

To Store or Not to Store: a graph theoretical approach for Dataset Versioning

Abstract—Dataset Versioning is extremely important for enterprises nowadays for ensuring reproducibility of results, tracking data changes over time, maintaining quality measures, enabling collaboration and ensuring legal compliance. Hence, cost-effective data management becomes necessary to reduce storage and reconstruction costs of datasets. In this work, we study the *cost efficient data versioning problem* where the goal is to optimize the storage and reconstruction costs of data versions, given a graph of datasets as nodes and edges capturing edit/reconstruction information. This problem (along with its variants) was introduced by Bhattacharjee et al. [1]. One central variant we study is MINSUM RETRIEVAL (MSR) where the goal is to minimize the total retrieval costs, keeping the storage costs bounded. While such problems are frequently encountered in collaborative tools, e.g., version control of source code, data analysis pipelines etc., to the best of our knowledge, there is not much research studying the theoretical aspects of these problems.

We show the best known heuristic, LMG (introduced in [1]) can perform arbitrarily badly in certain cases. Moreover, we show that it is hard to get $o(n)$ -approximation for MSR on general graphs even if we relax the storage constraints by an $O(\log n)$ factor. Similar hardness results are shown for other variants. We propose poly-time approximation schemes for tree-like graphs, motivated by the fact that the graphs arising in practice from typical edit operations are often not arbitrary. In fact, as version graphs typically have low treewidth, we further develop new algorithms for bounded treewidth graphs.

Furthermore, we propose two new heuristics and evaluate them empirically. First, we extend LMG by considering more potential “moves”, to propose a new heuristic LMG-All. LMG-All consistently outperforms LMG while having comparable run time on a wide variety of datasets, i.e., version graphs. Secondly, we apply our tree algorithms on the minimum-storage arborescence of an instance, yielding algorithms that are qualitatively better than all previous heuristics for MSR, as well as for another variant BOUNDEDMIN RETRIEVAL (BMR).

Index Terms—component, formatting, style, styling, insert

I. INTRODUCTION

The management and storage of data versions has become increasingly important to enterprises. As an example, the increasing usage of online collaboration tools allows many collaborators to edit an original dataset simultaneously, producing multiple versions of datasets to be stored daily. Large number of dataset versions also occur often in industry data lakes [2] where huge tabular datasets like product catalogs might require a few records (or rows) to be modified periodically, resulting in a new version for each such modification. Furthermore, in Deep Learning pipelines, multiple versions are generated from the same original data for the purpose of training and insight generation. At the scale of terabytes or even petabytes, storing and managing all the versions is extremely costly in the aforementioned situations [3]. Therefore, it is no surprise that data version control is emerging as a hot area in the industry [4, 5, 6, 7, 8, 9], and even popular cloud solution providers

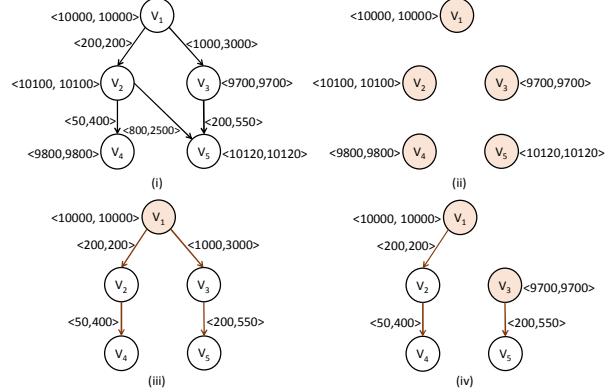


Fig. 1: (i) A version graph over 5 datasets – annotation $\langle a, b \rangle$ indicates a storage cost of a and a retrieval cost of b ; (ii, iii, iv) three possible storage graphs. The figure is taken from [1]

like Databricks are now capturing data lineage information, which helps in effective data version management [10].

In a pioneering paper, Bhattacharjee et al. [1] proposed a model capturing the trade-off between *storage* cost and *retrieval* (recreation) cost. The problems they studied can be defined as follows. Given dataset versions and a subset of the “*deltas*” between them, find a compact representation that minimizes the overall storage as well as the retrieval costs of the versions. This involves a decision for each version: either we *materialize* it (i.e., store it explicitly) or we store a “*delta*” and rely on edit operations to retrieve the version from another materialized version if necessary. The downside of the latter is that, to retrieve a version that was not materialized, we have to incur a computational overhead that we call *retrieval cost*.

Figure 1, taken from Bhattacharjee et al. [1], illustrates the central point through different storage options. (i) shows the input graph, with annotated storage and retrieval costs. If the storage size is not a concern, we should store all versions as in (ii). From (iii) to (iv), it is clear that, by materializing v_3 , we shorten the retrieval costs of v_3 and v_5 .

This retrieval/storage trade-off leads to the combinatorial problem of minimizing one type of cost, given a constraint on the other. There are variations of our objective function as well: retrieval cost of a solution can be measured by either the maximum or total (or equivalently average) retrieval cost of files. This yields four different optimization problems (Problems 3–6 in Table I).

There are some follow-up works on this model [11, 12, 13]. However, those either formulate new problems in different use cases [12, 13, 14] or implement a system incorporating the feature to store specific versions and deltas [13, 15, 16]. We will discuss this in more detail in Section I-B.

Problem Name	Storage	Retrieval
MINIMUM SPANNING TREE	min	$\mathcal{R}(v) < \infty, \forall v$
SHORTEST PATH TREE	$< \infty$	$\min \{\max_v R(v)\}$
MINSUM RETRIEVAL (MSR)	$\leq \mathcal{S}$	$\min \{\sum_v R(v)\}$
MINMAX RETRIEVAL (MMR)	$\leq \mathcal{S}$	$\min \{\max_v R(v)\}$
BOUNDED SUM RETRIEVAL (BSR)	min	$\sum_v R(v) \leq \mathcal{R}$
BOUNDED MAX RETRIEVAL (BMR)	min	$\max_v R(v) \leq \mathcal{R}$

TABLE I: Problems 1-6

A. Our Contributions

We provide the first set of *approximation algorithms* and *inapproximability results* for the aforementioned optimization problems under various conditions. In line with hardness results, we show in the full paper that a simple path structure causes LMG, the best previous heuristic for MSR [1], perform arbitrary poorly. Here, though, we focus on presenting new algorithms and refer to our full paper¹ for descriptions and proofs of the hardness results that are summarized below.

Problem	Graph type	Assumptions	Inapproximability
MSR	arborescence	Triangle inequality $r = s$ on edges ²	$\frac{1}{1 + \frac{1}{e} - \epsilon}$
	undirected		$\Omega(n)^3$
	general		$2 - \epsilon$
MMR	undirected		$\log^* n - \omega(1)^4$
	general		1
	arborescence		$(\frac{1}{2} - \epsilon) \log n$
BSR	undirected		$(1 - \epsilon) \log n$
BMR	undirected		

TABLE II: Hardness results

MMR and BMR. We prove that it is hard to approximate MMR within $\log^* n^4$ factor and BMR within $\log n$ factor on general inputs. Meanwhile, in Section III we give a polynomial-time dynamic programming (DP) algorithm for the two problems on *bidirectional trees*, i.e., digraphs whose underlying undirected graph is a tree. These inputs capture the cases where new versions are generated via edit operations.

We also briefly describe an FPTAS (defined below) for MMR, analogous to the main result for MSR in Section IV.

MSR and BSR. We prove that it is hard to design $(\Omega(n), \Omega(\log n))$ -bicriteria approximation⁵ for MSR or $\Omega(\log n)$ -approximation for MSR. It is also NP-hard to solve the two problems exactly on trees.

On the other hand, we again use DP to design a fully polynomial-time approximation scheme (FPTAS) for MSR on *bounded treewidth graphs*. These inputs capture many practical settings: bidirectional trees have width 1, series-parallel graphs have width 2, and the GitHub repositories we use in (Section VI) all have low treewidth.⁶

¹<https://github.com/Sooooffia/Graph-Versioning/blob/main/Full-Paper.pdf>

²Both are assumptions in previous work [1] that simplify the problems. We note that our algorithms function even without these assumptions.

³This is true even if we relax \mathcal{S} by $O(\log n)$.

⁴ $\log^* n$ is “iterated logarithm”, defined as the number of times we iteratively take logarithm before the result is at most 1.

⁵An (α, β) -bicriteria approximation refers to an algorithm that potentially exceeds the constraint by α times, in order to achieve a β -approximation of the objective. See Section II for an example.

⁶datasharing, styleguide, and leetcode have treewidth 2,3, and 6 respectively.

Graphs	Problems	Algorithm	Approx.
General Digraph	MSR	LMG-All	heuristic
	MSR & MMR	DP-BTW	$1 + \epsilon$
Bounded Treewidth	BSR & BMR		$(1, 1 + \epsilon)$
	MMR	DP-BMR	
Bidirectional Tree	BMR		exact

TABLE III: Algorithms summary. Here, \mathcal{R}_{\max} is the maximum retrieval cost between any pair of vertices in the tree.

New Heuristics. We improved LMG into a more general LMG-All algorithm for solving MSR. LMG-All outperforms LMG in all our experiments and runs faster on sparse graphs.

Inspired by our algorithms on trees, we also propose two DP heuristics for MSR and BMR. Both algorithms perform extremely well even when the input graph is not tree-like. Moreover, there are known procedures for parallelizing general DP algorithms [17], so our new heuristics are potentially more practical than previous ones, which are all sequential.

B. Related Works

1) Theory: There was little theoretical analysis on the exact problems we study. The optimization problems are first formalized in Bhattacherjee et al. [1], which also compared the effectiveness of several proposed heuristics on both real-world and synthetic data. Zhang et al. [11] followed up by considering a new objective that is a weighted sum of objectives in MSR and MMR. They also modified the heuristics to fit this objective. There are similar concepts, including *Light Approximate Shortest-path Tree (LAST)* [18] and *Shallow-light Tree (SLT)* [19, 20, 21, 22, 23, 24]. However, this line of work focuses mainly on undirected graphs and their algorithms do not generalize to the directed case. Among the two problems mentioned, SLT is closely related to MMR and BMR. Here, the goal is to find a tree that is **light** (minimize weight) and **shallow** (bounded depth). To our knowledge, there are only two works that give approximation algorithms for directed shallow-light trees. Chimani and Spoerhase [25] gives a bi-criteria $(1 + \epsilon, n^\epsilon)$ -approximation algorithm that runs in polynomial-time. Recently, Ghuge and Nagarajan [26] gave a $O(\frac{\log n}{\log \log n})$ -approximation algorithm for *submodular tree orienteering* that runs in quasi-polynomial time. Their algorithm can be adapted into $O(\frac{\log^2 n}{\log \log n})$ -approximation for BMR. For MSR, their algorithm gives $(O(\frac{\log^2 n}{\log \log n}), O(\frac{\log^2 n}{\log \log n}))$ -approximation. The idea is to run their algorithm for many rounds, where the objective of each round is to *cover as many nodes as possible*.

2) Systems: To implement a system captured by our problems, components spanning multiple lines of works are required. For example, to get a graph structure, one has to keep track of history of changes. This is related to the topic of data provenance [27, 28]. Given a graph structure, the question of modeling “deltas” is also of interest. There is a line of work dedicated to studying how to implement diff algorithms in different contexts [29, 30, 31, 32, 33].

In the more flexible case, one may think of creating deltas without access to the change history. However, computing all possible deltas is too wasteful, hence it is necessary to

utilize other approaches to identify similar versions/datasets. Such line of work is known as dataset discovery or dataset similarity [2, 34, 35, 36, 37].

Several follow-up works of Bhattacherjee et al. [1] have implemented systems with a feature that saves only selected versions to reduce redundancy. There are works focusing on version control for relational databases [13, 15, 16, 38, 39, 40, 41, 42] and works focusing on graph snapshots [14, 43, 44]. However, since their focus was on designing full-fledged systems, the algorithms they proposed are rather simple heuristics with no theoretical guarantees.

3) *Use cases*: In a version control system such as git, our problem is similar to what `git pack` command aims to do.⁷ The original heuristic for `git pack`, as described in an IRC log, is to sort objects in particular order and only create deltas between objects in the same window.⁸ It is shown in Bhattacherjee et al. [1] that git’s heuristic does not work well compared to other methods.

SVN, on the other hand, only stores the most recent version and the deltas to past versions [45]. Other existing data version management systems include [5, 6, 7, 8, 9], which offer git-like capabilities suited for different use cases, such as data science pipelines in enterprise setting, machine learning-focused, data lake storage, graph visualization, etc.

Though not directly related to our work, recently, there has been a lot of work exploring algorithmic and systems related optimizations for reducing storage and maintenance costs of data. For example, Mukherjee et al. [3] proposes optimal multi-tiering, compression and data partitioning, along with predicting access patterns for the same. Other works that exploit multi-tiering to optimize performance include e.g., [46, 47, 48, 49] and/or costs, e.g., [49, 50, 51, 52, 53, 54]. Storage and data placement in a workload aware manner, e.g., [49, 55, 56] and in a device aware manner, e.g., [57, 58, 59] have also been explored. [47] combine compression and multi-tiering for optimizing latency.

II. PRELIMINARIES

In this section, the definition of the problems, notations, simplifications, and assumptions will be formally introduced.

A. Problem Setting

In the problems we study, we are given a directed *version graph* $G = (V, E)$, where vertices represent *versions* and edges capture the *deltas* between versions. Every edge e is associated with two weights: storage cost s_e and retrieval cost r_e .⁹ The cost of storing e is s_e , and it takes r_e time to retrieve v once we retrieved u . Every vertex v is associated with only the storage cost, s_v , of storing (materializing) the version. Since there is usually a smallest unit of cost in the real world, we will assume $s_v, s_e, r_e \in \mathbb{N}$ for all $v \in V, e \in E$.

To retrieve a version v from a materialized version u , there must be some path $P = \{(u_{i-1}, u_i)\}_{i=1}^n$ with $u_0 = u, u_n = v$,

⁷<https://www.git-scm.com/docs/git-pack-objects>

⁸<https://github.com/git/git/blob/master/Documentation/technical/pack-heuristics.txt>

⁹If $e = (u, v)$, we may use $s_{u,v}$ in place of s_e and $r_{u,v}$ in place of r_e .

such that all edges along this path are stored. In such cases, we say that v is retrieved from materialized u with retrieval cost $R(v) = \sum_{i=1}^n r(u_{i-1}, u_i)$. In the rest of the paper, we say v is “retrieved from u ” if u is in the path to retrieve v , and v is “retrieved from materialized u ” if in addition u is materialized.

The general optimization goal is to select vertices $M \subseteq V$ and edges $F \subseteq E$ of *small size* (w.r.t. storage cost s), such that for each $v \in V \setminus M$, there is a *short path* (w.r.t retrieval cost r) from a materialized vertex to v . Formally, we want to minimize (a) total storage cost $\sum_{v \in M} s_v + \sum_{e \in F} s_e$, and (b) total (resp. maximum) retrieval cost $\sum_{v \in V} R(v)$ (resp. $\max_{v \in V} R(v)$).

Since the storage and retrieval objectives are negatively correlated, a natural problem is to constrain one objective and minimize the other. With this in mind, four different problems are formulated, as described by Problems 3-6 in Table I. These problems are originally defined in Bhattacherjee et al. [1], although we use different names for brevity. Since the first two problems are well studied, we do not discuss them further.

We note that MSR and BSR (MMR and BMR, resp.), are closely related. Given an algorithm for one, we can turn it into an algorithm for the other by binary-searching over the possible values of the constraint. Due to the somewhat exchangable nature of the storage and constraints in these problems, it’s worth considering (α, β) -bicriteria approximations, where we relax the constraint by some α factor in order to achieve a β -approximation. For example, an algorithm A is (α, β) -bicriteria approximation for MSR if it outputs a feasible solution with storage cost $\leq \alpha \cdot \mathcal{S}$ and retrieval cost $\leq \beta \cdot OPT$ where OPT is the retrieval cost of an optimal solution.

III. EXACT ALGORITHM FOR MMR AND BMR ON BI-DIRECTIONAL TREES

As discussed in introduction, we can use an algorithm for BMR to solve for MMR via binary search. Hence, it suffices to focus on BMR, namely, when we are given maximal retrieval constraint \mathcal{R} and want to minimize storage cost.

Let $T = (V, E)$ be a bidirectional tree (abbreviated “tree”) and let \mathcal{R} be the maximum retrieval cost constraint. We can pick any vertex v_0 as root, and orient the tree such that v_0 has no parent, while all other nodes have exactly one parent.

For each $v \in V$, let $T_{[v]}$ denote the subtree of T rooted at v . If v is retrieved from materialized u , we use p_v^u to denote the parent of v on the unique $u - v$ path to retrieve v . We write $p_v^v = v$. We now describe a dynamic programming (DP) algorithm DP-BMR that solves BMR exactly on T .

DP variables. For $u, v \in V$, let $DP[v][u]$ be the minimum storage cost of a *partial solution* on $T_{[v]}$, which satisfies the following: all descendants of v are retrieved from some node in $T_{[v]}$ (possibly itself), while v is retrieved from the materialized version u , which is *potentially outside the subtree* $T_{[v]}$. See Figure 2 for an illustration.

Importantly, note that when calculating the storage cost for $DP[v][u]$, if u is not a part of $T_{[v]}$, the incident edge (p_v^u, v) is involved in the calculation, while other edges in the $u - v$ path, or the cost to materialize u , are not involved in it.

Base case. We iterate from the leaves up. Let $R(u, v)$ denote the retrieval cost of the $u - v$ path. For a leaf v , we

Algorithm 1: DP-BMR

Input: T , a tree, and \mathcal{R} , the max retrieval constraint;
 Orient T arbitrarily. Sort V in reverse topological
 order;
 $DP[v][u] \leftarrow \infty$ for all $v, u \in V$;
for v in V **do**

- for** u in V such that $R(u, v) \leq \mathcal{R}$ **do**

 - if** $u = v$ **then**

 - $DP[v][u] \leftarrow s_v$;

 - else**

 - $DP[v][u] \leftarrow s_{p_v^u},$ where p_v^u is the node
 before v on the path from u to v ;

 - for** w that is a child of v **do**

 - if** w in the path from u to v **then**

 - $DP[v][u] \leftarrow DP[v][u] + DP[w][u]$;

 - else**

 - $DP[v][u] \leftarrow$
 $DP[v][u] + \min\{OPT[w], DP[w][u]\}$;

$OPT[v] \leftarrow \min\{DP[v][w] : w \in V(T_{[v]})\}$;

return $OPT[v_{root}]$;

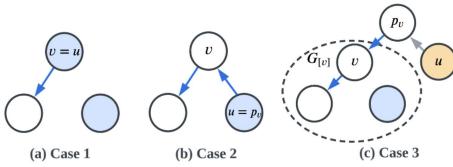


Fig. 2: 3 cases of DP-BMR, where $u = v, u \in V(T_{[v]})$, and $u \notin V(T_{[v]})$ respectively. The blue nodes and edges are stored in the partial solution.

set $DP[v][v] = s_v$, and $DP[v][u] = s_{(p_v^u, v)}$ for all $u \neq v$ with $R(u, v) \leq \mathcal{R}$. Here, p_v^u is just the parent of v in the tree structure. All choices of u, v such that $R(u, v) > \mathcal{R}$ are infeasible, and we therefore set $DP[v][u] = \infty$ in these cases.

Recurrence. For convenience, we define helper variable $OPT[v]$ to be the minimum storage cost on the subproblem $T_{[v]}$, such that v is either materialized or retrieved from one of its materialized descendants.¹⁰ In other words,

$$OPT[v] = \min\{DP[v][w] : w \in V(T_{[v]})\}$$

For recurrence on $DP[v][u]$ such that $R(v, u) \leq \mathcal{R}$, there are three possible cases of the relationship between v and u (see Figure 2). In each case, we outline what we add to $DP[v][u]$.

Case 1. If $u = v$, we materialize v , and each child w of v can be either materialized, or retrieved from their materialized descendants, or retrieved from the materialized $u = v$. Note that this is exactly $\min\{OPT[w], DP[w][u]\}$, and similar facts hold for the following two cases as well.

Case 2. If $u \in V(T_{[v]}) \setminus \{v\}$, we would store the edge (p_v^u, v) . Note that p_v^u is a child of v and hence is also retrieved from the materialized u , so we must add $DP[p_v^u][u]$. We then add $\min\{OPT[w], DP[w][u]\}$ for all other children w of v .

¹⁰Note that the case where v is retrieved from u outside of $T_{[v]}$, or case 3 in Figure 2, is not considered in this helper variable.

Case 3. If $u \notin V(T_{[v]})$, we add the edge (p_v^u, v) , where p_v^u is the parent of v in the tree structure. We then add $\min\{OPT[w], DP[w][u]\}$ for all children as before.

Output. We output $OPT[v_{root}]$, which is the storage cost of the optimal solution. To output the configuration achieving this optimum, we can use the standard procedure where we store the configuration in each DP variable.

Theorem 3.1: BOUNDEDMAX RETRIEVAL is solvable on bidirectional tree instances in $O(n^2)$ time.

The time complexity follows from the observation that each calculation of $DP[v][u]$ in the recurrence takes $O(\deg(v))$ time, and $\sum_u \sum_v \deg(v) = \sum_u O(n) = O(n^2)$. The optimality of this DP can be shown inductively from leaves up, and is omitted due to space limitations.

We note that by binary-searching the constraint value \mathcal{S} , this algorithm also solves MINMAX RETRIEVAL on trees.

IV. FPTAS FOR MSR VIA DYNAMIC PROGRAMMING

In this section we work on MINSUM RETRIEVAL and present a fully polynomial time approximation scheme (FPTAS) on digraphs whose *underlying undirected graph* has bounded treewidth. Similar techniques can be applied to MMR, but we will focus on MSR due to space constraints.

We start by describing a dynamic programming (DP) algorithm on trees in Section IV-A. In Section IV-B, we define all notations necessary for the latter subsection. Finally, in Section IV-C, we show how to extend our DP to the bounded treewidth graphs.

A. Warm-up: Bidirectional Trees

As a warm-up to the more general algorithm, we present an FPTAS for bidirectional tree instances of MSR via DP. This algorithm also inspired a practical heuristic DP-MSR, presented in Section V-B.

WLOG, we assume the tree has a designated root v_{root} and a parent-child hierarchy. We further assume that the tree is binary, via the standard trick of vertex splitting and adding edges of zero weight if necessary.

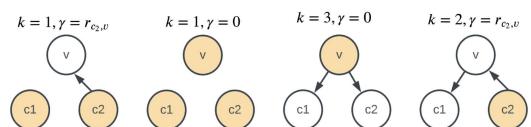


Fig. 3: An illustration of DP variables in Section IV-A

DP variables. We define $DP[v][k][\gamma][\rho]$ to be the minimum storage cost for the subproblem with constraints v, k, γ, ρ such that (with examples illustrated in Figure 3)

- 1) *Root for subproblem* $v \in V$ is a vertex on the tree; in each iteration, we consider the subtree rooted at v .
- 2) *Dependency number* $k \in \mathbb{N}$ is the number of versions retrieved from v (including v itself) in the subproblem solution. This is useful when calculating the extra retrieval cost incurred by retrieving v from its parent.
- 3) *Root retrieval* $\gamma \in \mathbb{N}$ represents the cost of retrieving the subtree root v , if it is retrieved from a materialized descendant. This is useful when calculating the extra

retrieval cost incurred by retrieving the parent of v from v . Note that the root retrieval cost will be discretized, as specified later.

- 4) Total retrieval $\rho \in \mathbb{N}$ represents the total retrieval cost of the subsolution. Similar to γ , ρ will also be discretized.

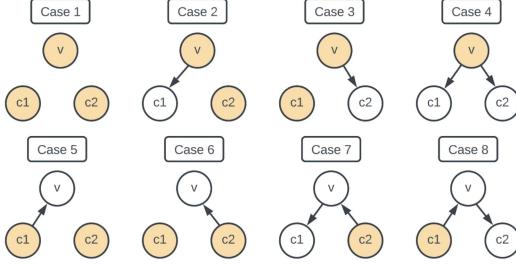


Fig. 4: Eight types of connections on a binary tree. A node is colored if it is materialized or retrieved via delta from outside the chart. Otherwise, an uncolored node is retrieved from another node as illustrated with the arrows.

Discretizing retrieval costs. Let $r_{max} = \max_{e \in E} \{r_e\}$. The possible total retrieval cost ρ is within range $\{0, 1, \dots, n^2 r_{max}\}$. To make the DP tractable, we partition this range further and define *approximated retrieval cost* $r'_{u,v}$ for edge $(u, v) \in E$ as follows:

$$r'_{u,v} = \lceil \frac{r_{u,v}}{l} \rceil \quad \text{where } l = \frac{n^2 r_{max}}{T(\epsilon)}, \quad T(\epsilon) = \frac{n^4}{\epsilon},$$

and $T(\epsilon)$ is the number of “ticks” we want to partition the retrieval range into. The choice for $T(\epsilon)$ will be specified in the proof for Theorem 4.2. We will work with r' in the rest of the subsection. However, by an abuse of notation, we still use r for discretized retrieval for the ease of representation.

Base case. For a leaf v , we let $DP[v][1][0][0] = s_v$.

Recurrence step. On each iteration at node v , we take the minimum over all possible situations as illustrated in Figure 4. For the DP recurrence, we want to restrict the possible solutions illustrated in Figure 4 on v to the corresponding compatible partial solutions on c_1 and c_2 . The recurrence relation for all cases is given in our full paper. Here, we select representative cases and explain the details of calculation below:

1) **Dealing with dependency:** When we decide to retrieve any child from v , like illustrated by case 4 of Figure 4, the children c_1, c_2 along with all their dependencies now become dependencies of v . The minimum storage cost in case 4 (given $v, k, \gamma = 0, \rho$) is:

$$S_4 = s_v + s_{v,c_1} + s_{v,c_2} - s_{c_1} - s_{c_2} \quad (1)$$

$$+ \min_{\substack{\rho_1 + \rho_2 = \rho \\ k_1 + k_2 = k-1}} \left\{ DP[c_1][k_1][0][\rho_1 - k_1 r_{v,c_1}] \right\} \quad (2a)$$

$$+ DP[c_2][k_2][0][\rho_2 - k_2 r_{v,c_2}] \quad (2b)$$

In Eq. (2a), v is required to have dependency number k and root retrieval 0. For each $k_1 + k_2 = k-1$, we must go through subproblems where c_1 has dependency number k_1 and c_2 has that of k_2 .

Also in Eq. (2a), the choice of ρ_1, ρ_2 determines how we are allocating retrieval costs budget ρ to c_1 and c_2 respectively. Specifically in Eq. (2a) and Eq. (2b), the total retrieval cost allocated to subproblem on $G_{[c_1]}$ is $\rho_1 - k_1 \cdot r_{v,c_1}$ since an extra $k_1 \cdot r_{v,c_1}$ cost is incurred by the edge (v, c_1) , as it is used k_1 times by all versions depending on c_1 . Similar applies to the subproblem on $G_{[c_2]}$.

Next, we highlight the idea of “invisible” dependency here: for case 2 on $G_{[v]}$, note the diff (v, c_1) and (v, c_2) was not available in any previous recurrence, since v has just been introduced. Therefore, the compatible solution for the subproblems on $G_{[c_1]}$ and $G_{[c_2]}$ have to materialize nodes c_1 and c_2 to ensure they can be retrieved. This explains the $-s_{c_1} - s_{c_2}$ terms in Eq. (1), as when calculating the retrieval cost on $G_{[v]}$, we need to subtract back the additional materialization cost.

When generalizing the DP onto graphs with bounded tree-width, similarly, restriction of a global solution does not always result in a feasible partial solution because of the existence of dependencies invisible to the subproblems. We resolve them using similarly ideas, as discussed in Section IV-C.

2) **Dealing with retrieval:** In contrast with dependencies, this refers to the case where v is retrieved from one of its children. We take case 5 as an example: given $v, k = 0, \gamma, \rho$,

$$\begin{aligned} S_5 &= s_{c_1,v} \\ &+ \min_{\rho_1 \leq \rho} \left\{ \min_{k_1} \{ DP[c_1][k_1][\gamma - r_{c_1,v}][\rho_1 - \gamma] \} \right. \\ &\quad \left. + \min_{k_2, \gamma'} \{ DP[c_2][k_2][\gamma'][\rho - \rho_1] \} \right\} \end{aligned}$$

We allocate the retrieval cost similar to case 2. We will care less about the dependency number, over which we will take minimum. The retrieval cost for c_1 now has to be $\gamma - r_{c_1,v}$ since v has to be retrieved from c_1 . Note importantly that now we are counting the retrieval cost for v in ρ_1 , and so the retrieval cost remaining for the left subproblem now is $\rho_1 - \gamma$. Notice that since only one way of retrieving v will be stored, this retrieval cost will not be over-counted in any cases.

Similarly, we take minimum on all other unused parameters to get the best storage for case 5.

3) **Combining the ideas:** We take case 8 as an example where both retrieval and dependencies are involved. In case 8, v is retrieved from child c_1 (retrieval), and child c_2 is retrieved from v (dependency). Given v, k, γ, ρ , we claim that:

$$\begin{aligned} S_8 &= s_{c_1,v} + s_{v,c_2} - s_{c_2} \\ &+ \min_{\rho_1 + \rho_2 = \rho} \left\{ \min_{k'} \{ DP[c_1][k'][\gamma - r_{c_1,v}][\rho_1 - \gamma] \} \right. \\ &\quad \left. + DP[c_2][k-1][0][\rho_2 - (k-1) \cdot (r_2 + \gamma)] \} \right\} \end{aligned}$$

Note that the c_1 side is identical to that for case 5. In combining both dependency and retrieval cases, there is slight adjustment in the dependency side: since v now might also depend on nodes further down c_1 side, the total extra retrieval cost created by adding edge (v, c_2) becomes $(k-1) \cdot (r_2 + \gamma)$ instead of $(k-1) \cdot (r_2)$.

Output. Finally, with storage constraint \mathcal{S} and root of the tree v_{root} , we output the configuration that outputs the minimum ρ which achieves the following

$$\exists k \leq n, \gamma \in \mathbb{N} \quad \text{s.t.} \quad DP[v_{root}][k][\gamma][\rho] \leq \mathcal{S}$$

We shall formally state and prove the FPTAS result below.

Lemma 4.1: The DP algorithm outputs a configuration with total retrieval cost at most $\text{OPT} + \epsilon r_{max}$ in $\text{poly}(n, 1/\epsilon)$ time.

Proof: By setting $T(\epsilon) = \frac{n^4}{\epsilon}$, we have $l = \frac{n^2 r_{max}}{T(\epsilon)} = \frac{\epsilon r_{max}}{n^2}$. Note that we can get an approximation of the original retrieval costs by multiplying each r'_e with l . This creates an estimation error of at most l on each edge. Note further that in the optimal solution, at most n^2 edges are materialized, so if ρ^* is the minimal discretized total retrieval cost, we have

$$\text{total retrieval of output} \leq l\rho^* \leq \text{OPT} + n^2l \leq \text{OPT} + \epsilon r_{max}.$$

■

Now we prove the main theorem of this subsection:

Theorem 4.2: For all $\epsilon > 0$, there is a $(1+\epsilon)$ -approximation algorithm for **MINSUM RETRIEVAL** on bidirectional trees that runs in $\text{poly}(n, \frac{1}{\epsilon})$ time.

Proof: Given parameter ϵ , we can use the DP algorithm as a black box and iterate the following for up to n times:

- 1) Run the DP for the given ϵ on the current graph. Record the output.
- 2) Let (u, v) be the most retrieval cost-heavy edge. We now set $r_{(u,v)} = 0$ and $s_{(u,v)} = s_v$. If the new graph is infeasible for the given storage constraint, or if all edges have already been modified, exit the loop.

At the end, we output the best out of all recorded outputs. This improves the previous bound when $r_{max} > \text{OPT}$: at some point we will eventually have $r_{max} \leq \text{OPT}$, which means the output configuration, if mapped back to the original input, is a feasible $(1+\epsilon)$ -approximation. ■

B. Treewidth-Related Definitions

We now consider a more general class of version graphs: any G whose *underlying undirected graph* G_0 has treewidth bounded by some constant k .

Definition 4.3 (Tree Decomposition [60]): A tree decomposition of undirected $G_0 = (V_0, E_0)$ is a tree $T = (V_T, E_T)$, where each $z \in V_T$ is associated with a subset (“bag”) S_z of V_0 . The bags must satisfy the following conditions:

- 1) $\bigcup_{z \in V_T} S_z = V_0$;
- 2) For each $v \in V_0$, the bags containing v induce a connected subtree of T ;
- 3) For each $(u, v) \in E_0$, there exists $z \in V_T$ such that S_z contains both u and v .

The *width* of a tree decomposition $T = (V_T, E_T)$ is $\max_{z \in V_T} |S_z| - 1$. The *treewidth* of G_0 is the minimum width over all tree decompositions of G_0 .

It follows that undirected forests have treewidth 1. We further note that there is also a notion of directed treewidth [61], but it is not suitable for our purpose.

We will WLOG assume a special kind of decomposition:

Definition 4.4 (Nice Tree Decomposition [62]): A nice tree decomposition is a tree decomposition with a designated root, where each node z is one of the following types:

- 1) A **leaf**, which has no children and whose bag has size 1;
 - 2) A **forget node**, which has one children c , and $S_z \subset S_c$ and $|S_c| = |S_z| + 1$.
 - 3) An **introduce node**, which has one children c , and $S_z \supset S_c$ and $|S_c| + 1 = |S_z|$.
 - 4) A **join**, which has children c_1, c_2 , and $S_z = S_{c_1} = S_{c_2}$.
- Given a bound k on the treewidth, there are multiple algorithms for calculating a tree decomposition of width k [63, 64, 65], or an approximation of k [66, 67, 68, 69].

For our case, the algorithm by Bodlaender [64] can be used to compute a tree decomposition in time $2^{\mathcal{O}(k^3)} \cdot O(n)$, which is linear if the treewidth k is constant. Given a tree decomposition, we can in $O(|V_0|)$ time find a nice tree decomposition of the same width with $O(k|V_0|)$ nodes [62].

C. Generalized Dynamic Programming

Here we outline the DP for MSR on graphs whose underlying undirected graph G_0 has treewidth at most k .

1) DP States: Similar to the warm-up, we will do the DP bottom-up on each $z \in V_T$ in the nice tree decomposition T . Before proceeding, let us define some additional notations. For any bag $z \in V_T$, let $T_{[z]}$ be the induced subtree of T rooted at z . We define $V_{[z]} = \bigcup_{z' \in V(T_{[z]})} S_{z'}$ be the set of vertices in the bags of $T_{[z]}$, including S_z . Following that, $G_{[z]}$ is the induced subgraph of G by vertices $V_{[z]}$.

We now define the *DP states*. At a high level, each state describes some number of *partial solutions* on the subgraph induced by $V_{[z]}, G_{[z]}$. When building a complete solution on G from the partial solutions, the state variables should give us *all* the information we need.

Each DP state on $z \in V_T$ consists of a tuple of functions

$$\mathcal{T}_z = (\text{Par}_z, \text{Dep}_z, \text{Ret}_z, \text{Anc}_z)$$

and a natural number ρ_z :

- (i) *Parent function* $\text{Par}_z : S_z \mapsto V_{[z]}$ describing the partial solution on $G_{[z]}$, restricted on S_z . If $\text{Par}_z(v) \neq v$ then v will be retrieved through the edge $(\text{Par}_z(v), v)$. If $\text{Par}_z(v) = v$ then v will be materialized.
- (ii) *Dependency function* $\text{Dep}_z : S_z \mapsto [n]$. Similar to the dependency parameter in the warm-up, $\text{Dep}_z(v)$ counts the number of nodes in $V_{[z]}$ retrieved through v .
- (iii) *Retrieval cost function* $\text{Ret}_z : S_z \mapsto \{0, \dots, nr_{max}\}$. Similar to the root retrieval parameter in the warm-up, $\text{Ret}_z(v)$ denotes the retrieval cost of version v in the partial solution on $G_{[z]}$.
- (iv) *Ancestor function* $\text{Anc}_z : S_z \mapsto 2^{S_z}$. If $u \in \text{Anc}_z(v)$, then u is retrieved in order to retrieve v in this partial solution, i.e., v is dependent on u . We need this extra information to avoid directed cycles.
- (v) ρ_z , the total retrieval cost of the subproblem according to the partial solution. Similar to its counterpart in the warm-up, all retrieval costs would be discretized by the same technique that makes the approximation an FPTAS.

A feasible state on $z \in V_T$ is a pair (\mathcal{T}_z, ρ_z) which correctly describes some partial solution on $G_{[z]}$ whose retrieval cost is exactly ρ_z . Each state is further associated with a storage value $\sigma(\mathcal{T}_z, \rho_z) \in \mathbb{Z}^+$, indicating the minimum storage needed to achieve the state (\mathcal{T}_z, ρ_z) on $G_{[z]}$.

We are now ready to describe how to compute the states.

2) **Recurrence on leaves:** For each leaf $z \in V_T$, the only feasible partial solution is to materialize the only vertex v in the leaf bag. We can easily calculate its state and storage cost.

3) **Recurrence on forget nodes:** This is also easy: for a forget node z with child c , we have $G_{[z]} = G_{[c]}$, and hence the states on z are simply the restrictions of states on c .

4) **Recurrence on introduce nodes:** At introduce node z with child c , we have $S_z = S_c \cup \{v_0\}$ for some ‘‘introduced’’ v_0 . Each feasible state (\mathcal{T}_z, ρ_z) on z must correspond to some state (\mathcal{T}_c, ρ_c) on c , which we can calculate as follows:

We first initialize \mathcal{T}_c to be the respective functions in \mathcal{T}_z restricted on S_c . For instance, $\text{Par}_c = \text{Par}_z|_{S_c}$.

If v_0 is retrieved through $u \in S_c$ according to \mathcal{T}_z ($\text{Par}_z(v_0) = u$), then we remove the dependencies related to v_0 and the retrieval cost incurred on edge (u, v_0) . Specifically:

- (i) Decrease the value of Dep_c by 1 on all vertices in $\text{Anc}_z(u)$.
- (ii) Decrease ρ_c by $\text{Dep}_z(v_0) \cdot \text{Ret}_z(v_0)$.
- (iii) Remove $\text{Anc}_z(u)$ from the ancestor functions of all descendants of z .

If v_0 has some child w according to \mathcal{T}_z (namely, $\text{Par}_z(w) = v_0$), then we reverse the *uprooting* process in the warm-up, such that vertex w , which was not a root in \mathcal{T}_z , is now a root in \mathcal{T}_c . Specifically:

- (i) Let $\text{Par}_c(w) = w$.
- (ii) Remove v_0 from the ancestor function of w and all its descendants.
- (iii) Decrease the retrieval cost function of w and its descendants by $\text{Ret}_z(w)$.
- (iv) Decrease ρ_c by $\text{Ret}_z(w) \times \text{Dep}_z(w)$.

Since v could have multiple children, the last procedure is potentially repeated multiple times.

5) **Recurrence on joins:** Suppose we are at a join z with children a, b , where $S_z = S_a = S_b$. On a high level, for each state (\mathcal{T}_z, ρ_z) on $G_{[z]}$, we want to find all pairs of states (\mathcal{T}_a, ρ_a) and (\mathcal{T}_b, ρ_b) such that the partial solutions they describe can combine into a partial solution on $G_{[z]}$, as described by (\mathcal{T}_z, ρ_z) . The pseudocode of the following functions can be found in the full version.

Compatibility. The algorithm COMPATIBILITY decides whether $\mathcal{T}_a, \mathcal{T}_b$ are indeed the ‘‘restrictions’’ of \mathcal{T}_z on $G_{[a]}$ and $G_{[b]}$ respectively. If the algorithm returns true, we would later proceed to calculate the correct value of $\rho_a + \rho_b$, based on this particular restriction.

Resolving external retrieval. COMPATIBILITY first deals with the vertices that are retrieved from outside S_z . For example, each $v \in S_z$ retrieved from $V_{[a]} \setminus S_z$, like the yellow node in (c) of Figure 5, is instead materialized from \mathcal{T}_b 's perspective. To check whether \mathcal{T}_a and \mathcal{T}_b resolve all such cases correctly, we define subroutine EXTERNAL-RETRIEVAL to loop through S_z topologically and calculate the correct $\text{Par}, \text{Ret}, \text{Anc}$ functions for both \mathcal{T}_a and \mathcal{T}_b .

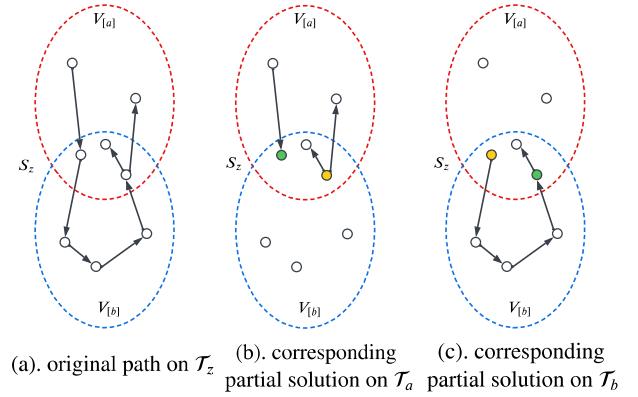


Fig. 5: Illustration for compatibility. Figure (b) and (c) show a pair of compatible configurations on \mathcal{T}_a and \mathcal{T}_b with the configuration on \mathcal{T}_z in (a). The configurations of yellow nodes and green nodes are analyzed in EXTERNAL-RETRIEVAL and EXTERNAL-DEPENDENCY respectively.

Resolving external dependency. The next step in COMPATIBILITY is to check whether the functions $\text{Dep}_a, \text{Dep}_b$ are compatible with Dep_z . Specifically, nodes in S_z could have *external dependencies* in $V_{[a]} \setminus S_z$ and $V_{[b]} \setminus S_z$, such as the green nodes in Figure 5 and Figure 6. The specific definition of $\text{ExtDep}_a(v)$ is the number of descendants that v have outside S_z , to whom v is the *closest* ancestor in S_z , according to \mathcal{T}_a . To see an example, note that only four red nodes are counted towards $\text{ExtDep}_a(A)$ in Figure 6. The functions ExtDep_b and ExtDep_z are defined similarly according to \mathcal{T}_b and \mathcal{T}_z .

We note that $\text{ExtDep}_a(v) + \text{ExtDep}_b(v) = \text{ExtDep}_z(v)$ for all $v \in S_z$ in order for $(\mathcal{T}_a, \mathcal{T}_b)$ to be compatible with \mathcal{T}_z . To check this, we call EXTERNAL-DEPENDENCY on $\mathcal{T}_z, \mathcal{T}_a, \mathcal{T}_b$ as a subroutine of COMPATIBILITY. We note that this is similar to distributing the dependency number k to the two children in case 4 of Figure 4.

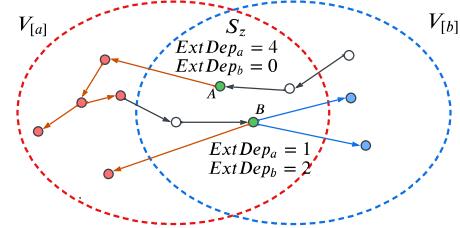


Fig. 6: Illustration for external dependency. Green nodes A and B both have non-zero external dependency, as labeled in the figure.

Calculating ρ . Given that $(\mathcal{T}_a, \mathcal{T}_b)$ are compatible with \mathcal{T}_z , we want to find the objective, $\sigma(\mathcal{T}_z, \rho_z)$, with the recurrence relation involving $\sigma(\mathcal{T}_a, \rho_a) + \sigma(\mathcal{T}_b, \rho_b)$ for suitable ρ_a and ρ_b . However, we cannot simply take $\rho_a + \rho_b = \rho_z$ due to the complicated procedure of combining $\mathcal{T}_a, \mathcal{T}_b$ into \mathcal{T}_z . We thus implement DISTRIBUTE RETRIEVAL to calculate ρ_Δ such that $\rho_a + \rho_b = \rho_z - \rho_\Delta$ and then iterate through all such ρ_a and ρ_b .

Recurrence relation. Finally, we have all we need for the recurrence relation. For each feasible (\mathcal{T}_z, ρ_z) , we take

$$\sigma(\mathcal{T}_z, \rho_z) = \min \{ \sigma(\mathcal{T}_a, \rho_a) + \sigma(\mathcal{T}_b, \rho_b) - \text{uproot} - \text{overcount} \}$$

where the minimum is taken over all $(\mathcal{T}_a, \mathcal{T}_b)$ that are compatible with \mathcal{T}_z and all $\rho_a + \rho_b = \rho_z - \rho_\Delta$, and where

$$\begin{aligned} \text{uproot} &= \sum_{v \in U_a} (s_v - s_{\text{Par}_z(v), v}) + \sum_{v \in U_b} (s_v - s_{\text{Par}_z(v), v}), \text{ and} \\ \text{overcount} &= \sum_{v \in S_a \cap S_b} s_{\text{Par}_z(v), v}. \end{aligned}$$

If k is constant, then the recurrence relation takes $\text{poly}(n)$ time. This is because there are $\text{poly}(n)$ many possible choices of \mathcal{T} and ρ on a, b, z , and it takes $\text{poly}(n)$ steps to check the compatibility of $(\mathcal{T}_a, \mathcal{T}_b)$ with \mathcal{T}_z and compute ρ_Δ .

Output The minimum retrieval cost of a global solution is just $\min\{\rho_z : \exists \mathcal{T}_z, \sigma(\mathcal{T}_z, \rho_z) \leq \mathcal{S}\}$ over all feasible (\mathcal{T}_z, ρ_z) , where z is the designated root of the nice tree decomposition.

We conclude this section with the following theorem.

Theorem 4.5: For a constant $k \geq 1$, on the set of graphs whose undelying undirected graph has treewidth at most k , MINSUM RETRIEVAL admits an FPTAS.

To see that our algorithm above is an FPTAS for MSR, the proof is almost identical to the proof of Theorem 4.2 (Section IV-A3) once we note that the number of partial solutions on each z is $\text{poly}(n)$.

An FPTAS for MMR arises from a similar procedure. When the objective becomes the maximum retrieval cost, we can use ρ_z to represent the maximum retrieval cost in the partial solution. We then modify $\text{Dep}_z(v)$ to represent the highest retrieval cost among all the nodes that are dependent on v . The recurrence relation is also changed accordingly. One can note that, like before, the new tuple \mathcal{T}_z contains all the information we need for a subsolution on $G_{[z]}$.

The same algorithms extend to $(1, 1 + \epsilon)$ bi-criteria approximation algorithms for BSR and BMR naturally, as the objective and constraint are reversed.

V. HEURISTICS ON MSR AND BMR

In this section, we propose three new heuristics that are inspired by empirical observations and theoretical results.

A. LMG-All: Improvement over LMG

Here we provide a brief description of the greedy heuristic LMG [1] and the improved LMG-All. We refer to the full paper for pseudocodes and formal definitions.

On a high level, LMG does the following:

- 1) Find a configuration that minimizes total storage cost.
- 2) Let V_{active} be the set of vertices not yet materialized, and, if materialized, does not cause the configuration to exceed storage constraint \mathcal{S} . If $V_{\text{active}} = \emptyset$, output the current configuration.
- 3) For each $v \in V_{\text{active}}$, calculate the cost and benefit of materializing v : storage cost increases by some $S(v)$, but total retrieval cost decreases by some $R(v)$.
- 4) From all such v , materialize the one that maximizes $\frac{R(v)}{S(v)}$. Go to step 2 and repeat.

Our improved heuristic LMG-All enlarges the scope of the search on each greedy step. Instead of searching for the most efficient version to *materialize*, we explore the payoff of *modifying any single edge*:

- 1) Find a configuration that minimizes total storage cost.
- 2) Let $\text{Par}(v)$ be the current parent of v on retrieval path. In addition to V_{active} , Define edge set E_{active} to be the edges that (a) does not cause the configuration to exceed storage constraint \mathcal{S} , and (b) does not form cycles, if $(u, v) \in E_{\text{active}}$ were to replace $(\text{Par}(v), v)$ in the current configuration. If $V_{\text{active}} = E_{\text{active}} = \emptyset$, output the current configuration.
- 3) Calculate cost and benefit of each $v \in V_{\text{active}}$ and $e \in E_{\text{active}}$. Materialize or store the most cost-effective node or edge. Go to step 2 and repeat.

While LMG-All considers more edges than LMG, it is not obvious that LMG-All always provides a better solution, due to its greedy nature.

B. DP on extracted bi-directional trees

We propose DP heuristics on both MSR and BMR, as inspired by algorithms in Sections III and IV. To ensure a reasonable running time, we only run the DP's on bi-directional trees (namely, with treewidth 1) extracted from our general input graphs, with the steps below:

- 1) Calculate a minimum spanning arborescence A of the graph G rooted at the first commit v_1 . We use the sum of retrieval and storage costs as weight.
- 2) Generate a bidirectional tree G' from A . Namely, we have $(u, v), (v, u) \in E(G')$ for each edge $(u, v) \in E(A)$.
- 3) Run the proposed DP for MSR and BMR on directed trees (see Section IV-A and Section III) with input G' .

In addition, we also implement the following modifications for MSR to further speed up the algorithm:

- 1) Total *storage* cost is discretized instead of retrieval cost, since the former generally has a smaller range.
- 2) Geometric discretization is used instead of linear discretization.
- 3) A pruning step is added, where the DP variable discards all subproblem solutions whose storage cost exceeds some bound.

All three original features are necessary in the proof for our theoretical results, but in practice, the modified implementations show comparable results but significantly improves the running time.

VI. EXPERIMENTS FOR MSR AND BMR

In this section, we discuss the experimental setup and results for empirical validation of the algorithms' performance, as compared to previous best-performing heuristic: LMG for MSR, and MP for BMR.¹¹

In all figures, the vertical axis (objective and run time) is presented in *logarithmic scale*. Run time is measured in *milliseconds*.

¹¹Our code can be found at <https://github.com/Sooooffia/Graph-Versioning>.

A. Datasets and Construction of Graphs

We use real-world GitHub repositories of varying sizes as datasets, from which we construct version graphs. Each commit corresponds to a node with its storage cost equal to its size in bytes. Between each pair of parent and child commits, we construct bidirectional edges. The storage and retrieval costs of the edges are calculated, in bytes, based on the actions (such as addition, deletion, and modification of files) required to change one version to the other in the direction of the edge. We use simple diff to calculate the deltas, hence the storage and retrieval costs are proportional to each other. Graphs generated this way are called “**natural graphs**” in the rest of the section.

In addition, we also aim to test (1) the cases where the retrieval and storage costs of an edge can greatly differ from each other, and (2) the effect of tree-like shapes of graphs on the performance of algorithms. Therefore, we also conduct experiments on modified graphs in the following two ways:

Random compression. We simulate compression of data by scaling storage cost with a random factor between 0.3 and 1, and increasing the retrieval cost by 20% (to simulate decompression).

ER construction. Instead of the naturally constructing edges between each pair of parent and child commits, we construct the edges as in an Erdős-Rényi random graph: between each pair (u, v) of versions, with probability p both deltas (u, v) and (v, u) are constructed, and with probability $1 - p$ neither are constructed. The resulting graphs are much less tree-like.¹²

Dataset	#nodes	#edges	avg. storage cost of edge
datasharing	29	74	395
styleguide	493	1250	8659
996.ICU	3189	9210	337038
freeCodeCamp	31270	2.5×10^7	14800
LeetCode	246	628	1.2×10^7
LeetCode 0.05	246	3032	1.0×10^8
LeetCode 0.2	246	11932	1.0×10^8
LeetCode 1	246	60270	1.0×10^8

TABLE IV: Natural and ER graphs overview.

B. Results in MSR

Figure 7, Figure 8, and Figure 9 demonstrate the performance of the three MSR algorithms on natural graphs, compressed natural graphs, and compressed ER graphs. The running times for the algorithms are shown in Figure 8 and Figure 9. Since run time for most non-ER graphs exhibit similar trends, many are omitted here due to space constraint. Also note that, since DP-MSR generates all data points in a single run, its running time is shown as a horizontal line over the full range for storage constraint.

We run DP-MSR with $\epsilon = 0.05$ on most graphs, except $\epsilon = 0.1$ for freeCodeCamp (for the feasibility of run time). The pruning value for DP variables is at twice the minimum storage for uncompressed graphs, and ten times the minimum storage for randomly compressed graphs.

¹²ER graphs have treewidth $\Theta(n)$ with high probability if the number of edges per vertex is greater than a small constant [70].

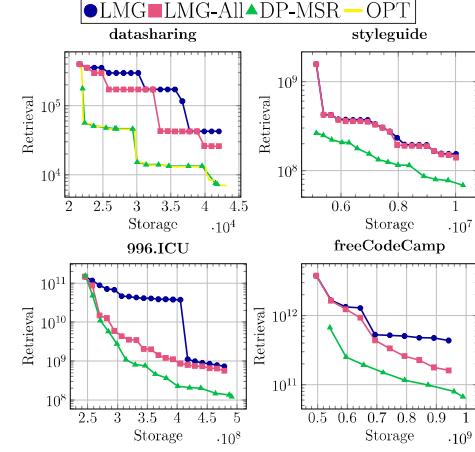


Fig. 7: Performance of MSR algorithms on natural graphs. OPT is obtained by solving an integer linear program (ILP) using Gurobi [71]. ILP takes too long to finish on all graphs except datasharing.

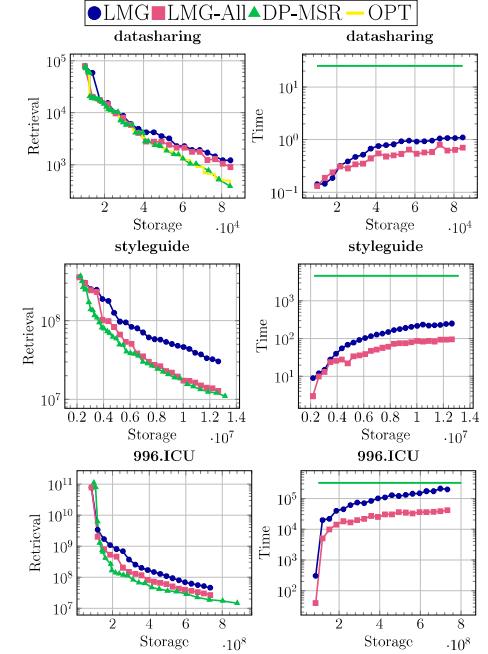


Fig. 8: Performance and run time of MSR algorithms on compressed graphs.

Performance analysis. On most graphs, DP-MSR outperforms LMG-All, which in turn outperforms LMG. This is especially clear on natural version graphs, where DP-MSR solutions are near 1000 times better than LMG solutions on 996.ICU in Figure 7. On datasharing, DP-MSR almost perfectly matches the optimal solution for all constraint ranges.

On naturally constructed graphs (Figure 7), LMG-All often has comparable performance with LMG when storage constraint is low. This is possibly because both algorithms can only iterate a few times when the storage constraint is almost tight. DP-MSR, on the other hand, performs much better on natural graphs even for low storage constraint.

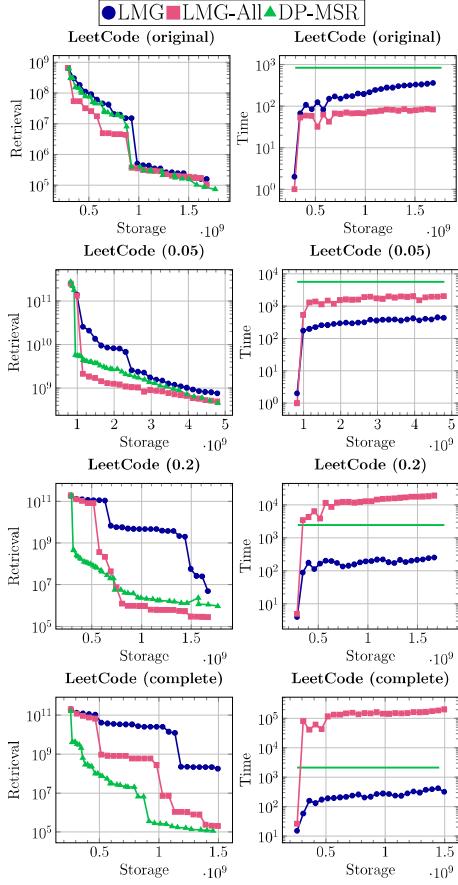


Fig. 9: Performance and run time of MSR algorithms on compressed ER graphs.

On graphs with random compression (Figure 8), the dominance of DP in performance over the other two algorithms become less significant. This is anticipated because of the fact that DP only runs on a subgraph of the input graph. Intuitively, most of the information is already contained in a minimum spanning tree when storage and retrieval costs are proportional. Otherwise, the dropped edges may be useful. (They could have large retrieval but small storage, and vice versa.)

Finally, LMG’s performance relative to our new algorithms is much worse on ER graphs. This may be due to the fact that LMG cannot look at non-auxiliary edges once the minimum arborescence is initialized, and hence losing most of the information brought by the extra edges. (Figure 9).

Run time analysis. For all natural graphs, we observe that LMG-All uses no more time than LMG (as shown in Figure 8). Moreover, LMG-All is significantly quicker than LMG on large natural graphs, which was unexpected considering that the two algorithms have almost identical structures in implementation. Possibly, this could be due to LMG making bigger, more expensive changes on each iteration (materializing a node with many dependencies, for instance) as compared to LMG-All.

As expected, though, LMG-All takes much more time than the other two algorithms on denser ER graphs (Figure 9), due to the large number of edges.

DP-MSR is often slower than LMG, except when ran on the natural construction of large graphs (Figure 8). However,

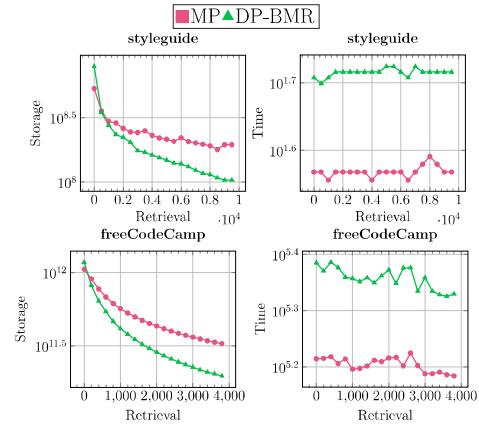


Fig. 10: Performance and run time of BMR algorithms on natural version graphs.

unlike LMG and LMG-All, the DP algorithm returns a whole spectrum of solutions at once, so it is difficult to make a direct comparison. We also note that the runtime of DP heavily depends on the choice of ϵ and the storage pruning bound. Hence, the user can trade-off the runtime with solution’s qualities by parameterize the algorithm with coarser configurations.

C. Results in BMR

As compared to MSR algorithms, the performance and run time of our BMR algorithms are much more predictable and stable. They exhibit similar trends across different ways of graph construction as mentioned in earlier sections - including the non-tree-like ER graphs, surprisingly.

Due to space limitation, we only present the results on natural graphs, as shown in Figure 10, to respectively illustrate their performance and run time.

Performance analysis. For every graph we tested, DP-BMR outperforms MP on most of the retrieval constraint ranges. As the retrieval constraint increases, the gap between MP and DP-BMR solution also increases. We also observe that DP-BMR performs worse than MP when the retrieval constraint is at zero. This is because the bidirectional tree have fewer edges than the original graph. (Recall that the same behavior happened for DP-MSR on compressed graphs)

We also note that, unlike MP, the objective value of DP-BMR solution monotonically decreases with respect to retrieval constraint. This is again expected since these are optimal solutions of the problem on the bidirectional tree.

Run time analysis. For all graphs, the runtimes of DP-BMR and MP are comparable within a constant factor. This is true with varying graph shapes and construction methods in all our experiments, and representative data is exhibited in Figure 10. Unlike LMG and LMG-All, their run times do not change much with varying constraint values.

Overall Evaluation For MSR, we recommend always using one of LMG-All and DP-MSR in place of LMG for practical use. On sparse graphs, LMG-All dominates LMG both in performance and run time. DP-MSR can also provide a frontier of better solutions in a reasonable amount of time, regardless of the input.

For BMR, DP-BMR usually outperforms MP, except when the retrieval constraint is close to zero. Therefore, we recommend using DP in most situations.

VII. CONCLUSION

In this paper, we developed fully polynomial time approximation algorithms for graphs with bounded treewidth. This often captures the typical manner in which edit operations are applied on versions. However, due to the high complexity of these algorithms, they are not yet practical to handle the size of real-life graphs. On the other hand, we extracted the idea behind this approach as well as previous LMG approach, and developed heuristics which significantly improved both the performance and run time in experiments, while potentially allowing for parallelization.

REFERENCES

- [1] S. Bhattacherjee, A. Chavan, S. Huang, A. Deshpande, and A. G. Parameswaran, “Principles of dataset versioning: Exploring the recreation/storage tradeoff,” *Proc. VLDB Endow.*, vol. 8, no. 12, pp. 1346–1357, 2015. [Online]. Available: <http://www.vldb.org/pvldb/vol8/p1346-bhattacherjee.pdf>
- [2] F. Nargesian, E. Zhu, R. J. Miller, K. Q. Pu, and P. C. Arocena, “Data lake management: Challenges and opportunities,” *Proc. VLDB Endow.*, vol. 12, no. 12, p. 1986–1989, aug 2019. [Online]. Available: <https://doi.org/10.14778/3352063.3352116>
- [3] K. Mukherjee, R. Shah, S. K. Saini, K. Singh, K. , H. Kesarwani, K. Barnwal, and A. Chauhan, “Towards optimizing storage costs on the cloud,” *IEEE 39th International Conference on Data Engineering (ICDE) (To Appear)*, 2023.
- [4] “Git,” <https://github.com/git/git>, 2005, last accessed: 13-Oct-22.
- [5] “Pachyderm,” <https://github.com/pachyderm/pachyderm>, 2016, last accessed: 13-Oct-22.
- [6] “DVC,” <https://github.com/iterative/dvc>, 2017, last accessed: 13-Oct-22.
- [7] “TerminusDB,” <https://github.com/terminusdb/terminusdb>, 2019, last accessed: 13-Oct-22.
- [8] “LakeFS,” <https://github.com/treeverse/lakeFS>, 2020, last accessed: 13-Oct-22.
- [9] “Dolt,” <https://github.com/dolthub/dolt>, 2019, last accessed: 13-Oct-22.
- [10] P. Roome, T. Feng, and S. Thakur, “Announcing the availability of data lineage with unity catalog,” <https://www.databricks.com/blog/2022/06/08/announcing-the-availability-of-data-lineage-with-unity-catalog.html>, 2022, last accessed: 13-Oct-22.
- [11] Y. Zhang, H. Liu, C. Jin, and Y. Guo, “Storage and recreation trade-off for multi-version data management,” in *Web and Big Data*, Y. Cai, Y. Ishikawa, and J. Xu, Eds. Cham: Springer International Publishing, 2018, pp. 394–409.
- [12] B. Derakhshan, A. Rezaei Mahdiraji, Z. Kaoudi, T. Rabl, and V. Markl, “Materialization and reuse optimizations for production data science pipelines,” ser. SIGMOD ’22. New York, NY, USA: Association for Computing Machinery, 2022, p. 1962–1976. [Online]. Available: <https://doi.org/10.1145/3514221.3526186>
- [13] S. Huang, L. Xu, J. Liu, A. J. Elmore, and A. Parameswaran, “ORPHEUSDB: bolt-on versioning for relational databases (extended version),” *The VLDB Journal*, vol. 29, no. 1, pp. 509–538, 2020.
- [14] N. N. Manne, S. Satpati, T. Malik, A. Bagchi, A. Gehani, and A. Chaudhary, “CHEX: multiversion replay with ordered checkpoints,” *Proceedings of the VLDB Endowment*, vol. 15, no. 6, pp. 1297–1310, 2022.
- [15] S. Wang, T. T. A. Dinh, Q. Lin, Z. Xie, M. Zhang, Q. Cai, G. Chen, B. C. Ooi, and P. Ruan, “Forkbase: An efficient storage engine for blockchain and forkable applications,” *Proc. VLDB Endow.*, vol. 11, no. 10, p. 1137–1150, jun 2018. [Online]. Available: <https://doi.org/10.14778/3231751.3231762>
- [16] M. E. Schule, L. Karnowski, J. Schmeißer, B. Kleiner, A. Kemper, and T. Neumann, “Versioning in Main-Memory Database Systems: From MusaeusDB to TardisDB,” *Proceedings of the 31st International Conference on Scientific and Statistical Database Management*, pp. 169–180, 2019.
- [17] A. Stivala, P. J. Stuckey, M. G. de la Banda, M. Hermenegildo, and A. Wirth, “Lock-free parallel dynamic programming,” *Journal of Parallel and Distributed Computing*, vol. 70, no. 8, pp. 839–848, 2010.
- [18] S. Khuller, B. Raghavachari, and N. E. Young, “Balancing minimum spanning trees and shortest-path trees,” *Algorithmica*, vol. 14, no. 4, pp. 305–321, 1995. [Online]. Available: <https://doi.org/10.1007/BF01294129>
- [19] G. Kortsarz and D. Peleg, “Approximating shallow-light trees,” in *Proceedings of the eighth annual ACM-SIAM symposium on Discrete algorithms*, 1997, pp. 103–110.
- [20] M. T. Hajiaghayi, G. Kortsarz, and M. R. Salavatipour, “Approximating buy-at-bulk and shallow-light k-steiner trees,” *Algorithmica*, vol. 53, no. 1, pp. 89–103, 2009.
- [21] B. Haeupler, D. E. Herskowitz, and G. Zuzic, “Tree embeddings for hop-constrained network design,” in *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, ser. STOC 2021. New York, NY, USA: Association for Computing Machinery, 2021, p. 356–369. [Online]. Available: <https://doi.org/10.1145/3406325.3451053>
- [22] M. V. Marathe, R. Ravi, R. Sundaram, S. Ravi, D. J. Rosenkrantz, and H. B. Hunt III, “Bicriteria network design problems,” *Journal of algorithms*, vol. 28, no. 1, pp. 142–171, 1998.
- [23] R. Ravi, “Rapid rumor ramification: approximating the minimum broadcast time,” in *Proceedings 35th Annual Symposium on Foundations of Computer Science*, 1994, pp. 202–213.
- [24] M. R. Khani and M. R. Salavatipour, “Improved approximations for buy-at-bulk and shallow-light k-steiner trees and (k,2)-subgraph,” *J. Comb. Optim.*, vol. 31, no. 2, p. 669–685, feb 2016. [Online]. Available: <https://doi.org/10.1007/s10878-014-9774-5>

- [25] M. Chimani and J. Spoerhase, “Network Design Problems with Bounded Distances via Shallow-Light Steiner Trees,” in *32nd International Symposium on Theoretical Aspects of Computer Science (STACS 2015)*, ser. Leibniz International Proceedings in Informatics (LIPIcs), E. W. Mayr and N. Ollinger, Eds., vol. 30. Dagstuhl, Germany: Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2015, pp. 238–248. [Online]. Available: <http://drops.dagstuhl.de/opus/volltexte/2015/4917>
- [26] R. Ghuge and V. Nagarajan, “Quasi-polynomial algorithms for submodular tree orienteering and directed network design problems,” *Mathematics of Operations Research*, vol. 47, no. 2, pp. 1612–1630, 2022.
- [27] P. Buneman, S. Khanna, and W.-C. Tan, “Data Provenance: Some Basic Issues,” *Lecture Notes in Computer Science*, pp. 87–93, 2000.
- [28] Y. L. Simmhan, B. Plale, D. Gannon *et al.*, “A survey of data provenance techniques.”
- [29] J. J. Hunt, K.-P. Vo, and W. F. Tichy, “Delta algorithms: an empirical analysis,” *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 7, no. 2, pp. 192–214, 1998.
- [30] R. C. Burns and D. D. E. Long, “In-place reconstruction of delta compressed files,” *Proceedings of the seventeenth annual ACM symposium on Principles of distributed computing - PODC '98*, pp. 267–275, 1998.
- [31] W. Xia, H. Jiang, D. Feng, L. Tian, M. Fu, and Y. Zhou, “Ddelta: A deduplication-inspired fast delta compression approach,” *Performance Evaluation*, vol. 79, pp. 258–272, 2014.
- [32] J. MacDonald, “File system support for delta compression,” Ph.D. dissertation, Masters thesis. Department of Electrical Engineering and Computer Science, University of California at Berkley, 2000.
- [33] D. T. N. M. T. Suel, “zdelta: An efficient delta compression tool,” 2002.
- [34] Y. Jayawardana and S. Jayarathna, “DFS: A Dataset File System for Data Discovering Users,” *2019 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, vol. 00, pp. 355–356, 2019.
- [35] R. C. Fernandez, Z. Abedjan, F. Koko, G. Yuan, S. Madden, and M. Stonebraker, “Aurum: A Data Discovery System,” *2018 IEEE 34th International Conference on Data Engineering (ICDE)*, pp. 1001–1012, 2018.
- [36] D. Brickley, M. Burgess, and N. Noy, “Google Dataset Search: Building a search engine for datasets in an open Web ecosystem,” *The World Wide Web Conference*, pp. 1365–1375, 2019.
- [37] A. Bogatu, A. A. A. Fernandes, N. W. Paton, and N. Konstantinou, “Dataset Discovery in Data Lakes,” *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, vol. 00, pp. 709–720, 2020.
- [38] J. Brown and N. Weber, “DSDB: An Open-Source System for Database Versioning & Curation,” *2021 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, vol. 00, pp. 299–307, 2021.
- [39] A. Chavan and A. Deshpande, “DEX: Query Execution in a Delta-based Storage System,” *Proceedings of the 2017 ACM International Conference on Management of Data*, pp. 171–186, 2017.
- [40] M. Maddox, D. Goehring, A. J. Elmore, S. Madden, A. Parameswaran, and A. Deshpande, “Decibel: The Relational Dataset Branching System,” *Proceedings of the VLDB Endowment. International Conference on Very Large Data Bases*, vol. 9, no. 9, pp. 624–635, 2016.
- [41] A. Bhardwaj, S. Bhattacherjee, A. Chavan, A. Deshpande, A. J. Elmore, S. Madden, and A. G. Parameswaran, “DataHub: Collaborative Data Science & Dataset Version Management at Scale,” *arXiv*, 2014.
- [42] A. Seering, P. Cudre-Mauroux, S. Madden, and M. Stonebraker, “Efficient Versioning for Scientific Array Databases,” *2012 IEEE 28th International Conference on Data Engineering*, vol. 1, pp. 1013–1024, 2012.
- [43] T. Ying, H. Chen, and H. Jin, “Pensieve: Skewness-aware version switching for efficient graph processing,” in *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD ’20. New York, NY, USA: Association for Computing Machinery, 2020, p. 699–713. [Online]. Available: <https://doi.org/10.1145/3318464.3380590>
- [44] U. Khurana and A. Deshpande, “Efficient Snapshot Retrieval over Historical Graph Data,” *arXiv*, 2012, graph database systems — stroing dynamic graphs so that a graph at a specific time can be queried. Vertices are marked with bits encoding information on which versions it belong to.
- [45] W. Nagel, “Subversion: not just for code anymore,” *Linux Journal*, vol. 2006, no. 143, p. 10, 2006.
- [46] A. Kougkas, H. Devarajan, and X.-H. Sun, “Hermes: a heterogeneous-aware multi-tiered distributed i/o buffering system,” in *Proceedings of the 27th International Symposium on High-Performance Parallel and Distributed Computing*, 2018, pp. 219–230.
- [47] H. Devarajan, A. Kougkas, L. Logan, and X.-H. Sun, “Hcompress: Hierarchical data compression for multi-tiered storage environments,” in *2020 IEEE IPDPS*. IEEE, 2020, pp. 557–566.
- [48] H. Devarajan, A. Kougkas, and X.-H. Sun, “Hfetch: Hierarchical data prefetching for scientific workflows in multi-tiered storage environments,” in *2020 IEEE IPDPS*. IEEE, 2020, pp. 62–72.
- [49] Y. Cheng, M. S. Iqbal, A. Gupta, and A. R. Butt, “Cast: Tiering storage for data analytics in the cloud,” in *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing*, 2015, pp. 45–56.
- [50] A. Erradi and Y. Mansouri, “Online cost optimization algorithms for tiered cloud storage services,” *Journal of Systems and Software*, vol. 160, p. 110457, 2020.
- [51] M. Liu, L. Pan, and S. Liu, “To transfer or not: An online cost optimization algorithm for using two-tier storage-as-a-service clouds,” *IEEE Access*, vol. 7, pp. 94 263–94 275, 2019.
- [52] ———, “Keep hot or go cold: A randomized online migration algorithm for cost optimization in staas clouds,” *IEEE Transactions on Network and Service Management*,

- vol. 18, no. 4, pp. 4563–4575, 2021.
- [53] W. Si, L. Pan, and S. Liu, “A cost-driven online auto-scaling algorithm for web applications in cloud environments,” *Knowledge-Based Systems*, vol. 244, p. 108523, 2022.
- [54] R. Kinoshita, S. Imamura, L. Vogel, S. Kazama, and E. Yoshida, “Cost-performance evaluation of heterogeneous tierless storage management in a public cloud,” in *2021 Ninth International Symposium on Computing and Networking (CANDAR)*. IEEE, 2021, pp. 121–126.
- [55] A. Anwar, Y. Cheng, A. Gupta, and A. R. Butt, “Taming the cloud object storage with mos,” in *Proceedings of the 10th Parallel Data Storage Workshop*, 2015, pp. 7–12.
- [56] ———, “Mos: Workload-aware elasticity for cloud object stores,” in *Proceedings of the 25th ACM International Symposium on High-Performance Parallel and Distributed Computing*, 2016, pp. 177–188.
- [57] L. Vogel, V. Leis, A. van Renen, T. Neumann, S. Imamura, and A. Kemper, “Mosaic: a budget-conscious storage engine for relational database systems,” *Proceedings of the VLDB Endowment*, vol. 13, no. 12, pp. 2662–2675, 2020.
- [58] R. Lasch, T. Legler, N. May, B. Scheirle, and K.-U. Sattler, “Cost modelling for optimal data placement in heterogeneous main memory,” *Proceedings of the VLDB Endowment*, vol. 15, no. 11, pp. 2867–2880, 2022.
- [59] R. Lasch, R. Schulze, T. Legler, and K.-U. Sattler, “Workload-driven placement of column-store data structures on dram and nvm,” in *Proceedings of the 17th International Workshop on Data Management on New Hardware (DaMoN 2021)*, 2021, pp. 1–8.
- [60] U. Bertelè and F. Brioschi, “On non-serial dynamic programming,” *Journal of Combinatorial Theory, Series A*, vol. 14, no. 2, pp. 137–148, 1973.
- [61] T. Johnson, N. Robertson, P. Seymour, and R. Thomas, “Directed tree-width,” *Journal of Combinatorial Theory. Series B*, vol. 82, no. 1, pp. 138–154, May 2001, funding Information: 1Partially supported by the NSF under Grant DMS-9701598. 2 Research partially supported by the DIMACS Center, Rutgers University, New Brunswick, NJ 08903. 3Partially supported by the NSF under Grant DMS-9401981. 4Partially supported by the ONR under Contact N00014-97-1-0512. 5Partially supported by the NSF under Grant DMS-9623031 and by the NSA under Contract MDA904-98-1-0517.
- [62] H. L. Bodlaender, “A partial k-arboretum of graphs with bounded treewidth,” *Theoretical Computer Science*, vol. 209, no. 1, pp. 1–45, 1998.
- [63] ———, “A linear time algorithm for finding tree-decompositions of small treewidth,” in *Proceedings of the Twenty-Fifth Annual ACM Symposium on Theory of Computing*, ser. STOC ’93. New York, NY, USA: Association for Computing Machinery, 1993, p. 226–234. [Online]. Available: <https://doi.org/10.1145/167088.167161>
- [64] S. Arnborg, D. G. Corneil, and A. Proskurowski, “Complexity of finding embeddings in a k-tree,” *SIAM Journal on Algebraic Discrete Methods*, vol. 8, no. 2, pp. 277–284, 1987. [Online]. Available: <https://doi.org/10.1137/0608024>
- [65] F. V. Fomin, I. Todinca, and Y. Villanger, “Large induced subgraphs via triangulations and cmso,” *SIAM Journal on Computing*, vol. 44, no. 1, pp. 54–87, 2015. [Online]. Available: <https://doi.org/10.1137/140964801>
- [66] M. Belbasi and M. Fürer, “Finding all leftmost separators of size $\leq k$,” in *Combinatorial Optimization and Applications: 15th International Conference, COCOA 2021, Tianjin, China, December 17–19, 2021, Proceedings*. Berlin, Heidelberg: Springer-Verlag, 2021, p. 273–287. [Online]. Available: https://doi.org/10.1007/978-3-030-92681-6_23
- [67] F. V. Fomin, D. Lokshtanov, S. Saurabh, M. Pilipczuk, and M. Wrochna, “Fully polynomial-time parameterized computations for graphs and matrices of low treewidth,” *ACM Trans. Algorithms*, vol. 14, no. 3, jun 2018. [Online]. Available: <https://doi.org/10.1145/3186898>
- [68] T. Korhonen, “A single-exponential time 2-approximation algorithm for treewidth,” in *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, 2022, pp. 184–192.
- [69] U. Feige, M. Hajiaghayi, and J. R. Lee, “Improved approximation algorithms for minimum-weight vertex separators,” in *Proceedings of the Thirty-Seventh Annual ACM Symposium on Theory of Computing*, ser. STOC ’05. New York, NY, USA: Association for Computing Machinery, 2005, p. 563–572. [Online]. Available: <https://doi.org/10.1145/1060590.1060674>
- [70] Y. Gao, “Treewidth of erdős–rényi random graphs, random intersection graphs, and scale-free random graphs,” *Discrete Applied Mathematics*, vol. 160, no. 4–5, pp. 566–578, 2012.
- [71] Gurobi Optimization, LLC, “Gurobi Optimizer Reference Manual,” 2022. [Online]. Available: <https://www.gurobi.com>