



Enumération Randomisée des Triangles dans des Graph à Grande Echelle à base de SQL

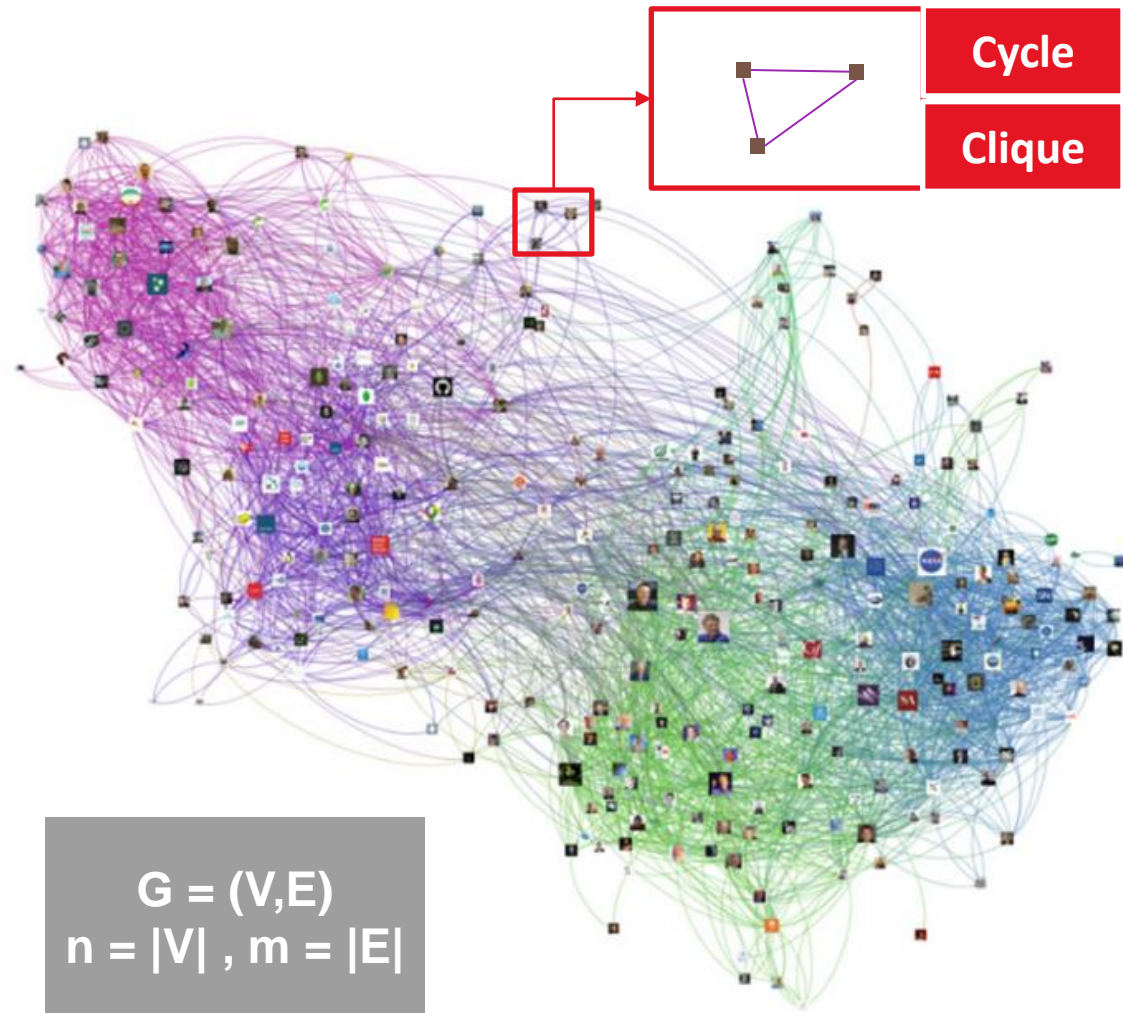
(Randomized Triangle Enumeration on Large Graph based on SQL)

Abir Farouzi^{1,3}, Ladjel Bellatreche¹, Carlos Ordonez², Gopal Pandurangan², Mimoun Malki³

¹ **LIAS/ISAE-ENSMA**, Poitiers, France
abir.farouzi@ensma.fr, bellatreche@ensma.fr

² Université de Houston, Texas, USA
carlos@central.uh.edu , gopal@cs.uh.edu

³ **ESI-SBA**, Sidi Bel Abbès, Algérie
a.farouzi@esi-sba.dz , m.malki@esi-sba.dz

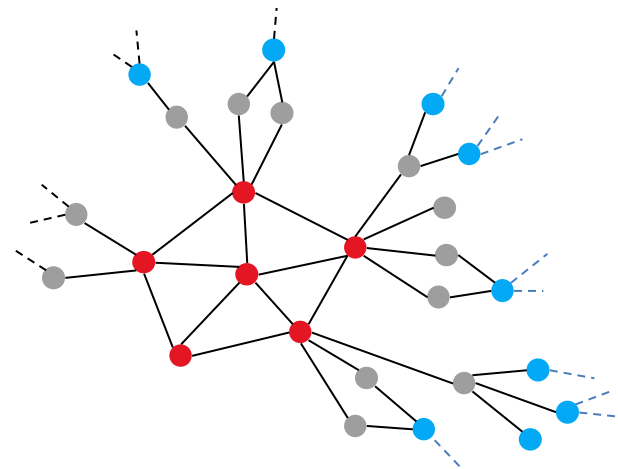


$G = (V, E)$
 $n = |V|, m = |E|$

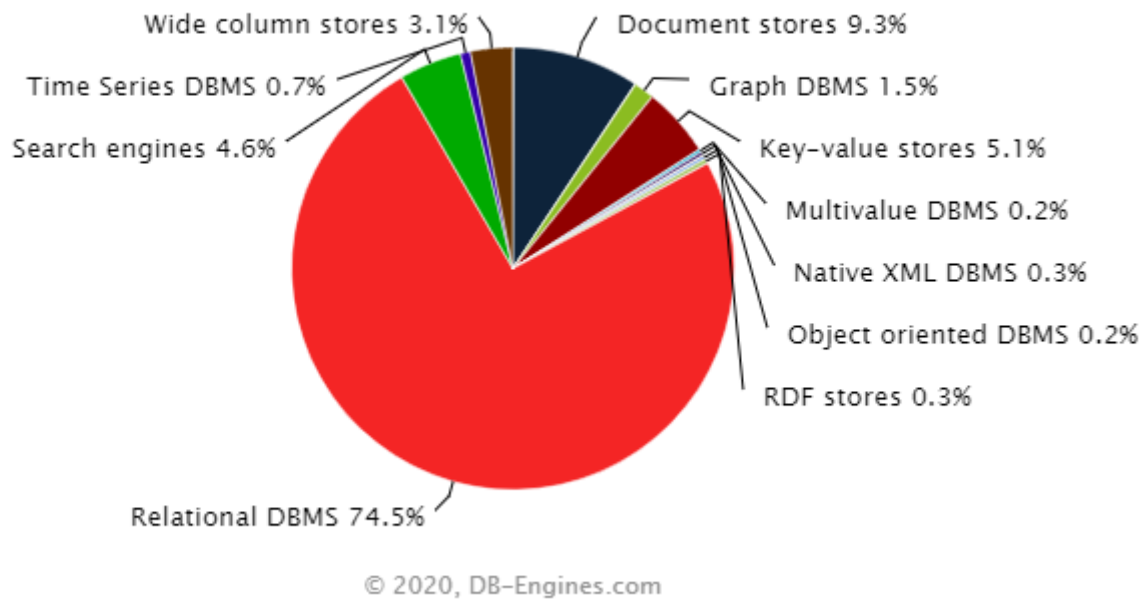
- Path Analytics
- Connectivity Analytics
- Community Analytics
- Centrality Analytics

- Social process analysis
- Dense sub-graph extraction
- Database joins
- Anomalies detection**
- Community age detection
- Recommandations

Spam filtering



- **Attacking node**
- **Victim node**
- **Safe node**



- Lot of data are stored in relational DBMS
- Less data stored as Graphs
- Export data from relational DBMS to graph DBMS causes performance issues
- Exploit RDBMS capabilities like security, concurrence management..etc

Parallel Randomized
Algorithm based on SQL
queries

Graph Analysis Inside
DBMS

Perfect load Balancing

Scalability

Motivations

Why triangles?

Why relational databases

Solution

Standard Algorithm

Randomized Algorithm

Experimental study

Configuration and Data sets

Results

Conclusion

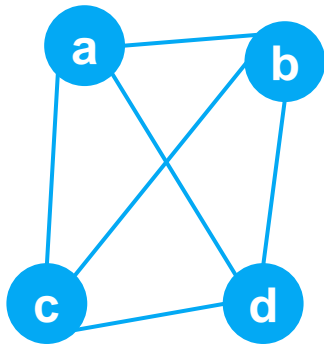
Conclusions

Perspectives



Solution

#standard_algorithm #randomized_algorithm #sql_queries



Undirected unweighted Graph

		j			
		a	b	c	d
i	a	0	1	1	1
	b	1	0	1	1
	c	1	1	0	1
	d	1	1	1	0

Adjacency matrix

i	j
a	b
a	c
a	d
b	a
b	c
b	d
c	a
c	b
c	d
d	a
d	b
d	c

Edge Table E
Adjacency list

Algorithm 1: Algorithme standard

```

Result: triangles (liste de tous les triangles dans le graphe)
chargement du graphe;
for  $e_1 \leftarrow 1$  to  $m$  do
  for  $e_2 \leftarrow 1$  to  $m$  do
    for  $e_3 \leftarrow 1$  to  $m$  do
      if  $(e_1.j == e_2.i \wedge e_2.j == e_3.i \wedge e_3.j == e_1.i)$  then
        else
          |  $triangles \leftarrow (e_1, e_2, e_3)$ 
        end
      end
    end
  end
end

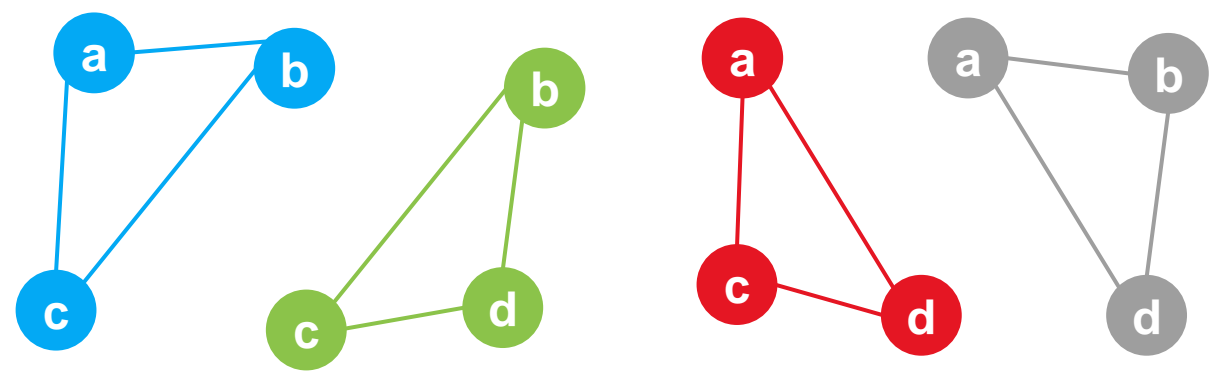
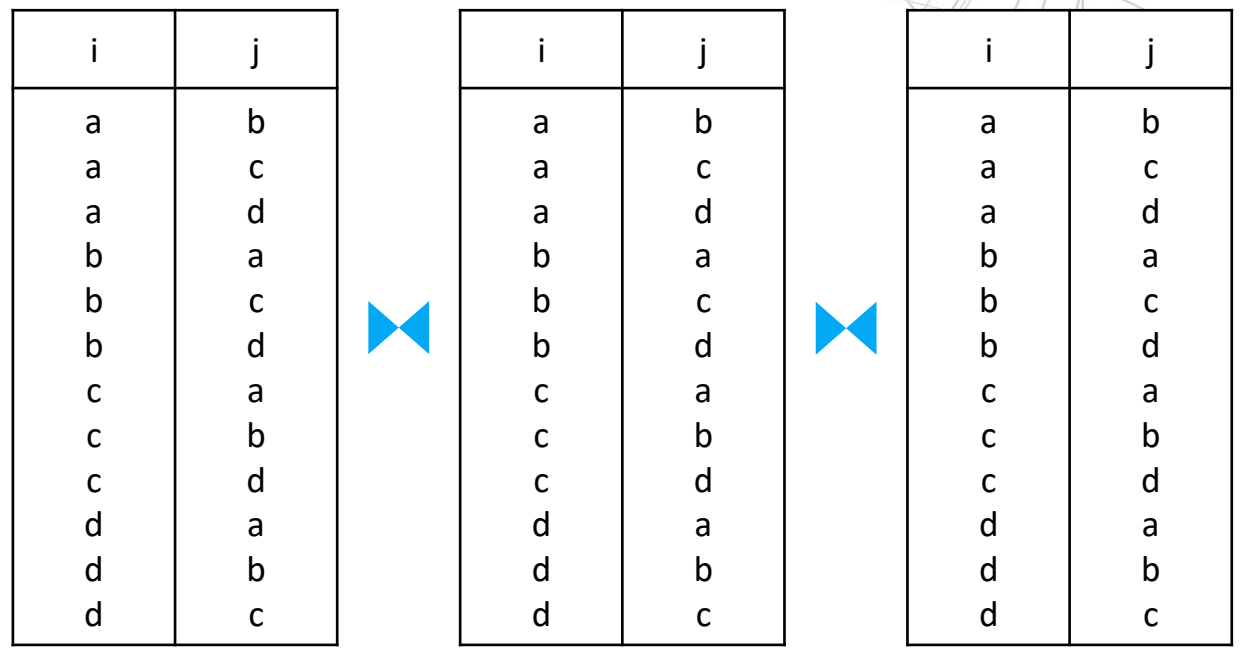
```

Complexity is $O(n^3)$

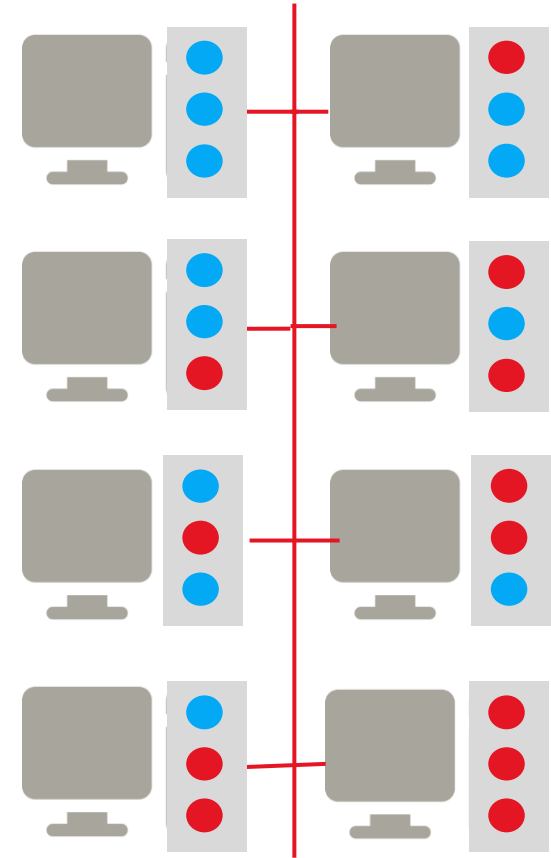
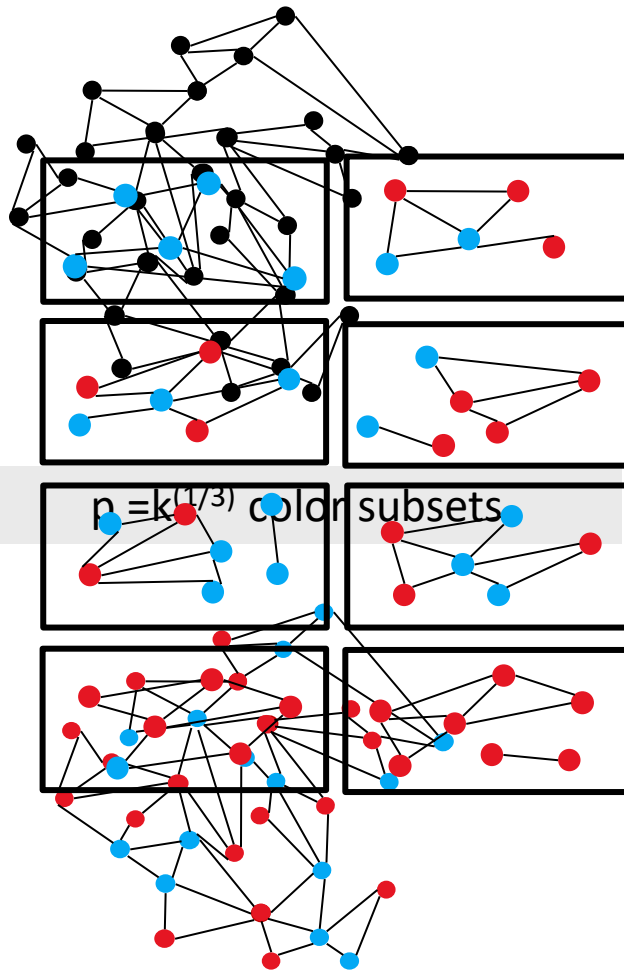
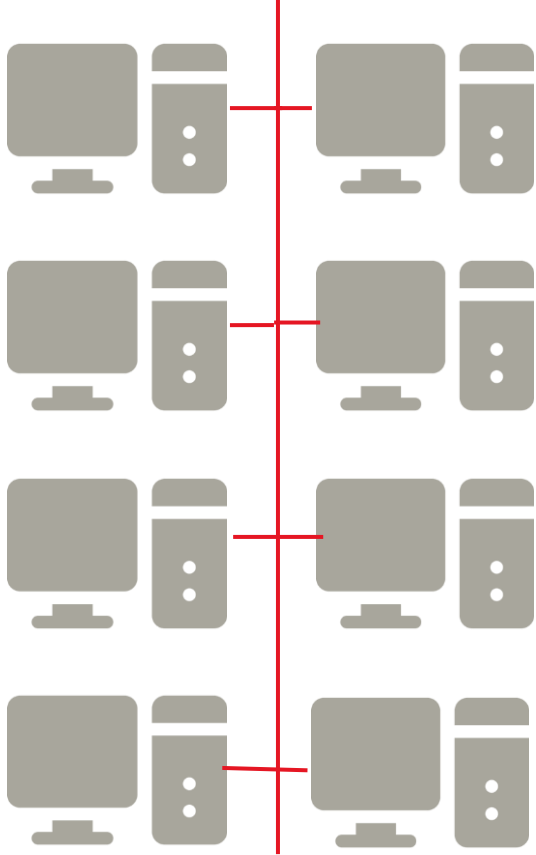
```

SELECT  $E_1.i$  AS  $v_1$ ,  $E_1.j$  AS  $v_2$ ,  $E_2.j$  AS  $v_3$ 
FROM
   $E$   $E_1$  JOIN  $E\_dup$   $E_2$  ON  $E_1.j = E_2.i$ 
  JOIN  $E$   $E_3$  ON  $E_2.j = E_3.i$  AND  $E_3.j = E_1.i$ 
WHERE  $E_1.i < E_1.j$  AND  $E_2.i < E_2.j$ ;

```

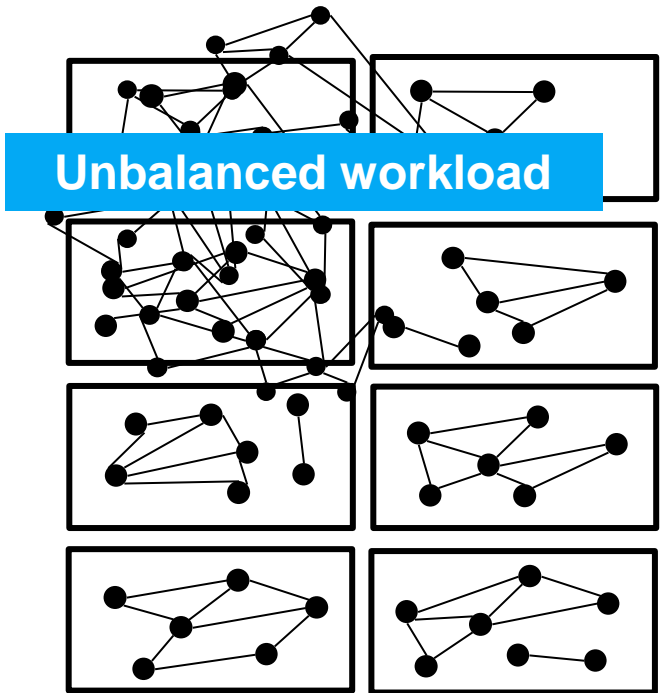


Complexity is $\tilde{O}(\max\{m/k^{2/3}, n/k^{1/3}\})$

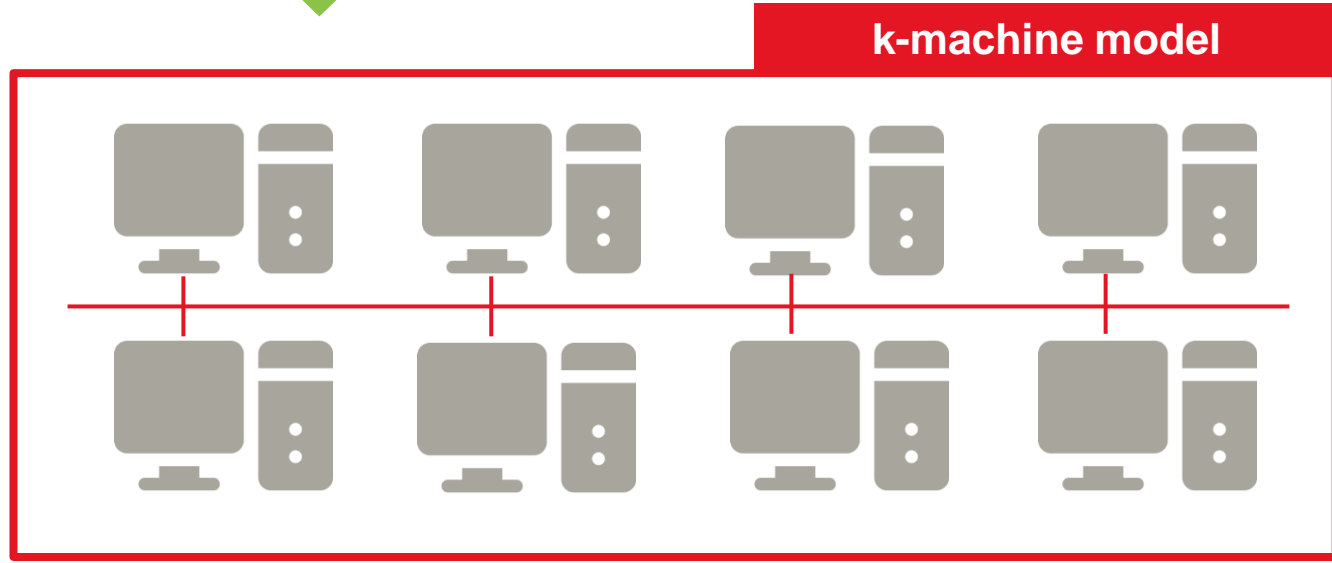
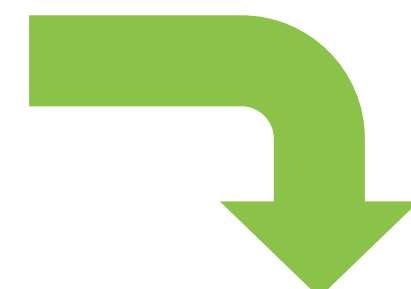


Randomized Algorithm

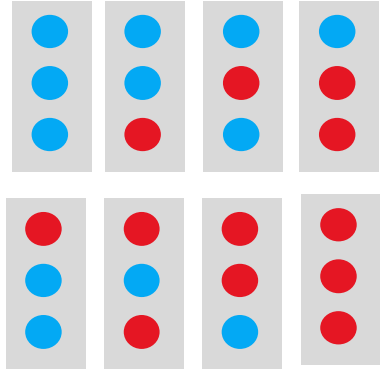
Graph loading



```
CREATE TABLE E_s (i int, j int);
COPY E_s FROM "Lien/vers/graph_data_set";
/*Si le graphe est non orienté */
INSERT INTO E_s SELECT j,i FROM E_s;
```

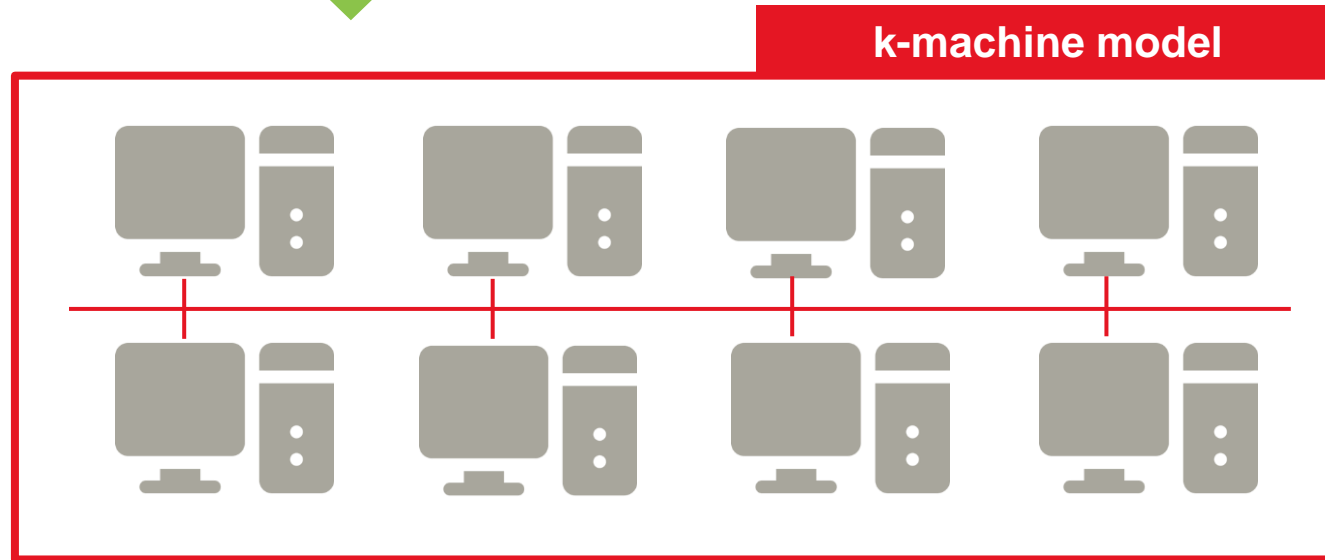


Color triplet assignment to machines



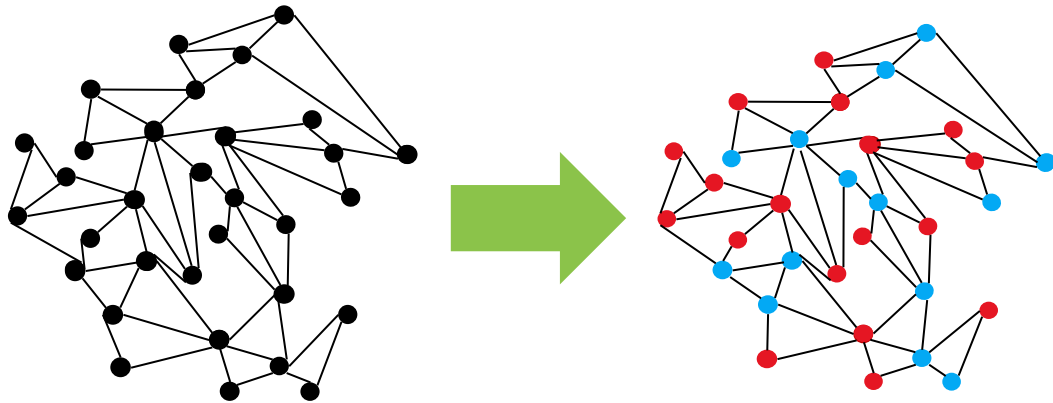
```
CREATE TABLE Triplet(machine int,color1 int,color2 int,color3 int)
UNSEGMENTED ALL NODES;
COPY Triplet FROM "lien/vers/triplet_file";
```

- 1,1,1,1
- 2,1,1,2
- 3,1,2,1
- 4,1,2,2
- 5,2,1,1
- 6,2,1,2
- 7,2,2,1
- 8,2,2,2



Send edges to proxies (1)

$p = k^{(1/3)}$ color subsets

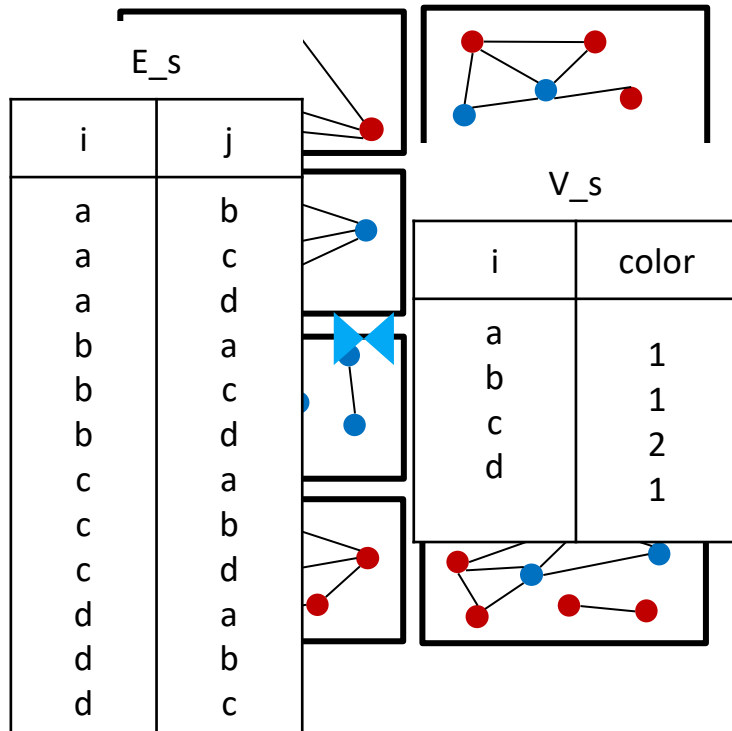


```

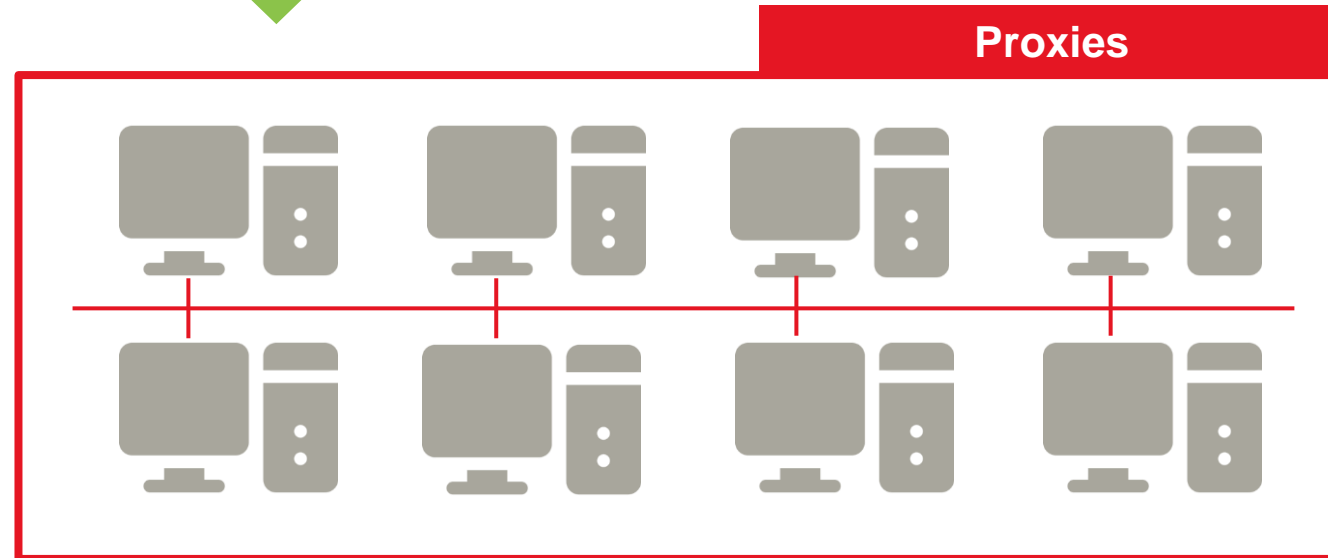
CREATE TABLE V_s(i int, color int);
INSERT INTO V_s
  SELECT i, randomint(p)+1
  FROM
    (SELECT DISTINCT i FROM E_s
    UNION
    SELECT DISTINCT j FROM E_s
    )V;

```

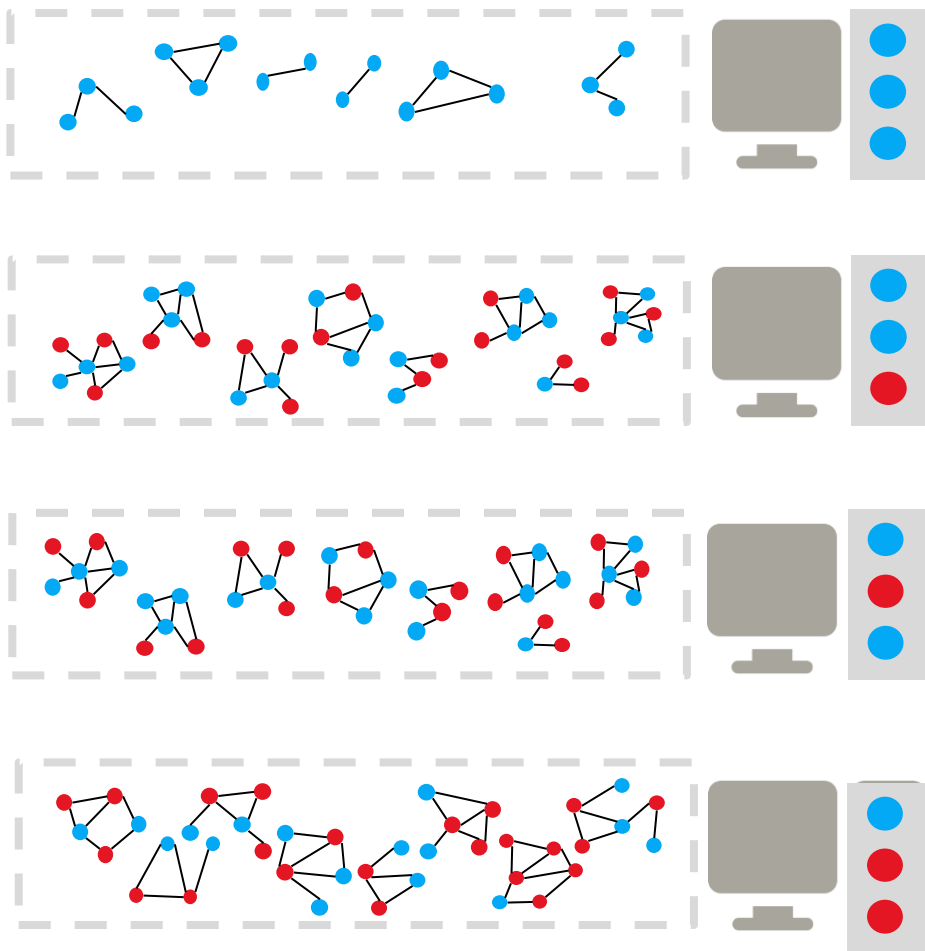
Send edges to proxies (2)



```
CREATE TABLE E_s_proxy (i_color int, j_color int, i int, j int);
INSERT INTO E_s_proxy
SELECT Vi.color, Vj.color, E.i, E.j
FROM
E_s E JOIN V_s Vi ON E.i = Vi.i
JOIN V_s Vj ON E.j = Vj.i;
```



Edges collecting by local machines



```

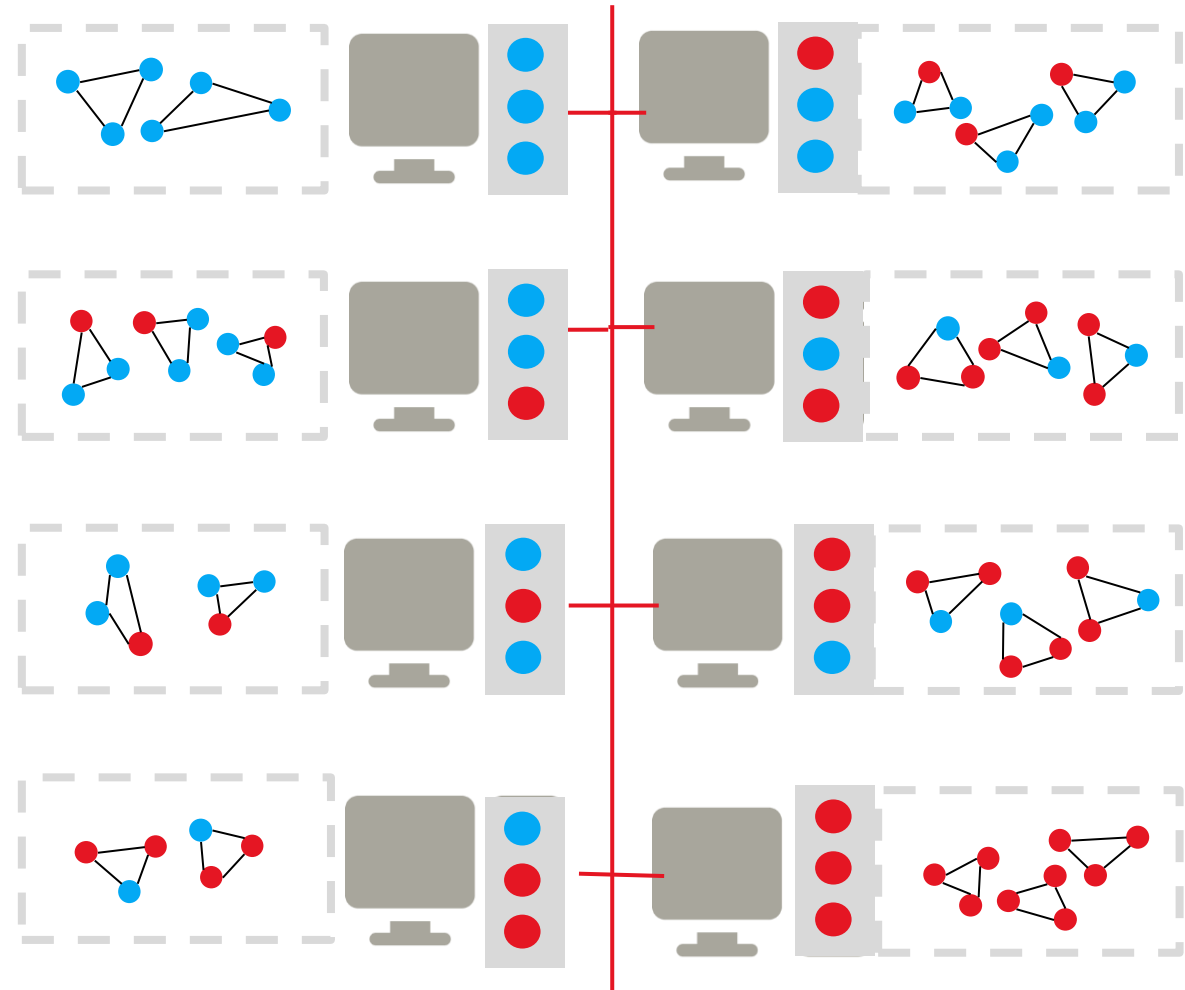
CREATE TABLE E_s_local(machine int,i int,j int,i_color int,j_color int);
INSERT INTO E_s_local
SELECT machine, i, j, i_color, j_color
FROM
    E_s_proxy E JOIN Triplet edge1 ON E.i_color=edge1.color1
    AND E.j_color=edge1.color2 WHERE E.i<E.j
UNION
SELECT machine, i, j, i_color, j_color
FROM
    E_s_proxy E JOIN Triplet edge2 ON E.i_color=edge2.color2
    AND E.j_color=edge2.color3 WHERE E.i<E.j
UNION
SELECT machine, i, j, i_color, j_color
FROM
    E_s_proxy E JOIN Triplet edge3 ON E.i_color=edge3.color3
    AND E.j_color=edge3.color1 WHERE E.i>E.j;
    
```

Local triangle enumeration

Balanced load

```

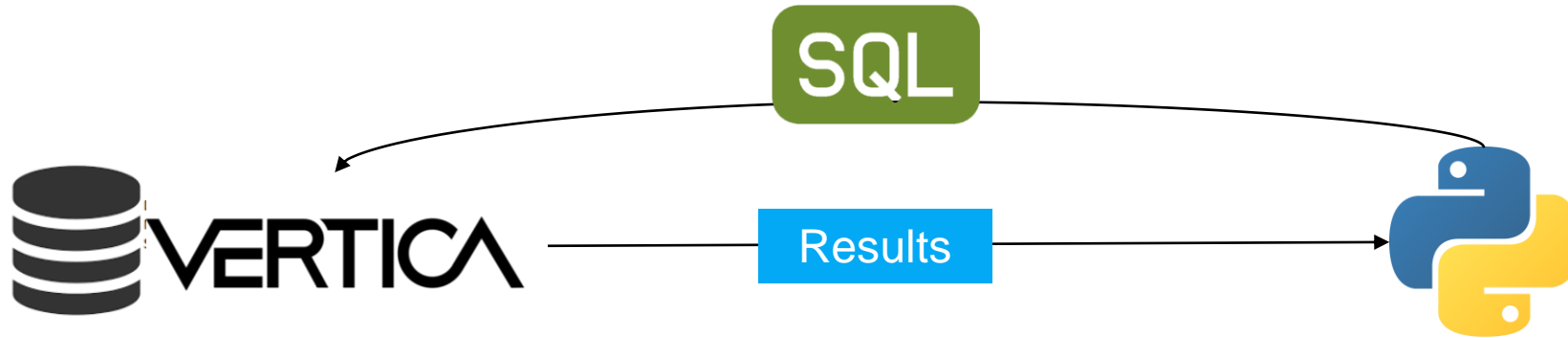
SELECT E1.machine, E1.i AS v1, E1.j AS v2, E2.j AS v3
FROM
  E_s_local E1 JOIN E_s_local E2
    ON E1.machine=E2.machine AND E1.j=E2.i
  JOIN E_s_local E3 ON E2.machine=E3.machine
    AND E2.j=E3.i
  JOIN Triplet T on T.machine = E3.machine
WHERE E1.i<E1.j AND E2.i<E2.j AND E1.i=E3.j
  AND E1.i_color=T.color_1 AND E1.j_color=T.color_2
  AND E2.j_color=T.color_3
  AND local_node_name()='node_name'
ORDER BY v1,v2,v3;
  
```





Experimental study

#data_sets #cluster #vertica_dbms



k-machine model

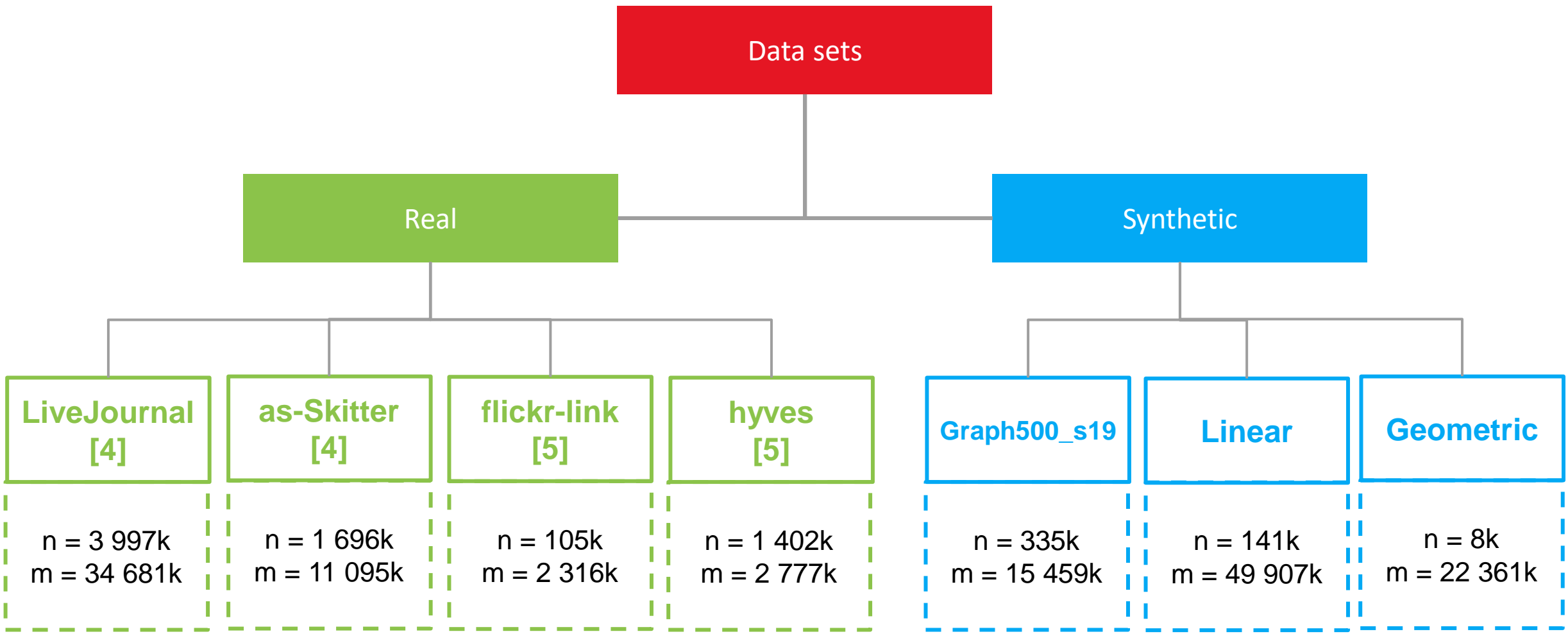
Operating System: **Ubuntu server 18.04**

RAM: **48 Go/host**

Storage: **1 To/ host**

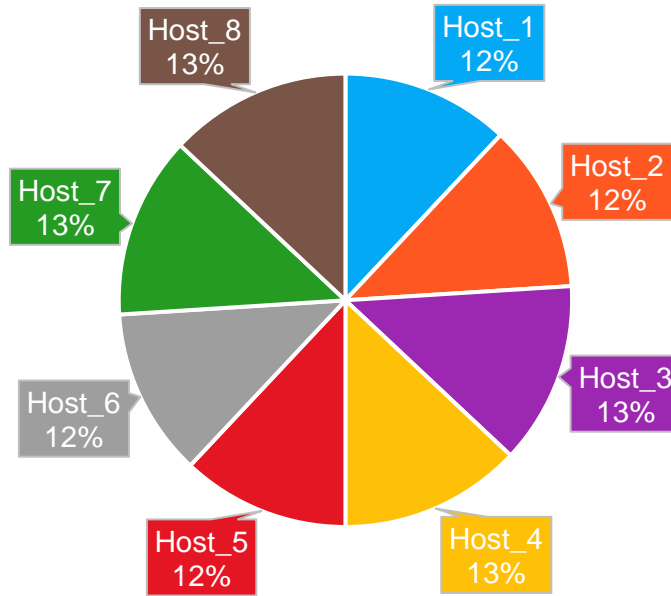
CPU: **Intel 4 cores/host**

Architecture: **Share-Nothing**

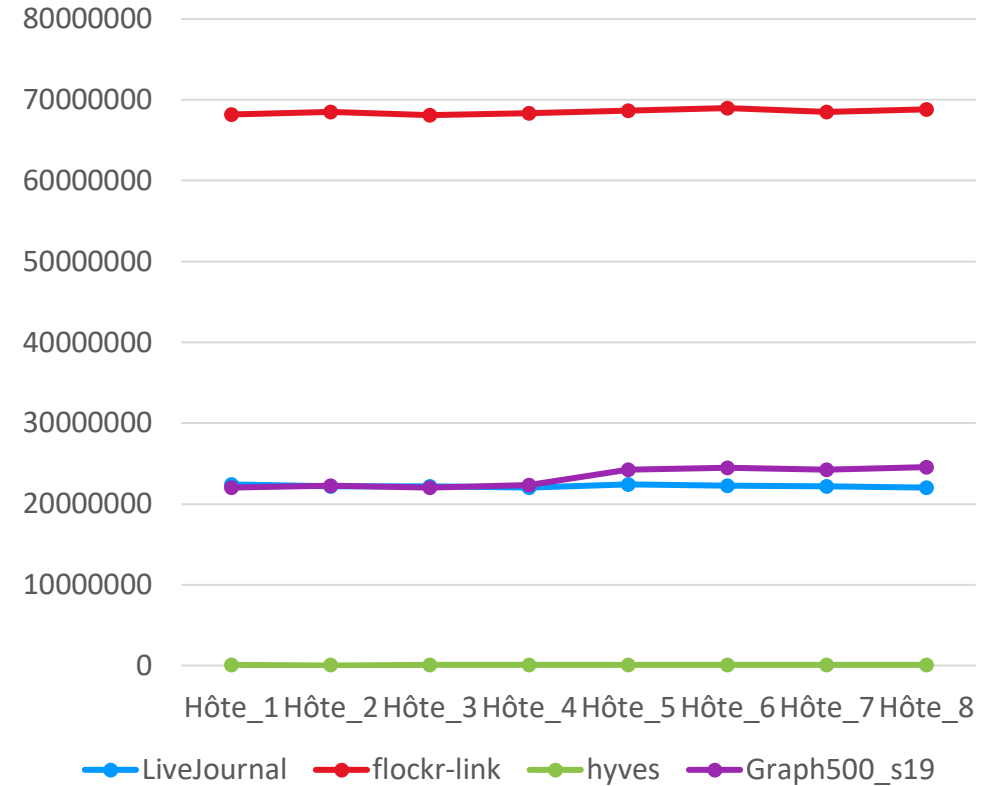


Randomized triangle enumeration by machine

Triangle count by machine
(Geometric data set)

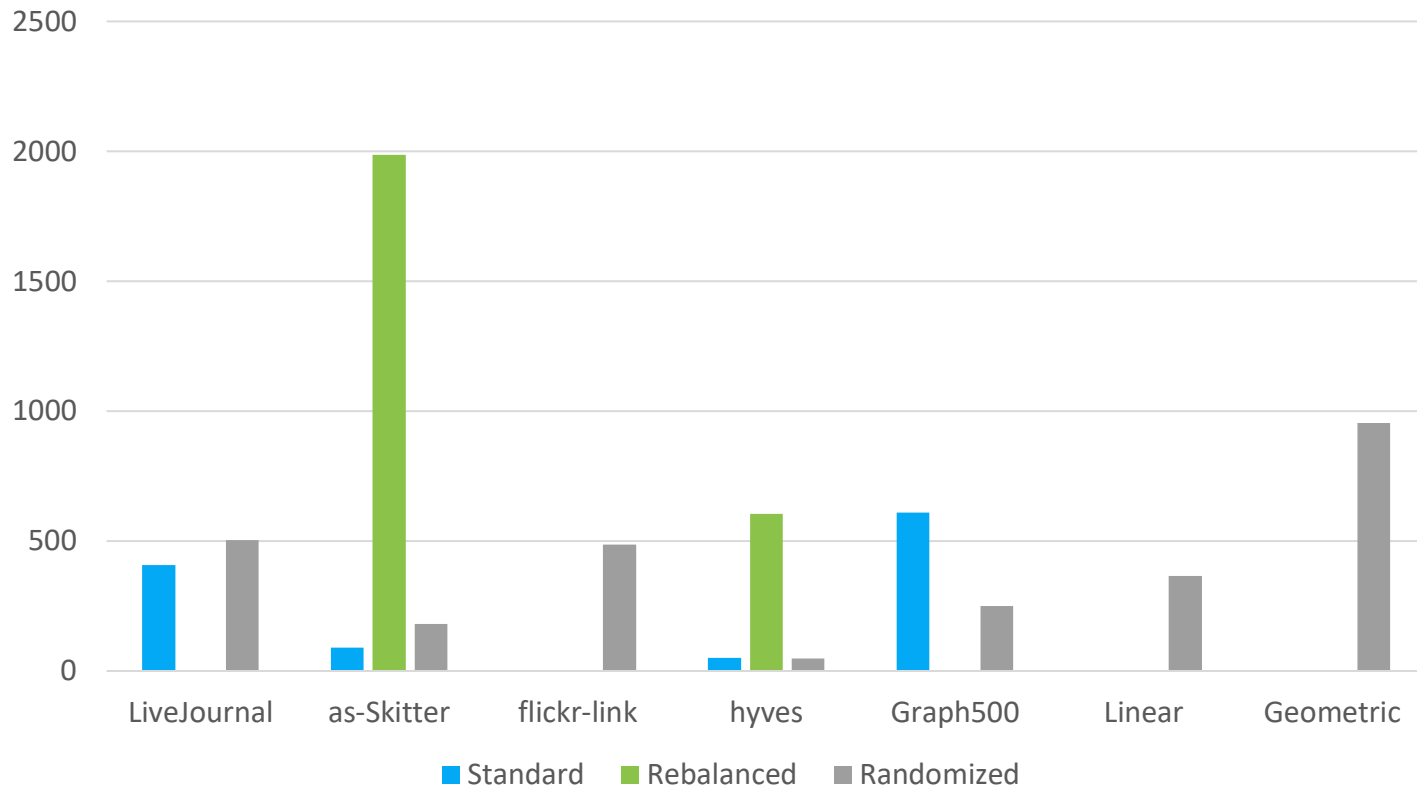


Triangle count by machine



Standard algorithm and randomized algorithm comparison

Execution time by data set



- Standard algorithm queries are **less efficient** with large-scale graphs and cause **performance problems when scaling**.
- Rebalancing is **expensive** and often **causes memory issues**.
- Randomized algorithm queries **perform better** with graphs having **skewed data distribution**.



Conclusions

#conclusion #perspectives

Parallel Randomized Algorithm for Triangles Enumeration on Large Graphs Using SQL Queries

Randomized algorithm for triangles enumeration ensures load balancing

The algorithm queries were able to process very large graphs and therefore they ensure scalability

- ◆ Deep study of Randomized Algorithms
- ◆ Comparaison with graph analysis engines
- ◆ Clique detection using triangles

- [1] https://db-engines.com/en/ranking_categories
- [2] Pandurangan, G., P. Robinson, et M. Scquizzato (2018). On the distributed complexity of large-scale graph computations. In SPAA, pp. 405–414.
- [3] Farouzi, A., L. Bellatreche, C. Ordonez, G. Pandurangan, et M. Malki (2020). A scalable randomized algorithm for triangle enumeration on graph based on SQL queries. In To Appear in DaWaK Conference.
- [4] Leskovec, J. et A. Krevl (2014). SNAP Datasets : Stanford large network dataset collection. <http://snap.stanford.edu/data>
- [5] Kunegis, J. (2013). Konekt : The koblenz network collection. Association for Computing Machinery.



Thank you for your attention

Stay safe!