Master's thesis

# A NEW ALGORITHMIC FRAMEWORK FOR ENUMERATING COMMUTABLE SET PROPERTIES

Supervisor:                                                   Author:
Prof. *Roberto Grossi*                                        *Luca Versari*

ACCADEMIC YEAR 2016–2017

**Abstract**

This thesis considers a new algorithmic framework for listing maximal sets satisfying a given property (e.g. being a clique, a cut, a cycle, etc.), which fall within the general framework of set systems. A set system $\mathcal{F}$ over a ground set $U$ (e.g. the network nodes) is a collection of subsets of $U$ for which there exists some function that checks if an arbitrary subset of $U$ belongs to $\mathcal{F}$. For all *maximal* subsets in $\mathcal{F}$ under inclusion to be listed, the ambitious goal is to cover a large class of set systems while preserving the efficiency of their enumeration algorithms at the same time. The best-known ones list the maximal subsets in time proportional to their number but may require exponential space. This thesis improves the state of the art in two directions by introducing an algorithmic framework that, under suitable conditions, simultaneously (i) extends the class that can be solved efficiently to *commutable set systems*, and (ii) reduces the additional space usage from exponential in $|E|$ to *stateless*, thus accounting for just $O(q)$ space, where $q \leqslant |E|$ is the largest size of a maximal set.

B

# Contents

# Chapter 1

# Introduction

Algorithms for graph listing have a long history and even if they were born in the 70s in the context of enumerative combinatorics and computational complexity [17, 26, 32, 38, 44], the interest has quickly broadened to a variety of other communities in computer science and not, massively involving algorithm design techniques.

In network analysis discovering special communities corresponds to finding all the subgraphs with a given property [1, 15, 27, 29, 30, 42]. In bioinformatics, listing all the solutions is desirable, as single or few solutions may not be meaningful due to noise of the data, noise of the models, or unclear objectives [9, 21, 25, 28, 40]. In graph databases, graph structures are used for semantic queries, where nodes, edges, and properties are used to represent and store data; retrieving information corresponds to find all the suitable subgraphs in these databases [2, 11, 46]. When dealing with incomplete information, it may be impossible to completely satisfy a query. Subgraph listing algorithms can find answers that *maximally* satisfy a partial query; for instance, there is a one-to-one correspondence between the results of a join or full disjunction query and certain subgraphs of the *assignment graph*, a special graph obtained by combining the relational database with the given query. Moreover, the kind of subgraphs to look for depend not only on the database, but also on the query [11].

In this scenario, graph enumeration has left the theoretical border [24] to meet more stringent requirements: not only a given listing problem must fit a given class of complexity, but its algorithms must be efficient also in real-world applications. Algorithm design has made a big effort to generalize the graph properties to be enumerated and to unify the corresponding approaches [3, 8, 10, 26, 43]. These generalizations allow the same algorithm to solve many different problems.

This thesis presents some of the work that has been done so far and new results that fit into this line of research: on one side, we want to obtain efficient listing algorithms able to deal with large networks; on the other side, we aim at designing an algorithmic framework which solves simultaneously many problems and leave the designer in charge of few core tasks depending on the specific application. In particular, we focus

on efficient enumeration algorithms for maximal subgraphs satisfying a given property (e.g. being a clique, a cut, a cycle, a matching, etc.), as they fall within the general framework of set systems [10, 26].

This thesis is organized as follows: Chapter 1 contains an overview of the thesis and some well-known definitions necessary to develop the content. Chapter 2 contains an introduction to set systems, that gives definitions, examples, known lower bounds and known algorithm for the enumeration of its maximal sets. Chapter 3 explains how the new algorithmic framework works, and Chapter 4 gives some applications of the framework to enumeration problems.

The original work in this thesis is mainly contained in Chapter 2 for the introduction of commutable set systems, Chapter 3 for the new algorithmic framework and Chapter 4 for its applications, except for Subsection 3.2.1 that gives an improved overview over already-existing work. Moreover, parts of the thesis build upon further work that was published in [12, 13].

## 1.1   Preliminaries

As graphs play a vital role in this thesis, we start by stating some standard definitions: an *undirected graph* $G$ is a pair $(V, E)$ with $E \subset V \times V$. The elements of $V$ are called *nodes* of the graph, while the elements of $E$ are called *edges*. The edge set of $G$ must be symmetrical, i.e. if $(a, b) \in E$ then $(b, a) \in E$, and must not contain self-loops (i.e. $(a, a) \notin E$ for every $a$ in $V$). We say that a node $v$ is a *neighbour* of $w$ if $(w, v) \in E$. Moreover, we denote with $N_G(v)$ the set of all neighbours of $v$ in $G$. When the graph is unambiguous, we use $N(v)$.

A *directed graph* $G$ is a pair $(V, A)$ with $A \subset V \times V$. The elements of $V$ are called *nodes*, while the elements of $A$ are called *arcs*. As with undirected graphs, $G$ must not contain self-loops. A node $v$ is a *in-neighbour* of $w$ if $(v, w) \in A$. Similarly, it is an *out-neighbour* of $w$ if $(w, v) \in A$. The sets of all in- and out-neighbours is denoted as $N^-(v)$ and $N^+(v)$, respectively.

When not further specified, *graph* means undirected graph.

A *path* in a graph $G$ is a sequence of nodes $v_1, \ldots, v_k$ such that $(v_i, v_{i+1})$ is an edge (or an arc) for every $i = 1 \ldots k - 1$. We say that the *length* of this path is $k - 1$ (i.e. the number of edges or arcs involved in it). A subset of nodes of $G$ is *connected* if, for every pair of nodes in it, there is at least a path that connects them that is made only of nodes of the subset. As being connected is an equivalence relation, $G$ may be partitioned into equivalence classes that are called *connected components*.

A *subgraph* of a graph $G$ is the graph that has as nodes a given subset $V'$ of $V$ and as edges some of the edges in $G$ that have as endpoints two members of $V'$. An *induced subgraph* is a subgraph that has all possible edges of $E$ for that choice of nodes.

A graph is *bipartite* if it has no odd cycles or, equivalently, if it can be partitioned into two sets $A$, $B$ such that there is no edge between any pair of nodes in $A$ nor in $B$.

We assume the nodes of a graph to be ordered, and denote them as $v_1, \ldots, v_n$, where $n = |V|$.

We will also use the concept of partial order:

**Definition 1.1.** A **partial order** is a binary relation $\sqsubseteq$ on a set $A$ that is

1. *reflexive*, i.e. $a \sqsubseteq a$ for every $a$ in $A$,

2. *transitive*, i.e. if $a \sqsubseteq b$ and $b \sqsubseteq c$, then $a \sqsubseteq c$,

3. *antisymmetric*, i.e. if $a \sqsubseteq b$ and $b \sqsubseteq a$, then $a = b$.

If $\sqsubseteq$ also satisfies that $\forall a, b \in A$ either $a \sqsubseteq b$ or $b \sqsubseteq a$, we call $\sqsubseteq$ a **total order** on $A$.

An important consequence of the fact that orders are transitive and antisymmetric is that an order contains no cycles, i.e. it is impossible to have a set of distinct elements $a_1, \ldots, a_n$ in $A$ such that $a_1 \sqsubseteq \cdots \sqsubseteq a_n \sqsubseteq a_1$.

If $\sqsubseteq$ is an order, we write $a \sqsubset b$ to mean that $a \sqsubseteq b$ and $a \neq b$.

Finally, we say that an element $m$ is *maximal* if there is no other element $v$ such that $m \sqsubset v$.

An interesting family of orders is the **lexicographical order**. If $A$ is a set with an order $\sqsubseteq$, we define an order on the set of all the tuples of elements of $A$ as follows. Let $P = (p_1, \ldots, p_k)$ and $Q = (q_1, \ldots, q_h)$ be two tuples of elements of $A$.

- If $P$ is a prefix of $Q$ (i.e. $k \leqslant h$ and $p_i = q_i$ for $1 \leqslant i \leqslant k$), then $P \prec Q$.

- Otherwise, let $i$ be the first index such that $p_i \neq q_i$. Then, $P \prec Q$ if and only if $p_i \sqsubset q_i$.

We may also define a lexicographical order on the power set $2^A$ of $A$. More precisely, the order between any two subsets of $A$ is the same as the order between the tuples that have as elements the same elements, in ascending order. As an example, take sets $B = \{9, 3, 7, 5\}$ and $C = \{7, 5, 3, 11\}$ in $\mathbb{N}$. The corresponding sorted tuples are $(3, 5, 7, 9)$ and $(3, 5, 7, 11)$; the former is lexicographically smaller than the latter, so we say that $B$ is lexicographically smaller than $C$.

Another interesting order is the **inclusion order** on the power set of any set $A$: given $B, C \in 2^A$, we say that $B$ is less than $C$ if and only if $B \subset A$.

## 1.2 Enumeration complexity measures

We call $\alpha$ the number of solutions of the enumeration problem. Since for many enumeration problems $\alpha$ may be exponential in the size of the input, traditional complexity definitions are not enough to provide all the information on the running time of an enumeration algorithm. For example, it was proven that in a graph there may be up to

$3^{n/3}$ maximal cliques [41], so the worst-case complexity of any enumeration algorithm for maximal cliques may not be any lower than that.

Considering this fact, complexity classes for enumeration algorithms have been defined [20] in a way that takes into account the effective number of solutions of the given problem. We will report the most useful among these classes here.

We say that an algorithm

- runs in *polynomial total time* if its running time is bounded by a polynomial in the input size and in $\alpha$.

- runs in *incremental polynomial time* if the time needed to generate solution number $X \leqslant \alpha$ is bounded by a polynomial in the input size and in $X$.

- has *polynomial cost per solution* if its running time is bounded by a polinomial in the input size times $\alpha$.

- produces *incremental output* if it outputs the first $X \leqslant \alpha$ within a time bounded by $X$ times a polynomial of the input size.

- has *bounded delay* if the running time between any two consecutive produced solutions is bounded by a polynomial in the input size.

# Chapter 2

# Set systems

This chapter introduces the notion of set systems, discussing some useful classes of them and giving an overview on what is known about enumeration algoritms that list all their maximal elements. Section 2.2 gives some examples of families that form various kinds of set systems.

## 2.1 Definitions

One of the most well-known enumeration problem is that of enumerating all the maximal cliques in a graph. A clique is a complete subgraph of a given graph $G$ (a more formal definition is given in Definition 2.7), and so it may be considered as a subset of the nodes of the graph. Families of solutions given by the subsets of a given set are a recurring situation in enumeration problems (as another example, feasible solutions of a knapsack problem may be considered subsets of a given set), so this situation has justified the following general definition.

**Definition 2.1.** A **set system** $\mathcal{F}$ over a *support set* $U$ is a nonempty family of subsets of $U$, i.e. $\mathcal{F} \subset 2^U$ and $\mathcal{F} \neq \varnothing$. A member of $\mathcal{F}$ is called a *feasible set*.

This thesis will focus mostly on set systems built on graphs. In this case, unless otherwise noted, the support set of the set system will be the set of the nodes of the graph.

As set systems are way too generic to provide any meaningful result, the following subclasses have been studied so far (see for example [6]).

**Definition 2.2.** An **accessible set system** $\mathcal{F}$ over $U$ is a set system that satisfies the following property: for any nonempty $X \in \mathcal{F}$, there is an element $x \in X$ such that $X \setminus \{x\} \in \mathcal{F}$.

**Definition 2.3.** A **strongly accessible set system** $\mathcal{F}$ over $U$ is a set system that satisfies the following property: for any nonempty $X \in \mathcal{F}$ and any $Y \subset X$ that belongs to $\mathcal{F}$, there
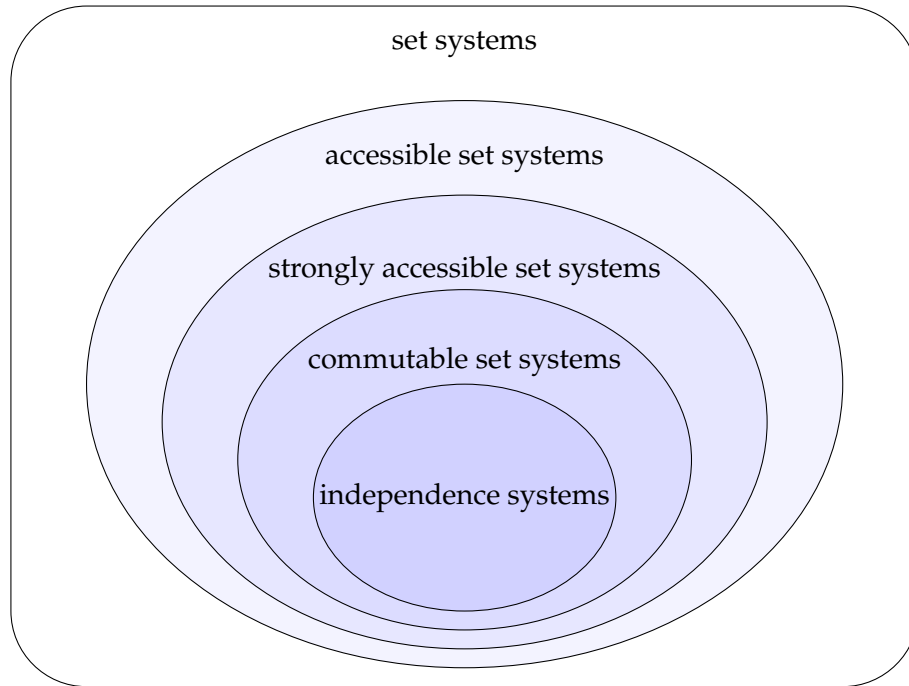
Figure 2.1: Relationships between the various kinds of set systems.

is an element $x \in X \setminus Y$ such that $X \setminus \{x\} \in \mathcal{F}$ or, equivalently, there is an element $y \in X \setminus Y$ such that $Y \cup \{y\} \in \mathcal{F}$.

**Definition 2.4.** An **independence system** $\mathcal{F}$ over $U$ is a set system that satisfies the following property: for any nonempty $X \in \mathcal{F}$, every $Y \subseteq X$ belongs to $\mathcal{F}$ too.

From the definitions, it clearly follows that any independence system is also a strongly accessible set system, and any strongly accessible set system is also an accessible set system. Our main focus will be on a new class that sits in the middle of independence systems and strongly accessible set systems. These relationships are represented in Figure 2.1.

**Definition 2.5.** A **commutable set system** $\mathcal{F}$ over $U$ is a strongly accessible system that satisfies the following property: for any $X \in \mathcal{F}$ and any $Y \subseteq X$ that belongs to $\mathcal{F}$, if $x, y \in X \setminus Y$ are such that $Y \cup \{x\} \in \mathcal{F}$ and $Y \cup \{y\} \in \mathcal{F}$ then $Y \cup \{x, y\} \in \mathcal{F}$.

Another possible definition is the following one:

**Definition 2.6.** A **commutable set system** $\mathcal{F}$ over $U$ is a strongly accessible system that satisfies the following property: for any $X \in \mathcal{F}$ and any $Y \subseteq X$ that belongs to $\mathcal{F}$, if $A, B \subseteq X$ are members of $\mathcal{F}$ that properly contain $Y$, then $A \cup B \in \mathcal{F}$.

**Lemma 2.1.** *Definition 2.5 and Definition 2.6 are equivalent.*

*Proof.* Definition 2.6 trivially implies Definition 2.5.

If $|A \setminus X| = 0$ then there is nothing to prove. We will first prove the equivalence in the case in which $A = X \cup \{a\}$.

If $B = X \cup \{b\}$, then the thesis is exactly Definition 2.5. Otherwise let $b \in B \setminus X$ be such that $X \cup \{b\} \in \mathcal{F}$ (such a b exists because $\mathcal{F}$ is strongly accessible). Then we have that $X \cup \{a, b\} \in \mathcal{F}$ because of Definition 2.5. We may repeat this reasoning to the pair of sets $X \cup \{a, b\} = A \cup \{b\}$ and B, both of which contain $X \cup \{b\}$ and whose union is still $A \cup B$. Since now the number of elements that should be added to the common part of the two sets to obtain B is reduced, this completes the proof by induction.

Let us now consider the general case, and let $a \in A \setminus X$ be chosen such that $X \cup \{a\} \in \mathcal{F}$. Thanks to what we just proved, we know that $B \cup \{a\} \in \mathcal{F}$. As before, we may repeat this process on A and $B \cup \{a\}$, both of which contain $X \cup \{a\}$ and whose union is still $A \cup B$. By induction, this completes the proof.

□

As a set system may contain an exponential number of members, we will assume that the set system is not given explicitly, but through a **membership oracle**, a function $f_{\mathcal{F}}$ that takes a subset of U and returns true if that subset is a member of $\mathcal{F}$.

We will focus on the task of enumerating all the *maximal* (under inclusion) members of $\mathcal{F}$. This is a good compromise between the need of knowing the whole $\mathcal{F}$ and the issue of its size: all members of $\mathcal{F}$ are included in a maximal one and, for set systems that are at least strongly accessible, any element of $\mathcal{F}$ may be found from a maximal one by iteratively removing an element.

**Notation.** Throughout the thesis, n will denote the cardinality of U, q will denote the maximum size of any member of $\mathcal{F}$ and $\alpha$ will denote the number of maximal elements in $\mathcal{F}$.

## 2.2 Some accessible set systems

One of the most studied examples of independent systems are maximal cliques. This problem has been studied in a lot of contexts, including bioinformatics [14, 35], computational chemistry [33, 23, 4] and social network analysis [18, 45]. The family may be described as follows:

**Definition 2.7** (Clique). Given a graph $G = (V, E)$, a subset C of its vertices is called a clique if, for every $a \neq b \in C$ there is an edge between a and b.

Another studied example is the k-plex, which is a generalization of a clique that relaxes the condition that all pairs of nodes must be connected:
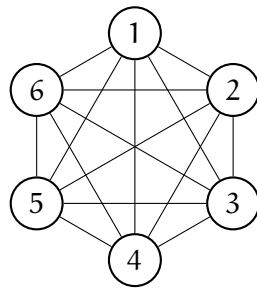
**Definition 2.8** (k-plex)**.** Given a graph $G = (V, E)$, a subset $C$ of its vertices is called a k-plex if, for every $a$ in $C$,

$$|C \cap N(a)| \geqslant |C| - k$$

holds.[1]

This definition of k-plex has some practical issues, such as every set of size at most k being a k-plex. To reduce the impact of this issue, a minor variant of this family has been considered [5]:

**Definition 2.9** (Connected k-plex)**.** Given a graph $G = (V, E)$, a subset $C$ of its vertices is called a connected k-plek if it is a k-plex and it is connected.



2.2a: Clique with six nodes                    2.2b: 4-plex with six nodes

2.2c: Connected 2-plex with six nodes

An example of a clique, a 4-plex and a connected 2-plex can be found in Figure 2.2a, Figure 2.2b and Figure 2.2c respectively. In particular, Figure 2.2b shows a degenerate example of a k-plex that is not even connected, which should not be considered interesting as the objective of finding k-plexes is to find dense structures.

It is easy to see that any k-plex $C$ of size at least $2k - 1$ is also a connected k-plex: indeed, consider any two nodes $a$ and $b$ in $C$. If they are connected by an edge, then there is a path between them and we are done. Otherwise, suppose that their neighbourhoods are disjoint in $C$. Then:

$$|C| - 2 = |C \setminus \{a, b\}| \geqslant |C \cap (N(a) \cup N(b))| = |C \cap N(a)| + |C \cap N(b)| \geqslant 2|C| - 2k$$

---

[1]This reduces to the definition of clique when $k = 1$.

Figure 2.3: Black-connected clique with six nodes. White edges are dashed, black edges are represented with a continuous line.

which implies that $2k - 2 \geqslant |C|$, which is a contraddiction. So their neighbourhoods must intersect in some node $v \in C$, providing a valid path $a, v, b$ between the two nodes.

This implies that considering just connected $k$-plexes does not alter the results for bigger elements of the family, while removing a good amount of small sets.
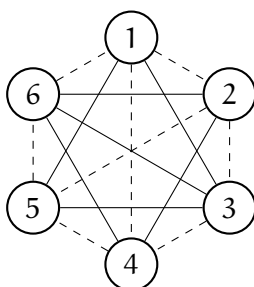
Unfortunately, connected $k$-plexes do not form an independent system (as can easily be seen by considering the subset $\{1, 4\}$ of the connected 2-plex of Figure 2.2c: indeed, this subset is not a connected 2-plex because it is not connected), but they still are commutable set system. This can be easily seen by using Definition 2.6, as any subset of a $k$-plex is still a $k$-plex and the union of any two connected sets that share at least a node is still connected.

Another interesting family that is not an independent system is the family of black-connected cliques. In this context, we have a colored graph that is composed of white and black edges:

**Definition 2.10** (Black-connected clique)**.** Given a graph $G = (V, E)$ in which every edge has an associated color (either black or white), a subset $C$ of its vertices is a black-connected clique if it is a clique and if it is connected in the graph where all white edges are removed.

Figure 2.3 gives an example. Moreover, from the example it is clear that black-connected cliques do not form an independence system, as the subset $\{1, 4\}$ is not a black-connected clique.

This problem arises naturally while studying molecule similarity in bioinformatics and computational chemistry[22, 12]. As with connected $k$-plexes, this family is a commutable set system.

Another family that may be interesting to study for network analysis is the family of connected bipartite (induced) subgraphs of a given graph $G$:

**Definition 2.11** (Connected bipartite (induced) subgraph)**.** Given a graph $G = (V, E)$, a subset $C$ of its edges (vertices) is a connected bipartite (induced) subgraph if the subgraph it defines is both bipartite and connected.

The previous examples of set systems that are not independence systems all share one common trait: they are obtained from independence systems by adding a connectivity constraint. These kinds of independence systems have been called *connected hereditary*, and an algorithm for enumerating them (that uses exponential memory) was given in [10].

It is natural, then, to wonder if any commutable set system is, in fact, a connected hereditary property. This conjecture can be proven false by considering the family $\mathcal{F}$ of cliques that all share a given node $v$ in common: indeed, they are a commutable set system, but they are not an independent set system nor they can be obtained as the family of all cliques that satisfy a connectivity constraint (because in that case any singleton would be a member of the family, while in this case $\{v\}$ is the only singleton that belongs to $\mathcal{F}$).

The techniques used for the enumeration of set systems may be used to solve enumeration problems that one would not ordinarily associate with set systems. For example, one might enumerate all the s-t paths by enumerating the maximal elements of the set system made of all paths that begin in a certain node $s$ and may be extended to reach a certain destination t.

In a similar way, the set of all sub-forests in a graph form an independence system, whose maximal elements are spanning trees: so it is possible to enumerate spanning trees with techniques that enumerate independence systems.

Another family that may be considered an independence system is the (complement of) s-t (vertex-)cuts in a graph:

**Definition 2.12** (s-t (vertex-)cuts)**.** Given a connected graph $G = (V, E)$ and two vertices s, t, a subset C of E $(V \setminus \{s, t\})$ is an s-t (vertex-)cut if and only if the graph obtained by removing C from G has s and t in two different connected components.

This family is clearly not even accessible, but the family formed by its complement is. It may be defined directly as the family of all subsets of E (or $V \setminus \{s, t\}$) that induce connected components such that s and t belong to different components.

A similar family to the previous one is given by the feedback vertex (arc) set problem:

**Definition 2.13** (Feedback vertex (arc) set)**.** Given a graph (possibly directed) $G = (V, E)$, a feedback vertex (arc) set is a subset F of V (E) such that the graph obtained from G by removing the elements in F contains no cycles.

As before, this is not an accessible family, but considering its complement we obtain an independence system. If we furthermore reduce ourselves to consider subset that leave the remaining graph connected, we obtain a commutable set system. This problem has been studied before, but known algorithms use an exponential amount of memory [47, 36].

## 2.3 Lower bounds

It was proven in [26] that it is not reasonable to try to find an algorithm that is able to find the maximal sets of **any** independence system in polynomial total time, as shown by the following reduction:

**Theorem 2.1** (Lawler et al [26]). *There is no algorithm that enumerates all the maximal sets in any independence system in polynomial total time, unless* $P = NP$.

*Proof.* Let $F(X_1, \ldots, X_n)$ be a boolean expression in conjunctive normal form. We will define an independence system on the ground set $E = \{T_1, F_1, \ldots, T_n, F_n\}$ as follows. Let $J$ be any subset of $U$ and $x_i(J)$ be defined as:

$$x_i(J) = \begin{cases} \texttt{true} & \text{if } T_i \in J, \ F_i \notin J \\ \texttt{false} & \text{if } F_i \in J, \ T_i \notin J \\ \texttt{undefined} & \text{if } T_i \notin J, \ F_i \notin J \\ \texttt{overdefined} & \text{otherwise} \end{cases}$$

We say that $J \in \mathcal{F}$ if and only if one of the following holds:

- The assignment $X_i = x_i(J)$ satisfies the boolean formula $F$, or

- We have $T_i \notin J$ and $F_i \notin J$.

This is clearly an independence system, as any subset of a $J$ that satisfies the second condition still satisfies it and any subset of a $J$ that satisfies the first condition has at least one undefined literal.

Moreover, any solution to $F$ provides a maximal set for $\mathcal{F}$; maximal sets not obtained this way must have at least an undefined literal, so they are of the form $E \setminus \{T_i, F_i\}$ for some $i$. We call these maximal sets trivial. Note that there are $n$ trivial maximal sets.

Let us now suppose that there is an algorithm which is able to enumerate all the maximal elements of $\mathcal{F}$ in $P(n, \alpha)$ time. Then, if we let this algorithm run for $P(n, n)$ time, one of the following situations may happen:

- The algorithm terminates and finds no non-trivial maximal set: we know that $F$ has no solution.

- The algorithm terminates and finds a non-trivial maximal set: we know that $F$ has a solution.

- The algorithm does not terminate: this implies that $\mathcal{F}$ has more than $n$ maximal solutions, so at least one of them must be non-trivial and thus $F$ must have a solution.

This reasoning gives an algorithm that takes $P(n, n)$ (i.e. polynomial) time to check if $F$ has a solution or not, thus proving the reduction.

$\square$

Moreover, from the same proof another condional lower bound follows, based on the well-known Strong Exponential Time Hypothesis of Impagliazzo and Paturi[19].

**Corollary 2.1.** *There is no algorithm that enumerates all the maximal sets in any independence system in* $O(2^{\frac{q}{2}-\epsilon}P(\alpha, n))$, *unless SETH is false.*

To circumvent this issue, any general technique that promises to enumerate maximal elements in strongly accessible set systems relies on some form of problem-specific insight. In particular, *certificates* and *restricted problems* have been considered.

**Definition 2.14** (Certificate). A certificate for a given enumeration problem is some kind of information that can be computed in polynomial time and guarantees the existence of at least one solution for the given subproblem that must be solved during the recursive enumeration.

**Definition 2.15** (Restricted problem). Given a strongly accessible set system $\mathcal{F}$ on $U$, a maximal feasible set $F$ and an element $v \in E \setminus F$, the restricted problem $\mathcal{P}(F, v)$ is the problem of enumerating all the maximal elements of the family

$$\mathcal{G}_F^v = \{A \in \mathcal{F} : A \subseteq F \cup \{v\}\}$$

Figure 2.4 represents the solutions of the restricted problem that is obtained for the connected k-plex problem from the k-plex of Figure 2.2c by adding node 7. The solution on top is the trivial solution (i.e. the solution that generated the restricted problem), while the other four figures are obtained by choosing 7 and all the possible subsets of non-neighbours of 7, and choosing the other nodes accordingly.

## 2.4  Previous algorithms

Previous works on the enumeration of maximal elements of set systems are mainly based on two different techniques: binary partition and reverse search. See for example [34] for an overview of enumeration algorithms, where its Chapter 9 deals with enumeration of subgraphs satisfying certain properties.

### 2.4.1  Binary partition

Binary partition is a general, recursive scheme that works by splitting the problem into the following subproblems: "enumerate all maximal sets that contain all the elements that belong to set $S$, no elements that belong to set $X$ and possibly some of the elements that do not belong to either of them". The recursion starts with $S = X = \varnothing$, and

Figure 2.4: Example of a restricted problem using the connected 2-plex of Figure 2.2c and the extra node 7. The figures below highlight the non-trivial solutions, the one on the top the original k-plex that defines the restricted problem.

proceeds by choosing, at every recursive step, some node $v$ in $V \setminus (S \cup X)$ and then by producing one or two recursive children, with $v$ added to $X$ or to $S$. This second option is only possible if we have that $S \cup \{v\} \in \mathcal{F}$, as otherwise the algorithm would produce incorrect solutions. When $S \cup X = V$, the recursion ends and the current solution $S$ is generated as output.

This scheme guarantees a polynomial time per solution whenever we have a certificate for the problem we are considering. Indeed, the presence of a certificate allows us to immediately stop the recursion in nodes where we know there will not be any solution, producing a recursion tree that has at most $\alpha n$ nodes. As the cost of the computations done in every node of the recursion tree is polynomial in the input size, the total cost is polynomial in the input and the number of solutions. Further optimizations (dependent on the problem) usually allow to give algorithms that are extremely

efficient in these cases.

When no certificate is available, no polynomial upper bound in the input and output size can usually be given. Nonetheless, algorithms that are very efficient in practice may be obtained by using some sort of heuristic that allows to eliminate a good percentage of the "dead ends" of the recursion tree.

As an example, binary partition can easily be used to enumerate spanning trees or s-t paths in polynomial time per solution. For the first case, a simple certificate is given by the check that the graph obtained by removing all the excluded edges is still connected, as that guarantees the existence of a spanning tree. As for s-t path, a certificate can be obtained by checking that the endpoint of path defined by the current set of "taken" edges is still able to reach t, without using any excluded edge.

The Bron-Kerbosch algorithm [7], used for enumerating maximal cliques and the basis of a lot of algorithms for the enumeration of cliques, k-plexes and black-connected cliques [22], is a straightforward application of binary partition, with a smart pivoting rule that allows to remove most "unwanted" recursive nodes.

### 2.4.2   Reverse search

Reverse search is a general technique in enumeration which was initially introduced by Avis and Fukuda [3]. It works by implicitly constructing a directed graph that has as nodes all the solutions to the enumeration problem, with an arc from a certain solution $P$ to another solution $Q$ if the second can be obtained in a certain way from the first. As long as the graph is properly connected, a simple graph visit starting from any solution will then enumerate all the solutions to the problem.

As this graph may contain multiple paths to a single solution, we either need to keep track of all the solutions already found so far or we need some way to remove edges from the graph, making it a directed forest whose roots are known.

In the first case, there is no performance penalty, but the resulting algorithm takes exponential memory to run. In the second case, some more computations may be required to discard unneeded edges.

Regarding set system enumeration, this technique is usually applied by building a graph over all the maximal elements of $\mathcal{F}$. Given a solution $S$, outgoing edges are obtained by using a node from $V \setminus S$ as a "guide" to find new maximal elements. In particular, this is usually done by solving the restricted problem on $(S, v)$ and applying some form of post-processing to the found solutions. Examples and further details are provided in Section 3.2.

This idea is the one at the basis of most algorithms for maximal clique enumeration that guarantee polynomial time per solution, as in [13]. Moreover, it was also employed in [10] to give an algorithm for the enumeration of maximal elements in an independence system. In the same paper, the authors give an algorithm for the enumeration of maximal sets that satisfy "connected hereditary" properties, using expo-

nential space. As the sets that satisfy connected hereditary properties are a special case of commutable set systems, this thesis improves their work by removing the need for exponential space.

# Chapter 3

# Framework description

This chapter will explain how the new framework for the enumeration of maximal elements in commutable set systems works. We will start by giving some necessary definitions, then we will proceed to explain the algorithm itself. Throughout this section, we assume the elements of $U$ to be ordered with an arbitrary total order.

## 3.1 Core concepts

We will now introduce some concepts that will be fundamental while explaining how the algorithm works. As they are not very intuitive, most definitions will have an example on the black-connected clique $C$ of Figure 2.3.

**Definition 3.1.** Given a commutable set system $\mathcal{F}$ and one of its feasible sets $S$, we say that $s$ is a **seed** of $S$ if $s \in S$ and $\{s\} \in \mathcal{F}$. The **canonical seed** of $S$, denoted by $seed(S)$, is the smallest possible seed according to the ordering of the elements of $U$.

Any singleton in $C$ is a seed, since it is both a clique and connected with black edges. So, its canonical seed is 1.

The concept of level, that is introduced with the next definition, will be fundamental to define a $complete$ function (one of the main ingredients in the reverse search algorithm) that satisfies the property stated in Lemma 3.2 but is not NP-hard to compute (see Lemma 3.1). Moreover, the properties of the $complete$ function obtained using this definition will allow us to prove Theorem 3.1.

**Definition 3.2.** Given a commutable set system $\mathcal{F}$, one of its feasible sets $S$ and a seed $s$ of $S$, we define the **level** of an element $v$ with respect to $s$ as follows:

- if $v = s$, then the level of $v$ is 0.

- let $k$ be the smallest integer such that there is a set $S' \subset S$ that is composed of elements of level $\leqslant k$ and such that $S' \cup \{v\} \in \mathcal{F}$ and $s \in S'$. Then the level of $v$ is $k + 1$.

Figure 3.1: Black-connected clique of Figure 2.3 with nodes partitioned into groups according to their level. Only black edges are represented, and the canonical seed 1 is used.

- if there is no such subset, we say that the level of $v$ is $\infty$.

We will use $level_S^s(v)$ to denote the level of $v$ in $S$ with respect to seed $s$. When the seed is not specified, it is assumed to be the canonical seed.

Figure 3.1 represents C according to its levels. It follows immediately from the definition of level that the level of a node $v$ with respect to a given seed $s$ is given by the distance between $v$ and $s$ (i.e. the length of the shortest path) according to black edges. This actually holds true for any connected hereditary family.

These two definitions are crucial to allow us to define the following order between two feasible sets.

**Definition 3.3.** The **level order** $\prec$ between any two solutions P, Q of a commutable set system $\mathcal{F}$ is defined as follows:

- Let $lP = \{(level_P(v), v) \ \forall v \in P\}$, the set pairs made by the level of an element and the element itself.

- Let $lQ$ be defined in the same way.

- We say that $P \prec Q$ if and only if $lP$ is smaller than $lQ$ using the lexicographical order between sets defined in Section 1.1.

According to this definition, the set corresponding to C, i.e. the set of pairs of the form $(dist(s, v), v)$, is $\{(0, 1), (1, 3), (1, 5), (2, 6), (3, 2), (3, 4)\}$.

We define $complete(S, s)$ as the maximal solution that is obtained from $S$ by iteratively adding the element $v \in E \setminus S$ such that adding $v$ keeps the current set in $\mathcal{F}$ and that, among all those possible choices, minimizes $(level_S^s(v), v)$. The definition of $complete(S)$ takes into account that, while we are adding elements to $S$, we might add an element $s'$ that is smaller than $seed(S)$. In that case, subsequent iterations consider the levels to be relative to the new canonical seed.

Figure 3.2: A graph with black and white edges that is used in the proof of Lemma 3.1. Only black edges are represented, white edges are implicitly assumed to be present between any two nodes, except for the ones that already have a black edge or that have a crossed-out red edge. Assuming $k = 3$ and $n = 5$, the figure represents the formula $(x_1 \lor x_3) \land (x_1 \lor x_3 \lor x_4) \land (x_2 \lor \neg x_3 \lor x_4 \lor \neg x_5)$

In previous works, $\texttt{complete}(S)$ was simply defined as the lexicographically minimum maximal solution that contains $S$. Unfortunately, for a generic commutable set system this function is NP-hard to compute, as stated by the following lemma.

**Lemma 3.1.** *Given a graph $G$ whose edges are either black or white and a non-maximal black-connected clique $X$ of $G$, it is NP-hard to find the lexicographically minimum among the maximal black-connected cliques containing $X$.*

*Proof.* We prove that a $\texttt{complete}(X)$ function that returns the lexicographically mini-

mum black-connected clique containing $X$ can be used to solve a SAT problem in polynomial time, by building a graph with a number of nodes linear in the amount of the clauses and variables in the formula.

Given a SAT formula in conjunctive normal form with $n$ variables $x_1 \ldots x_n$ and $k$ clauses $d_1 \ldots d_k$, we build the the graph in Figure 3.2, whose nodes are $C_1 \ldots C_k$, $T_1 \ldots T_n$, $F_1 \ldots F_n$ and $Y_1 \ldots Y_n$, labelled increasingly in this order (i.e., nodes $C_1 \ldots C_k$ have smaller label than all other nodes). Each $Y_i$ is connected with a *black* edge to $T_i$ and $F_i$, and, except for $Y_1$, also with $T_{i-1}$ and $F_{i-1}$. Each $C_i$, which corresponds to $d_i$, is connected with *black* edge to $T_j$ (resp. $F_j$) whether $d_i$ contains a positive (resp. negative) occurrence of $x_j$. Hence, nodes in $C_1 \ldots C_k$ are connected with black edge to an arbitrary amount of $T_i$ and $F_i$ nodes, but not to any $Yi$ node. All other pairs of nodes are connected with a *white* edge, except for the pairs $T_i, F_i$ (symbolized by the crossed-out red edge in Figure 3.2).

It is straightforward to see that any maximal black-connected clique in this graph will contain exactly one of $T_i$ and $F_i$, for any $i$, and that any maximal black-connected clique containing all nodes in $C_1 \ldots C_k$ will be lexicographically smaller than any that does not contain all of them (as they have the smallest labels).

Consider `complete({Y1})`, the lexicographically smallest maximal black-connected clique containing $Y_1$. Any black-connected clique containing $Y_1$ and all $C_i$ nodes represents a satisfying assignment for the formula at hand. Indeed, for each $C_i$ node to be reachable from $Y_1$ with black edges, at least one of the $T_j$ or $F_j$ nodes connected to $C_i$ must be in the black-connected clique; the set of $T_i$ and $F_i$ nodes in the black-connected clique will thus give us the value (true or false) of the corresponding variable $x_i$ (recall that we cannot have the pair $T_j$-$F_j$ in the same black connected clique). Hence, in order to verify that the formula is satisfiable, we only need to compute `complete({Y1})` and check whether this contains all $C_i$ nodes.

$\square$

Analogously, we say that $P$ is a *prefix* of $S$ if it is obtained from $S$ by iteratively removing the element with the the highest level and, in case of ties, the highest element among those with maximal level. Note that the definition of levels, combined with the properties of commutative set systems, ensures that removing a node with maximal level will not cause $S$ to become unfeasible. We also say that $P$ is a *prefix of $S$ with respect to $s$* to explicitly specify that $s$ is the seed to be used.

Given the set form of $C$, it follows that its prefixes are given by $\{1\}$, $\{1, 3\}$, $\{1, 3, 5\}$, $\{1, 3, 5, 6\}$, $\{1, 3, 5, 6, 2\}$ and $\{1, 3, 5, 6, 2, 4\}$, that are all black-connected cliques.

Note that, if $P$ is a subset of $S$, $\text{level}_P^s(v) \geqslant \text{level}_S^s(v)$. If $P$ is a prefix of $S$, then the equality holds for any $v$ such that both levels are finite.

By the definition of `complete`, any prefix of $S$ may be computed by executing a limited number of steps of `complete(P)`, where $P$ is any smaller prefix of $S$, while limiting the selection of elements to add to the current solution to only elements that are part of $S$.

Thanks to the definition of $\mathtt{complete(S)}$ and the definition of level ordering, we obtain the following important lemma.

**Lemma 3.2.** *Given a feasible set* S, *if* Q *is any maximal feasible set of* $\mathcal{F}$ *such that* $S \subseteq Q$ *and* $\mathtt{seed(S)} = \mathtt{seed(Q)}$, *we have* $\mathtt{complete(S)} \preceq Q$.

*Proof.* If $\mathtt{seed(complete(S))} \neq \mathtt{seed(S)}$, then the thesis clearly follows, as the seed is the only level 0 element. Moreover, if $Q = \mathtt{complete(S)}$, there is nothing to prove.

Let now $i > 0$ be the smallest level such that the set of elements of level $i$ in $Q$ is different from the set of elements of level $i$ in $\mathtt{complete(S)}$. Moreover, let $v$ be the smallest element that belongs to exactly one of these two sets.

We will now prove that $v$ does, in fact, belong to level $i$ of $\mathtt{complete(S)}$. Let X be the union of all levels smaller than $i$ in $Q$, as well as of all elements on level $i$ that are smaller than $Q$. Thanks to the definitions of levels and to the alternative definition of commutable set systems (Definition 2.6), we have that $X \cup \{v\}$ is feasible, as it may be obtained as the union of feasible sets all containing $\mathtt{seed(Q)}$ and all contained in $Q$. Moreover, since S is feasible, is contained in $Q$ and contains $\mathtt{seed(Q)}$, we have that $X \cup \{v\} \cup S$ is feasible. Since it can be easily seen that $X \cup S$ is one of the sets obtained during the iterations that build $\mathtt{complete(S)}$, we have that $v$ was considered for adding during that iteration. As $v$ was chosen as the smallest node in the difference between level $i$ of the two sets, we know that at that iteration no node smaller than $v$ was viable, and so $v$ was added to $X \cup S$ by $\mathtt{complete}$.

Now, if $v$ is bigger than all the elements in level $i$ of $Q$, we have that level $i$ of $\mathtt{complete(S)}$ is a superset of level $i$ of $Q$. Since $Q$ is maximal, it cannot be a proper subset of $\mathtt{complete(S)}$, so it must have at least one element on level $i + 1$. Otherwise, level $i$ of $\mathtt{complete(S)}$ compares smaller than level $i$ of $Q$, so by the definition of level order it follows that $\mathtt{complete(S)} \prec Q$.

$\square$

All this machinery allows us to definte the concepts used for executing the reverse search.

**Definition 3.4.** If a maximal solution S is such that $\mathtt{complete(seed(S))} = S$, we call it a $\mathtt{root}$. If S is not a $\mathtt{root}$, we call one of its nonempty prefixes $\mathtt{core(S)}$ if it is the longest prefix such that $\mathtt{complete(core(S))} \neq S$. Moreover, in that case the smallest element of $S \setminus \mathtt{core(S)}$ (according to the level order) is called parent index of S (denoted by $\mathtt{parind(S)}$) and $\mathtt{complete(core(S))}$ is called $\mathtt{parent(S)}$.

Note that Lemma 3.2 shows that, for any non-root maximal solution, $\mathtt{parent(S)} \prec S$. This property will be crucial when proving that the solution graph defined by $\mathtt{parent}$ does not contain any directed cycles.

## 3.2   Algorithm

This section will first outline the reverse search algorithm employed by the framework and will then provide further details on how the graph of the reverse search is generated.

### 3.2.1   Reverse search explained

We will now describe the graph that is explored by the reverse search algorithm used by the framework. While doing this, we use the definition of parent from the Section 3.1, but the results are valid for any definition of parent that satisfies $parent(S) \prec S$ for some order $\prec$.

Consider the graph $\mathcal{R}$ that has all the maximal sets of $\mathcal{F}$ as nodes. The arcs of $\mathcal{R}$ are defined as follows: for every maximal non-root solution S, we have an arc from $parent(S)$ to S.

Thanks to the fact that $parent(S) \prec S$ and to the properties of orders, we know that the arcs of $\mathcal{R}$ cannot contain any directed cycle. Moreover, as every node has at most one ingoing edge, there can be no undirected cycle either, as otherwise one node would have at least two ingoing edges.

Since $\mathcal{R}$ contains no cycles, it must be a forest. Let us now consider any maximal subtree in this forest. Its root has no ingoing edge, so it is, in fact, a maximal solution with no parent (i.e. a root). Moreover, any other solution in that subtree is reachable from its root: since a root is the only node in that tree that has no parent, walking back from S to $parent(S)$ as long as it is possible must eventually reach the root.

Let now S be any maximal solution. We define $children(S)$ to be the family of maximal sets such that, for any maximal solution Q, we have that the parent of Q is in S if and only if Q belongs to $children(S)$.

It is now easy to define an algorithm that enumerates all the maximal solutions in a commutable set system $\mathcal{F}$: simply start the explorations from the root of every forest, exploring every tree in a depth-first manner. Pseudocode for this algorithm is shown in Algorithm 1, which is a simple recursive visit on a directed forest. A more complicated algorithm, that avoids using a recursion stack (thus potentially improving memory usage), can be found in Algorithm 2. More in detail, this algorithm performs a depth-first search that jumps from a node to the next without storing anywhere what the last node was. So, to find the next node we must first check if the current node has some child that is yet unvisited (this can be obtained by visiting the children in some deterministic order, thus allowing to skip all the children that have already been visited just by knowing the last explored node). If that is not the case, then the algorithm has to climb up in the tree, by executing the parent function.

Note that in some cases the parent of a given maximal solution, or the exact set of children, may be hard to compute. We suppose that, instead of $children$, we have a similar $candidates$ function available, that satisfies $candidates(S) \supseteq children(S)$. In

---

**Algorithm 1:** Reverse search algorithm

---

**for** *every solution* S *that is a* root **do**
  | spawn(S)
**end**
**Function** spawn(P)
  | **for** S ∈ children(P) **do**
  |   | spawn(S)
  | **end**

---

this case, we may define a graph $\mathcal{S}$ that has the same nodes of $\mathcal{R}$ but different arcs. In particular, the arcs of $\mathcal{S}$ are defined as follows: for any pair of maximal solutions S, Q, we have an arc from S to Q if and only if Q ∈ candidates(S). Clearly $\mathcal{R}$ is a subgraph of $\mathcal{S}$, so an exploration of this graph starting from all the roots will report all the maximal solutions. The issue of avoiding duplication may be solved by keeping a global set of all solutions discovered so far. The pseudocode for this approach (that corresponds to a "traditional" implementation of a depth-first search on a generic graph) can be found in Algorithm 3 and is the one used to enumerate maximal sets that satisfy connected hereditary properties in [10].

Finally, note that by using the technique of *alternative output*, as described in [39], that is by outputting a solution when going down in the computational tree if the current height is even and when going up otherwise, we obtain algorithms whose maximum **delay** is given by the cost of running the algorithm for children.

### 3.2.2 Computing children(S)

We will now explain how to compute the children of a given maximal solution S of a commutable set system $\mathcal{F}$. As this cannot be done in the generic case in sub-exponential time (since doing so would give a polynomial total time algorithm for the enumeration of any commutable set system, in violation of Theorem 2.1), we will use the concept of restricted problem from Definition 2.15 as a blackbox to aid us in the search of the children of S. Moreover, when doing so some care will be necessary to avoid generating the same child multiple times.

The following key fact will give us an the algorithm that generates children(S), as it gives us a reasonably-sized set of candidates that is assured to contain all the children of S.

**Theorem 3.1.** *Given a maximal solution* C *such that* parent(C) = S, *we have that* core(C) *is the prefix ending just before* parind(C) *of a solution of* $\mathcal{P}(S, \text{parind}(C))$.

*Proof.* Let $v$ be parind(C) and let C$v$ be core(C) ∪ {$v$}. We will first prove that $v \notin$ S.

Suppose by contraddiction $v \in$ S. Since we know parent(C) = complete(core(C)), let y be the first element that is added to core(C) by complete. Clearly, it cannot

---

**Algorithm 2:** Reverse search algorithm without recursion

---

**for** *every solution* S *that is a* root **do**

    $\text{Last} \leftarrow \varnothing$;

    **loop**

        $\text{Child} \leftarrow \varnothing$;

        **for** $C \in \text{children}(S)$ *that comes after* Last **do**

            $\text{Child} \leftarrow C$;

            **break**;

        **end**

        **if** $\text{Child} = \varnothing$ **then**

            **if** S *is a* root **then**

                **break**;

            **end**

            $\text{Last} \leftarrow S$;

            $S \leftarrow \text{parent}(S)$;

        **else**

            $\text{Last} \leftarrow \varnothing$;

            $S \leftarrow \text{Child}$;

        **end**

    **end**

**end**

---

be $v$, as otherwise we would have $\text{parent}(C) = \text{complete}(Cv) = C$, that contraddicts the definition of parent. By the definition of $\text{complete}$, we immediately know that $(\text{level}_{\text{core}(C)}(y), y) < (\text{level}_{\text{core}(C)}(v), v)$ and that, as $\text{core}(C)$ is a prefix of $C$, $(\text{level}_C(y), y) < (\text{level}_C(v), v)$. Thus, $y \notin C$, as that would mean that $v$ is not the element immediately after $\text{core}(C)$ in $C$.

By Definition 2.6, we know that, since both $Cv$ and $\text{core}(C) \cup \{y\}$ belong to $\mathcal{F}$, $\text{core}(C) \cup \{v, y\}$ belongs to $\mathcal{F}$ too. Thus, $y$ is a viable choice for $\text{complete}$ when it is expanding the set $Cv$. As any element that is a viable candidate for $Cv$ but not for $\text{core}(C)$ must have a level of $\text{level}_{Cv}(v) = \text{level}_{\text{core}(C)}(v) + 1$, then $(\text{level}_{Cv}(y), y)$ is still the lowest level-value pair among the viable nodes for $\text{complete}(Cv)$, and so it will be chosen as the next element. This implies that $\text{complete}(Cv)$ contains $y$, which contraddicts $y \notin C$ since, by the definition of $\text{core}$, we know $\text{complete}(Cv) = C$. This proves that $v \notin S$.

As $Cv \in \mathcal{F}$, we have that $Cv \in \mathcal{G}_S^v$. So $Cv$ must be contained in a maximal solution of $\mathcal{G}_S^v$, i.e. in a solution of $\mathcal{P}(S, v)$. We will call this solution R.

We now need to prove that $Cv$ is indeed a prefix of R with respect to $s$, the canonical seed of $C$. If $C = Cv$, then there is nothing to prove, as then $R \supseteq C$ and so $R = C$ (since

---

**Algorithm 3:** Reverse search algorithm with exponential memory

$\mathcal{S} = \varnothing$;
**for** *every solution S that is a* root **do**
  | $\mathrm{spawn}(S)$
**end**
**Function** $\mathrm{spawn}(P)$
  | **if** $S \in \mathcal{S}$ **then**
  |   | **return**;
  | **end**
  | add $S$ to $\mathcal{S}$;
  | **for** $S \in \mathrm{children}(P)$ **do**
  |   | $\mathrm{spawn}(S)$
  | **end**

---

$C$ is maximal).

Otherwise, $R$ strictly contains $Cv$. Let $r$ be the smallest element in $R \setminus Cv$ and suppose by contraddiction that $(\mathrm{level}_R^s(r), r) < (\mathrm{level}_R^s(v), v)$. Then we also have that $(\mathrm{level}_R^s(r), r) < (\mathrm{level}_{Cv}^s(v), v)$, thanks to the fact that the level may not decrease when taking subsets.

Let us denote by $Rp$ the prefix of $R$ that ends just before $r$: note that $Rp$ is a prefix of $C$. Moreover, by Definition 2.6, we know that $Cv \cup \{r\}$ is in $\mathcal{F}$ too, as both $Cv$ and $Rp \cup \{r\}$ are and they both contain $Rp$.

Let $y$ be the element of $C$ that immediately follows $Cv$ in $C$. Then, by the definition of prefix, we know that $(\mathrm{level}_C^s(v), v) < (\mathrm{level}_C^s(y), y)$. Since the level of a node does not change by taking prefixes (as long as it does not become infinite), it follows that $(\mathrm{level}_{Cv}^s(v), v) < (\mathrm{level}_{Cv}^s(y), y)$. Moreover, we also have the following equality:

$$\mathrm{level}_R^s(r) = \mathrm{level}_{Rp}^s(r) = \mathrm{level}_{Cv}^s(r)$$

Putting all this together allows us to conclude that

$$(\mathrm{level}_{Cv}^s(r), r) < (\mathrm{level}_{Cv}^s(y), y)$$

that imples that the first step of $\mathrm{complete}(Cv)$ chooses $r$, so that $\mathrm{complete}(Cv) \neq C$, contraddicting the definition of $\mathrm{core}$. This proves that $Cv$ is indeed a prefix of $R$.  □

We will denote with $R(C)$ the solution of $\mathcal{P}(\mathrm{parent}(C), \mathrm{parind}(C))$ (according to level order) that is obtained by running $\mathrm{complete}(\mathrm{core}(C) \cup \{\mathrm{parind}(C)\})$ using only candidates from $\mathrm{parent}(C) \cup \{\mathrm{parind}(C)\}$.

We now know that any child $C$ of $S$ may be found by first chosing the node $v$ that should be $\mathrm{parind}(C)$, then trying all the possible solutions $R$ of $\mathcal{P}(S, v)$ as candidates for $R(C)$, then checking if the prefix with respect to any seed of $C$ ending just before

$v$ is, in fact, the core of a child. To avoid duplication, we check that the tuple $(v, R, s)$ choosen to generate the child was indeed the one that can be found from $C$ itself as $(\text{parind}(C), R(C), \text{seed}(C))$. Since we have proven that when that specific tuple is chosen then $C$ is obtained by the previous procedure, this argument shows the correctness of the algorithm, whose pseudocode can be found in Algorithm 4.

---

**Algorithm 4:** Children generation algorithm

---

**Function** children(S)

    **for** $v \in E \setminus S$ **do**

        **for** R *solution of* $\mathcal{P}(S, v)$ *different from* S **do**

            **for** s *possible seed of* R **do**

                $Cv \leftarrow$ the prefix of R with respect to s ending with $v$;

                $C \leftarrow$ complete(Cv);

                **if** $\text{parind}(C) = v$ *and* $\text{parent}(C) = S$ *and* $R = R(C)$ *and* $\text{seed}(C) = s$ **then**

                    **yield** C;

                **end**

            **end**

        **end**

    **end**

---

Note that some speed-ups may be obtained by choosing $v$ from a smaller set, as long as any $v$ such that $\{v\}$ is not the only non-trivial solution to $\mathcal{P}(S, v)$ belongs to that set. This is because complete($\{v\}$) is either a root or does not satisfy $\text{seed}(\text{complete}(\{v\})) = v$.

### 3.2.3 An algorithm that does not require the restricted problem

In Subsection 3.2.2, the restricted problem is only used as a tool to find the cores of the children of the current node. Another way to reach the same objective is trying all the possible subsets of $S \cup \{v\}$. This algorithm would be slower in the majority of cases, but requires no knowledge of the problem to be applied. Moreover, with slight modifications to the definition of level it can be used on set systems that are not commutable but only strongly accessible.

More precisely, we may define the level of a node as follows:

- If $v \notin S$, then $\text{level}(v) = |S|$.

- If $v \in S$, then run complete($\{s\}$) while considering only candidates from S. Let $S'$ be the solution such that the next iteration of complete adds $v$. Then $\text{level}(v) = |S'|$.

When replacing Definition 3.2 with this one, both `complete` and `parent` are still well-defined, and the same proof for Lemma 3.2 still holds. We thus obtain that the algorithms of Subsection 3.2.1 may be used to enumerate any strongly accessible set system.

### 3.2.4 Analysis

The running time of Algorithm 4 critically depends on the cost of solving $\mathcal{P}$. The following lemma will tie together $\alpha$ and the number of solutions of $\mathcal{P}$.

**Lemma 3.3.** *The number of solutions of* $\mathcal{P}(S, \nu)$ *is at most* $(|S| + 1)\alpha$.

*Proof.* Let us fix a node $s$ of $S \cup \{\nu\}$ and consider all the solutions of $\mathcal{P}(S, \nu)$ that have $s$ as a seed. We will prove that there are at most $\alpha$ such solutions: as $s$ may be chosen in $|S| + 1$ ways, this proves the thesis.

Indeed, let us suppose by contraddiction that there are more that $\alpha$ solutions of $\mathcal{P}(S, \nu)$ that contain s. As any solution of $\mathcal{P}(S, \nu)$ is contained in a maximal solution of $\mathcal{F}$, by the pidgeonhole principle, at least two such solutions (say Q and R) must be contained in the same maximal solution B of $\mathcal{F}$. From Definition 2.6 then follows that, since Q and R belong to $\mathcal{F}$, are contained in an element of $\mathcal{F}$ and both contain $\{s\}$, another viable solution, $Q \cup R \in \mathcal{F}$ too. But this contraddicts the maximality of Q and R in $\mathcal{G}_S^\nu$, so the thesis is proven. $\square$

This fact, combined with the algorithm explained above, gives us the following theorem:

**Theorem 3.2.** *Let $\mathcal{F}$ be a commutable set system over a set $U$. If the restricted problem can be solved in polynomial time (resp. polynomial total time), then the maximal solutions of $\mathcal{F}$ may be enumerated with polynomial delay (resp. in polynomial total time).*

Note that algorithms obtained with this framework have a slightly higher computational cost with respect to the corresponding ones given by the framework in [10]. However, algorithms obtained with our framework have an huge advantage in a parallelized or distributed settings, as they do not require a centralized set that stores all the maximal solutions generated so far. Such a set would create a huge bottleneck, as it would be the source of a lot of inherently sequential work, thus limiting the maximum speed-up that can be obtained. In contrast, our framework requires almost no coordination, with the exception of the work required to split the forest in roughly equal-sized chunks to be assigned to each worker.

Let us now consider the variant of the algorithm that does an exaustive search of the possible cores of the children. That algorithm gives us a delay bounded by $O(P(n)2^q)$, which is the square of the lower bound given in Corollary 2.1.

This bound gives us the following corollary:

**Corollary 3.1.** *Let $\mathcal{F}$ be a commutable set system over a set $U$ such that any element $S \in \mathcal{F}$ has size at most $q = O(\log(n))$. Then all the maximal solutions of $\mathcal{F}$ may be enumerated with polynomial delay.*

Note that the trivial algorithm (try all the sets with size up to $q$) does not guarantee this, as its running time is $O(n^q) = O(n^{c \log n})$, that is super-polynomial.

Finally, we will spend a couple words on (possible) implementations of complete and parent and their runtime costs. Suppose that determining the level of a node takes $\mathcal{L}(q)$ time, and that we know that the viable elements to be added belong to a set $A$ of size at most $\mathcal{A}$ that can be listed in $O(\mathcal{A})$ time. Then, complete may be easily implemented by scanning $A$, finding the level of each element, adding to the current solution the one with the lowest level-value pair and repeating. The total cost of each of these steps is $O(\mathcal{A}\mathcal{L}(q))$, and there is a total of $O(q)$ steps, so the total running time of complete is $O(q\mathcal{A}\mathcal{L}(q))$.

As for parent, we can proceed by removing an element $v$ from the current solution and executing one step of complete. If this step adds $v$ back, then $v$ was not the parind and we need to remove another element, otherwise what we have left is the core of the solution and we just need to end the execution of complete. This algorithm has the same running time as the one above, i.e. $O(q\mathcal{A}\mathcal{L}(q))$.

We may summarize these results with the following lemma.

**Lemma 3.4.** *If we may find the level of a node in $\mathcal{L}(q)$ time, and we know that the elements that may be added to a partial solution can be iterated on using at most $O(\mathcal{A})$ total time, then both complete and parent may be implemented to run in $O(q\mathcal{A}\mathcal{L}(q))$ time, while using only $O(q)$ memory.*

# Chapter 4

# Applications of the framework

This chapter will show some applications of the algorithmic framework of Section 3.2 to problems that were presented in Section 2.2. In particular, the main focus will be on solving the restricted problem, as in many cases this is not trivial.

## 4.1 Cliques and (connected) $k$-plexes

Applying the framework to cliques is quite straightforward, since the restricted problem is very simple: if $C$ is a clique, then $C \cup \{v\}$ has $C$ and $C \cap N(v) \cup \{v\}$ as the only maximal cliques. As the restricted problem may be solved in polynomial time, the maximal clique problem may be solved with polynomial delay.

This is a well-known result, as seen for example in [13].

Solving the restricted problem for maximal $k$-plexes is not so simple: in fact, it may have an exponential number of solution. In [5], the authors find a polynomial-delay algorithm (assuming $k$ to be a constant) to generate all the solutions of the restricted problem.

Applying this solution of the restricted problem to our framework gives a polynomial total time algorithm to find all the maximal $k$-plexes in a graph, a weaker result than the one in [5]. Because the delay between two solutions of the restricted problem is not very useful in our framework, we will give an upper bound on the number of solutions of the restricted problem and an algorithm to compute them.

**Lemma 4.1.** *Assuming $k$ to be a constant, if $S$ is a $k$-plex and $v$ is a node in $V \setminus S$, then there are at most $1 + f(k)|S|^{k-1}$ maximal $k$-plexes in $S \cup \{v\}$, with $f(k) = (k-1)^{2k}$ for $k > 1$ and $f(1) = 1$. Moreover, they can be computed in $O(kf(k)|S|^k)$ time using only $O(kq)$ memory.*

*Proof.* During the proof, we will use the following upper bound on the sum of certain binomial coefficients, that can easily be derived from the fact that $\binom{n}{i} \leqslant n^i$:

$$\sum_{i=0}^{k} \binom{n}{i} \leqslant kn^k \qquad \forall n > 0$$

S is clearly a maximal k-plex in $S \cup \{v\}$. From now on, we will only consider maximal k-plexes that include $v$. Any of those maximal k-plexes may contain up to $k-1$ elements of S that are not neighbours of $v$. So, we can partition the maximal k-plexes that contain $v$ according to the set $\tilde{N}$ of non-neighbours of $v$ that they contain.

Consider now $K(\tilde{N}) = \{v\} \cup \tilde{N} \cup (N(v) \cap S)$. If it is not a k-plex, then it cannot be because of any element in $N(v)$: indeed, the number of non-neighbours that these elements have in $K(\tilde{N})$ is at most the number of non-neighbours that they have in S, and so it is at most k. Moreover, it cannot be due to $v$ either, as we know $|\tilde{N}| \leqslant k - 1$.

So any element x that breaks the k-plex constraint must be in $\tilde{N}$. Call $\tilde{N}_b$ the set of such xs. The constraint may only be broken by having $k + 1$ non-neighbours, as $v$ is the only element in $K(\tilde{N})$ that was not in S. Since we are only interested in k-plexes that contain $v$, the only way to fix this issue for any $x \in \tilde{N}_b$ is by removing from $K(\tilde{N})$ some element that belongs to $B(x) = \tilde{N}(v) \cap S \setminus N(x)$. As we need to do this for all the x in $\tilde{N}_b$, we need to find some set X such that:

- $X \subseteq \bigcup_{x \in \tilde{N}_b} B(x)$,

- $X \cup B(x) \neq \varnothing$ for any x in $\tilde{N}_b$,

- X is minimal among all the sets that satisfy the first two properties, as otherwise the resulting k-plex would not be maximal.

We will now prove that any minimal X satisfies $|X| \leqslant |\tilde{N}_b| \leqslant k - 1$. Let $f(y)$ be the function that maps any element of X into the set of $B(x)$ such that $B(x) \cap X = \{y\}$. We will now prove that if $|X| > |N_b|$ then there is an y for which $f(y) = \varnothing$. Indeed, if it were not the case, we would have that, since the various $f(y)$ are disjoint, $|\bigcup_{y \in X} f(y)| > |\tilde{N}_b|$. But this is a contraddiction, as the number of $B(x)$ is smaller than $|\tilde{N}_b|$ because of their definition. So there is an y such that $f(y) = \varnothing$. Removing such an y would not alter the number of sets covered by X, and so we have proven that X is not minimal.

So, we have that X must be a set of size at most $|\tilde{N}_b|$ contained in a set of size at most $|N_b| \times (k - 1) \leqslant (k - 1)^2$.

Taking these results together, we have that the number of maximal k-plexes in $S \cup \{v\}$ containing $v$ is at most

$$\sum_{i=0}^{k-1} \binom{|S|}{i} \sum_{i=0}^{k-1} \binom{(k-1)^2}{i} \leqslant |S|^{k-1}(k-1)^{2k}$$

Finally, since these sets can be easily enumerated in $O(|S|)$ time per set, it is enough to show that we may check if they are a k-plex and their maximality in $O(k|S|)$ time.

Indeed, given a preprocessing time of $O(|S|^2)$ (that reduces to $O(|S|)$ if $k = 1$), due to the computation of the set of non-neighbours for every node in S, we may check if a set is a k-plex in $O(k|S|)$ time; moreover we can check if extending a given k-plex with one extra node still gives a k-plex in $O(k)$ time per node (since we only need to check that node and its neighbours). This completes the proof.

$\square$

Unfortunately, there is no easy bound on the size of the candidate set for k-plexes (either for `complete` or for the choice of $v$ in Algorithm 4), so the resulting enumeration algorithm still has a time complexity of $O(n^2 q^{k+2} f(k)\alpha)$, which is obtained as the n (for the choice of $v$) times the number of solutions of the restricted problem $O(f(k)q^{k-1})$ times q (the number of choices for a seed in every solution) times the cost of running `complete` or `parent`, which according to Lemma 3.4 is given by the cost of checking the level of a node ($O(q)$ by updating it from one iteration to the next) times q times $A = O(n)$. Thus running `complete` takes $O(nq^2)$. Multiplying all these values together gives the cost of computing `children` as $O(n^2 q^{k+2} f(k))$. As this is basically the only cost in the reverse search, we obtain a delay of $O(n^2 q^{k+2} f(k))$ and thus the given total time.

A more interesting problem is the one of enumerating connected k-plexes. This problem was studied in [5] too, where a polynomial-delay algorithm for the restricted problem was given. The resulting algorithm using their framework takes incremental polynomial time but uses additional exponential memory. With our framework, we obtain a polynomial-delay algorithm (for fixed k) that only takes $O(kq)$ memory:

**Theorem 4.1.** *All the connected* k*-plexes in a graph* G *can be enumerated with* $O(q^{k+4}\Delta^2 f(k))$ *delay, using only* $O(kq)$ *extra memory (other than the memory used to store* G*).*

The theorem holds because we can restrict our choice of candidates to nodes that only have at least one neighbour in the current solution, thus replacing all the instances of n in the running time of the enumeration of k-plexes with $q\Delta$. Moreover, we may re-use Lemma 4.1 with the following simple observation: any connected k-plex must be contained in a "normal" k-plex, and may be obtained from it by extracting the connected component containing $v$. Moreover, duplicates may be avoided by choosing an easily-computable canonical k-plex that should generate a given maximal connected k-plex (for example, the one obtained by adding iteratively adding the smallest element in the input to the restricted problem that would keep the current set a k-plex). So, when we generate a connected k-plex C from a given k-plex K, we check if K is the canonical k-plex associated to C and, if it is not, we discard C, as we know that it will be (or has already been) generated when K is the canonical k-plex.

## 4.2 Black-connected cliques

In the case of black-connected cliques, the restricted problem $\mathcal{P}(S, v)$ is easy to solve: it is enough to find the only maximal clique in $S \cup \{v\}$ that is not $S$ and to remove any vertex from it that is not in the same black-connected component of $v$. This takes $O(q)$ time and gives exactly one solution.

Let us now consider `complete` and `parent`. We will denote by $\Delta_b$ the maximum "black degree" of a node in G. To execute a step of `complete`, it is enough to consider the elements of G that are black neighbours of a node in the current solution. There are at most $q\Delta_b$ such elements, and it is easy to check if one such element may be added to the solution in $O(q)$ time. Thanks to Lemma 3.4, the execution of both `complete` and `parent` takes at most $O(q^3\Delta_b)$ time.

Moreover, note that in Algorithm 4 we only need to consider nodes that are black neighbours of a node in the current solution as a possible $v$, as otherwise the only non-trivial solution of the restricted problem would be $\{v\}$ itself. With a calculation similar to the one that gave us the bound for connected k-plexes, this gives the following result:

**Theorem 4.2.** *All the maximal black-connected cliques in a graph* G *may be enumerated with* $O(q^5\Delta_b^2)$ *delay, using only* $O(q)$ *extra memory (other than the memory used to store* G*).*

## 4.3 Connected bipartite (induced) subgraphs

A bipartite subgraph can be equivalently seen as a graph that is two-colorable, or as a graph that contains no cycle with odd length.

We will first consider non-induced subgraphs. In this case, when solving the restricted problem we have a bipartite graph with one extra edge $e$ that connects two nodes of the same color, as otherwise the solution that generates the restricted problem would not be maximal. Any non-trivial subset of the edges in the restricted problem that gives a solution must then be missing enough edges to disconnect these two nodes in the original bipartite graph, i.e. the set of removed edge must be a cut in the bipartite graph. Moreover, since we are interested in maximal solutions, the cut must be minimal. Note that, since minimal cuts divide the graph into two connected components, and that adding $e$ to the graph connects these two components, the connectivity constraint is not an issue. Since enumeration of minimal cuts can be done in polynomial total time, as in [31], the restricted problem may be solved in polynomial total time.

For induced subgraphs, the extra node $v$ will have some neighbours of one color and some of the other. We will call these two sets of neighbours $B$ and $W$ respectively. As before, we need to break the odd cycles involving $v$. To do so, we need to remove one or more nodes from the original bipartite graph in such a way that either

- $B$ becomes empty, or

- *W* becomes empty, or

- there is no path from any node in B to any node in *W*

This can be done by enumerating the minimal vertex-cuts between two fake nodes b and *w*, where b (*w*) have all B (*W*) as neighbours (respectively). As before, removing minimal vertex-cuts leaves the graph split into two connected components, and since *v* is the node obtained by identifying b with *w*, we have that the resulting graph is connected. Moreover, minimal vertex-cuts can be enumerated in polynomial total time, as in [37].

Putting it all together, we obtain the following result.

**Theorem 4.3.** *All the maximal connected bipartite (induced) subgraphs in* G *may be enumerated in polynomial total time, using only polynomial space.*

## 4.4 Feedback vertex (arc) sets

As with bipartite subgraphs, the restricted problem requires one to break all the cycles that are formed in a graph with no cycles by adding a single node or a single arc. Thus, it may be solved by enumerating maximal cuts or vertex-cuts in the original graph. As vertex cuts in undirected graphs may be enumerated in polynomial total time (see [37]), and edge cuts may be enumerated in polynomial total time both in undirected and directed graphs (see [31]), we have the following result.

**Theorem 4.4.** *All the minimal feedback vertex sets in an undirected graph may be enumerated in polynomial total time, using only polynomial space. Moreover, all the minimal feedback arc sets in both directed and undirected graphs may be enumerated in polynomial total time, using only polynomial space.*

# Chapter 5

# Conclusions

This thesis analyzed known results regarding the enumeration of set systems, discussing various classes of set systems (independence systems, commutable set systems, (strongly) accessible set systems).

Despite the hardness result given in Theorem 2.1, many techniques have succeded in producing efficient algorithms, in practice if not in theory (for example [16]). In particular, we focused on binary partition, which was presented in Section 2.4, and on reverse search, initially presented in Section 2.4 and presented more in-depth in Section 3.2.1.

We then explained our new algorithmic framework, based on the technique of reverse search, that achieves under suitable conditions polynomial total time enumeration of commutable set systems while using a limited amount of memory. More precisely, our framework allows us to transform, in a space-efficient way, an algorithm for the enumeration of maximal solutions to a restricted problem that consists in enumerating the maximal solutions in a set obtained by adding an extra element to a given maximal solution.

The runtime performance of this algorithm influences heavily the running time of the algorithm obtained by applying our framework: a polynomial algorithm for the restricted problem gives a polynomial delay for the general enumeration, while a polynomial total time algorithm for the restricted problem translates to a polynomial total time algorithm for the general problem (this last implication clearly also holds in the other direction).

Thus, we studied in Chapter 4 how to solve the restricted problem in various commutable set systems. This analysis is by no means exhaustive, and a possible direction for future work would be to solve the restricted problem for more set systems, or even to find some general techniques related to restricted problems.

Other interesting extensions to this work could be related to the performance of the algorithm obtained. This could be done both on a theoretical aspect, by studying further properties of these families in order to lower the time complexity of the algorithm,

and on a practical aspect, as the algorithm should be easy to implement in a distributed manner. Indeed, we expect that the ability to easily implement a distributed version of the algorithm will make the biggest difference in practice between this framework and previous work.

# Bibliography

[1] N. K. Ahmed, J. Neville, R. A. Rossi, and N. Duffield. Efficient graphlet counting for large networks. In *ICDM*, pages 1–10. IEEE, 2015.

[2] R. Angles and C. Gutierrez. Survey of graph database models. *CSUR*, 40(1):1, 2008.

[3] D. Avis and K. Fukuda. Reverse search for enumeration. *Discrete Applied Mathematics*, 65(1-3):21–46, 1996.

[4] H. G. Barrow and R. M. Burstall. Subgraph isomorphism, matching relational structures and maximal cliques. *Information Processing Letters*, 4(4):83–84, 1976.

[5] D. Berlowitz, S. Cohen, and B. Kimelfeld. Efficient enumeration of maximal k-plexes. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, pages 431–444. ACM, 2015.

[6] M. Boley, T. Horváth, A. Poigné, and S. Wrobel. Listing closed sets of strongly accessible set systems with applications to data mining. *Theoretical Computer Science*, 411(3):691–700, 2010.

[7] C. Bron and J. Kerbosch. Algorithm 457: finding all cliques of an undirected graph. *Communications of the ACM*, 16(9):575–577, 1973.

[8] N. Chiba and T. Nishizeki. Arboricity and subgraph listing algorithms. *SIAM J Comput*, 14(1):210–223, 1985.

[9] G. Ciriello and C. Guerra. A review on models and algorithms for motif discovery in protein–protein interaction networks. *Brief Funct genomics & proteomics*, 7(2):147–156, 2008.

[10] S. Cohen, B. Kimelfeld, and Y. Sagiv. Generating all maximal induced subgraphs for hereditary and connected-hereditary graph properties. *Journal of Computer and System Sciences*, 74(7):1147–1159, 2008.

[11] S. Cohen and Y. Sagiv. An abstract framework for generating maximal answers to queries. In *ICDT*, pages 129–143. Springer, 2005.

[12] A. Conte, R. Grossi, A. Marino, L. Tattini, and L. Versari. A fast algorithm for large common connected induced subgraphs. *Algorithms for Computational Biology*, page 62.

[13] A. Conte, R. Grossi, A. Marino, and L. Versari. Sublinear-space bounded-delay enumeration for massive network analytics: Maximal cliques. In *43rd International Colloquium on Automata, Languages, and Programming (ICALP 2016)*, volume 148, pages 1–148, 2016.

[14] W. H. Day and D. Sankoff. Computational complexity of inferring phylogenies by compatibility. *Systematic Biology*, 35(2):224–229, 1986.

[15] N. Du, B. Wu, X. Pei, B. Wang, and L. Xu. Community detection in large-scale social networks. In *Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 workshop on Web mining and social network analysis*, pages 16–25. ACM, 2007.

[16] D. Eppstein and D. Strash. Listing all maximal cliques in large sparse real-world graphs. *Experimental Algorithms*, pages 364–375, 2011.

[17] L. A. Goldberg. *Efficient algorithms for listing combinatorial structures*, volume 5. Cambridge University Press, 1992.

[18] F. Harary and I. C. Ross. A procedure for clique detection using the group matrix. *Sociometry*, 20(3):205–215, 1957.

[19] R. Impagliazzo and R. Paturi. Complexity of k-SAT. In *Computational Complexity, 1999. Proceedings. Fourteenth Annual IEEE Conference on*, pages 237–240. IEEE, 1999.

[20] D. S. Johnson, M. Yannakakis, and C. H. Papadimitriou. On generating all maximal independent sets. *Inform Process Lett*, 27(3):119 – 123, 1988.

[21] C. Klein, A. Marino, M.-F. Sagot, P. V. Milreu, and M. Brilli. Structural and dynamical analysis of biological networks. *Brief Funct Genomics*, 11(6):420–433, 2012.

[22] I. Koch. Enumerating all connected maximal common subgraphs in two graphs. *Theoretical Computer Science*, 250(1-2):1–30, 2001.

[23] F. S. Kuhl, G. M. Crippen, and D. K. Friesen. A combinatorial algorithm for calculating ligand binding. *Journal of Computational Chemistry*, 5(1):24–34, 1984.

[24] W. Kunhiro. Enumeration of enumeration algorithms and its complexity.

[25] V. Lacroix, L. Cottret, P. Thébault, and M.-F. Sagot. An introduction to metabolic networks and their structural analysis. *IEEE ACM T Comput Bi*, 5(4):594–617, 2008.

[26] E. L. Lawler, J. K. Lenstra, and A. Rinnooy Kan. Generating all maximal independent sets: NP-hardness and polynomial-time algorithms. *SIAM Journal on Computing*, 9(3):558–565, 1980.

[27] V. E. Lee, N. Ruan, R. Jin, and C. Aggarwal. A survey of algorithms for dense subgraph discovery. In *Managing and Mining Graph Data*, pages 303–336. Springer, 2010.

[28] P. V. Milreu, C. C. Klein, L. Cottret, V. Acuña, E. Birmelé, M. Borassi, C. Junot, A. Marchetti-Spaccamela, A. Marino, L. Stougie, et al. Telling metabolic stories to explore metabolomics data: a case study on the yeast response to cadmium exposure. *Bioinformatics*, 30(1):61–70, 2014.

[29] N. Modani and K. Dey. Large maximal cliques enumeration in sparse graphs. In *CIKM*, pages 1377–1378. ACM, 2008.

[30] R. J. Mokken. Cliques, clubs and clans. *Quality & Quantity*, 13(2):161–173, 1979.

[31] J. S. Provan and D. R. Shier. A paradigm for listing (s, t)-cuts in graphs. *Algorithmica*, 15(4):351–372, Apr 1996.

[32] R. C. Read. A survey of graph generation techniques. In *Combinatorial mathematics VIII*, pages 77–89. Springer, 1981.

[33] N. Rhodes, P. Willett, A. Calvet, J. B. Dunbar, and C. Humblet. Clip: similarity searching of 3d databases using clique detection. *Journal of chemical information and computer sciences*, 43(2):443–448, 2003.

[34] F. Ruskey. Combinatorial generation. *Preliminary working draft. University of Victoria, Victoria, BC, Canada*, 11:20, 2003.

[35] R. Samudrala and J. Moult. A graph-theoretic algorithm for comparative modeling of protein structure. *Journal of molecular biology*, 279(1):287–302, 1998.

[36] B. Schwikowski and E. Speckenmeyer. On enumerating all minimal solutions of feedback problems. *Discrete Applied Mathematics*, 117(1):253 – 265, 2002.

[37] H. Shen and W. Liang. Efficient enumeration of all minimal separators in a graph. *Theoretical Computer Science*, 180(1):169 – 180, 1997.

[38] D. Stanton and D. White. *Constructive combinatorics*. Undergraduate Texts in Mathematics (Springer-Verlag), 1986.

[39] U. Takeaki. Two general methods to reduce delay and change of enumeration algorithms. 2003.

[40] A. Tanay, R. Sharan, and R. Shamir. Discovering statistically significant biclusters in gene expression data. *Bioinformatics*, 18(suppl 1):S136–S144, 2002.

[41] E. Tomita, A. Tanaka, and H. Takahashi. The worst-case time complexity for generating all maximal cliques and computational experiments. *Theoretical Computer Science*, 363(1):28 – 42, 2006. Computing and Combinatorics.

[42] C. Tsourakakis, F. Bonchi, A. Gionis, F. Gullo, and M. Tsiarli. Denser than the densest subgraph: extracting optimal quasi-cliques with quality guarantees. In *SIGKDD*, pages 104–112. ACM, 2013.

[43] S. Tsukiyama, M. Ide, H. Ariyoshi, and I. Shirakawa. A new algorithm for generating all the maximal independent sets. *SIAM J Comput*, 6(3):505–517, 1977.

[44] L. G. Valiant. The complexity of computing the permanent. *Theor Comput Sci*, 8(2):189–201, 1979.

[45] S. Wasserman and K. Faust. *Social network analysis: Methods and applications*, volume 8. Cambridge university press, 1994.

[46] D. W. Williams, J. Huan, and W. Wang. Graph database indexing using structured graph decomposition. In *ICDE*, pages 976–985. IEEE, 2007.

[47] S. Yau. Generation of all hamiltonian circuits, paths, and centers of a graph, and related problems. *IEEE Transactions on Circuit Theory*, 14(1):79–81, 1967.

# Appendix A

# Code for the framework

This appendix contains a simple implementation of the enumeration framework described in this thesis. It contains the generic implementations of the functions described in Chapter 3, plus some specific implementations of the functions needed to apply the framework to black-connected cliques in graphs obtained by computing the product of two connected graphs.

The implementation provided here is not as memory-efficient as the one described in the rest of the work, but improves a bit the running time of the algorithm. More in detail, for black-connected cliques, we obtain an algorithm that has a delay of $O(q^4\Delta_b^2)$ while using $O(q\Delta_b)$ memory per recursive node. For simpicity, only the recursive version was implemented.

## A.1 Generic algorithm and graph data structures

### A.1.1 `framework.hpp`

This file contains the implementation of the functions defined in Chapter 3.

```
1   #ifndef FRAMEWORK_HPP
2   #define FRAMEWORK_HPP
3   #include <vector>
4   #include <cstdint>
5   #include <functional>
6   #include <algorithm>
7   #include <memory>
8   #include <set>
9   #include <unordered_map>
10  #include <unordered_set>
11
12  template <typename univ_t>
13  class CommutableSystem {
14  public:
15      const univ_t& e;
```

```cpp
16
17          CommutableSystem(const univ_t& e): e(e) {}
18
19      typedef typename univ_t::elem_t elem_t;
20      /**
21       * Checks if a given subset is a solution.
22       */
23      virtual bool is_good(const std::vector<elem_t>& s) = 0;
24
25      /**
26       * Solves the restricted problem
27       */
28      virtual void restricted_problem(
29          const std::vector<elem_t>& s,
30          elem_t v,
31          const std::function<bool(std::vector<elem_t>)>& cb
32      ) = 0;
33
34      /**
35       * Reports a solution
36       */
37      virtual void report_solution(const std::vector<elem_t>& s) = 0;
38
39
40      /**
41       * Checks if we can add a given element to a solution
42       */
43      virtual bool can_add(const std::vector<elem_t>& s, elem_t v) {
44          auto cnd = s;
45          cnd.push_back(v);
46          return is_good(cnd);
47      }
48
49      /**
50       * Returns true if the resticted problem may have at least two solutions.
51       */
52      virtual bool restr_multiple() {
53          return true;
54      }
55
56      /**
57       * Checks if the given element can be a valid seed of a solution,
58       * or a root if NULL is specified.
59       */
60      virtual bool is_seed(elem_t v, const std::unordered_set<elem_t>* s) {
61          return is_good({v});
62      }
63
64      /**
65       * Iterates over all the possible new elements that could be added
66       * because of a single new element in a solution.
```

```cpp
67          */
68      virtual void complete_cands(
69          const std::vector<elem_t>* ground_set,
70          elem_t new_elem,
71          const std::function<bool(elem_t)>& cb
72      ) {
73          if (!ground_set) {
74              for (elem_t i=0; i<e.size(); i++) {
75                  if (cb(i))
76                      break;
77              }
78          } else {
79              for (auto i: *ground_set) {
80                  if (cb(i))
81                      break;
82              }
83          }
84      }
85
86      /**
87       * Iterates over all the possible new elements that could be used
88       * for the restricted problem
89       */
90      virtual void restricted_cands(
91          const std::vector<elem_t>& s,
92          const std::vector<int32_t>& level,
93          const std::function<bool(elem_t)>& cb
94      ) {
95          auto ss = s;
96          std::sort(ss.begin(), ss.end());
97          for (elem_t i=0; i<e.size(); i++) {
98              if (std::binary_search(ss.begin(), ss.end(), i)) continue;
99              if (cb(i))
100                 break;
101         }
102     }
103
104     /**
105      * Checks if complete of a given element is a root.
106      */
107     virtual bool get_root(
108         elem_t v,
109         std::vector<elem_t>& s,
110         std::vector<int32_t>& level
111     ) {
112         if (!is_seed(v, nullptr)) return false;
113         s.clear();
114         level.clear();
115         s.push_back(v);
116         level.push_back(0);
117         auto ret = !complete(s, level, true);
```

```cpp
118          return ret;
119      }
120
121      /**
122       * Update candidate list when a new element is added to the solution.
123       */
124      virtual void update_step(
125          std::vector<elem_t>& s,
126          elem_t v,
127          int32_t level,
128          std::set<std::pair<int32_t, elem_t>>& candidates,
129          std::unordered_map<elem_t, int32_t>& cand_level,
130          const std::vector<elem_t>* ground_set
131      ) {
132          complete_cands(ground_set, v, [&](elem_t cnd) {
133              if (!can_add(s, cnd)) return false;
134              if (cand_level.count(cnd)) return false;
135              cand_level[cnd] = level+1;
136              candidates.emplace(level+1, cnd);
137              return false;
138          });
139      }
140
141      /**
142       * Extracts the next valid cand from candidates
143       */
144      virtual std::pair<elem_t, int32_t> next_cand(
145          const std::vector<elem_t>& s,
146          std::set<std::pair<int32_t, elem_t>>& candidates
147      ) {
148          while (!candidates.empty()) {
149              auto p = *candidates.begin();
150              candidates.erase(candidates.begin());
151              if (!can_add(s, p.second)) continue;
152              return {p.second, p.first};
153          }
154          return {e.size(), -1};
155      }
156
157      /**
158       * Recomputes the order and the level of the elements in s with another
        seed.
159       */
160      virtual void resort(
161          std::vector<elem_t>& s,
162          std::vector<int32_t>& level,
163          elem_t seed
164      ) {
165          std::vector<elem_t> sn{seed};
166          std::vector<int32_t> ln{0};
167          complete_inside(sn, ln, s, false);
```

```
168            s = sn;
169            level = ln;
170        }
171
172        /**
173         * Complete function. Returns true if there was a seed change, false
    ↪   otherwise
174         */
175        virtual bool complete(
176            std::vector<elem_t>& s,
177            std::vector<int32_t>& level,
178            bool stop_on_seed_change = false
179        ) {
180            if (s.empty()) throw std::runtime_error("??");
181            std::set<std::pair<int32_t, elem_t>> candidates;
182            std::unordered_map<elem_t, int32_t> cand_level;
183            for (uint32_t i=0; i<s.size(); i++) {
184                update_step(s, s[i], level[i], candidates, cand_level, nullptr);
185            }
186            bool seed_change = false;
187            while (true) {
188                elem_t n;
189                int32_t l;
190                std::tie(n, l) = next_cand(s, candidates);
191                if (n == e.size()) break;
192                unsigned pos = s.size();
193                while (pos > 0 && (l < level[pos-1] || (l==level[pos-1] && n <
                    ↪   s[pos-1]))) pos--;
194                s.insert(s.begin()+pos, n);
195                level.insert(level.begin()+pos, l);
196                if (n < s[0]) { // Seed change
197                    if (stop_on_seed_change) return true;
198                    seed_change = true;
199                    resort(s, level, n);
200                    cand_level.clear();
201                    candidates.clear();
202                    for (uint32_t i=0; i<s.size(); i++) {
203                        update_step(s, s[i], level[i], candidates, cand_level,
                            ↪   nullptr);
204                    }
205                } else {
206                    update_step(s, n, l, candidates, cand_level, nullptr);
207                }
208            }
209            return seed_change;
210        }
211
212        /**
213         * Runs complete inside a given set.
214         */
215        virtual void complete_inside(
```

```
216            std::vector<elem_t>& s,
217            std::vector<int32_t>& level,
218            const std::vector<elem_t>& inside,
219            bool change_seed = true
220        ) {
221            if (s.empty()) throw std::runtime_error("??");
222            std::set<std::pair<int32_t, elem_t>> candidates;
223            std::unordered_map<elem_t, int32_t> cand_level;
224            for (uint32_t i=0; i<s.size(); i++) {
225                update_step(s, s[i], level[i], candidates, cand_level, &inside);
226            }
227            while (true) {
228                elem_t n;
229                int32_t l;
230                std::tie(n, l) = next_cand(s, candidates);
231                if (n == e.size()) break;
232                unsigned pos = s.size();
233                while (pos > 0 && (l < level[pos-1] || (l==level[pos-1] && n <
             ↪  s[pos-1]))) pos--;
234                s.insert(s.begin()+pos, n);
235                level.insert(level.begin()+pos, l);
236                if (n < s[0] && change_seed) { // Seed change
237                    resort(s, level, n);
238                    cand_level.clear();
239                    candidates.clear();
240                    for (uint32_t i=0; i<s.size(); i++) {
241                        update_step(s, s[i], level[i], candidates, cand_level,
                 ↪  &inside);
242                    }
243                } else {
244                    update_step(s, n, l, candidates, cand_level, &inside);
245                }
246            }
247        }
248
249        /**
250         * Computes the prefix of the solution with a given seed and ending with
    ↪  v
251         */
252        virtual void get_prefix(
253            std::vector<elem_t>& s,
254            std::vector<int32_t>& level,
255            elem_t seed,
256            elem_t v
257        ) {
258            resort(s, level, seed);
259            std::size_t i;
260            for (i=0; i<s.size(); i++)
261                if (s[i] == v)
262                    break;
263            s.resize(i+1);
```

```
264                  level.resize(i+1);
265          }
266
267      /**
268       * Parent function, returns the parent index.
269       */
270      virtual elem_t parent(
271          const std::vector<elem_t>& s,
272          const std::vector<int32_t>& level,
273          std::vector<elem_t>& parent,
274          std::vector<int32_t>& parent_level
275      ) {
276          for (unsigned parind_pos = s.size()-1; parind_pos > 0; parind_pos--)
             ↪ {
277              parent = s;
278              parent_level = level;
279              parent.resize(parind_pos);
280              parent_level.resize(parind_pos);
281              complete(parent, parent_level);
282              if (parent != s) {
283                  return s[parind_pos];
284              }
285          }
286          parent.clear();
287          return e.size();
288      }
289
290      /**
291       * Computes the children of a given solution. Returns true if we stopped
     ↪ generating
292       * them because the callback returned true.
293       */
294      virtual bool children(
295          const std::vector<elem_t>& s,
296          const std::vector<int32_t>& level,
297          const std::function<bool(const std::vector<elem_t>&, const
             ↪ std::vector<int32_t>&)>& cb
298      ) {
299          bool done = false;
300          restricted_cands(s, level, [&] (elem_t cand) {
301              restricted_problem(s, cand, [&](const std::vector<elem_t>& sol) {
302                  std::unordered_set<elem_t> sol_set(sol.begin(), sol.end());
303                  for (auto seed: sol) {
304                      if (!is_seed(seed, &sol_set)) continue;
305                      if (cand <= seed) continue;
306                      std::vector<elem_t> core = sol;
307                      std::vector<int32_t> clvl = level;
308                      get_prefix(core, clvl, seed, cand);
309                      std::vector<elem_t> child = core;
310                      std::vector<int32_t> lvl = clvl;
311                      // There was a seed change
```

```cpp
312                    if (complete(child, lvl, true)) continue;
313                    std::vector<elem_t> p;
314                    std::vector<int32_t> plvl;
315                    elem_t parind = parent(child, lvl, p, plvl);
316                    // Not the parent of this child
317                    if (p != s) continue;
318                    // Wrong parent index
319                    if (parind != cand) continue;
320                    if (restr_multiple()) {
321                        p.push_back(parind);
322                        complete_inside(core, clvl, p);
323                        // Wrong restricted problem solution
324                        if (core != sol) continue;
325                    }
326                    if (cb(child, lvl)) {
327                        done = true;
328                        break;
329                    }
330                }
331                return done;
332            });
333            return done;
334        });
335        return done;
336    }
337 };
338
339 template<typename CS>
340 class ReverseSearch {
341 protected:
342     std::unique_ptr<CS> cs;
343     typedef typename CS::elem_t elem_t;
344     virtual void handle_solution(const std::vector<elem_t>& s, const
      ↪  std::vector<int32_t>& level) {
345         cs->report_solution(s);
346         cs->children(s, level, [this](const std::vector<elem_t>& sn, const
           ↪  std::vector<int32_t>& leveln) {
347             handle_solution(sn, leveln);
348             return false;
349         });
350     }
351 public:
352     template <typename... Args>
353     ReverseSearch(const Args&... args): cs(std::make_unique<CS>(args...)) {}
354
355     void run() {
356         std::vector<elem_t> s, p;
357         std::vector<int32_t> level, pl;
358         for (elem_t i=0; i<cs->e.size(); i++) {
359             if (cs->get_root(i, s, level)) {
360                 handle_solution(s, level);
```

```
361                     }
362                 }
363           }
364     };
365
366     #endif
```

### A.1.2 `graph.hpp`

This file contains the implementation of the graph data structures, including the product graph as defined in [12]. This product graph is the graph with colored edges in which black connected cliques are enumerated.

```
1    #ifndef GRAPH_HPP
2    #define GRAPH_HPP
3    #include <stdint.h>
4    #include <assert.h>
5    #include "dynarray.hpp"
6    #include "cuckoo.hpp"
7    #include "binary_search.hpp"
8    #include "common.hpp"
9
10   #include <vector>
11   #include <algorithm>
12   #include <unordered_map>
13   #include <functional>
14
15   template <typename node_t_ = uint32_t, bool lowmem = false>
16   class graph_t {
17   public:
18       typedef node_t_ node_t;
19       typedef node_t elem_t;
20   private:
21       node_t N_;
22       dynarray<binary_search_t<node_t>> edges;
23       // Structures for the "fast" version
24       dynarray<cuckoo_hash_set<node_t>> edges_fast;
25       dynarray<typename binary_search_t<node_t>::iterator> fwd_iter;
26
27   protected:
28       static int64_t nextInt(FILE* in) {
29           int64_t n = 0;
30           int64_t ch = getc_unlocked(in);
31           while (ch != EOF && (ch < '0' || ch > '9')) ch = getc_unlocked(in);
32               if (ch == EOF) return EOF;
33               while (ch >= '0' && ch <= '9') {
34               n = 10*n + ch - '0';
35               ch = getc_unlocked(in);
36           }
37           return n;
```

```cpp
38          }
39      static void read_edges(FILE* in, bool directed,
        ↪  std::vector<std::vector<node_t>>& graph) {
40          while(true) {
41              int a = nextInt(in);
42              int b = nextInt(in);
43              if (a == EOF || b == EOF) return;
44              if(a == b) continue;
45              graph[a].push_back(b);
46              if (!directed) graph[b].push_back(a);
47          }
48      }
49  public:
50      graph_t(node_t N, const std::vector<std::vector<node_t>>& edg, bool
        ↪  sorted = false): N_(N) {
51          edges.resize(N);
52          for (node_t i=0; i<N; i++) {
53              edges[i].init(edg[i]);
54              if (!sorted) std::sort(edges[i].data().begin(),
                ↪  edges[i].data().end());
55          }
56          if (!lowmem) {
57              edges_fast.resize(N);
58              fwd_iter.resize(N);
59              for (node_t i=0; i<N; i++) {
60                  for (auto x: edg[i])
61                      edges_fast[i].insert(x);
62                  fwd_iter[i] = edges[i].upper_bound(i);
63              }
64          }
65      }
66
67      static graph_t read_oly(FILE* in = stdin, bool directed = false) {
68          node_t N = nextInt(in);
69          nextInt(in);
70          std::vector<std::vector<node_t>> graph(N);
71          read_edges(in, directed, graph);
72          for (node_t i=0; i<N; i++) {
73              sort(graph[i].begin(), graph[i].end());
74              graph[i].erase(unique(graph[i].begin(), graph[i].end()),
                ↪  graph[i].end());
75          }
76          return {N, graph, true};
77      }
78
79      static graph_t read_nde(FILE* in = stdin, bool directed = false) {
80          node_t N = nextInt(in);
81          std::vector<std::vector<node_t>> graph(N);
82          for(node_t i=0; i<N; i++) {
83              int a = nextInt(in);
84              int b = nextInt(in);
```

```
85          graph[a].reserve(b);
86      }
87      read_edges(in, directed, graph);
88      for (node_t i=0; i<N; i++) {
89          sort(graph[i].begin(), graph[i].end());
90          graph[i].erase(unique(graph[i].begin(), graph[i].end()),
            ↪  graph[i].end());
91      }
92      return {N, graph, true};
93  }
94
95  /**
96   *  Node new_order[i] will go in position i.
97   */
98  template <typename G>
99  static graph_t permute(G g, const std::vector<node_t>& new_order) {
100     assert(new_order.size() == (size_t) g.size());
101     std::vector<node_t> new_pos(g.size(), -1);
102     for (node_t i=0; i<g.size(); i++) new_pos[new_order[i]] = i;
103     std::vector<std::vector<node_t>> new_edges(g.size());
104     for (node_t i=0; i<g.size(); i++) {
105         for (auto x: g.neighs(i)) {
106             new_edges[new_pos[i]].push_back(new_pos[x]);
107         }
108     }
109     return {g.size(), new_edges};
110 }
111
112 node_t size() const {
113     return N_;
114 }
115
116 node_t degree(node_t i) const {
117     return edges[i].size();
118 }
119
120 const binary_search_t<node_t>& neighs(node_t i) const {
121     return edges[i];
122 }
123
124 class fwd_neighs_t {
125     node_t n;
126     const graph_t* g;
127 public:
128     fwd_neighs_t(const graph_t* g, node_t n): g(g), n(n) {}
129     typename binary_search_t<node_t>::iterator begin() const {
130         if (lowmem) return g->edges[n].upper_bound(n);
131         else return g->fwd_iter[n];
132     }
133     typename binary_search_t<node_t>::iterator end() const {
134         return g->edges[n].end();
```

```cpp
135            }
136        node_t size() const {
137            return end() - begin();
138        }
139        };
140
141        friend class fwd_neighs_t;
142
143        const fwd_neighs_t fwd_neighs(node_t n) const {
144            return {this, n};
145        }
146
147        node_t fwd_degree(node_t n) const {
148            return fwd_neighs(n).size();
149        }
150
151        bool are_neighs(node_t a, node_t b) const {
152            if (lowmem) return edges[a].count(b);
153            else return edges_fast[a].count(b);
154        }
155    };
156
157    template <typename label_t_ = uint32_t, typename node_t_ = uint32_t, bool
        ↪   lowmem = false>
158    class labeled_graph_t: public graph_t<node_t_, lowmem> {
159        std::vector<label_t_> labels;
160        using ugraph_t = graph_t<node_t_, lowmem>;
161    public:
162        typedef label_t_ label_t;
163        typedef typename ugraph_t::node_t node_t;
164        typedef typename ugraph_t::elem_t elem_t;
165
166        labeled_graph_t(
167            node_t N, std::vector<label_t> labels,
168            const std::vector<std::vector<node_t>>& edg,
169            bool sorted = false
170        ): ugraph_t(N, edg, sorted), labels(labels) {}
171
172        static labeled_graph_t read_oly(FILE* in = stdin, bool directed = false)
        ↪   {
173            node_t N = ugraph_t::nextInt(in);
174            ugraph_t::nextInt(in);
175            std::vector<std::vector<node_t>> graph(N);
176            std::vector<label_t> labels(N);
177            for (node_t i=0; i<N; i++) {
178                labels[i] = ugraph_t::nextInt(in);
179            }
180            ugraph_t::read_edges(in, directed, graph);
181            for (node_t i=0; i<N; i++) {
182                sort(graph[i].begin(), graph[i].end());
```

```
183              graph[i].erase(unique(graph[i].begin(), graph[i].end()),
                 ↪  graph[i].end());
184          }
185          return {N, labels, graph, true};
186      }
187
188      label_t get_label(node_t node) const {
189          return labels[node];
190      }
191  };
192
193  template <typename label_t_ = uint32_t, typename node_t_ = uint32_t, bool
     ↪  lowmem = false>
194  class product_graph_t {
195      typedef label_t_ label_t;
196      using ugraph_t = graph_t<node_t_, lowmem>;
197      using lgraph_t = labeled_graph_t<label_t, node_t_, lowmem>;
198      lgraph_t g1;
199      lgraph_t g2;
200      std::vector<std::pair<node_t_, node_t_>> nds;
201      std::unordered_map<std::pair<node_t_, node_t_>, node_t_, pair_hash> rmp;
202      void gen_node_list() {
203          std::unordered_map<label_t, std::vector<node_t_>> g2_nodes;
204          for (node_t_ i=0; i<g2.size(); i++)
205              g2_nodes[g2.get_label(i)].push_back(i);
206          for (node_t_ i=0; i<g1.size(); i++) {
207              for (auto second: g2_nodes[g1.get_label(i)]) {
208                  nds.emplace_back(i, second);
209                  rmp[nds.back()] = nds.size()-1;
210              }
211          }
212      }
213  public:
214      typedef typename ugraph_t::node_t node_t;
215      typedef typename ugraph_t::elem_t elem_t;
216      product_graph_t(lgraph_t&& g1, lgraph_t&& g2): g1(g1), g2(g2) {
217          gen_node_list();
218      }
219      product_graph_t(const lgraph_t& g1, const lgraph_t& g2): g1(g1), g2(g2) {
220          gen_node_list();
221      }
222
223      static product_graph_t read_oly(FILE* in1, FILE* in2, bool directed =
         ↪  false) {
224          return {lgraph_t::read_oly(in1, directed), lgraph_t::read_oly(in2,
             ↪  directed)};
225      }
226
227      node_t size() const {
228          return nds.size();
229      }
```

```
230
231        std::pair<node_t, node_t> to_pair(node_t node) const {
232            return nds[node];
233        }
234
235        void black_neighs(node_t node, const std::function<bool(node_t)>& cb)
         ↪  const {
236            for (auto a: g1.neighs(nds[node].first)) {
237                for (auto b: g2.neighs(nds[node].second)) {
238                    auto p = std::make_pair(a, b);
239                    if (rmp.count(p) == 0) continue;
240                    if (cb(rmp.at(p))) break;
241                }
242            }
243        }
244
245        bool are_neighs(node_t a, node_t b) const {
246            return nds[a].first != nds[b].first && nds[a].second != nds[b].second
247                && g1.are_neighs(nds[a].first, nds[b].first) ==
248                   g2.are_neighs(nds[a].second, nds[b].second);
249        }
250
251        bool are_black_neighs(node_t a, node_t b) const {
252            bool ans = g1.are_neighs(nds[a].first, nds[b].first) &&
253                   g2.are_neighs(nds[a].second, nds[b].second);
254            return ans;
255        }
256 };
257 #endif
```

## A.2   Black connected cliques

### A.2.1  bccliques.cpp

This file contains the implementation of the main function.

```
1  #include "graph.hpp"
2  #include "permute.hpp"
3  #include "bccliques.hpp"
4
5  int main(int argc, char** argv) {
6      if (argc < 3) {
7          fprintf(stderr, "Usage: %s g1 g2\n", argv[0]);
8          return 1;
9      }
10     FILE* f1 = fopen(argv[1], "r");
11     FILE* f2 = fopen(argv[2], "r");
12     auto tmp = product_graph_t<>::read_oly(f1, f2);
13     auto rs = ReverseSearch<BlackConnectedCliques<>>(tmp);
```

```
14        rs.run();
15    }
```

### A.2.2 `bccliques.hpp`

This file contains the functions that are needed by the framework to enumerate black-connected cliques.

```cpp
1   #ifndef BCCLIQUES_HPP
2   #define BCCLIQUES_HPP
3   #include "framework.hpp"
4   #include "graph.hpp"
5   #include "cuckoo.hpp"
6   #include <queue>
7
8   template <typename node_t = uint32_t>
9   class BlackConnectedCliques: public CommutableSystem<product_graph_t<node_t>>
    ↪ {
10  public:
11      using CommutableSystem<product_graph_t<node_t>>::CommutableSystem;
12      typedef typename CommutableSystem<product_graph_t<node_t>>::elem_t
        ↪ elem_t;
13      /**
14       * Checks if a given subset is a solution.
15       */
16      virtual bool is_good(const std::vector<elem_t>& s) override {
17          throw std::runtime_error("This function should never be called!");
18      }
19
20      /**
21       * Solves the restricted problem
22       */
23      virtual void restricted_problem(
24          const std::vector<elem_t>& s,
25          elem_t v,
26          const std::function<bool(std::vector<elem_t>)>& cb
27      ) override {
28          cuckoo_hash_set<elem_t> ok;
29          ok.insert(v);
30          for (auto n: s)
31              if (this->e.are_neighs(v, n))
32                  ok.insert(n);
33          std::vector<elem_t> sol;
34          cuckoo_hash_set<elem_t> visited;
35          std::queue<elem_t> q;
36          q.push(v);
37          while (!q.empty()) {
38              auto t = q.front(); q.pop();
39              if (!ok.count(t)) continue;
40              if (visited.count(t)) continue;
```

```cpp
41                 visited.insert(t);
42                 sol.push_back(t);
43                 this->e.black_neighs(t, [&](node_t n) -> bool {
44                     q.push(n);
45                     return false;
46                 });
47             }
48             cb(sol);
49         }
50
51         /**
52          * Reports a solution
53          */
54         virtual void report_solution(const std::vector<elem_t>& s) override {
55             printf("{");
56             for (auto n: s) printf("%u, ", this->e.to_pair(n).first);
57             printf("\b\b} -> {");
58             for (auto n: s) printf("%u, ", this->e.to_pair(n).second);;
59             printf("\b\b}\n");
60             fflush(stdout);
61         }
62
63         /**
64          * Checks if we can add a given element to a solution
65          */
66         virtual bool can_add(const std::vector<elem_t>& s, elem_t v) override {
67             uint32_t black_cnt = 0;
68             uint32_t neigh_cnt = 0;
69             for (auto n: s) {
70                 if (this->e.are_neighs(v, n)) neigh_cnt++;
71                 if (this->e.are_black_neighs(v, n)) black_cnt++;
72             }
73             return black_cnt > 0 && neigh_cnt == s.size();
74         }
75
76         /**
77          * Returns true if the resticted problem may have at least two solutions.
78          */
79         virtual bool restr_multiple() override {
80             return false;
81         }
82
83         /**
84          * Checks if the given element can be a valid seed of a solution,
85          * or a root if NULL is specified.
86          */
87         virtual bool is_seed(elem_t v, const std::unordered_set<elem_t>* s)
        ↪ override {
88             bool can_be = true;
89             this->e.black_neighs(v, [&] (elem_t e) {
90                 if (e > v) return true;
```

```cpp
 91             if (s == nullptr || s->count(e)) {
 92                 can_be = false;
 93                 return true;
 94             }
 95             return false;
 96         });
 97         return can_be;
 98     }
 99
100     /**
101      * Iterates over all the possible new elements that could be added
102      * because of a single new element in a solution.
103      */
104     virtual void complete_cands(
105         const std::vector<elem_t>* ground_set,
106         elem_t new_elem,
107         const std::function<bool(elem_t)>& cb
108     ) override {
109         if (!ground_set) {
110             this->e.black_neighs(new_elem, cb);
111         } else {
112             for (auto i: *ground_set) {
113                 if (cb(i))
114                     break;
115             }
116         }
117     }
118
119     /**
120      * Iterates over all the possible new elements that could be used
121      * for the restricted problem
122      */
123     virtual void restricted_cands(
124         const std::vector<elem_t>& s,
125         const std::vector<int32_t>& level,
126         const std::function<bool(elem_t)>& cb
127     ) override {
128         std::set<elem_t> els;
129         for (auto i: s) {
130             this->e.black_neighs(i, [&](elem_t v) {
131                 els.insert(v);
132                 return false;
133             });
134         }
135         for (auto i: s) els.erase(i);
136         for (auto i: els) {
137             if (cb(i))
138                 break;
139         }
140     }
141 };
```

```
142
143   #endif
```

## A.3   Data structures and other common things

### A.3.1   binary_search.hpp

This file contains a fast implementation of binary search.

```
1    #ifndef _BINARY_SEARCH_T
2    #define _BINARY_SEARCH_T
3    #include "dynarray.hpp"
4    #include <vector>
5
6    template<typename T = uint32_t>
7    class binary_search_t {
8    private:
9        dynarray<T> support;
10   public:
11       typedef dynarray<T>& data_type;
12       typedef const T* iterator;
13
14       void init(const std::vector<T>& v) {
15           support.resize(v.size());
16           unsigned cnt = 0;
17           while (cnt != v.size()) {
18               support[cnt] = v[cnt];
19               cnt++;
20           }
21       }
22
23       iterator begin() const {
24           return support.begin();
25       }
26
27       iterator it_at(size_t p) const {
28           return begin() + p;
29       }
30
31       iterator end() const {
32           return support.end();
33       }
34
35       size_t size() const {return support.size();}
36
37       T get_at(size_t idx) const {
38           return support[idx];
39       }
40
41       bool count(T v) const {
```

```cpp
42            int64_t n = support.size();
43            const T* arr = &support[0];
44            while (n > 1) {
45                const int64_t half = n/2;
46                arr = (arr[half] < v)?(arr+half):arr;
47                n -= half;
48            }
49            const T* tmp = (*arr < v)+arr;
50            return tmp < support.end() && *tmp == v;
51        }
52
53        iterator lower_bound(T v) const {
54            return std::lower_bound(support.begin(), support.end(), v);
55        }
56
57        iterator upper_bound(T v) const {
58            return std::upper_bound(support.begin(), support.end(), v);
59        }
60
61        data_type data() {
62            return support;
63        }
64    };
65    #endif
```

### A.3.2  `common.hpp`

This file contains an hash function for pairs.

```cpp
1    #ifndef COMMON_HPP
2    #define COMMON_HPP
3    #include <unordered_map>
4
5    struct pair_hash {
6        template <class T1, class T2>
7        std::size_t operator () (const std::pair<T1,T2> &p) const {
8            auto h1 = std::hash<T1>{}(p.first);
9            auto h2 = std::hash<T2>{}(p.second);
10           h1 ^= h2 + 0x9e3779b9 + (h1<<6) + (h1>>2);
11           return h1;
12       }
13   };
14   #endif
```

### A.3.3  `cuckoo.hpp`

This file contains a fast implementation of a cuckoo hash set.

```cpp
1    #ifndef _CUCKOO_H
2    #define _CUCKOO_H
```

```cpp
#include <vector>
#include <assert.h>
#include <immintrin.h>
#include <stdlib.h>
#include <stdint.h>

template<typename T, T missing = -1,
#ifdef __KNC__
    int bucket_size = 64/sizeof(T)
#else
    int bucket_size = 16/sizeof(T)
#endif
>
class cuckoo_hash_set {
public:
    typedef T value_type;
    typedef value_type& reference;
    typedef const value_type& const_reference;
    typedef value_type* pointer;
    typedef const value_type* const_pointer;
    typedef std::ptrdiff_t difference_type;
    typedef size_t size_type;
private:
    pointer ht;
    size_t mask;
    size_t sz;

    size_t hash_1(const value_type& k) const {
        return  k & mask;
    }

    size_t hash_2(const value_type& k) const {
        return ~k & mask;
    }

    void insert(value_type& k, value_type*& table) {
        int h1 = hash_1(k);
        for (int pos=0; pos<bucket_size; pos++)
            if (table[h1*bucket_size+pos] == missing) {
                table[h1*bucket_size+pos] = std::move(k);
                return;
            }
        int h2 = hash_2(k);
        for (int pos=0; pos<bucket_size; pos++)
            if (table[h2*bucket_size+pos] == missing) {
                table[h2*bucket_size+pos] = std::move(k);
                return;
            }
        bool use_hash_1 = true;
        for (unsigned i=0; i<mask; i++) {
            value_type cuckooed;
```

```
54              size_t hash;
55              if (use_hash_1) hash = hash_1(k);
56              else hash = hash_2(k);
57              int pos = 0;
58              for (; pos<bucket_size; pos++)
59                  if (table[hash*bucket_size+pos] == missing)
60                      break;
61              if (pos == bucket_size) {
62                  cuckooed = std::move(table[hash*bucket_size]);
63                  pos = 1;
64                  for (; pos<bucket_size; pos++)
65                      table[hash*bucket_size+pos-1] =
                        ↪   std::move(table[hash*bucket_size+pos]);
66                  table[hash*bucket_size+pos-1] = std::move(k);
67              } else {
68                  cuckooed = std::move(table[hash*bucket_size+pos]);
69                  table[hash*bucket_size+pos] = std::move(k);
70              }
71              use_hash_1 = hash == hash_2(cuckooed);
72              k = std::move(cuckooed);
73              if (k == missing) return;
74          }
75          rehash(table);
76          insert(k, table);
77      }
78
79      void rehash(value_type*& table) {
80          auto oldmask = mask;
81          if (mask == 0) mask = 1;
82          else mask = (mask<<1) | mask;
83          pointer newt = 0;
84          posix_memalign((void**)&newt, sizeof(T)*bucket_size,
            ↪   sizeof(T)*capacity());
85          std::fill(newt, newt+capacity(), missing);
86          for (size_t i=0; i<(oldmask+1)*bucket_size; i++)
87              if (table[i] != missing)
88                  insert(table[i], newt);
89          std::swap(table, newt);
90          free(newt);
91      }
92  public:
93      class const_iterator {
94      private:
95          const cuckoo_hash_set& container;
96          size_type offset;
97      public:
98          typedef cuckoo_hash_set::value_type value_type;
99          typedef cuckoo_hash_set::const_reference const_reference;
100         typedef cuckoo_hash_set::const_pointer const_pointer;
101         typedef cuckoo_hash_set::difference_type difference_type;
102         typedef std::forward_iterator_tag iterator_category;
```

```
103
104         const_iterator(const cuckoo_hash_set& container, size_type offset_):
            ↪  container(container), offset(offset_) {
105             while (offset != container.capacity() && container.ht[offset] ==
                ↪  missing)
106                 ++offset;
107         }
108         bool operator==(const const_iterator& other) const {
109             return &container == &other.container && offset == other.offset;
110         }
111         bool operator!=(const const_iterator& other) const {
112             return !(*this == other);
113         }
114
115         const_iterator& operator++() {
116             ++offset;
117             while (offset != container.capacity() && container.ht[offset] ==
                ↪  missing)
118                 ++offset;
119             return *this;
120         }
121         const_iterator operator++(int) {
122             const_iterator tmp = *this;
123             ++*this;
124             return tmp;
125         }
126
127         const_reference operator*() const {
128             return container.ht[offset];
129         };
130     };
131
132     typedef const_iterator iterator;
133     friend class const_iterator;
134
135     cuckoo_hash_set& operator=(const cuckoo_hash_set& other) = delete;
136
137     cuckoo_hash_set(const cuckoo_hash_set& other): ht(nullptr) {
138         posix_memalign((void**)&ht, sizeof(T)*bucket_size,
            ↪  sizeof(T)*other.capacity());
139         for (uint64_t i=0; i<other.capacity(); i++)
140             ht[i] = other.ht[i];
141         mask = other.mask;
142         sz = other.sz;
143     };
144
145     ~cuckoo_hash_set() {
146         free(ht);
147     }
148
149     cuckoo_hash_set(): mask(0), sz(0) {
```

```cpp
150         posix_memalign((void**)&ht, sizeof(T)*bucket_size,
     ↪   sizeof(T)*bucket_size);
151         std::fill(ht, ht+bucket_size, missing);
152     }
153
154     const_iterator begin() const {
155         return const_iterator(*this, 0);
156     }
157     const_iterator end() const {
158         return const_iterator(*this, capacity());
159     }
160
161     bool operator==(const cuckoo_hash_set<T>& oth) {
162         if (oth.size() != size()) return false;
163         for (const auto& x: oth)
164             if (!count(x))
165                 return false;
166         return true;
167     }
168
169     bool operator!=(const cuckoo_hash_set<T>& oth) {
170         return !(*this == oth);
171     }
172
173     void insert(value_type k) {
174         if (count(k)) {
175             return;
176         }
177         insert(k, ht);
178         sz++;
179     }
180     bool count(const value_type& k) const {
181             int h1 = hash_1(k);
182         int h2 = hash_2(k);
183 #ifndef __KNC__
184         if (bucket_size == 4 && sizeof(T) == 4) {
185             __m128i cmp = _mm_set1_epi32(k);
186             __m128i b1 = _mm_load_si128((__m128i*)&ht[bucket_size*h1]);
187             __m128i b2 = _mm_load_si128((__m128i*)&ht[bucket_size*h2]);
188             __m128i flag = _mm_or_si128(_mm_cmpeq_epi32(cmp, b1),
     ↪   _mm_cmpeq_epi32(cmp, b2));
189             return _mm_movemask_epi8(flag);
190         }
191 #else
192         if (bucket_size == 16 && sizeof(T) = 4) {
193             __m512i cmp = _mm512_set1_epi32(k);
194             __m512i b1 = _mm512_load_epi32(&ht[bucket_size*h1]);
195             __m512i b2 = _mm512_load_epi32(&ht[bucket_size*h2]);
196             return _mm512_cmpeq_epi32_mask(b1, cmp) ||
     ↪   _mm512_cmpeq_epi32_mask(b2, cmp);
197         }
```

```cpp
198   #endif
199           bool result = false;
200           for (unsigned i=0; i<bucket_size; i++)
201               result |= (ht[(bucket_size*h1)|i] == k || ht[(bucket_size*h2)|i]
                  ↪   == k);
202           return result;
203       }
204       void reserve(size_type sz) {
205           if (sz <= capacity()) return;
206           mask++;
207           while (mask <= sz/bucket_size) mask <<= 1;
208           free(ht);
209           posix_memalign((void**)&ht, sizeof(T)*bucket_size,
              ↪   sizeof(T)*capacity());
210           std::fill(ht, ht+capacity(), missing);
211           mask--;
212       }
213       size_type size() const {
214           return sz;
215       }
216       size_type capacity() const {
217           return (mask+1)*bucket_size;
218       }
219       bool empty() const {
220           return sz == 0;
221       }
222       void erase(const value_type& k) {
223           int h1 = hash_1(k);
224           for (int pos=0; pos<bucket_size; pos++)
225               if (ht[h1*bucket_size+pos] == k) {
226                   ht[h1*bucket_size+pos] = missing;
227                   sz--;
228                   return;
229               }
230           int h2 = hash_2(k);
231           for (int pos=0; pos<bucket_size; pos++)
232               if (ht[h2*bucket_size+pos] == k) {
233                   ht[h2*bucket_size+pos] = missing;
234                   sz--;
235                   return;
236               }
237       }
238       void clear() {
239           std::fill(ht, ht+bucket_size*capacity(), missing);
240       }
241
242       int front() const {
243           return *begin();
244       }
245   };
246   #endif
```

### A.3.4 `dynarray.hpp`

This file contains a simple, fixed size array whose size is decided at construction time.

```cpp
#ifndef _DYNARRAY_H
#define _DYNARRAY_H
#include <cstddef>
#include <stdexcept>
#include <algorithm>
#include <memory>

template< class T >
struct dynarray {
    // types:
    typedef       T                                value_type;
    typedef       T&                               reference;
    typedef const T&                               const_reference;
    typedef       T*                               iterator;
    typedef const T*                               const_iterator;
    typedef std::reverse_iterator<iterator>        reverse_iterator;
    typedef std::reverse_iterator<const_iterator>  const_reverse_iterator;
    typedef size_t                                 size_type;
    typedef ptrdiff_t                              difference_type;

    // fields:
private:
    T*        store;
    size_type count;

    // helper functions:
    void check(size_type n) {
        if (store == 0) throw std::out_of_range("dynarray");
        if (n >= count) throw std::out_of_range("dynarray");
    }
    T* alloc(size_type n) {
        if (n>std::numeric_limits<size_type>::max()/sizeof(T))
            throw std::out_of_range("dynarray");
        return reinterpret_cast<T*>(malloc(n*sizeof(T)));
    }

    void init(size_t n) {
        size_type i;
        try {
            for (size_type i=0; i<count; ++i)
                new (store+i) T;
        } catch (...) {
            for (; i>0; --i)
                (store+(i-1))->~T();
            throw;
        }
    }
```

```
48
49      void init(size_t n, const value_type& v) {
50          size_type i;
51          try {
52              for (size_type i=0; i<count; ++i)
53                  new (store+i) T(v);
54          } catch (...) {
55              for (; i>0; --i)
56                  (store+(i-1))->~T();
57              throw;
58          }
59      }
60  public:
61      // construct and destruct:
62      dynarray(): store(nullptr), count(0) {};
63      const dynarray operator=(const dynarray&) = delete;
64
65      explicit dynarray(size_type c): store(alloc(c)), count(c) {
66          init(c);
67      }
68
69      explicit dynarray(size_type c, const value_type& v): store(alloc(c)),
      ↪   count(c) {
70          init(c, v);
71      }
72
73      dynarray(const dynarray& d): store(alloc(d.count)), count(d.count) {
74          try {
75              std::uninitialized_copy(d.begin(), d.end(), begin());
76          } catch (...) {
77              free(store);
78              throw;
79          }
80      }
81
82      ~dynarray() {
83          if (store == 0) return;
84          for (size_type i = 0; i<count; ++i)
85              (store+i)->~T();
86          free(store);
87      }
88
89      void resize(size_type n) {
90          this->~dynarray();
91          store = alloc(n);
92          count = n;
93          init(n);
94      }
95
96      void resize(size_type n, const value_type& v) {
97          this->~dynarray();
```

```
98          store = alloc(n);
99          count = n;
100         init(n, v);
101     }
102
103     // iterators:
104     iterator       begin()          { return store; }
105     const_iterator begin()   const { return store; }
106     const_iterator cbegin()  const { return store; }
107     iterator       end()             { return store + count; }
108     const_iterator end()     const { return store + count; }
109     const_iterator cend()    const { return store + count; }
110
111     reverse_iterator       rbegin()          { return reverse_iterator(end());
        ↪   }
112     const_reverse_iterator rbegin()  const { return reverse_iterator(end());
        ↪   }
113     reverse_iterator       rend()            { return
        ↪   reverse_iterator(begin()); }
114     const_reverse_iterator rend()     const { return
        ↪   reverse_iterator(begin()); }
115
116     // capacity:
117     size_type size()     const { return count; }
118     size_type max_size() const { return count; }
119     bool      empty()    const { return count == 0; }
120
121     // element access:
122     reference        operator[](size_type n)       { return store[n]; }
123     const_reference  operator[](size_type n) const { return store[n]; }
124
125     reference       front()       { return store[0]; }
126     const_reference front() const { return store[0]; }
127     reference       back()        { return store[count-1]; }
128     const_reference back()  const { return store[count-1]; }
129
130     const_reference at(size_type n) const { check(n); return store[n]; }
131     reference       at(size_type n)       { check(n); return store[n]; }
132
133     // data access:
134     T*       data()       { return store; }
135     const T* data() const { return store; }
136
137     friend void swap(dynarray<value_type>& a, dynarray<value_type>& b) {
138         std::swap(a.store, b.store);
139         std::swap(a.count, b.count);
140     }
141 };
142
143 #endif
```