

# Finding the Hierarchy of Dense Subgraphs using Nucleus Decompositions

A. Erdem Sarıyüce<sup>\*</sup>, C. Seshadhri<sup>+</sup>, Ali Pınar<sup>#</sup>, Ümit V. Çatalyürek<sup>\*</sup>

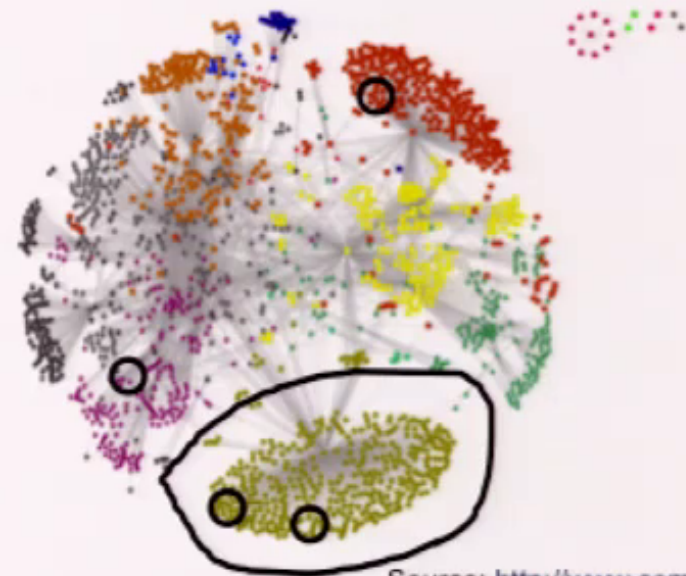
<sup>\*</sup> The Ohio State University

<sup>+</sup> University of California, Santa Cruz

<sup>#</sup> Sandia National Labs, Livermore

# Graphs are globally sparse... yet locally dense.

- Graphs in real world are SPARSE
  - Number of vertices = millions
  - Number of edges  $\approx 10 \times$  vertices
    - Two random vertices unlikely to be connected (prob =  $10^{-5}$ )
- But they contain many dense substructures
  - Within dense region, two random vertices highly likely to be connected (prob = 0.4)



Source: <http://www.complexworld.net/virhulab/ongoing-projects-main>



Community detection:  
label most/all vertices



Dense subgraph discovery:  
Regions with lots of "activity"

# Many applications find dense subgraphs

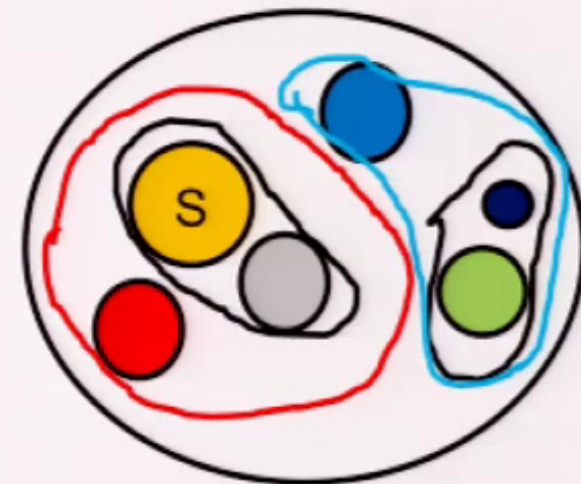
List is long, time is short. Why don't you just trust me?

- Finding communities, spam link farms [Gibson et al., 2005]
- Graph visualization [Alvarez-Hamelin et al., 2006]
- Real-time story identification [Angel et al., 2012]
- DNA motif detection [Fratkin et al., 2006]
- Finding correlated genes [Zhang and Horvath, 2005]
- Finding price value motifs in financial data [Du et al., 2009]
- Graph compression [Buehrer and Chellapilla, 2008]
- Distance query indexing [Jin et al., 2009]
- Throughput of social networking sites [Gionis et al., 2013]
- To name a few...

# Dense subgraphs

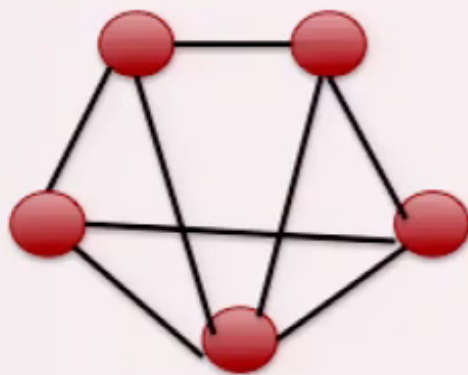
Concept is intuitive, yet formalizations are tricky

- Factors for consideration: size, density of internal edges, density of external edges
- Many formalizations lead to NP-hard problems, and heuristics are used.
- Hard to distinguish, whether an observation is an artifact of the heuristic or not.
- **Our goal:**
  - Can we formulate the problem such that the result is well-defined?
  - Can we find all dense graph not just the densest?
  - Is there a “natural” hierarchy of dense subgraphs?
  - Can we design efficient, provable algorithms and minimize heuristics/approximations?



# K-cores in graphs

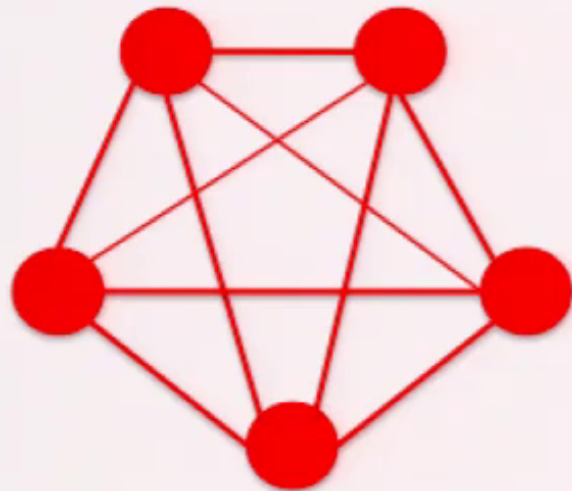
*Unit of observation: vertex*  
*Witness: Edge*



- $k$ -core of a graph is its largest induced subgraph, where degree of each vertex is at least  $k$ .
- Introduced by [Matula and Beck, 1983]
- Algorithm
  - Compute degrees
  - Iteratively remove in increasing order
    - Assign  $K$  value during removal
- **$O(|E|)$  complexity**

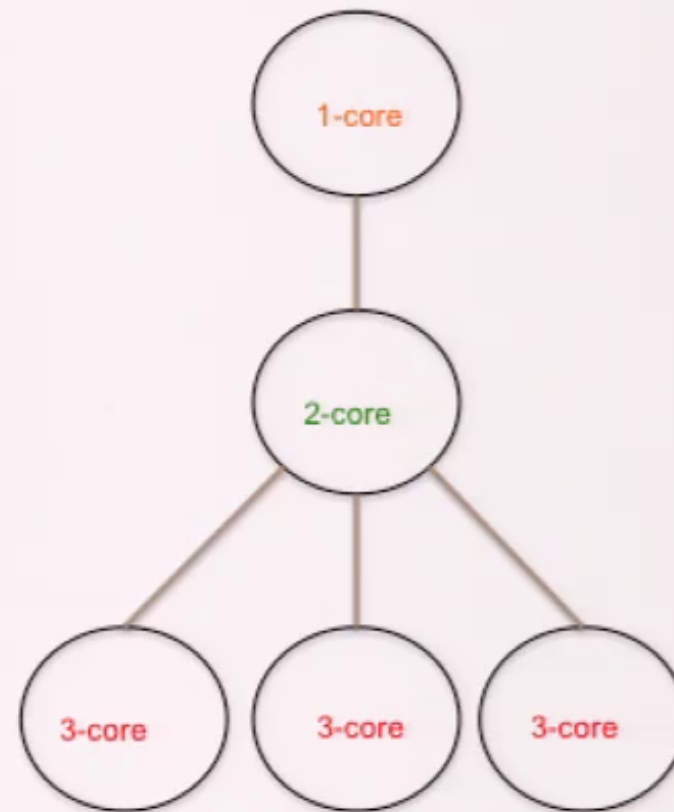
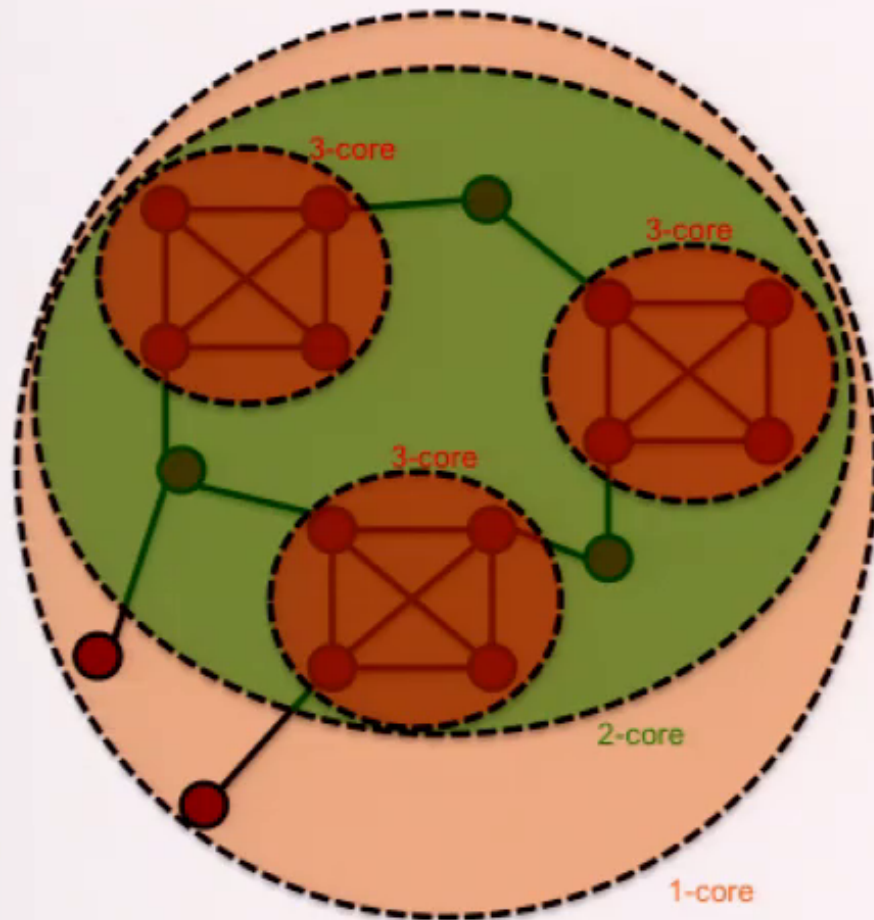
# K-truss decompositions go a step further

Unit of observation: Edge  
Witness: Triangle



- $K$ -truss of a graph is its largest induced subgraph, where each edge participates in at least  $k$  triangles.
- Introduced by Cohen and Parthasarathy independently.
- Applied to visualization and dense graph finding

# Decompositions lead to hierarchies



Caveat:  $k$ -core decomposition typically leads to long chains as opposed to well-branched trees.

# Let us go a step further: Nucleus Decomposition

DEFINITION 1. Let  $r < s$  be positive integers and  $\mathcal{S}$  be a set of  $K_s$ s in  $G$ .

- $K_r(\mathcal{S})$  the set of  $K_r$ s contained in some  $S \in \mathcal{S}$ .
- The number of  $S \in \mathcal{S}$  containing  $R \in K_r(\mathcal{S})$  is the  $\mathcal{S}$ -degree of that  $K_r$ .
- Two  $K_r$ s  $R, R'$  are  $\mathcal{S}$ -connected if there exists a sequence  $R = R_1, R_2, \dots, R_k = R'$  in  $K_r(\mathcal{S})$  such that for each  $i$ , some  $S \in \mathcal{S}$  contains  $R_i \cup R_{i+1}$ .

DEFINITION 2. Let  $k, r$ , and  $s$  be positive integers such that  $r < s$ . A  $k$ - $(r, s)$ -nucleus is a maximal union  $\mathcal{S}$  of  $K_s$ s such that:

- The  $\mathcal{S}$ -degree of any  $R \in K_r(\mathcal{S})$  is at least  $k$ .
- Any  $R, R' \in K_r(\mathcal{S})$  are  $\mathcal{S}$ -connected.

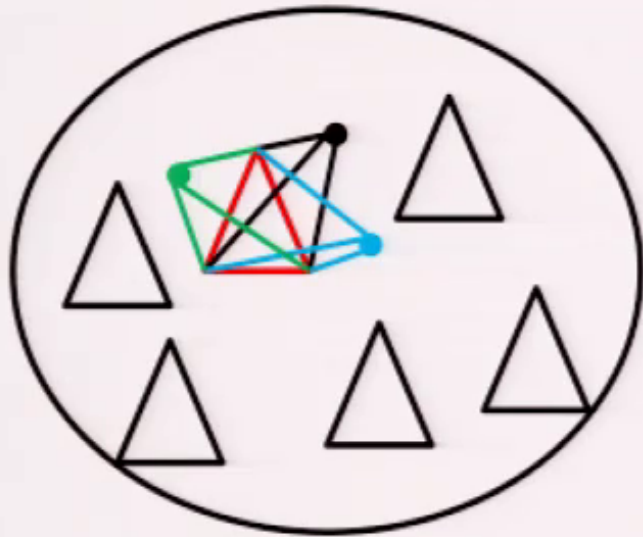
**$r$**  refers to the size of the unit of observation

**$s$**  refers to the size of the witness + unit

**$k$**  is the number of witnesses; not a parameter we sweep through  $k$

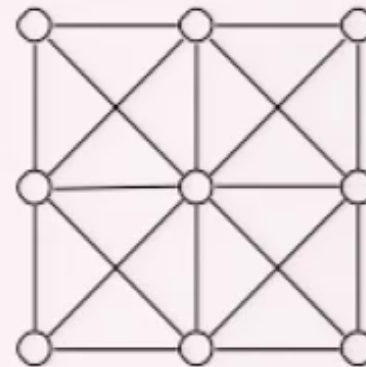


# Examples of nuclei



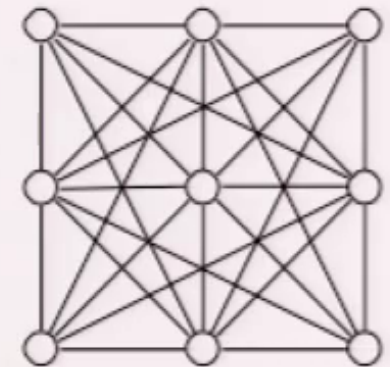
- $k$ -(3,4) nucleus: subgraph formed by maximal union of triangles. Every triangle in at least  $k$  four-cliques
- $k$ -(1,2) is core decomposition
- $k$ -(2,3) is truss decomposition

Edge (2-clique) and  
3-clique interaction

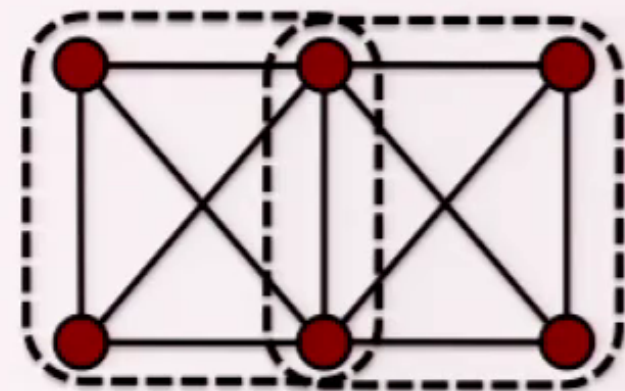


2-(2,3)  
nucleus

Edge (2-clique) and  
4-clique interaction



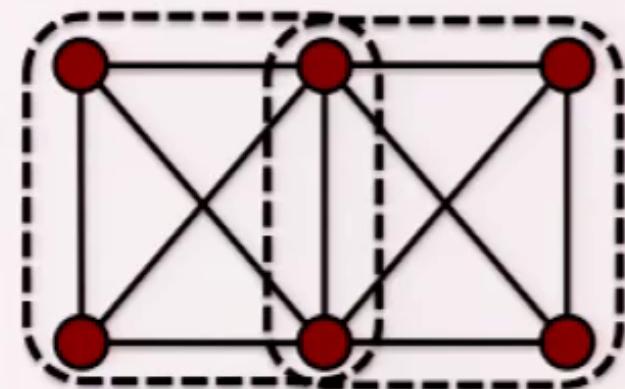
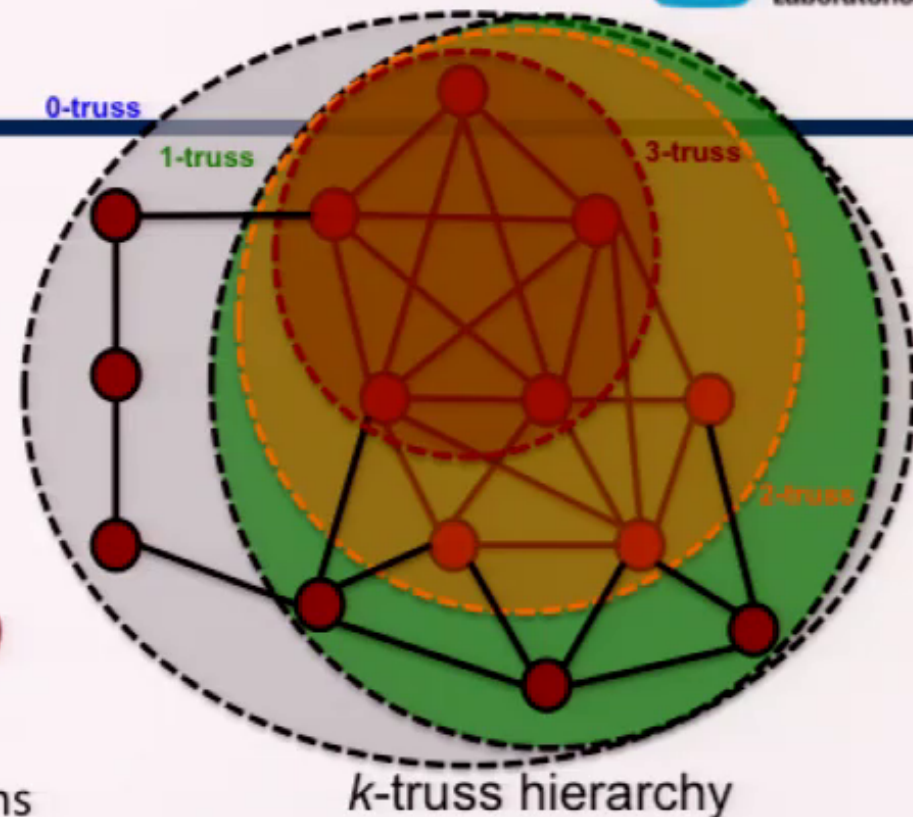
2-(2,4)  
nucleus



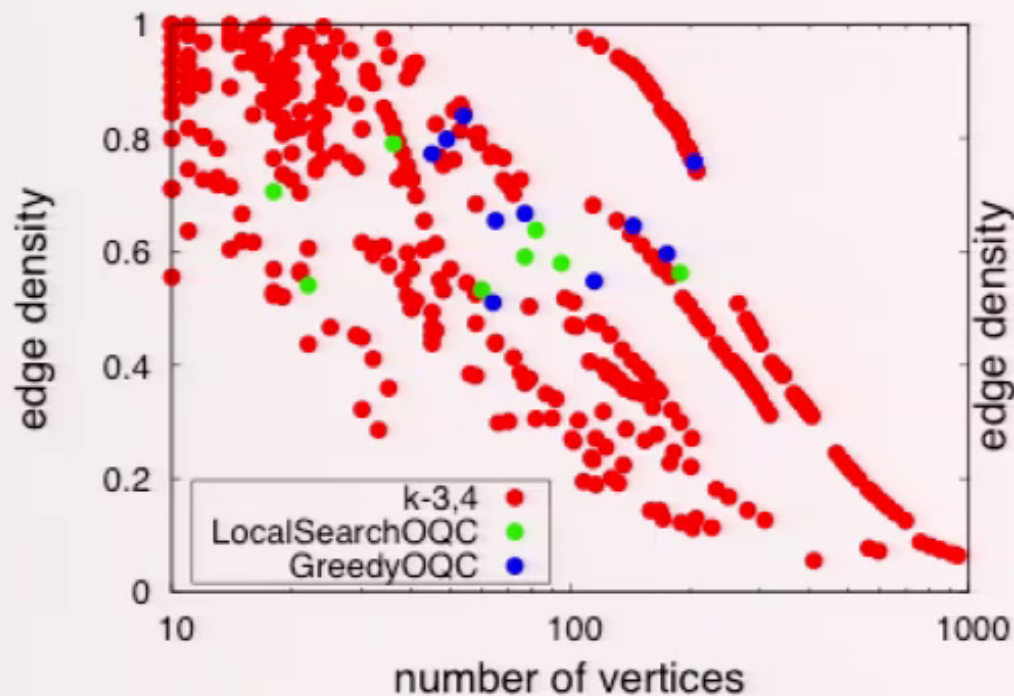
Two 1-(3,4) nuclei

# Properties of nuclei decomposition

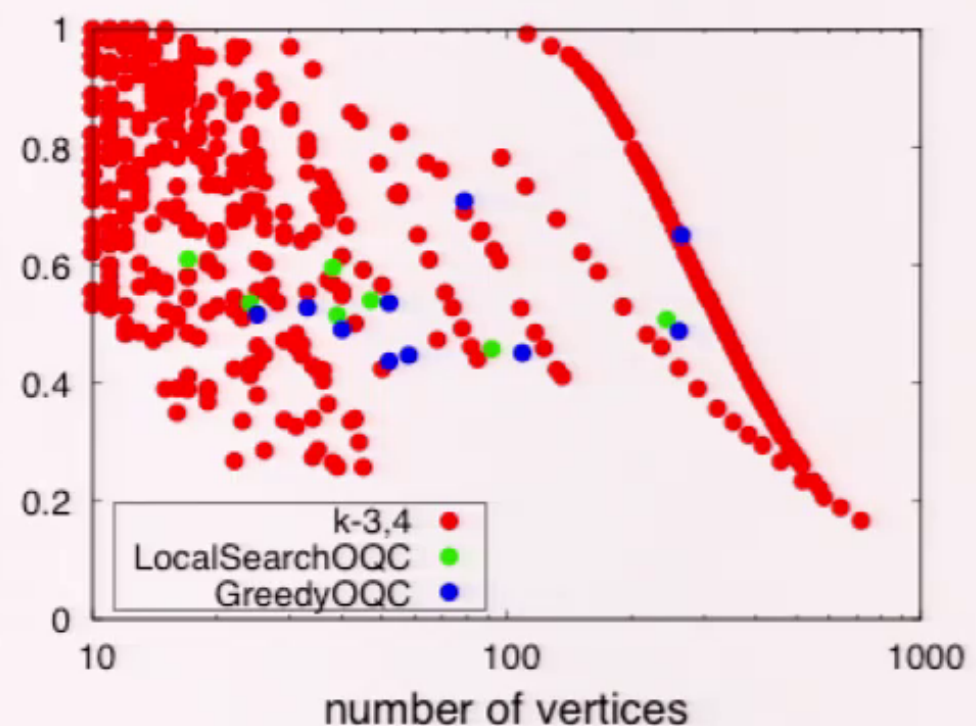
- Well-defined property of the graph
  - Not heuristic
  - No optimization
  - Deterministic**
- Forest of nuclei
  - Smaller  $k$ -( $r,s$ ) contained in larger  $k$ -( $r,s$ )**
  - Hierarchy of dense subgraphs
    - Finding many and understanding relations
- Overlaps of nuclei**
  - For  $r \geq 2$ ,** lower order structures can be shared among nuclei
  - No overlaps for  $k$ -cores!



# Nucleus decomposition finds dense subgraphs



Facebook

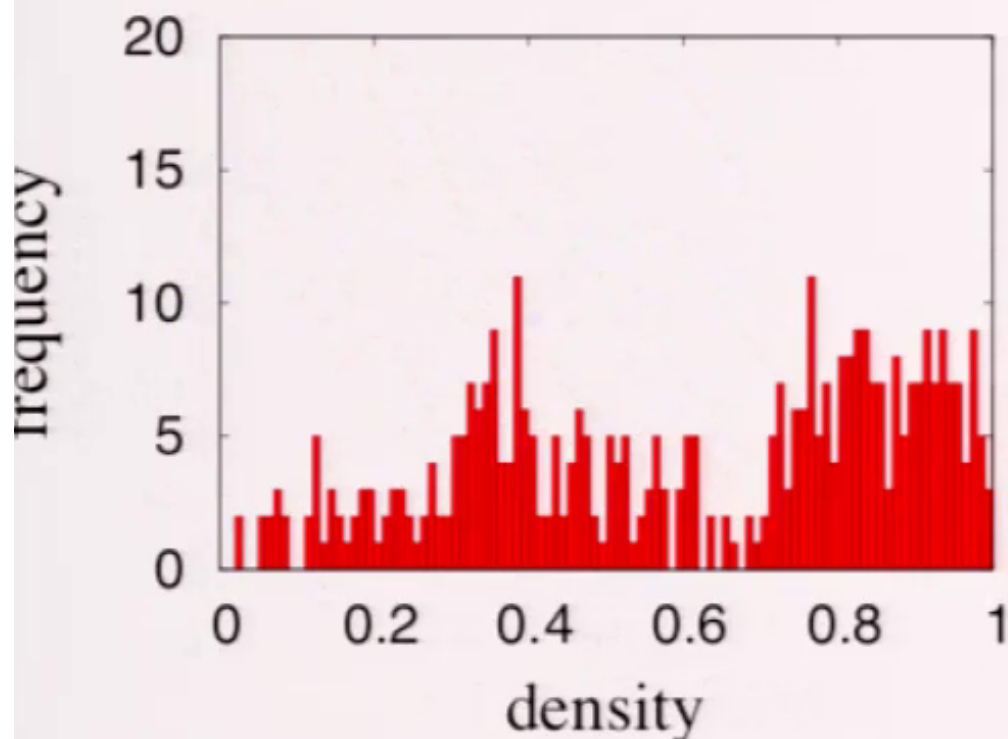


Soc-Epinions

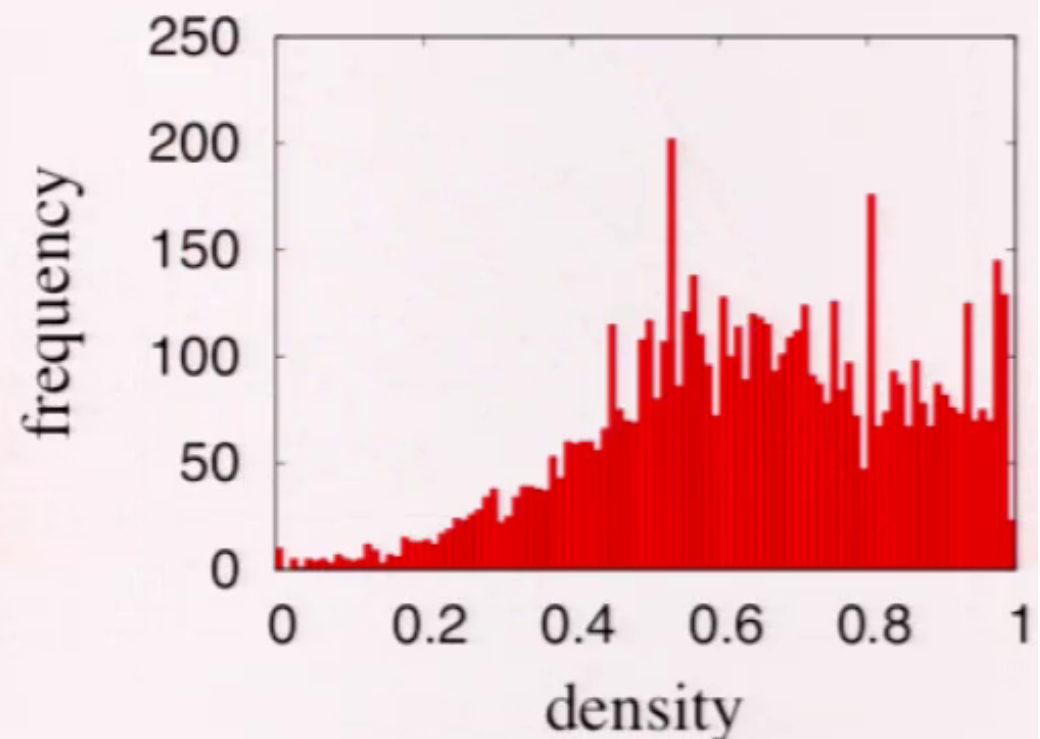
$$\text{Density of } S = E(S,S) / (|S|^2)$$

- We can find many dense subgraphs, not just one at the same time.
- Solution qualities can match the state of the art tools.

# Distributions of dense structures



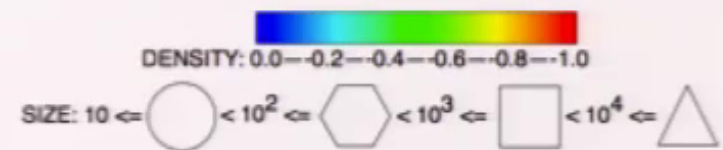
Facebook



Wikipedia

- Finding many dense structures enables producing a density structure profile.

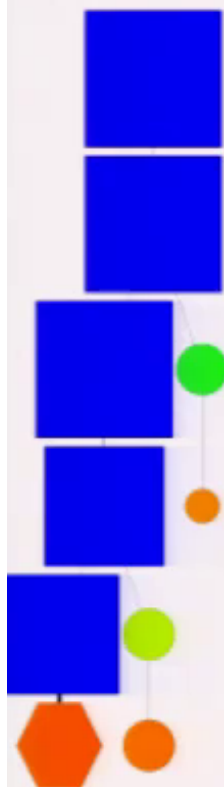
# Hierarchy reveals structure among communities.



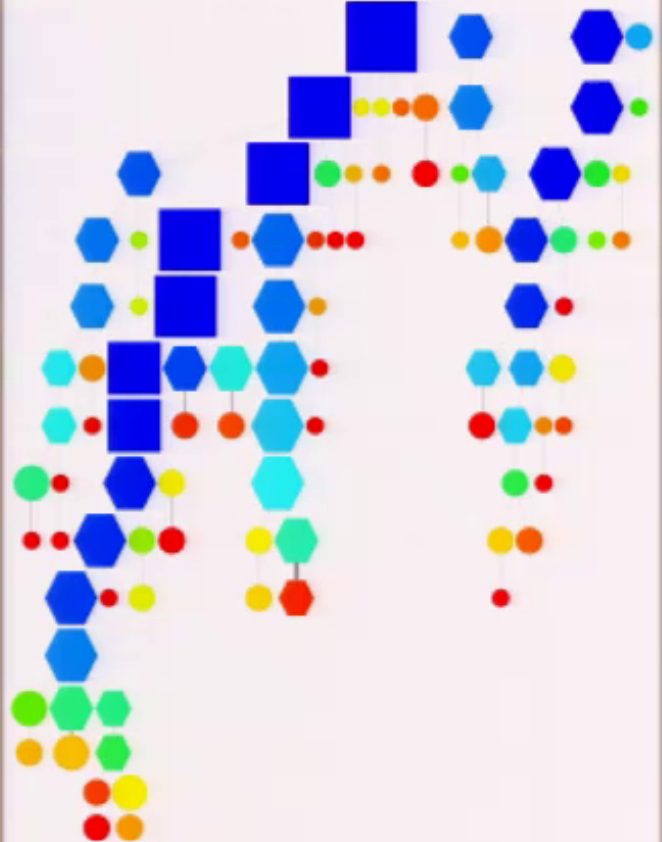
- Results on experimental protein interaction data from Baylor College of Medicine.
- More than 50K vertices, 400K edges, but only few hundred nuclei, with tree of size 50

# Hierarchies (facebook $|V|: 4K, |E|: 88K$ )

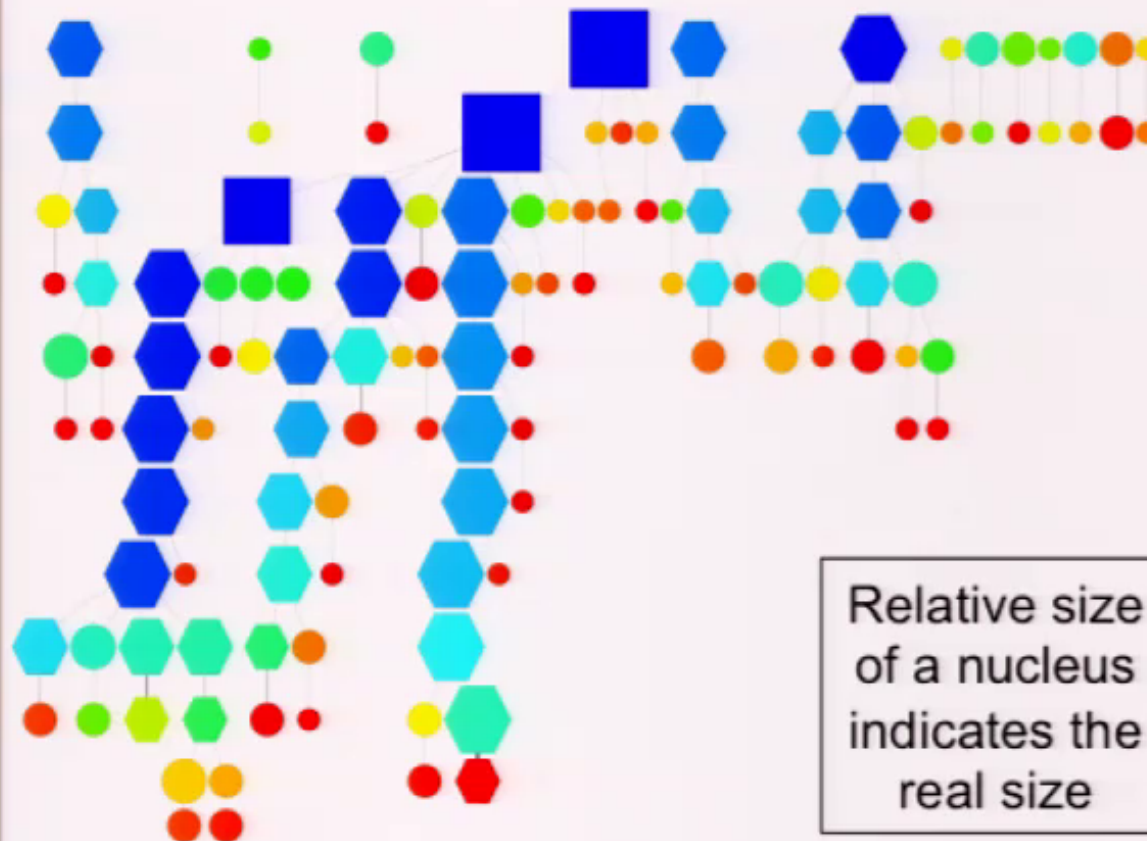
$k-(1,2)$   
( $k$ -core)



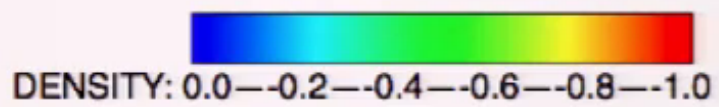
$k-(2,3)$   
( $k$ -truss)



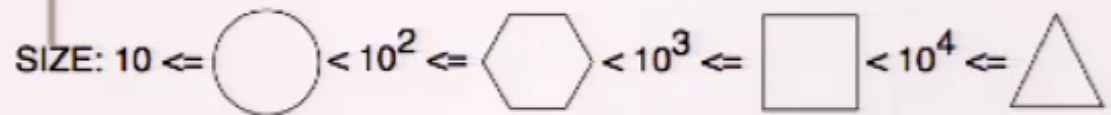
$k-(3,4)$



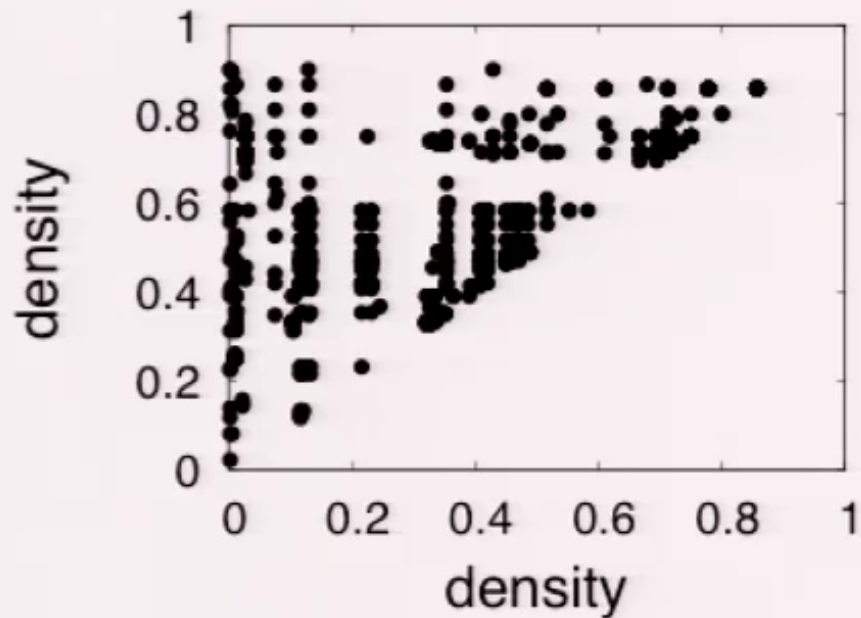
Relative size  
of a nucleus  
indicates the  
real size



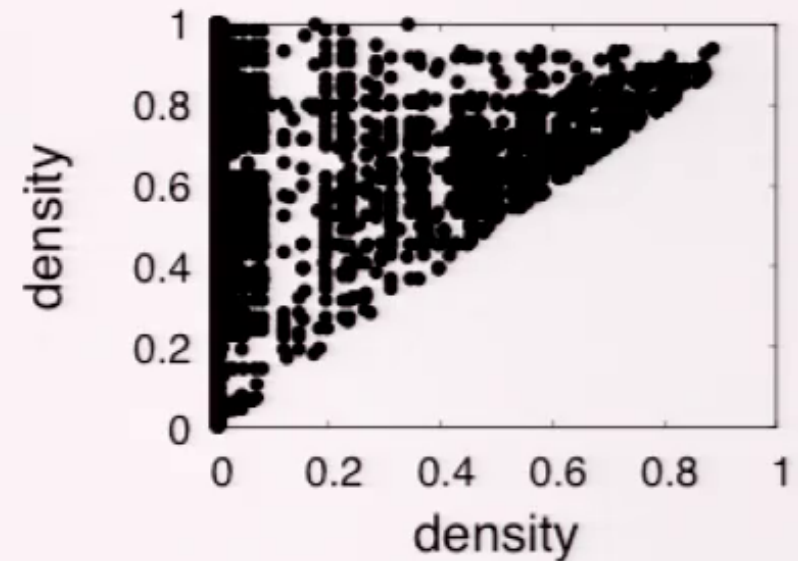
Any nucleus  $\geq 10$  vertices



# Many dense structures overlap

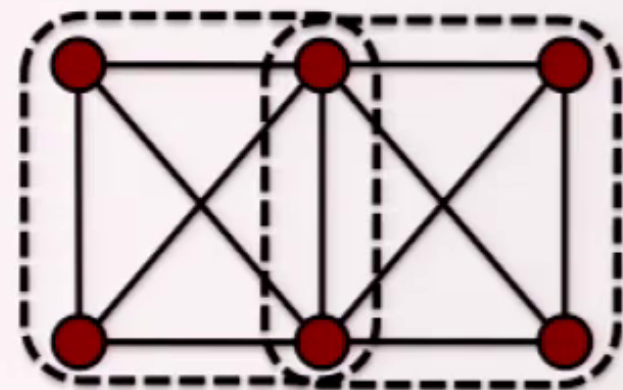


Web-NotreDame



Wikipedia

- Overlap size is at least 5



# How to compute nucleus decomposition

- Given  $r$  and  $s$ , **find all  $k$ - $(r,s)$  nuclei**
- Just like  $k$ -core decomposition
- Find  $K$  values of all  $K_r$ s**

Two ways to implement:

- Enumerate all  $K_r$ s and  $K_s$ s
  - Not feasible for large  $r, s$
  - Huge space complexity
- Construct adj. lists of  $K_r$ s online (only enumerate  $K_r$ s)
  - Better space complexity**
  - Time complexity is**

---

**Algorithm 1:** set- $k(G, r, s)$


---


```


1 Enumerate all  $K_r$ s and  $K_s$ s in  $G(V, E)$ ;
2 For every  $K_r$   $R$ , initialize  $\delta(R)$  to be the number of  $K_s$ s
  containing  $R$ ;
3 Mark every  $K_r$  as unprocessed;
4 for each unprocessed  $K_r$   $R$  with minimum  $\delta(R)$  do
5    $\kappa(R) = \delta(R)$ ;
6   Find set  $S$  of  $K_s$ s containing  $R$ ;
7   for each  $S \in \mathcal{S}$  do
8     if any  $K_r$   $R' \subset S$  is processed then
9       Continue;
10    for each  $K_r$   $R' \subset S$ ,  $R' \neq R$  do
11      if  $\delta(R') > \delta(R)$  then
12         $\delta(R') = \delta(R) - 1$ ;
13    Mark  $R$  as processed;
14 return array  $\kappa(\cdot)$  ;
  
```

---

$$O(RT_r(G) + \sum_v ct_r(v)d(v)^{s-r})$$

  
 Total num of  $K_r$ s

  
 Num of  $K_r$ s of  $v$

  
 Degree of  $v$



# Future Directions

- **Applications of nucleus decomposition**
  - Protein-protein and protein-gene interaction networks
  - Ongoing collaboration
- Larger values of  $r$  and  $s$ 
  - Computational cost of increasing  $r$  and  $s$  is significant.
  - **Preliminary experimentation for (4,5)**
    - Very little quality benefit
  - **Is (3,4) a sweet spot?**
- Faster  $k$ -(3,4)
  - Clique enumeration
  - **Parallel algorithms**
    - GPU implementation of  $k$ -core [Jiang et al., 2014]
    - Pregel algorithm for  $k$ -truss [Shao et al., 2014]