

d2k: Scalable Community Detection in Massive Networks via Small-Diameter k-Plexes

Alessio Conte
National Institute of
Informatics, Japan
conte@nii.ac.jp

Tiziano De Matteis
University of Pisa
dematteis@di.unipi.it

Daniele De Sensi
University of Pisa
desensi@di.unipi.it

Roberto Grossi
University of Pisa
grossi@di.unipi.it

Andrea Marino
University of Pisa
marino@di.unipi.it

Luca Versari
University of Pisa
luca.versari@di.unipi.it

ABSTRACT

This paper studies k -plexes, which are arguably the most popular pseudo-clique model for network communities. In a k -plex, each node can miss at most $k - 1$ links. Our goal is to detect large communities in today's real-world graphs which can have hundreds of millions of edges. While many have tried, this task has been elusive so far due to its computationally challenging nature: k -plexes and other pseudo-cliques are harder to find and more numerous than cliques, a well known hard problem. We present d2k, which is the first algorithm able to find large k -plexes of very large graphs in just a few minutes. The good performance of our algorithm follows from a combination of graph-theoretical concepts, careful algorithm engineering and a high-performance implementation. In particular, we exploit the low degeneracy of real-world graphs, and the fact that large enough k -plexes have diameter 2. We validate a sequential and a parallel/distributed implementation of d2k on real graphs with up to half a billion edges.

KEYWORDS

k -plexes, graph enumeration, community discovery, parallel programming

ACM Reference Format:

Alessio Conte, Tiziano De Matteis, Daniele De Sensi, Roberto Grossi, Andrea Marino, and Luca Versari. 2018. d2k: Scalable Community Detection in Massive Networks via Small-Diameter k -Plexes. In *Proceedings of ACM SIGKDD conference (KDD'18)*. ACM, New York, NY, USA, Article 4, 9 pages.

1 INTRODUCTION

Finding communities and clusters is one of the most fundamental tasks when analyzing any form of data, albeit computationally demanding. In networks, communities are generally associated with densely interconnected subgraphs [13, 19, 21]: a clique represents the ideal situation, where nodes are pairwise linked, and is the earliest and arguably most studied community model. Many modern approaches exist to find cliques in networks, based on the seminal work by Bron and Kerbosch [5]. Great effort has been dedicated to

optimize various goals, such as worst-case running time on general graphs [17] and sparse, real-world, graphs [12], usage of main memory [7], and running time as a function of the output [10].

In real networks, where data can be noisy or faulty, large and closely linked communities hardly appear as ideal cliques, so other forms of subgraphs are sought for. A natural answer to this question is a relaxed notion of *pseudo-clique*, such as k -core, k -plex, n -clan, n -club, s -clique, dense subgraph [13, 15, 16, 21], but there is another side to this coin: the number of pseudo-cliques grows exponentially, at an even faster pace than that of cliques; moreover, the cost of detecting the former ones is higher due to their more complex structure. These issues are well known and investigated in the literature.

In this paper, we focus on k -plexes, arguably the most popular notion of pseudo-cliques [3, 9, 16, 22, 23]: the requirement of each node being linked to all others is loosened to each node missing no more than $k - 1$ links or, equivalently, missing k links including the one to itself (absence self-loops is assumed). Some examples are shown in Figure 1. A clique is a 1-plex according to this definition, assuming that there are no self-loops, and as k grows, the number of k -plexes increases exponentially with respect to that of cliques: Figure 2 testifies how striking this is even on small networks (see Table 1 for reference). To make k -plexes effective models of communities, we want to focus on interesting configurations that are larger and more densely connected by fulfilling three constraints (trivially satisfied by cliques when $k = 1$): the first two appeared in the literature, and the last one is introduced and motivated here.

- The k -plexes should be *connected*, otherwise they cannot be regarded as one community.
- Their size should be *large*, so as to involve as many entities as possible. In particular, they should be *maximal* under inclusion, namely, adding one more node violates the constraints given here.
- Their *diameter* (i.e. maximum pairwise distance among their nodes) should be *small*, as a node cannot go too far in reaching the nodes in its community.

In particular, we require that the diameter should be at most 2 for the following reasons: on one hand, it is equivalent to saying that every two nodes in the community are directly linked or at least share a common neighbor. Indeed, any k -plex with diameter greater than 2 has nodes that are not linked and do not share common neighbors: it could be argued that such nodes should not belong to the same community (see the examples in Figure 1). On the

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

KDD'18, August 2018, London, United Kingdom, UK

© 2018 Copyright held by the owner/author(s).

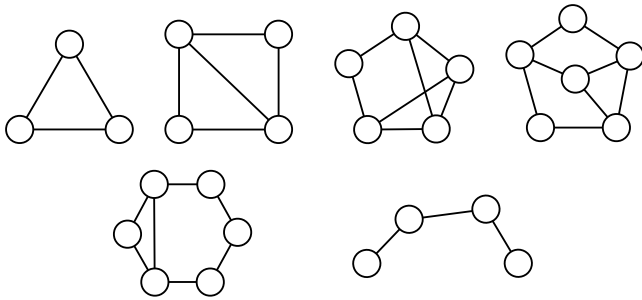


Figure 1: From left to right, top: a 1-, 2-, 3-, 4-plex. Bottom: a 4-plex and a 3-plex with diameter larger than 2.

other hand, if the diameter is less than 2, we obtain exactly a clique. Asking the diameter to be *at most 2* is the right balance between these extremes. Furthermore, this does not miss any significant communities (the meaning of significance is discussed in Section 6) in comparison with standard k -plexes. Specifically, *all* connected 1-plexes and 2-plexes have diameter at most 2, and the only connected 3-plex of diameter more than 2 is the induced path of 4 nodes, which can hardly be considered a community. For larger values ($k \geq 4$), we still guarantee to find at least all k -plexes with $2k - 1$ or more nodes, as it is well known that they have diameter 2 (e.g., see Lemma 5.1 in [9]). As k is a small constant in practice, this guarantees to hit all large (or even medium-sized) communities.

From a computational point of view, the fact that the diameter is at most 2 gives us a powerful handle to efficiently pull out the most interesting k -plexes from massive networks. A practical example can be seen in Section 2, where we briefly describe the biggest communities in some large real networks, whose detection only took few minutes with $D2K$.

Our main contribution is the algorithm $D2K$ (for diameter-2 k -plexes), which shows that significant k -plexes can be efficiently discovered in large networks. $D2K$ features a recursive backtracking structure, and a decomposition which is particularly suitable for parallelization and distribution. Furthermore, we introduce some novel pruning techniques for the problem, based on simple yet effective insights on both the problem and the structure of real data, which dramatically reduce the search space and consequently the computation time.

While computation of cliques on large and sparse real-world networks is feasible with ad-hoc algorithms[10, 12], computing k -plexes on such networks has been so far a task out of reach: [9] finds the largest 2-plex on a graph with 1.8 million edges, and [22] scales up to a graph with 22 million edges (see *pokec* in Table 1), but only processes a subgraph made by 10 of its nodes and the surrounding areas.

While these results greatly improve upon previous approaches (e.g., [3, 23]), it is not yet enough for today's data: large real world graphs can have hundreds of millions of nodes, or even billions. We argue that our approach can give a satisfying solution to the sought-after problem of finding pseudo-cliques in very large graphs. This is motivated in Section 5, where we show how our approach can process real world graphs with hundreds of millions of nodes

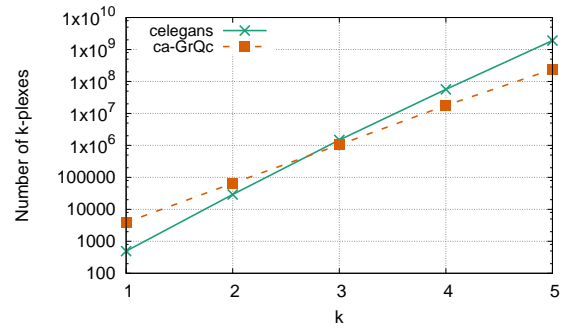


Figure 2: Number of k -plexes for increasing values of k .

in very reasonable times. Moreover, further discussion with related work is given in Section 6.

2 DATA ANALYSIS

For the sake of presentation, we computed the largest 4-plexes on three networks, whose content is briefly described below. Although performing data analysis is out of the scope of this work, we believe that the following toy examples can help to grasp its flavour.

pokec. This dataset is a snapshot of the most popular Slovakian social network. The largest 4-plex is a group of 32 users, aged between 15 and 20 years. Remarkably, the community is mainly Czech, with 81% of the users from Czech Republic (whose language, should be remarked, is closely related to Slovak), and is 94% female. Furthermore, data suggests it may have been a pre-existing group of friends who decided to join together, as 97% registered within 5 months of each other (January to May 2012). The users also seem to have similar music interests, with 56% selecting Slovak actress and singer *Lucia Molnárová* as favorite singer.

it-2004. This is a 2004 crawl of the Italian web (*.it* domain) made by [4]. The largest 4-plex consists in a collection of 3210 pages from the website *www.cuoko.it*. The reason behind such a structure is not evident from the data.

uk-2005. This is a 2005 crawl of the British web (*.co.uk* domain, but including *bbc.com*) again made by [4]. The largest 4-plex consists in a collection of 587 pages from a series of similarly named housing websites, plus the two domains *gibbinsrichards.co.uk* and *doorkeys.co.uk*. The following is an excerpt.

<i>gibbinsrichards.co.uk</i>	<i>homes-for-sale-oxford.co.uk</i>
<i>doorkeys.co.uk</i>	<i>homes-for-sale-paisley.co.uk</i>
...	...
<i>estate-agent-oldham.co.uk</i>	<i>house-for-sale-paisley.co.uk</i>
<i>estate-agent-oxbridge.co.uk</i>	<i>house-for-sale-perth.co.uk</i>
...	...
<i>home-for-sale-oxbridge.co.uk</i>	<i>houses-for-sale-perth.co.uk</i>
<i>home-for-sale-oxford.co.uk</i>	<i>houses-for-sale-peterborough.co.uk</i>
...	...

Among these pages, only *gibbinsrichards.co.uk* is currently active and indeed corresponds to a housing company. While a *whois* lookup does not reveal any information on the similarly named websites, one may conjecture that they were bogus websites set up for the benefit of the first two domains, e.g., for improving their rank scores on well known search engines.

3 ALGORITHM

In this section we show our proposed algorithm for listing maximal k -plexes of diameter at most 2, on a graph $G = (V, E)$.

First we will explain the core structure of our approach, then some crucial algorithm engineering and optimization which allow us to reduce the search space.

3.1 Main structure

Let v_1, \dots, v_n be any ordering of the nodes, from small to large. For a node v_i , let $N(v_i)$ be the set of nodes adjacent to v_i (its neighbors), and $N_{>}(v_i)$ be the set of *forward neighbors* of v_i , i.e., $N(v_i) \cap v_{i+1}, \dots, v_n$. Furthermore, let $N_{>}^2(v_i)$ be the set of *forward cousins* of v_i , that are neighbors of a node in $N_{>}(v_i)$ which are also larger than v_i (and not already forward neighbors of v_i).

We decompose the original problem into enumerating, for each node v , all k -plexes whose smallest node is v itself. To do so, we will exploit the following key property:

OBSERVATION 3.1. *When v_i is the smallest node in a diameter two k -plex, all other nodes of the k -plex must be either forward neighbors or forward cousins of v_i .*

This allows us to look for solutions with smallest node v in the subgraph containing just v , its forward neighbors, and its forward cousins, which we call $\text{BLOCK}(v)$. To improve the performance of our approach, we want $\text{BLOCK}(v)$ to be as small as possible, as will be discussed later.

In order to process each $\text{BLOCK}(v)$, we use a binary partition scheme similar to the Bron-Kerbosch clique enumeration algorithm [5]. We want to recursively solve the following subproblem.

PROBLEM 1. *Given a connected k -plex K and a set of nodes EXCL , find all maximal connected k -plexes K' such that $K \subseteq K'$ and $K' \cap \text{EXCL} = \emptyset$.*

Solving Problem 1 recursively is easy: for a given $v \notin K \cup \text{EXCL}$ such that $K \cup \{v\}$ is a connected k -plex, we split the problem into finding all solutions containing v by adding v to K , and all those who do not by adding it to EXCL instead, using two nested recursive calls. For efficiency, for each k -plex K we keep a set CAND of all v that should be tested, and directly generate all nested calls that would add one node of CAND to K and the previously considered one to EXCL as children of the same recursive call.

Whenever CAND is empty, it is easy to see that K is a maximal connected k -plex if and only if EXCL is also empty, since any node left in EXCL may be added to K to make a larger k -plex, and if both sets are empty no node can be used to enlarge K .

Finally, for any k -plex maximal in $\text{BLOCK}(v)$, our algorithm checks whether it is maximal in G , in which case it is output, or it is a subset of a larger k -plex whose smallest node is some $w < v$, in which case it is discarded as it is not a solution.

The resulting algorithm is as follows:

For each child, CAND and EXCL are updated using the function $\text{update}()$ to retain just the nodes x such that $K \cup \{x\}$ is still a connected k -plex.

To speed up the $\text{update}()$, we can compute CAND' and EXCL' by difference from the current ones: the addition of c to K may only shrink EXCL . As for CAND , we may have to add some nodes which became connected to K , namely, neighbors x of c which had no

Algorithm 1: Structure of our algorithm

```

Input : A graph  $G = (V(G), E(G))$ , an integer  $k$ 
Output: All diameter 2  $k$ -plexes of  $G$ 
1 foreach  $v$  in an ordering  $\{v_1, \dots, v_n\}$  of  $V(G)$  do
2    $H = (V(H), E(H)) \leftarrow \text{BLOCK}(v)$ 
3    $\text{enum}(H, \{v\}, \emptyset)$ 
4 Function  $\text{enum}(H, K, \text{CAND}, \text{EXCL})$ 
5   // if  $K$  is maximal in  $H$ 
6   if  $\text{CAND} \cup \text{EXCL} = \emptyset$  then
7     if  $K$  is maximal in  $G$  then
8       output  $K$ 
9     foreach  $c \in \text{CAND}$  do
10       $\text{CAND}', \text{EXCL}' \leftarrow \text{update}(H, K, c, \text{EXCL})$ 
11       $\text{enum}(H, K \cup \{c\}, \text{CAND}', \text{EXCL}')$ 
12       $\text{CAND} \leftarrow \text{CAND} \setminus \{c\}$ 
13       $\text{EXCL} \leftarrow \text{EXCL} \cup \{c\}$ 
14 Function  $\text{update}(H, K, c, \text{EXCL})$ 
15  $\text{CAND}' \leftarrow \{v \in V(H) \setminus \text{EXCL} : v \neq c \text{ and } K \cup \{c\} \cup \{v\} \text{ is a}$ 
    $\text{connected } k\text{-plex}\}$ 
    $\text{EXCL}' \leftarrow \{v \in \text{EXCL} : K \cup \{c\} \cup \{v\} \text{ is a connected } k\text{-plex}\}$ 

```

neighbor in K . Whenever $|K| \geq k$, however, our algorithm can skip this check, as any such x would have at least k non-neighbors in K , meaning that $K \cup \{x\} \cup \{c\}$ would not be a k -plex.

A strong point of this algorithm is the suitability for parallel and distributed computation, since each $\text{BLOCK}(v)$ can be processed independently. Furthermore, we will now present how the computation on each block can be further optimized.

3.2 Node ordering

The first and most surprising cut is simply ordering $V(G)$ in a *degeneracy ordering* [12].

If we call Δ the maximum degree of a node in G , note that the number of nodes in $\text{BLOCK}(v)$ is bounded by Δ^2 . However, a degeneracy ordering minimizes the maximum number of *forward* neighbors of a node in G : this number is called the *degeneracy* of the graph, d . Using such ordering reduces the size of $\text{BLOCK}(v)$ from Δ^2 to $d \cdot \Delta$.

Extensive experimental evidence has shown the degeneracy to be a small number on most real-world networks, even large ones [10–12], and can be noticed in Table 1. The practical effect of this ordering is striking. Horizontal lines *baseline* and d in Figure 3 show the maximum block sizes produced by our algorithm, using respectively a *random* or a *degeneracy* node ordering: the latter ordering produced blocks always less than half the size of the former, and up to a factor 7. (Lines sp and $sp + d$ will be covered later in Section 3.4)

3.3 Pivoting generalization

The Bron-Kerbosch algorithm for clique enumeration uses an effective technique called *pivoting* to cut useless branches from the computational trees. This is based on the principle that “each maximal clique must contain either the node u or a non-neighbor of

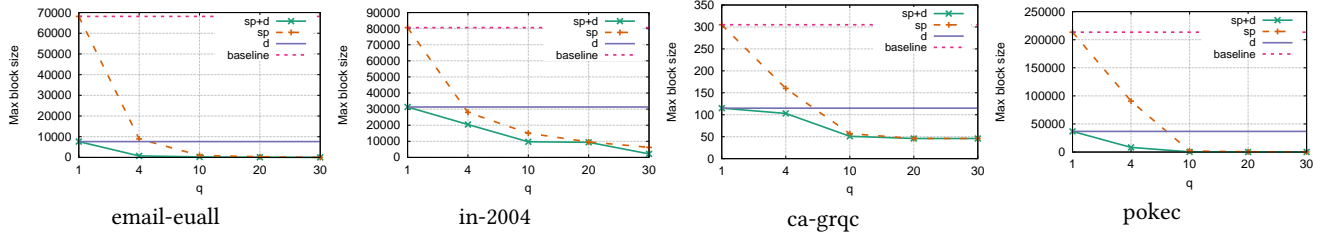


Figure 3: Effectiveness of the pruning techniques described in Sections 3.2 and 3.4 on the size of the generated BLOCKS

u . In this work, we present of a generalization of this principle applied to k -plexes:

LEMMA 3.2 (k -PLEX PIVOTING). *Let K be a k -plex, $CAND = \{v \notin K : K \cup \{v\} \text{ is a } k\text{-plex}\}$, and u a node in $CAND$. Any maximal k -plex containing K contains either u , a non-neighbor of u , or a neighbor w of u such that v and u have a common non-neighbor in K .*

PROOF. Assume by contradiction that K' is a maximal k -plex that violates the above constraint: u has at most $k-1$ non-neighbors in K' since $K \cup \{u\}$ is a k -plex and $K' \setminus K$ is made by only neighbors of u . Furthermore, any non-neighbor w of u in K is a neighbor of all nodes in $K' \setminus K$, thus w may not have more than $k-1$ non-neighbors in K' as it did not in K . Thus $K' \cup \{u\}$ is a k -plex and K' is not maximal, a contradiction. \square

We say that this is a generalization of the pivoting for clique enumeration as for $k=1$ we obtain exactly the pivoting of the Bron-Kerbosch algorithm. Furthermore, as for the Bron-Kerbosch algorithm, the same applies if u is chosen in EXCL rather than CAND. The cut of the search space is in that, when considering K , we can skip the recursive call for nodes which are both neighbors of u and do not share any non-neighbors in K with u . In practice, this translates to replacing Line 8 in Algorithm 1 with the following

```
foreach { $c \in CAND : c \notin N(u)$  or  $K \setminus (N(u) \cup N(c)) \neq \emptyset$ } do
```

where u is the chosen pivot. In order to maximize the effectiveness of this cut, we adopt the philosophy of Tomita et al. [17], and pick at each step the u which maximizes the number of prevented recursive calls.

3.4 Soluzion size pruning

Finally, we show some cuts that leverage our interest in k -plexes with minimum size q .

The first and most obvious cut follows from the fact that all the k -plexes generated from a given recursive call will be a subset of $K \cup CAND$. This means that, whenever $|K \cup CAND|$ is less than q , we can cut the search as no interesting solutions will be produced.

Another less obvious, yet essential cut, is obtained from the following lemma:

LEMMA 3.3 (SIZE PRUNING). *Any two nodes u and v in a k -plex K of size q , have at least $q - 2k + 2$ common neighbors in K .¹*

¹If $u \in N(v)$, we consider u as common neighbor of u and v

PROOF. As each node has at most $k-1$ non-neighbors in K , the number of node in K that are not neighbors of at least one of u and v is at most $2k-2$. From this follows that there are at least $q - 2k + 2$ nodes that have both of them as neighbors. \square

As when processing $BLOCK(v)$ we are only interested in k -plexes of size at least q containing v , this means we can immediately (and recursively) remove from $BLOCK(v)$ any node that does not share $q - 2k + 2$ neighbors with v .

The effect of this cut is not just heuristic: we give no proof for space reasons, but simple calculations show that this reduces the maximum size of $BLOCK(v)$ from $d \cdot \Delta$ to $d \cdot \Delta / (q - 2k + 2)$.

In practice, as q grows, size pruning may reduce the size of the subgraphs processed by up to orders of magnitude, as shown in Figure 3: line sp shows the size maximum size of a block generated on the corresponding graph with the specified q , using the size pruning. Compared to the *baseline* (i.e., no optimization) one can see how the difference is already important for $q=4$, and becomes even larger for increasing q , so much so that for $q=30$ on most graph there is hardly anything left to process.

Finally, line $sp+d$ in Figure 3 represent the maximum block size generated by $d2k$, i.e., using both the degeneracy ordering and the size pruning: one can see how this technique takes the best of both reduction, obtaining much smaller sizes than the baseline for all values of q . Notably, the effectiveness of the cuts seems to be independent from the size of the graph, but strongly influenced by the degeneracy (relatively high on *ca-grqc* and *in-2004*, and small on *email-euall* and *pokec*).

This allows us to process even huge graphs in a short time, provided a large enough q is chosen (see Section 5).

4 PARALLELIZATION

Concerning the parallelization of the algorithm, we first describe how we parallelized it for the execution on a single shared memory multi-core machine. Then, we extend our parallelization to support a cluster of multi-core computing machines with a distributed memory.

4.1 Shared Memory, Single Machine

The parallel implementation of $d2k$ is characterized by the presence of multiple threads, each one performing part of the processing. For doing this, we resort to the observation that each $BLOCK(v)$ can be processed independently from the others (see Section 3.1). Accordingly, the first step of our solution is to keep all the nodes v in a queue shared among the threads (in the following *vqueue*).

Each thread extracts a node from the queue and begins to process the BLOCK associated with that node. When a thread terminates the processing of a BLOCK, a new node is extracted from $vqueue$, until there are no more nodes (and thus BLOCKS to be processed). Since the time required to process different BLOCKS may be different, with this solution we could experience some load unbalancing effects. In particular, it could happen that most threads have terminated the available BLOCKS while few threads are still processing some heavier BLOCKS. To mitigate this effect, we use a *work requesting* approach. If a thread detects that there are no more BLOCKS to be processed but some thread is still processing something, it will require the offloading of additional work to improve the load balancing and reduce the execution time.

The offloading is performed by means of a globally shared queue $wqueue$. We keep an atomic counter indicating how many threads have finished all the BLOCKS and are available for additional work. During the processing of a BLOCK, each thread periodically checks if the size of $wqueue$ is lower than the value of the counter. If this is the case, the thread will put a part of the BLOCK in the queue, instead of processing it directly. The requesting worker will then be able to dequeue the sub-BLOCK from $wqueue$ and to process it. To implement this, D2K can generate an object representing a nested recursive call (i.e., the sub-BLOCK), then hand it over to the requesting worker, who will process it, and skip directly to the next recursive call.

4.2 Distributed Memory, Multiple Machines

We adopt a similar approach for a cluster. Since we are in a distributed memory environment, we adopt a *master-worker* solution. One computing machine will act as the *master*, dealing with the distribution of the BLOCKS to be processed to a set of *workers*, each one executed on a different multi-core computing machine.

When the computation starts, the master sends a chunk of nodes to each worker. A generic worker, after receiving a chunk, will insert all the nodes in its shared queue, leveraging on the parallelization scheme used for the single computing machine to process all the assigned nodes. Differently from the single machine case, when there are no more BLOCKS to be processed, instead of terminating the execution, the worker will ask the master for a new chunk.

When the master receives a request for a new BLOCK from a worker $wreq$ and there are no more BLOCKS to assign, the master tries to balance the work by requesting a BLOCK (or a part of it) from an overloaded worker $wvictim$ and by redirecting it to $wreq$. Each worker periodically checks for pending stealing requests from the master. When this is the case, and there are other nodes in the current chunk which were not yet processed, the worker sends some of those nodes to the master; otherwise (i.e., it is processing the last node of the assigned chunk), it sends a sub-BLOCK to the master. In both cases, the master will redirect the work received from $wvictim$ to $wreq$. This process continues until all the workers have nothing left compute.

To select $wvictim$ we adopt a simple heuristic. We store the timestamp of the last chunk request received by each worker. When selecting a victim for the stealing, we pick the worker characterized by the lower timestamp, since the worker who has not requested

data for the longest time is likely to be the most loaded worker, i.e., the most suitable victim for stealing.

5 EXPERIMENTS

The evaluations have been performed on a cluster of 16 homogeneous machines. Each machine is a dual CPU Intel Xeon E5-2640 v4, Broadwell based architecture, composed of 20 cores operating at 2.40GHz. The HyperThreading feature was not used.

Each core has a private L1 (32KB) and L2 (256KB) cache, and each CPU is equipped with a shared L3 cache of 25MB. The machine has 128 GB of RAM and the machines are interconnected using a OmniPath network. The program is compiled with `gcc-6.4.0`, using the `-O3` optimization flag. Our code is publicly available at https://github.com/veluca93/parallel_enum. In the following discussion, the sequential and parallel evaluations are performed on a single machine, while the distributed tests use a different number of machines. For all tests, we considered 12 hours of execution time as a hard limit.

Dataset. Our graphs have been taken from SNAP (<http://snap.stanford.edu/>) and LAW (<http://law.di.unimi.it/>) and are reported in Table 1. We divided our networks in three classes depending on their size: in the first class, shown in Table 1(a), there are networks with less than about 50 thousands of edges, in the second one (b) there are the remaining networks with less than 10 millions of edges, while other graphs (c) are in the third class. As the size and the number of solutions of the networks change considerably, different settings are needed and the three classes of networks reflect the corresponding different performance evaluations. We use the former classes of networks to compare our approach with the existing ones to alleviate their computational cost setting $q = 4, 10, 20$ for the 2-plexes and $q = 10, 20$ for the 3-plexes in a sequential setting. We use the bigger networks to study the behaviour of our approach in parallel and distributed settings. Since these networks are much bigger, and hence they contain many k -plexes, the value of q changes from graph to graph. We will show next how to choose suitable values for q .

For each network we report some statistics and its type.² In particular, we report the number of 2-plexes and 3-plexes greater than q for different values of q . For some of these graphs, we know just a lower bound on the number of these k -plexes (the ones reporting \geq) since our computations exceeded 12 hours (which was the time limit used in our experiments). This happened in the cases where there was a huge number of solutions.

It is worth observing that, even if graphs have similar size in terms of nodes, their number of k -plexes can be very different. Sometimes this can be related to an higher density, and, it seems that a small increase on the average degree can correspond to a huge increase in terms of 2- and 3-plexes. The most striking example is the biological graph `interdom`, relatively small compared to the others but with average degree 92.6, that has millions of 2-plexes and billions of 3-plexes. On the other hand, the collaboration network `ca-grqc` has a similar number of nodes and a smaller average degree (about $1/17th$ w.r.t. that of `interdom`), and indeed the number of 2- and 3-plexes is orders of magnitude smaller. In some other cases,

²as: autonomous systems, collab: collaboration, bio: biological, prod: product co-purchasing, comm: communication, soc: social, web: web crawl.

NETWORK	TYPE	n	m	Δ	d	2-plexes			3-plexes	
						$q = 4$	$q = 10$	$q = 20$	$q = 10$	$q = 20$
caida	as	26 475	53 381	2 628	22	1 337 044	23 314	0	1 531 876	0
ca-grqc	collab	5 241	14 484	81	43	12 038	377	118	13 352	1 568
jazz	collab	198	2 742	100	29	26 172	8 059	2	257 233	2
celegans	bio	354	1 501	186	10	12 814	5	0	248	0
homo-sapiens	bio	1 027	1 166	443	17	929 795	576	0	20 301	0

(a) SMALL NETWORKS

NETWORK	TYPE	n	m	Δ	d	2-plexes			3-plexes	
						$q = 4$	$q = 10$	$q = 20$	$q = 10$	$q = 20$
interdom	bio	1 706	78 983	728	129	60 858 742	54 562 092	34 424 664	536 552 584 954	378 852 643 230
amazon0505	prod	410 236	2 439 436	2 760	10	8 996 909	22 483	0	248 433	0
email-euAll	comm	265 214	365 569	7 636	37	11 489 357	1 042 929	0	75 751 394	2 637
epinions1	soc	75 888	405 739	3 044	67	176 891 842	73 518 387	3 286 364	11 947 340 789	538 275 804
slashdot090221	soc	82 144	500 480	2 548	54	62 130 214	28 862 926	11 411 028	2 993 676 468	1 303 148 522
wiki-vote	soc	8 298	100 761	1 065	53	42 757 442	9 162 660	52	1 337 634 391	156 727

(b) MEDIUM NETWORKS

NETWORK	TYPE	n	m	Δ	d	2-plexes		3-plexes	
						$q = 50$	$q = 100$	$q = 50$	$q = 100$
as-skitter	as	1 696 415	11 095 298	35 455	111	47 969 775	0	21 070 497 438	0
in-2004	web	1 353 703	13 126 172	21 869	488	25 855 779	9 978 037	29 045 783 792	4 257 410 159
						$q = 20$	$q = 30$	$q = 20$	$q = 30$
pokec	soc	1 632 803	22 301 964	14 854	47	94 184	3	5 911 456	5
						$q = 250$	$q = 500$	$q = 250$	$q = 500$
uk-2005	web	36 022 725	549 996 068	1 372 171	584	106 243 475	256 406	$\geq 18 336 111 409$	28 199 814
						$q = 2000$	$q = 3000$	$q = 2000$	$q = 3000$
it-2004	web	35 051 939	473 649 719	1 243 927	3 209	2 727 030	6 304	$\geq 165 823 718$	2 722 875

(c) BIG NETWORKS

Table 1: Considered networks and their properties, including number of k -plexes of different sizes

the number of k -plexes is more deeply related to the topology of the networks. This is the case of `as-skitter` and `in-2004` which have roughly the same number of nodes and edges but the number of 2-plexes greater than 100 is very different: almost 10 millions for `in-2004` and 0 for `as-skitter`. Surprisingly, relatively small networks, like `interdom`, having this huge size of k -plexes seem to be not common, meaning that, despite the worst case upper bound on the number of solutions, in practice the number of k -plexes in real-word networks is relatively small for suitable values of q .

Choosing q . Due to the effectiveness of the pruning techniques, `D2K` finishes quickly when q is too high to find any result. This feature can be exploited to quickly find the largest k -plexes in a graph. The philosophy is similar to that used in [9], but with the important difference that we do not require listing all maximal cliques first, something which is trivial on small graphs but far from it on larger ones. For the larger graphs in Table 1(c), we thus adaptively found values of q corresponding to their larger communities.

5.1 Performance Evaluation

In this section we evaluate the performance of our method `D2K` with respect to the competitors `GP` [22] and `LP` [9] using the graphs in Table 1(a) and (b) when generating 2-plexes and 3-plexes of size at least q . All the methods were run in a sequential setting. Moreover, note that, as q is always greater than k^2 , all the k -plexes to be listed have diameter at most 2, which means that all the algorithms return the same set of solutions.

To give a general picture of the running times for the different values of q , we reported time performances of all the approaches in

Table 2 using the small graphs in Table 1(a). It is worth observing that `LP` and `GP` show their best behaviour in opposite settings: the former is running out of time (12 hours) for smaller values of q and improves when q grows, the latter works pretty well for smaller values of q and it gets worse growing q even if there are no solutions to list. Indeed, sometimes there could be no k -plexes whose size is greater than q (like for the values reporting * in Table 2). Looking at the number of k -plexes in Table 1(a), it seems that `GP` spends less time whenever there are many solutions and much more time whenever there are no or few solutions, while `LP` spends less time in the latter case. On top of the two methods, `D2K` spends less than one second or few seconds to conclude its computation for all the values of q and for both 2- and 3-plexes, always improving the competitors.

In order to quantify our improvement with respect to the state of art, we have used the medium sized graphs in Table 1(b). We report the time needed by `D2K` (indicated as T) and the speedup S of `D2K` with respect to the best running time achieved by both `LP` and `GP`. In particular, we divide the best running time of the competitors by T to obtain S . Table 3 reports T and S for both 2- and 3-plexes, for different values of q . OOT means that `D2K` ran out of the maximum time allowed (12 hours). The speedup S is set to \star whenever both our competitors ran out of the maximum time or memory allowed. For the great majority of the graphs the values of S is always greater than 20, meaning that we spent less than 1/20 of the running time of our competitors. In the case of the 2-plexes with $q = 10$ and $q = 20$, our improvement is even more evident, as `D2K` is much faster and often outperforms the competitors of orders of magnitude. In the case of 3-plexes with $q = 20$ and $q = 30$, for the great majority of the cases both our competitors ran out of time

NETWORK	2-plexes									3-plexes					
	$q = 4$			$q = 10$			$q = 20$			$q = 10$			$q = 20$		
	D2K	GP	LP	D2K	GP	LP	D2K	GP	LP	D2K	GP	LP	D2K	GP	LP
caida	18.29	449.89	OOT	0.72	357.35	OOT	*0.01	*321.49	*0.10	50.86	OOT	OOT	*0.01	*OOT	*OOT
ca-grqc	0.07	3.98	OOT	0.01	3.19	2 964.91	0.01	2.50	12.49	0.07	17.27	OOT	0.01	48.42	471.90
jazz	0.25	6.47	OOT	0.12	5.28	OOT	0.01	2.40	OOT	2.89	29.12	OOT	0.01	41.26	OOT
celegans	0.08	3.27	OOT	0.01	1.93	1 186.19	*0.03	*1.65	*0.04	0.01	48.20	OOT	*0.01	*50.48	*0.01
homo-sapiens	4.89	123.17	OOT	0.02	101.91	OOT	*0.01	*91.71	*0.10	1.27	OOT	OOT	*0.01	*24 383.34	*0.01

Table 2: Time comparison among D2K, GP, LP when generating all 2- and 3-plexes of size greater than q (time is in seconds) for the small graphs in Table 1(a). OOT means the process did not terminate within 12 hours, while * means there were no solutions.

NETWORK	2-plexes						3-plexes					
	$q = 4$		$q = 10$		$q = 20$		$q = 10$		$q = 20$		$q = 30$	
	T(s)	S	T(s)	S	T(s)	S	T(s)	S	T(s)	S	T(s)	S
interdom	1 032.47	16.24x	1 040.76	15.45x	989.02	14.39x	OOT	*	OOT	*	OOT	*
amazon0505	56.76	28.68x	0.33	4 183.13x	*0.01	124 200x	4.16	*	0.01	*	0.01	*
email-euAll	157.35	56.53x	31.46	228.09x	*0.65	9384x	2 995.51	*	76.96	*	2.99	*
epinions1	2 945.71	OOT	1 545.27	OOT	283.93	67.75x	OOT	*	OOT	*	3 720.38	*
slashdot090221	989.99	39.81x	408.99	88.74x	260.39	90.29x	OOT	*	30 278.48	*	3 313.04	*
wiki-vote	532.22	22.16x	262.72	27.01x	16.79	76.59x	OOT	*	2 763.42	*	341.17	*

Table 3: Time (τ) of D2K and speedup (s) D2K obtains with respect to GP for medium sized graphs. OOT: D2K ran out of time. *: GP ran out of time.

or memory while D2K finished in time. On the other hand, when setting $q = 10$, D2K and both the competitors went out of time. We will be able to deal with these graphs using our parallel version of our algorithm.

Setup time. Since D2K needs the degeneracy ordering of the nodes of the network to be processed, for the sake of completeness, we have reported the time needed to compute this ordering on the networks of our dataset. These times (in seconds) are reported in the following, showing that this setup time is completely negligible with respect to the time needed to list k -plexes (see Table 2 and Table 3).

NETWORKS	TIME	NETWORKS	TIME
ca-grqc	0.01	email-euAll	0.12
celegans	0.01	slashdot090221	0.13
jazz	0.01	amazon0505	0.73
caida	0.02	as-skitter	2.86
homo-sapiens	0.02	in-2004	1.95
interdom	0.02	pokec	5.45
wiki-vote	0.04	uk-2005	79.21
soc-epinions1	0.10	it-2004	94.04

5.2 Evaluation of the parallel and distributed implementations

For the parallel evaluation of D2K, we executed the program on a single machine of the cluster with a variable number of threads, ranging from 1 to 20, the number of cores available in the machine.

To evaluate the effectiveness of the parallelization, we consider the speedup obtained with respect to the sequential version, i.e. the ratio between the execution time of the sequential implementation over the execution time of the parallel version with a given number of threads. Figure 4 shows three notable scenarios with the related speedups. As it can be noticed, in both situations (i.e. 2-plexes and 3-plexes) we are able to obtain good speedups, close to the ideal one.

Table 4 reports the parallel execution time obtained with 20 threads and the corresponding speedup with respect to the sequential implementation for the graphs in Table 1(b). We obtained good

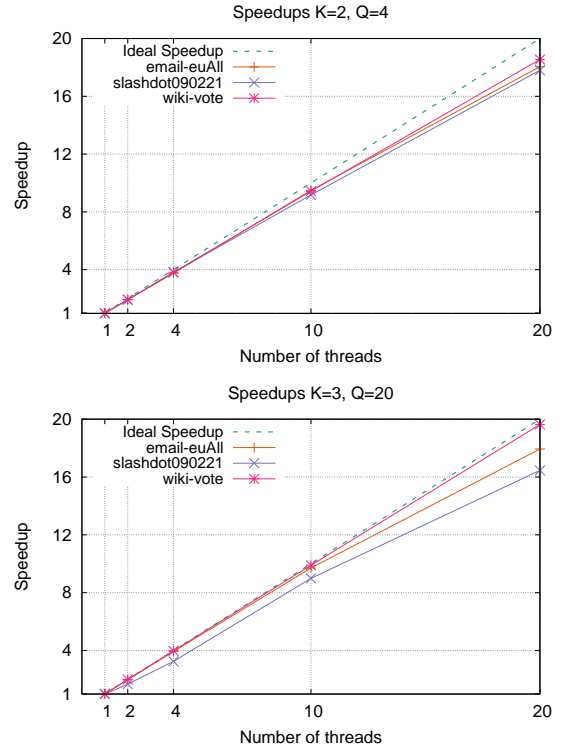


Figure 4: Speedup over the sequential version

values of speedup for all the situations in which the sequential running time is not negligible (i.e. more than 1 seconds). Reasonably, if this is not the case, the overhead of the parallelization reduces (or eliminates) the benefits of the parallel execution.

In general, the parallel implementation gives an almost linear (in the number of used threads) reduction of the execution time. This allowed us to complete the execution within the deadline of a

NETWORK	2-plexes						3-plexes					
	$q = 4$		$q = 10$		$q = 20$		$q = 10$		$q = 20$		$q = 30$	
	T (s)	S	T (s)	S	T (s)	S	T (s)	S	T (s)	S	T (s)	S
interdom	79.91	2.92x	61.34	16.97x	57.23	17.28x	OOT	★	OOT	★	OOT	★
amazon0505	3.57	15.90x	0.09	3.81x	0.01	1.00x	0.347	11.98x	*0.063	0.21x	*0.05	0.17x
email-euAll	8.72	18.05x	1.84	17.14x	0.10	6.33x	164.47	18.88x	4.29	17.9x	*0.04	0.16x
epinions1	165.82	17.76x	90.57	17.06x	15.72	18.07x	18496.36	★	2869.65	★	186.02	19.3x
slashdot090221	55.63	17.80x	24.06	17.00x	15.61	16.68x	2800.90	★	1513.92	16.46x	165.65	18.7x
wiki-vote	28.68	18.55x	14.04	18.72x	0.91	18.37x	2224.84	★	140.79	19.63x	0.01	4.3x

Table 4: Parallel execution time and speedup with respect to sequential version. ★: sequential algorithm ran out of time.

NETWORK	2-plexes		3-plexes	
	q	T (s)	q	T (s)
as-skitter	50	3531.87	50	OOT
	100	0.24	100	0.25
in-2004	50	43.51	50	OOT
	100	18.209	100	9720.47
pokec	30	0.63	30	0.82
	50	*0.48	50	*0.6
uk-2005	250	1567.26	250	OOT
	500	16.14	500	190.6
it-2004	2000	125.53	2000	OOT
	3000	18.59	3000	407.49

Table 5: Parallel execution time (s) over big graphs with 20 threads.

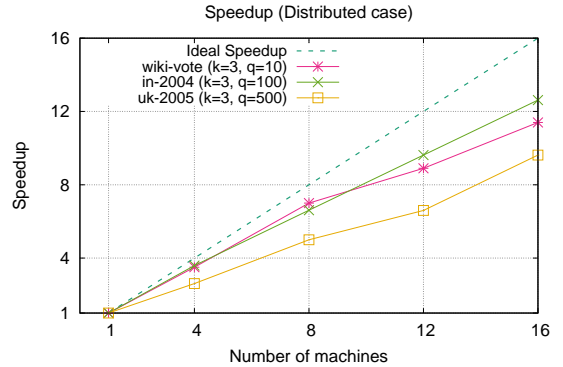


Figure 5: Speedup over the single computing machine

higher number of scenarios with respect to the case of the sequential implementation (i.e. all the 2-plexes cases and all the 3-plexes with the exception of interdom). Achieving such a good speedup has been possible thanks to a careful design of the sequential algorithm. For example, due to the efficient pruning performed on the BLOCKS, the size of the data to keep in the processors' caches is reduced, decreasing contention effects on the last level caches and improving the scalability of the parallel implementation.

Furthermore, the parallel implementation enables the enumeration of meaningful k -plexes on the biggest graphs of Table 1. Table 5 shows the obtained execution times.

Concerning the support for distributed memory clusters of multicore machines, we show in Figure 5 the speedup we achieve for different numbers of machines with respect to a single computing machine. In this case the speedup is defined as the execution time when using one machine (with 20 cores), divided by the execution time when using a given number of machines (each one with 20 cores). For this test we selected some of the biggest graphs, since it would not be very meaningful to use a cluster for problems which takes few seconds to be computed on a single parallel machine. As depicted by the Figure, the speedup differs according to the graph.

In some cases (e.g. wiki-vote) by using 16 multicore machines we achieve a speedup of 12 with respect to the single machine. Considering that on the single machine we had a speedup of 19 with respect to the sequential version, we reduced the execution time by 230 times with respect to the sequential implementation.

An interesting study case involves the computation of 3-plexes with $q = 50$ for in-2004. While a parallel implementation on a single machine was not able to terminate in 12 hours, our distributed implementation using 16 machines was able to find all the k -plexes in less than 3 hours. As the experiments suggest a close to ideal scaling factor, we may conjecture that a sequential implementation

of the algorithm would have required several weeks to find the same result.

6 RELATED WORK

The papers by Bron and Kerbosch [6] and Tsukiyama et al. [20] are currently at the heart of many algorithms for enumerating cliques and their relaxed versions such as k -plexes due to their effectiveness. The Bron-Kerbosch algorithm relies on a backtracking scheme that is adopted in several efficient algorithms due to its simplicity and good practical performance [12, 18, 24]. While the original version [6] does not provide any guarantee, the version in [18] guarantees a total running time of $O(3^{n/3})$, which worst-case-optimal, and the one in [12] further improves the work for sparse graphs, which may have up to $(n-d)3^{d/3}$ maximal cliques, by producing an algorithm with $O(d(n-d)3^{d/3})$ time. This strategy has been adapted to the enumeration of maximal k -plexes in [23], and inspired the similar backtracking structure in [22].

The algorithm by Tsukiyama et al. has been originally conceived for the equivalent problem of enumerating maximal independent sets, and has been subsequently adapted to maximal clique enumeration by Chiba-Nishizeki [8]. Makino-Ueno [14] has reinterpreted some of its ideas in the paradigm of reverse search introduced by Avis and Fukuda [1]. A space efficient algorithm with bounded delay has been given in [10]. A simplified version of this strategy has been applied to the enumeration of maximal k -plexes in [3], which proposed the first output-sensitive (i.e., whose time is proportional to the number of solutions found) k -plex enumeration algorithm, and then further exploited by [9] for large k -plexes.

As k -plexes generalize cliques, their discover takes more time than finding cliques. For this reason several algorithms for listing

all k -plexes [3, 23] do not scale well with large graphs. For example, [3] needs to maintain a large database of all maximal k -plexes found so far to avoid duplicates, which can create contention and excessive space usage. A new trend has emerged in recent years to simultaneously speed up the computation and avoid finding solutions of little interest for the final user: that is looking for a reduced number of just *significant* results. For cliques, an example is the work by Zhou et al. for finding the top- k cliques [26]. Zhai et al. [25], rather than aiming for large sized solutions, define a new pseudo-clique structure, similar to the k -plex, where small solutions are allowed as long as they are dense enough. Behar and Cohen [2] also aim at finding large pseudo cliques using a different model, that is connected s -cliques. As for k -plexes, examples are [22] and [9]. Both algorithms include a *minimum size* parameter which cuts off small solutions from the search, with hard-coded lower bounds in their structure to improve their efficiency: [9] finds k -plexes at least larger than k^2 , and after a careful reading of the paper (and as confirmed by the authors) it can be seen that [22] finds k -plexes at least larger than $k + 1$.

7 CONCLUSION

We proposed $\mathfrak{D}2\mathfrak{K}$, the first algorithm that can find k -plexes (and to the best of our knowledge, pseudo-cliques) of very large real world networks, by an effective combination of algorithm design and insight on the problem. To maximize this benefit, we proposed a parallel and distributed implementation which scales up nicely to tens of machines, with tens of cores each, further extending the reach of $\mathfrak{D}2\mathfrak{K}$ to networks whose sequential processing time would otherwise be massive. This allowed to compute for the first time k -plexes in large real-world networks.

$\mathfrak{D}2\mathfrak{K}$ moves an important chunk of today's real data within the reach of community detection, improving the applicability of existing network analysis methods. As a toy example, we provide some insight on the largest communities in two networks with half a billion edges each, and further validate our algorithm on various real networks of different types.

We believe that this work, paired with its open source implementation, sets a milestone for dense subgraph enumeration, and that it will open new directions for community discovery in very large graphs.

ACKNOWLEDGEMENTS

We wish to thank authors of [3, 9, 22] for providing their source code, and prof. Marco Danelutto and the people at *ITC Center* of the University of Pisa for providing the machines used in the experiments. The work was partially supported by JST CREST, Grant Number JPMJCR1401, Japan.

REFERENCES

- [1] David Avis and Komei Fukuda. Reverse search for enumeration. *Discrete Applied Mathematics*, 65(1-3):21–46, 1996.
- [2] Rachel Behar and Sara Cohen. Finding all maximal connected s -cliques in social networks. In *Proceedings of the 21th International Conference on Extending Database Technology, EDBT 2018. (to appear)*, pages –, 2018.
- [3] Devora Berlowitz, Sara Cohen, and Benny Kimelfeld. Efficient enumeration of maximal k -plexes. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, SIGMOD '15*, pages 431–444, New York, NY, USA, 2015. ACM.
- [4] Paolo Boldi, Bruno Codenotti, Massimo Santini, and Sebastiano Vigna. Ubcrawler: A scalable fully distributed web crawler. *Software: Practice & Experience*, 34(8):711–726, 2004.
- [5] Coen Bron and Joep Kerbosch. Algorithm 457: finding all cliques of an undirected graph. *Communications of the ACM*, 16(9):575–577, 1973.
- [6] Coenraad Bron and Joep Kerbosch. Finding all cliques of an undirected graph (algorithm 457). *Commun. ACM*, 16(9):575–576, 1973.
- [7] James Cheng, Linhong Zhu, Yiping Ke, and Shumo Chu. Fast algorithms for maximal clique enumeration with limited memory. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD*, pages 1240–1248, 2012.
- [8] Norishige Chiba and Takao Nishizeki. Arboricity and subgraph listing algorithms. *SIAM Journal on Computing*, 14(1):210–223, 1985.
- [9] Alessio Conte, Donatella Firmani, Caterina Mordente, Maurizio Patrignani, and Riccardo Torlone. Fast enumeration of large k -plexes. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017*, pages 115–124, 2017.
- [10] Alessio Conte, Roberto Grossi, Andrea Marino, and Luca Versari. Sublinear-space bounded-delay enumeration for massive network analytics: Maximal cliques. In *43rd International Colloquium on Automata, Languages, and Programming, ICALP 2016, July 11-15, 2016, Rome, Italy*, pages 148:1–148:15, 2016.
- [11] Alessio Conte, Roberto De Virgilio, Antonio Maccioni, Maurizio Patrignani, and Riccardo Torlone. Finding all maximal cliques in very large social networks. In *Proceedings of the 19th International Conference on Extending Database Technology, EDBT 2016, Bordeaux, France, March 15-16, 2016, Bordeaux, France, March 15-16, 2016*, pages 173–184, 2016.
- [12] David Eppstein, Maarten Löffler, and Darren Strash. Listing all maximal cliques in large sparse real-world graphs. *ACM J of Experimental Algorithmics*, 18, 2013.
- [13] Santo Fortunato. Community detection in graphs. *Phys rep*, 486(3):75–174, 2010.
- [14] Kazuhisa Makino and Takeaki Uno. New algorithms for enumerating all maximal cliques. In *Algorithm Theory-SWAT 2004*, pages 260–272. Springer, 2004.
- [15] Robert J. Mokken. Cliques, clubs and clans. *Quality and quantity*, 13(2):161–173, 1979.
- [16] Stephen B Seidman and Brian L Foster. A graph-theoretic generalization of the clique concept. *Journal of Mathematical sociology*, 6(1):139–154, 1978.
- [17] Etsuji Tomita, Akira Tanaka, and Haruhisa Takahashi. The worst-case time complexity for generating all maximal cliques and computational experiments. *Theoretical Computer Science*, 363(1):28–42, 2006.
- [18] Etsuji Tomita, Akira Tanaka, and Haruhisa Takahashi. The worst-case time complexity for generating all maximal cliques and computational experiments. *Theor. Comput. Sci.*, 363(1):28–42, 2006.
- [19] Charalampos E. Tsourakakis, Francesco Bonchi, Aristides Gionis, Francesco Gullo, and Maria A. Tsiarli. Denser than the densest subgraph: extracting optimal quasi-cliques with quality guarantees. In *The 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2013, Chicago, IL, USA, August 11-14, 2013*, pages 104–112, 2013.
- [20] Shuji Tsukiyama, Mikio Ide, Hiromu Ariyoshi, and Isao Shirakawa. A new algorithm for generating all the maximal independent sets. *SIAM J. Comput.*, 6(3):505–517, 1977.
- [21] Takeaki Uno. An efficient algorithm for solving pseudo clique enumeration problem. *Algorithmica*, 56(1):3–16, 2008.
- [22] Zhuo Wang, Qun Chen, Boyi Hou, Bo Suo, Zhanhui Li, Wei Pan, and Zachary G. Ives. Parallelizing maximal clique and k -plex enumeration over graph data. *Journal of Parallel and Distributed Computing*, 106:79–91, 2017.
- [23] Bin Wu and Xin Pei. A parallel algorithm for enumerating all the maximal k -plexes. In Takashi Washio, Zhi-Hua Zhou, Joshua Zhexue Huang, Xiaohua Hu, Jinyan Li, Chao Xie, Jieyue He, Deqing Zou, Kuan-Ching Li, and Mário M. Freire, editors, *Emerging Technologies in Knowledge Discovery and Data Mining*, pages 476–483, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg.
- [24] Yanyan Xu, James Cheng, Ada Wai-Chee Fu, and Yingyi Bu. Distributed maximal clique computation. In *Big Data (BigData Congress), 2014 IEEE International Congress on*, pages 160–167. IEEE, 2014.
- [25] Hongjie Zhai, Makoto Haraguchi, Yoshiaki Okubo, and Etsuji Tomita. A fast and complete algorithm for enumerating pseudo-cliques in large graphs. *International Journal of Data Science and Analytics*, 2(3):145–158, Dec 2016.
- [26] Zhaonian Zou, Jianzhong Li, Hong Gao, and Shuo Zhang. Finding top- k maximal cliques in an uncertain graph. In *Data Engineering (ICDE), 2010 IEEE 26th International Conference on*, pages 649–652. IEEE, 2010.