Sorting With Forbidden Intermediates

Carlo Comin^{1,2}, Anthony Labarre¹, Romeo Rizzi³, and Stéphane Vialette¹

¹ Université Paris-Est, LIGM (UMR 8049), UPEM, CNRS, ESIEE, ENPC, F-77454,

Marne-la-Vallée, France {Anthony.Labarre,Stephane.Vialette}@u-pem.fr

² Department of Mathematics, University of Trento, Italy

³ Department of Knatheniatics, University of Trento, Italy

Carlo.CominGunitn.it, Romeo.RizziGunivr.it

Abstract. A wide range of applications, most notably in comparative genomics, involve the computation of a shortest sorting sequence of operations for a given permutation, where the set of allowed operations is fixed beforehand. Such sequences are useful for instance when reconstructing potential scenarios of evolution between species, or when trying to assess their similarity. We revisit those problems by adding a new constraint on the sequences to be computed: they must avoid a given set of *forbidden intermediates*, which correspond to species that cannot exist because the mutations that would be involved in their creation are lethal. We initiate this study by focusing on the case where the only mutations that can occur are exchanges of any two elements in the permutations, and give a polynomial time algorithm for solving that problem when the permutation to sort is an involution.

key Guided Sorting, Lethal Mutations, Forbidden Vertices, Permutation Sorting, Hypercube Graphs, st-Connectivity.

1 Introduction

Computing distances between permutations, or sequences of operations that transform them into one another, are two generic problems that arise in a wide range of applications, including comparative genomics [7], ranking [5], and interconnection network design [16]. Those problems are well-known to reduce to constrained sorting problems of the following form: given a permutation π and a set S of allowed operations, find a sequence of elements from S that sorts π and is as short as possible. In the context of comparative genomics, the sequence to be reconstructed yields a possible scenario of evolution between the genomes represented by π and the target identity permutation ι , where all permutations obtained inbetween are successive descendants of π (and ancestors of ι). The many possible choices that exist for S, as well as other constraints or cost functions with which they can be combined, have given rise to a tremendous number of variants whose algorithmic and mathematical aspects have now been studied for decades [7]. Specific issues that biologists feel need to be addressed to improve the applicability of these results in a biological context include:

- 1. the oversimplicity of the model (permutations do not take duplications into account),
- 2. the rigid definition of allowed operations, which fails to capture the complexity of evolution, and
- 3. the complexity of the resulting problems, where algorithmic hardness results abound even for deceivingly simple problems.

A large body of work has been devoted to addressing those issues, namely by proposing richer models for genomes, encompassing several operations with different weights [7]. Some approaches for increasing the reliability of rearrangement methods by adding additional biologically motivated constraints have been investigated (for instance, Bergeron et al. [2] consider conserved intervals, Figeac and Varré [8] restrict the set of allowed inversions and Bérard et al. [1] take into account the number of inversions in the wanted scenario which commute with all *common intervals*). However, another critical issue has apparently been overlooked: to the best of our knowledge, no model takes into account the fact that the solutions it produces may involve allele mutations that are lethal to the organism on which they act. Lethals are usually a result of mutations in genes that are essential to growth or development [10]; they have been known to occur for more than a century [4], dating back to the works of Cuénot in 1905 who was studying the inheritance of coat colour in mice. As a consequence, solutions that may be perfectly valid from a mathematical point of view should nonetheless be rejected on the grounds that some of the intermediate ancestors they produce are nonviable and can therefore not have had any descendants. We revisit the family of problems mentioned above by adding a natural constraint which, as far as we know, has not been previously considered in this form (see e.g. [2, 8, 1] for connected attempts): namely, the presence of a set of forbidden intermediate permutations, which the sorting sequence that we seek must avoid. We refer to this family of problems as GUIDED SORTING problems, since they take additional guidance into account. In this paper, we focus our study on the case where only exchanges (i.e., algebraic transpositions) are allowed; furthermore, we simplify the problem by demanding that the solutions we seek be *optimal* in the sense that no shorter sorting sequence of exchanges exists even when no intermediate permutation is forbidden. We choose to focus on exchanges because of their connection to the underlying *disjoint cycle structure* of permutations, which plays an important role in many related sorting problems where a similar cycle-based approach, using this time the ubiquitous breakpoint graph, has proved extremely fruitful [15]. Therefore, we believe that progress on this particular variant will be helpful when attempting to solve related variants based on more complex operations.

1.1 Contribution

Our main contribution in this work is a polynomial time algorithm for solving GUIDED SORTING by exchanges when the permutation to sort is an *involution*. We show that, in that specific case, the space of all feasible sorting sequences admits a suitable description in terms of directed (s, t)-paths in hypercube graphs.

We achieve this result by reducing GUIDED SORTING to the problem of finding directed (s, t)-paths that avoid a prescribed set $\mathcal{F} \subseteq V$ of *forbidden vertices*. Our main contribution, therefore, consists in solving this latter problem in time polynomial in just the encoding length of \mathcal{F} , if G is constrained to be a *hypercube* graph; which is a novel algorithmic result that may be of independent interest. Specific properties that will be described later on [11, 17] allow us to avoid the full construction of that graph, which would lead to an exponential time algorithm.

1.2 Related Works

We should mention that constrained variants of the (s, t)-connectivity problem have been studied already to some extent. For instance, in the early '70s, motivated by some problems in the field of automatic software testing and validation, Krause et al. [14] introduced the *path avoiding forbidden pairs* problem, namely, that of finding a directed (s, t)-path in a graph G = (V, E) that contains at most one vertex from each pair in a prescribed set $\mathcal{P} \subseteq V \times V$ of *forbidden pairs* of vertices. Gabow et al. [9] proved that the problem is NP-complete on DAGs. A number of special cases were shown to admit polynomial time algorithms, e.g. Yinnone [19] studied the problem in directed graphs under a *skew-symmetry* condition. However, the involved techniques and the related results do not extend to our problem, for which we are aware of no previously known algorithm that runs in time polynomial in just the encoding length of \mathcal{F} .

A preliminary version of this article appeared in the proceedings of the 3rd International Conference on Algorithms for Computational Biology (AlCoB 2016), see [3]. Here, the previous results are improved and the presentation is extended:

(1) The time complexity of Algorithm 2 is improved by a factor of $d_{S,T} \cdot n$ (see Theorem 1 for the actual time bound).

(2) Subsection 3.5 is extended by presenting Algorithm 1 plus all the details of its correctness and running time analysis.

(3) Subsection 3.6 is extended by including a detailed correctness and complexity analysis of Algorithm 2.

(4) Fig. 1, Fig. 2 and Fig. 5 have been added to support some of the more technical constructions with an illustration.

1.3 Organization

The remainder of the article is organized as follows. Section 2 provides some background notions and notation on which the rest of this work relies. The main contribution is offered in Section 3. In Subsection 3.1, the problem HY-STCON is formulated. The reduction from GUIDED SORTING for exchanges (and adjacent exchanges) to HY-STCON is offered in Subsection 3.2 (and 3.3, respectively). The formal statement of our main algorithmic contribution is detailed in Subsection 3.4. Next, Subsection 3.5 concerns the specific properties [11, 17] that allow us to avoid the full construction of the hypercube search space. Subsection 3.6 presents the polynomial-time algorithm for solving HY-STCON. In

Subsection 3.7, it is shown how to speed up the algorithm in the case in which one is interested just in the decision task of HY-STCON. The correctness analysis of the main algorithm is carried on in Subsection 3.8, while the complexity is analyzed in Subsection 3.9. We conclude in Section 4 with a discussion of several open problems.

2 Background and Notation

We use the notation $\pi = \langle \pi_1 \ \pi_2 \ \cdots \ \pi_k \rangle$ when viewing permutations as sequences, i.e. $\pi_i = \pi(i)$ for $i \in [k] = \{1, 2, \dots, k\}$. Our aim is to sort a given permutation π , i.e. to transform it into the identity permutation $\iota = \langle 1 \ 2 \ \cdots \ k \rangle$, using a predefined set of allowed operations specified as a generating set S of the symmetric group \mathfrak{S}_k . We seek a sorting sequence that uses only elements from S and:

- 1. avoids a given set \mathcal{F} of forbidden permutations, i.e. no intermediary permutation produced by applying the operations specified by the sorting sequence belongs to \mathcal{F} , and
- 2. is *optimal*, i.e. no shorter sorting sequence exists for π even if $\mathcal{F} = \emptyset$.

We use standard notions and notation from graph theory (see e.g. Diestel [6] for undefined concepts), using $\{u, v\}$ (resp. (u, v)) to denote the edge (resp. arc) between vertices u and v of an undirected (resp. directed) graph G = (V, E). All graphs we consider are *simple*: they contain neither loops nor parallel edges / arcs. If $\mathcal{F} \subseteq V$, a directed path $\mathbf{p} = v_0 v_1 \cdots v_n$ avoids \mathcal{F} when $v_i \notin \mathcal{F}$ for every i. If $S \subseteq V$ and $\mathcal{T} \subseteq V$, we say that \mathbf{p} goes from S to \mathcal{T} in G if $v_0 \in S$ and $v_n \in \mathcal{T}$. When G is directed, we partition the neighbourhood N(u) of a vertex u into the sets $N^{\text{out}}(u) = \{v \in V \mid (u, v) \in E\}$ and $N^{\text{in}}(u) = \{v \in V \mid (v, u) \in E\}$. Some of our graphs may be vertex-labelled, using any injective mapping $\ell : V \to \mathbb{N}$. For any $n \in \mathbb{N}$, $\wp_n = \wp([n])$ denotes the power set of [n]. The hypercube graph on ground set [n], denoted by \mathcal{H}_n , is the graph with vertex set \wp_n and in which the arc (U, V) connects vertices $U, V \subseteq [n]$ if there exists some $q \in [n]$ such that $U = V \setminus \{q\}$. If $S, T \in \wp_n$ and $|S| \leq |T|$, then $d_{S,T} = |T| - |S|$ is the distance between S and T. Finally, $\mathcal{H}_n^{(i)}$ denotes the family of all subsets of \wp_n of size i.

3 Solving GUIDED SORTING For Involutions

The Cayley graph $\Gamma(\mathfrak{S}_n, S)$ of \mathfrak{S}_n for a given generating set S of \mathfrak{S}_k contains a vertex for each permutation in \mathfrak{S}_k and an edge between any two permutations that can be obtained from one another using one element from S. A naïve approach for solving any variant of the GUIDED SORTING problem would build the part of $\Gamma(\mathfrak{S}_k, S)$ that is needed (i.e. without the elements of \mathcal{F}), then run a shortest path algorithm to compute an optimal sequence that avoids all elements of \mathcal{F} . This is highly impractical, since the size of Γ is exponential in k.

We describe in this section a polynomial time algorithm for the case where S is the set of all exchanges and π is an *involution*, i.e. a permutation such that for each $1 \leq i \leq |\pi|$, either $\pi_i = i$ or there exists an index j such that $\pi_i = j$ and $\pi_j = i$. From our point of view, involutions reduce to collections of disjoint pairs of elements that each need to be swapped by an exchange until we obtain the identity permutation, and the only forbidden permutations that could be produced by an optimal sorting sequence are involutions whose pairs of unsorted elements all appear in π . Therefore, we can reformulate our GUIDED SORTING problem in that setting as that of finding a directed (π, ι) -path in \mathcal{H}_n that avoids all vertices in \mathcal{F} , where the permutation to sort π corresponds to the bottom vertex \emptyset of \mathcal{H}_n . We give more details on the reduction in Section 3.2.

3.1 Problem Formulation

We shall focus on the following problem from here on.

INPUT: the size $n \in \mathbb{N}$ of the underlying ground set [n], a family of forbidden vertices $\mathcal{F} \subseteq \wp_n$, a source set $S \in \wp_n$ and a target set $T \in \wp_n$.

DECISION-TASK: Decide whether there exists a directed path \mathbf{p} in \mathcal{H}_n that goes from source S to target T avoiding \mathcal{F} ; SEARCH-TASK: Compute a directed path \mathbf{p} in \mathcal{H}_n that goes from source

S to target T avoiding \mathcal{F} , provided that at least one such path exists.

We examine in this section specific instances of GUIDED SORTING which can be solved through a reduction to HY-STCON. We say that permutations that may occur in an optimal sorting sequence for a given permutation π are *relevant*, and all others are *irrelevant*. The distinction will matter when sorting a particular permutation since, as we shall see, the structure of π (however it is measured) will have implications on that of relevant permutations and will allow us to simplify the set of forbidden permutations by discarding irrelevant ones. For a fixed set S of operations, we let $R_S(\pi)$ denote the set of permutations that are relevant to π . Undefined terms and unproven properties of permutations below are well-known, and details are in standard references, such as Knuth [13].

3.2 GUIDED SORTING For Exchanges

Recall that every permutation π in \mathfrak{S}_k decomposes in a single way into *disjoint* cycles (up to the ordering of cycles and of elements within each cycle). This

decomposition corresponds to the cycle decomposition of the directed graph $G(\pi) = (V, A)$, where V = [k] and $A = \{(i, \pi_i) \mid 1 \leq i \leq k\}$. The *length* of a cycle of π is then simply the number of elements it contains, and the number of cycles of π is denoted by $c(\pi)$.

The Cayley distance of a permutation π is the length of an optimal sorting sequence of exchanges for π , and its value is $|\pi| - c(\pi)$. Therefore, when searching for an optimal sorting sequence, we may restrict our attention to exchanges that split a particular cycle into two smaller ones.

Let (π, \mathcal{F}, S, K) be an instance of GUIDED SORTING such that S is the set of all exchanges and where the permutation π to sort is an *involution*, i.e. a permutation whose cycles have length at most two. It is customary to omit cycles of length 1, and to write a permutation $\pi = \langle \pi_1 \ \pi_2 \ \cdots \ \pi_k \rangle$ with n cycles of length 2 as $c_1c_2 \cdots c_n$. Since we are looking for an optimal sorting sequence, we may assume that all permutations in \mathcal{F} are relevant, which in this case means that every permutation ϕ in \mathcal{F} is an involution and its 2-cycles form a proper subset of those of π . Our instance of GUIDED SORTING then translates to the following instance of HY-STCON:

- $-\pi \mapsto [n]$ in the following way: $c_i \mapsto i$ for $1 \leq i \leq n$;
- each permutation ϕ in \mathcal{F} is mapped onto a subset of [n] by replacing its cycles with the indices obtained in the first step; let \mathcal{F}' denote the collection of subsets of [n] obtained by applying that mapping to each ϕ in \mathcal{F} .

The resulting HY-STCON instance is then $\langle [n], \emptyset, \mathcal{F}', n \rangle$, and a solution to instance (π, \mathcal{F}, S, K) of GUIDED SORTING exists if and only if a solution to instance $\langle [n], \emptyset, \mathcal{F}', n \rangle$ of HY-STCON exists; the translation of the solution from the latter formulation to the former is straightforward.

3.3 GUIDED SORTING For Adjacent Exchanges

Recall that an *inversion* in a permutation π in \mathfrak{S}_k is a pair (π_i, π_j) with $1 \leq i < j \leq k$ and $\pi_i > \pi_j$. Let (π, \mathcal{F}, S, K) be an instance of GUIDED SORTING where S is the set of all *adjacent* exchanges, i.e. exchanges that act on consecutive positions. It is well-known that in this case, any optimal sorting sequence for π has length equal to the number of inversions of π , which means that in the search for an optimal sorting sequence, we may restrict our attention to adjacent exchanges that act on inversions that consist of adjacent elements.

Let us now assume that all n inversions of π are made of adjacent elements, and denote $\pi = i_1 i_2 \cdots i_n$, where each i_j is an inversion. Since we are looking for an optimal sorting sequence, we may assume that all permutations in \mathcal{F} are relevant, which in this case means that all inversions of any permutation ϕ in \mathcal{F} form a proper subset of those of π . The reduction to HY-STCON in that setting is very similar to that given in the case of exchanges:

- $-\pi \mapsto [n]$ in the following way: $i_j \mapsto j$ for $1 \le j \le n$;
- each permutation ϕ in \mathcal{F} is mapped onto a subset of [n] by replacing its inversions with the indices obtained in the first step; let \mathcal{F}' be the collection of subsets of [n] obtained by applying that mapping to each ϕ in \mathcal{F} .

The resulting HY-STCON instance is then $\langle [n], \emptyset, \mathcal{F}', n \rangle$, and a solution to instance (π, \mathcal{F}, S, K) of GUIDED SORTING exists if and only if a solution to instance $\langle [n], \emptyset, \mathcal{F}', n \rangle$ of HY-STCON exists; the translation of the solution from the latter formulation to the former is straightforward.

3.4 Main Result

In the rest of this section, we will show how to solve HY-STCON in time polynomial in $|\mathcal{F}|$ and n. The algorithm mainly consists in the continuous iteration of two phases:

- 1. Double-BFS. This phase explores the outgoing neighbourhood of the source S by a breadth-first search denoted by BFS_{\uparrow} going from lower to higher levels of \mathcal{H}_n while avoiding the vertices in \mathcal{F} . BFS_{\uparrow} collects a certain (polynomially bounded) amount of visited vertices. Symmetrically, the incoming neighbourhood of the target vertex T is also explored by another breadth-first search BFS_{\downarrow} going from higher to lower levels of \mathcal{H}_n while avoiding the vertices in \mathcal{F} , also collecting a certain (polynomially bounded) amount of visited vertices.
- 2. Compression. If a valid solution has not yet been determined, then a compression technique is devised in order to shrink the size of the remaining search space. This is possible thanks to some nice regularities of the search space and to certain connectivity properties of hypercube graphs [11, 17]. This allows us to reduce the search space in a suitable way and, therefore, to continue with the Double-BFS phase in order to keep the search towards valid solutions going.

Our main contribution is summarized in the following theorem. We devote the rest of this section to an in-depth description of the algorithms it mentions.

Theorem 1. Concerning the HY-STCON problem, the following propositions hold on any input (S, T, \mathcal{F}, n) , where $d_{S,T}$ is the distance between S and T.

1. There exists an algorithm for solving the DECISION-TASK of HY-STCON whose time complexity is:

$$O(\min(\sqrt{|\mathcal{F}| \, d_{S,T} \, n}, |\mathcal{F}|) \, |\mathcal{F}|^2 \, d_{S,T}^3 \, n).$$

2. There exists an algorithm for solving the SEARCH-TASK of HY-STCON whose time complexity is:

$$O(\min(\sqrt{|\mathcal{F}| d_{S,T} n}, |\mathcal{F}|) |\mathcal{F}|^2 d_{S,T}^3 n + |\mathcal{F}|^{5/2} n^{3/2} d_{S,T}).$$

3.5 On Vertex-Disjoint Paths in Hypercube Graphs

The proof of Theorem 1 relies on connectivity properties of hypercube graphs [11]. The next result, which proves the existence of a family of certain vertex-disjoint paths in \mathcal{H}_n that are called *Lehman-Ron paths*, will be particularly useful.

Theorem 2 (Lehman, Ron [17]). Given $n, m \in \mathbb{N}$, let $\mathcal{R} \subseteq \mathcal{H}_n^{(r)}$ and $\mathcal{S} \subseteq \mathcal{H}_n^{(s)}$ with $|\mathcal{R}| = |\mathcal{S}| = m$ and $0 \leq r < s \leq n$. Assume there exists a bijection $\varphi: \mathcal{S} \to \mathcal{R}$ such that $\varphi(S) \subset S$ for every $S \in \mathcal{S}$. Then, there exist m vertexdisjoint directed paths in \mathcal{H}_n whose union contains all subsets in \mathcal{S} and \mathcal{R} .

We call tuples $\langle \mathcal{R}, \mathcal{S}, \varphi, n \rangle$ that satisfy the hypotheses of Theorem 2 Lehman-Ron tuples, and we refer to the quantity d = s - r as the distance between $\mathcal{R} \subseteq \mathcal{H}_n^{(r)}$ and $\mathcal{S} \subseteq \mathcal{H}_n^{(s)}$. Lehman and Ron [17] give an elementary inductive proof of Theorem 2; also, they showed that Theorem 2 does not hold if one requires that the disjoint chains exactly correspond to the given bijection φ . Anyway, a careful and in-depth analysis of their proof, from the algorithmic perspective, yields a polynomial time algorithm for computing all the Lehman-Ron paths.

Theorem 3. There exists an algorithm for computing all the Lehman-Ron paths within time $O(m^{5/2}n^{3/2}d)$ on any Lehman-Ron input $\langle \mathcal{R}, \mathcal{S}, \varphi, n \rangle$ with $|\mathcal{R}| = |\mathcal{S}| = m$, where d is the distance between \mathcal{R} and \mathcal{S} and n is the size of the underlying ground set.

Now we provide all the details of the algorithm sketched above as well as a proof of the time complexity stated in Theorem 3, in which Menger's vertexconnectivity theorem [6] and Hopcroft-Karp's algorithm [12] for maximum cardinality matching in *undirected* bipartite graphs play a major role.

As mentioned, our proof of Theorem 1 relies on certain connectivity properties of hypercube graphs, and in particular the existence of a family of certain vertex-disjoint paths in \mathcal{H}_n that we call "Lehman-Ron paths", which is guaranteed by Theorem 2.

Although Theorem 2 was initially proved and applied in the specific area of testing monotonicity [11], it is of independent interest and related results could be useful in the context of packet routing on the hypercube network. Lehman and Ron provided an elegant inductive proof of that result [17]. In the present work, we point out that a careful analysis of their proof allows us to "extract" a simple recursive algorithm for computing all Lehman-Ron paths in polynomial time. We now describe that algorithm, whose correctness follows from the arguments used by Lehman and Ron in their original proof of Theorem 2 (see [17] for more details). Its time complexity can be derived by taking into account Hopcroft-Karp's algorithm for computing maximum cardinality matchings in bipartite graphs [12], and is analyzed in detail at the end of this section.

The algorithm that we are going to describe is named compute_Lehman-Ron_paths(). The intuition underlying it is simply to follow the structure of Lehman and Ron's proof of Theorem 2 and to analyze it from the algorithmic standpoint. Its pseudocode is given in Algorithm 1.

The algorithm takes as input a Lehman-Ron tuple $\langle \mathcal{R}, \mathcal{S}, \varphi, n \rangle$, and outputs a family $\mathbf{p}_1, \mathbf{p}_2, \ldots, \mathbf{p}_m$ of Lehman-Ron paths joining \mathcal{R} to \mathcal{S} . Recall that Lehman-Ron tuples satisfy the following properties:

1. the families of sets $\mathcal{R} \subseteq \mathcal{H}_n^{(r)}$ and $\mathcal{S} \subseteq \mathcal{H}_n^{(s)}$ are such that $|\mathcal{S}| = |\mathcal{R}| = m$,

Algorithm 1: computing Lehman-Ron's paths.

Procedure compute_Lehman-Ron_paths($\mathcal{R}, \mathcal{S}, \varphi, n$) **Input**: a Lehman-Ron tuple $\langle \mathcal{R}, \mathcal{S}, \varphi, n \rangle$. **Output**: a family of *m* vertex-disjoint directed paths p_1, \ldots, p_m in \mathcal{H}_n such that $\mathcal{R} \cup \mathcal{S} \subseteq \bigcup_{i=1}^{m} \mathsf{p}_i$. if s = r + 1 then 1 return compute_paths_from_bijection(S, φ, n); $\mathbf{2}$ $m \leftarrow |\mathcal{S}|;$ // assume $|\mathcal{S}| = |\mathcal{R}|$ 3 $\mathcal{Q} \leftarrow \texttt{compute}_Q(\mathcal{S});$ 4 $\mathcal{K} \leftarrow \texttt{compute_auxiliary_network}(\mathcal{R}, \mathcal{Q}, \mathcal{S});$ $\mathbf{5}$ $\langle \mathsf{p}'_1, \mathsf{p}'_2, \dots, \mathsf{p}'_m \rangle \leftarrow \texttt{compute_vertex_disjoint_paths}(\mathcal{K});$ 6 $\langle \mathcal{Q}', \varphi', \varphi'' \rangle \leftarrow \texttt{compute_auxiliary_bjcts}(\langle \mathsf{p}_1', \mathsf{p}_2', \dots, \mathsf{p}_m' \rangle, m)$ 7 $\langle \mathsf{p}_1'',\mathsf{p}_2'',\ldots,\mathsf{p}_m''\rangle \gets \texttt{compute_Lehman-Ron_paths}(\mathcal{R},\mathcal{Q}',(\varphi')^{-1},n);$ 8 $\langle \mathsf{p}_1, \mathsf{p}_2, \dots, \mathsf{p}_m \rangle \leftarrow \texttt{extend_paths}(\langle p_1'', \mathsf{p}_2'', \dots, p_m'' \rangle, \mathcal{Q}', \varphi'', m);$ 9 return $\langle \mathsf{p}_1, \mathsf{p}_2, \ldots, \mathsf{p}_m \rangle$; 10

2. r, s and $n \in \mathbb{N}$ are such that $0 \leq r < s \leq n$, and 3. $\varphi : S \to \mathcal{R}$ is a bijection such that $\forall S \in S : \varphi(S) \subset S$.

As a base case of the algorithm, if s = r + 1 (line 1), then the sought family of directed paths $\mathbf{p}_1, \mathbf{p}_2, \ldots, \mathbf{p}_m$ is simply a set of *m* pairwise vertex-disjoint arcs oriented from S to \mathcal{R} , which are already given by the input bijection φ (line 2).

We now focus on the general case s > r + 1. To begin with, we introduce the following proposition, which was already implicit in [17] and which is actually a straightforward consequence of Theorem 2.

Proposition 1. [17] Given $n, m \in \mathbb{N}$, consider two families of sets $\mathcal{R} \subseteq \mathcal{H}_n^{(r)}$ and $\mathcal{S} \subseteq \mathcal{H}_n^{(s)}$ where $|\mathcal{R}| = |\mathcal{S}| = m$ and $0 \leq r < s \leq n$. Let \mathcal{Q} (resp. \mathcal{P}) be the set of vertices in $\mathcal{H}_n^{(s-1)}$ (resp. $\mathcal{H}_n^{(r+1)}$) that lie on any directed path from some vertex in \mathcal{R} to some vertex in \mathcal{S} .

Then, $|\mathcal{Q}| \geq m$ and $|\mathcal{P}| \geq m$.

The algorithm first computes the set Q of all vertices in $\mathcal{H}_n^{(s-1)}$ that lie on any directed path from some vertex in \mathcal{R} to some vertex in \mathcal{S} . This step is encoded by compute_Q() (line 4). The algorithm then invokes (at line 5) a procedure called compute_auxiliary_network(), which constructs a directed auxiliary network $\mathcal{K} = (V_{\mathcal{K}}, A_{\mathcal{K}})$ which will be useful in the following steps and is defined by:

 $- V_{\mathcal{K}} = \{\mathbf{s}, \mathbf{t}\} \cup \mathcal{R} \cup \mathcal{Q} \cup \mathcal{S}, \text{ where } \mathbf{s} \text{ (resp. } \mathbf{t}) \text{ is an auxiliary source (resp. target)} \\ \text{vertex, i.e. } \{\mathbf{s}, \mathbf{t}\} \cap (\mathcal{R} \cup \mathcal{Q} \cup \mathcal{S}) = \emptyset;$

 $-A_{\mathcal{K}}$ is defined as follows:

- the source vertex \mathbf{s} is joined to every vertex in \mathcal{R} ;
- for each $R \in \mathcal{R}$ and $Q \in \mathcal{Q}$, R is joined to Q if and only if $R \subset Q$;
- similarly, for each $Q \in \mathcal{Q}$ and $S \in \mathcal{S}$, Q is joined to S if and only if $Q \subset S$;



Fig. 1: The auxiliary network $\mathcal{K} = (V_{\mathcal{K}}, A_{\mathcal{K}})$ and bijection φ .

• finally, every vertex in S is joined to \mathbf{t} .

Fig. 1 shows an example of an auxiliary network. We remark that, as shown in [17], the following proposition holds on \mathcal{K} .

Proposition 2. [17] The minimum (s, t)-vertex-separator of \mathcal{K} has size m.

As a corollary, and by applying Menger's vertex-connectivity theorem (which is recalled below), the existence of m internally-vertex-disjoint directed (\mathbf{s}, \mathbf{t}) -paths, denoted $\mathbf{p}'_1, \mathbf{p}'_2, \ldots, \mathbf{p}'_m$, is thus guaranteed.

Theorem 4 (Menger [6]). Let G = (V, A) be a directed graph, and let u and v be nonadjacent vertices in V. Then the maximum number of internally-vertexdisjoint directed (u, v)-paths in G equals the minimum number of vertices from $V \setminus \{u, v\}$ whose deletion destroys all directed (u, v)-paths in G.

How to compute $\mathbf{p}'_1, \mathbf{p}'_2, \ldots, \mathbf{p}'_m$ We argue that it is possible to compute efficiently the family of directed paths $\mathbf{p}'_1, \mathbf{p}'_2, \ldots, \mathbf{p}'_m$ in \mathcal{K} by finding a maximum cardinality matching in an auxiliary, *undirected* bipartite graph \mathcal{K}' . This reduction is performed by compute_vertex_disjoint_paths() at line 6. The undirected graph $\mathcal{K}' = (V_{\mathcal{K}'}, E_{\mathcal{K}'})$ is obtained from the directed graph \mathcal{K} as follows: first, the set family \mathcal{Q} gets split into two (disjoint) twin set families $\mathcal{Q}^{(in)}$ and $\mathcal{Q}^{(out)}$, i.e. $\mathcal{Q}^{(in)} = \{Q^{(in)} \mid Q \in \mathcal{Q}\}$ and $\mathcal{Q}^{(out)} = \{Q^{(out)} \mid Q \in \mathcal{Q}\}$ where $\mathcal{Q}^{(in)} \cap \mathcal{Q}^{(out)} = \emptyset$ and $|\mathcal{Q}^{(in)}| = |\mathcal{Q}^{(out)}| = |\mathcal{Q}|$. Thus, the vertex set of \mathcal{K}' is:

$$V_{\mathcal{K}'} = \mathcal{R} \cup \mathcal{Q}^{(\mathbf{in})} \cup \mathcal{Q}^{(\mathbf{out})} \cup \mathcal{S}.$$

The edge set $E_{\mathcal{K}'}$ is obtained as follows:



Fig. 2: The undirected bipartite graph $\mathcal{K}' = (V_{\mathcal{K}'}, E_{\mathcal{K}'})$ and a perfect matching \mathcal{M} (thick edges).

- for each $R \in \mathcal{R}$ and $Q \in \mathcal{Q}$, R is joined to $Q^{(in)}$ if and only if $R \subset Q$; similarly, for each $Q \in \mathcal{Q}$ and $S \in \mathcal{S}$, $Q^{(out)}$ is joined to S if and only if $Q \subset S;$
- finally, $Q^{(in)}$ is joined to $Q^{(out)}$ for every $Q \in Q$.

Fig. 1 shows an example of \mathcal{K}' . The next proposition derives some useful properties of \mathcal{K}' .

Proposition 3. The graph $\mathcal{K}' = (V_{\mathcal{K}'}, E_{\mathcal{K}'})$, as defined above, is bipartite and it admits a perfect matching.

Proof. The bipartiteness of \mathcal{K}' follows from the bipartition $(\mathcal{R} \cup \mathcal{Q}^{(out)}, \mathcal{Q}^{(in)} \cup$ S). To see that \mathcal{K}' admits a perfect matching, recall that, by Proposition 2 and by Theorem 4, there exist m internally-vertex-disjoint directed (\mathbf{s}, \mathbf{t}) -paths p'_1, p'_2, \ldots, p'_m in \mathcal{K} . Then, for every $i \in [m]$, let $p'_i = \mathbf{s} R_i Q_i S_i \mathbf{t}$ for some $R_i \in \mathcal{R}_i$ $\mathcal{R}, Q_i \in \mathcal{Q}, S_i \in \mathcal{S}$. Finally, let us define $\hat{\mathcal{Q}} = \mathcal{Q} \setminus \{Q \mid \exists i \in [m] \text{ s.t. } \mathbf{p}'_i = \mathcal{Q} \}$ $\mathbf{s}R_iQS_i\mathbf{t}$. At this point, let us consider the following matching \mathcal{M} of \mathcal{K}' :

$$\mathcal{M} = \left\{ \{R_i, Q_i^{(\mathbf{in})}\}, \{Q_i^{(\mathbf{out})}, S_i\} \mid \exists i \in [m] \text{ s.t. } \mathsf{p}'_i = \mathsf{s}R_iQ_iS_i\mathsf{t} \right\} \\ \cup \left\{ \{Q^{(\mathbf{in})}, Q^{(\mathbf{out})}\} \mid Q \in \hat{\mathcal{Q}} \} \right\}.$$

Since $m = |\mathcal{R}| = |\mathcal{S}|$ and p'_1, p'_2, \dots, p'_m are internally-vertex-disjoint, it follows that \mathcal{M} is a perfect matching of \mathcal{K}' .

We are in position to show how to compute p'_1, p'_2, \ldots, p'_m based on \mathcal{K} . Firstly, the procedure compute_vertex_disjoint_paths() constructs \mathcal{K}' as explained above and computes a maximum cardinality matching \mathcal{M} of \mathcal{K}' (e.g. with Hopcroft-Karp's algorithm [12]), which is perfect by Proposition 3. Therefore, the following property holds: for every $Q \in \mathcal{Q}$, there exists $R \in \mathcal{R}$ such that $\{R, Q^{(in)}\} \in \mathcal{M}$

if and only if there exists $S \in S$ such that $\{Q^{(\text{out})}, S\} \in \mathcal{M}$. We can then proceed as follows: for each $R_i \in \mathcal{R}$, the algorithm finds $Q_i \in \mathcal{Q}$ such that $\{R_i, Q_i^{(\text{in})}\} \in \mathcal{M}$ and then it finds $S_i \in S$ such that $\{Q_i^{(\text{out})}, S_i\} \in \mathcal{M}$. Then, compute_vertex_disjoint_paths() returns the family of paths p'_1, p'_2, \ldots, p'_m defined as: $p'_i = \mathbf{s}R_iQ_iS_i\mathbf{t}$ for every $i \in [m]$. Since \mathcal{M} is a perfect matching of \mathcal{K}' , the paths p'_1, p'_2, \ldots, p'_m are internally-vertex-disjoint.

Let $Q' = \{Q \mid \exists i \in [m] \text{ s.t. } \mathbf{p}'_i = \mathbf{s}R_iQS_i\mathbf{t}\}$. Once we have computed $\mathbf{p}'_1, \mathbf{p}'_2, \ldots, \mathbf{p}'_m$, we can deduce two bijections that will be helpful in obtaining the wanted paths:

$$\varphi' : \mathcal{R} \to \mathcal{Q}' \text{ and } \varphi'' : \mathcal{Q}' \to \mathcal{S}.$$

The first bijection is defined for any $R \in \mathcal{R}$ as $\varphi'(R) = Q$ (where $Q \in Q'$) provided there exists some p'_i joining R to Q; similarly, the second bijection is defined for any $Q \in Q'$ as $\varphi''(Q) = S$ (where $S \in S$) provided there exists some p'_i joining Q to S. These bijections are computed by compute_auxiliary_bjcts() at line 7.

At this point, since the distance between \mathcal{R} and \mathcal{Q}' equals s - 1, a recursive call to compute_Lehman-Ron_paths () on input $\langle \mathcal{R}, \mathcal{Q}', (\varphi')^{-1}, n \rangle$ yields, at line 8, a family of Lehman-Ron paths $p''_1, p''_2, \ldots, p''_m$ joining \mathcal{R} to \mathcal{Q}' .

Indeed, we argue that it is possible to construct, starting from $p''_1, p''_2, \ldots, p''_m$, the sought family of Lehman-Ron paths p_1, p_2, \ldots, p_m that join \mathcal{R} to \mathcal{S} . Actually, this can be done just by taking into account the bijection φ'' : since φ'' joins \mathcal{Q}' to \mathcal{S} , it suffices to perform the following steps in practice:

- 1. consider the last vertex Q_i of p''_i (i.e. the unique vertex $Q_i \in \mathcal{Q}'$ such that $Q_i \in \mathsf{p}''_i \cap \mathcal{Q}'$);
- 2. let $S_i = \varphi''(Q_i);$
- 3. concatenate S_i at the end of \mathbf{p}''_i (i.e. $\mathbf{p}_i = \mathbf{p}''_i S_i$).

This construction is performed by the extend_paths() procedure at line 9. Since $p_1'', p_2'', \ldots, p_m''$ are vertex-disjoint and $\varphi'' : Q' \to S$ is a bijection, p_1, p_2, \ldots, p_m is the sought family of Lehman-Ron paths joining \mathcal{R} to \mathcal{S} .

Complexity Analysis (Proof of Theorem 3) We now turn to the time complexity analysis of Algorithm 1, going through each line in detail.

- line 2: compute_paths_from_bijection() (line 2) takes time at most O(m), which corresponds to the time needed to inspect the input bijection φ .
- line 4: compute_Q() takes time at most O(mn): for each $S \in S$, the procedure inspects the predecessors $N^{\text{in}}(S)$, and the time bound follows from the fact that |S| = m and $|N^{\text{in}}(S)| \leq n$.
- line 5: we argue that $|V_{\mathcal{K}}| = O(mn)$ and $|A_{\mathcal{K}}| = O(m^2n)$. Indeed, recall that $|\mathcal{R}| = |\mathcal{S}| = m$ by hypothesis; and since every vertex of \mathcal{S} has at most n neighbours in \mathcal{Q} , we have $|\mathcal{Q}| \leq mn$. This in turn implies that $|V_{\mathcal{K}}| \leq 2 + 2m + mn$; moreover, each of the m vertices in \mathcal{R} has at most mn neighbours, which all lie in \mathcal{Q} . Therefore, $|A_{\mathcal{K}}| \leq 2m + m^2n + mn$, and the procedure compute_auxiliary_network() takes time at most $O(|V_{\mathcal{K}}| + |A_{\mathcal{K}}|) = O(m^2n)$.

Algorithm 2: Solving the HY-STCON problem.

Procedure *solve*_HY-STCON (S, T, \mathcal{F}, n) **Input**: an instance $\langle S, T, \mathcal{F}, n \rangle$ of HY-STCON. **Output**: a pair $\langle YES, p \rangle$ where the path p is a solution to HY-STCON if such a path exists, NO otherwise. // let $d_{\boldsymbol{S},T}$ be the distance between \boldsymbol{S} and T $d_{S,T} \leftarrow |T| - |S|;$ 1 $\mathcal{S} \leftarrow \{S\}; \ell_{\uparrow} \leftarrow 0; // \text{ init the frontier } \mathcal{S} \text{ and its level counter } \ell_{\uparrow}$ $\mathbf{2}$ $\mathcal{T} \leftarrow \{T\}; \ell_{\perp} \leftarrow 0; // \text{ init the frontier } \mathcal{T} \text{ and its level counter } \ell_{\perp}$ 3 $\mathbf{4}$ while TRUE do $\langle \mathcal{S}, \mathcal{T}, \ell_{\uparrow}, \ell_{\downarrow} \rangle \leftarrow \texttt{double-bfs_phase}(\mathcal{S}, \mathcal{T}, \mathcal{F}, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n);$ 5 if $S = \emptyset$ OR $T = \emptyset$ OR $(\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$ AND $S \cap T = \emptyset$) then 6 return NO; 7 if $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$ AND $S \cap \mathcal{T} \neq \emptyset$ then 8 $p \leftarrow \text{reconstruct_path}(\mathcal{S}, \mathcal{T}, n);$ 9 return $\langle YES, p \rangle;$ 10 returned_val \leftarrow compression_phase($\mathcal{S}, \mathcal{T}, \mathcal{F}, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n$); 11 12 if $returned_val = \langle YES, p \rangle$ then return p; else $\mathcal{T} \leftarrow \texttt{returned_val};$ 13

- line 6: compute_vertex_disjoint_paths() takes time at most $O(m^{5/2}n^{3/2})$. Indeed, let us consider the auxiliary (undirected) bipartite graph $\mathcal{K}' = (V_{\mathcal{K}'}, E_{\mathcal{K}'})$ defined above. Since $|V_{\mathcal{K}}| = O(mn)$ and $|A_{\mathcal{K}}| = O(m^2n)$, we have $|V_{\mathcal{K}'}| = O(mn)$ and $|E_{\mathcal{K}'}| = O(m^2n)$ by construction. A maximum cardinality matching \mathcal{M} of \mathcal{K}' can be computed with Hopcroft-Karp's algorithm [12] within time $O(\sqrt{|V_{\mathcal{K}'}|}|E_{\mathcal{K}'}|) = O(m^{5/2}n^{3/2})$, which yields the claimed time bound.
- finally, lines 7 (compute_auxiliary_bjcts()) and 9 (extend_paths()) take time at most O(m).

To obtain the total time complexity of compute_Lehman-Ron_paths(), it is sufficient to observe that the depth of the recursion stack (originating from line 8) equals the distance d = s - r between the families of sets that were originally given as input, \mathcal{R} and \mathcal{S} , and that the most expensive computation at each step of the recursion is clearly the maximum cardinality matching computation that is performed on the auxiliary bipartite graph \mathcal{K}' . Therefore, we conclude that the worst-case time complexity of compute_Lehman-Ron_paths() is $O(m^{5/2}n^{3/2}d)$.

3.6 A Polynomial Time Algorithm For Solving HY-STCON

We now describe a polynomial time algorithm for solving HY-STCON, called solve_HY-STCON(), which takes as input an instance $\langle S, T, \mathcal{F}, n \rangle$ of HY-STCON, and returns a pair $\langle \text{YES}, \mathbf{p} \rangle$ where \mathbf{p} is a directed path in \mathcal{H}_n that goes from source S to target T avoiding \mathcal{F} if such a path exists (otherwise, the algorithm simply returns NO). Algorithm 2 shows the pseudocode for that procedure. The

Algorithm 3: Breadth-First-Search phases.

Procedure double-bfs_phase($S, T, F, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n$) $\langle \mathcal{S}^*, \ell_{\uparrow}^* \rangle \leftarrow \texttt{bfs_phase}(\mathcal{S}, \mathcal{F}, \ell_{\uparrow}, \ell_{\downarrow}, \texttt{out}, d_{S,T}, n); // \mathsf{BFS}_{\uparrow}$ 1 $\langle \mathcal{T}^*, \ell_{\downarrow}^* \rangle \leftarrow \texttt{bfs_phase}(\mathcal{T}, \mathcal{F}, \ell_{\downarrow}, \ell_{\uparrow}^*, \texttt{in}, d_{S,T}, n); // \mathsf{BFS}_{\downarrow}$ 2 return $\langle \mathcal{S}^*, \mathcal{T}^*, \ell^*_{\uparrow}, \ell^*_{\downarrow} \rangle$; 3 **SubProcedure** $bfs_phase(\mathcal{X}, \mathcal{F}, \ell_x, \ell_y, drt, d_{S,T}, n)$ while $1 \leq |\mathcal{X}| \leq |\mathcal{F}| d_{S,T}$ AND $\ell_x + \ell_y < d_{S,T}$ do 1 $\mathcal{X} \leftarrow \texttt{next_step_bfs}(\mathcal{X}, \mathcal{F}, \texttt{drt}, n);$ $\mathbf{2}$ $\ell_x \leftarrow \ell_x + 1;$ 3 return $\langle \mathcal{X}, \ell_x \rangle;$ 4 SubProcedure $next_step_bfs(\mathcal{X}, \mathcal{F}, drt, n)$ 1 $\mathcal{X}' \leftarrow \emptyset$: foreach $v \in \mathcal{X}$ do 2 $\begin{tabular}{ll} L $\mathcal{X}' \leftarrow \mathcal{X}' \cup N^{\mathtt{drt}}(v) \setminus \mathcal{F}; $$ $// $N^{\mathtt{drt}}$ is $N^{\mathtt{in}}$ if $\mathtt{drt} = \mathtt{in}$, otherwise it is $N^{\mathtt{out}}$ end{tabular} \end{tabular} \end{tabular}$ 3 return \mathcal{X}' ; 4

rationale at the base of $solve_HY-STCON()$ consists in the continuous iteration of two major phases: double-bfs_phase() (line 5) and compression_ phase() (line 11). Throughout computation, both phases alternate repeatedly until a final state of termination is eventually reached (either at line 7, line 10 or line 12). At that point, the algorithm either returns a pair $\langle YES, p \rangle$ where p is the sought directed path, or a negative response NO instead. We now describe both phases in more detail, and give the corresponding pseudocode.

Breadth-First Search phases The first search BFS_{\uparrow} starts from the source vertex S and moves upward, from lower to higher levels of \mathcal{H}_n . Meanwhile, it collects a certain (polynomially bounded) amount of vertices that do not lie in \mathcal{F} . In particular, at the end of any BFS_{\uparrow} phase, the number of collected vertices will always lie between $|\mathcal{F}| d_{S,T} + 1$ and $|\mathcal{F}| d_{S,T} n$ (see line 1 of $bfs_phase()$). The set S of vertices collected at the end of BFS_{\uparrow} is called the *(source) frontier* of BFS_{\uparrow} . All vertices within S have the same cardinality, i.e. $|X_1| = |X_2|$ for every $X_1, X_2 \in S$. Also, the procedure keeps track of the highest level of depth ℓ_{\uparrow} that is reached during BFS_{\uparrow} . Thus, ℓ_{\uparrow} corresponds to the distance between the source vertex S and the current frontier S, formally, $\ell_{\uparrow} = |X| - |S|$ for every $X \in S$. Since at the beginning of the computation BFS_{\uparrow} starts from the source vertex S, solve_HY-STCON() initializes S to $\{S\}$ and ℓ_{\uparrow} to 0 at line 2.

Similarly, the second search BFS_{\downarrow} starts from the target vertex T and moves downward, from higher to lower levels of \mathcal{H}_n , also collecting a certain (polynomially bounded) amount of vertices that do not lie in \mathcal{F} . As in the previous case, this amount will always lie between $|\mathcal{F}| d_{S,T} + 1$ and $|\mathcal{F}| d_{S,T} n$. The set \mathcal{T} of vertices collected at the end of BFS_{\downarrow} is called the *(target) frontier* of BFS_{\downarrow} . All vertices within \mathcal{T} have the same cardinality. Also, the procedure keeps track of the *lowest* level of depth ℓ_{\downarrow} that BFS_{\downarrow} has reached. Thus, ℓ_{\downarrow} corresponds to the distance between the target vertex T and the frontier \mathcal{T} , so that $\ell_{\downarrow} = |T| - |X|$ for every $X \in \mathcal{T}$. Since at the beginning of the computation, BFS_{\downarrow} starts from the target vertex T, solve_HY-STCON() initializes $\mathcal{T} = \{T\}$ and $\ell_{\downarrow} = 0$ at line 3. Fig. 3 provides an illustration of the behaviour of double-bfs_phase().

In summary, after any round of double-bfs_phase(), we are left with two (possibly empty) frontier sets S and T. In Algorithm 2, whenever $S = \emptyset$ or $T = \emptyset$ holds at line 6, then at least one frontier set could not proceed one level further in \mathcal{H}_n while avoiding \mathcal{F} , and thus the procedure halts by returning NO at line 7. Similarly, whenever $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$ and $S \cap T = \emptyset$ holds at line 6, the computation halts by returning NO at line 7 — the underlying intuition being that S and Thave finally reached one another's level of depth without intersecting each other, which means that \mathcal{H}_n contains no directed path from S to T that avoids \mathcal{F} .



Fig. 3: A double_bfs_phase() on \mathcal{H}_3 that starts from $S = \emptyset$ and $T = \{1, 2, 3\}$. The forbidden vertices are $\mathcal{F} = \{\{2\}, \{3\}, \{1, 2\}, \{2, 3\}\}$, while the edges explored by BFS₁ and BFS₁ are $(\emptyset, \{1\})$ and $(\{1, 2, 3\}, \{1, 3\})$ (respectively).

On the other hand, if both $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$ and $S \cap T \neq \emptyset$ hold at line 8, then we can prove that for every $S' \in S$, there exists at least one directed path in \mathcal{H}_n that goes from the source S to S' avoiding \mathcal{F} . Similarly, for every $T' \in \mathcal{T}$, there exists at least one directed path in \mathcal{H}_n that goes from T' to target T avoiding \mathcal{F} . Therefore, whenever $S \cap T \neq \emptyset$, the algorithm is in the right position to reconstruct a directed path p in \mathcal{H}_n that goes from source S to $S \cap \mathcal{T}$ and from $S \cap \mathcal{T}$ to target T avoiding \mathcal{F} (line 9). In practice, the reconstruction can be implemented by maintaining a map throughout the computation, which associates to every vertex v (possibly visited during the BFSs) the *parent vertex*, parent(v), which led to discover v first. As soon as p gets constructed, $solve_HY$ -STCON() returns (YES, p) at line 10, and the computation halts.

Compression Phase After double-bfs_phase() has completed, the procedure solve_HY-STCON() also needs to handle the case where $S, T \neq \emptyset$ and $\ell_{\downarrow} + \ell_{\uparrow} < d_{S,T}$. The phase that starts at that point is named compression_phase() (see Algorithm 4). This procedure takes as input a tuple $\langle S, T, F, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n \rangle$, where S and T are the current frontier sets. Recall that $|T| > |F| d_{S,T}$ holds

Algorithm 4: Compression phase.

F	$\textbf{Procedure compression_phase}(\mathcal{S}, \mathcal{T}, \mathcal{F}, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n)$
1	$\mathcal{T}' \leftarrow \emptyset;$
2	while TRUE do
3	$\mathcal{G} \leftarrow \texttt{construct_bipartite_graph}(\mathcal{S}, \mathcal{T}, n);$
4	$\mathcal{M} \gets \texttt{compute_max_matching}(\mathcal{G}, \mathcal{F} + 1);$
5	$ ext{ if } \mathcal{M} > \mathcal{F} ext{ then }$
6	$\mathcal{M}_{\mathcal{S}} \leftarrow \{ X \in \mathcal{S} \mid \exists Y \in \mathcal{T} \text{ s.t } (X, Y) \in \mathcal{M} \};$
7	$\mathcal{M}_{\mathcal{T}} \leftarrow \{ Y \in \mathcal{T} \mid \exists X \in \mathcal{S} \text{ s.t. } (X, Y) \in \mathcal{M} \};$
8	$\{p_1, \dots, p_{ \mathcal{M} }\} \leftarrow \texttt{compute_Lehman-Ron_paths}(\mathcal{M}_{\mathcal{S}}, \mathcal{M}_{\mathcal{T}}, \mathcal{M}, n);$
9	$p \leftarrow \texttt{reconstruct_path}(\mathcal{S}, \mathcal{T}, \{p_i\}_{i=1}^{ \mathcal{M} }, n);$
10	return $\langle YES, p \rangle;$
11	$\mathcal{X} \leftarrow \texttt{compute_min_vertex_cover}(\mathcal{G}, \mathcal{M});$
12	$\mathcal{X}_{\mathcal{S}} \leftarrow \mathcal{X} \cap \mathcal{S}; \mathcal{X}_{\mathcal{T}} \leftarrow \mathcal{X} \cap \mathcal{T};$
13	$\mathcal{T}' \leftarrow \mathcal{T}' \cup \mathcal{X}_{\mathcal{T}};$
14	$\langle \mathcal{S}, \mathcal{T}, \ell_{\uparrow}, \ell_{\downarrow} angle \leftarrow \texttt{double-bfs_phase}(\mathcal{X}_{\mathcal{S}}, \mathcal{T}, \mathcal{F}, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n);$
15	$ \ \ {\bf if} \ \ {\cal S}=\emptyset \ \ {\it OR} \ \ (\ell_{\downarrow}+\ell_{\uparrow}=d_{S,T} \ \ {\it AND} \ \ {\cal S}\cap {\cal T}=\emptyset \) \ {\bf then} \ \ $
16	$\mathbf{return} \ \mathcal{T}';$
17	$\mathbf{if} \hspace{0.1in} \ell_{\uparrow} + \ell_{\downarrow} = d_{S,T} \hspace{0.1in} \textit{AND} \hspace{0.1in} \mathcal{S} \cap \mathcal{T} \neq \emptyset \hspace{0.1in} \mathbf{then}$
18	$p \gets \texttt{reconstruct_path}(\mathcal{S}, \mathcal{T}, n);$
19	return $\langle Y\!E\!S, p \rangle;$

due to line 1 of **bfs_phase()**. Also, $\mathcal{F} \subseteq \wp_n$ is the set of forbidden vertices; ℓ_{\uparrow} is the level counter of \mathcal{S} and ℓ_{\downarrow} is that of \mathcal{T} ; finally $d_{S,T}$ is the distance between the source S and the target T, and n is the size of the ground set. The output returned by compression_phase() is either a path **p** going from source S to target T avoiding \mathcal{F} or a subset $\mathcal{T}' \subset \mathcal{T}$ such that the following two basic properties hold:

(1) $|\mathcal{T}'| \leq |\mathcal{F}| d_{S,T}$, and (2) if **p** is any directed path in \mathcal{H}_n going from \mathcal{S} to \mathcal{T} avoiding \mathcal{F} , then **p** goes from \mathcal{S} to \mathcal{T}' .

This frontier set \mathcal{T}' is dubbed the *compression* of \mathcal{T} . The underlying rationale goes as follows. On one hand, because of (1), it is possible to keep the search going on by applying yet another round of double-bfs_phase() on input \mathcal{S} and \mathcal{T}' (in fact, the size of \mathcal{T} has been compressed down to $|\mathcal{T}'| \leq |\mathcal{F}| d_{S,T}$, thus matching the threshold condition " $|\mathcal{X}| \leq |\mathcal{F}| d_{S,T}$ " checked at line 1 of bfs_phase()). On the other hand, because of (2), it is indeed sufficient to seek for a directed path in \mathcal{H}_n that goes from \mathcal{S} to \mathcal{T}' avoiding \mathcal{F} , namely, the search can actually forget about $\mathcal{T} \setminus \mathcal{T}'$ because it leads to a dead end. We now describe compression_phase() in more details, and give a graphical summary in Fig. 4. The procedure repeatedly builds an undirected bipartite graph $\mathcal{G} = (V_{\mathcal{G}}, E_{\mathcal{G}})$, where $V_{\mathcal{G}} = \mathcal{S} \cup \mathcal{T}$ and every vertex $U \in \mathcal{S}$ is adjacent to a vertex $V \in \mathcal{T}$ if and only if $U \subset V$. It then uses the procedure compute_max_matching() to find a matching \mathcal{M} of size $|\mathcal{M}| = \min(m^*, |\mathcal{F}| + 1)$, where m^* denotes the size of



Fig. 4: The frontier sets of the compression_phase().

a maximum cardinality matching of \mathcal{G} . Notice that the following holds due to line 1 of bfs_phase():

$$|V_{\mathcal{G}}| = |\mathcal{S}| + |\mathcal{T}| \le 2 |\mathcal{F}| d_{S,T} n,$$

thus, we have the following bound on the size of its edge set:

$$|E_{\mathcal{G}}| \le |V_{\mathcal{G}}|^2 \le 4 |\mathcal{F}|^2 d_{S,T}^2 n^2$$

The fact is that, given that we are content with a cardinality matching of size at most $k = |\mathcal{F}| + 1$, it is worth applying the following recursive self-reduction(\mathcal{G}, k) (Algorithm 5), on input $(\mathcal{G}, |\mathcal{F}| + 1)$, in order to shrink the upper bound on the size of $|E_{\mathcal{G}}|$ from $|V_{\mathcal{G}}|^2$ down to $|V_{\mathcal{G}}| \cdot |\mathcal{F}|$: at line 1, $\mathcal{M} \leftarrow \emptyset$ is initialized to the empty set. At line 2, if k = 0, the empty matching $\mathcal{M} = \emptyset$ is returned. Then, at line 3, let $\hat{v} \in V$ be some vertex having maximum degree $\delta(\hat{v})$ in \mathcal{G} . If $\delta(\hat{v}) < k$ at line 4, the Hopcroft-Karp's algorithm [12] is invoked at line 5 to compute a matching \mathcal{M} of \mathcal{G} such that $|\mathcal{M}| = \min(m^*, k)$, where m^* is the maximum cardinality of any matching in \mathcal{G} . In practice, this step can be implemented in the same manner as a maximum cardinality matching procedure, e.g. as Hopcroft-Karp's algorithm [12], although with the following basic variation: if the size of the augmenting matching \mathcal{M} eventually reaches the cutoff value k, then compute_max_matching() returns \mathcal{M} and halts (i.e. even if $m^* > k$). Otherwise, $\delta(\hat{v}) \ge k$ holds at line 6. So, at line 7, let \mathcal{G}' be the graph obtained from \mathcal{G} by removing \hat{v} and all of its adjacent edges; next, it is invoked self-reduction($\mathcal{G}', k-1$) at line 8, recursively; and, then, the returned matching is assigned to \mathcal{M}' . Since $\delta(\hat{v}) \geq k$, there must be at least one edge $\{u, \hat{v}\} \in E_{\mathcal{G}}$ such that u is not matched in \mathcal{M}' , therefore, $\{u, \hat{v}\}$ is added to \mathcal{M}' ; and the corresponding matching is assigned to \mathcal{M} , at line 9. Finally, \mathcal{M}

Algorithm 5: Self-Reduction for computing \mathcal{M} .

Procedure self-reduction(\mathcal{G}, k) $\mathcal{M} \leftarrow \emptyset$: 1 if k = 0 then return \mathcal{M} ; 2 $\hat{v} \leftarrow \text{pick one vertex } \hat{v} \in V \text{ having maximum degree } \delta(\hat{v}) \text{ in } \mathcal{G};$ 3 if $\delta(\hat{v}) < k$ then 4 $\mathcal{M} \leftarrow \text{compute a matching } \mathcal{M} \text{ of } \mathcal{G} \text{ s.t. } |\mathcal{M}| = \min(m^*, k), \text{ with the}$ 5 Hopcroft-Karp's algorithm [12]; if $\delta(\hat{v}) \geq k$ then 6 $\mathcal{G}' \leftarrow$ remove \hat{v} from \mathcal{G} ; and call the resulting graph \mathcal{G}' ; 7 8 $\mathcal{M}' \leftarrow \texttt{self-reduction}(\mathcal{G}', k-1);$ $\mathcal{M} \leftarrow$ there must be at least one edge $\{u, \hat{v}\} \in E_{\mathcal{G}}$ such that u is not 9 matched in \mathcal{M}' , therefore, add $\{u, \hat{v}\}$ to \mathcal{M}' ; and assign the resulting matching to \mathcal{M} ; return \mathcal{M} ; $\mathbf{10}$

is returned at line 10. In so doing, as shown in Lemma 12, the complexity of compute_max_matching(), at line 4 of compression_phase() (Algorithm 4), is going to improve by a factor $n \cdot d_{S,T}$.

The course of the next actions depends on $|\mathcal{M}|$:

- 1. If $|\mathcal{M}| = |\mathcal{F}| + 1$, then the procedure relies on Theorem 3 to compute a family $p_1, p_2, \ldots, p_{|\mathcal{M}|}$ of $|\mathcal{M}|$ vertex-disjoint directed paths in \mathcal{H}_n that go from \mathcal{S} to \mathcal{T} . In order to do that, the procedure considers the subset $\mathcal{M}_{\mathcal{S}} \subseteq \mathcal{S}$ (resp. $\mathcal{M}_{\mathcal{T}} \subseteq \mathcal{T}$) of all vertices in \mathcal{S} (resp. in \mathcal{T}) that are incident to some edge in \mathcal{M} (lines 6 and 7). Notice that the matching \mathcal{M} can be viewed as a bijection between $\mathcal{M}_{\mathcal{S}}$ and $\mathcal{M}_{\mathcal{T}}$. Then, the algorithm underlying Theorem 3 gets invoked on input $\langle \mathcal{M}_{\mathcal{S}}, \mathcal{M}_{\mathcal{T}}, \mathcal{M}, n \rangle$ (line 8). Once all the Lehman-Ron paths $p_1,p_2,\ldots,p_{|\mathcal{M}|}$ have been found, it is then possible to reconstruct the sought directed path p in \mathcal{H}_n that goes from source S to target T avoiding \mathcal{F} (line 9). In fact, since $|\mathcal{M}| > |\mathcal{F}|$ by hypothesis, and since $\mathsf{p}_1, \mathsf{p}_2, \ldots, \mathsf{p}_{|\mathcal{M}|}$ are distinct and pairwise vertex-disjoint, there must exist at least one path \mathbf{p}_i that goes from \mathcal{S} to \mathcal{T} avoiding \mathcal{F} . It is therefore sufficient to find such a path $\mathbf{p}_i = v_0 v_1 \cdots v_k$ by direct inspection. At that point, it is possible to reconstruct a path **p** going from S to v_0 (because $v_0 \in S$), as well as a path going from v_k to T (because $v_k \in \mathcal{T}$). As already mentioned, in practice, the reconstruction can be implemented by maintaining a map that associates to every vertex v (eventually visited during the BFSs) the parent vertex that had led to discover v first. Then, $\langle YES, p \rangle$ is returned at line 10.
- 2. If $|\mathcal{M}| \leq |\mathcal{F}|$, then the compression_phase() aims to compress the size of \mathcal{T} down to $|\mathcal{T}'| \leq |\mathcal{F}| d_{S,T}$ as follows. Notice that in this case \mathcal{M} is a maximum cardinality matching of \mathcal{G} , because $|\mathcal{M}| \leq |\mathcal{F}|$. So, the algorithm computes a minimum cardinality vertex-cover \mathcal{X} of \mathcal{G} at line 11, whose size is $|\mathcal{M}|$ by König's theorem [6]. The algorithm then proceeds at line 12 by considering

the set $\mathcal{X}_{\mathcal{S}} = \mathcal{X} \cap \mathcal{S}$ (resp. $\mathcal{X}_{\mathcal{T}} = \mathcal{X} \cap \mathcal{T}$) of all vertices that lie both in the vertex-cover \mathcal{X} and in the frontier set \mathcal{S} (resp. \mathcal{T}). Here, it is crucial to notice that both $|\mathcal{X}_{\mathcal{S}}| \leq |\mathcal{F}|$ and $|\mathcal{X}_{\mathcal{T}}| \leq |\mathcal{F}|$ hold, because $|\mathcal{X}| = |\mathcal{M}| \leq |\mathcal{F}|$. The fact that, since \mathcal{X} is a vertex-cover of \mathcal{G} , any directed path in \mathcal{H}_n that goes from \mathcal{S} to \mathcal{T} must go either from $\mathcal{X}_{\mathcal{S}}$ to \mathcal{T} or from $\mathcal{S} \setminus \mathcal{X}_{\mathcal{S}}$ to $\mathcal{X}_{\mathcal{T}}$ plays a pivotal role. Stated otherwise, there exists no directed path in \mathcal{H}_n that goes from $\mathcal{S} \setminus \mathcal{X}_{\mathcal{S}}$ to $\mathcal{T} \setminus \mathcal{X}_{\mathcal{T}}$, simply because \mathcal{X} is a vertex cover of \mathcal{G} . At that point, the compression \mathcal{T}' gets enriched with $\mathcal{X}_{\mathcal{T}}$ at line 13.

Then, compression_phase() seeks a directed path in \mathcal{H}_n that eventually goes from \mathcal{X}_S to \mathcal{T} . This is done at line 14 by running double-bfs_phase() on $\langle \mathcal{X}_S, \mathcal{T}, \mathcal{F}, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n \rangle$. Since $|\mathcal{X}_S| \leq |\mathcal{F}|$, that execution results into an update of both the frontier set S and of its level counter ℓ_{\uparrow} . Let $\mathcal{S}^{(i+1)}$ be the updated value of S and let $\ell_{\uparrow}^{(i+1)}$ be that of ℓ_{\uparrow} . Note that, since $|\mathcal{T}| > |\mathcal{F}| d_{S,T}$ holds as a pre-condition of compression_phase(), neither \mathcal{T} nor ℓ_{\downarrow} are ever updated at line 14. Upon completion of this supplementary double-bfs_phase(), if $\mathcal{S}^{(i+1)} = \emptyset$ or both $\ell_{\uparrow}^{(i+1)} + \ell_{\downarrow} = d_{S,T}$ and $\mathcal{S}^{(i+1)} \cap \mathcal{T} = \emptyset$ at line 15, then \mathcal{T}' is returned at line 16 of compression_phase().

Otherwise, if $\ell_{\uparrow}^{(i+1)} + \ell_{\downarrow} = d_{S,T}$ and $\mathcal{S}^{(i+1)} \cap \mathcal{T} \neq \emptyset$ at line 17, the sought directed path \mathbf{p} in \mathcal{H}_n that goes from source S to target T avoiding \mathcal{F} can be reconstructed from $\mathcal{S}^{(i+1)}$ and \mathcal{T} at line 18, so that compression_phase() returns $\langle \text{YES}, \mathbf{p} \rangle$ and halts soon after at line 19.

Otherwise, if $S^{(i+1)} \neq \emptyset$ and $\ell_{\uparrow}^{(i+1)} + \ell_{\downarrow} < d_{S,T}$, the next iteration will run on the novel frontier set $S^{(i+1)}$ and its updated level counter $\ell_{\uparrow}^{(i+1)}$. It is not difficult to prove that each iteration increases ℓ_{\uparrow} by at least one unit, so that the while-loop at line 2 of compression_phase() can be iterated at most $d_{S,T}$ times overall. In particular, this fact implies that $|\mathcal{T}'| \leq |\mathcal{F}| d_{S,T}$ always holds at line 16 of compression_phase().

Fig. 4 illustrates the family of all frontier sets considered throughout compression_phase(), where the following notation is assumed: max_i is the total number of iterations of the while-loop at line 2 of compression_phase(), $\mathcal{X}^{(i)}$ is the vertex-cover computed at the *i*th iteration of line 11, $\mathcal{X}_{S}^{(i)}$ and $\mathcal{X}_{T}^{(i)}$ are the sets computed at the *i*th iteration of line 12, and $\mathcal{S}^{(i)}$ is the frontier set computed at the *i*th iteration of line 14. The compression of \mathcal{T} (possibly returned at line 16) is $\mathcal{T}' = \bigcup_{i=1}^{\max_i} \mathcal{X}_{T}^{(i)}$.

3.7 A Remark On Decision Versus Search

Algorithm 2 tackles the SEARCH-TASK of HY-STCON. If we merely want to answer the DECISION-TASK instead, we can simplify the algorithm by immediately returning YES if $|\mathcal{M}| > |\mathcal{F}|$ at line 5 of compression_phase(). This is because in that case, Theorem 2 guarantees the existence of a family of $|\mathcal{M}| > |\mathcal{F}|$ vertexdisjoint paths in \mathcal{H}_n that go from the current source frontier \mathcal{S} to the target frontier \mathcal{T} , which suffices to conclude that at least one of those paths avoids \mathcal{F} . This simplification improves the time complexity of our algorithm for solving the DECISION-TASK by a polynomial factor over that for the SEARCH-TASK.

3.8 Correctness Analysis of Algorithm 2

The present subsection aims to show that the procedure **solve_Hy-stCon()** is correct. A formal statement of that is provided in the next theorem.

Theorem 5. Let $\mathcal{I} = \langle S, T, \mathcal{F}, n \rangle$ be any instance of HY-STCON. Given \mathcal{I} as input, the procedure solve_HY-STCON() halts within a finite number of steps. Moreover, it returns as output a directed path \mathbf{p} in \mathcal{H}_n that goes from source S to target T avoiding \mathcal{F} , provided that at least one such path exists; otherwise, the output is simply NO.

We are going to show a sequence of results that shall ultimately lead us to prove Theorem 5. Hereafter, it is assumed that $\langle S, T, \mathcal{F}, n \rangle$ is an instance (of HY-STCON) given as input to the solve_HY-STCON() procedure. Lemmas 1 to 3 below show that procedures double-bfs_phase() and compression_phase(), which are called by solve_HY-STCON(), halt within a finite number of steps.

Lemma 1. Any invocation of double-bfs_phase() halts within a finite number of steps. In particular, the while-loop at line 1 of the bfs_phase() iterates at most $d_{S,T}$ times.

Proof. Consider the while-loop at line 1 of bfs_phase(). At each iteration of line 3, the level counter ℓ_x gets incremented. Notice that this is the only line at which ℓ_x may be modified, and also notice that ℓ_y is never modified. Therefore, $\ell_x + \ell_y$ can only increase and not decrease. Since the while-loop at line 1 of bfs_phase() halts as soon as $\ell_x + \ell_y = d_{S,T}$, the thesis follows.

Lemma 2. Each iteration of the while-loop at line 2 of compression_phase() increases $\ell_{\uparrow} + \ell_{\downarrow}$ by at least one unit, either until $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$ or until the procedure halts by reaching either line 10, line 16 or line 19.

Proof. Consider any iteration of the while-loop at line 2 of compression_phase(). Let \mathcal{G} be the bipartite graph computed at line 3, and let \mathcal{M} be the matching of \mathcal{G} computed at line 4. If $|\mathcal{M}| > |\mathcal{F}|$, then line 10 gets executed, so the procedure halts within a finite number of steps by virtue of our discussion in Section 3.5. Otherwise $|\mathcal{M}| \leq |\mathcal{F}|$. Recall that, since $|\mathcal{M}| \leq |\mathcal{F}|$, then \mathcal{M} is a maximum matching of \mathcal{G} ; also recall that $\mathcal{X}_{\mathcal{S}} = \mathcal{X} \cap \mathcal{S}$ where \mathcal{X} is a minimum vertex cover of \mathcal{G} (line 12). Since $|\mathcal{X}| = |\mathcal{M}|$, then $|\mathcal{X}_{\mathcal{S}}| \leq |\mathcal{X}| = |\mathcal{M}| \leq |\mathcal{F}|$. Moreover, since $|\mathcal{M}| \leq |\mathcal{F}|$, double-bfs_phase() gets invoked at line 14 on input $\langle \mathcal{X}_{\mathcal{S}}, \mathcal{T}, \mathcal{F}, \ell_{\downarrow}, \ell_{\uparrow}, d_{S,T}, n \rangle$ and halts within a finite number of steps by Lemma 1. Let us analyze its behavior with respect to $\mathcal{X}_{\mathcal{S}}$. If $\mathcal{X}_{\mathcal{S}} = \emptyset$, then double-bfs_phase() returns an empty frontier set \mathcal{S} as output, which leads to the termination of compression_phase() at line 16. Moreover, if $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$, then compression_phase() halts either at line 16 or at line 19. Otherwise, we

must have $1 \leq |\mathcal{X}_S| \leq |\mathcal{F}|$ and $\ell_{\uparrow} + \ell_{\downarrow} < d_{S,T}$, in that case the condition for entering the while-loop at line 1 of the bfs_phase() is satisfied; therefore, at line 3 of bfs_phase(), the level counter ℓ_{\uparrow} gets incremented. This implies the thesis.

Lemma 3. Any invocation of compression_phase() halts within a finite number of steps. In particular, the while-loop at line 2 of the compression_phase() iterates at most $d_{S,T}$ times.

Proof. Firstly, recall Lemma 2. Then, notice that as soon as $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$ the compression_phase() then halts either at line 16 (if $S \cap T = \emptyset$) or at line 19 (if $S \cap T \neq \emptyset$). This implies that the while-loop at line 2 of compression_phase() iterates at most $d_{S,T}$ times.

We now prove some useful properties of compression_phase() and solve_HY-STCON().

Lemma 4. The following invariant is maintained at each line of solve_HY-STCON() and at each line of compression_phase(). For every $S' \in S$ there exists a directed path in \mathcal{H}_n that goes from S to S' avoiding \mathcal{F} ; similarly, for every $T' \in \mathcal{T}$ there is a directed path in \mathcal{H}_n that goes from T' to T avoiding \mathcal{F} .

Proof. At the beginning of the procedure $S = \{S\}$ and $T = \{T\}$, so the thesis holds. At each subsequent step, the only way in which a novel vertex can be added either to S or T is by invoking the double_bfs_phase(), which preserves connectivity and avoids \mathcal{F} by construction at line 3 of next_step_bfs().

Lemma 5. Assume that any invocation of compression_phase() halts by returning $\langle YES, p \rangle$. Then p is a directed path in \mathcal{H}_n that goes from source S to target T avoiding \mathcal{F} .

Proof. If compression_phase() returns p as output, then the last iteration of the while-loop at line 2 must reach either line 10 or line 19:

1. Assume that line 10 is reached at the last iteration. Then, during that iteration, the matching \mathcal{M} (computed at line 4 on input \mathcal{G}) has size $|\mathcal{M}| > |\mathcal{F}|$. Recall that \mathcal{G} is a bipartite graph on bipartition $(\mathcal{S}, \mathcal{T})$. Let $\mathcal{M}_{\mathcal{S}}$ (resp. $\mathcal{M}_{\mathcal{T}}$ be the subset of all vertices in \mathcal{S} (resp. \mathcal{T}) that belong to some edge in \mathcal{M} . Then, by Theorem 2, there exist $|\mathcal{M}|$ vertex-disjoint directed paths in \mathcal{H}_n , say $\mathsf{p}_1, \mathsf{p}_2, \ldots, \mathsf{p}_{|\mathcal{M}|}$, whose union contains all the vertices in $\mathcal{M}_{\mathcal{S}}$ and $\mathcal{M}_{\mathcal{T}}$. Since $|\mathcal{M}| > |\mathcal{F}|$, at least one of those paths — say, $\mathsf{p}_i = v_0 \cdots v_k$ — must avoid \mathcal{F} . By Proposition 4, the procedure reconstruct_path() (invoked at line 9) is able to compute a directed path p_{S,v_0} in \mathcal{H}_n that goes from S to v_0 avoiding \mathcal{F} (because $v_0 \in \mathcal{S}$, being the first step of p_i), and it is also able to compute a directed path $\mathsf{p}_{v_k,T}$ that goes from v_k to T avoiding \mathcal{F} (because $v_k \in \mathcal{T}$, being the last step of p_i). Let $\mathsf{p} = \mathsf{p}_{S,v_0}\mathsf{p}_i\mathsf{p}_{v_k,T}$ be the directed path obtained by concatenation. compression_phase() then returns p at line 10. 2. Assume that line 19 is reached at the last iteration. Then, at that iteration, the condition checked at line 17 of compression_phase() must be satisfied; that is, we have $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$ and $S \cap \mathcal{T} \neq \emptyset$. Let X be an arbitrary vertex in $\mathcal{S} \cap \mathcal{T}$. By Lemma 4, there exists at least one directed path $\mathbf{p}_{S,X}$ in \mathcal{H}_n that goes from S to X avoiding \mathcal{F} (because $X \in \mathcal{S}$); similarly, there exists at least one directed path $p_{X,T}$ in \mathcal{H}_n that goes from X to T avoiding \mathcal{F} (because $X \in \mathcal{T}$). Therefore, during that iteration, the procedure reconstruct_path() (invoked at line 18) is able to compute a path $\mathbf{p} = \mathbf{p}_{S,X} \mathbf{p}_{X,T}$ that goes from S to X, and then from X to T, which is the result returned by compression_phase() at line 19.

The following result shows two useful properties of the frontier set returned by compression_phase(), for which we will need additional notation. Denote by \max_i be the number of times that the while-loop at line 2 gets iterated throughout the whole execution of the compression_phase().

Also, let us introduce the following notation, for each index $i \in [\max_i]$:

- let $\mathcal{X}^{(i)}$ be the vertex cover that is computed during the *i*-th iteration of line 11;
- let $\mathcal{X}_{\mathcal{S}}^{(i)}$ and $\mathcal{X}_{\mathcal{T}}^{(i)}$ be the sets computed during the *i*-th iteration of line 12; let $\mathcal{S}^{(i)}$ be the novel frontier set that is computed during the *i*-th iteration of line 14:

Moreover, we assume the notation $\mathcal{S}^{(0)} = \mathcal{S}$, so that $\mathcal{X}^{(i)}_{\mathcal{S}} = \mathcal{S}^{(i-1)} \cap \mathcal{X}^{(i)}$ holds for each iteration $i \in [\max_i]$. Notice that, since $|\mathcal{T}| > |\mathcal{F}| d_{S,T}$ holds by hypothesis, then \mathcal{T} is not modified, at line 14, by the invocation of double-bfs_phase(). Indeed, \mathcal{T} is never modified throughout the compression_phase(). Nevertheless, a novel set $\mathcal{T}' \subset \mathcal{T}$ gets constructed and possibly returned.

Proposition 4. Assume that the procedure compression_phase() is invoked on input $\langle S, T, F, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n \rangle$, where $|T| > |F| d_{S,T}$ is required to hold as a pre-condition. Also, assume that the procedure halts at line 16, returning a novel frontier set $\mathcal{T}' \subset \mathcal{T}$. Then, the following properties hold:

- 1. $|\mathcal{T}'| \leq |\mathcal{F}| d_{S,T};$
- 2. if p is any directed path in \mathcal{H}_n that goes from S to T avoiding F, then p goes from S to T'.

Proof. Firstly notice that, if an invocation of the compression_phase() halts at line 16 by returning a novel frontier set $\mathcal{T}' \subset \mathcal{T}$, this means that neither line 10 nor line 19 are ever reached throughout that invocation. In particular this implies that, at each iteration i of the while-loop at line 2, the maximum matching $\mathcal{M}^{(i)}$ (computed at line 4) has size $|\mathcal{M}^{(i)}| \leq |\mathcal{F}|$; this fact is assumed throughout the whole proof.

1. Proof of (1). At each iteration $i \in [\max_i]$, the minimum vertex cover $\mathcal{X}^{(i)}$ has size:

$$|\mathcal{X}^{(i)}| = |\mathcal{M}^{(i)}| \le |\mathcal{F}|.$$

Since $\mathcal{X}_{\mathcal{T}}^{(i)} = \mathcal{X}^{(i)} \cap \mathcal{T}$ at line 12, then $|\mathcal{X}_{\mathcal{T}}^{(i)}| \leq |\mathcal{X}^{(i)}| \leq |\mathcal{F}|$. Moreover, recall that \mathcal{T}' gets enriched by $\mathcal{X}^{(i)}$ at each iteration of line 13, so that the following holds at the termination of the compression_phase():

$$\mathcal{T}' = \bigcup_{i=1}^{\max_i} \mathcal{X}_{\mathcal{T}}^{(i)}$$

Also recall that, by Lemma 3, the while-loop at line 2 can be iterated at most $d_{S,T}$ times, so that $\max_i \leq d_{S,T}$. Therefore, when compression_phase() terminates, we have $|\mathcal{T}'| \leq |\mathcal{F}| d_{S,T}$.

2. Proof of (2). In order to prove (2), we exhibit a number of invariants which hold for each iteration of the while-loop at line 2 of compression_phase(). In what follows, we assume that the procedure compression_phase() gets invoked on input $\langle S, T, F, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n \rangle$, and that $S^{(0)} = S$ holds by notational convention.

Lemma 6. Let $i \in [\max_i]$ be any iteration of the while-loop at line 2 of compression_phase(). Let p be any directed path in \mathcal{H}_n that goes from $\mathcal{S}^{(i-1)}$ to \mathcal{T} . Then p goes either from $\mathcal{X}_S^{(i)}$ to \mathcal{T} or from $\mathcal{S}^{(i-1)} \setminus \mathcal{X}_S^{(i)}$ to $\mathcal{X}_{\mathcal{T}}^{(i)}$. In other words, there exists no directed path in \mathcal{H}_n that goes from $\mathcal{S}^{(i-1)} \setminus \mathcal{X}_S^{(i)}$ to $\mathcal{T} \setminus \mathcal{X}_{\mathcal{T}}^{(i)}$.

Proof. Recall that $\mathcal{X}^{(i)}$ is a vertex cover of the bipartite graph defined as $\mathcal{G}^{(i)} = ((\mathcal{S}^{(i-1)}, \mathcal{T}), \subset)$, which is constructed during the *i*-th iteration of line 3 within the procedure compression_phase(). Also, $\mathcal{X}^{(i)}_{\mathcal{S}} = \mathcal{X}^{(i)} \cap \mathcal{S}^{(i-1)}$ and $\mathcal{X}^{(i)}_{\mathcal{T}} = \mathcal{X}^{(i)} \cap \mathcal{T}$, so that the existence of any directed path in \mathcal{H}_n going from $\mathcal{S}^{(i-1)} \setminus \mathcal{X}^{(i)}_{\mathcal{S}}$ to $\mathcal{T} \setminus \mathcal{X}^{(i)}_{\mathcal{T}}$ would imply the existence of some edge of $\mathcal{G}^{(i)}$ that would be uncovered by $\mathcal{X}^{(i)}$, contradicting the fact that $\mathcal{X}^{(i)}$ is vertex cover of $\mathcal{G}^{(i)}$.

Fig. 5 illustrates the intuition underlying Lemma 6.

Lemma 7. Let $i \in [\max_i]$ be any iteration of the while-loop at line 2 of compression_phase(). Let U be any subset of $S^{(i-1)}$ and let V be any subset of \mathcal{T} . Let p be any directed path in \mathcal{H}_n that goes from U to V. Then p goes from S to V in \mathcal{H}_n .

Proof. Induction on $i \in [\max_i]$.

- Base Case. If i = 1, recall that $\mathcal{S}^{(0)} = \mathcal{S}$. Then $U \subseteq \mathcal{S}$, which implies the base case.
- Inductive Step. Let us assume, by induction hypothesis, that the claim holds for some $i \in [\max_i -1]$ and let us prove it for i + 1. So, let $U \subseteq S^{(i)}$, and let **p** by any directed path in \mathcal{H}_n that goes from U to V. Recall that $S^{(i)}$ is the frontier set that is returned by an invocation of double-bfs_phase() on input $\mathcal{X}_{S}^{(i)}$, at the *i*-th iteration of line 14, within compression_phase(). This amounts to saying that all vertices in $S^{(i)}$ have been discovered by a BFS starting



Fig. 5: The undirected bipartite graph $\mathcal{G}^{(i)} = ((\mathcal{S}^{(i-1)}, \mathcal{T}), \subset)$, and vertex cover $\mathcal{X}^{(i)} = (X_{\mathcal{S}}^{(i)}, X_{\mathcal{T}}^{(i)})$ (doubly-circular nodes).

from $\mathcal{X}_{\mathcal{S}}^{(i)}$. Recall that $\mathcal{X}_{\mathcal{S}}^{(i)} = \mathcal{X}^{(i)} \cap \mathcal{S}^{(i-1)}$ so that $\mathcal{X}_{\mathcal{S}}^{(i)} \subseteq \mathcal{S}^{(i-1)}$. Therefore, **p** is indeed a directed path in \mathcal{H}_n that goes from $\mathcal{S}^{(i-1)}$ to V in \mathcal{H}_n . By induction hypothesis, the thesis follows.

Lemma 8. Let $i \in [\max_i]$ be any index of iteration of the while-loop at line 2 of compression_phase(). Let p be a directed path in \mathcal{H}_n that goes from S to \mathcal{T} avoiding \mathcal{F} . Then, p goes either from $\mathcal{X}_{\mathcal{S}}^{(i)}$ to \mathcal{T} or from S to $\bigcup_{j=1}^{i} \mathcal{X}_{\mathcal{T}}^{(j)}$.

Proof. Induction on $i \in [\max_i]$.

- Base Case. If i = 1, recall that $\mathcal{S}^{(0)} = \mathcal{S}$. Then, by Lemma 6, we have that \mathbf{p} either goes from $X_{\mathcal{S}}^{(1)}$ to \mathcal{T} or from $\mathcal{S} \setminus \mathcal{X}_{\mathcal{S}}^{(1)}$ to $\mathcal{X}_{\mathcal{T}}^{(1)}$. If \mathbf{p} goes from $\mathcal{S} \setminus \mathcal{X}_{\mathcal{S}}^{(1)}$ to $\mathcal{X}_{\mathcal{T}}^{(1)}$, then clearly \mathbf{p} goes from \mathcal{S} to $\mathcal{X}_{\mathcal{T}}^{(1)}$. This implies the base case.
- Inductive Step. Let us assume, by induction hypothesis, that the claim holds for some $i \in [\max_i -1]$, and let us prove it for i + 1. By induction hypothesis, **p** either goes from $\mathcal{X}_{\mathcal{S}}^{(i)}$ to \mathcal{T} or from \mathcal{S} to $\bigcup_{i=1}^{i} \mathcal{X}_{\mathcal{T}}^{(j)}$ in \mathcal{H}_{n} .

If **p** goes from $\mathcal{X}_{\mathcal{S}}^{(i)}$ to \mathcal{T} avoiding \mathcal{F} in \mathcal{H}_n , then **p** must go from $\mathcal{S}^{(i)}$ to \mathcal{T} : in fact, recall that $\mathcal{S}^{(i)}$ is the frontier set that is returned by

the invocation of double-bfs_phase() on input $\mathcal{X}_{S}^{(i)}$, at line 14 of the compression_phase().

If p goes from $\mathcal{S}^{(i)}$ to \mathcal{T} then, by Lemma 6, we also have that p goes either from $\mathcal{X}_{\mathcal{S}}^{(i+1)}$ to \mathcal{T} or from $\mathcal{S}^{(i)} \setminus \mathcal{X}_{\mathcal{S}}^{(i+1)}$ to $\mathcal{X}_{\mathcal{T}}^{(i+1)}$ in \mathcal{H}_n . If **p** goes from $\mathcal{S}^{(i)} \setminus \mathcal{X}_{\mathcal{S}}^{(i+1)}$ to $\mathcal{X}_{\mathcal{T}}^{(i+1)}$, then **p** goes from \mathcal{S} to $\mathcal{X}_{\mathcal{T}}^{(i+1)}$

by Lemma 7.

Since p either goes from $\mathcal{X}_{\mathcal{S}}^{(i+1)}$ to \mathcal{T} , or from \mathcal{S} to $\mathcal{X}_{\mathcal{T}}^{(i+1)}$, or from \mathcal{S} to $\bigcup_{i=1}^{i} \mathcal{X}_{\mathcal{T}}^{(j)}$ in \mathcal{H}_{n} , we have that **p** either goes from $\mathcal{X}_{S}^{(i+1)}$ to \mathcal{T} , or from \mathcal{S} to $\bigcup_{j=1}^{i+1} \mathcal{X}_{\mathcal{T}}^{(j)}$ in \mathcal{H}_n , thus concluding the induction and the proof of Lemma 8.

We now have everything we need to prove (2). Let $i = \max_i$ be the last iteration of the while-loop at line 2 of compression_phase(). Moreover, assume that **p** is a directed path in \mathcal{H}_n that goes from \mathcal{S} to \mathcal{T} avoiding \mathcal{F} . By Lemma 8, **p** either goes from $\mathcal{X}_{\mathcal{S}}^{(\max_i)}$ to \mathcal{T} or from \mathcal{S} to $\bigcup_{i=1}^{\max_i} \mathcal{X}_{\mathcal{T}}^{(i)}$. We argue that **p** cannot go from $\mathcal{X}_{\mathcal{S}}^{(\max_i)}$ to \mathcal{T} in \mathcal{H}_n . In fact, any such path must first visit $\mathcal{S}^{(\max_i)}$ in order to reach \mathcal{T} . Then, it is sufficient to show that there exists no path that goes from $\mathcal{S}^{(\max_i)}$ to \mathcal{T} . Recall that \max_i is the last iteration of the while-loop at line 2, and by hypothesis the compression_phase() halts by returning \mathcal{T}' at line 16. Therefore, at line 15, it must hold that $\mathcal{S}^{(\max_i)} = \emptyset$ or that both $\ell_{\perp}^{(\max_i)} + \ell_{\uparrow} = d_{S,T}$ and $\mathcal{S}^{(\max_i)} \cap \mathcal{T} = \emptyset$. Thus, there exists no directed path in \mathcal{H}_n that goes from $\mathcal{S}^{(\max_i)}$ to \mathcal{T} .

Since **p** does not go from $\mathcal{X}_{\mathcal{S}}^{(\max_i)}$ to \mathcal{T} , it must go from \mathcal{S} to $\bigcup_{i=1}^{\max_i} \mathcal{X}_{\mathcal{T}}^{(i)}$ instead; and since $\mathcal{T}' = \bigcup_{i=1}^{\max_i} \mathcal{X}_{\mathcal{T}}^{(i)}$, **p** must therefore go from \mathcal{S} to \mathcal{T}' , which concludes the proof of (2).

Now that we have established the correctness of the procedures it uses, we go back to establishing the correctness of solve_HY-STCON().

Lemma 9. Each iteration of the while-loop at line 4 of solve_HY-STCON() increases $\ell_{\uparrow} + \ell_{\downarrow}$ by at least one unit; until $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$ or until the procedure halts by reaching either line 7, line 10 or line 12.

Proof. Induction on the index i of iteration of the while-loop at line 4.

- Base Case. Consider the first iteration of the while-loop at line 4. We have $\mathcal{S} = \{S\}, \mathcal{T} = \{T\}$, and $\ell_{\uparrow} = \ell_{\downarrow} = 0$. Therefore, if $d_{S,T} = 0$, then the procedure halts immediately, either at line 7 (if $S \neq T$) or at line 10 (if S = T). If $d_{S,T} > 0$, then a first execution of double-bfs_phase() is invoked at line 5, which halts after a finite number of steps by Lemma 1. Notice that the condition for entering the while-loop at line 1 of bfs_phase() is satisfied, so $\ell_{\uparrow} + \ell_{\downarrow}$ gets incremented at line 3 of bfs_phase().
- Inductive Step. Assume that at the *i*-th iteration of the while-loop at line 4, we have $\ell_{\uparrow} + \ell_{\downarrow} < d_{S,T}$. Furthermore, assume that none of the conditions checked by solve_HY-STCON() at line 6, line 8 and line 12 are

satisfied. Then, the procedure does not halt at *i*-th iteration. Recall that double-bfs_phase(), which is invoked at line 5, halts within finite time by Lemma 1; also, recall that compression_phase(), which is invoked at line 11, halts within finite time by Lemma 3. Thus, at the end of the *i*-th iteration, line 13 gets finally executed. At line 13, the current frontier \mathcal{T} gets replaced by the value \mathcal{T}' , previously returned by compression_phase() at line 11. Notice that $|\mathcal{T}'| \leq |\mathcal{F}| d_{S,T}$ holds by Proposition 4. The (i + 1)-th iteration of the while-loop at line 4 starts at this point. Then, at line 5, another round of double-bfs_phase() is executed. If $\mathcal{T} \neq \emptyset$ and $\ell_{\uparrow} + \ell_{\downarrow} < d_{S,T}$, the condition for entering the while-loop at line 3. If $\mathcal{T} = \emptyset$ or $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$, then the procedure halts at line 7. This implies that the invariant is maintained for each iteration *i*.

Proposition 5. The procedure solve_HY-STCON () halts within a finite number of steps. In particular, the while-loop at line 4 iterates at most $d_{S,T}$ times.

Proof. Recall the statement of Lemma 9. As soon as $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$, then solve_HY-STCON() halts either at line 7 (if $S \cap T = \emptyset$) or at line 10 (if $S \cap T \neq \emptyset$). In particular, this implies that the while-loop at line 4 of the solve_HY-STCON() can be iterated at most $d_{S,T}$ times.

Proposition 6. Assume that solve_HY-STCON() halts by returning the pair $\langle YES, p \rangle$. Then p is a directed path in \mathcal{H}_n that goes from S to T avoiding \mathcal{F} .

Proof. Observe that solve_HY-STCON() can return $\langle \text{YES}, \mathsf{p} \rangle$ as output only at line 10 or at line 12. In the latter case, p gets constructed at line 11 by invoking compression_phase(), so the thesis follows by Lemma 5. Otherwise, assume that p is returned at line 10. Therefore, at the last iteration of line 8, it must hold that $S \cap \mathcal{T} \neq \emptyset$. Then, let $X \in S \cap \mathcal{T}$. By Lemma 4 there exists a directed path $p_{S,X}$ in \mathcal{H}_n that goes from S to X avoiding \mathcal{F} (because $X \in S$), and there exists another directed path $\mathsf{p}_{X,T}$ in \mathcal{H}_n that goes from X to T avoiding \mathcal{F} (because $X \in \mathcal{T}$). Therefore, reconstruct_path() at line 9, is able to compute a directed path $\mathsf{p} = \mathsf{p}_{S,X}\mathsf{p}_{X,T}$ in \mathcal{H}_n that goes from S to T avoiding \mathcal{F} , which gets returned at line 12.

Lemma 10. The following invariant is maintained at each line of solve_HY-STCON(). If p is any directed path in \mathcal{H}_n that goes from S to T avoiding \mathcal{F} , then p goes from S to \mathcal{T} .

Proof. Induction on the index *i* of iteration of the while-loop at line 2.

- Base Case. Before entering the first iteration, since $S = \{S\}$ and $T = \{T\}$, the thesis holds.
- Inductive Step. Assume that the thesis holds at the end of the *i*-th iteration. So, let $\mathcal{S}^{(i)}$ and $\mathcal{T}^{(i)}$ be the frontier sets at the end of the *i*-th iteration. When i = 0, just recall that $\mathcal{S}^{(0)} = \{S\}$ and $\mathcal{T}^{(0)} = \{T\}$. Now, at the

beginning of the (i + 1)-th iteration, in particular at line 5 of solve_HY-STCON(), let S and T be the frontier sets returned by the invocation of double-bfs_phase(). If p is any directed path in \mathcal{H}_n that goes from S to T avoiding \mathcal{F} , then p goes from $S^{(i)}$ to $\mathcal{T}^{(i)}$ by induction hypothesis. It is not difficult to see that if p goes from $S^{(i)}$ to $\mathcal{T}^{(i)}$ avoiding \mathcal{F} , then p must go from S to T as well: at this point, the reader can check that this is a direct consequence of double-bfs_phase()'s construction. If the (i + 1)-th iteration doesn't halt, then the compression_phase() at line 11 gets invoked. Then, let \mathcal{T}' be the value returned by compression_phase() at line 11. By Proposition 4, if p is a directed path in \mathcal{H}_n that goes from S to \mathcal{T} avoiding \mathcal{F} , then p goes from S to \mathcal{T}' . Thus, it is indeed correct to update \mathcal{T} by \mathcal{T}' at line 13 of solve_HY-STCON(). This implies that the thesis holds for each iteration of the while-loop at line 2, until termination.

Proposition 7. Assume that solve_HY-STCON() halts by returning NO. Then there is no directed path in \mathcal{H}_n that goes from S to T avoiding \mathcal{F} .

Proof. Since solve_HY-STCON() returns NO, the condition checked at line 6 must be satisfied: if $S = \emptyset$ or $\mathcal{T} = \emptyset$, then there exists no directed path in \mathcal{H}_n that goes from S to \mathcal{T} ; similarly, if $\ell_{\uparrow} + \ell_{\downarrow} = d_{S,T}$ and $S \cap \mathcal{T} = \emptyset$, then there exists no directed path in \mathcal{H}_n that goes from S to \mathcal{T} . By Lemma 10, there exists no directed path in \mathcal{H}_n that goes from S to \mathcal{T} avoiding \mathcal{F} .

Theorem 5 follows, at this point, from Propositions 5 to 7.

3.9 Complexity Analysis

We now analyze the time complexity of **solve_HY-STCON()**, starting with that of the procedures it relies on.

Lemma 11. The double-bfs_phase() always halts within $O(|\mathcal{F}| d_{S,T}^2 n)$ time.

Proof. It is sufficient to prove that bfs_phase() always halts within $O(|\mathcal{F}| d_{S,T}^2 n)$ time. Recall that, by Lemma 1, the while-loop at line 1 of bfs_phase() iterates at most $d_{S,T}$ times. At each iteration, next_step_bfs() gets invoked on some input set $\mathcal{X} \in \wp_n$ and flag variable drt $\in \{in, out\}$ (see line 2 of bfs_phase()).

We argue that each of these invocations takes at most $O(|\mathcal{F}| d_{S,T} n)$ time. Assume that N^{drt} is N^{in} when drt = in, and that it is N^{out} otherwise. Then, each invocation of next_step_bfs() takes $O(|\mathcal{X}| \max_{v \in \mathcal{X}} \{|N^{drt}(v)|\})$ time, because it involves visiting $N^{drt}(v)$ for each $v \in \mathcal{X}$; still, in order to enter the while-loop at line 1 of bfs_phase(), we must have $|\mathcal{X}| \leq |\mathcal{F}| d_{S,T}$, and moreover we have $|N^{drt}(v)| = O(n)$ for every $v \in \mathcal{X}$. Since the total number of iterations is bounded above by $d_{S,T}$, the bound follows.

Lemma 12. Assume that compression_phase() gets invoked at line 11 of the procedure solve_HY-STCON(). If compression_phase() halts without ever executing the procedure compute_Lehman-Ron_paths() at line 8, then it halts within

the following time bound:

$$O\left(\min\left(\sqrt{|\mathcal{F}|\,d_{S,T}\,n},|\mathcal{F}|\right)|\mathcal{F}|^2\,d_{S,T}^2\,n\right)\tag{1}$$

Otherwise, if compression_phase() executes compute_Lehman-Ron_paths() at line 8, then it halts within the following time bound:

$$O\left(\min\left(\sqrt{|\mathcal{F}|d_{S,T}n}, |\mathcal{F}|\right)|\mathcal{F}|^2 d_{S,T}^2 n + |\mathcal{F}|^{5/2} n^{3/2} d_{S,T}\right)$$
(2)

Proof. We start with some preliminary observations that will be useful in proving time bounds (1) and (2). Let us assume that compression_phase() is invoked on the following input $\langle S, T, F, \ell_{\uparrow}, \ell_{\downarrow}, d_{S,T}, n \rangle$ at line 11 of solve_HY-STCON(). We argue that the following bounds hold on the size of S and T:

$$|\mathcal{S}| \le |\mathcal{F}| \, d_{S,T} \, n \quad \text{and} \quad |\mathcal{T}| \le |\mathcal{F}| \, d_{S,T} \, n. \tag{3}$$

In fact, notice that S and T were computed during a previous invocation of double-bfs_phase(), at line 5 of solve_HY-STCON(). Therefore, it suffices to consider the set \mathcal{X} which is computed by passing through the while-loop at line 1 of bfs_phase(). The condition for entering that while-loop requires $|\mathcal{X}| \leq |\mathcal{F}| d_{S,T}$. Therefore, as soon as bfs_phase() exits that while-loop, we must have $|\mathcal{X}| \leq |\mathcal{F}| d_{S,T} n$. This implies the bounds specified by (3).

Consider the bipartite graph $\mathcal{G} = (V_{\mathcal{G}}, E_{\mathcal{G}}) = ((\mathcal{S}, \mathcal{T}), \subset)$, which is constructed at line 3 of compression_phase(). Since we have:

$$|V_{\mathcal{G}}| = |\mathcal{S}| + |\mathcal{T}| \le 2 |\mathcal{F}| \, d_{S,T} \, n,$$

we also have the following bound on the size of its edge set:

$$|E_{\mathcal{G}}| \le |V_{\mathcal{G}}|^2 \le 4 |\mathcal{F}|^2 d_{S,T}^2 n^2$$

We can now proceed with the proof of the two time bounds.

1. In the case where compute_Lehman-Ron_paths() never gets executed, recall that, at line 4, the compression_phase() computes a matching \mathcal{M} of \mathcal{G} such that $|\mathcal{M}| = \min(m^*, |\mathcal{F}|+1)$, where m^* is the size of a maximum cardinality matching of \mathcal{G} . At this point, the self-reduction($\mathcal{G}, |\mathcal{F}|+1$) (Algorithm 5), allows us to shrink the upper bound on the size of $|E_{\mathcal{G}}|$ from $|V_{\mathcal{G}}|^2$ down to:

$$|V_{\mathcal{G}}| \cdot |\mathcal{F}| \le 2 |\mathcal{F}|^2 d_{S,T} n.$$

The total overhead introduced by self-reduction() is only $O(|V_{\mathcal{G}}| + |E_{\mathcal{G}}|)$, because there are at most $|V_{\mathcal{G}}|$ recursive calls, each one inspecting the neighbourhood of some node of \mathcal{G} . So, \mathcal{M} is computed within the following time bound $t_{\mathcal{M}}$:

$$t_{\mathcal{M}} = O\left(\min(\sqrt{|V_{\mathcal{G}}|}, |\mathcal{F}|) |E_{\mathcal{G}}|\right)$$
$$= O\left(\min\left(\sqrt{|\mathcal{F}| d_{S,T} n}, |\mathcal{F}|\right) |\mathcal{F}|^2 d_{S,T} n\right)$$

At this point, let us observe that the time complexity of compute_min_vertex_cover(), which is invoked at line 11 of compression_phase(), is bounded above by the time complexity of computing \mathcal{M} at line 4. Also, by Lemma 11, the time complexity of the double-bfs_phase(), which is invoked at line 14 of compression_phase(), is bounded above by the same quantity.

If compute_Lehman-Ron_paths() never gets executed at line 8, then during each iteration of the while-loop at line 2 of compression_phase(), the most expensive task is that of computing the matching \mathcal{M} at line 4. Recall that, according to Lemma 3, the while-loop at line 2 iterates at most $d_{S,T}$ times. We conclude that, in this case, the compression_phase() halts within the following time bound:

$$t_{\mathcal{M}} d_{S,T} = O\left(\min\left(\sqrt{|\mathcal{F}| d_{S,T} n}, |\mathcal{F}|\right) |\mathcal{F}|^2 d_{S,T}^2 n\right).$$

2. In the case where compute_Lehman-Ron_paths() gets executed, which happens whenever $|\mathcal{M}| = |\mathcal{F}| + 1$, we must now take its time complexity into account, which we analyze below.

First, consider the set $\mathcal{M}_{\mathcal{S}}$ computed at line 6 of compression_phase(). The following bound holds on its size:

$$|\mathcal{M}_{\mathcal{S}}| = |\mathcal{M}| = |\mathcal{F}| + 1.$$

The same bound holds for the set $\mathcal{M}_{\mathcal{T}} \subseteq \mathcal{T}$ which is computed at line 7 namely: $|\mathcal{M}_{\mathcal{T}}| = |\mathcal{M}| = |\mathcal{F}| + 1$. By Theorem 3, provided that we consider the parameter $m = |\mathcal{M}| = O(|\mathcal{F}|)$, invoking compute_Lehman-Ron_paths() on input $\langle \mathcal{M}_{\mathcal{S}}, \mathcal{M}_{\mathcal{T}}, \mathcal{M}, n \rangle$ takes time at most t_{LR} , where:

$$t_{\rm LR} = O\left(m^{5/2}n^{3/2}d_{S,T}\right) = O\left(|\mathcal{F}|^{5/2}n^{3/2}d_{S,T}\right).$$

Recall that, by Lemma 3, the while-loop at line 2 iterates at most $d_{S,T}$ times. At each of such iterations, a brand new matching \mathcal{M} gets computed at line 4. Finally, at the very last of such iterations, provided that $|\mathcal{M}| > |\mathcal{F}|$, then the procedure compute_Lehman-Ron_paths() is invoked at line 8. Therefore, we conclude that whenever compression_phase() executes compute_Lehman-Ron_paths() at line 8, then it halts within the following time bound:

$$t_{\mathcal{M}} d_{S,T} + t_{\mathrm{LR}} = O\left(\min\left(\sqrt{|\mathcal{F}| d_{S,T} n}, |\mathcal{F}|\right) |\mathcal{F}|^2 d_{S,T}^2 n + |\mathcal{F}|^{5/2} n^{3/2} d_{S,T}\right)$$

Proposition 8. The DECISION-TASK of HY-STCON can be solved within the following time bound on any input (S, T, \mathcal{F}, n) :

$$O\left(\min\left(\sqrt{|\mathcal{F}|\,d_{S,T}\,n},|\mathcal{F}|\right)|\mathcal{F}|^2\,d_{S,T}^3\,n\right).$$

Proof. Let us consider the procedure $solve_HY$ -STCON() of Algorithm 2. By Proposition 5, the while-loop at line 4 iterates at most $d_{S,T}$ times. At each iteration, double-bfs_phase() is invoked at line 5, and compression_phase() is invoked soon after at line 11. By Lemma 11, the most expensive one between the two procedures is clearly compression_phase(). Recall that, if we are content with solving the DECISION-TASK of HY-STCON, then the compression_phase() can be implemented so that it always halts without ever executing the procedure compute_Lehman-Ron_paths() at line 8. Therefore, by Lemma 12, each invocation of compression_phase() takes time at most

$$O\left(\min\left(\sqrt{|\mathcal{F}|\,d_{S,T}\,n},|\mathcal{F}|\right)|\mathcal{F}|^2\,d_{S,T}^2\,n\right)$$

Since we have at most $d_{S,T}$ of such invocations, then the thesis follows.

Proposition 9. The SEARCH-TASK of HY-STCON can be solved within the following time bound on any input (S, T, \mathcal{F}, n) :

$$O\left(\min\left(\sqrt{|\mathcal{F}|d_{S,T}n}, |\mathcal{F}|\right)|\mathcal{F}|^2 d_{S,T}^3 n + |\mathcal{F}|^{5/2} n^{3/2} d_{S,T}\right)$$

Proof. Let us consider the procedure solve_HY-STCON() of Algorithm 2. By Proposition 5, the while-loop at line 4 iterates at most $d_{S,T}$ times. At each iteration, double-bfs_phase() is invoked at line 5, and compression_phase() is invoked shortly after at line 11. By Lemma 11, the most expensive step between the two is clearly the compression_phase(). Recall that, if we aim to solve the SEARCH-TASK of HY-STCON, then the compression_phase() possibly executes the compute_Lehman-Ron_paths() procedure at line 8. Nevertheless, whenever compression_phase() executes compute_Lehman-Ron_paths() at line 8, then the procedure **solve_Hy-stCon()** halts shortly after at line 12. This means that the only invocation of compression_phase() that possibly executes compute_Lehman-Ron_paths() is the very last invocation. Then, each invocation of compression_phase(), except the very last one, halts within the following time bound by Lemma 12: $O(\min(\sqrt{|\mathcal{F}| d_{S,T} n}, |\mathcal{F}|) |\mathcal{F}|^2 d_{S,T}^2 n)$. Since the very last invocation of compression_phase() possibly executes the procedure compute_Lehman-Ron_paths() at line 8, the following time bound holds on the last invocation of compression_phase() by Lemma 12:

$$O\left(\min\left(\sqrt{|\mathcal{F}|d_{S,T}n}, |\mathcal{F}|\right)|\mathcal{F}|^2 d_{S,T}^2 n + |\mathcal{F}|^{5/2} n^{3/2} d_{S,T}\right)$$

Since there are at most $d_{S,T}$ invocations of the compression_phase(), the thesis follows.

4 Conclusion

With the intention of integrating more biologically relevant constraints into classical genome rearrangement problems, we introduced in this paper the GUIDED SORTING problem. We broadly define it as the problem of transforming two genomes into one another using as few operations as possible from a given fixed set of allowed operations while avoiding a set of nonviable genomes. We gave a polynomial time algorithm for solving this problem in the case where genomes are represented by permutations, under the assumptions that 1) permutations can only be modified by exchanging any two elements, 2) the sequence to seek must be optimal, and 3) the permutation to sort is an involution.

Many questions remain open, most notably that of the computational complexity of the GUIDED SORTING problem, whether under assumptions (1) and (2) or in a more general setting (i.e., using structures other than permutations, operations other than exchanges, or allowing sequences to be "as short as possible" instead of optimal). One could also investigate "implicit" representations for the set of forbidden intermediate permutations, e.g. all permutations that avoid a given (set of) pattern(s), or that belong to a specific conjugacy class. Aside from complexity issues, future work shall also focus on extending the approach we proposed to other families of instances of the GUIDED SORTING problem, and identifying other tractable (or intractable) cases or variants of it; for instance, we plan to extend our algorithmic results to the family of graphs satisfying the *shadow-matching* [18] condition.

Bibliography

- S. BÉRARD, A. BERGERON, AND C. CHAUVE, Conservation of combinatorial structures in evolution scenarios, in RECOMB'04, vol. 3388 of LNCS, Springer, 2004, pp. 16–19.
- [2] A. BERGERON, M. BLANCHETTE, A. CHATEAU, AND C. CHAUVE, Reconstructing ancestral gene orders using conserved intervals, in WABI'04, vol. 3240 of LNCS, Springer, 2004, pp. 14–25.
- [3] C. COMIN, A. LABARRE, R. RIZZI, AND S. VIALETTE, Sorting with Forbidden Intermediates, Algorithms for Computational Biology: Third International Conference, AlCoB 2016, Trujillo, Spain, June 21-22, 2016, Proceedings, Springer, 2016, pp. 133–144.
- [4] L. CUÉNOT, Les races pures et leurs combinaisons chez les souris, Archives de Zoologie Experimentale, 3 (1905), pp. cxxiii–cxxxii.
- [5] M. DEZA AND T. HUANG, Metrics on permutations, a survey, Journal of Combinatorics, Information and System Sciences, 23 (1998), pp. 173–185.
- [6] R. DIESTEL, Graph Theory (Graduate Texts in Mathematics), Springer, 2005.
- [7] G. FERTIN, A. LABARRE, I. RUSU, E. TANNIER, AND S. VIALETTE, Combinatorics of Genome Rearrangements, The MIT Press, 2009.
- [8] M. FIGEAC AND J.-S. VARRÉ, Sorting by reversals with common intervals, in WABI'04, vol. 3240 of LNCS, Springer, 2004, pp. 26–37.
- [9] H. GABOW, S. MAHESHWARI, AND L. OSTERWEIL, On two problems in the generation of program test paths, IEEE Trans. Software Eng., (1976), pp. 227–231.
- [10] S. GLUECKSOHN-WAELSCH, Lethal genes and analysis of differentiation, Science, 142 (1963), pp. 1269–1276.
- [11] O. GOLDREICH, S. GOLDWASSER, E. LEHMAN, D. RON, AND A. SAMORODNITSKY, *Testing monotonicity*, Combinatorica, 20 (2000), pp. 301–337.
- [12] J. HOPCROFT AND R. KARP, An n^{5/2} algorithm for maximum matchings in bipartite graphs, SIAM Journal on Computing, 2 (1973), pp. 225–231.
- [13] D. E. KNUTH, The Art of Computer Programming, Volume III: Sorting and Searching, Addison-Wesley, 1973.
- [14] K. KRAUSE, R. SMITH, AND M. GOODWIN, Optimal software test planning through automated network analysis, in Proceedings 1973 IEEE Symposium Computer Software Reliability, IEEE, 1973, pp. 18–22.
- [15] A. LABARRE, Lower bounding edit distances between permutations, SIAM Journal on Discrete Mathematics, 27 (2013), pp. 1410–1428.
- [16] S. LAKSHMIVARAHAN, J.-S. JWO, AND S. K. DHALL, Symmetry in interconnection networks based on Cayley graphs of permutation groups: A survey, Parallel Computing, 19 (1993), pp. 361–407.
- [17] E. LEHMAN AND D. RON, On disjoint chains of subsets, Journal of Combinatorial Theory, Series A, 94 (2001), pp. 399–404.

- [18] M. LOGAN AND S. SHAHRIARI, A new matching property for posets and existence of disjoint chains, Journal of Combinatorial Theory, Series A, 108 (2004).
- [19] H. YINNONE, On paths avoiding forbidden pairs of vertices in a graph, Discrete Applied Mathematics, 74 (1997), pp. 85–92.