


On the Maximum Colorful Arborescence Problem and Color Hierarchy Graph Structure

Guillaume Fertin¹

LS2N UMR CNRS 6004, Université de Nantes, Nantes, France

guillaume.fertin@univ-nantes.fr

 <https://orcid.org/0000-0002-8251-2012>

Julien Fradin²


LS2N UMR CNRS 6004, Université de Nantes, Nantes, France

julien.fradin@univ-nantes.fr

Christian Komusiewicz³

Fachbereich Mathematik und Informatik, Philipps-Universität Marburg, Marburg, Germany

komusiewicz@informatik.uni-marburg.de

 <https://orcid.org/0000-0003-0829-7032>

Abstract

Let $G = (V, A)$ be a vertex-colored arc-weighted directed acyclic graph (DAG) rooted in some vertex r . The color hierarchy graph $\mathcal{H}(G)$ of G is defined as follows: the vertex set of $\mathcal{H}(G)$ is the color set \mathcal{C} of G , and $\mathcal{H}(G)$ has an arc from c to c' if G has an arc from a vertex of color c to a vertex of color c' . We study the MAXIMUM COLORFUL ARBORESCENCE (MCA) problem, which takes as input a DAG G such that $\mathcal{H}(G)$ is also a DAG, and aims at finding in G a maximum-weight arborescence rooted in r in which no color appears more than once. The MCA problem models the *de novo* inference of unknown metabolites by mass spectrometry experiments. Although the problem has been introduced ten years ago (under a different name), it was only recently pointed out that a crucial additional property in the problem definition was missing: by essence, $\mathcal{H}(G)$ must be a DAG. In this paper, we further investigate MCA under this new light and provide new algorithmic results for this problem, with a focus on fixed-parameter tractability (FPT) issues for different structural parameters of $\mathcal{H}(G)$. In particular, we develop an $\mathcal{O}^*(3^{x_{\mathcal{H}}})$ -time algorithm for solving MCA, where $x_{\mathcal{H}}$ is the number of vertices of indegree at least two in $\mathcal{H}(G)$, thereby improving the $\mathcal{O}^*(3^{|\mathcal{C}|})$ -time algorithm of Böcker et al. [Proc. ECCB '08]. We also prove that MCA is W[2]-hard with respect to the treewidth $t_{\mathcal{H}}$ of the underlying undirected graph of $\mathcal{H}(G)$, and further show that it is FPT with respect to $t_{\mathcal{H}} + \ell_{\mathcal{C}}$, where $\ell_{\mathcal{C}} := |V| - |\mathcal{C}|$.

2012 ACM Subject Classification F.2.2 Nonnumerical Algorithms and Problems, G.2.1 Combinatorics, G.2.2 Graph Theory

Keywords and phrases Subgraph problem, computational complexity, algorithms, fixed-parameter tractability, kernelization

Digital Object Identifier 10.4230/LIPIcs.CPM.2018.17

1 Introduction

Motivated by *de novo* inference of metabolites from mass spectrometry experiments, Böcker et al. [4] introduced the MAXIMUM COLORFUL SUBTREE problem. This optimization prob-

¹ GF was partially supported by PHC PROCOPE number 37748TL.

² JF was partially supported by PHC PROCOPE number 37748TL.

³ CK was partially supported by the DFG, project MAGZ (KO 3669/4-1).



© Guillaume Fertin, Julien Fradin and Christian Komusiewicz;
licensed under Creative Commons License CC-BY

29th Annual Symposium on Combinatorial Pattern Matching (CPM 2018).

Editors: Gonzalo Navarro, David Sankoff, and Binhai Zhu; Article No. 17; pp. 17:1–17:15

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

lem takes as input a vertex-colored arc-weighted directed acyclic graph $G = (V, A)$ rooted in some vertex r , and asks for a maximum-weight colorful arborescence in G with root r . Herein, a vertex-colored graph or a vertex set is called *colorful* if the vertices have pairwise different colors and a directed graph G is an *arborescence* with root r if the underlying undirected graph of G is a tree and there is a path from r to every vertex in G . In this model, the root r in G represents the sought metabolite, any vertex in G represents a molecule obtained from r after (possibly several) fragmentation(s), and vertices are colored according to their masses. An arc connects two molecules (vertices) u and v when v can be obtained from u by fragmentation, and is assigned a weight that indicates the (possibly negative) degree of confidence that the fragmentation from u to v actually occurs. A maximum-weight colorful arborescence from G with root r thus represents a most plausible fragmentation scenario from r . Let $\mathcal{H}(G)$ be the following graph built from G : $V(\mathcal{H}(G))$ is the set \mathcal{C} of colors used to color $V(G)$, and there is an arc from c to c' in $\mathcal{H}(G)$ if there is an arc in G from a vertex of color c to a vertex of color c' . We call $\mathcal{H}(G)$ the *color hierarchy graph* of G . Observe that $\mathcal{H}(G)$ must be a DAG since colors represent masses and fragmenting a molecule gives new molecules with lower mass. As recently pointed out [14], the initial definition of MAXIMUM COLORFUL SUBTREE omits this crucial property of G . This led Fertin et al. [14] to reformulate the initial MAXIMUM COLORFUL SUBTREE problem as follows.

Maximum Colorful Arborescence (MCA)

Input: A DAG $G = (V, A)$ rooted in some vertex r , a set \mathcal{C} of colors, a coloring function $\text{col} : V \rightarrow \mathcal{C}$ such that $\mathcal{H}(G)$ is a DAG and an arc weight function $w : A \rightarrow \mathbb{R}$.

Output: A colorful arborescence $T = (V_T, A_T)$ rooted in r of maximum weight $w(T) := \sum_{a \in A_T} w(a)$.

The study of MCA initiated in [14] essentially focused on the particular case where G is an arborescence and showed for example that MCA is NP-hard even for very restricted such instances. This work was also the first one to explicitly exploit that $\mathcal{H}(G)$ is a DAG. In particular, it was shown that if $\mathcal{H}(G)$ is an arborescence, then MCA is polynomially solvable. This latter promising result is the starting point of the present paper, in which we aim at better understanding the structural parameters of $\mathcal{H}(G)$ that could lead to fixed-parameter tractable (FPT), *i.e.* exact and moderately exponential, algorithms. As pointed out in a recent study [12], obtaining exact solutions instead of approximate ones is indeed preferable for MCA. Hence, improved exact algorithms are truly desirable for this problem.

Related Work

The MCA problem is NP-hard and highly inapproximable even when G is an arborescence and every arc weight is equal to 1 [14]. Moreover, MCA is NP-hard even if $\ell_{\mathcal{C}} = 0$ where $\ell_{\mathcal{C}} := |V(G)| - |\mathcal{C}|$ [14] (a consequence of the proof of [19, Theorem 1]). On the positive side, MCA can be solved in $\mathcal{O}^*(3^{|\mathcal{C}|})$ time by dynamic programming [4]. Moreover, as previously mentioned, MCA is in P when $\mathcal{H}(G)$ is an arborescence [14]. This result can be extended to some arborescence-like color hierarchy graphs as MCA can be solved by a branching algorithm in time $\mathcal{O}^*(2^s)$ where s is the minimum number of arcs of \mathcal{H} whose removal turns \mathcal{H} into an arborescence [14].⁴ Finally, a solution of MCA of order k can be computed in $\mathcal{O}^*((3e)^k)$ time using the color-coding technique [1] in combination with dynamic programming [7].

⁴ The notation $\mathcal{O}^*(\cdot)$ does not take polynomial factors into account.

■ **Table 1** Overview of the results for the MCA problem presented in this paper. Here, $x_{\mathcal{H}}$ is the number of vertices of indegree at least two in \mathcal{H} , $t_{\mathcal{H}}$ is the treewidth of the underlying undirected graph of \mathcal{H} , $\ell_c := |V(G)| - |\mathcal{C}|$ and $\ell \geq \ell_c$ is the number of vertices that are not part of the solution.

Parameter	FPT status	Kernel status
$x_{\mathcal{H}}$	$\mathcal{O}^*(3^{x_{\mathcal{H}}})$ (Thm. 2.2)	No poly. kernel (Thm. 2.4)
ℓ	W[1]-hard (from [19])	
$x_{\mathcal{H}} + \ell_c$	FPT (from Thm. 2.2)	No poly. kernel (Thm. 2.7)
$x_{\mathcal{H}} + \ell$	Poly. kernel (Thm. 2.8)	
$t_{\mathcal{H}}$	W[2]-hard (Thm. 3.3)	
$t_{\mathcal{H}} + \ell_c$	$\mathcal{O}^*(2^{\ell_c} \cdot 4^{t_{\mathcal{H}}})$ (Thm. 3.7)	No poly. kernel (Cor. 3.8)

A related pattern matching problem in graphs is GRAPH MOTIF where, in its simplest version, we are given an undirected vertex-colored graph and ask whether there is a connected subgraph containing one vertex of each color [18, 13, 2, 3]. In contrast to MCA, GRAPH MOTIF is fixed-parameter tractable for the parameter ℓ_c [2, 15].

Our Contribution

Our results are summarized in Table 1. We focus on two parameters from $\mathcal{H}(G)$, namely its number $x_{\mathcal{H}}$ of vertices of indegree at least two, and the treewidth $t_{\mathcal{H}}$ of its underlying undirected graph. This choice is motivated by the fact that when $\mathcal{H}(G)$ is an arborescence, each of these two parameters is constant (namely, $x_{\mathcal{H}} = 0$ and $t_{\mathcal{H}} = 1$) and MCA is in P. Thus, our parameters measure the distance from this trivial case [16]. In addition, we consider the parameter $\ell_c := |V(G)| - |\mathcal{C}|$ and the parameter ℓ which is the number of vertices that are not part of a solution with a maximum number of vertices. More precisely, whenever we refer to the parameter ℓ we consider the problem variant where we are constrained to report the best arborescence among those with at least $|V| - \ell$ vertices. Intuitively, ℓ_c is the number of vertices that we need to delete just to obtain a colorful subgraph of G , and hence $\ell \geq \ell_c$. Observe that MCA is W[1]-hard parameterized by ℓ [14]; this is a consequence of the proof of [19, Theorem 1].

Together with FPT issues, we also address the (in)existence of polynomial problem kernels for these parameters. In a nutshell, we provide a complete dichotomy for fixed-parameter tractability and problem kernelization for these parameters.

Preliminaries

In the following, let $G = (V, A)$ be the input graph of MCA, with $n_G := |V(G)|$. For any integer p , we let $[p] := \{1, \dots, p\}$. For any vertex $v \in V$, $N^+(v)$ is the set of outneighbors of v . We say that a vertex v is *reachable* from another vertex $v' \in V(G)$ in a directed graph G if there exists a path from v' to v in G . The color hierarchy graph of G is denoted $\mathcal{H}(G) := (\mathcal{C}, A_{\mathcal{C}})$, or, when clear from the context, simply \mathcal{H} .

We briefly recall the relevant notions of parameterized algorithmics (see e.g. [8]). A parameterized problem is a subset of $\Sigma \times \mathbb{N}$ where the second component is the parameter. A parameterized problem is *fixed-parameter tractable* if every instance (x, k) can be solved in $f(k) \cdot |x|^{\mathcal{O}(1)}$ time. A *reduction to a problem kernel*, or *kernelization*, is an algorithm that takes as input an instance (x, k) of a parameterized problem Q and produces in polynomial time an equivalent (*i.e.*, having the same solution) instance (x', k') of Q such that (i) $|x'| \leq g(k)$, and (ii) $k' \leq k$. The instance (x', k') is called *problem kernel*, and g is called the *size of*

the problem kernel. If g is a polynomial function, then the problem admits a *polynomial-size kernel*. Classes $W[1]$ and $W[2]$ are classes of presumed fixed-parameter intractability: if a parameterized problem is $W[1]$ -hard or $W[2]$ -hard, then it is generally assumed that it is not fixed-parameter tractable.

This paper is organized as follows. In Section 2, we study in detail the impact of $x_{\mathcal{H}}$ on the parameterized complexity of the MCA problem, while in Section 3, the same type of study is realized with parameter $t_{\mathcal{H}}$.

2 Parameterizing the MCA Problem by $x_{\mathcal{H}}$

Two main reasons lead us to be particularly interested in $x_{\mathcal{H}}$, the number of vertices with indegree at least two in \mathcal{H} . First, MCA is in P when \mathcal{H} is an arborescence [14], thus when $x_{\mathcal{H}} = 0$. Second, MCA can be solved in $\mathcal{O}^*(3^{|\mathcal{C}|})$ time [4]. Since by definition $x_{\mathcal{H}} \leq |\mathcal{C}|$, determining whether MCA is FPT with respect to $x_{\mathcal{H}}$ is of particular interest. We answer this question positively in Theorem 2.2. We first need some additional definitions.

Let X be the set of vertices of indegree at least two in \mathcal{H} (thus $|X| = x_{\mathcal{H}}$) and call X the set of *difficult colors*. For any $V' \subseteq V(G)$, let $\text{col}(V')$ denote the set of colors used by col on the vertices in V' . Moreover, for any vertex $v \in V$ that has at least one outneighbor in G , assume that $\text{col}(N^+(v))$ has an arbitrary but fixed ordering. Therefore, for any $i \in [|\text{col}(N^+(v))|]$, we may let $\text{col}^+(v, i)$ denote the i th color in $\text{col}(N^+(v))$. Finally, for any arborescence T in G or in \mathcal{H} , let $X(T) := X \cap \text{col}(V(T))$ denote the set of difficult colors in T . We have the following lemma.

► **Lemma 2.1.** *Let T_1 and T_2 be two arborescences in \mathcal{H} such that T_1 is rooted in c_1 , T_2 is rooted in $c_2 \neq c_1$, and $c_1, c_2 \in N^+(c)$ for some $c \in \mathcal{C}$. If $X(T_1)$ and $X(T_2)$ are disjoint, then $V(T_1)$ and $V(T_2)$ are disjoint.*

Proof. Assume without loss of generality that c_1 is not reachable from c_2 in \mathcal{H} . If $V(T_1)$ and $V(T_2)$ are not disjoint, then there exists a color $c^* \in \mathcal{C}$ that belongs to T_1 and to T_2 . In order to prove that such a color c^* cannot exist, let τ_1 (resp. τ_2) be the set of colors on the path from c_1 (resp. c_2) to c^* including c_1 in T_1 (resp. c_2 in T_2). Then, either $\tau_2 \subset \tau_1$ or $c_2 \notin \tau_1$. First, if $\tau_2 \subset \tau_1$, then there exists a vertex $c' \in \tau_1$ such that $c' \neq c$ with an arc (c', c_2) . Since \mathcal{H} contains the arc (c, c_2) , the color c_2 is thus difficult. This contradicts the assumption that $X(T_1)$ and $X(T_2)$ are disjoint. Second, if $c_2 \notin \tau_1$, then $|\tau_1 \cap \tau_2| \geq 1$ since $c^* \in \tau_1 \cap \tau_2$. Therefore, let $\bar{c} \in \tau_1 \cap \tau_2$ such that there exists a path from \bar{c} to any other color of $\tau_1 \cap \tau_2$. By definition, the father of \bar{c} in τ_1 is different from the father of \bar{c} in τ_2 , which means that \bar{c} is a difficult color. This contradicts the assumption that $X(T_1)$ and $X(T_2)$ are disjoint. ◀

► **Theorem 2.2.** *MCA can be solved in $\mathcal{O}^*(3^{x_{\mathcal{H}}})$ time and $\mathcal{O}^*(2^{x_{\mathcal{H}}})$ space.*

Proof. We propose a dynamic programming algorithm which makes use of two tables. The first one, $A[v, X', i]$, is computed for all $v \in V(G)$, $X' \subseteq X$ and $i \in \{0\} \cup [|\text{col}(N^+(v))|]$ and stores the weight of a maximum colorful arborescence $T_A(v, X', i)$ in G such that

- $T_A(v, X', i)$ is rooted in v ,
- $(X(T_A(v, X', i)) \setminus \{\text{col}(v)\}) \subseteq X'$, and
- $T_A(v, X', i)$ contains an arc (v, u) only if $\text{col}(u) = \text{col}^+(v, j)$ for some $j \leq i$.

The second one, $B[v, X', i]$, is computed for all $v \in V$, $X' \subseteq X$ and $i \in [|\text{col}(N^+(v))|]$ and stores the weight of a maximum colorful arborescence $T_B(v, X', i)$ in G such that

- $T_B(v, X', i)$ is rooted in v ,

Algorithm 1 COMPUTING THE ENTRIES IN TABLES A AND B

```

for all  $v \in V$  from last to first in some topological ordering of  $G$  do
  for all  $X' \subseteq X$  do
    for all  $i \in \{1, \dots, |\text{col}(N^+(v))|\}$  do
      Compute  $B[v, X', i]$ 
    end for
  end for
  for all  $X' \subseteq X$  do
    for all  $i \in \{0, \dots, |\text{col}(N^+(v))|\}$  do
      Compute  $A[v, X', i]$ 
    end for
  end for
end for

```

- $(X(T_B(v, X', i)) \setminus \{\text{col}(v)\}) \subseteq X'$, and
- $T_B(v, X', i)$ contains an arc (v, u) only if $\text{col}(u) = \text{col}^+(v, i)$.

In a nutshell, $T_A(v, X', i)$ and $T_B(v, X', i)$ share the same root v and the same allowed set of difficult colors X' (disregarding $\text{col}(v)$), but $T_A(v, X', i)$ contains outneighbors of v up to color $\text{col}^+(i)$ and $T_B(v, X', i)$ contains at most one outneighbor of v which is of color $\text{col}^+(v, i)$. Hence, there is no $u \in N^+(v)$ such that $(v, u) \in T_A(v, X', i-1)$ and $(v, u) \in T_B(v, X', i)$. We now show how to compute the two abovementioned tables.

$$A[v, X', i] = \begin{cases} 0 & \text{if } i = 0, \\ \max_{X'' \subseteq X'} \{A[v, X'', i-1] + B[v, X' \setminus X'', i]\} & \text{otherwise.} \end{cases}$$

For an entry $A[v, X', i]$ with $i = 0$ note that $T_A(v, X', i)$ can only contain v . For $i > 0$, by definition there cannot exist any $u \in N^+(v)$ such that u belongs both to $T_A(v, X'', i-1)$ and $T_B(v, X' \setminus X'', i)$. Therefore, Lemma 2.1 shows that $\text{col}(v)$ is the only color occurring in $T_A(v, X'', i-1)$ and $T_B(v, X' \setminus X'', i)$. Thus, the union of $T_A(v, X'', i-1)$ and $T_B(v, X' \setminus X'', i)$ is a colorful arborescence. Finally, testing every possible $X'' \subseteq X'$ ensures the correctness of the formula.

$$B[v, X', i] = \begin{cases} 0 & \text{if } \text{col}^+(v, i) \in X \setminus X', \\ \max_{\substack{u \in N^+(v): \\ \text{col}(u) = \text{col}^+(v, i)}} \{0, w(v, u) + A[u, X', |\text{col}(N^+(u))|\]\} & \text{otherwise.} \end{cases}$$

For an entry of type $B[v, X', i]$, if $\text{col}^+(v, i)$ is a difficult color which does not belong to X' , then $V(T_B(v, X', i)) = \{v\}$, and hence $B[v, X', i] = 0$. Otherwise, recall that $B[v, X', i]$ stores the weight of a maximum colorful arborescence rooted in v containing at most one further vertex $u \in N^+(v)$ of color $\text{col}^+(v, i)$. Therefore, computing the maximum colorful arborescences for any such u and only keeping the best one if it is positive ensures the correctness of the formula.

Recall that any DAG has a topological ordering of its vertices, *i.e.* a linear ordering of its vertices such that for every arc (u, v) , u appears before v in this ordering. In Algorithm 1, we show how to compute all the entries of both dynamic programming tables. For this, we consider the entries from last to first according to some topological ordering of G . The total running time derives from the fact that our algorithm needs at most 3^{x_n} steps to compute $A[v, X', i]$ since a difficult color can be in X'' , $X' \setminus X''$ or in $X \setminus X'$. ◀

Recall that a parameterized problem Q is FPT with respect to a parameter k if and only if it has a kernelization algorithm for k [11], but that such a kernel is not necessarily polynomial. In Theorem 2.4, we prove that although MCA parameterized by $x_{\mathcal{H}}$ is FPT (as proved by Theorem 2.2), MCA is unlikely to admit a polynomial kernel for $x_{\mathcal{H}}$. For this, we use the or-composition technique which, roughly speaking, is a reduction that combines many instances of a problem into one instance of the problem Q . We first recall the definition of or-compositions.

► **Definition 2.3.** ([5]) An *or-composition* for a parameterized problem $Q \in \Sigma \times \mathbb{N}$ is an algorithm that receives as input a sequence $(x_1, k), (x_2, k), \dots, (x_t, k)$ with $(x_i, k) \in \Sigma \times \mathbb{N}$ for each $1 \leq i \leq t$, takes polynomial time in $\sum_{i=1}^t |x_i| + k$, and outputs $(y, k') \in \Sigma \times \mathbb{N}$ with $(y, k') \in Q$ if and only if $\exists_{1 \leq i \leq t} (x_i, k) \in Q$ and k' is polynomial in k .

If an NP-hard parameterized problem Q admits an or-composition, then Q does not admit any polynomial-size problem kernel (unless $\text{NP} \subseteq \text{coNP/Poly}$) [5]. Our or-composition actually shows that MCA is unlikely to admit a polynomial kernel for the parameter $|\mathcal{C}|$.

► **Theorem 2.4.** *Unless $\text{NP} \in \text{coNP/Poly}$, MCA does not admit a polynomial kernel for parameter $|\mathcal{C}|$, and consequently for parameter $x_{\mathcal{H}}$, even if G is an arborescence.*

Proof. In the following, let t be a positive integer. For any $i \in [t]$, let $G_i = (V_i, A_i)$ be the graph of an instance of MCA which is rooted in a vertex r_i and assume that the t instances are built on the same color set $\mathcal{C}' = \{c_1, \dots, c_{|\mathcal{C}'|}\}$, otherwise colors can be relabeled suitably.

We now compose the t instances of MCA into a new instance of MCA. Let $G = (V, A)$ be the graph of such a new instance with $V = \{r\} \cup \{r'_i : i \in [t]\} \cup \{v \in V_i : i \in [t]\}$ and $A = \{(r, r'_i) : i \in [t]\} \cup \{(r'_i, r_i) : i \in [t]\} \cup \{(u, v) \in A_i : i \in [t]\}$. Here, r is a vertex not contained in any of the t MCA instances and which has a path of length 2 towards the root r_i of any graph G_i ; thus G is clearly a DAG. Let \mathcal{C} be the color set of G , and let us define the coloring function on $V(G)$ as follows: the root r is assigned a unique color $c_r \notin \mathcal{C}'$; all vertices of type r'_i are assigned the same color $c_{r'} \notin (\mathcal{C}' \cup \{c_r\})$; all arcs of type (r'_i, r_i) and (r, r'_i) are given a weight of 0; the color (resp. weight) of all other vertices (resp. arcs) is the same in the new instance as in their initial instance. Clearly, $(G, \mathcal{C}, \text{col}, w, r)$ is a correct instance of MCA and $|\mathcal{C}| = |\mathcal{C}'| + 2$. Moreover, if G_i is an arborescence for every $i \in [t]$, then G is also an arborescence. We now prove that there exists $i \in [t]$ such that G_i has a colorful arborescence $T = (V_T, A_T)$ rooted in r_i of weight $W > 0$ if and only if G has a colorful arborescence $T' = (V_{T'}, A_{T'})$ rooted in r and of weight $W > 0$.

(\Rightarrow) If there exists $i \in [t]$ such that G_i has a colorful arborescence $T = (V_T, A_T)$ rooted in r_i and of weight $W > 0$, then let $T' = (V_{T'}, A_{T'})$ with $V_{T'} = V_T \cup \{r, r'_i\}$ and $A_{T'} = A_T \cup \{(r, r'_i), (r'_i, r_i)\}$. Clearly, T' is connected, colorful and of weight W .

(\Leftarrow) Suppose G contains a colorful arborescence $T' = (V_{T'}, A_{T'})$ with root r and weight $W > 0$. Since T' is colorful and all vertices of type r'_i share the same color, there cannot exist i and j in $[t]$, $v_i \in V_i$ and $v_j \in V_j$ such that both v_i and v_j belong to T' . Thus, let i^* be the only index in $[t]$ such that $V_{i^*} \cap V_{T'} \neq \emptyset$ and let $T = (V_T, A_T)$ with $V_T = V_{T'} \setminus \{r, r'_{i^*}\}$ and $A_T = A_{T'} \setminus \{(r, r'_{i^*}), (r'_{i^*}, r_{i^*})\}$. Clearly, T is connected, colorful and of weight W .

Now, recall that $|\mathcal{C}| = |\mathcal{C}'| + 2$ and thus that we made a correct composition of MCA into MCA. Moreover, recall that MCA is NP-hard [14] and that $x_{\mathcal{H}} \leq |\mathcal{C}|$. As a consequence, MCA does not admit a polynomial kernel for the parameter $|\mathcal{C}|$, and hence for the parameter $x_{\mathcal{H}}$, even in arborescences, unless $\text{NP} \subseteq \text{coNP/Poly}$. ◀

Recall that MCA can be solved in time $\mathcal{O}^*(2^s)$ where s is the minimum number of arcs needed to turn \mathcal{H} into an arborescence [14]. Since $s < |\mathcal{C}|^2$, we have the following.

► **Corollary 2.5.** *Unless $\text{NP} \in \text{coNP}/\text{Poly}$, MCA parameterized by s does not admit a polynomial kernel, even if G is an arborescence.*

In the following, we use a different technique, called polynomial parameter transformation [6], to show that MCA is also unlikely to admit a polynomial kernel for the parameter $x_{\mathcal{H}} + \ell_{\mathcal{C}}$, where $\ell_{\mathcal{C}} = n_G - |\mathcal{C}|$.

► **Definition 2.6.** ([6, 10, 9]) Let P and Q be two parameterized problems. We say that P is *polynomial parameter reducible* to Q if there exists a polynomial-time computable function $f : \Sigma^* \times \mathbb{N} \rightarrow \Sigma^* \times \mathbb{N}$ and a polynomial p , such that for all $(x, k) \in \Sigma^* \times \mathbb{N}$ the following holds: $(x, k) \in P$ if and only if $(x', k') = f(x, k) \in Q$, and $k' \leq p(k)$. The function f is called a *polynomial parameter transformation*.

If P is an NP-hard problem and Q belongs to NP, then a polynomial parameter transformation from P parameterized by k to Q parameterized by k' has the following consequence: if Q parameterized by k' admits a polynomial kernel, then P parameterized by k admits a polynomial kernel [6]. Using such a transformation, we obtain the following result.

► **Theorem 2.7.** *MCA parameterized by $x_{\mathcal{H}}$ does not admit a polynomial kernel unless $\text{NP} \subseteq \text{coNP}/\text{Poly}$ even when restricted to the special case where $\ell_{\mathcal{C}} = 0$.*

Proof. We reduce from SET COVER, which is defined as follows.

Set Cover

Input: A universe $\mathcal{U} = \{u_1, u_2, \dots, u_q\}$, a family $\mathcal{F} = \{S_1, S_2, \dots, S_p\}$ of subsets of \mathcal{U} , an integer k .

Output: A k -sized subfamily $\mathcal{S} \subseteq \mathcal{F}$ of sets whose union is \mathcal{U} .

The reduction is as follows: for any instance of SET COVER, we create a three-levels DAG $G = (V = V_1 \cup V_2 \cup V_3, A)$ with $V_1 = \{r\}$, $V_2 = \{v_i : i \in [p]\}$ and $V_3 = \{z_j : j \in [q]\}$. We call V_2 the *second level* of G and V_3 the *third level* of G . Informally, we associate one vertex at the second level to each set of \mathcal{F} and one vertex at the third level to each element of \mathcal{U} . There is an arc of weight -1 from r to each vertex at level 2 and an arc of weight p from v_i to z_j , for all $i \in [p]$ and $j \in [q]$ such that the element u_j is contained in the set S_i . Now, our coloring function col is as follows: give a unique color to each vertex of G . Notice that \mathcal{H} is also a three-levels DAG with $\text{col}(V_1)$, $\text{col}(V_2)$, and $\text{col}(V_3)$ at the first, second, and third levels, respectively. Therefore, the above construction is a correct instance of MCA. We now prove that there exists a k -sized subfamily $\mathcal{S} \subseteq \mathcal{F}$ of sets whose union is \mathcal{U} if and only if there exists a colorful arborescence T in G of weight $w(T) = pq - k$.

(\Rightarrow) Suppose there exists a k -sized subfamily $\mathcal{S} \subseteq \mathcal{F}$ of sets whose union is \mathcal{U} and let $\text{True} = \{i \in [p] : S_i \in \mathcal{S}\}$. Then, we set $V_T = \{r\} \cup \{v_i : i \in \text{True}\} \cup \{z_j : j \in [q]\}$. Necessarily, $G[V_T]$ is connected: first, r is connected to every level-2 vertex; second, a vertex z_j corresponds to an element u_j which is contained in some set $S_i \in \mathcal{S}$. Now, let T be a spanning arborescence of $G[V_T]$. Clearly, T is colorful and of weight $pq - k$.

(\Leftarrow) Suppose there exists a colorful arborescence $T = (V_T, A_T)$ in G of weight $w(T) = pq - k$. Notice that any arborescence T' in G which contains r and at least one vertex from V_3 must contain at least one vertex from V_2 in order to be connected. Therefore, if such an arborescence T' does not contain one vertex of type z_j , then $w(T') < pq - p - 1$ and $w(T') < w(T)$. Hence, if $w(T) = pq - k$ then T contains each vertex of the third level and T contains exactly k vertices at the second level. Now, let $\mathcal{S} = \{S_i : i \in [p] \text{ s.t. } v_i \in V_T\}$ and notice that \mathcal{S} is a k -sized subfamily of \mathcal{F} whose union is \mathcal{U} as all vertices of the third level belong to T . Our reduction is thus correct.

Now, recall that \mathcal{H} is a three-levels DAG with $\text{col}(V_1)$, $\text{col}(V_2)$, and $\text{col}(V_3)$ at the first, second and third levels, respectively. By construction of G , if there exists $c \in V(\mathcal{H})$ such that $d^-(c) \geq 1$, then $c \in \text{col}(V_3)$. Moreover, recall that $|\text{col}(V_3)| = |\mathcal{U}|$ and thus $x_{\mathcal{H}} \leq |\mathcal{U}|$. Thus we provided a correct polynomial parameter transformation from SET COVER parameterized by $|\mathcal{U}|$ to MCA parameterized by $x_{\mathcal{H}}$. Now, recall that SET COVER does not admit a polynomial kernel for $|\mathcal{U}|$ unless $\text{NP} \subseteq \text{coNP}/\text{Poly}$ [10] and that SET COVER is NP-hard [17]. Moreover, the decision version of MCA, which asks for a solution of weight at least k , clearly belongs to NP. Finally, observe that $\ell_{\mathcal{C}} = 0$ as G is colorful. As a consequence, MCA does not admit any polynomial kernel for $x_{\mathcal{H}}$ unless $\text{NP} \subseteq \text{coNP}/\text{Poly}$ even if $\ell_{\mathcal{C}} = 0$. ◀

Since $\ell \geq \ell_{\mathcal{C}}$, and in light of Theorem 2.7, we aim at determining whether a polynomial kernel exists for MCA parameterized $x_{\mathcal{H}} + \ell$. We have the following theorem.

► **Theorem 2.8.** *MCA admits a problem kernel with $\mathcal{O}(x_{\mathcal{H}} \cdot \ell^2)$ vertices.*

To show this result we provide three data reduction rules. To formulate the rules, we introduce some notation first.

For any vertex $v \in V(G)$, we define $G^+(v)$ as the subgraph of G that is induced by the set of vertices that are reachable from v in G (including v). Similarly, for any color $c \in V(\mathcal{H})$, we define $\mathcal{H}^+(c)$ as the subgraph of \mathcal{H} that is induced by the set of vertices that are reachable from c in \mathcal{H} (including c). We call a color c *autonomous* if (i) $\mathcal{H}^+(c)$ is an arborescence, and (ii) there does not exist an arc from a color $c_1 \notin \mathcal{H}^+(c)$ to a color $c_2 \in \mathcal{H}^+(c)$ in \mathcal{H} . For a vertex v , let T_v denote a maximum colorful arborescence in G that is rooted at v . Finally, for a color $c \in \mathcal{C}$, let $V_c := \{v \in V : \text{col}(v) = c\}$ denote the set of vertices with color c .

► **Reduction Rule 1.** *If an instance $(G, \mathcal{C}, \text{col}, w, r)$ of MCA contains an autonomous color c such that $\mathcal{H}^+(c)$ contains at least two vertices, then do the following.*

- For each vertex $v \in V_c$, compute the value $w(T_v)$ of T_v , and add $w(T_v)$ to the weight of each incoming arc of v .
- Remove from G all vertices that are reachable from a vertex in V_c , except the vertices of V_c .

► **Lemma 2.9.** *Reduction Rule 1 is correct and can be performed exhaustively in polynomial time.*

Proof. Consider a vertex $v \in V_c$. Since c is autonomous, $\mathcal{H}^+(c)$ is an arborescence and thus we may compute T_v which contains only colors from $\mathcal{H}^+(c)$ in polynomial time [14].

Now, we prove the correctness of the rule, that is, the original instance $(G, \mathcal{C}, \text{col}, w, r)$ has a colorful arborescence $T = (V_T, A_T)$ of weight at least W if and only if the new instance $(G', \mathcal{C}', \text{col}', w', r')$ has a colorful arborescence $T' = (V_{T'}, A_{T'})$ of weight at least W . We only show the forward direction of the equivalence; the converse can be seen by symmetric arguments. First, recall that c is autonomous. Therefore, if T does not contain any vertex of color c , then T does not contain any vertex whose color belongs to $V(\mathcal{H}^+(c))$ and we can trivially set $T' = T$. Otherwise, if T contains a vertex v of color c , then let $S_c \subseteq V_T$ be the set of vertices that are reachable from v in T . We now set $V_{T'} := (V_T \setminus S_c) \cup \{v\}$ and let $A_{T'}$ contain all the arcs from A_T that are not in $\mathcal{H}^+(c)$. Now, recall that we computed the weight $w(T_v)$ of the maximum colorful arborescence in G that was rooted in v and that $w'(v^-, v) = w(v^-, v) + w(T_v)$ where v^- is the inneighbor of v in T . This ensures that $w(T) \leq w'(T')$. ◀

In the following, for any vertices $v, v' \in V(G)$ such that v' is reachable from v in G , we denote $\pi(v, v')$ as the length of the maximum weighted path from v to v' in G .

► **Reduction Rule 2.** *If an instance $(G, \mathcal{C}, \text{col}, w, r)$ of MCA contains a triple $\{c_1, c_2, c_3\} \subseteq \mathcal{C}$ such that (i) c_1 is the unique inneighbor of c_2 , (ii) c_2 is the unique inneighbor of c_3 and (iii) c_3 is the unique outneighbor of c_2 , then do the following.*

- For any $v_1 \in V_{c_1}$ and $v_3 \in V_{c_3}$ such that there exists a path from v_1 to v_3 in G , create an arc (v_1, v_3) and set $w(v_1, v_3) := \pi(v_1, v_3)$.
- Add a vertex v^* of color c_3 and, for any vertex $v_1 \in V_{c_1}$ that has at least one outneighbor of color c_2 in G , add the arc (v_1, v^*) and set $w(v_1, v^*)$ to the highest weighted outgoing arc from v_1 to any vertex of color c_2 in G .
- Remove all vertices of V_{c_2} from G' .

► **Lemma 2.10.** *Reduction Rule 2 is correct and can be performed exhaustively in polynomial time.*

Proof. We first prove that our transformation is correct. We show only the direction that an arborescence of weight at least W in the original instance $(G, \mathcal{C}, \text{col}, w, r)$ implies an arborescence of weight at least W in the new instance $(G', \mathcal{C}', \text{col}', w', r')$; the converse direction can be shown by symmetric arguments. Let $T = (V_T, A_T)$ be a colorful arborescence of weight W in the original instance. First, if T does not contain a vertex of color c_2 , then T is an arborescence of the new instance. Second, if T contains a vertex v_2 of color c_2 whose inneighbor is v_1 in T and if T does not contain any vertex of color c_3 , then setting $V_{T'} := V_T \setminus \{v_2\} \cup \{v^*\}$ and $A_{T'} := A_T \setminus \{(v_1, v_2)\} \cup \{(v_1, v^*)\}$ gives an arborescence $T' = (V_{T'}, A_{T'})$ of the new instance. Moreover, $w(T) = w'(T')$ since $w(v_1, v_2) = w'(v_1, v^*)$. Third, if T contains a vertex v_2 of color c_2 whose inneighbor is v_1 in T and if T contains a vertex v_3 of color c_3 (whose inneighbor is necessarily v_2), then setting $V_{T'} := V_T \setminus \{v_2\}$ and $A_{T'} := A_T \setminus \{(v_1, v_2), (v_2, v_3)\} \cup \{(v_1, v_3)\}$ gives an arborescence $T' = (V_{T'}, A_{T'})$ of the new instance. Moreover, $w(T) = w'(T')$ since $w(v_1, v_2) + w(v_2, v_3) = w'(v_1, v_3)$.

The polynomial running time follows from the fact that $\pi(v_1, v_3)$ can be computed in polynomial time. ◀

To describe the final rule, let $N_{\bar{v}}^-(v)$ denote the set of unique colors in the inneighborhood of v in G , where a color c is *unique* if $|V_c| = 1$. Recall also that ℓ is the maximum number of vertices that do not belong to T in G .

► **Reduction Rule 3.** *If an instance $(G, \mathcal{C}, \text{col}, w, r)$ of MCA contains a vertex $v \in V$ such that $|N_{\bar{v}}^-(v)| > \ell + 1$, then delete the $|N_{\bar{v}}^-(v)| - \ell - 1$ least-weighted arcs from $N_{\bar{v}}^-(v)$ to v .*

► **Lemma 2.11.** *Reduction Rule 3 is correct and can be performed exhaustively in polynomial time.*

Proof. Since $|N_{\bar{v}}^-(v)| > \ell + 1$, T has to contain at least two vertices from $N_{\bar{v}}^-(v)$. Now, let v_1 be a vertex from $N_{\bar{v}}^-(v)$ such that (v_1, v) is the least-weighted incoming arc from a unique color to v in G . Even if v_1 belongs to T , there will always exist at least one other vertex v_2 that will also belong to T and such that $w(v_1, v) \leq w(v_2, v)$. Thus, we may assume that T does not contain the arc (v_1, v) and safely delete it. The correctness of the rule now follows from repeated application of this argument. ◀

We are now ready to prove Theorem 2.8.

Proof. The kernelization algorithm consists of the exhaustive application of Reduction Rules 1–3 in polynomial time. Let $(G, \mathcal{C}, \text{col}, w, r)$ denote the resulting equivalent instance and let $T = (V_T, A_T)$ be a solution of this instance. It remains to show that G has $\mathcal{O}(x_{\mathcal{H}} \cdot \ell^2)$ vertices. First, we show that the indegree of any color in \mathcal{H} is at most $(\ell + 1)^2 + \ell$. This will allow us to show, subsequently, the claimed bound on n_G .

Let us first bound the indegree of any color in \mathcal{H} . Since T is colorful and since $|V_T| = n_G - \ell$, there exist at most ℓ non-unique colors in \mathcal{C} and hence the inneighborhood of any color $c \in V(\mathcal{H})$ cannot contain more than ℓ non-unique colors in \mathcal{H} . Moreover, since the instance is reduced with respect to Reduction Rule 3, the inneighborhood of any vertex $v \in V(G)$ contains at most $\ell + 1$ vertices of unique color in G . Furthermore, we may assume $|V_c| \leq \ell + 1$ for any any color $c \in V(\mathcal{H})$ as T cannot be colorful if there exists more than $\ell + 1$ occurrences of c in G . As a consequence, for any color $c \in V(\mathcal{H})$, the inneighborhood of c cannot contain more than $|V_c| \cdot (\ell + 1) = (\ell + 1)^2$ unique colors in \mathcal{H} , and hence c has at most $(\ell + 1)^2 + \ell$ inneighbors.

Now, let F be the forest whose vertex set is $\mathcal{C}_F = \mathcal{C} \setminus X$ and which contains each arc (c, c') of \mathcal{H} such that $\{c, c'\} \subseteq \mathcal{C}_F$. In the following, we successively bound the maximum number of leaves of F , the maximum number of vertices of F , of $V(\mathcal{H})$ and finally of $V(G)$ in a function of ℓ and $x_{\mathcal{H}}$. First, recall that there does not exist any autonomous color $c \in \mathcal{C}$ to which Reduction Rule 1 applies. Thus, each leaf c of \mathcal{H} is in fact a difficult color. Consequently, every leaf of F is in \mathcal{H} an inneighbor of a difficult color. Since the maximum indegree of any color in \mathcal{H} is at most $(\ell + 1)^2 + \ell$, the number of leaves in F is at most $x_{\mathcal{H}}((\ell + 1)^2 + \ell)$. Now, by Lemma 2.10, \mathcal{H} does not contain any color which has a unique inneighbor and a unique outneighbor. As a consequence, F has no internal vertices of degree two that are not inneighbors of a difficult color. Hence, the number of nonleaves of F that are not inneighbors of a difficult color is $\mathcal{O}(x_{\mathcal{H}} \cdot \ell^2)$, and thus $|V(F)| = \mathcal{O}(x_{\mathcal{H}} \cdot \ell^2)$. Moreover, since $\mathcal{C}_F = \mathcal{C} \setminus X$, we have that $|\mathcal{C}| \leq x_{\mathcal{H}} + \mathcal{O}(x_{\mathcal{H}} \cdot \ell^2)$. Finally, the number of vertices in G can exceed the number of colors in \mathcal{H} by at most ℓ . Therefore, $n_G = \mathcal{O}(x_{\mathcal{H}} \cdot \ell^2)$ as claimed. ◀

3 Parameterizing the MCA Problem by the Treewidth of the Color Hierarchy Graph

Let $U(\mathcal{H})$ denote the underlying undirected graph of \mathcal{H} . In this section, we are interested in parameter $t_{\mathcal{H}}$, defined as the treewidth of $U(\mathcal{H})$. Indeed, since MCA is in P whenever \mathcal{H} is an arborescence [14], it is natural to study whether MCA parameterized by $t_{\mathcal{H}}$ is FPT. To do so, we first introduce some definitions.

► **Definition 3.1.** Let $G = (V, E)$ be a undirected graph. A tree decomposition of G is a pair $\langle \{X_i : i \in I\}, \mathcal{T} \rangle$, where \mathcal{T} is a tree whose vertex set is I , and each X_i is a subset of V , called a *bag*. The following three properties must hold:

1. $\cup_{i \in I} X_i = V$.
2. For every edge $(u, v) \in E$, there is an $i \in I$ such that $\{u, v\} \subseteq X_i$.
3. For all $i, j, k \in I$, if j lies on the path between i and k in \mathcal{T} , then $X_i \cap X_k \subseteq X_j$.

The *width* of $\langle \{X_i : i \in I\}, \mathcal{T} \rangle$ is defined as $\max\{|X_i| : i \in I\} - 1$, and the *treewidth* of G is the minimum k such that G admits a tree decomposition of width k .

► **Definition 3.2.** A tree decomposition $\langle \{X_i : i \in I\}, \mathcal{T} \rangle$ is called *nice* if the following conditions are satisfied:

1. Every node of \mathcal{T} has at most two children.

2. If a node i has two children j and k , then $X_i = X_j = X_k$ and in this case, X_i is called a JOIN NODE.
3. If a node i has one child j , then one of the following situations must hold:
 - a) $|X_i| = |X_j| + 1$ and $X_j \subset X_i$ and in this case, X_i is called an INTRODUCE NODE, or
 - b) $|X_i| = |X_j| - 1$ and $X_i \subset X_j$ and in this case, X_i is called a FORGET NODE.
4. If a node i has no child, then $|X_i| = 1$ and in this case, X_i is called a LEAF NODE.

We first show that MCA is unlikely to be FPT with respect to parameter $t_{\mathcal{H}}$.

► **Theorem 3.3.** *MCA parameterized by $t_{\mathcal{H}}$ is $W[2]$ -hard.*

Proof. We reduce from the k -MULTICOLORED SET COVER problem, which is defined below.

k -Multicolored Set Cover

Input: A universe $\mathcal{U} = \{u_1, u_2, \dots, u_q\}$, a family $\mathcal{F} = \{S_1, S_2, \dots, S_p\}$ of subsets of \mathcal{U} , a set of colors Λ with a coloring function $\text{col}' : \mathcal{F} \rightarrow \Lambda$, an integer k .

Output: A subfamily $\mathcal{S} \subseteq \mathcal{F}$ of sets whose union is \mathcal{U} , and such that (i) $|\mathcal{S}| = k$ and (ii) \mathcal{S} is colorful, i.e. $\text{col}'(S_i) \neq \text{col}'(S_j)$ for any $i \neq j$ such that $S_i, S_j \in \mathcal{S}$.

The reduction is as follows: for any instance of k -MULTICOLORED SET COVER, we create a three-level DAG $G = (V = V_1 \cup V_2 \cup V_3, A)$ with $V_1 = \{r\}$, $V_2 = \{v_i : i \in [p]\}$ and $V_3 = \{z_j : j \in [q]\}$. Informally, we associate a vertex at the second level to each set of \mathcal{F} and a vertex at the third level to each element of \mathcal{U} . We then add an arc of weight -1 from r to each vertex at level 2 and an arc of weight p from v_i to z_j , for all $i \in [p]$ and $j \in [q]$ such that $u_j \in S_i$. Now, our coloring function col is as follows: we give a unique color to each vertex in $V_1 \cup V_3$, while at the second level (thus in V_2), two vertices of type v_i are assigned the same color if and only if their two associated sets are assigned the same color by col' . Notice that \mathcal{H} is also a three-levels DAG with $\text{col}(V_1)$, $\text{col}(V_2)$, and $\text{col}(V_3)$ at the first, second and third levels, respectively. Therefore, $(G, \mathcal{C}, \text{col}, w, r)$ is a correct instance of MCA. We now prove that there exists a colorful set $\mathcal{S} \in \mathcal{F}$ of size k whose union is \mathcal{U} if and only if there exists a colorful arborescence T in G of weight $w(T) = pq - k$.

(\Rightarrow) Suppose there exists a colorful set $\mathcal{S} \in \mathcal{F}$ of size k whose union is \mathcal{U} and let $\text{True} = \{i \in [p] : S_i \in \mathcal{S}\}$. Let $V_T = \{r\} \cup \{v_i : i \in \text{True}\} \cup \{z_j : j \in [q]\}$. Necessarily, $G[V_T]$ is connected: first, r is connected to every level-2 vertex; second, a vertex z_j corresponds to an element u_j which is contained in some set $S_i \in \mathcal{S}$. Now, let T be a spanning arborescence of $G[V_T]$. Clearly, T is colorful and of weight $pq - k$.

(\Leftarrow) Suppose there exists a colorful arborescence $T = (V_T, A_T)$ in G of weight $w(T) = pq - k$. Notice that any arborescence T' in G which contains r and at least one vertex from V_3 must contain at least one vertex from V_2 in order to be connected. Therefore, if such an arborescence T' does not contain one vertex of type z_j , then $w(T') < pq - p - 1$ and $w(T') < w(T)$. Hence, if $w(T) = pq - k$ then T necessarily contains each vertex from V_3 , and thus contains exactly k vertices from V_2 . Now, let $\mathcal{S} = \{S_i : i \in [p] \text{ s.t. } v_i \in V_T\}$ and notice that \mathcal{S} is a colorful subfamily of size k whose union is \mathcal{U} as all vertices of the third level belong to T . Our reduction is thus correct.

Now, recall that \mathcal{H} is a three-levels DAG with resp. $\text{col}(V_1)$, $\text{col}(V_2)$ and $\text{col}(V_3)$ at the first, second and third levels. Thus, there exists a trivial tree decomposition $\langle \{X_i : i \in [|\text{col}(V_3)| + 2]\}, \mathcal{T} \rangle$ of $U(\mathcal{H})$ which is as follows: the bag $X_0 = \{\text{col}(r)\}$ has an arc towards the bag $X_1 = \{\{\text{col}(r)\} \cup \text{col}(V_2)\}$ and, for any $i \in [|\text{col}(V_3)|]$, there exists an arc from X_1 to X_i where each X_i contains $\text{col}(V_2)$ and a different vertex of $\text{col}(V_3)$. Consequently, the width of $\langle \{X_i : i \in [|\text{col}(V_3)| + 2]\}, \mathcal{T} \rangle$ is k , and hence MCA is $W[2]$ -hard parameterized by $t_{\mathcal{H}}$ as k -MULTICOLORED SET COVER is well-known to be $W[2]$ -hard parameterized by k . ◀

We now use the above proof to show that MCA is unlikely to admit FPT algorithms relatively for different further parameters related to \mathcal{H} . The vertex cover number of $U(\mathcal{H})$ is the size of a smallest subset $S \subseteq V(\mathcal{H})$ such that at least one incident vertex of any arc of \mathcal{H} belongs to S . Notice that $\text{col}(V_2)$ is a vertex cover of $U(\mathcal{H})$ and thus $U(\mathcal{H}) \leq k$. The *feedback vertex set* number is the size of a smallest subset $S \subseteq \mathcal{H}$ whose removal makes $U(\mathcal{H})$ acyclic. The size of such a subset S is an interesting parameter as $x_{\mathcal{H}} = 0$ in $\mathcal{H}[V(\mathcal{H}) \setminus S]$ and any vertex cover of $U(\mathcal{H})$ is also a feedback vertex set of $U(\mathcal{H})$ – hence, $\text{col}(V_2)$ is also a feedback vertex set of $U(\mathcal{H})$. Altogether, we thus obtain the following corollary.

► **Corollary 3.4.** *MCA parameterized by the vertex cover number of $U(\mathcal{H})$ or the feedback vertex set number of $U(\mathcal{H})$ is W[2]-hard.*

Next, recall that in the proof of Theorem 3.3 each color from the third level of \mathcal{H} is a leaf. Hence, the number of colors of outdegree at least two in \mathcal{H} is $|\text{col}(V_1)| + |\text{col}(V_2)| = k + 1$. Although Theorem 2.2 showed that MCA is FPT relatively to $x_{\mathcal{H}}$, we obtain the following.

► **Corollary 3.5.** *MCA parameterized by the number of colors of outdegree at least two in \mathcal{H} is W[2]-hard.*

By Theorem 3.3, MCA parameterized by $t_{\mathcal{H}}$ is W[2]-hard; thus, one may look for a parameter whose combination with $t_{\mathcal{H}}$ may lead to MCA being FPT. Here, we focus on parameter $\ell_{\mathcal{C}} = n_G - |\mathcal{C}|$. We know that MCA parameterized by $\ell_{\mathcal{C}}$ is W[1]-hard, but the problem can be solved in $\mathcal{O}^*(2^{\ell_{\mathcal{C}}})$ time when G is an arborescence [14]. Recall also that MCA is in P when \mathcal{H} is an arborescence [14], and hence when $t_{\mathcal{H}} = 1$. In the following, a *fully-colorful subgraph* of G is a subgraph of G that contains *exactly* one occurrence of each color $c \in \mathcal{C}$.

► **Lemma 3.6.** *Any graph G with $|\mathcal{C}|$ colors has at most $2^{\ell_{\mathcal{C}}}$ fully-colorful subgraphs.*

Proof. Let n_c be the number of vertices of color $c \in \mathcal{C}$ and notice that $\prod_{c \in \mathcal{C}} n_c$ is the number of fully-colorful subgraphs of G . Then, observe that $n_c \leq 2^{n_c - 1}$ for all $n_c \in \mathbb{N}$, which implies $\prod_{c \in \mathcal{C}} n_c \leq 2^{\sum_{c \in \mathcal{C}} n_c - 1}$ and thus $\prod_{c \in \mathcal{C}} n_c \leq 2^{\ell_{\mathcal{C}}}$. ◀

► **Theorem 3.7.** *MCA can be solved in $\mathcal{O}^*(2^{\ell_{\mathcal{C}}} \cdot 4^{t_{\mathcal{H}}})$ time and $\mathcal{O}^*(3^{t_{\mathcal{H}}})$ space.*

Proof. In the following, let $(\{X_i : i \in I\}, \mathcal{T})$ be a nice tree decomposition of $U(\mathcal{H})$. In this proof, we provide a dynamic programming algorithm that makes use of $(\{X_i : i \in I\}, \mathcal{T})$ in order to compute a solution to MCA in any fully-colorful subgraph $G' \subseteq G$, to which we remove all vertices that are not accessible from r . First, observe that $(\{X_i : i \in I\}, \mathcal{T})$ is also a correct nice tree decomposition for the (undirected) color hierarchy graph of any subgraph of G . Second, as any colorful graph is equivalent to its color hierarchy graph, notice that $(\{X_i : i \in I\}, \mathcal{T})$ is also a correct nice tree decomposition of any fully-colorful subgraph $G' \subseteq G$. Therefore, we assume without loss of generality that any bag X_i contains vertices of such graph G' instead of colors, and that $X_0 = \{r\}$ is the root of $(\{X_i : i \in I\}, \mathcal{T})$.

Now, for any $i \in I$ and for any subsets L_1, L_2, L_3 that belong to X_i such that $L_1 \oplus L_2 \oplus L_3 = X_i$, let $T_i[L_1, L_2, L_3]$ store the weight of a *partial solution* of MCA in G' , which is a collection of $|L_1|$ disjoint arborescences such that:

- each $v \in L_1$ is the root of exactly one such arborescence,
- each $v \in L_2$ is contained in exactly one such arborescence,
- no vertex $v \in L_3$ belongs to any of these arborescences,
- any vertex $v \in V$ whose color is forgotten below X_i can belong to any such arborescence,

- there does not exist another collection of arborescences with a larger sum of weights under the same constraints.

Besides, let us define an entry of type $D_i[L_1, L_2, L_3]$ which stores the same partial solution as entry $T_i[L_1, L_2, L_3]$, except for the vertices $v \in V$ whose colors are forgotten below X_i which cannot belong to any arborescence of the partial solution. We now detail how to compute each entry of $T_i[L_1, L_2, L_3]$. We stress that each entry of $D_i[L_1, L_2, L_3]$ is filled exactly as an entry of type $T_i[L_1, L_2, L_3]$, apart from the case of forget nodes which we detail below.

- If X_i is a leaf node: $T_i[L_1, L_2, L_3] = 0$

Notice that leaf nodes are base cases of the dynamic programming algorithm as $\langle \{X_i : i \in I\}, \mathcal{T} \rangle$ is a nice tree decomposition. Moreover, recall that leaf nodes have size 1 and thus that the only partial solution for such nodes has a weight of zero.

- If X_i is an introduce node having a child X_j and if v^* is the introduced vertex:

$$T_i[L_1, L_2, L_3] = \begin{cases} A) \max_{\forall S \subseteq L_2} \left\{ \sum_{v \in S} w(v^*, v) + T_j[L_1 \cup S \setminus \{v^*\}, L_2 \setminus S, L_3] \right\} & \text{if } v^* \in L_1 \\ B) \max_{\forall u \in (L_1 \cup L_2)} \left\{ w(u, v^*) + \max_{\forall S \subseteq (L_2 \setminus \{u\})} \left\{ \sum_{v \in S} w(v^*, v) + T_j[L_1 \cup S \setminus \{v^*\}, L_2 \setminus S, L_3] \right\} \right\} & \text{if } v^* \in L_2 \\ C) T_j[L_1, L_2, L_3 \setminus \{v^*\}] & \text{if } v^* \in L_3 \end{cases}$$

where we set $w(u, v) := -\infty$ when there is no arc from u to v in G' . There are three cases: v^* is the root of an arborescence in a partial solution (case A)), an internal vertex of such a solution (case B)) or v^* does not belong to such a solution (case C)). In case A), S corresponds to the set of outneighbors of v^* in the partial solution, thus the vertices of S do not have any other inneighbor in the partial solution. Therefore, in the corresponding entry T_j , the vertices of S are roots. Now, notice that B) is very similar to A). In addition to a given set S of outneighbors, v^* being in L_2 implies that v^* has an inneighbor $u \in (L_1 \cup L_2)$ in the partial solution. Since the inneighbor u cannot be an outneighbor at the same time, u is not contained in S . Exhaustively trying all possibilities for both S and u ensures the correctness of the solution. Finally, by definition of L_3 , observe that v^* does not belong to the partial solution of $T_i[L_1, L_2, L_3]$ if $v^* \in L_3$.

- If X_i is a forget node having a child X_j and if v^* is the forgotten vertex:

$$T_i[L_1, L_2, L_3] = \max\{T_j[L_1, L_2 \cup \{v^*\}, L_3], T_j[L_1, L_2, L_3 \cup \{v^*\}]\}$$

Informally, the above formula determines whether the collection of arborescences that is stored in $T_i[L_1, L_2, L_3]$ had a higher weight with or without v^* as an internal vertex. Observe that we do not consider the case where v^* is the root of an arborescence as such an arborescence could not be connected to the rest of the partial solution via an introduced vertex afterwards. Besides, notice that $D_i[L_1, L_2, L_3] = D_j[L_1, L_2, L_3 \cup \{v^*\}]$ as the partial solution in $D_i[L_1, L_2, L_3]$ does not contain any forgotten vertex by definition.

- If X_i is a join node having two children X_j and X_k :

$$T_i[L_1, L_2, L_3] = T_j[L_1, L_2, L_3] + T_k[L_1, L_2, L_3] - D_i[L_1, L_2, L_3]$$

Informally, the partial solution in $T_i[L_1, L_2, L_3]$ can contain both the forgotten vertices of the partial solution in $T_j[L_1, L_2, L_3]$ and those of the partial solution in $T_k[L_1, L_2, L_3]$. Recall that the partial solution in $D_i[L_1, L_2, L_3]$ does not contain any forgotten vertices and therefore that any arc of the partial solution in $T_i[L_1, L_2, L_3]$ is only counted once.

We fill the tables from the leaves to the root for all $i \in I$ until T_0 and any entry of type $T_i[L_1, L_2, L_3]$ is directly computed after the entry of type $D_i[L_1, L_2, L_3]$. If $T' = (V_{T'}, A_{T'})$ is a solution of MCA in a fully-colorful subgraph $G' \subseteq G$, then $w(T') = T_0[\{r\}, \emptyset, \emptyset]$. Thus, for each fully-colorful subgraph we can compute the solution by filling the tables T and D . The table has $3^{t_{\mathcal{H}}}$ entries which implies the upper bound on the space consumption. The most expensive recurrences in terms of running time are the one of cases A) and B) for introduce nodes X_i where we consider altogether $\mathcal{O}(4^{t_{\mathcal{H}}})$ cases: each term corresponds to a partition of X_i into four sets L_1 , $L_2 \setminus S$, $L_2 \cap S$, and L_3 . Finally, the solution of MCA in G is also the solution of at least one fully-colorful subgraph $G' \subseteq G$. Therefore, computing the solution of MCA for any such subgraph G' ensures the correctness of the algorithm and hence, by Lemma 3.6, adding a factor $\mathcal{O}(2^{\ell_C})$ to the complexity of the above algorithm proves our theorem. \blacktriangleleft

We now use the proof of Theorem 2.7 to show that MCA parameterized by $t_{\mathcal{H}} + \ell_C$ is unlikely to admit a polynomial kernel. Recall that the proof shows a polynomial parameter transformation from SET COVER to MCA and notice that $(\text{col}(V_1) \cup \text{col}(V_3))$ is a vertex cover of $U(\mathcal{H})$ that is of size $x_{\mathcal{H}} + 1$. Moreover, recall that the size of a minimum vertex cover of a graph is lower-bounded by its treewidth. As a consequence, MCA does not admit any polynomial kernel for $t_{\mathcal{H}}$ unless $\text{NP} \subseteq \text{coNP}/\text{Poly}$ even if $\ell_C = 0$.

► **Corollary 3.8.** *MCA parameterized by $t_{\mathcal{H}}$ does not admit a polynomial kernel unless $\text{NP} \subseteq \text{coNP}/\text{Poly}$, even when restricted to the special case where $\ell_C = 0$.*

4 Conclusion

In this paper, we obtained an $\mathcal{O}^*(3^{x_{\mathcal{H}}})$ time algorithm for MCA, which improves upon the $\mathcal{O}^*(3^{|C|})$ of Böcker *et al.* [4]. We also showed that MCA parameterized by $x_{\mathcal{H}} + \ell_C$ is unlikely to admit a polynomial kernel and then that the problem admits such a kernel for the parameter $x_{\mathcal{H}} + \ell$. Furthermore, we proposed an FPT algorithm for MCA relatively to $t_{\mathcal{H}} + \ell_C$ and showed that MCA is W[2]-hard relatively to $t_{\mathcal{H}}$. Moreover, we showed that MCA parameterized by $\ell_C + t_{\mathcal{H}}$ does not admit a polynomial kernel. In light of these results, we ask the following question: does MCA parameterized by the larger parameter $\ell + t_{\mathcal{H}}$ admit a polynomial kernel?

A further issue that is not addressed by our algorithm and previous algorithms is that parameterization by ℓ or k essentially constrains the cardinality of the arborescences that are considered to be solutions. In other words, to make use of these parameters we need to know the number of vertices in an optimal solution in advance. Can we obtain fixed-parameter algorithms also when we do not know the number of vertices in the optimal solution?

References

- 1 Noga Alon, Raphael Yuster, and Uri Zwick. Color-coding. *J. ACM*, 42(4):844–856, 1995.
- 2 Nadja Betzler, René van Bevern, Michael R. Fellows, Christian Komusiewicz, and Rolf Niedermeier. Parameterized algorithmics for finding connected motifs in biological networks. *IEEE/ACM Trans. Comput. Biology Bioinform.*, 8(5):1296–1308, 2011.

- 3 Andreas Björklund, Petteri Kaski, and Lukasz Kowalik. Constrained multilinear detection and generalized graph motifs. *Algorithmica*, 74(2):947–967, 2016.
- 4 Sebastian Böcker and Florian Rasche. Towards *de novo* identification of metabolites by analyzing tandem mass spectra. In *Proceedings of the 7th European Conference on Computational Biology (ECCB '08)*, volume 24(16), pages i49–i55, 2008.
- 5 Hans L. Bodlaender, Rodney G. Downey, Michael R. Fellows, and Danny Hermelin. On problems without polynomial kernels. *J. Comput. Syst. Sci.*, 75(8):423–434, 2009.
- 6 Hans L. Bodlaender, Stéphan Thomassé, and Anders Yeo. Kernel bounds for disjoint cycles and disjoint paths. *Theor. Comput. Sci.*, 412(35):4570–4578, 2011.
- 7 Sharon Bruckner, Falk Hüffner, Richard M. Karp, Ron Shamir, and Roded Sharan. Topology-free querying of protein interaction networks. *J. Comput. Biol.*, 17(3):237–252, 2010.
- 8 Marek Cygan, Fedor V. Fomin, Lukasz Kowalik, Daniel Lokshtanov, Dániel Marx, Marcin Pilipczuk, Michal Pilipczuk, and Saket Saurabh. *Parameterized Algorithms*. Springer, 2015.
- 9 Marek Cygan, Marcin Pilipczuk, Michal Pilipczuk, and Jakub Onufry Wojtaszczyk. Kernelization hardness of connectivity problems in d -degenerate graphs. *Discr. Appl. Math.*, 160(15):2131–2141, 2012.
- 10 Michael Dom, Daniel Lokshtanov, and Saket Saurabh. Kernelization lower bounds through colors and IDs. *ACM Trans. Algorithms*, 11(2):13:1–13:20, 2014.
- 11 Rodney G. Downey and Michael R. Fellows. *Fundamentals of Parameterized Complexity*. Texts in Computer Science. Springer, 2013.
- 12 Kai Dührkop, Marie Anne Lataretu, W. Timothy J. White, and Sebastian Böcker. Heuristic algorithms for the maximum colorful subtree problem. *arXiv*, 2018.
- 13 Michael R. Fellows, Guillaume Fertin, Danny Hermelin, and Stéphane Vialette. Upper and lower bounds for finding connected motifs in vertex-colored graphs. *J. Comput. Syst. Sci.*, 77(4):799–811, 2011.
- 14 Guillaume Fertin, Julien Fradin, and Géraldine Jean. Algorithmic aspects of the maximum colorful arborescence problem. In *Proceedings of the 14th Annual Conference of Theory and Applications of Models of Computation (TAMC '17)*, volume 10185 of LNCS, pages 216–230, 2017.
- 15 Guillaume Fertin and Christian Komusiewicz. Graph motif problems parameterized by dual. In *Proceedings of the 27th Annual Symposium on Combinatorial Pattern Matching (CPM '16)*, volume 54 of LIPIcs, pages 7:1–7:12. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2016.
- 16 Jiong Guo, Falk Hüffner, and Rolf Niedermeier. A structural view on parameterizing problems: Distance from triviality. In *Proceedings of the First International Workshop on Parameterized and Exact Computation (IWPEC '04)*, volume 3162 of LNCS, pages 162–173. Springer, 2004.
- 17 Richard M. Karp. Reducibility among combinatorial problems. In *Proceedings of a Symposium on the Complexity of Computer Computations*, The IBM Research Symposia Series, pages 85–103. Plenum Press, New York, 1972.
- 18 Vincent Lacroix, Cristina G. Fernandes, and Marie-France Sagot. Motif search in graphs: Application to metabolic networks. *IEEE/ACM Trans. Comput. Biology Bioinform.*, 3(4):360–368, 2006.
- 19 Imran Rauf, Florian Rasche, Francois Nicolas, and Sebastian Böcker. Finding maximum colorful subtrees in practice. *J. Comput. Biol.*, 20(4):311–321, 2013.