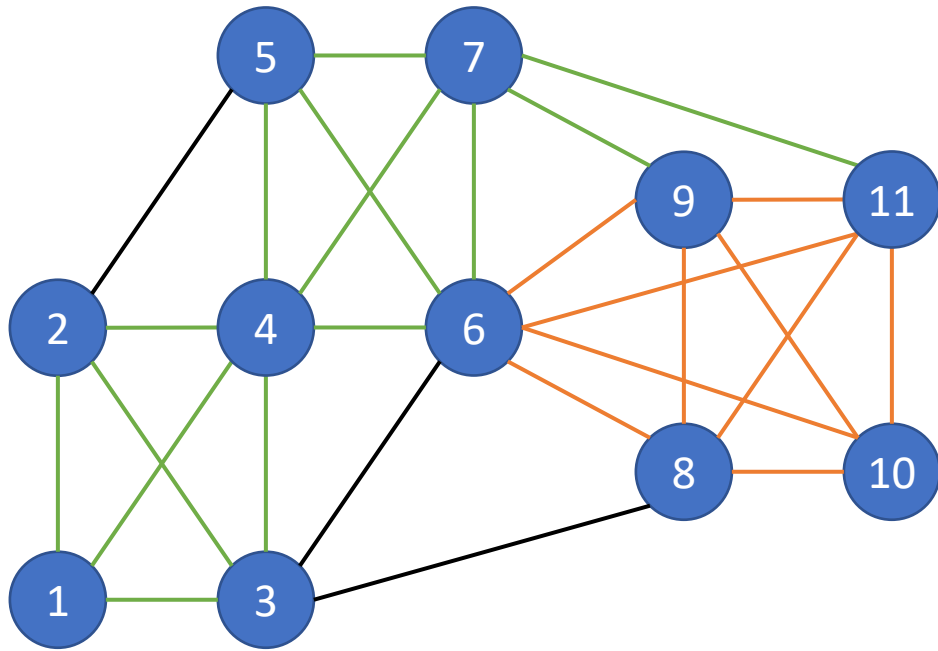
Algorithm 2: Minimum Search based Support Updating

```

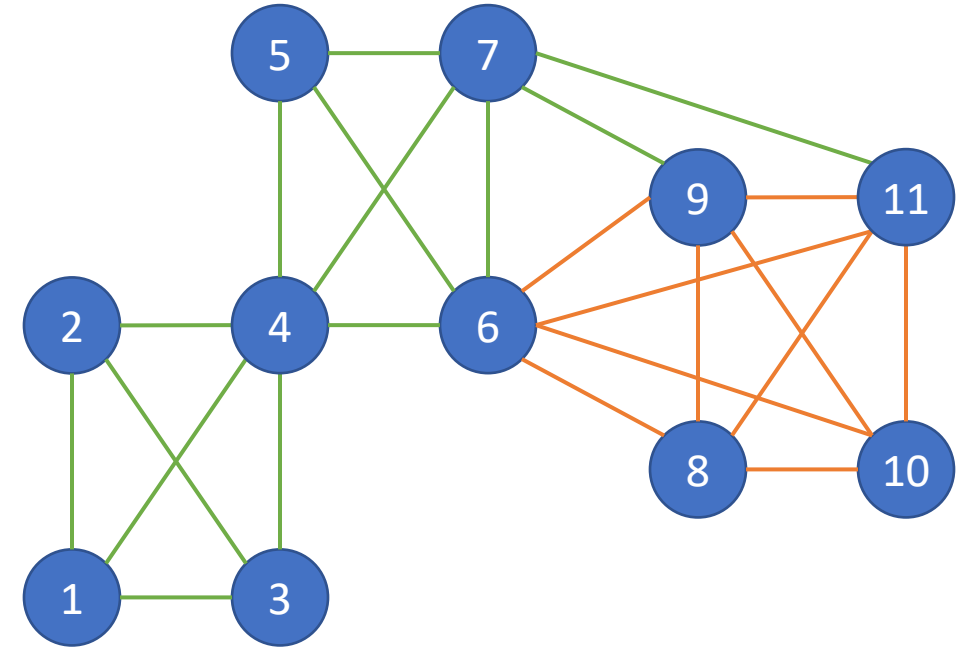
1  SupportUpdate( $G, EDel, JustDelEBag, AtoSup$ )
2  coforall loc in Locales do
3    forall ( $e_1 = \langle u, v \rangle \in JustDelEBag$ ) && ( $e$  is local)
4      do
5        /* We assume that  $|Adj(u)| < |Adj(v)|$  */
6        forall ( $e_2 = \langle u, w \rangle, w \in Adj(u) - \{v\}$ ) &&
7          (EDel[ $e_2$ ] < 0) do
8          Search  $e_3 = \langle v, w \rangle$  or  $e_3 = \langle w, v \rangle$ ;
9          if ( $e_3$  exists) then
10            if ( $e_2$  and  $e_3$  are undeleted) then
11              AtoSup[ $e_2$ ].sub(1);
12              AtoSup[ $e_3$ ].sub(1);
13            end
14            else
15              if ( $e_2$  is undeleted) && ( $e_1 < e_3$ ) then
16                AtoSup[ $e_2$ ].sub(1);
17              end
18              if ( $e_3$  is undeleted) && ( $e_1 < e_2$ ) then
19                AtoSup[ $e_3$ ].sub(1);
20              end
21            end
22          end
23        end
24      end
25    end
26  end

```

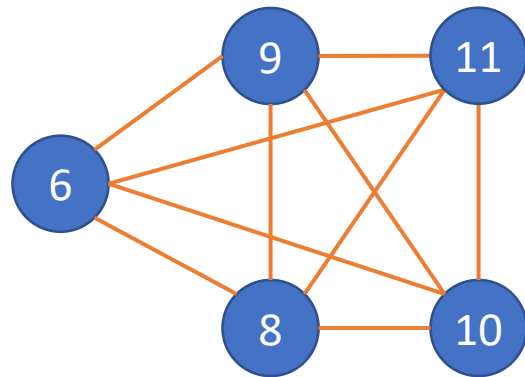
fun fact: beginning iteration (+ reduce AtoSup)/3 = tricount



Delete black edges.
Update supports.



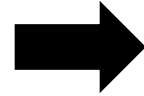
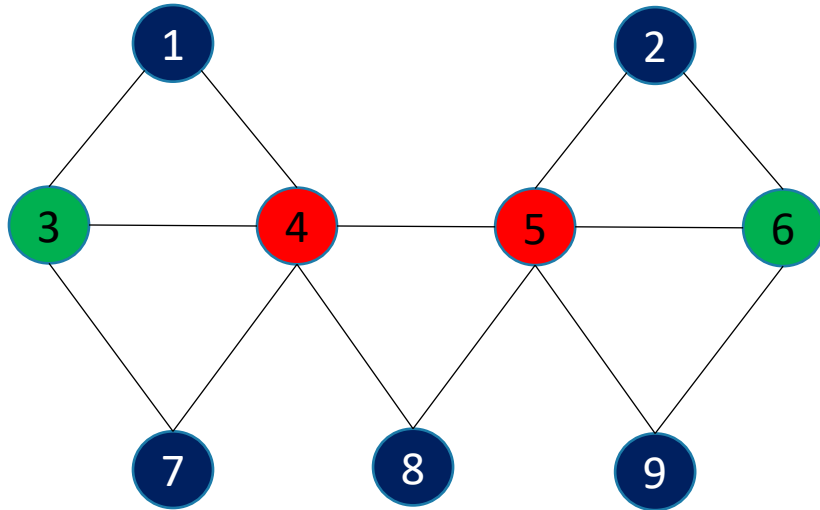
Delete green edges.
Update supports.



Output

Triangle Centrality Example

Input:



$$TC(v) = \frac{\frac{1}{3} \sum_{u \in N_{\Delta}^+(v)} \Delta(u) + \sum_{w \in \{N(v) \setminus N_{\Delta}(v)\}} \Delta(w)}{\Delta(G)}.$$

[Burkhardt, 2021]

Output: D = [0.4, 0.4, 0.47, 0.73, 0.73, 0.47, 0.4, 0.4, 0.4]

What kind(s) of questions can each algorithm answer?

- **Breadth-first search**
 - How far is vertex u from a vertex v ? Is a vertex v reachable from a vertex u ?
 - What is the diameter of a graph?
 - How are the nodes distributed amongst depths based off a source vertex?
- **Connected components**
 - What is the largest connected component in a graph?
 - What is the distribution of the sizes of all the connected components?
 - Is the graph connected?
- **Triangle counting**
 - Most obvious: how many triangles (3-cliques) in the graph?
 - How many triangles are adjacent to a vertex u ?
- **Truss analytics**
 - Is there a subgraph where every edge belongs to $[2,3,+]$ triangles?
 - What is the maximum number of triangles every edge belongs to?
- **Triangle centrality**
 - How important is a vertex u dependent on how many triangles are around it?
 - How does this match with other metrics such as betweenness centrality?

Jupyter Notebook Time 😊

Thank You 😊
Questions?