

Fast and Flexible Flow Decompositions in General Graphs via Dominators

Francisco Sena¹ and Alexandru I. Tomescu¹

¹Department of Computer Science, University of Helsinki, Finland,
`{alexandru.tomescu, francisco.sena}@helsinki.fi`

Abstract

Multi-assembly methods rely at their core on a flow decomposition problem, namely, decomposing a weighted graph into weighted paths or walks. However, most results over the past decade have focused on decompositions over directed acyclic graphs (DAGs). This limitation has lead to either purely heuristic methods, or in applications transforming a graph with cycles into a DAG via preprocessing heuristics. In this paper we show that flow decomposition problems can be solved in practice also on general graphs with cycles, via a framework that yields fast and flexible Mixed Integer Linear Programming (MILP) formulations.

Our key technique relies on the graph-theoretic notion of *dominator tree*, which we use to find all *safe sequences of edges*, that are guaranteed to appear in some walk of any flow decomposition solution. We generalize previous results from DAGs to cyclic graphs, by showing that maximal safe sequences correspond to extensions of common leaves of two dominator trees, and that we can find all of them in time linear in their size.

Using these, we can accelerate MILPs for *any* flow decomposition into walks in general graphs, by setting to (at least) 1 suitable variables encoding solution walks, and by setting to 0 other walks variables non-reachable to and from safe sequences. This reduces model size and eliminates costly linearizations of MILP variable products.

We experiment with three decomposition models (Minimum Flow Decomposition, Least Absolute Errors and Minimum Path Error), on four bacterial datasets. Our pre-processing enables up to thousand-fold speedups and solves even under 30 seconds many instances otherwise timing out. We thus hope that our dominator-based MILP simplification framework, and the accompanying software library can become building blocks in multi-assembly applications.

Keywords: directed graph, walk cover, dominator tree, graph algorithm, graph reachability, integer linear programming, flow decomposition, optimization problem, multi-assembly

Co-funded by the European Union (ERC, SCALEBIO, 101169716). Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council. Neither the European Union nor the granting authority can be held responsible for them. This work has also received funding from the Research Council of Finland grants No. 346968, 358744.



1 Introduction

Background. Decomposing a weighted graph into weighed paths is a basic computational problem at the core of several types of *multi-assembly problems*, for example RNA transcript assembly [47, 24, 23, 55, 36, 53, 37, 43, 12, 10, 8, 54, 52, 34], metagenomic assembly [29, 31, 22, 30, 21, 18, 51, 38] or the assembly of viral strains [49, 3, 4, 5, 26, 46, 13, 25]. In this kind of problems, multiple (similar) genomic sequences, in different abundances in a sample, need to be recovered from mixed reads sequenced from all of them.

In popular models of this problem, the reads give rise to a weighted graph, whose vertices encode partial genomic sequences, and whose edges encode observed connections between them (e.g. reads spanning two such genomic sequences). In addition, vertices or edges (or both) have weights corresponding to their read coverage. The genomic sequences to be recovered correspond to *weighted paths* in this graph, the weight being their abundance in the sample. Therefore, to solve the initial multi-assembly problem, one needs to “optimally” decompose the weighted graph into weighted paths.

For directed acyclic graphs (DAGs), there is a wide body of work on this problem. For example, various formulations of optimality for a set of weighted paths exist [47, 24, 23, 55, 36, 53, 37, 43, 12, 10, 8, 9, 44]. Among these, the most basic is the Minimum Flow Decomposition (MFD) problem for error-free data [33, 41, 20], where the input is a DAG with a weight on every edge (its *flow value*), satisfying *flow conservation* at every vertex except sources or sinks (i.e. the sum of flow values on incoming edges equals that on outgoing edges). The MFD problem asks for a minimum number of weighted paths such that the flow value of every edge equals the sum of weights of paths passing through it [50]. MFD is NP-hard [50], even on simple DAG classes [17, 16], and hence most multi-assembly formulations for real, imperfect data are NP-hard [47].

Motivation. The acyclicity of the graph is a powerful structural property that has been exploited in numerous practical solution for MFD, despite its hardness. For example, on DAGs, MFD admits a fixed-parameter tractable algorithm [20], tailored heuristics [41, 28, 50], approximation algorithms on some classes of DAGs [7], or integer linear programs (ILPs) [10], which further admit optimizations to significantly reduce the runtime of ILP solvers [15].

The situation is radically different for general graphs that may contain cycles, where there are far fewer methods to approach MFD: two heuristics appear in [50], an approximation algorithm working on a certain class of graphs was proposed by [7], and an exact ILP-based method was proposed by [11]. Moreover, all these work only for perfect data, and thus cannot be applied to real data. In fact, some multi-assembly methods employ various strategies or construction parameters to make a graph built from real data acyclic [32, 4].

Moreover, we have shown that on DAGs, *any* flow decomposition problem for erroneous data admits fast ILP-based solvers [39], not only MFD. This methodology works by (1) abstracting the decomposition problem as one whose solution consists of a set of source-to-sink paths that cover (a given subset of) the edges (they form a *walk cover* of the DAG); (2) characterizing the sequences of edges that must appear in a path of *every* such walk cover (*safe sequences*); (3) developing a linear-time algorithm finding all maximal safe sequences in a DAG, based on the well-known notion of *dominator* (see e.g. [14]); and (4) reducing the search space of the ILP by fixing to 1 suitable binary variables corresponding to incompatible safe sequences.

Contributions. In this paper we show that this methodology is not limited to DAGs, but can be adapted also to general graphs that may contain cycles. For (1), we consider decomposition problems whose solution consists of a set of source-to-sink *walks* that cover (a given subset of) the edges (a *walk cover*—walks may repeat vertices and edges). For (2) and (3) we show that the techniques developed for the DAG case can be extended to the general case of graphs with cycles; we also obtain a linear-time algorithm outputting all maximal safe sequences of walk covers, also based on dominators. For (4) we use the condensation DAG of strongly connected components to find maximal safe sequences that can be guaranteed to appear in different solution walks. Based on them, we use again the condensation DAG to generalize the fixing of walk variables of the ILP. As a novelty, in this paper we also fix to 0 some integer variables corresponding to edges that cannot be traversed by a solution path. Additionally, by this fixing to 0 or 1, we can also avoid the expensive linearization of products between two integer variables, leading to smaller models that are faster to solve.

As a technical novelty compared to [39], note that the characterization of maximal safe sequences from [39] (in DAGs) is based on first compressing paths in the input graph where every vertex but the first has in-degree one and every vertex but the last one has out-degree one. This leads to involved characterizations and analysis for the case when the solution paths are required to cover only a subset of vertices/edges. Instead, we propose to collapse particular paths on the dominator trees, without affecting the dominance relation of G . This turns out to give a unified view on standard walk covers and on walk covers where a subset of vertices/edges has to be covered, for general directed graphs. Moreover, the techniques used in [39] for DAGs can thus be seen as a special application of the techniques developed here.

As base ILP formulation of source-to-sink walks in general graphs, we use the one from [11], which was proposed for the MFD problem. Here we slightly simplify it, and we extend it also to decomposition problems applicable to imperfect data. As two concrete ILP models for such data, in this paper we consider the LeastAbsErrors model [3, 48, 6] (where we need to minimize the total sum of the absolute differences between the weight of every edge and the sum of the weights of the walks using it), and a more complex and more accurate “robust” model, MinPathError [9] (where paths also get a slack variable, whose total sum needs to be minimized, and this absolute difference for every edge must be below the sum of the slacks of the walks passing through it). We define these models formally in Section 2.

Moreover, in applications of flow decomposition problems, more information is often available than just the graph structure and the edge weights. For example, long reads overlapping more than one edge indicate that these edges must co-appear in some solution path or walk [5, 54]. For acyclic graphs, these have been modeled previously as *subpath constraints* [52, 35, 15, 10], namely given paths that must appear as a subpath of at least one solution path of the decomposition problem. For general graphs with cycles, we introduce here a weaker form of this constraint (but which can be directly added to the MILP models above), in the form of a *subset constraint*, namely a set S of edges that must appear in at least one walk of any solution.

We implemented all our ILP models and their optimizations inside the `flowpaths` Python package for flow decomposition problems (<https://algbio.github.io/flowpaths>). This package uses by default the popular open source MILP solver HiGHS, see www.hiighs.dev, which is installable without a license, thus maximizing the potential for applications in bioinformatics software. The architecture of the `flowpaths` package is designed so that the formulation of paths is inside an abstract class, which is then instantiated by any flow decomposition model adding its own problem-specific constraints. We followed the same approach for walks: the optimizations proposed here and the subset constraints are all implemented in an abstract class, and thus applicable to any other decomposition model into walks that might be developed in the future based on `flowpaths`. Given the generality of this software design, we focused our experiments on measuring the speed-ups that our optimizations can give on input graphs created from different bacterial datasets. We observe speed-ups that grow with the model complexity (from MinFlowDecomp, to LeastAbsErrors, to MinPathError), and likewise grow with the number of bacterial genomes that are in the graph. On the hardest datasets that we tested we obtain speed-ups of $393\times$ for MinFlowDecomp, $560\times$ for LeastAbsoluteErrors, and $1465\times$ for MinPathError, and solving under 30 seconds on average many instances otherwise timing out.

2 Notation and preliminaries

Graphs. An s - t graph is a directed graph without parallel edges (but possibly with self-loops), with a unique source s (a vertex with no in-coming edges) and a unique sink t (a vertex with no out-going edges). Let $G = (V, E)$ be an s - t graph. A *walk* in G is a sequence of vertices $v_1 v_2 \dots v_t$ such that $v_i v_{i+1} \in E$ for all $1 \leq i < t$; note that a walk can revisit vertices and edges. For designated vertices $u, v \in V$, we say that a walk W in G is an u - v walk if the first vertex is u and the last vertex is v . If the walk does not repeat vertices (and thus also edges), we say that it is a *path*. For an edge $uv \in E$, we denote by $W(uv)$ the number of times the edge uv is traversed by the walk W ; if W does not traverse uv then we set $W(uv) = 0$. The *distance* between two vertices is the number of edges in a shortest path between them.

Sequences and safety. Let C be a subset of vertices/edges. A C -walk cover of G is a set P of s - t walks such that for every vertex/edge $u \in C$ there is a walk in P containing u . We denote a *sequence* X

through vertices/edges u_1, \dots, u_ℓ as $X = u_1, \dots, u_\ell$ or $X = (u_1, \dots, u_\ell)$ if there is a $u_i - u_{i+1}$ path for all $1 \leq i \leq \ell - 1$. Any path is a sequence. We say that X *contains* a vertex v if $v = u_i$ for some $i \in \{1, \dots, \ell\}$. The internal vertices of a path are the vertices contained in the path except its first and last vertices. Let $X' = v_1, v_2, \dots, v_{\ell'}$ be a sequence. If X can be obtained from X' by deleting any number of vertices/edges then X is a *subsequence* of X' (and X' is a *supersequence* of X). Suppose that there is a path from the last vertex u_ℓ of X to the first vertex v_1 of X' . The *concatenation* of X and X' is the sequence XX' obtained as X followed by X' ; if $u_\ell = v_1$, then the concatenation is $XX' := u_1, \dots, u_\ell, v_2, \dots, v_{\ell'}$, i.e. the repeated occurrence of $u_\ell = v_1$ is removed. For sequences of edges, the definitions given above for vertices adapt in the obvious way.

Let $G = (V, E)$ be an s - t graph, let C be a subset of V (resp. E) and let S be a sequence of vertices (resp. edges) of G . We say that X is a *C-safe sequence* if for any C -walk cover P , there is a walk W in P such that X is a subsequence of W (when clear from the context, we omit the prefix C -). See Figure 1a for an illustration of four safe sequences for $C = V$.

Dominators. We say that vertex u *s-dominates* vertex v if every s - v path contains u . Analogously, vertex u *t-dominates* vertex v if every v - t path contains u . Every vertex s - and t -dominates itself by definition. Vertex u *strictly s-dominates* (*strictly t-dominates*) v if u *s-dominates* (*t-dominates*) v and $u \neq v$. For $v \neq s$, the *immediate s-dominator* (*immediate t-dominator*) of v is the vertex $u \neq v$ that *s-dominates* (*t-dominates*) v and is dominated by all the vertices that strictly dominate v ; we write $idom_s(v) = u$ ($idom_t(v) = u$). The domination relation on G with respect to s can be represented as a tree rooted at s with the same vertex set as G , called the *s-dominator tree* of G . In this tree, every vertex has as its parent its immediate s -dominator, and every vertex is s -dominated by all its ancestors and s -dominates all its descendants and itself. A vertex is said to be an *s-dominator* if it strictly s -dominates some vertex. The domination relation on G with respect to t is defined analogously. Sequences on dominator trees are to be understood as subsequences of vertex-to-root paths of the tree, so that sequences in the s - and t -dominator trees are in correspondence with sequences of s - and t -dominators. See Figure 1 for an illustration.

We also define a function $dom : V \times \mathbb{N}^+ \rightarrow V$ on dominator trees as follows. Let $dom_s(v, k) := idom_s(v)$ when $k = 1$ and $dom_s(v, k) := dom_s(idom_s(v, 1), k - 1)$ when $k > 1$. If k is larger than the number of *strict* dominators of a given vertex, then it defaults to s (resp. t) in the s -dominator tree (resp. t -dominator tree). Essentially, dom gives the k -th ancestor of a vertex in a rooted tree, if it is well defined, otherwise it returns the root of the tree. By $extension(v)$ we mean the sequence of vertices obtained by concatenating the path from s to v in the s -dominator tree with the path from v to t in the t -dominator tree. The *depth* of a vertex v in a rooted tree is the distance between v and the root of the tree. See [14] for a more in-depth presentation of dominators.

Flow decomposition models. Next, we formally define the three flow decomposition problem variants that we address in the paper.

Definition 1 (Flow decomposition variants). Let $G = (V, E)$ be an s - t graph, let $E' \subseteq E$ be a subset of edges, let $f : E \rightarrow \mathbb{Z}_+ \cup \{0\}$ be a function assigning a non-negative weight to every edge of G , and let $k \geq 1$ be an integer. We define below three problems that require finding k s - t walks W_1, \dots, W_k in G , with associated weights $w_1, \dots, w_k \in \mathbb{Z}_+$, respectively, with the following objectives:

k -Flow Decomposition (k -FD) problem For every edge $uv \in E'$, we have: $\sum_{i=1}^k W_i(uv)w_i = f(uv)$.

k -Least Absolute Errors (k -LAE) problem The walks and their associated weights minimize:

$$\sum_{uv \in E'} \left| f(uv) - \sum_{i=1}^k W_i(uv)w_i \right|.$$

k -Min Path Error (k -MPE) problem Find also an associated slack $\rho_i \in \mathbb{Z}_+ \cup \{0\}$ for each walk W_i ,

$$\text{such that: } \left| f(uv) - \sum_{i=1}^k W_i(uv)w_i \right| \leq \sum_{i=1}^k W_i(uv)\rho_i \text{ holds for all } uv \in E', \text{ and minimizing } \sum_{i=1}^k \rho_i.$$

Note that in the above we assumed that all walks start in the same vertex s and end in the same vertex t . However, in practice we might have a *set* S of vertices where the walks are allowed to start, especially if the graph has more sources, and likewise for a *set* T where the walks are allowed to end. However, allowing for such S - T walks is easy, because we considered the generalization in which the objectives are applied only to the subset of edges E' . We can then simply add a new vertex s to such an input, together with edges from s to every vertex in S , and a new vertex t together with edges from every vertex in T to t . These new edges incident to s or t are not added to E' , hence they can get weight 0 and they do not interfere with the problem objectives (which are applied to E').

Due to lack of space, we describe the MILP formulations for the problems from Definition 1 in Section A. For the rest of the paper, it suffices to know that every walk W_i , with $i \in \{1, \dots, k\}$ is modeled by variables $x_{uv,i} \in \mathbb{Z}_+$ for every edge $uv \in E$, such that $x_{uv,i} = W_i(uv)$.

3 Safe sequences and dominators

Note that the flow decomposition models from Definition 1 are defined so that their solution walks cover the edges of the graph. As such, in our applications, we are interested in safe sequences for $C \subseteq E$. However, dominator trees are usually defined in terms of vertices. Thus, we state the results in terms of $C \subseteq V$, as they may also be of independent interest outside of our bioinformatics application. We first solve the case $C = V$, and then show how to generalize for $C \subseteq V$ (the former results serve as preliminary steps toward understanding the more general theorems). Further, we do not give proofs of the statements $C = V$ as they follow as a particular case of the results for $C \subseteq V$. The case $C \subseteq E$ can be handled in two manners: either (1) introduce a new vertex v_e in the middle of every edge $e \in C$, add all such v_e to C' , and compute all C' -safe sequences; or (2) obtain direct analogues of the characterizations and enumeration algorithms in terms of edges using the line graph of G and using dominator trees of *edges* (for a running example of this approach, observe that Figure 1a is the line graph of Figure 3a and compare the safe sequences shown with matching colors, the former figure w.r.t. vertices and the latter w.r.t. edges). The relevant missing proofs from this section can be found in Section B.

Theorem 1. *Let $G = (V, E)$ be an s - t graph. A sequence X of vertices is safe for walk covers if and only if there exists a vertex $v \in V$ such that X is a subsequence of $\text{extension}(v)$.*

As every safe sequence is a subsequence of the extension of some vertex v , every maximal safe sequence corresponds to the extension of some vertex, i.e., the corresponding paths in the dominator trees do not skip any vertices on their way to the roots. This implies that in a graph with n vertices there are at most n maximal safe sequences. Moreover, it also implies that any vertex u appears at most twice in any safe sequence X of G , since any vertex appears at most once in the sequence of the s -dominators of u , and likewise for t -dominators.

The next lemma was first showed in [39], and essentially relates the s - and t -dominance relations so as to generalize the immediate dominance relation. The result is merely technical but is required for the characterization of maximal safe sequences.

Lemma 1. *Let $G = (V, E)$ be an s - t graph, let $u, v \in V$ be vertices, and let $k \in \mathbb{N}^+$. If u is the k -th ancestor of v in the s -dominator tree, then $\text{dom}_t(u, k)$ t -dominates v , and v is not a t -dominator of u unless $v = \text{dom}_t(u, k)$.*

In order to obtain a characterization of maximal safe sequences on DAGs, [39] compresses paths in the input graph where every vertex but the first has in-degree one and every vertex but the last one has out-degree one. (In bioinformatics, this kind of paths are known as unitigs, and [39] refers to them as unitary paths.) This causes the characterization of maximal safe sequences with respect to C -walk covers to be more involved than necessary. Instead, we propose to collapse particular paths on the dominator trees. As we will show, this operation preserves the set of maximal safe sequences of G . Moreover, the paths compressed in [39] always correspond to these particular paths in the trees. However, the converse does not hold due to the presence of cycles in our graph, and hence our approach indeed generalizes that of [39].

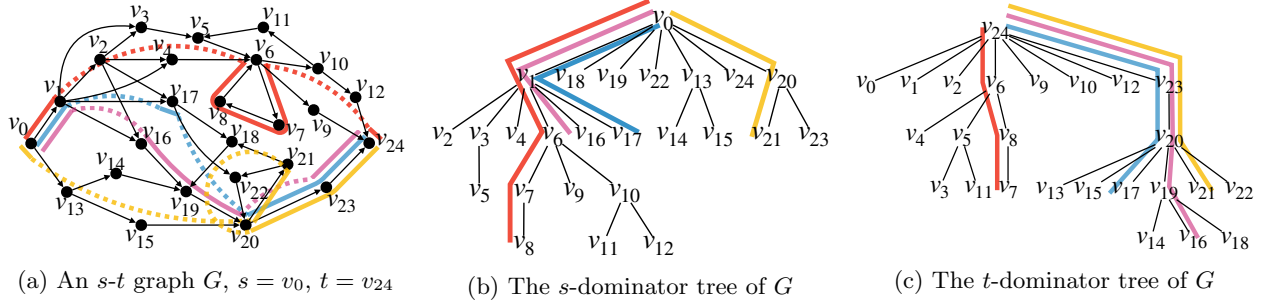


Figure 1: Example of a graph G , its s - and t -dominator trees, four C -safe sequences in G , for $C = V$. For illustration purposes, in a sequence we draw univocal edges between nodes as solid lines, and other connections as dashed: for example, the yellow safe sequence is $(v_0, v_{20}, v_{21}, v_{20}, v_{23}, v_{24})$. Because of Theorem 2, maximal safe sequences are obtained by concatenating the paths in the s - and t -dominator trees: for example, the yellow sequence is obtained by concatenating the path $v_0 v_{20} v_{21}$ in the s -dominator tree with the path $v_{21} v_{20} v_{23} v_{24}$ in the t -dominator tree.

Definition 2 (Univocal path). Let G be an s - t graph. Let $p = v_1 \dots v_k$ be a path in the t -dominator tree such that v_{i+1} has a unique child v_i for $i = 1, \dots, k-1$ ($k \geq 1$). If $v_k \dots v_1$ is a path in the s -dominator tree where each vertex but the deepest one has exactly one child, then p is called a *univocal path* of G .

Notice that univocal paths are syntactically described with respect to the order in which their vertices appear in the t -dominator tree. Trivially, a single vertex forms a univocal path.

Let $p = u_1 \dots u_k$ be a nontrivial maximal univocal path ($k \geq 2$) and let T_s and T_t denote the s -dominator and t -dominator trees, respectively. By definition of univocal path, the extensions of the vertices of p all yield the same sequence. Importantly, the extension of any vertex v not in p either contains p or does not contain any vertex in p : if v is a descendant of v_1 in T_t or of v_k in T_s then $\text{extension}(v)$ clearly contains p ; otherwise, v is not a descendant of v_1 in T_t and is not a descendant of v_k in T_s , and thus the lowest common ancestor of v and v_1 (resp. v_k) in T_t (resp. T_s) is a strict ancestor of v_k (resp. v_1) in T_t (resp. T_s), and thus $\text{extension}(v)$ does not contain any vertex in p . We claim that there is a unique set of maximal safe sequences. For this, it is enough to show that two distinct maximal univocal paths are vertex disjoint.

Lemma 2. Let G be an s - t graph and let p_1, p_2 be two distinct maximal univocal paths. Then p_1 and p_2 are vertex disjoint.

Collapsing each maximal univocal path into a single vertex then yields two trees with essentially the same set of extensions of the trees before the collapse; we choose the deepest vertex in the t -dominator tree to represent the respective collapsed path. To ensure that there is no loss of information during the process, every vertex arising from a collapse operation should store additionally the ordered sequence of vertices of the respective collapsed path (and it is not hard to see that it suffices to do this in just one of the trees). In this way, whenever an extension passes through a collapsed vertex, the sequence stored at that vertex is considered for the extension, thus mimicking the extension in the original trees. Therefore the set of maximal safe sequences before and after the collapse is identical.

For example, in Figure 2, the paths uv and bc are maximal univocal paths (and in fact these paths do not correspond to unitary paths in the sense of [39]). Thus, these should be collapsed. On the other hand, the path df should not be collapsed even though $\text{extension}(d) = \text{extension}(f)$, since otherwise this would mistakenly encode the fact that vertex f s -dominates vertex e (further, while this collapse preserves the t -dominance relation, it does not preserve the s -dominance relation). After collapse, the path uv in the t -dominator tree becomes vertex u with the additional information “ v ”, and in the s -dominator tree the path vu becomes simply vertex u . We are now ready to give a characterization of maximal safe sequences.

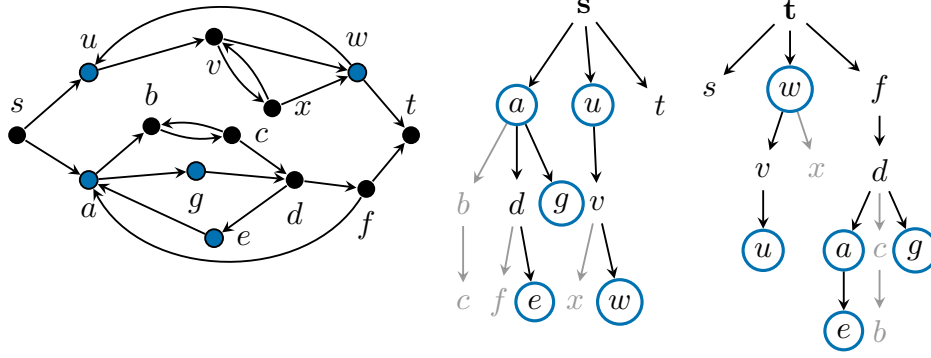


Figure 2: Example of a graph G , its blue-dominator trees with respect to $C = \{a, e, g, u, w\}$. Lower opacity vertices are not part of the blue-dominator trees. Underlying the blue-dominator trees are the dominator trees of G for $C = V$. Vertex e is a blue-child of vertex a in the s -dominator tree (and vertex a is its closest blue-ancestor). The path uw is maximal C -univocal and its collapsed into vertex u , which stores the sequence (v, w) . The path ea is not C -univocal because of the position of vertex g relative to $\{a, e\}$ in the s -dominator tree. The maximal C -safe sequences are: $\text{extension}(u) = \text{extension}(w) = (s, u, v, w, t)$, $\text{extension}(e) = (s, a, d, e, a, d, f, t)$, and $\text{extension}(g) = (s, a, g, d, f, t)$.

Theorem 2 (Characterization of maximal safe sequences). *Let $G = (V, E)$ be an s - t graph and suppose that the dominator trees of G are collapsed. Let u be a vertex in the collapsed trees. Then $\text{extension}(u)$ is a maximal safe sequence with respect to G if and only if u is a leaf in both dominator trees.*

Linear-time algorithms to build dominator trees are known, see e.g. [2], so we get the next result.

Theorem 3 (Optimal enumeration and representation of safe sequences). *Let G be an s - t graph with n vertices and m edges. There is an $O(m + o)$ time algorithm outputting the set of all maximal safe sequences with no duplicates, where o denotes the total length of all the maximal safe sequences.*

Moreover, Theorem 2 also tells us that the dominator trees encode every maximal safe sequence of G via extensions. Hence, dominator trees represent all the maximal safe sequences of G while requiring $O(n)$ -space. We dedicate the rest of this section in generalizing the previous results for walk covers where only a given subset of vertices has to be covered.

C -safe sequences. The characterization of C -safe sequences is completely analogous to Theorem 1.

Theorem 4 (Characterization of safe sequences for C -walk covers). *Let $G = (V, E)$ be an s - t graph, $C \subseteq V$ a set of vertices. A sequence X of vertices is safe for C -walk covers if and only if there exists a vertex $v \in C$ such that X is a subsequence of $\text{extension}(v)$.*

Our goal is to establish analogues of Theorem 2 and Theorem 3 for C -walk covers. We begin by describing univocal paths with respect to C -walk covers (or C -univocal paths). Mark with blue the vertices of C in the both dominator trees of G .

A vertex is *blue* if it is marked. A *blue-child* of a vertex u in the s -dominator tree is a blue vertex that is a descendant of u and no internal vertex on the (unique) v - u path in the s -dominator tree is blue (and analogously for the t -dominator tree). By Theorem 4, vertices that do not appear in the extension of any blue vertex can be ignored/removed from the trees. This produces the *blue-dominator trees*, which essentially encode the dominance relation of G restricted to C while maintaining relevant non blue vertices (i.e., vertices not in C which appear in the extension of a blue vertex). See Figure 2.

A C -univocal path in the blue-dominator trees generalizes in the natural way: if $v_1 \dots v_k$ is a sequence of blue vertices in the blue t -dominator tree, then v_1, \dots, v_k is C -univocal if each vertex v_{i+1} has v_i as its *unique* blue-child for $i = 1, \dots, k - 1$ ($k \geq 1$) (i.e., in the blue t -dominator tree, no internal vertex on

the v_i-v_{i+1} path is blue and each vertex on the this path except v_i has a unique child) and v_k, \dots, v_1 is a sequence in the s -dominator tree with the analogous properties. So Lemma 2 generalizes for C -univocal paths, and further, since Lemma 1 suitably explains how the s - and t -dominance relation interact in arbitrary “ancestry-dominance” relations, the characterization of maximal C -safe sequences is a direct generalization of Theorem 2. We remark that in order to prove Theorem 2, application of the Lemma 1 is required only for $k = 1$ (i.e., immediate-dominance relation).

Theorem 5 (Characterization of maximal safe sequences for C -walk covers). *Let $G = (V, E)$ be an s - t graph, let $C \subseteq V$, and suppose that the dominator trees of G are collapsed with respect to C . Let $u \in C$ be a vertex in the collapsed trees. Then $\text{extension}(u)$ is a maximal C -safe sequence if and only if u is a leaf in both dominator trees.*

To compute maximal C -safe sequences we can proceed essentially as when $C = V$. The description of this procedure is given in the proof of the next theorem. Moreover, this result also implies that blue-dominator trees represent every maximal C -safe sequence of G in $O(n)$ space.

Theorem 6 (Maximal safe sequence enumeration for C -walk covers). *Let $G = (V, E)$ be an s - t graph with n vertices and m arcs, and let $C \subseteq V$. There is an $O(m + o)$ time algorithm outputting the set of all maximal C -safe sequences for C -walk covers with no duplicates, where o denotes the total length of all the maximal C -safe sequences.*

4 Simplifying flow decomposition MILP models via safe sequences

In [39], maximal safe sequences in acyclic s - t graphs were used to fix MILP variables encoding s - t paths in the following manner. For every edge e of the graph, one stores a weight $w(e)$ computed as the length of a longest safe sequence containing e , and stores one such sequence as $S(e)$. Then, one computes an antichain of edges, i.e. a set of edges $A = \{e_1, \dots, e_t\}$ that are pairwise unreachable, i.e. there is no path in the graph from some $e_i \in A$ to some $e_j \in A$. Since the graph is acyclic, we then have that also the sequences $S(e_1), \dots, S(e_t)$ are pairwise unreachable. As such, they must appear in different s - t paths in any solution. Thus, for acyclic graphs, for each edge $e_i = u_i v_i \in A$ one can fix $x_{uv,i} = 1$, for all edges uv in sequence $S(e_i)$. For a large number of x variables to be set in this manner, [39] found an antichain with the property that the sum of the lengths of the edges in A , i.e. $\sum_{e \in A} w(e)$, is maximized. This problem of computing a *maximum-weight antichain* can be solved in polynomial time with a reduction to a min-flow problem [27].

The reasoning behind such maximum-weight antichain is that (i) safe sequences containing edges of such an antichain must appear in different solution paths (call them *incompatible*), as we also argue below for walks; and (ii) maximizing the total weight maximizes the number of path variables $x_{uv,i}$ that can be fixed to (at least) 1 in a preprocessing step.

To adapt this approach to general s - t graphs, possibly with cycles, we first show that we can analogously compute a maximum-weight antichain in a general graphs with cycles, with the same complexity as for acyclic graphs (proof in Section B):

Theorem 7. *Let $G = (V, E)$ be a graph, and let $w : E \rightarrow \mathbb{Z}_+ \cup \{0\}$. Computing a maximum-weight antichain of edges of G , namely a set $A = \{e_1, \dots, e_t\}$ of pairwise unreachable edges maximizing $\sum_{e \in A} w(e)$, can be reduced in time $O(|V| + |E|)$ to computing a maximum-weight antichain of edges of a DAG $G' = (V', E')$, with edge weights $w' : E' \rightarrow \mathbb{Z}_+ \cup \{0\}$, where $|V'| \leq 2|V|$ and $|E'| \leq |E|$.*

Next, we show how to use such an antichain $A = \{e_1, \dots, e_t\}$ to fix $x_{uv,i}$ variables encoding s - t walks. Let $S(e_1), \dots, S(e_t)$ be sequences of edges passing through e_1, \dots, e_t . We claim that there can be no s - t walk containing two distinct sequences, say $S(e_i)$ and $S(e_j)$. Indeed, if this were so, then both e_i and e_j would appear on a same s - t walk, and thus either e_i reaches e_j or e_j reaches e_i , which contradicts the fact that A is an antichain.

Therefore, we have that any solution to Definition 1 made up of s - t walk W_1, \dots, W_k is such that: (i) each $S(e_i)$ belongs to a distinct s - t walk; without loss of generality, we can assume that $S(e_i)$ belongs to walk W_i ;

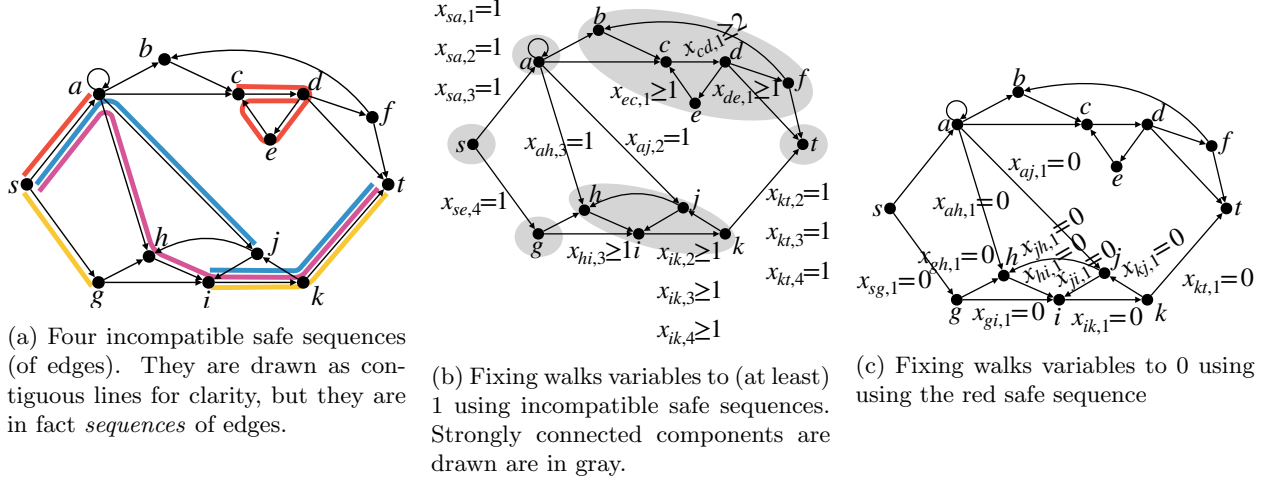


Figure 3: Fixing walk variables using incompatible safe sequences of edges. For the i -th safe sequence (assumed here in the order red, blue, violet, yellow), and for every edge uv in it, we set $x_{uv,i} = 1$ if uv connects different SCCs (it cannot be traversed more than once by the i -th walk, since G^{sc} is acyclic). If uv lies within an SCC, we set $x_{uv,i}$ at least to the number of times uv appears in the sequence (it must be traversed at least this many times, but possibly more). When fixing variables to 0, note that for the red sequence (1st), edge (j, h) does not reach s , is not reached by t , and while j is reachable from a , h does not reach c —thus it cannot be used by the first solution walk.

and (ii), as a consequence, $k \geq t$. In analogy to the walk notation, for a sequence $S(e_i)$, let $S(e_i)(uv)$ denote the number of times uv appears in $S(e_i)$. Based on these observations, we can fix the following variables.

Fixing variables to (at least) 1. We can constrain the walk variables $x_{uv,i}$ of the MILPs by setting (see also Figure 3b):

$$\forall i \in \{1, \dots, t\}, \forall uv \in S(e_i), \quad x_{uv,i} \begin{cases} = 1, & \text{if } uv \text{ is an edge between different SCCs of } G, \\ \geq S(e_i)(uv), & \text{if } uv \text{ is an edge inside some SCC of } G. \end{cases} \quad (1)$$

This is correct because edges between different SCCs cannot be traversed more than once by any s - t walk, while edges inside SCC can be repeated by a walk. Moreover, we can use this observation also to transform *any* integer variable $x_{uv,i}$ (for all $uv \in E$, and all $i \in \{1, \dots, k\}$) to a binary variable belonging to $\{0, 1\}$, whenever uv is an edge between different SCCs of the graph.

Fixing variables to 0. Once we fix the edges of a safe sequence $S(e_i)$, we can further infer that some other edges uv *cannot* be part of the i -th solution walk, in which case we can set $x_{uv,i} = 0$. These are the edges such that *none* of the following conditions hold (see Figure 3c for an illustration):

1. v reaches the first vertex of the safe sequence, i.e. the i -th walk could pass through uv before passing through $S(e_i)$;
2. the last vertex of the safe sequence reaches u , i.e. the i -th walk could pass through uv after passing through $S(e_i)$;
3. (3) for two consecutive edges ab and cd in the sequence, b reaches u and v reaches c , i.e. i -th walk could pass through uv between two consecutive edges ab and cd in $S(e_i)$.

5 Experimental results

Data. In this paper we experiment with graphs created from four bacterial datasets: **ecoli**, containing 30 E.coli strains from the dataset [1], **medium20** and **complex32**, containing 20 and 32 bacterial genomes, across 19 and 26 genera, respectively, from [40], and **JGI**, a mock community containing 23 bacterial strains and 3 archaeal strains, across 22 genera, from [42]. For the latter three, we use the genome abundances provided with the datasets. For **ecoli**, where there are no abundances, we simulate them following the lognormal distribution with mean 1 and variance 1, multiplied by 10 and rounded up to the nearest integer.

To create datasets of increasing complexity, we create the graphs from an increasing number of genomes (5, 10, 15, ...), until including all genomes in the dataset. Further, to create graphs of a feasible size for the ILP models, we partition these genomes into strings (windows) of a fixed length ℓ . Because of the dataset complexity, in **medium20** and **JGI** we use $\ell = 50,000$, in **complex32** we use $\ell = 10,000$, and in **ecoli** we use $\ell = 5,000$. For the i -th window in each genome, we construct the i -th graph as an edge-centric de Bruijn graph, as follows. The nodes are $(k - 1)$ -mers, and the edges are k -mers, for a fixed value k (in all datasets we set $k = 63$). For each occurrence of a k -mer in a genome, we increase the flow value of the corresponding edge by the abundance value of the genome. The genome windows then correspond to walks in this graph, and moreover they, with the abundances used to create the graph, form a flow decomposition. In this graph, we compact all unitigs (i.e. paths whose internal nodes have in-degree and out-degree equal to 1). Finally, we keep only those graphs that contain at least one cycle.

The above creates perfect edge weights where flow conservation holds at every node. For imperfect data, we follow [9] and replace each error-free edge weight x by a sample from the Poisson distribution $Pois(x)$. To provide also subset constraints to the models, when creating the i -th graph, we also simulate 5 reads from each genome window, from a location chosen uniformly at random, by reading (at most) 1,000 bases until the end of the genome window. This string is then mapped to a walk in the de Bruijn graph, by traversing each edge in the order of its k -mers. All graphs are available at <https://doi.org/10.5281/zenodo.17549958>.

Implementation. We implemented the ILPs for the three problems from Definition 1, and described in Section A, in the **flowpaths** Python package, see <https://algbio.github.io/flowpaths>. This package uses by default the popular HiGHS MILP solver, see www.hiogs.dev and [19], ensuring that our implementations can be immediately used without procuring a license. We used the HiGHS Python API (v1.11.0) that gets installed automatically with **flowpaths**. See Section D for minimal code examples of these models. We ran our experiments with version v0.2.12 of **flowpaths**, on a machine with an AMD Threadripper processor and 512GB RAM, and with a timeout of 300 seconds for HiGHS.

To fix a value of k in Definition 1, for k -FD we use the smallest value for which k -FD exists, i.e. we solve the minimum flow decomposition problem (we call the resulting model MinFlowDecomp). For k -LAE and k -MPE, k is the minimum number of s - t walks needed to cover the edges in E' , computed with Theorem 7.

The implementations for MinFlowDecomp, LeastAbsErrors and MinPathError support setting the set of edges E' from Definition 1, support arbitrary sets S and T of start and end nodes for the solution walks (recall the discussion after Definition 1), and support subset constraints. For the MinFlowDecomp and MinPathError models, any solution is also a walk cover of the edges in E' , by Definition 1. Hence, for them we use safe sequences with respect to walk covers of $C = E'$. For LeastAbsErrors, we cannot easily guarantee which edges will be traversed by the solution walks. For it, we add to E' only those edges in the subset constraints, since by definition, they must be covered by the solution walks, thus guaranteeing optimality.

To test different experimental scenarios, for MinFlowDecomp we use $E' = E$ and no subset constraints, for LeastAbsErrors we use $E' = E$ and the subset constraints from our simulation (and use C -safe sequence with C as the edges of the subset constraints). For MinPathError we set E' as the set of those edges whose flow value is above the 25-th percentile of all edge values, and no subset constraints. Here, by definition, all edges in E' must appear in a solution walk, thus C -safe sequences for $C = E'$ still guarantee optimality.

Discussion. The experimental results for MinFlowDecomp are in Table 1. The results for LeastAbsErrors and MinPathError are in Tables 2 and 3 in Section C. We generally observe that, without optimizations,

Table 1: **Experimental results on perfect data with MinFlowDecomp.** Column “#gen” contains the number of genomes used to create the graphs, “#graphs” contains the number of graphs (with cycles) obtained from all windows of that number of genomes, “avg n / m ” contain the average number of nodes / edges of the graphs, while in parentheses we list the maximum number of nodes / edges of the graphs. Column “prep (s)” contains the average time taken by the safety optimizations described in Section 4, “#solved” contains the number of solved instances within the 300s timeout, “time (s)” contains the average time taken by the solved instances (without and with the safety optimizations, from beginning to end), and “speed-up (\times)” contains the average speed-up obtained with the safety optimizations, where if an instance timed out, we considered it to run in 300s.

	#gen	#graphs	avg n (max n)	avg m (max m)	prep (s)	#solved, time (s)		speed-up (\times)
						no safety	safety	
complex32	5	63	26 (144)	37 (194)	0.009	63, 1.00	63, 0.08	13.8
	10	162	28 (151)	44 (209)	0.015	162, 4.34	162, 0.14	28.5
	15	114	33 (156)	54 (219)	0.024	111, 9.13	114, 0.20	58.7
	20	116	42 (226)	69 (318)	0.041	113, 18.61	116, 0.30	65.3
	25	118	70 (284)	111 (406)	0.087	106, 42.04	118, 0.66	94.2
	30	122	77 (289)	124 (416)	0.106	97, 50.26	122, 0.83	111.3
	32	124	80 (291)	131 (420)	0.115	99, 56.74	124, 0.87	115.4
ecoli	5	127	11 (108)	17 (157)	0.004	127, 0.49	127, 0.04	13.1
	10	209	30 (249)	46 (370)	0.016	208, 7.84	209, 0.20	55.0
	15	271	53 (334)	82 (502)	0.036	255, 24.97	270, 0.45	107.5
	20	369	84 (398)	128 (614)	0.072	311, 45.91	367, 0.87	138.7
	25	436	99 (438)	151 (672)	0.107	318, 54.53	432, 2.26	143.8
	30	487	139 (466)	210 (721)	0.194	288, 78.27	476, 1.99	128.6
JGI	5	47	16 (86)	25 (127)	0.005	47, 1.93	47, 0.05	20.7
	10	41	34 (122)	52 (175)	0.019	37, 16.53	41, 0.17	165.8
	15	44	41 (132)	66 (194)	0.032	37, 44.84	44, 0.25	306.3
	20	42	53 (151)	87 (224)	0.048	30, 76.25	42, 0.40	345.7
	26	38	67 (181)	111 (272)	0.081	22, 88.11	38, 0.68	315.8
medium20	5	52	16 (75)	24 (108)	0.005	52, 2.06	52, 0.04	27.5
	10	39	26 (86)	42 (131)	0.015	39, 15.52	39, 0.11	96.6
	15	40	36 (113)	60 (170)	0.026	36, 37.60	40, 0.22	216.1
	20	40	50 (125)	84 (200)	0.042	28, 76.88	40, 0.34	393.2

the running times (and the proportion of instances where the MILP solver times out) increase with the complexity of the MILP models (in the order MinFlowDecomp, LeastAbsErrors, MinPathError), and with the number of genomes used to create the graphs (column “#gen”).

Overall, we observe that the time taken to preprocess the MILP as described in Section 4 is negligible (often less than 0.1s on average), confirming efficiency. Applying these optimizations significantly increases the number of instances solved within the time limit, and presents significant speed-ups.

For MinFlowDecomp all instances are solved, except for a small number in the **ecoli** dataset. The average running times of the optimized MILPs from beginning to end is generally under 1 second on average, presenting speed-ups of up to $393\times$ (on **medium20**, 20 genomes). For LeastAbsErrors, note that we also add subset constraints, and the edges E' used for safety are only from the subset constraints, and not all edges of the graph. The number of solved instances is less than for MinFlowDecomp, especially when more genomes are used to construct the graphs. We observe speed-ups of up to $560\times$ (on **complex32**, 15 genomes); the speed-up values decrease for graphs constructed from more genomes because the number of instances timing out increases. For MinPathError, recall that E' is computed by selecting edges whose flow value is larger than the 25-th percentile. This means that safe sequences are smaller than if setting $E' = E$ (because less edges are guaranteed to appear in a solution walk), but we still observe speedups of more than $1000\times$ for graphs created from 10 genomes. Overall, our optimized MinPathError solves most instances in under one minute on average, and solve most instances within the timeout (except for **ecoli** graphs). We also observe more solved instances compared to LeastAbsErrors, because of the different set E' used to compute safety.

References

- [1] Jarno N. Alanko. 3682 E. coli assemblies from NCBI, May 2022. <https://doi.org/10.5281/zenodo.6577997>.
- [2] Stephen Alstrup, Dov Harel, Peter W Lauridsen, and Mikkel Thorup. Dominators in linear time. *SIAM Journal on Computing*, 28(6):2117–2132, 1999.
- [3] Jasmijn A Baaijens, Leen Stougie, and Alexander Schönhuth. Strain-aware assembly of genomes from mixed samples using flow variation graphs. In *International Conference on Research in Computational Molecular Biology*, pages 221–222. Springer, 2020.
- [4] Jasmijn A Baaijens, Bastiaan Van der Roest, Johannes Köster, Leen Stougie, and Alexander Schönhuth. Full-length de novo viral quasispecies assembly through variation graph construction. *Bioinformatics*, 35(24):5086–5094, 2019.
- [5] Niko Beerenwinkel, Stefano Beretta, Paola Bonizzoni, Riccardo Dondi, and Yuri Pirola. Covering pairs in directed acyclic graphs. *Comput. J.*, 58(7):1673–1686, 2015.
- [6] Elsa Bernard, Laurent Jacob, Julien Mairal, and Jean-Philippe Vert. Efficient rna isoform identification and quantification from rna-seq data with network flows. *Bioinformatics*, 30(17):2447–2455, 2014.
- [7] Manuel Cáceres, Massimo Cairo, Andreas Grigorjew, Shahbaz Khan, Brendan Mumey, Romeo Rizzi, Alexandru I. Tomescu, and Lucia Williams. Width helps and hinders splitting flows. *ACM Trans. Algorithms*, 20(2), March 2024.
- [8] Fernando H C Dias, Manuel Cáceres, Lucia Williams, Brendan Mumey, and Alexandru I Tomescu. A safety framework for flow decomposition problems via integer linear programming. *Bioinformatics*, 39(11):btad640, 10 2023.
- [9] Fernando H. C. Dias and Alexandru I. Tomescu. Accurate flow decomposition via robust integer linear programming. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, pages 1–14, 2024.
- [10] Fernando H. C. Dias, Lucia Williams, Brendan Mumey, and Alexandru I. Tomescu. Fast, flexible, and exact minimum flow decompositions via ILP. In Itsik Pe’er, editor, *Research in Computational Molecular Biology - 26th Annual International Conference, RECOMB 2022, San Diego, CA, USA, May 22-25, 2022, Proceedings*, volume 13278 of *Lecture Notes in Computer Science*, pages 230–245. Springer, 2022.
- [11] Fernando H. C. Dias, Lucia Williams, Brendan Mumey, and Alexandru I. Tomescu. Minimum flow decomposition in graphs with cycles using integer linear programming. *Journal of Global Optimization*, 2025.
- [12] Jianxing Feng, Wei Li, and Tao Jiang. Inference of isoforms from short sequence reads. *Journal of Computational Biology*, 18(3):305–321, 2011.
- [13] Borja Freire, Susana Ladra, José R Paramá, and Leena Salmela. Viquif: de novo viral quasispecies reconstruction using unitig-based flow networks. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 20(2):1550–1562, 2022.
- [14] Loukas Georgiadis and Nikos Parotsidis. Dominators in directed graphs: a survey of recent results, applications, and open problems. *Proc. of ISCIM*, 13:15–20, 2013.

- [15] Andreas Grigorjew, Fernando H. C. Dias, Andrea Cracco, Romeo Rizzi, and Alexandru I. Tomescu. Accelerating ILP solvers for minimum flow decompositions through search space and dimensionality reductions. In Leo Liberti, editor, *22nd International Symposium on Experimental Algorithms, SEA 2024, July 23-26, 2024, Vienna, Austria*, volume 301 of *LIPIcs*, pages 14:1–14:19. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2024.
- [16] Andreas Grigorjew, Wanchote Jiamjitrak, Brendan Mumey, and Alexandru I. Tomescu. Parameterised approximation and complexity of minimum flow decompositions, 2024.
- [17] Tzvika Hartman, Avinatan Hassidim, Haim Kaplan, Danny Raz, and Michal Segalov. How to split a flow? In *2012 Proceedings IEEE INFOCOM*, pages 828–836. IEEE, 2012.
- [18] Steffen Heber, Max Alekseyev, Sing-Hoi Sze, Haixu Tang, and Pavel A Pevzner. Splicing graphs and est assembly problem. *Bioinformatics*, 18(suppl_1):S181–S188, 2002.
- [19] Qi Huangfu and JA Julian Hall. Parallelizing the dual revised simplex method. *Mathematical Programming Computation*, 10(1):119–142, 2018.
- [20] Kyle Kloster, Philipp Kunke, Michael P O’Brien, Felix Reidl, Fernando Sánchez Villaamil, Blair D Sullivan, and Andrew van der Poel. A practical FPT algorithm for flow decomposition and transcript assembly. In *ALENEX 2018 – Proceedings of the Twentieth Workshop on Algorithm Engineering and Experiments*, pages 75–86. SIAM, 2018.
- [21] Mikhail Kolmogorov, Derek M Bickhart, Bahar Behsaz, Alexey Gurevich, Mikhail Rayko, Sung Bong Shin, Kristen Kuhn, Jeffrey Yuan, Evgeny Pevnikov, Timothy PL Smith, et al. metaflye: scalable long-read metagenome assembly using repeat graphs. *Nature methods*, 17(11):1103–1110, 2020.
- [22] Dinghua Li, Chi-Man Liu, Ruibang Luo, Kunihiko Sadakane, and Tak-Wah Lam. Megahit: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de bruijn graph. *Bioinformatics*, 31(10):1674–1676, 2015.
- [23] Yen-Yi Lin, Phuong Dao, Faraz Hach, Marzieh Bakhshi, Fan Mo, Anna Lapuk, Colin Collins, and S Cenk Sahinalp. CLIIQ: Accurate comparative detection and quantification of expressed isoforms in a population. In *WABI 2012 – 12th International Workshop on Algorithms in Bioinformatics*, pages 178–189. Springer, 2012.
- [24] Juntao Liu, Ting Yu, Zengchao Mu, and Guojun Li. TransLiG: a de novo transcriptome assembler that uses line graph iteration. *Genome Biology*, 20:1–9, 2019.
- [25] Runpeng Luo and Yu Lin. Vstrains: De novo reconstruction of viral strains via iterative path extraction from assembly graphs. In *International Conference on Research in Computational Molecular Biology*, pages 3–20. Springer, 2023.
- [26] Xiao Luo, Xiongbao Kang, and Alexander Schönhuth. Strainline: full-length de novo viral haplotype reconstruction from noisy long reads. *Genome biology*, 23(1):29, 2022.
- [27] Ralf Möhring. Algorithmic Aspects of Comparability Graphs and Interval Graphs. In Ivan Rival, editor, *Graphs and Order: the role of graphs in the theory of ordered sets and its applications*. D. Reidel Publishing Company, 1984.
- [28] Brendan Mumey, Samareh Shahmohammadi, Kathryn McManus, and Sean Yaw. Parity balancing path flow decomposition and routing. In *2015 IEEE Globecom Workshops (GC Wkshps)*, pages 1–6. IEEE, 2015.
- [29] Toshiaki Namiki, Tsuyoshi Hachiya, Hideaki Tanaka, and Yasubumi Sakakibara. Metavelvet: an extension of velvet assembler to de novo metagenome assembly from short sequence reads. In *Proceedings of the 2nd ACM conference on bioinformatics, computational biology and biomedicine*, pages 116–124, 2011.

- [30] Sergey Nurk, Dmitry Meleshko, Anton Korobeynikov, and Pavel A Pevzner. metaspades: a new versatile metagenomic assembler. *Genome research*, 27(5):824–834, 2017.
- [31] Yu Peng, Henry CM Leung, Siu-Ming Yiu, and Francis YL Chin. Meta-idba: a de novo assembler for metagenomic data. *Bioinformatics*, 27(13):i94–i101, 2011.
- [32] Alexander J Petri and Kristoffer Sahlin. isonform: reference-free transcriptome reconstruction from oxford nanopore data. *Bioinformatics*, 39(Supplement_1):i222–i231, 2023.
- [33] Mihai Pop. Genome assembly reborn: recent computational challenges. *Briefings in bioinformatics*, 10(4):354–366, 2009.
- [34] Andrey D Prjibelski, Alla Mikheenko, Anoushka Joglekar, Alexander Smetanin, Julien Jarroux, Alla L Lapidus, and Hagen U Tilgner. Accurate isoform discovery with isoquant using long reads. *Nature Biotechnology*, 41(7):915–918, 2023.
- [35] Romeo Rizzi, Alexandru I. Tomescu, and Veli Mäkinen. On the complexity of minimum path cover with subpath constraints for multi-assembly. *BMC Bioinform.*, 15(S-9):S5, 2014.
- [36] Zhaleh Safikhani, Mehdi Sadeghi, Hamid Pezeshk, and Changiz Eslahchi. SSP: An interval integer linear programming for de novo transcriptome assembly and isoform discovery of RNA-seq reads. *Genomics*, 102(5-6):507–514, 2013.
- [37] Palash Sashittal, Chuanyi Zhang, Jian Peng, and Mohammed El-Kebir. Jumper enables discontinuous transcript assembly in coronaviruses. *Nature Communications*, 12(1):6728, 2021.
- [38] Francisco Sena, Eliel Ingervo, Shahbaz Khan, Andrey Prjibelski, Sebastian Schmidt, and Alexandru Tomescu. Flowtigs: Safety in flow decompositions for assembly graphs. *iScience*, 27(12), 2024.
- [39] Francisco Sena, Romeo Rizzi, and Alexandru I Tomescu. Safe sequences via dominators in dags for path-covering problems. In *33rd Annual European Symposium on Algorithms (ESA 2025)*, pages 55–1. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2025.
- [40] Daria Shafranskaya, Varsha Kale, Rob Finn, Alla L Lapidus, Anton Korobeynikov, and Andrey D Prjibelski. Metagt: A pipeline for de novo assembly of metatranscriptomes with the aid of metagenomic data. *Frontiers in Microbiology*, 13:981458, 2022.
- [41] Mingfu Shao and Carl Kingsford. Theory and a heuristic for the minimum path flow decomposition problem. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 16(2):658–670, 2017.
- [42] Esther Singer, Bill Andreopoulos, Robert M Bowers, Janey Lee, Shweta Deshpande, Jennifer Chiniquy, Doina Ciobanu, Hans-Peter Klenk, Matthew Zane, Christopher Daum, et al. Next generation sequencing data of a defined microbial mock community. *Scientific data*, 3(1):1–8, 2016.
- [43] Li Song, Sarven Sabuncian, and Liliana Florea. CLASS2: accurate and efficient splice variant annotation from RNA-seq reads. *Nucleic Acids Research*, 44(10):e98–e98, 2016.
- [44] Moritz Stinzendörfer, Philine Schiewe, and Fabricio Oliveira. A robust optimization approach to flow decomposition. *arXiv preprint arXiv:2410.21140*, 2024.
- [45] Robert Tarjan. Depth-first search and linear graph algorithms. *SIAM journal on computing*, 1(2):146–160, 1972.
- [46] Jia Tian, Ziyu Gao, Minghao Li, Ergude Bao, and Jin Zhao. Accurate assembly of full-length consensus for viral quaspecies. *BMC bioinformatics*, 26(1):36, 2025.

- [47] A. I. Tomescu, T. Gagie, A. Popa, R. Rizzi, A. Kuosmanen, and V. Makinen. Explaining a Weighted DAG with Few Paths for Solving Genome-Guided Multi-Assembly. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 12(06):1345–1354, 2015.
- [48] Alexandru I. Tomescu, Anna Kuosmanen, Romeo Rizzi, and Veli Mäkinen. A novel min-cost flow method for estimating transcript expression with RNA-Seq. *BMC Bioinformatics*, 14(5):S15, 2013.
- [49] Armin Töpfer, Tobias Marschall, Rowena A Bull, Fabio Luciani, Alexander Schönhuth, and Niko Beerenwinkel. Viral quasispecies assembly via maximal clique enumeration. *PLoS computational biology*, 10(3):e1003515, 2014.
- [50] Benedicte Vatinlen, Fabrice Chauvet, Philippe Chrétienne, and Philippe Mahey. Simple bounds and greedy algorithms for decomposing a flow into a minimal set of paths. *European Journal of Operational Research*, 185(3):1390–1401, 2008.
- [51] Riccardo Vicedomini, Christopher Quince, Aaron E Darling, and Rayan Chikhi. Strainberry: automated strain separation in low-complexity metagenomes using long reads. *Nature Communications*, 12(1):4485, 2021.
- [52] Lucia Williams, Alexandru I. Tomescu, and Brendan Mumey. Flow decomposition with subpath constraints. *IEEE ACM Trans. Comput. Biol. Bioinform.*, 20(1):360–370, 2023.
- [53] Zheng Xia, Jianguo Wen, Chung-Che Chang, and Xiaobo Zhou. NSMAP: a method for spliced isoforms identification and quantification from RNA-Seq. *BMC Bioinformatics*, 12:1–13, 2011.
- [54] Qimin Zhang, Qian Shi, and Mingfu Shao. Accurate assembly of multi-end rna-seq data with scallop2. *Nature computational science*, 2(3):148–152, 2022.
- [55] Jin Zhao, Haodi Feng, Daming Zhu, and Yu Lin. MultiTrans: an algorithm for path extraction through Mixed Integer Linear Programming for transcriptome assembly. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 19(1):48–56, 2021.

A MILP formulations for flow decomposition problems into set of weighted walks

In this section we describe the MILP models for the three problems defined in Definition 1. These models use the same main ideas as the MILP model introduced in [11] for the k -FD problem in cyclic graphs, but improve upon it by having fewer constraints. The formulation follows the same general methodology as for flow decomposition in acyclic graphs: first formulate k s - t walks (now walks instead of paths)—as in the first part of Definition 1; second, add constraints for each specific problem objective, as in the second part of Definition 1. The main difficulty here with respect to acyclic graphs is to model s - t walks via MILP.

A.1 Formulation of s - t walks

We first introduce suitable variables and constraints that model k s - t walks in an s - t graph. As in [11], for each walk W_i , and every edge $uv \in E$, we introduce an integer variable $x_{uv,i} \in \mathbb{Z}_+ \cup \{0\}$ representing the number of times the edge uv is traversed by the walk W_i . To ensure that the variables $x_{uv,i}$ correctly model a walk, we use the following lemma from [11].

Lemma 3 ([11]). *Let $G = (V, E)$ be an s - t graph, and let W be a multiset of edges of G , where for every edge $uv \in W$, we denote by $W(uv)$ the number of times the edge appears in the multiset W . It holds that the edges of W can be ordered to form an s - t walk passing $W(uv)$ times through each edge $uv \in E$ if and only if the following conditions hold:*

1. For every $v \in V$, $\sum_{uv \in E} W(uv) - \sum_{vw \in E} W(vw) = \begin{cases} -1 & \text{if } v = s, \\ 1 & \text{if } v = t, \\ 0 & \text{otherwise.} \end{cases}$
2. For every vertex v appearing in some edge of the multiset W , there is an s - v path using only edges in W (i.e. v is reachable from s using edges in W).

Condition 1 of Lemma 3 can be directly modeled with the following constraints, for each vertex $v \in V$:

$$\sum_{uv \in E} x_{uv,i} - \sum_{vw \in E} x_{vw,i} = \begin{cases} -1 & \text{if } v = s, \\ 1 & \text{if } v = t, \\ 0 & \text{otherwise,} \end{cases} \quad \forall v \in V \setminus \{s, t\}, \forall i \in \{1, \dots, k\}. \quad (2)$$

To model condition 2 of Lemma 3, the main idea from [11] is to model a *reachability tree* from s to every vertex v appearing in some edge uv such that $x_{uv,i} > 0$, and made up only of edges with positive x values. That is, we want to ensure that if $x_{uv,i} > 0$, then there is a path from s to v using only edges in the tree.

In this paper we offer a simpler formulation of such a reachability tree than in [11], as follows. First, as in [11], we introduce binary variables $y_{uv,i} \in \{0, 1\}$ with the meaning that $y_{uv,i} = 1$ iff the edge uv is part of this tree. We first state that an edge uv can be in the tree (i.e. $y_{uv,i} = 1$) only if $x_{uv,i} \geq 1$:

$$y_{uv,i} \leq x_{uv,i} \quad \forall uv \in E, \forall i \in \{1, \dots, k\}. \quad (3)$$

Next, we state that if a vertex v appears in some edge uv such that $x_{uv,i} \geq 1$ (call such a vertex *selected*), then there is one, and exactly one, incoming edge wv such that $y_{wv,i} = 1$ (call such an edge *selected*):

$$\sum_{wv \in E} x_{wv,i} \leq M_1 \sum_{wv \in E} y_{wv,i} \quad \forall v \in V, \forall i \in \{1, \dots, k\}, \quad (4)$$

$$\sum_{wv \in E} y_{wv,i} = 1 \quad \forall v \in V, \forall i \in \{1, \dots, k\}, \quad (5)$$

where M_1 is a sufficiently large upper-bound for the sum on the left-hand side of (4), for example the in-degree of v , which we denote $d^-(v)$, multiplied by the maximum number of times that a walk in the solution can visit an edge.

At this point, we have that for every selected vertex v , there is exactly one incoming selected edge wv (i.e. such that $y_{wv,i} = 1$). This means that such edges induce a collection of disjoint cycles, with trees “hanging” from them, and whose edges are oriented downwards. To ensure that they induce a *single tree*, rooted at s , we introduce a *distance-like* variable $d_{v,i} \in \mathbb{Z}_+$ for each vertex v and every walk i . We set $d_{s,i} = 0$ and for every other vertex v , we have the following constraint:

$$d_{v,i} \geq d_{u,i} + 1 - M_2(1 - y_{uv,i}) \quad \forall uv \in E, \forall i \in \{1, \dots, k\}, \quad (6)$$

where M_2 is a sufficiently large upper-bound for $d_{v,i}$, i.e., the number of vertices in the graph. The interpretation of this constraint is as follows. If $y_{uv,i} = 0$ (i.e. the edge uv is not in the tree), then the constraint is vacuous since the right-hand side is non-positive. However, if $y_{uv,i} = 1$, then the constraint becomes $d_{v,i} \geq d_{u,i} + 1$, i.e. the distance of v is at least one more than that of u . This ensures that there are no cycles induced by the edges with $y_{uv,i} = 1$, since otherwise they would induce a cycle of unsatisfiable inequalities. This ensures that they induce a forest, but because of eq. (2), then they induce a single tree, rooted at s .

The key difference between our model for s - t walks and that of [11] lies in how the properties enforced by eq. (6) are modeled. In [11], these properties were expressed through a product of the y and d variables. Although such products can be linearized (as we will discuss in Section A.2), doing so introduces extra variables and constraints, which in turn slows down the MILP solver. In contrast, our approach avoids this complication, since eq. (6) is already a linear constraint.

We summarize the correctness of our model in the following theorem, whose proof is in Section B.

Theorem 8. *Let $G = (V, E)$ be an s - t graph, and let $x_{uv,i} \in \mathbb{Z}_+ \cup \{0\}$ be variables satisfying eq. (2). It holds that there exists an s - t walk W_i such that $W_i(uv) = x_{uv,i}$ for every edge $uv \in E$, if and only if the MILP model made up of eqs. (3) to (6) is feasible.*

Proof. (\Rightarrow) Suppose that there exists an s - t walk W_i such that $W_i(uv) = x_{uv,i}$ for every edge $uv \in E$. Consider the subgraph G' of G consisting only of edges that are present in W , and consider a breath-first search tree T in this graph, rooted at s . For every edge $uv \in E$, set $y_{uv,i} = 1$ if uv is also an edge of T , and $y_{uv,i} = 0$ otherwise. Moreover, for every vertex $v \in V$, set $d_{v,i} = 0$ if $v = s$ or v is not in G' , and otherwise set $d_{v,i}$ as the length of a shortest path from s to v in G' (in terms of number of edges). We claim that these assignments of $y_{uv,i}$ and $d_{v,i}$ variables satisfy eqs. (3) to (6):

- Equation (3) is satisfied because $y_{uv,i} = 1$ only if uv is also an edge of T , which implies $x_{uv,i} \geq 1$.
- For Equation (4), if the sum $\sum_{wv \in E} x_{wv,i}$ is zero, then the condition trivially holds. Otherwise, observe that $\sum_{wv \in E} x_{wv,i} \leq d^-(v) \max_{uv \in E} x_{uv,i}$ and we choose $M_1 \geq d^-(v) \max_{uv \in E} x_{uv,i}$. Since $\sum_{wv \in E} x_{wv,i}$ is not zero, then v is a vertex of G' and thus it has an incoming edge in the breath-first search tree T , say uv . Since we set $y_{uv,i} = 1$ for the edges of the tree, we have that $\sum_{wv \in E} y_{wv,i} \geq 1$ and thus Equation (4) holds.
- Equation (5) holds by the fact that T is a tree, and for every vertex v in G' (except s) we set $y_{uv,i} = 1$ for exactly one edge uv in-coming to v (i.e. the edge of the tree that is in-coming to v).
- For Equation (6), note that if $y_{uv,i} = 0$, the condition holds, as M_2 is chosen larger than the number of the vertices of the graph, and the values of $d_{v,i}$ are distances from s , which are at most the number of vertices of the graph. If $y_{uv,i} = 1$, since then uv is an edge on a shortest path from s to v , and thus the distance $d_{v,i}$ from s to v is equal the distance $d_{u,i}$ from s to u , plus one, and thus Equation (6) holds.

(\Leftarrow) Suppose that the MILP model made up of eqs. (3) to (6) is feasible. We need to show that there exists an s - t walk W_i such that $W_i(uv) = x_{uv,i}$ for every edge $uv \in E$. Thanks to Lemma 3, it suffices to show

that for every selected vertex v (i.e. for which $x_{uv,i} \geq 1$ for some $uv \in E$), there is an s - v path using only edges zw having positive $x_{zw,i}$ value.

Consider the selected edges uv (i.e. for which $y_{uv,i} = 1$). By eqs. (4) and (5), we have that every selected vertex v has exactly one in-coming selected edge. Let G' be the subgraph of G made up of selected edges; we claim that G' has no cycles. Indeed, if there were a cycle made up of selected edges $v_1v_2, v_2v_3, \dots, v_{t-1}v_t, v_tv_1$, then eq. (6) applied to them would imply $v_2 > v_1, v_3 > v_2, \dots, v_t > v_{t-1}, v_1 > v_t$, thus $v_1 > v_1$, a contradiction. Thus, G' is a forest, made up of trees directed away from their roots.

Let s' be an arbitrary root in some tree in G' and note that each tree in the forest contains at least one edge (recall that we created G' from a set of edges). Thus, there is some edge $s'v$ such that $y_{s'v,i} = 1$. By eq. (3), it follows that $x_{s'v,i} \geq y_{s'v,i} \geq 1$. Since the x variables satisfy eq. (2), if $s' \neq s$, then there is at least one in-coming edge of s' , say us' , such that $x_{us'} \geq 1$. Therefore, s' is a selected vertex and therefore $y_{us'} = 1$ for some in-coming edge $u's'$ of s' (possibly $u' = u$), by eqs. (4) and (5). Therefore, also the edge $u's'$ is in G' , which contradicts the fact that s' is the root of a tree of G' .

Therefore, it follows that $s' = s$, and thus G' is made up of a single tree rooted at s , whose edges are directed away from s , and which contains all selected vertices (i.e. those vertices v for which $x_{uv,i} \geq 1$ for some $uv \in E$). Thus, for every selected vertex v we can obtain an s - v path as the path in G' from s to v . All edges uw on such a path have $y_{uw,i} = 1$, by definition of G' , and thus also $x_{uw,i} \geq 1$ (by eq. (3)). \square

A.2 Formulation of problem objectives

Having a formulation for k s - t walks, encoded as variables $x_{uv,i} \in \mathbb{Z}_+ \cup \{0\}$, for all $uv \in E$, and $i \in \{1, \dots, k\}$, we now show how also the problem objectives from Definition 1 can be modeled in MILP.

For each walk $i \in \{1, \dots, k\}$ we introduce a weight variable $w_i \in \mathbb{Z}_+ \cup \{0\}$ (and also a slack variable ρ_i for k -MPE). The problem objectives then require modeling the product $x_{uv,i}w_i$ (and also $x_{uv,i}\rho_i$ for k -MPE). Even though these are not linear terms, as required by a linear program, they can be linearized using a standard power-of-two technique technique, see [11, Remark 7]. For completeness, we present it also here. Suppose we have a product between an integer variable x and an arbitrary variable y , and that \bar{x} is an upper-bound for x , and let $t = \lceil \log_2(\bar{x}) \rceil$. Then we can introduce $t + 1$ binary variables x_0, \dots, x_t together with the linear constraint $x = \sum_{j=0}^t 2^j x_j$. Then, we have that $xy = \sum_{j=0}^t 2^j x_j y$. Each of $x_j y$ is a product between a binary variable x_j and the variable y , which can in turn be linearized in a standard manner, as follows. Assuming that y can take values only in the interval $[\underline{y}, \bar{y}]$ we introduce a new variable z_j to model this product, together with the constraints:

$$x_j \underline{y} \leq z_j \leq x_j \bar{y}, \quad (7)$$

$$y - \bar{y}(1 - x_j) \leq z_j \leq y - \underline{y}(1 - x_j). \quad (8)$$

As such, the objective for k -FD is (assuming we further linearize non-linear terms as described above):

$$\sum_{i=1}^k x_{uv,i} w_i = f(uv), \quad \forall uv \in E, \forall i \in \{1, \dots, k\}. \quad (9)$$

For k -LAE we introduce a new error variable $z_{uv} \in \mathbb{Z}_+ \cup \{0\}$ for every edge $uv \in E'$, together with the constraints:

$$-z_{uv} \leq f(uv) - \sum_{i=1}^k x_{uv,i} w_i \leq z_{uv}, \quad \forall uv \in E'. \quad (10)$$

Then, we require from the MILP to minimize the sum of errors, namely $\min \sum_{uv \in E'} z_{uv}$.

For k -MPE (assuming again we handle non-linear terms as above) it suffices to add the constraints:

$$-\sum_{i=1}^k x_{uv,i} \rho_i \leq f(uv) - \sum_{i=1}^k x_{uv,i} w_i \leq \sum_{i=1}^k x_{uv,i} \rho_i, \quad \forall uv \in E', \quad (11)$$

and require minimizing the sum of slacks $\min \sum_{i=1}^k \rho_i$.

A.3 Formulation of subset constraints

For general graphs with cycles, we introduce here a weaker form of a set of subpath constraints, but one which can be added to the MILP models developed in the previous subsections, i.e. can be expressed as further constraints on the $x_{uv,i}$ variables modeling walks.

More specifically, assume that we are also given a family $\mathcal{S} = \{S_1, \dots, S_\ell\}$ of subsets of edges, namely $S_j \subseteq E$, for all $j \in \{1, \dots, \ell\}$. The decomposition problems defined in Definition 1 are then generalized to additionally require that for every set S_j , there is at least one walk W_i such that W_i contains all edges of S_j . Formally:

$$\forall j \in \{1, \dots, \ell\}, \exists i \in \{1, \dots, k\}, \text{ such that } W_i(uv) \geq 1, \forall uv \in S_j.$$

Note that these constraints just apply to the walks, not to the objectives of the decomposition problems. As such, it suffices to model them on the $x_{uv,i}$ variables and they will be valid for all problems.

We proceed as in [10], with the difference that we now need to handle solution walks and subset constraints (instead of solution paths and subpath constraints). First, we introduce binary variables $p_{uv,i} \in \{0, 1\}$ which will indicate whether edge $uv \in E$ is present in walk $i \in \{1, \dots, k\}$, namely $p_{uv,i} = \min(1, x_{uv,i})$. We enforce this equality with the following two constraints:

$$p_{uv,i} \leq x_{uv,i}, \quad \forall uv \in E, \forall i \in \{1, \dots, k\}, \quad (12)$$

$$x_{uv,i} \leq M_3 p_{uv,i}, \quad \forall uv \in E, \forall i \in \{1, \dots, k\}, \quad (13)$$

where M_3 is a suitably large constant, for example the maximum number of times an edge can appear in a walk.

For every set S_j and every $i \in \{1, \dots, k\}$, we then introduce a binary indicator variable $s_{i,j} \in \{0, 1\}$ which is allowed to equal 1 only if the walk W_i contains all the edges of set S_j . To enforce this, we add the constraint:

$$\sum_{uv \in S_j} p_{uv,i} \geq |S_j| s_{i,j}, \quad \forall i \in \{1, \dots, k\}, \forall j \in \{1, \dots, \ell\}. \quad (14)$$

Finally, we require that every set S_j is present in at least one solution walk W_i :

$$\sum_{i=1}^k s_{i,j} \geq 1, \quad \forall j \in \{1, \dots, \ell\}. \quad (15)$$

B Missing proofs

This section contains the necessary proofs missing from the main text.

Proof of Lemma 1. Note that $u \neq t$ since u is an s -dominator. Let $w = \text{dom}_t(u, k)$ and let y_1, \dots, y_{k-1} be the vertices between u and v in the s -dominator tree. We can assume $w \neq v$.

Observe that $w \neq y_i$ for every $i \in \{1, \dots, k-1\}$. Suppose otherwise. Then $w \neq t$ because t is not an s -dominator, and so there are $k-1$ vertices x_1, \dots, x_{k-1} between w and u in the t -dominator tree. Some x_j is not between u and v in the s -dominator tree, hence there is a $u-w$ path avoiding x_j . Further, since x_j is strictly t -dominated by w there is a $w-t$ path avoiding x_j . Concatenating these paths at w results in a $u-t$ path avoiding x_j , a contradiction.

Suppose now for a contradiction that v is not strictly t -dominated by w . From the point above, we know that there is a $u-v$ path avoiding w . This path concatenated with a $v-t$ path avoiding w gives a $u-t$ path also avoiding w , contradicting the fact that w t -dominates u .

It remains to show that v cannot t -dominate u . Suppose otherwise. Then we have $v = x_i$ for some $i \in \{1, \dots, k-1\}$, since v is strictly t -dominated by w from the point above and t -dominates u by hypothesis. Thus, some y_j is not between v and u in the t -dominator tree, implying that there is a u - v path avoiding y_j . Since y_j is s -dominated by u , there is an s - u path avoiding y_j , which concatenated with the previous path at u gives an s - v path also avoiding y_j , contradicting the fact that y_j s -dominates v . \square

Proof of Lemma 2. Since p_1 and p_2 are distinct, we can assume without loss of generality that p_1 has a vertex not contained in p_2 . Suppose for a contradiction that there is a vertex v in p_1 and p_2 .

Let u be the vertex in p_1 that is not contained in p_2 and that is closest to v in p_1 (notice that the distance between two vertices of the same univocal path is identical in either tree). Then vertex u is either the immediate s -dominator of v or it is immediately s -dominated by v . Suppose that u is the immediate s -dominator of v , so in particular the vertex of p_2 closest to s in the s -dominator tree is v . Since p_1 is univocal, u immediately s -dominates only v and v immediately t -dominates only u . Thus, the path u concatenated with p_2 is univocal because p_2 is univocal, contradicting the maximality of p_2 . For the remaining case when u is immediately s -dominated by v , identically we can deduce that p_2 concatenated with u is univocal, a contradiction. \square

Proof of Theorem 4. (\Rightarrow) Suppose for a contradiction that X is safe and that for every vertex $v \in C$, X is not a subsequence of $\text{extension}(v)$. Then we can build a C -walk cover of G by adding to it an s - t walk avoiding X for every vertex in C . This contradicts the C -safety of X .

(\Leftarrow) If there is a vertex $v \in C$ such that X is a subsequence of $\text{extension}(v)$, then X is safe. To see why, notice that v has to be covered by some s - t walk in any C -walk cover of G . This walk contains the sequence of s -dominators of v followed by the sequence of t -dominators of v by definition of dominance relation. Thus, X is contained in such a walk and therefore it is safe. \square

Proof of Theorem 5. (For brevity, we say “dominator trees” instead of “blue-dominator trees”.)

(\Leftarrow) If u is a leaf in both dominator trees then $\text{extension}(u)$ is clearly C -maximal: no other vertex can be added to this sequence by the equivalence of safe sequences and paths in the dominator trees.

(\Rightarrow) Let T_s and T_t denote the s - and t - dominator trees. Suppose for a contradiction that $\text{extension}(u)$ is a maximal C -safe sequence but u is not a leaf in both dominator trees. Moreover, suppose without loss of generality that u is not a leaf in T_s and let v be a blue-child of u in T_s . Let $k \in \mathbb{N}^+$ be the smallest number such that $\text{dom}_s(v, k) = u$, so u is the k -th ancestor of v in T_s . Notice that $u \neq t$ since u is an s -dominator. Let $w := \text{dom}_t(u, k)$ (note that w may or may not be in C). First we argue that u has a blue-child z in T_s such that z is strictly t -dominated by w and u is not t -dominated by u . We consider two cases. If $v \neq w$ then we can take $z = v$: applying Lemma 1 to u, v, k we have that w t -dominates v and moreover v is not a t -dominator of u since $v \neq w$, so w strictly t -dominates v . Otherwise we have $v = w$.

Vertex v has a blue-child different from u in T_t or u has a blue-child distinct from v in T_s , otherwise uv is a C -univocal path, contradicting the fact that the trees are collapsed. So suppose without loss of generality that u has a blue-child v' in T_s distinct from v . Let $k' \in \mathbb{N}^+$ be the smallest number such that $\text{dom}_s(v', k') = u$, so u is the k' -th ancestor of v' in T_t . Let $w' := \text{dom}_t(u, k')$. Applying Lemma 1 to u, v', k' we have that w' t -dominates v' and v' is not a t -dominator of u unless $v' = w'$. Our goal now is to show that $v' \neq w'$. Suppose for a contradiction that $w' = v'$. We proceed by cases.

If $w = w'$ then we have the following equalities $v = w = w' = v'$, which contradicts the fact that $v \neq v'$.

Otherwise we have $w \neq w'$ and hence w' is proper ancestor or proper descendant of w in T_t , since they both t -dominate u . Suppose without loss of generality that w' is a proper descendant of w in T_t . Notice that v' does not s -dominate v , since otherwise v is not a blue-child of u . Hence, there is a u - v path in G avoiding v' . Moreover, since $w = v$ and $w' = v'$ and w' is a proper descendant of w in T_t , v' is a proper descendant of v and so there is a v - t path in G avoiding v' . These paths concatenated at v result in a u - t path avoiding v' , contradicting the fact that u is t -dominated by v' . Therefore $v' \neq w'$. Recall that w' t -dominates v' and that v' is not a t -dominator of u (from previous application of Lemma 1). Since now we have $v' \neq w'$, we conclude that w' strictly t -dominates v' . So u has a blue-child v' , w' t -dominates u , and w' strictly t -dominates v' and moreover v' is not a t -dominator of u . We can thus put $w = w'$ and take $z = v'$.

Now we argue that $\text{extension}(z)$ properly contains $\text{extension}(u)$. The sequence $\text{extension}(u)$ consists of the sequence of s -dominators of u (call it Y_1) concatenated with the sequence of t -dominators of u (call it Y_2). Analogously, the sequence $\text{extension}(z)$ consists of the sequence of s -dominators of z (call it X_1) concatenated with the sequence of t -dominators of z (call it X_2). Since z is a blue-child of u in T_s , the sequence Y_1 is *strictly* contained in X_1 , and since z is a proper blue-descendant of w in T_t , the sequence Y_2 is contained in X_2 (possibly $Y_2 = X_2$). Therefore, $\text{extension}(u)$ is a proper subsequence of $\text{extension}(z)$. Since z is a blue-child of u , z is blue and so its extension is a C -safe sequence. A contradiction to the maximality of $\text{extension}(u)$ follows. \square

Proof of Theorem 6. Build the dominator trees of G . From these, build the blue-dominator trees of G . In these, collapse C -univocal paths (while storing the ordered sequence of vertices of the collapsed path in the resulting node in one of the trees). These three steps take $O(m + n)$ time.

For every vertex v that is a leaf in both dominator trees, output $\text{extension}(v)$. By Theorem 5 we find every maximal safe sequence during this process. Further, no safe sequence is computed twice since different leaves common to both trees have different extensions (as they differ in at least the extended vertex). Therefore the algorithm runs in time equal to the length of all the maximal C -safe sequences plus the time required for preprocessing described above, so $O(m + o)$ altogether. \square

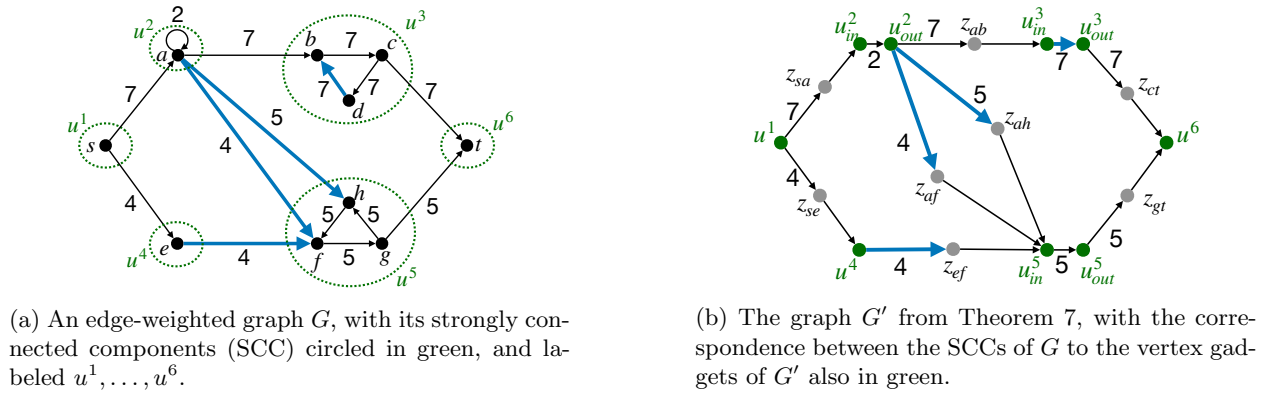


Figure 4: The four blue edges form a maximum-weight edge antichain of G , of weight 20. In G' , every non-trivial SCC u^i of G (i.e. with at least one edge) is replaced by a gadget made up of a single edge $u_{in}^i u_{out}^i$ with weight equaling the maximum weight of an edge in that SCC. Every trivial SCC of G (i.e. with no edges) is kept as a single vertex. All edges between SCCs are replaced by a path of length 2, with a private vertex (in gray) for every edge; the first edge of this length-2 path has the same weight as the original edge, and the second edge of this length-2 path has weight 0 (not shown). The maximum-weight edge antichain of G' , corresponding to the one shown for G , is also shown in blue.

Proof of Theorem 7. We build the following reduction (see also Figure 4 for an illustration). From G , we compute the acyclic graph $G^{scc} = (V^{scc}, E^{scc})$ of strongly connected components (SCCs) of G (the *condensation* of G), in time $O(|V| + |E|)$. In this graph, the strongly connected components correspond to the vertices V^{scc} , and two vertices v' and u' in V^{scc} are joined by an edge if there is an edge between a vertex in the SCC represented by v' to a vertex in the SCC represented by u' . Moreover, we say that an SCC (or a vertex in G^{scc}) is non-trivial if the SCC contains at least one edge (note that a non-trivial SCC can also be made up of a single vertex with a self-loop). We build the reduction graph G' from G^{scc} as follows. If v' is a non-trivial SCC, we represent v' as an edge $v'_{in} v'_{out}$ in G' . In G' we set the weight $w'(v'_{in} v'_{out})$ as the maximum weight of an edge in the SCC represented by v' . Otherwise, we simply add v' to G' . For uniformity of future notation, in such a case we let $v'_{in} := v'$ and $v'_{out} := v'$. For every edge uv in G such that u belongs to some SCC u' and v belongs to a different SCC v' , we add the 3-vertex path $u'_{out} z_{uv} v'_{in}$,

where z_{uv} is a new vertex private to the edge uv , which we add to avoid the creation of parallel edges. In G' , we set $w'(u'_{out}z_{uv}) = w(uv)$, and $w'(z_{uv}v'_{in}) = 0$; namely, the weight of uv is carried over to the weight of the first edge of the newly created path in G' . The graph G' can be computed in linear time, since G^{scc} is computable in linear time [45], and the construction of G' involves only a linear pass over V and E . Note that the number of vertices of G' is at most $2|V|$ and the number of edges is at most $|E|$.

To prove the correctness of the reduction, we prove first that for any antichain A of G , there exists an antichain A' of G' whose weight is at least the weight of A . Let A' be a set of edges of G' obtained from A as follows:

- If $e \in A$ is an edge inside an SCC v' , then v' is a non-trivial SCC and we take in A' the edge $v'_{in}v'_{out}$. Since $w'(v'_{in}v'_{out})$ was chosen as the maximum weight of an edge in the SCC represented by v' , we have that $w'(v'_{in}v'_{out}) \geq w(e)$.
- Otherwise, $e = uv \in A$ is an edge between two SCCs, say u' and v' , and we take in A' the first edge of the corresponding 3-vertex path between u' and v' , namely $u'_{out}z_{uv}$. By construction, its weight satisfies $w'(u'_{out}z_{uv}) = w(uv) = w(e)$.

By the above construction, we have that the weight of A' is at least the weight of A . It remains to show that A' is an antichain in G' . Suppose for a contradiction this is not the case, and suppose that some edge $e'_i \in A'$ (created from an edge $e_i \in A$) reaches some edge $e'_j \in A'$ (created from an edge $e_j \in A$). However, by construction, we have that also e_i reaches e_j , which contradicts the fact that A is an antichain.

Second, it remains to prove that for any antichain A' of G' , there exists an antichain A of G whose weight is at least that of A' . As above, from A' we construct A in a symmetric manner. First, we remove from A' all edges whose weight is 0, as they do not change the weight of A . Next, we construct A as follows.

- If $e' \in A'$ is an edge of the type $e' = v'_{in}v'_{out}$, then we include in A an arbitrary edge e inside the SCC corresponding to v' of maximum weight inside that SCC. By construction, we have that $w(e) = w'(v'_{in}v'_{out}) = w'(e')$.
- If e' is an edge of the type $e' = u'_{out}z_{uv}$, then we include in A the edge uv . By construction, we have that $w(uv) = w'(u'_{out}z_{uv}) = w'(e')$.

Note that edges of G' of the type $z_{uv}v'_{in}$ have weight 0, so they cannot belong to A' anymore. Thus, the weight of A equals the weight of A' . To see that A is an antichain in G , note that if edges e_i and e_j of A are such that there is a path from v_i to u_j in G , then also by construction there is a path in G' between the corresponding edges e'_i and e'_j , which contradicts the fact that A' is an antichain in G' . \square

C Additional experimental results

Table 2: **Experimental results on imperfect data with LeastAbsErrors.** The header is as in Table 1.

	#gen	#graphs	avg n (max n)	avg m (max m)	prep (s)	#solved, time (s)		speed-up (\times)
						no safety	safety	
complex32	5	63	26 (144)	37 (194)	0.012	45, 44.66	59, 12.93	95.8
	10	162	28 (151)	44 (209)	0.020	106, 69.09	152, 14.17	249.3
	15	114	33 (156)	54 (219)	0.031	26, 185.00	105, 15.17	560.2
	20	116	42 (226)	69 (318)	0.046	4, 191.93	107, 17.39	411.2
	25	118	70 (284)	111 (406)	0.074	3, 144.66	91, 17.72	168.8
	30	122	77 (289)	124 (416)	0.100	2, 145.15	93, 27.22	102.1
	32	124	80 (291)	131 (420)	0.110	2, 226.23	90, 28.33	97.6
ecoli	5	127	11 (108)	17 (157)	0.006	104, 27.47	126, 5.48	111.1
	10	209	30 (249)	46 (370)	0.014	53, 135.49	169, 17.74	401.5
	15	271	53 (334)	82 (502)	0.031	4, 137.49	195, 31.57	364.4
	20	369	84 (398)	128 (614)	0.058	1, 9.37	236, 38.81	185.3
	25	436	99 (438)	151 (672)	0.070	0, -	180, 75.12	28.6
	30	487	139 (466)	210 (721)	0.106	0, -	163, 68.33	21.9
JGI	5	47	16 (86)	25 (127)	0.007	29, 63.04	45, 9.13	123.8
	10	41	34 (122)	52 (175)	0.018	3, 112.00	33, 16.05	387.6
	15	44	41 (132)	66 (194)	0.033	0, -	35, 36.96	259.5
	20	42	53 (151)	87 (224)	0.057	1, 9.79	34, 38.19	111.3
	26	38	67 (181)	111 (272)	0.089	1, 92.37	27, 50.20	47.1
	26	38	67 (181)	111 (272)	0.089	1, 92.37	27, 50.20	47.1
medium20	5	52	16 (75)	24 (108)	0.007	36, 43.28	49, 10.15	113.4
	10	39	26 (86)	42 (131)	0.014	7, 159.33	32, 19.05	450.8
	15	40	36 (113)	60 (170)	0.028	1, 247.94	31, 37.66	262.5
	20	40	50 (125)	84 (200)	0.069	0, -	25, 49.29	159.3
	20	40	50 (125)	84 (200)	0.069	0, -	25, 49.29	159.3

Table 3: **Experimental results on imperfect data with MinPathError.** The header is as in Table 1.

	#gen	#graphs	avg n (max n)	avg m (max m)	prep (s)	#solved, time (s)		speed-up (\times)
						no safety	safety	
complex32	5	63	26 (144)	37 (194)	0.012	56, 15.81	63, 0.61	143.0
	10	162	28 (151)	44 (209)	0.017	41, 120.60	156, 2.85	1465.6
	15	114	33 (156)	54 (219)	0.029	6, 128.93	109, 4.09	965.3
	20	116	42 (226)	69 (318)	0.048	2, 193.06	112, 8.32	619.5
	25	118	70 (284)	111 (406)	0.085	0, -	110, 14.02	277.5
	30	122	77 (289)	124 (416)	0.116	0, -	110, 19.70	189.5
	32	124	80 (291)	131 (420)	0.122	0, -	106, 18.28	173.0
ecoli	5	127	11 (108)	17 (157)	0.005	118, 14.85	127, 0.53	126.1
	10	209	30 (249)	46 (370)	0.017	21, 131.66	207, 4.07	1228.5
	15	271	53 (334)	82 (502)	0.039	1, 250.74	260, 7.25	694.4
	20	369	84 (398)	128 (614)	0.072	0, -	336, 17.70	361.4
	25	436	99 (438)	151 (672)	0.094	0, -	377, 19.75	177.4
	30	487	139 (466)	210 (721)	0.148	0, -	390, 23.22	127.4
JGI	5	47	16 (86)	25 (127)	0.010	33, 32.78	46, 1.10	330.4
	10	41	34 (122)	52 (175)	0.021	0, -	39, 6.79	840.5
	15	44	41 (132)	66 (194)	0.043	0, -	39, 2.43	487.7
	20	42	53 (151)	87 (224)	0.057	0, -	39, 5.91	277.6
	26	38	67 (181)	111 (272)	0.102	0, -	33, 7.50	174.2
	26	38	67 (181)	111 (272)	0.102	0, -	33, 7.50	174.2
medium20	5	52	16 (75)	24 (108)	0.007	42, 22.17	52, 1.38	170.4
	10	39	26 (86)	42 (131)	0.019	0, -	39, 9.27	1129.3
	15	40	36 (113)	60 (170)	0.033	0, -	39, 5.32	537.0
	20	40	50 (125)	84 (200)	0.052	0, -	37, 18.03	273.7
	20	40	50 (125)	84 (200)	0.052	0, -	37, 18.03	273.7

D Flowpaths example code

The models in this paper were implemented in the `flowpaths` Python package version 0.2.11, which is installable as:

```
pip install flowpaths
```

Below, we give minimal usage examples of our implementation, to illustrate their use in potential applications employing these models. Note that the decomposition models in `flowpaths` accept graph object created using the popular NetworkX Python library. Moreover, the `flowpaths` package uses the `graphviz` Python package to visualize graphs, which can be useful in applications. We also show these visualizations for the code fragments below. The complete documentation and further examples of the classes implementing our models is available in the `flowpaths` documentation, at algbio.github.io/flowpaths.

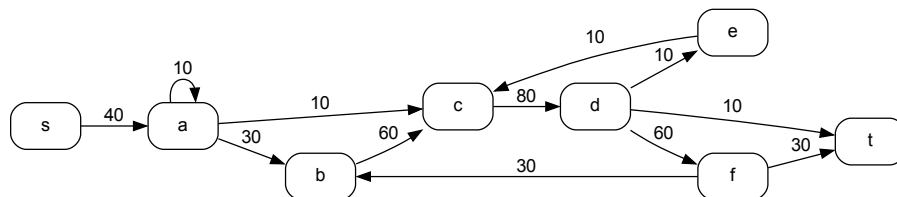
D.1 Minimum Flow Decomposition

```
import flowpaths as fp
import networkx as nx

# We create a NetworkX graph with a flow value on each edge
graph = nx.DiGraph()
graph.add_edge('s', 'a', flow=4)
graph.add_edge('a', 'a', flow=1)
graph.add_edge('a', 'b', flow=3)
graph.add_edge('a', 'c', flow=1)
graph.add_edge('b', 'c', flow=6)
graph.add_edge('c', 'd', flow=8)
graph.add_edge('d', 'e', flow=1)
graph.add_edge('d', 'f', flow=6)
graph.add_edge('d', 't', flow=1)
graph.add_edge('e', 'c', flow=1)
graph.add_edge('f', 'b', flow=3)
graph.add_edge('f', 't', flow=3)

# We visualize the input graph
fp.utils.draw(
    G = graph,
    filename = 'input-mfd.pdf',
    flow_attr = 'flow',
    draw_options={'show_edge_weights': True})
```

The input graph rendered as a PDF file by the above code is the following:



Next, we proceed to creating and solving the model:

```
# We create the model
model = fp.MinFlowDecompCycles(
    G = graph,
    flow_attr = 'flow')

# We solve the model
model.solve()
```

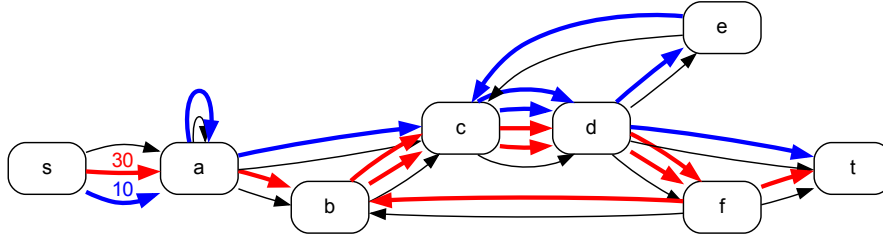


```

# If the model was solved, we retrieve the solution
if model.is_solved():
    solution = model.get_solution()
    print('Solution walks:', solution['walks'])
    print('Solution weights:', solution['weights'])
    fp.utils.draw( # We visualize the solution
        G = graph,
        filename = 'solution-mfd.pdf',
        paths=solution['walks'],
        weights=solution['weights'])

```

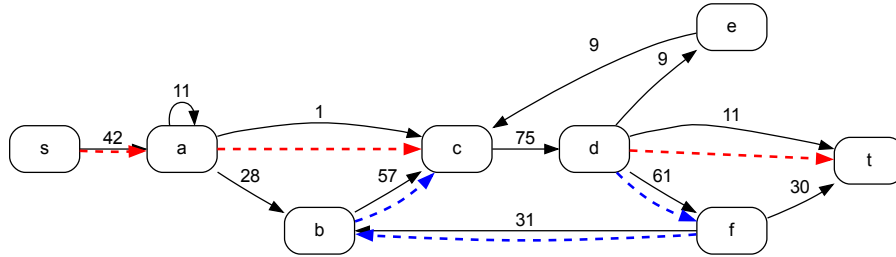
The solution walks visualized by the above code are the blue and red walks in the figure below, of weights 10 and 30, respectively. With the default drawing options, the weight is printed on the first edge of the walk. Note also that e.g. the blue walk traverses the edge (c, d) two times.



D.2 Least Absolute Errors

For the LeastAbsErrors model, we illustrate the use of both subset constraints, and of enforcing the objective from Definition 1 on a proper subset $E' \subsetneq E$ of edges.

Assume that the input graph is the one illustrated below, where there are two subset constraints $[('s', 'a'), ('a', 'c'), ('d', 't')]$ and $[('f', 'b'), ('b', 'c'), ('d', 'f')]$, visualized as dotted colored edges below.



We can pass such subset constraints via the parameter `subset_constraints`. Moreover, since the flow value on edge (a, c) is much smaller than the other values, one could remove it from the set E' . This is achieved by passing such edges to the parameter `elements_to_ignore`.

```

# We create the model
model = fp.kLeastAbsErrorsCycles(
    G = graph,
    weight_type = int,
    flow_attr = 'flow',
    subset_constraints = [
        [('s', 'a'), ('a', 'c'), ('d', 't')],
        [('b', 'c'), ('d', 'f')]
    ],
    elements_to_ignore = [('a', 'c')]
)
# We solve the model

```

```

model.solve()

# If the model was solved, we retrieve the solution
if model.is_solved():
    solution = model.get_solution()
    print('Solution walks:', solution['walks'])
    print('Solution weights:', solution['weights'])

```

The solution walks are visualized in the PDF figure below.

