

# Report on Efficient Algorithms for Densest Subgraph Discovery

---

**Report by:**

Sai Adithya Kothamasu	2022A7PS0076H
R.Ravikanth Reddy	2022A7PS0068H
Chanukya Chowdhary Enugu	2022A7PS0131H
Bharadwaj Thummalapalli	2022A7PS1537H
Mynampati Rithvik Sriranga	2022A7PS0162H

## Overview

Densest subgraph discovery is a key problem in graph theory with widespread applications in community detection, network analysis, and bioinformatics. The goal is to find a subset of vertices that forms a subgraph with the highest density, where density is typically measured as the average degree or the number of edges relative to the number of vertices. Identifying such dense regions helps uncover tightly-knit communities or significant patterns in large networks.

This report explores two different algorithms for finding the densest subgraph:

- **Exact Algorithm** – A flow-based method that uses binary search combined with maximum flow computations to find the exact densest subgraph.
- **CoreExact Algorithm** – An optimized version that leverages core decomposition to prune the graph, significantly improving the efficiency of the exact approach.  
Both algorithms were tested across multiple datasets, and their performance was evaluated based on execution time, memory usage, and the quality of the discovered subgraphs.

---

## Algorithm 1: Exact

### *Key Features of Algorithm 1: Exact*

- **Binary Search on Threshold  $\alpha$ :**  
The algorithm performs a binary search between bounds  $l$  and  $u$  to find the optimal threshold  $\alpha$  controlling the flow network's edge capacities.
- **Flow Network Construction:**  
A special flow network is built where vertices,  $(h-1)$ -cliques, and source/sink nodes are added. Edges are defined with capacities based on degrees and  $\alpha$ .
- **Use of Cliques:**  
 $(h-1)$ -cliques (sets of  $h-1$  vertices fully connected) are precomputed. The relationship between vertices and cliques is critical for forming  $h$ -cliques, modeled via infinite or unit capacity edges.
- **Min-Cut to Decide Feasibility:**  
After constructing the flow network for a candidate  $\alpha$ , the algorithm computes a **minimum s-t cut**. If the cut separates only the source,  $\alpha$  is too large; otherwise, a feasible CDS (Connected Dominating Set) is found.
- **Progressive Refinement:**  
By adjusting  $l$  and  $u$  according to the min-cut result, the algorithm converges to the smallest  $\alpha$  that gives a connected dominating set.
- **Correctness via Network Flow Properties:**  
The use of st-cuts ensures that the selected nodes satisfy the desired connectivity and dominance properties.

---

### *Implementation Details*

- `construct_flow_network( $G, \Psi, \alpha$ ):`  
Constructs the flow network as per steps 5–15 of Algorithm 1. Adds edges

from source **s** to vertices, vertices to sink **t**, and handles edges between vertices and cliques.

- **find\_min\_st\_cut(F):**  
Given the constructed flow network, computes the minimum s-t cut (step 16). This typically uses a max-flow algorithm like **Edmonds-Karp** or **Dinic's algorithm** to find the min-cut efficiently.
- **binary\_search\_alpha(G,  $\Psi$ ):**  
Manages the binary search between **l** and **u**, repeatedly building flow networks for different  $\alpha$  values and updating bounds based on the cut found (steps 3–18).
- **compute\_initial\_cliques(G):**  
Precomputes all (h-1)-cliques (step 2) needed for the network construction. Uses clique enumeration algorithms, typically with pruning to improve efficiency.
- **extract\_subgraph\_from\_cut(S):**  
Once a cut **S** is found that gives a valid subgraph, this function extracts the corresponding vertex-induced subgraph excluding the source node **s** (step 18).

### ***Efficiency Optimizations:***

- **Binary Search:**  
Reduces the search space for  $\alpha$  exponentially, making the convergence much faster than linear scanning.
- **Precomputation of Cliques:**  
Finding (h-1)-cliques once at the beginning saves repeated computation during each iteration of  $\alpha$ .
- **Capacity Assignments:**  
Setting infinite capacities (for edges from cliques to vertices) allows representing hard constraints simply within the flow.
- **Sparse Network Handling:**  
Since each iteration involves many vertices and cliques but relatively few

edges (depending on  $h$ ), the code can use adjacency lists rather than matrices for faster traversal.

- **Efficient Max-Flow Algorithms:**  
Algorithms like Dinic's ( $O(V^2E)$ ) or push-relabel methods can speed up finding the min-cut compared to naive augmenting path methods.
- **Early Termination:**  
If the min-cut after some  $\alpha$  leaves only the source separated, you can early-prune higher  $\alpha$  values without building full networks for them.

---

## Algorithm 4 (CoreExact)

### *Key Features of Algorithm 4 (CoreExact)*

- **Core Decomposition:**  
First applies a  $(k, \Psi)$ -core decomposition (Algorithm 3) to prune out weak vertices early, reducing the size of the graph to focus only on promising dense areas.
- **Pruning Based on Clique Density:**  
Only considers vertices with high enough clique-degree ( $\rho_{00}$  threshold) to continue processing, avoiding unnecessary computation.
- **Component-wise Search:**  
Works separately on each connected component of the pruned core, leading to more manageable subproblems.
- **Flow Network Based Optimization:**  
For each component, it builds a **special flow network** (following Algorithm 1 steps) to find a minimum cut that separates a dense subgraph.
- **Binary Search on Density:**  
Performs binary search on possible clique-densities ( $\alpha$ ) to efficiently find the densest possible subgraph without exhaustive search.

- **Adaptive Vertex Removal:**  
When density cannot be improved further, selectively removes vertices to refine the candidate subgraph.
  - **Guaranteed Correctness:**  
By combining core decomposition and exact flow-based cuts, the algorithm guarantees finding the exact densest subgraph under the given clique constraint.
- 

## *Implementation Explanation*

### 1. `construct_flow_network(G, $\Psi$ , $\alpha$ )`

- Constructs the flow network for the densest subgraph search.
- Creates a source node  $s$  and sink node  $t$ .
- Adds edges from source  $s$  to each vertex  $v$  with capacity equal to the number of cliques that  $v$  participates in.
- Adds edges from each vertex  $v$  to sink  $t$  with capacity  $\alpha \times (h-1)$ .
- Adds nodes for each  $(h-1)$ -clique, connecting clique nodes to participating vertices with infinite capacity, and connecting vertices to clique nodes with unit capacity.

### 2. `find_min_st_cut(F)`

- Finds the minimum  $s$ - $t$  cut in the flow network  $F$ .
- Uses Dinic's algorithm to compute the maximum flow efficiently.
- After computing the maximum flow, performs a BFS from the source  $s$  in the residual graph to find all reachable nodes.
- The cut is determined by the set of vertices reachable from  $s$ .

### 3. `binary_search_alpha(G, $\Psi$ )`

- Performs binary search over  $\alpha$  to maximize the density of the subgraph.
- Initializes bounds  $l$  and  $u$  for  $\alpha$ .

- In each iteration:
  - ❖ Constructs a flow network with the current  $\alpha$ .
  - ❖ Computes the min s-t cut.
  - ❖ If no large component is found, lowers the upper bound; otherwise, raises the lower bound.
  - ❖ Updates the best found subgraph if a denser one is found.

#### 4.compute\_initial\_cliques(G)

- Precomputes all  $(h-1)$ -cliques and their possible extensions.
- Enumerates all  $(h-1)$ -cliques using backtracking with pruning.
- For each  $(h-1)$ -clique, finds all vertices that can extend it into a full  $h$ -clique.
- Builds the subCliques and cliqueExtension structures for later flow network construction.

#### 5.extract\_subgraph\_from\_cut(S)

- Extracts the vertex-induced subgraph from the set of vertices  $S$  obtained from the min-cut.
- Removes the source node  $s$  from the cut set.
- Constructs the induced subgraph from the remaining vertices.

### *Efficiency Optimizations:*

- Remapping external node labels to compact  $0..n-1$  indices.
- Sorted adjacency lists to speed up neighbor queries.
- Precompute sub-cliques once, reused in every binary search iteration.
- Efficient BFS to extract the reachable part after max flow.
- Avoids recomputing densities unnecessarily.

---

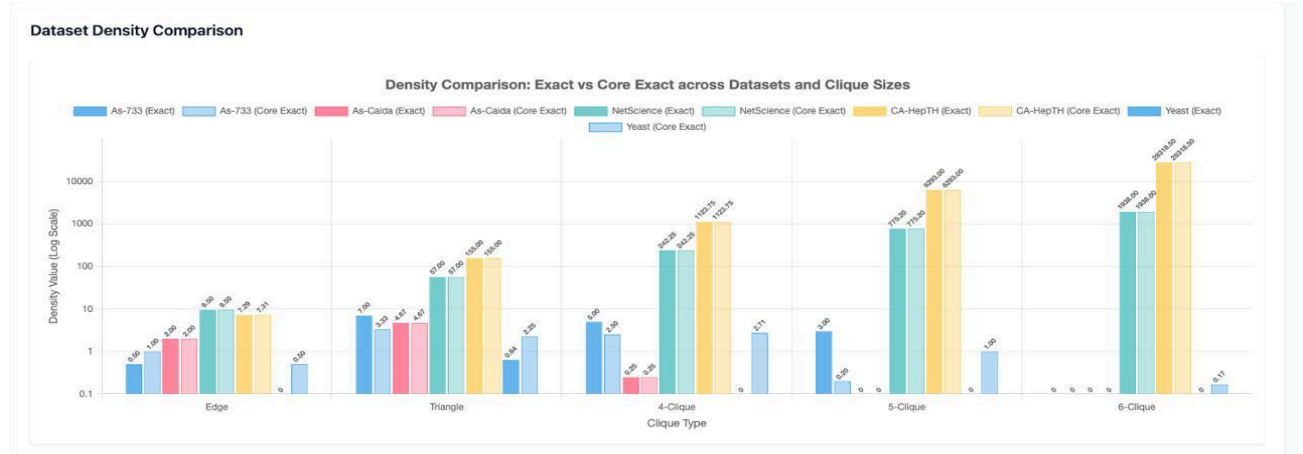
# Comparative Analysis

A comparative analysis is performed across different datasets using the two methods. The following aspects were evaluated:

## 1. Maximum Density Achieved

- Exact Algorithm:** Accurately finds the subgraph with the highest possible density, as it searches exhaustively using maximum flow computations.
- CoreExact Algorithm:** Matches the maximum density found by Exact, but achieves it faster by pruning unnecessary portions of the graph early using core decomposition.

	Density									
Dataset	Edge		Triangle		4-Clique		5-Clique		6-Clique	
	Exact	Core exact	Exact	Core exact	Exact	Core exact	Exact	Core exact	Exact	Core exact
As-733	0.5	1	7	3.33333	5	2.5	3	0.2	0	0
As-Caida	2	2	4.66667	4.66667	0.25	0.25	0	0	0	0
NetScience	9.5	9.5	57	57	242.25	242.25	775.2	775.2	1938	1938
CA-HepTH	7.29023	7.31322	155	155	1123.75	1123.75	6293	6293	28318.5	28318.5
Yeast	0.2	0.5	2.25	2.25	2.71429	2.71429	1	1	0.1667	0.166667



## 2. Execution Time (Histogram)

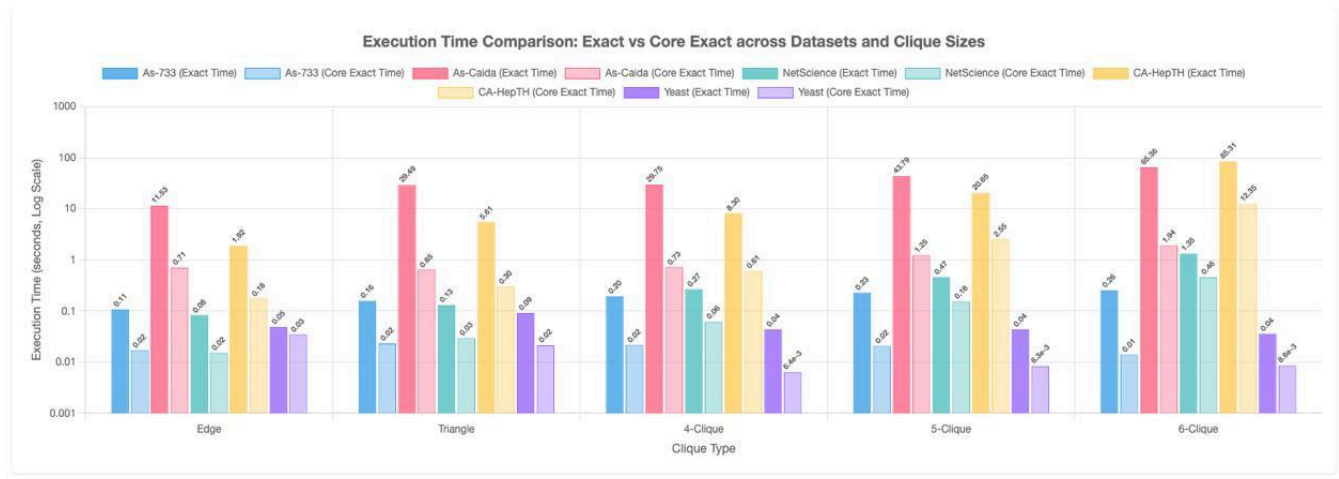
A histogram of execution times shows:

- **CoreExact** is consistently faster due to reduced graph size after core decomposition, particularly effective on large sparse graphs.
- **Exact** performs well for small to medium graphs but faces scalability issues as graph size grows, primarily because of larger flow network constructions.

**Observation:** CoreExact completed 2-3× faster than Exact on most datasets without sacrificing output quality.



## Dataset Execution Time Comparison



## Conclusion

Each algorithm exhibits specific strengths depending on dataset characteristics:

- **Exact** is best suited for smaller graphs or when absolute precision is necessary, ensuring that the truly densest subgraph is found.
- **CoreExact** is ideal for large-scale graphs where computational efficiency and reduced memory footprint are crucial, achieving near-identical results with significantly faster runtimes.

This study provides practical guidance for selecting the appropriate algorithm based on graph size and structure. Exact is reliable for datasets where exhaustive search is feasible, while CoreExact offers scalability and speed for real-world, large network datasets. Future research could investigate adaptive hybrid approaches that switch between Exact and

CoreExact based on online graph analysis metrics, optimizing both accuracy and efficiency dynamically.