

Scalable Interactive Dynamic Graph Clustering on Multicore CPUs

SON et. al., IEEE Trans.KDE 2018

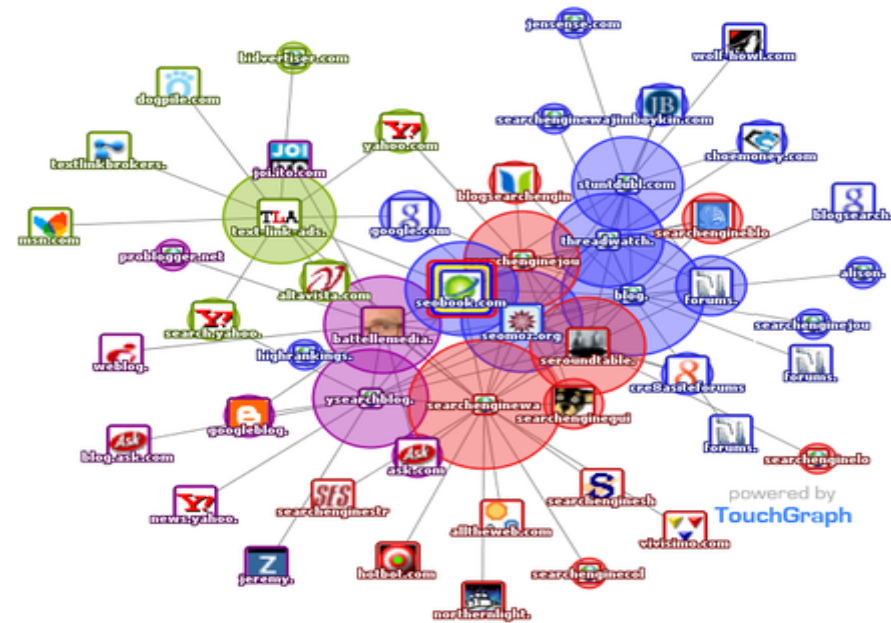
Wissem Inoubli
FST, LIPAH

Graphs

Graphs are ubiquitous and can model complex relationships



Social network



Web

Graph clustering

- Group vertices into clusters;
 - e.g: detect the strong connected sub-graphs;
- Optimization problem:
 - **Maximize** the distance between **sub-graphs**;
 - **Minimize** the distance between the vertices in the same **sub-graphe**.

Applications

- **Digital marketing:** deliver intelligently the messages or offers;
- **Bio-informatic:** identify the target proteins or the functions of the protein groups;
- **Social network:** community detection, user profiling;
- **Scientific Research :** identify the scientific community .

challenges

- **High volume of the graph:** Graph storage;
- **Scalability:** Graph processing;
- **Velocity:** Fast updating of graph.



1.49 trillion user/ month



500 Million Tweets/ dy
320 Million user/month

Scalable Interactive Dynamic Graph Clustering on Multicore CPUs

SON et. al., IEEE Trans. KDE 2018

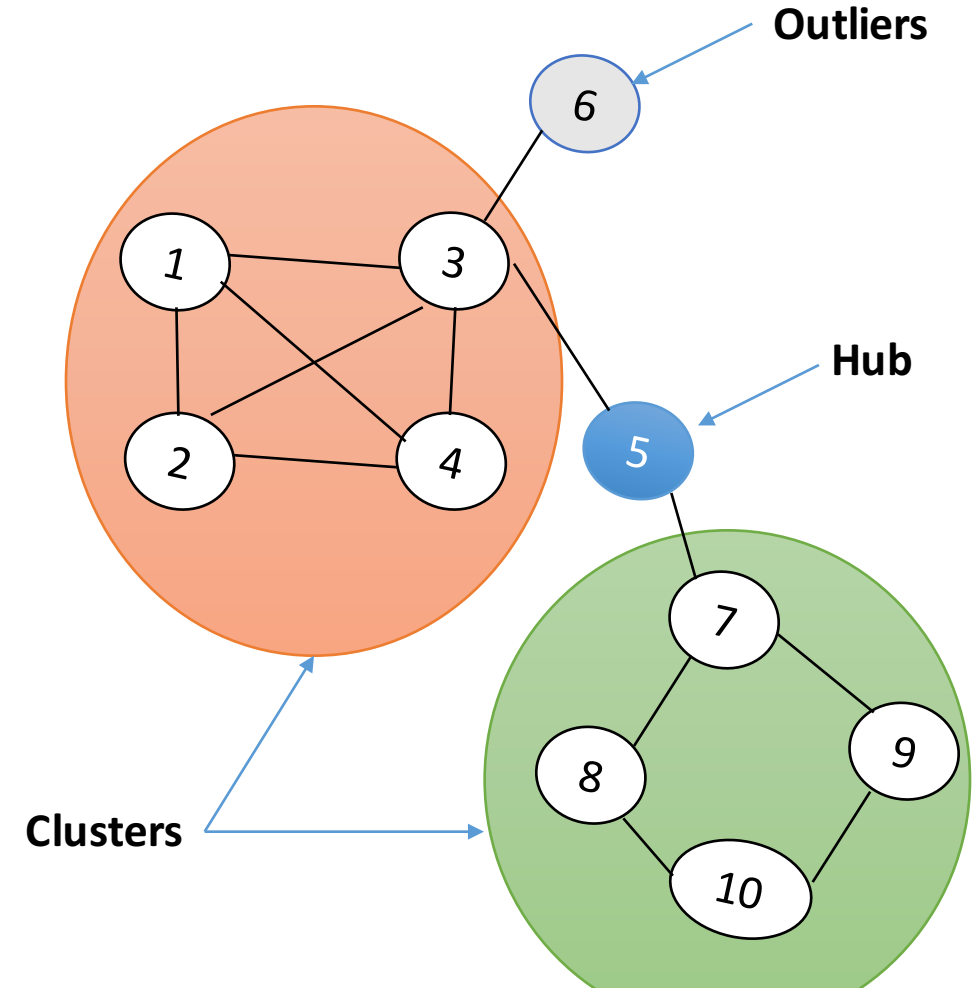
Contributions

1. Structural clustering: SCAN;
2. Interactive clustering in real time: intermediate results;
3. Parallel processing on multi-core and shared memory;
4. Dynamic graph clustering.

Scan: a structural clustering algorithm for networks

Xu et. al., KDD'07.

- Based on common neighbors between the vertices;
- Uses the structural similarity;
- Identifies Clusters, Outliers and Hubs;



The main steps of the basis SCAN algorithm

1. Define the neighbors of each vertex;
2. Calculate the *structural similarity* in the graph \mathbf{G} ;
3. Define the **core** vertices;
4. Build the **clusters**;
5. Define the **Hub** and **outliers** vertices.

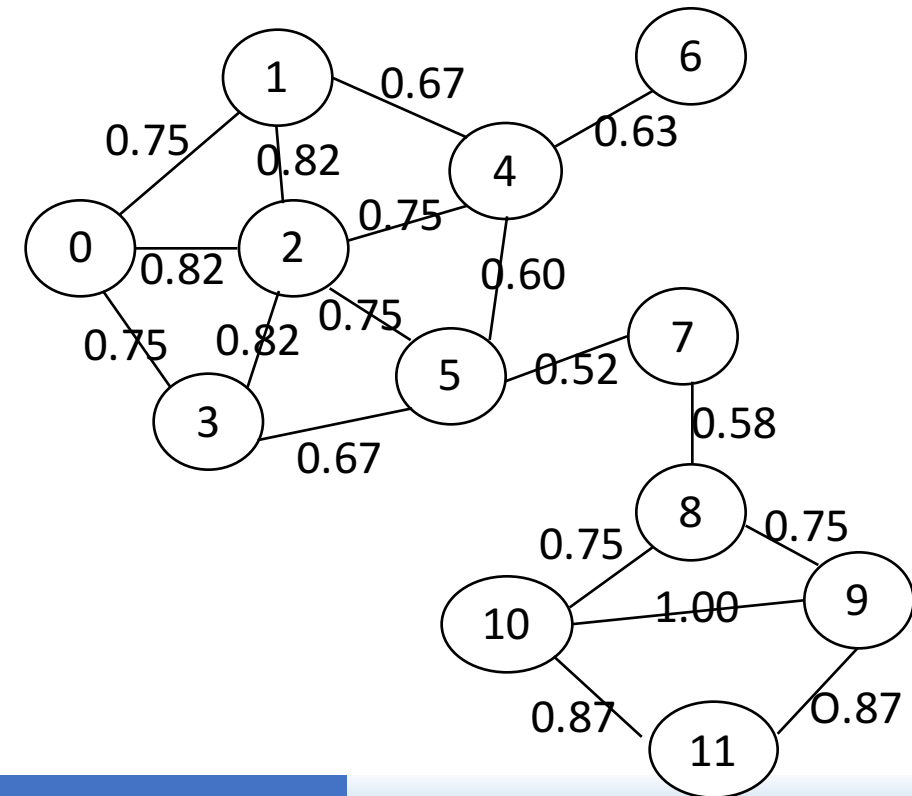
1. Define the neighbors of each vertex;

List of neighbors for each vertex

$$\Gamma(v) = \{w \in V \mid (v,w) \in E\} \cup \{v\}$$

2. Structural similarity of the graph **G**

$$\sigma(v,w) = \frac{|\Gamma(v) \cap \Gamma(w)|}{\sqrt{|\Gamma(v)| |\Gamma(w)|}}$$



3. Cores detection

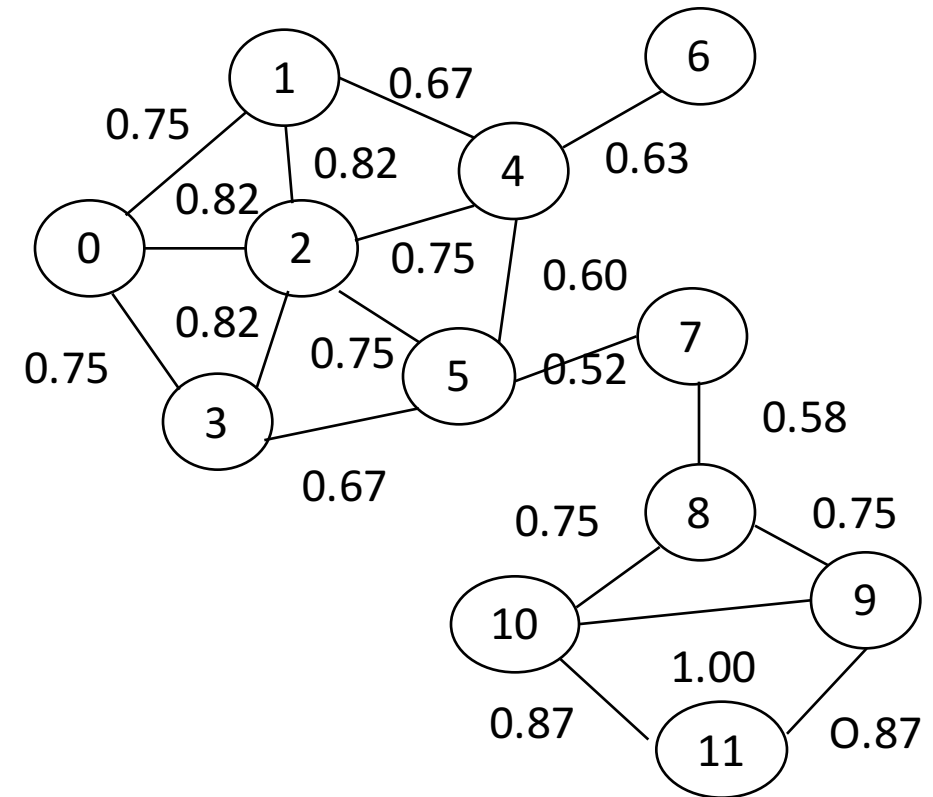
Uses a threshold ϵ to determine a dense connections:

$$N_\epsilon(v) = \{w \in \Gamma(v) \mid \sigma(v, w) \geq \epsilon\}$$

A **core** vertex shares structural similarity of at least ϵ with at least μ neighbors:

$$CORE_{\epsilon, \mu}(v) \Leftrightarrow |N_\epsilon(v)| \geq \mu$$

$\epsilon = 0.7$
 $\mu = 3$



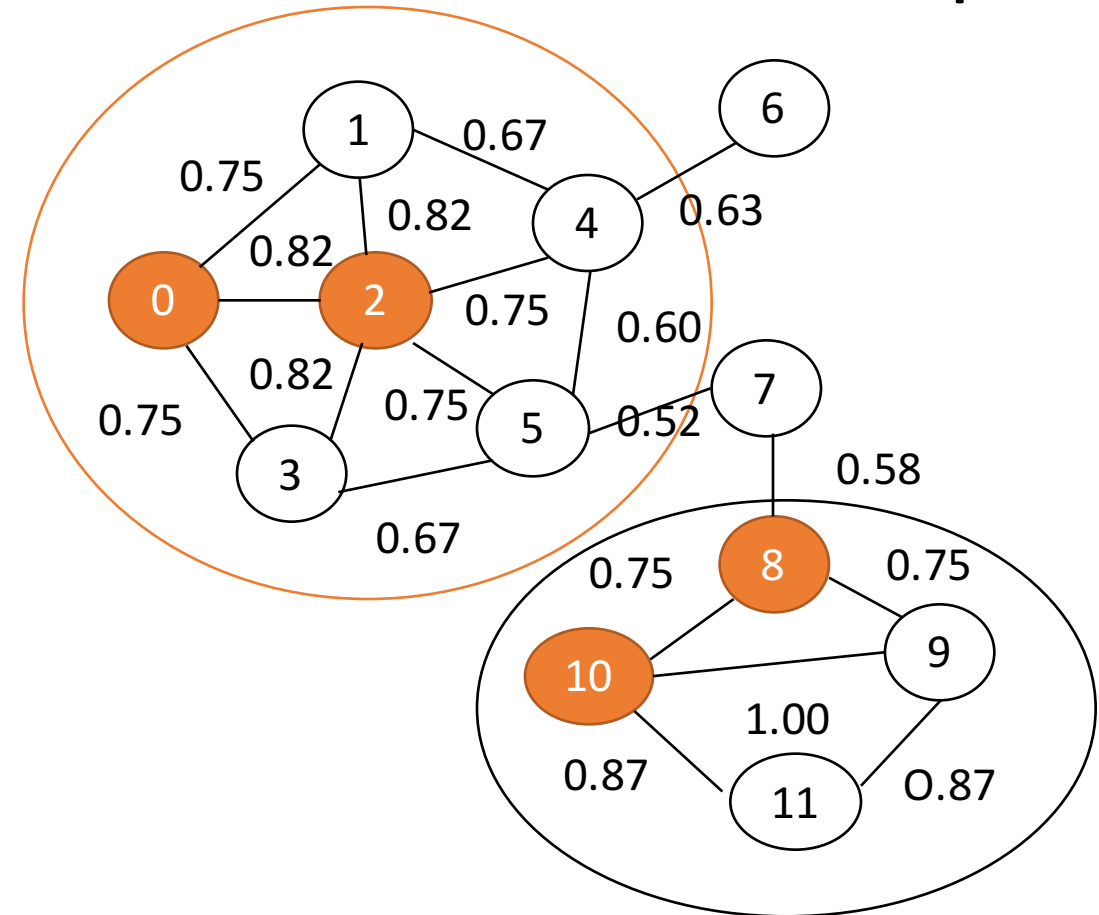
$\epsilon = 0.7$
 $\mu = 3$

4. Build the Clusters

Direct structure reachability:

IF vertex is in ϵ -**Neighborhood** of a **core** vertex,
 they should be **IN** the **same cluster**

$$\text{DirREACH}_{\epsilon, \mu}(v, w) \Leftrightarrow \text{CORE}_{\epsilon, \mu}(v) \wedge w \in N_{\epsilon}(v)$$



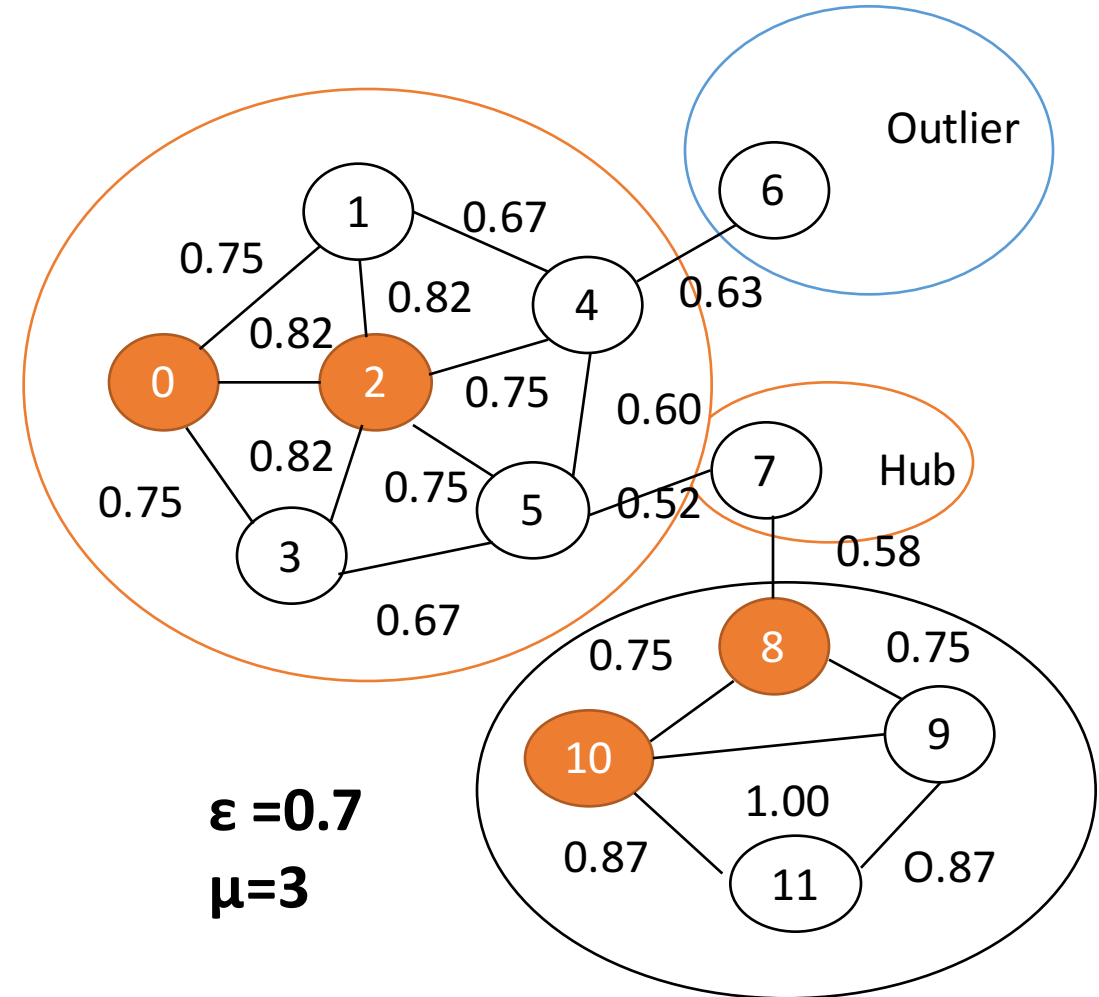
5. Define the **Hub** and **outliers** vertices

Hub:

Is an *isolated vertex* that's neighbors *belonging* of two or more different clusters.

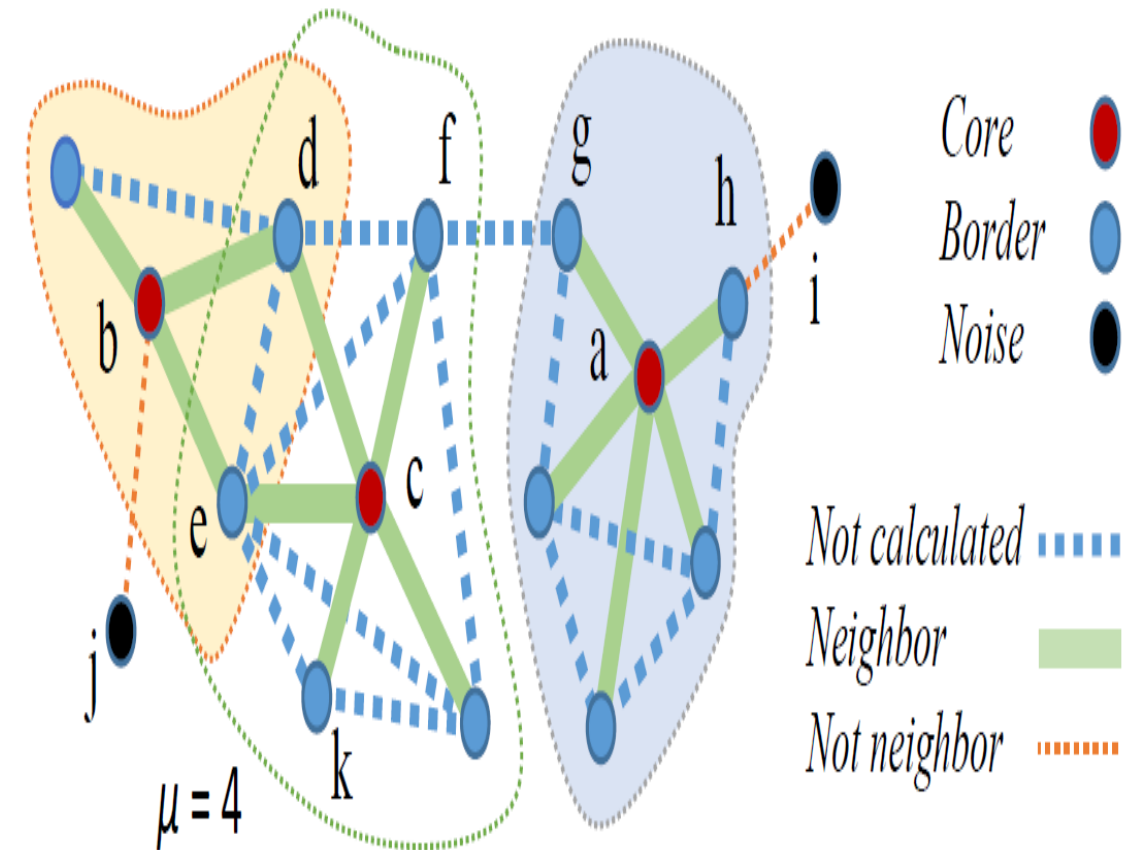
Outlier:

Is an *isolated vertex* that *do not belong* to any cluster.



AnyTimeSCAN method: an extension of SCAN

- AnyTimeSCAN: extension of the SCAN;
- Gives the same result as SCAN;
- (+) **Reduce** the calculation operation of the structural similarity



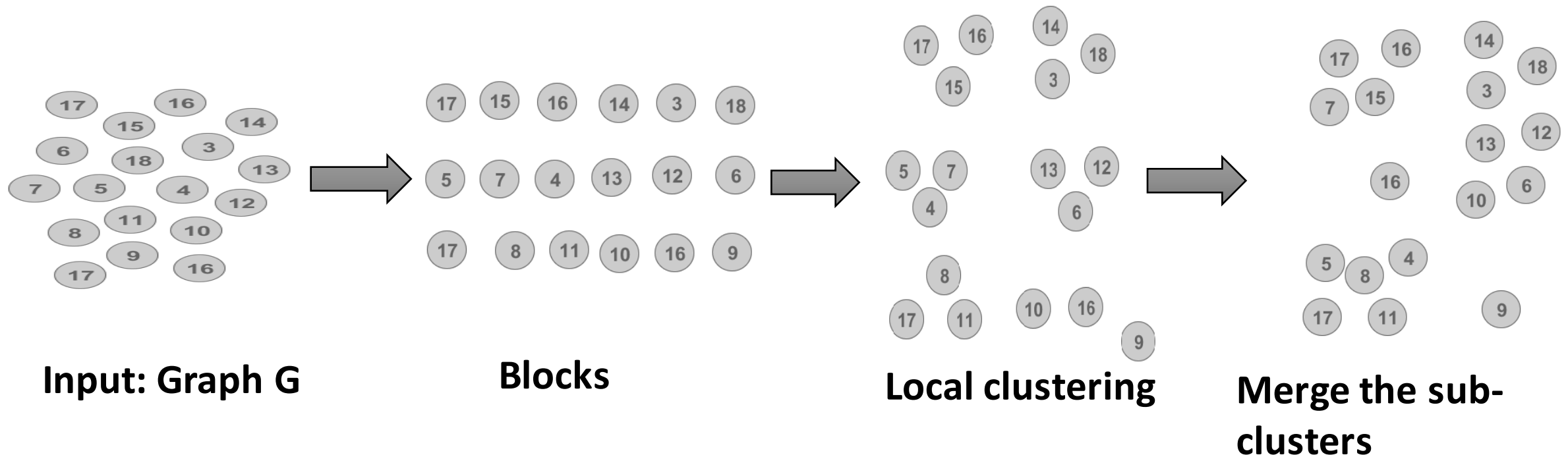
a, b and c: **Core**

i and j: **Noise**

AnyTimeSCAN

The steps of the AnyTime SCAN algorithm

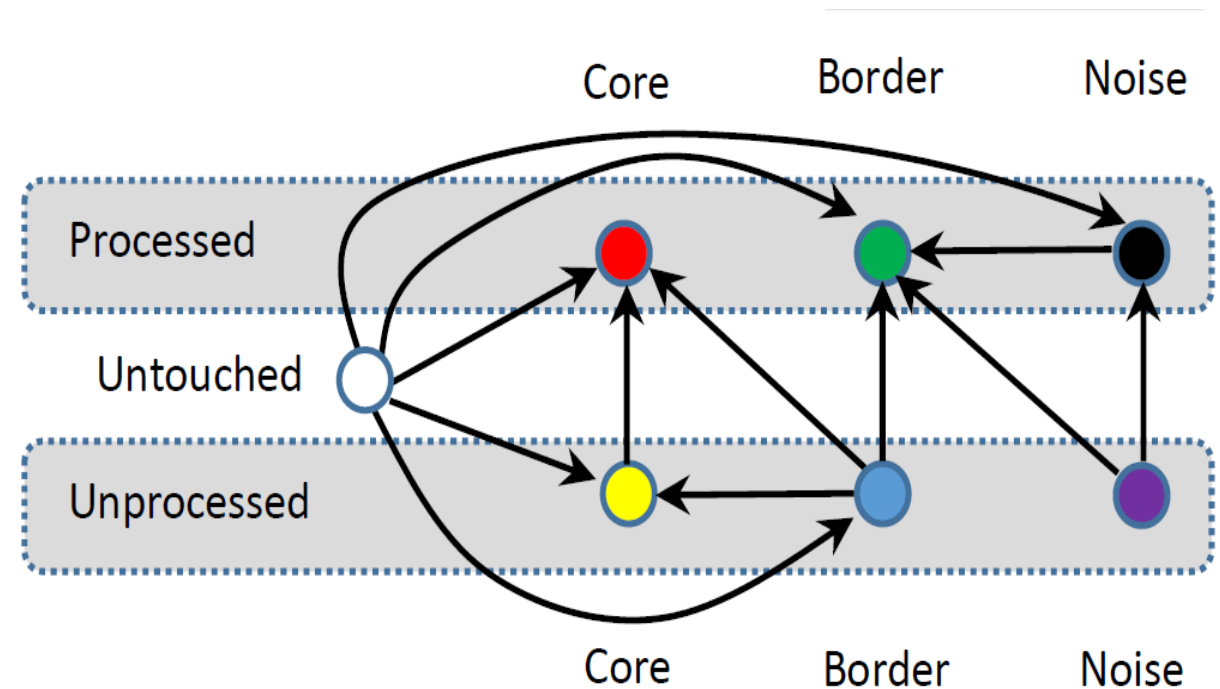
1. Summarization: decomposes the set of vertices into equal blocks
2. Combines the sub-clusters;
3. Determines the borders.



AnyTimeSCAN

AnyTimeSCAN: Summarization step

- At $t = 0$, all vertices \mathbf{V} are marked as untouched vertices.
- For each vertex \mathbf{v} in a bloc \mathbf{b} , we determine its state according to its neighbors.

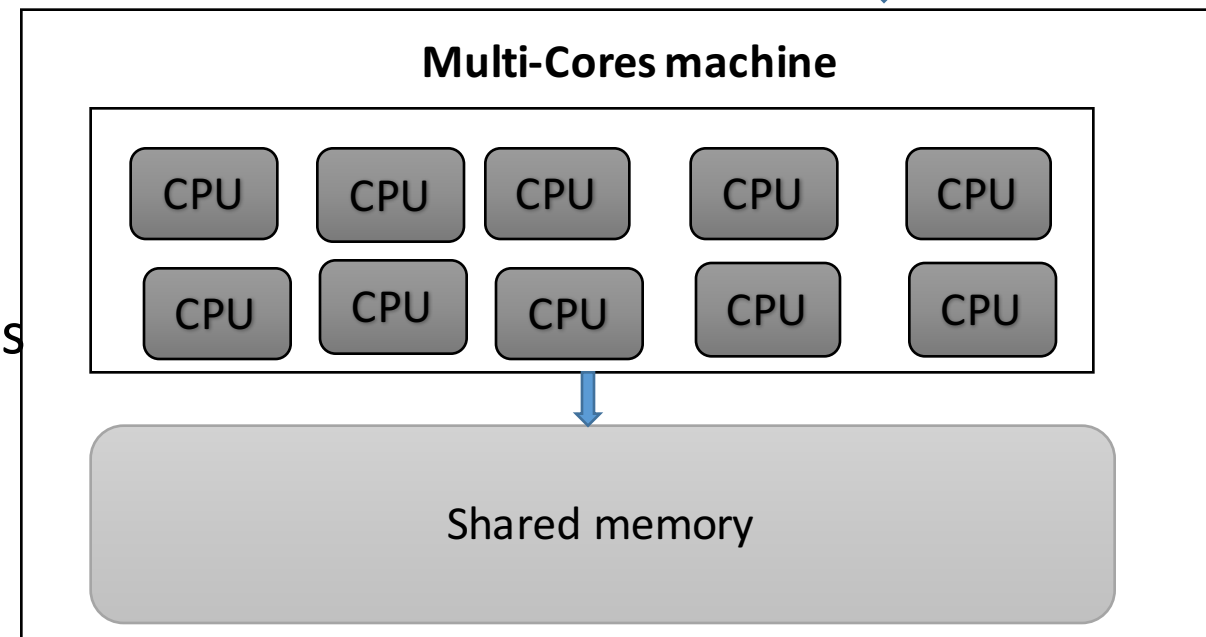
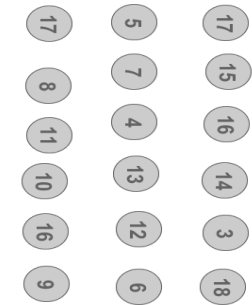


The state transition schema for vertices

AnyTimeSCAN: Parallel processing

- Implementation based on OpenMP API
- Degree of parallelism according to the number of blocks or the parameter $\alpha \gg 1$
- Shared memory:
 - Group all the blocks;
 - Keep the links between the sub-graphs
 - Keep the shared data.

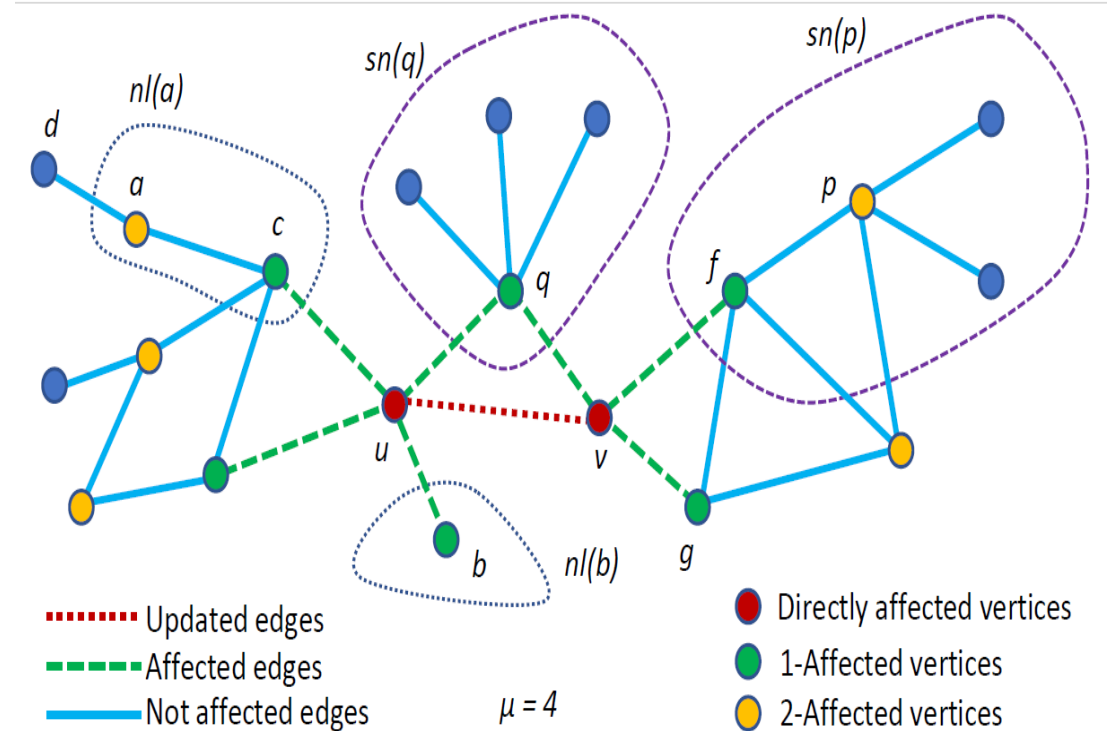
Blocs of the Graph G



AnyTimeSCAN

DanySCAN: Dynamic Clustering

- Incremental approach;
- Each update in the graph, a set of vertices and edges will be affected by this change;



Experimental study

Dataset:

Graph	Vertices	Edges	\bar{d}	c
Ego-Gplus (GR01)	107,614	13,673,453	127.06	0.4901
Soc-LiveJournal1 (GR02)	4,847,571	68,993,773	14.23	0.2742
Soc-Poket (GR03)	1,632,803	30,622,564	18.75	0.1094
Com-Orkut (GR04)	3,072,441	117,185,083	38.14	0.1666
Kron_g500-logn21 (GR05)	2,097,152	182,082,942	86.82	0.1649
Twitter (GR06)	41,652,230	1,369,000,750	32.8	0.0730

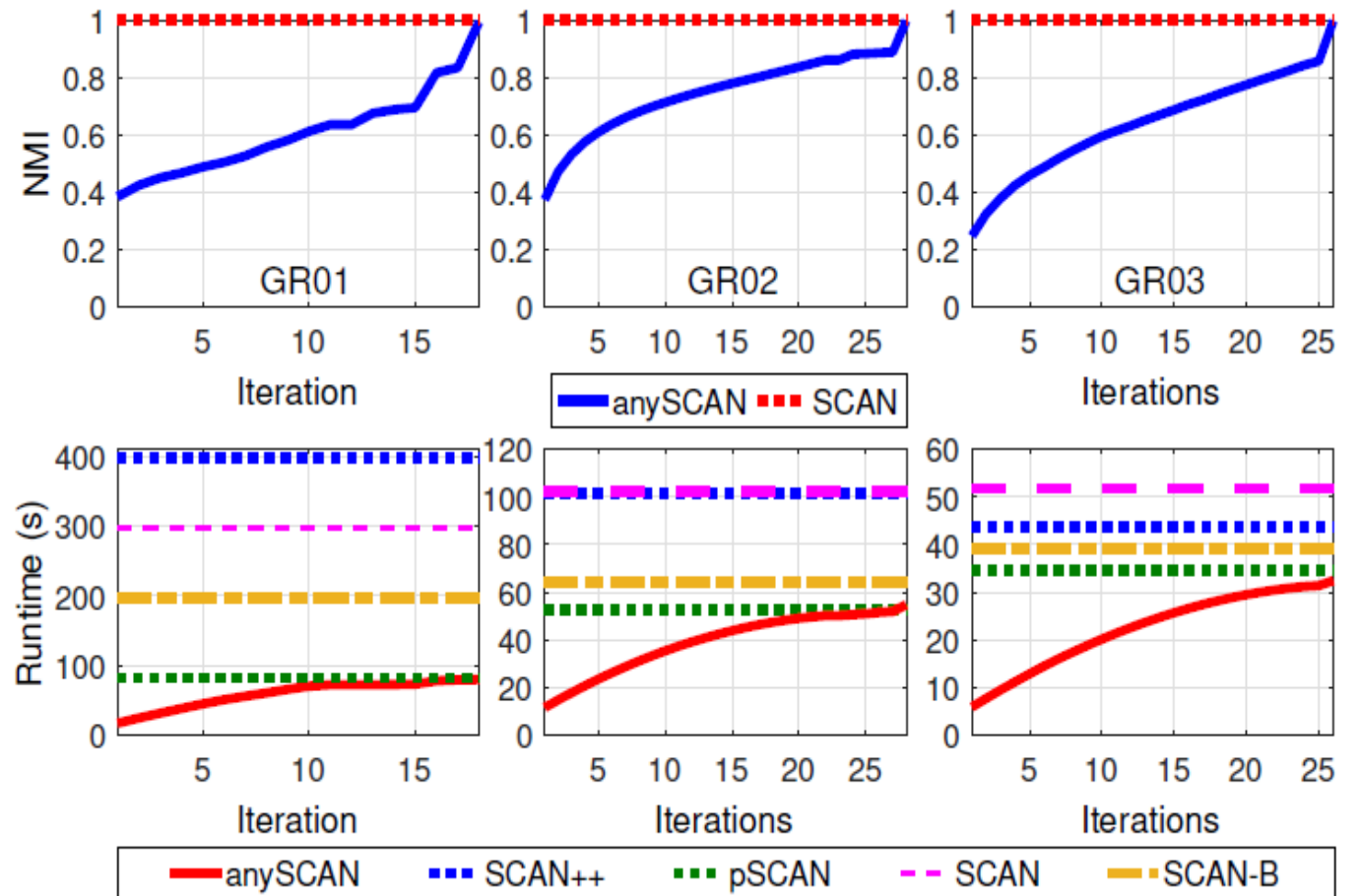
Environment:

- 2 * 3.1 GHz Intel Xeon CPUs with 64 GB local RAM
- OpenMP 3.2

AnyTimeSCAN

AnyTimeSCAN: experimental results

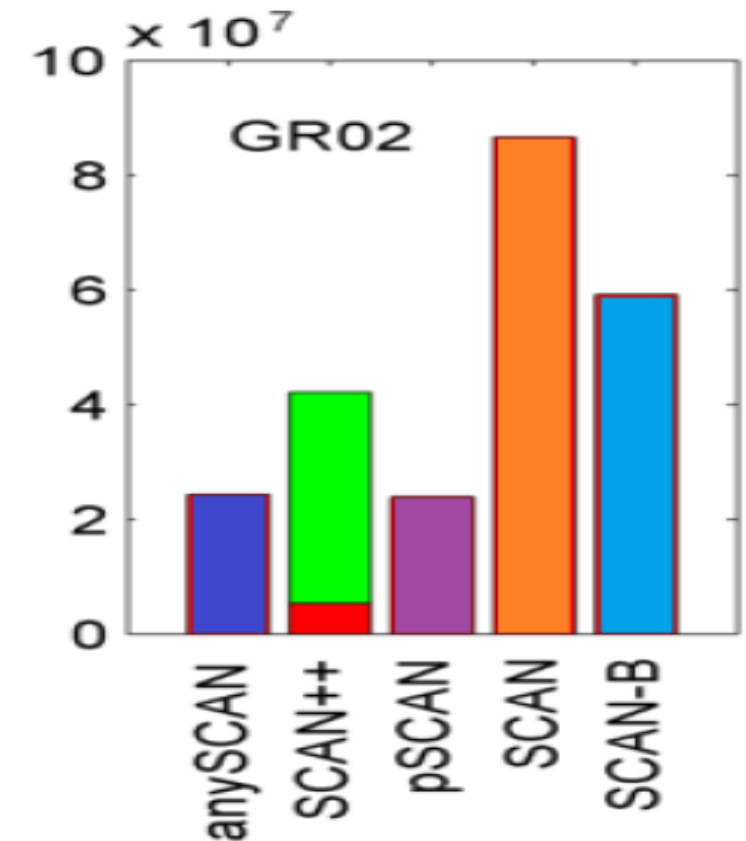
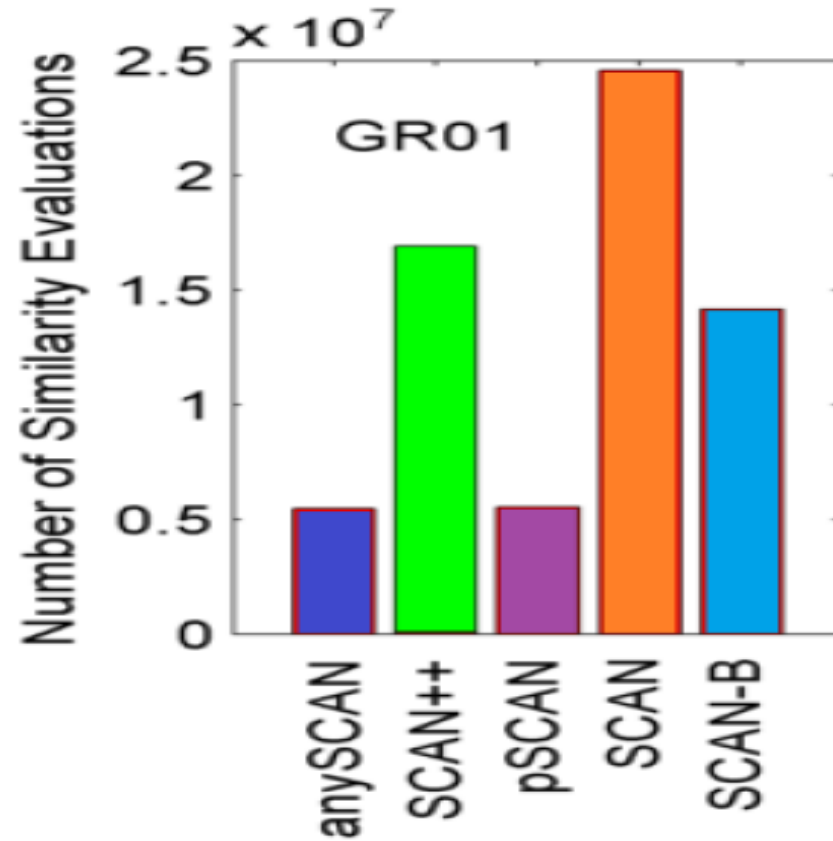
- Give an intermediate results, which has not been guaranteed by the others approaches;
- Runtimes comparison with others approaches.



AnyTimeSCAN

AnyTimeSCAN: experimental results

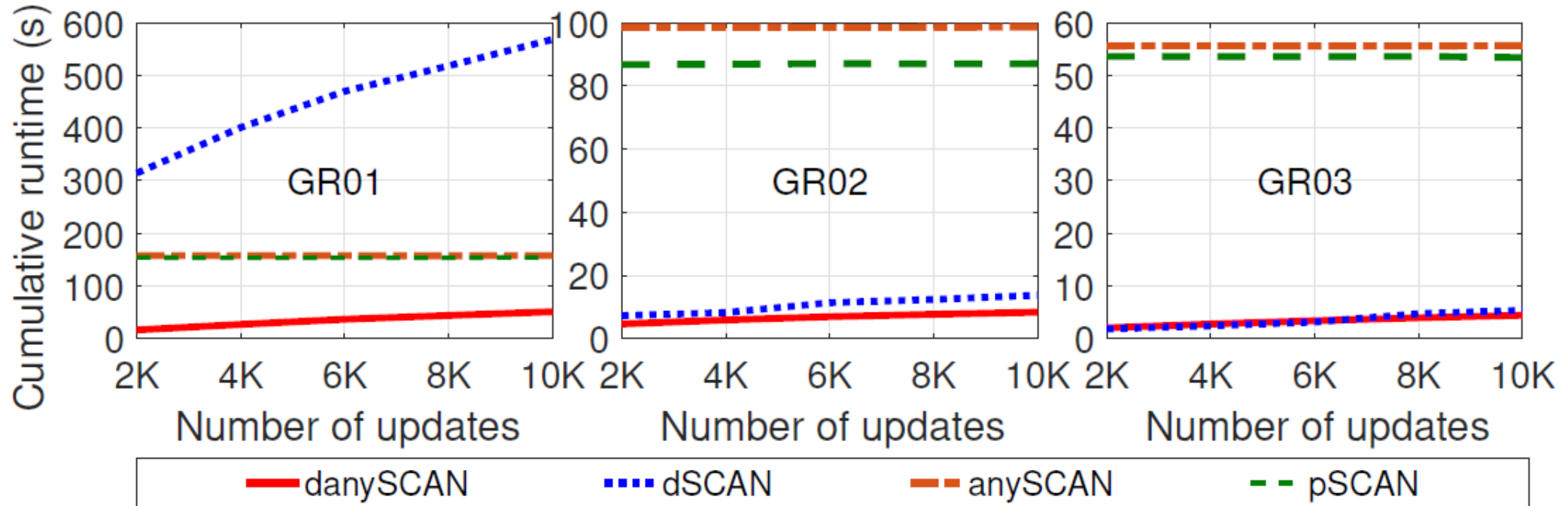
Reduce the structural similarity operation between a pairs of vertices in the graph.



Comparison in terms of the number of structural similarity operation

DanySCAN

DanySCAN: experimental results



Performance of danySCAN according the number of updates for different datasets

Conclusion

- A new paradigm to reduce the similarity calculation operation;
- An incremental approach for dynamic graphs clustering ;
- Interactive model to provide an intermediate results during the execution of the algorithm;
- A parallel approach based on shared memory architecture;
- A strong experimental study;