

Dynamic Path-Decomposed Tries

SHUNSUKE KANDA, RIKEN Center for Advanced Intelligence Project, Japan

DOMINIK KÖPPL, Kyushu University, Japan and Japan Society for Promotion of Science

YASUO TABEL, RIKEN Center for Advanced Intelligence Project, Japan

KAZUHIRO MORITA, Tokushima University, Japan

MASAO FUKETA, Tokushima University, Japan

A keyword dictionary is an associative array whose keys are strings. Recent applications handling massive keyword dictionaries in main memory have a need for a space-efficient implementation. When limited to static applications, there are a number of highly-compressed keyword dictionaries based on the advancements of practical succinct data structures. However, as most succinct data structures are only efficient in the static case, it is still difficult to implement a keyword dictionary that is *space efficient* and *dynamic*. In this article, we propose such a keyword dictionary. Our main idea is to embrace the path decomposition technique, which was proposed for constructing cache-friendly tries. To store the path-decomposed trie in small memory, we design data structures based on recent compact hash trie representations. Experiments on real-world datasets reveal that our dynamic keyword dictionary needs up to 68% less space than the existing smallest ones, while achieving a relevant space-time tradeoff.

1 INTRODUCTION

An associative array is called a *keyword dictionary* if its keys are strings. In this article, we study the problem to maintain a keyword dictionary in main memory efficiently. When storing words extracted from text collections written in natural or computer languages, the size of a keyword dictionary A is not of major concern. This is because, after carefully polishing the extracted strings with natural language processing tools like stemmers, the size of A grows sublinearly as $O(N^\beta)$ for some $\beta \approx 0.5$ over a text of N words due to Heaps' Law [10, 32]. However, as reported in [45], some natural language applications such as web search engines and machine translation systems need to handle large datasets that are not under Heaps' Law. Also, other recent applications as in Semantic Web graphs and in bioinformatics handle massive string databases with keyword dictionaries [45, 46]. Although common implementations like hash tables are fast, their memory consumption is a severe drawback in such scenarios. Here, a *space-efficient* implementation of the keyword dictionary is important. In this paper, we focus on the *practical* side of this problem.

In the static setting, omitting the insertion and deletion of keywords, a number of compressed keyword dictionaries have been developed for a decade, some of which we highlight in the following. We start with Martínez-Prieto et al. [45], who proposed and evaluated a number of compressed keyword dictionaries based on techniques like hashing, front-coding, full-text indexes, and tries. They demonstrated that their implementations use up to 5% space of the original dataset size, while also supporting searches of prefixes and substrings of the keywords. Subsequently, Grossi and Ottaviano [30] proposed a cache-friendly keyword dictionary through path decomposition of tries. Arz and Fischer [5] adapted the LZ78 compression to devise a keyword dictionary. Finally, Kanda et al. [39] proposed a keyword dictionary based on a compressed double-array trie. As we can see from these representations, space-efficient static keyword dictionaries have been well studied because of the advancements of practical (yet static) succinct data structures collected in well maintained libraries such as SDSL [24] and Succinct [29].

Authors' addresses: Shunsuke Kanda, shunsuke.kanda@riken.jp, RIKEN Center for Advanced Intelligence Project, Japan; Dominik Köppl, dominik.koepl@inf.kyushu-u.ac.jp, Kyushu University, Japan, Japan Society for Promotion of Science; Yasuo Tabei, yasuo.tabei@riken.jp, RIKEN Center for Advanced Intelligence Project, Japan; Kazuhiro Morita, kam@is.tokushima-u.ac.jp, Tokushima University, Japan; Masao Fuketa, fuketa@is.tokushima-u.ac.jp, Tokushima University, Japan.

Under the dynamic setting, however, only a few space-efficient keyword dictionaries have been realized, probably due to the implementation difficulty. Although HAT-trie [7] and Judy [11] are representative space-efficient dynamic implementations as demonstrated in previous experiments¹, they still waste memory by maintaining many pointers. The Cedar trie [64] is a space-efficient implementation embracing heavily 32-bit pointers to address memory, and therefore cannot be applied to massive datasets. Its implementation makes it hard to switch to 64-bit pointers, but we expect that doing so will increase its space consumption considerably. Although several practical dynamic succinct data structures [48–50] have been recently developed, modern dynamic keyword dictionaries are heavily based on pointers, consuming a large fraction of the entire space requirement. Nonetheless, there are some applications that need dynamic keyword dictionaries for massive datasets such as search engines [15, 16], RDF stores [46], or Web crawler [59]. Consequently, realizing a practical space-efficient dynamic keyword dictionaries is an important open challenge.

1.1 Space-Efficient Dynamic Tries

Common keyword dictionary implementations represent the keywords in a trie, supporting the retrieval of keywords with trie navigation operations. In this subsection, we summarize space-efficient dynamic tries.

Theoretical Discussion. We consider a dynamic trie with t nodes over an alphabet of size σ . Arroyuelo et al. [4] introduced succinct representations that require almost optimal $2t + t \log \sigma + o(t \log \sigma)$ bits of space, while supporting insertion and deletion of a leaf in $O(1)$ amortized time if $\sigma = O(\text{polylog}(t))$ and in $O(\log \sigma / \log \log \sigma)$ amortized time otherwise.² Jansson et al. [36] presented a dynamic trie representation that uses $O(t \log \sigma)$ bits of space, while supporting insertion and deletion of a leaf in $O(\log \log t)$ expected amortized time.

Hash Tries. On the practical side, Poyias et al. [49] proposed the *m-Bonsai* trie, a practical dynamic compact trie representation. It is a variant of the Bonsai trie [19] that represents the trie nodes as entries in a compact hash table. It takes $O(t \log \sigma)$ bits of space, while supporting update and some traversal operations in $O(1)$ expected time. Fischer and Köppl [22] presented and evaluated a number of dynamic tries for LZ78 [65] and LZW [60] factorization. They also proposed an efficient hash-based trie representation in a similar way to *m-Bonsai*, which is referred to as *FK-hash*.³ Although *FK-hash* uses $O(t \log \sigma + t \log t)$ bits of space, its update algorithm is simple and practically fast. However, we are not aware of any space-efficient approach using them as keyword dictionaries.

Compacted Tries. Another line of research focuses on limiting the space of the trie in relation to the number of keywords. Suppose that we want to maintain a set of n strings with a total length of N on a machine, where $\beta = \log_{\sigma} N$ characters fit into a single machine word w . In this setting, Belazzougui et al. [12] proposed the (dynamic) *z-fast* trie, which takes $N \log \sigma + O(n \log N)$ bits of space and supports retrieval, insertion and deletion of a string S in $O(|S|/\beta + \log |S| + \log \log \sigma)$ expected time. Takagi et al. [52] proposed the packed compact trie, which takes $N \log \sigma + O(nw)$ bits of space and supports the same operations in $O(|S|/\beta + \log \log N)$ expected time. Recently, Tsuruta et al. [58] developed a hybrid data structure of the *z-fast* trie and the packed compact trie, which also takes $N \log \sigma + O(nw)$ bits of space, but improves each of these operations to run in $O(|S|/\beta + \log \beta)$ expected time.

¹Such as <http://www.tkl.iis.u-tokyo.ac.jp/~ynaga/cedar/#perf> and <https://github.com/Tessil/hat-trie/blob/master/README.md#benchmark>.

²Throughout this paper, the base of the logarithm is 2, whenever not explicitly indicated.

³The representation is referred to as *hash* or *cht* in their paper [22]. To avoid confusion, we name it *FK-hash* by using the initial letters of the proposers, Fischer and Köppl.

1.2 Our Contribution

We propose a novel space-efficient dynamic keyword dictionary, called the *dynamic path-decomposed trie* (abbreviated as *DynPDT*). DynPDT is based on a trie formed by *path decomposition* [21]. The path decomposition is a trie transformation technique, which was proposed for constructing cache-friendly trie dictionaries. It was up to now utilized only in static applications [30, 35]. Here, we adapt this technique for the dynamic construction of DynPDT, which gives DynPDT two main advantages over other known keyword dictionaries.

- (1) The first is that the data structure is cache efficient because of the path decomposition. During the retrieval of a keyword, most parts of the keyword can be scanned in a cache-friendly manner without node-to-node traversals based on random accesses.
- (2) The second is that the path decomposition allows us to plug in any dynamic trie representation for the path-decomposed trie topology. For this job, we choose the hash-based trie representations m-Bonsai and FK-hash as these are fast and memory efficient in the setting when all trie nodes have to be represented explicitly (which is the case for the nodes of the path-decomposed trie).

Based on these advantages, DynPDT becomes a fast and space-efficient dynamic keyword dictionary.

From experiments using massive real-world datasets, we demonstrate that DynPDT is more space efficient compared to existing keyword dictionaries while achieving a relevant space-time tradeoff. For example, to construct a keyword dictionary from a large URI dataset of 13.8 GiB, DynPDT needs only 2.5 GiB of working space, while a HAT-trie and a Judy trie need 9.5 GiB and 7.8 GiB, respectively. The time performance is competitive in many cases thanks to the path decomposition. The source code of our implementation is available at <https://github.com/kampersanda/poplar-trie>.

1.3 Paper Structure

In Section 2, we introduce the keyword dictionary, and review the trie data structure and the path decomposition in our preliminaries. We introduce our new data structure DynPDT in Section 3. Subsequently, we present our DynPDT representations based on m-Bonsai and FK-hash in Sections 4 and 5, respectively. In Section 6, we provide our experimental results. Finally, we conclude the paper in Section 7.⁴

2 PRELIMINARIES

A *string* is a (finite) sequence of characters over a finite alphabet. Our strings always start at position 0. Given a string S of length n , $S[i, j]$ denotes the *substring* $S[i], S[i+1], \dots, S[j-1]$ for $0 \leq i \leq j \leq n$. Particularly, $S[0, j]$ is a *prefix* of S and $S[i, n]$ is a *suffix* of S . Let $|S| := n$ denote the length of S . The same notation is also applied to *arrays*. The cardinality of a set A is denoted by $|A|$.

Our model of computation is the transdichotomous word RAM model of word size $w = \Theta(\log N)$, where N is the total length of all keywords of a given problem, i.e., the size of the problem. We can read and process $\mathcal{O}(w)$ bits in constant time.

2.1 Keyword Dictionary

A *keyword* is a string over an alphabet \mathcal{A} that is terminated with a special character $\$ \notin \mathcal{A}$ at its end. In a *prefix-free* set of strings, no string is a prefix of another string. A set of keywords is

⁴A preliminary version of this work appeared in our conference paper [40] and the first author's Ph.D. thesis [37]. This paper contains the significant differences as follows: (1) a fast variant of m-Bonsai was incorporated in Section 4.1; (2) an efficient implementation of the bijective hash function in m-Bonsai was incorporated in Section 4.2; (3) a growing algorithm of m-Bonsai was presented in Section 4.3; (4) FK-hash was also considered in addition to m-Bonsai in Section 5; (5) the experimental results in Section 6 and all descriptions were significantly enhanced.

always prefix-free due to the character \$. A *keyword dictionary* is a dynamic associative array that maps a dynamic set of n keywords $\mathcal{S} = \{K_1, K_2, \dots, K_n\} \subset \mathcal{A}^*$ to values x_1, x_2, \dots, x_n , where x_i belongs to a finite set \mathcal{X} . It supports the retrieval, the insertion, and the deletion of keywords while maintaining the *key-value mapping*. In detail, it supports the following operations:

- $\text{lookup}(K)$ returns the value associated with the keyword K if $K \in \mathcal{S}$ or \perp otherwise.
- $\text{insert}(K, x)$ inserts the keyword K in \mathcal{S} , i.e., $\mathcal{S} \leftarrow \mathcal{S} \cup \{K\}$, and associates the value x with K .
- $\text{delete}(K)$ removes the keyword K from \mathcal{S} , i.e., $\mathcal{S} \leftarrow \mathcal{S} \setminus \{K\}$.

2.2 Tries

A trie [23, 41] is a rooted labeled tree $\mathcal{T}_{\mathcal{S}}$ representing a set of keywords \mathcal{S} . Each edge in $\mathcal{T}_{\mathcal{S}}$ is labeled by a character. All outgoing edges of a node are labeled with a distinct character. The label c of the edge (u, v) between a node v and its parent u is called the *branching character* of v . The parent u and branching character c uniquely determine v . Each keyword $K \in \mathcal{S}$ is represented by exactly one path from the root to a leaf u , i.e., the keyword K can be extracted by concatenating the edge labels on the path from the root to u . Since \mathcal{S} is prefix-free ($\$$ is a unique delimiter of each keyword), there is a 1-to-1 correlation between leaves and keywords.

Given a keyword K of length m , $\mathcal{T}_{\mathcal{S}}$ retrieves K by traversing nodes from the root to a leaf while matching the characters of K with the edge labels of the traversed path. In representations storing all trie nodes explicitly, we visit m nodes during this traversal. However, this traversal suffers poor locality of reference since it needs to access pointers usually addressing non-consecutive memory. In practice, this cache inefficiency is a critical bottleneck especially for long strings such as URLs. Grossi and Ottaviano [30] successfully solved this problem through *path decomposition* [21] in practice (but, in static settings).

2.3 Path Decomposition

The *path decomposition* [21] of a trie $\mathcal{T}_{\mathcal{S}}$ is a recursive procedure that first chooses an arbitrary root-to-leaf path π in $\mathcal{T}_{\mathcal{S}}$, then compactifies the path π to a single node, and subsequently repeats the procedure in each subtree hanging off the path π . As a result, $\mathcal{T}_{\mathcal{S}}$ is partitioned into a set of n node-to-leaf paths because there are n leaves in $\mathcal{T}_{\mathcal{S}}$. This decomposition produces the *path-decomposed* trie $\mathcal{T}_{\mathcal{S}}^c$, which is composed of n compactified nodes.

For explaining the properties of $\mathcal{T}_{\mathcal{S}}^c$, we call the concatenation of the labels of all edges of a node-to-leaf path π in $\mathcal{T}_{\mathcal{S}}$ the *path string* of π . The path strings of the compactified paths of $\mathcal{T}_{\mathcal{S}}$ are the node labels of $\mathcal{T}_{\mathcal{S}}^c$. In detail, each node u in $\mathcal{T}_{\mathcal{S}}^c$ is associated with a node-to-leaf path π of $\mathcal{T}_{\mathcal{S}}$ and is labeled by the path string of π , denoted by $L_u \in \mathcal{A}^*$. Each edge in $\mathcal{T}_{\mathcal{S}}^c$ is labeled by a pair consisting of a branching character and an integer, which are defined as follows (see also Figure 1): Take a node u in $\mathcal{T}_{\mathcal{S}}^c$ and one of its children v . Suppose that u and v are associated with the paths π_u and π_v in $\mathcal{T}_{\mathcal{S}}$, respectively, such that L_u and L_v are the path labels of π_u and π_v . The edge (u, v) has the label (b, i) if, in $\mathcal{T}_{\mathcal{S}}$, the first node on the path π_v is the node

- whose branching character is b , and
- whose parent is the i -th node⁵ visited on the path π_u .

The edge labels of $\mathcal{T}_{\mathcal{S}}^c$ are characters drawn from the alphabet $\mathcal{B} := \mathcal{A} \times \{0, 1, \dots, \Lambda - 1\}$, where Λ is the longest length of all node labels.

Example 2.1 (Path-Decomposed Trie). Figure 2 illustrates a root-to-leaf path π in $\mathcal{T}_{\mathcal{S}}$ and the corresponding root r in $\mathcal{T}_{\mathcal{S}}^c$ after compactifying π to r . The root r is labeled by the path string of π ,

⁵Throughout this paper, we start counting from zero.

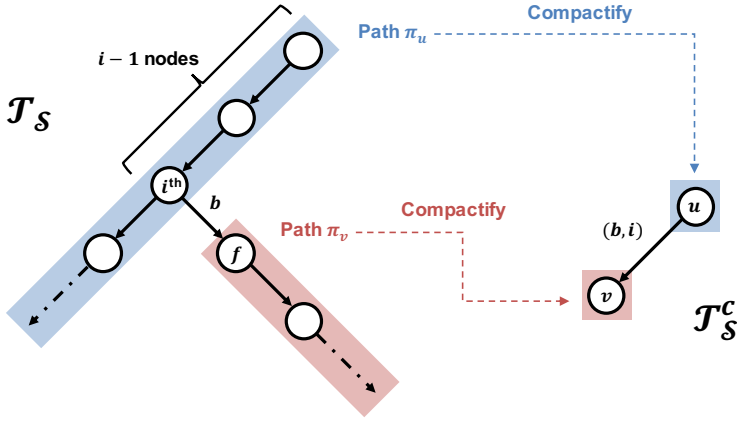


Fig. 1. Illustration of the path decomposition of the path π_v whose first node f is a child of the i -th node on the path π_u represented by node u in \mathcal{T}_S^c . Since the branching character of f is b , u and v are connected with an edge with label (b, i) .

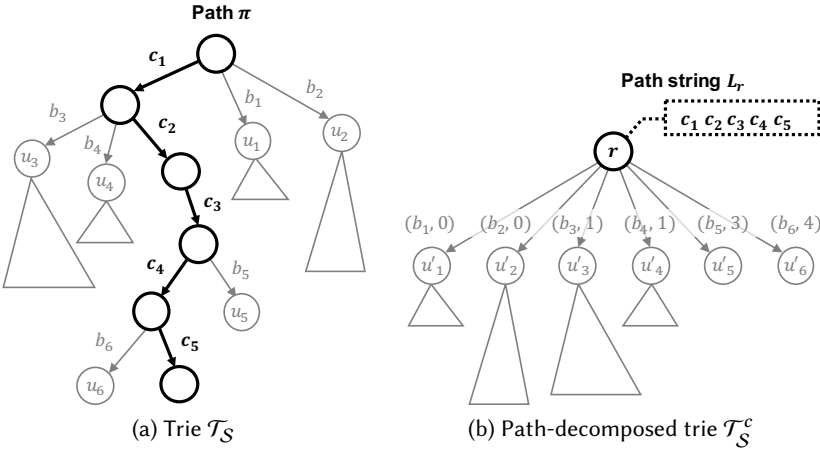


Fig. 2. Illustration of the first transformation of the path decomposition.

which is $L_r = c_1 c_2 c_3 c_4 c_5$. The branching character of u'_5 in \mathcal{T}_S^c is $(b_5, 3)$ because u_5 in \mathcal{T}_S is the child of the third node on the path π with branching character b_5 . Also for the subtrees rooted at the nodes u_1, u_2, \dots, u_6 in \mathcal{T}_S , the decomposition is recursively applied to produce the children of the root in \mathcal{T}_S^c .

Given a keyword K , the retrieval on \mathcal{T}_S can be simulated with a traversal of \mathcal{T}_S^c starting at its root: Let u denote the currently visited node in \mathcal{T}_S^c . On visiting u , we compare the path string L_u with the characters of K . If we find a mismatch at $L_u[i]$ with $b := K[i] \neq L_u[i]$, we descend to the child with branching character (b, i) and drop the first $i + 1$ characters of K .

When storing the characters of each path string L_u in consecutive memory locations, the number of random accesses involved in the retrieval on \mathcal{T}_S^c is bounded by $O(h)$, where h is the height of \mathcal{T}_S^c . The following property regarding the height is satisfied by construction.

Property 1. The height of \mathcal{T}_S^c cannot be larger than that of \mathcal{T}_S .

Centroid Path Decomposition. A way to improve this height bound in the static case is the *centroid path decomposition* [21]. Given an inner node u in \mathcal{T}_S , the *heavy child* of u is the child whose subtree has the most leaves (ties are broken arbitrarily). Given a node u , the *centroid path* is the path from u to a leaf obtained by descending only to heavy children. The centroid path decomposition yields the following property by always choosing centroid paths in the decomposition.

Property 2 ([21]). Through the centroid path decomposition, the height of \mathcal{T}_S^c is bounded by $O(\log n)$.

Key-Value Mapping. We can implement the key-value mapping through \mathcal{T}_S^c because there is a 1-to-1 correlation between nodes in \mathcal{T}_S^c and keywords in \mathcal{S} . A simple approach is to store the associated values in an array A such that $A[u]$ stores the value associated with node u . If we assign each of the n nodes in \mathcal{T}_S^c a unique id from the range $[0, n)$, then A has no vacant entry (i.e. $|A| = n$). Another approach is to embed the value of K_i at the end of L_u , where the node u corresponds to the keyword K_i . This approach can be used without considering the assignment of node ids. In our experiments, we used the latter approach.

3 DYNAMIC PATH-DECOMPOSED TRIE

Although the centroid path decomposition gives a logarithmic upper bound on the height of \mathcal{T}_S^c (cf. Section 2), it can be adapted only in static settings because we have to know the complete topology of \mathcal{T}_S *a priori* to determine the centroid paths. As a matter of fact, previous data structures embracing the path decomposition [21, 30, 35] consider only static applications.

In this section, we present the *incremental path decomposition*, which is a novel procedure to construct a *dynamic path-decomposed trie*, which we call DynPDT in the following. Our procedure incrementally chooses⁶ a node-to-leaf path in \mathcal{T}_S and directly updates the DynPDT \mathcal{T}_S^c on inserting a new keyword of \mathcal{S} . This incrementally chosen path is not a centroid path in general. Thus, the incremental path decomposition does not necessarily satisfy Property 2 but always satisfies Property 1.

In this section, we drop the technical detail of storing the values to ease the explanation of DynPDT, for which we omit the second argument in the insert operation $\text{insert}(K)$ of a new keyword K .

3.1 Incremental Path Decomposition

In the following, we simulate a dynamic trie \mathcal{T}_S by DynPDT \mathcal{T}_S^c . Suppose that \mathcal{T}_S is non-empty. On inserting a new keyword $K \notin \mathcal{S}$ into \mathcal{T}_S , we proceed as follows:

- (1) First traverse \mathcal{T}_S from the root by matching characters of K until reaching the deepest node u whose string label X is a prefix of K .
- (2) Decompose K into $K = XbY$ for $b \in \mathcal{A}$ and $Y \in \mathcal{A}^*$, which is possible since $K \notin \mathcal{S}$ and $K[|K| - 1] = \$$.
- (3) Finally, insert a new child v of u with branching character b and append, from node v , new nodes corresponding to the suffix Y .

In other words, the task of $\text{insert}(K)$ on \mathcal{T}_S is to create a new node-to-leaf path π representing the suffix Y . We call that path π the *incremental path* of the keyword K . We simulate $\text{insert}(K)$ by creating a new node in \mathcal{T}_S^c whose label is the path label of this incremental path π :

- If $\mathcal{S} = \emptyset$, create the root u_1 and associate the keyword K with u_1 by $L_{u_1} \leftarrow K$.

⁶We actually do not construct \mathcal{T}_S , but represent it with the DynPDT \mathcal{T}_S^c

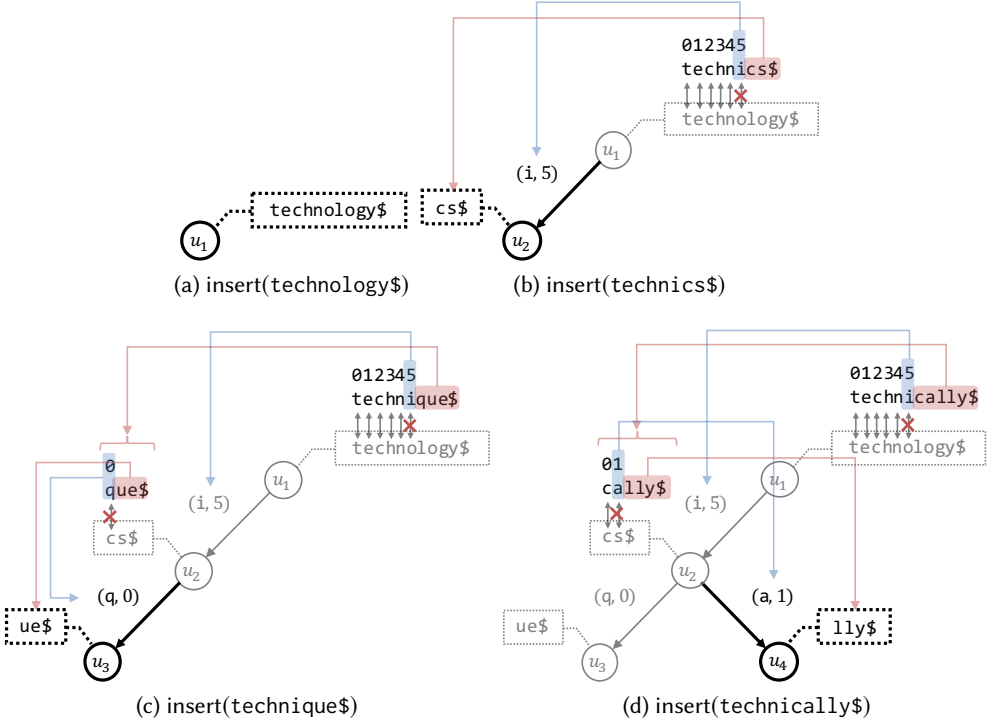


Fig. 3. Process of incremental path decomposition for keywords technology\$, technics\$, technique\$ and technically\$ in this order.

- Otherwise ($S \neq \emptyset$), retrieve the keyword K from the root u_1 in three steps after setting variables $u \leftarrow u_1$ and $S \leftarrow K$:
 - (1) Compare S with L_u . If $S = L_u$, terminate because K is already inserted; otherwise, proceed with Step 2.
 - (2) Find i such that $S[0, i) = L_u[0, i)$ and $S[i] \neq L_u[i]$ (i exists since $K \notin S$ and $K[|K| - 1] = \$$), and search the child of u with branching character $(S[i], i)$. If found, go back to Step 1 after setting the variable u to this child and S to the remaining suffix $S[i + 1, |S|)$; otherwise, proceed with Step 3.
 - (3) Insert K into S by creating a new child v of u with branching character $(S[i], i)$, and store the remaining suffix in v by $L_v \leftarrow S[i + 1, |S|)$.

Example 3.1 (Construction). Figure 3 illustrates the construction process of DynPDT \mathcal{T}_S^c when inserting the keywords $K_1 = \text{technology}\$, K_2 = \text{technics}\$, K_3 = \text{technique}\$, and K_4 = \text{technically}\$$ in this order, where the i -th created node is denoted by u_i . The process begins with an empty trie \mathcal{T}_S^c .

- In the first insertion $\text{insert}(K_1)$, we create the root u_1 and associate K_1 with L_{u_1} , that is, L_{u_1} becomes $\text{technology}\$$. The resulting \mathcal{T}_S^c for $S = \{K_1\}$ is shown in Figure 3a.
- In the second insertion $\text{insert}(K_2)$, we define a string variable S initially set to $S \leftarrow K_2$. We try to retrieve K_2 in \mathcal{T}_S^c by comparing S with L_{u_1} , but fail as there is a mismatching character i at position 5 with $S[0, 5) = L_{u_1}[0, 5) = \text{techn}$ and $S[5] = i \neq o = L_{u_1}[5]$. Based on this mismatch

result, we search the child of u_1 with branching character $(i, 5)$. However, since there is no such child, we add a new child u_2 to u_1 with branching character $(i, 5)$ and associate the remaining suffix $S[6, |S|) = cs\$$ with L_{u_2} . The resulting \mathcal{T}_S^c for $\mathcal{S} = \{K_1, K_2\}$ is shown in Figure 3b.

- (c) In the third insertion $\text{insert}(K_3)$, we initially set the string variable S to $S \leftarrow K_3$ and then compare S with L_{u_1} in the same manner as the second insertion. Since $S[0, 5) = L_{u_1}[0, 5) = \text{techn}$ and $S[5] = i \neq o = L_{u_1}[5]$, we descend to child u_2 with branching character $(i, 5)$. After updating $S \leftarrow S[6, |S|) = \text{que}\$,$ we subsequently compare S with L_{u_2} to obtain the mismatch character q at position 0 with $S[0] = q \neq c = L_{u_2}[0]$. We search the child with branching character $(q, 0)$, but there is no such child; thus, we create the child u_3 and set L_{u_3} to be the remaining suffix $S[1, |S|) = \text{ue}\$$. The resulting \mathcal{T}_S^c for $\mathcal{S} = \{K_1, K_2, K_3\}$ is shown in Figure 3c.
- (d) The fourth insertion $\text{insert}(K_4)$ is also conducted in the same manner. The final trie \mathcal{T}_S^c is shown in Figure 3d.

3.2 Dictionary Operations

It is left to define the operations lookup and delete to make DynPDT a keyword dictionary. Similar to insert, the operation lookup can be performed by traversing \mathcal{T}_S^c from the root. After matching all the characters of K , $\text{lookup}(K)$ returns the value associated with the last visited node. It returns \perp on a mismatch.

Example 3.2 (Retrieval). We provide an example for a successful and an unsuccessful search. Both examples are similar to the construction described in Example 3.1.

- (1) We consider $\text{lookup}(\text{technically}\$)$ for the \mathcal{T}_S^c in Figure 3d. We define a string variable S initially set to $S \leftarrow \text{technically}\$,$ and compare S with L_{u_1} to retrieve (a part of) the keyword from the root. Since $S[0, 5) = L_{u_1}[0, 5) = \text{techn}$ and $S[5] = i \neq o = L_{u_1}[5]$, we descend to child u_2 with branching character $(i, 5)$. Subsequently, we update S to be the remaining suffix as $S \leftarrow S[6, |S|) = \text{cally}\$$ and descend to child u_4 with branching character $(a, 1)$ since $S[0, 1) = L_{u_4}[0, 1) = c$ and $S[1] = a \neq s = L_{u_2}[1]$. Finally, we update $S \leftarrow S[2, |S|) = \text{lly}\$$ and compare S with L_{u_4} . As both match, we return the value stored in u_4 .
- (2) We consider $\text{lookup}(\text{technical}\$)$ for the \mathcal{T}_S^c in Figure 3d. In the same manner as in the above case, we reach node u_4 with the prefix technica and subsequently compare $S = \text{l}\$$ and L_{u_4} . Since $S[0, 1) = L_{u_4}[0, 1) = c$ and $S[1] = \$ \neq L_{u_4}[1] = l$, we search a child with branching character $(\$, 1)$; however, there is no such child. As a result, $\text{lookup}(\text{technical}\$)$ returns \perp .

The operation delete can be implemented by introducing deletion flags for each node (i.e., for each keyword), a trick that is also used in hashing with open addressing [41, Chapter 6.4, Algorithm L]. In other words, $\text{delete}(K)$ retrieves K and sets the deletion flag for the node corresponding to K . However, this approach additionally needs one bit for each node. Another approach is to set the value associated with the deleted keyword to \perp as an invalid value. This approach does not need additional space for the deletion flags. Although these approaches do not free up space after deletion, the space is reused for keywords inserted subsequently if the new keywords share sufficiently long prefixes with the deleted ones.

3.3 Fixing the Alphabet

In practice, a critical problem of DynPDT is that the domain of the edge labels \mathcal{B} in \mathcal{T}_S^c and the longest length of all node labels Λ are not constant in general. We tackle this problem by limiting the size of \mathcal{B} . To this end, we introduce a new parameter λ to forcibly fix the alphabet as $\mathcal{B} = \mathcal{A} \times \{0, 1, \dots, \lambda - 1\}$ in advance. Within this limitation, suppose that we want to create an

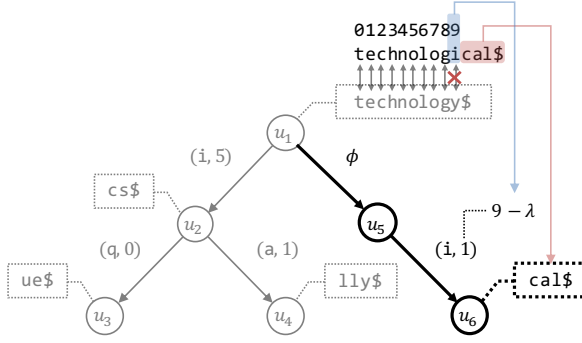


Fig. 4. Process of insert(technological\$) when $\lambda = 8$.

edge labeled (c, i) from node u with $i \geq \lambda$. As this label is not in \mathcal{B} , we create dummy nodes called *step nodes* with a special character ϕ by repeating the following procedure until i becomes less than λ : add a new child v of u with branching character ϕ and recursively set $u \leftarrow v$ and $i \leftarrow i - \lambda$. L_u is the empty string if u is a step node.

Example 3.3 (Step Node). We consider insert(technological\$) for \mathcal{T}_S^c in Figure 3d with $\lambda = 8$. We set $S \leftarrow \text{technological\$}$ and compare S with L_{u_1} . Since $S[0, 9) = L_{u_1}[0, 9) = \text{technolog}$ and $S[9] = i \neq y = L_{u_1}[9]$, we try to create the edge label $(i, 9)$; however, as $i \geq \lambda$, we instead create a step child u_5 with branching character ϕ , descend to this child, and set $i \leftarrow i - \lambda = 1$. Since i becomes less than λ , we define a child u_6 of the step node u_5 with branching character $(i, 1)$ and associate the remaining suffix $S[10, |S|) = \text{cal\$}$ with L_{u_6} . The resulting DynPDT \mathcal{T}_S^c is depicted in Figure 4.

This solution creates additional nodes depending on λ . When λ is too small, many step nodes are created and extra node traversals are involved. When λ is too large, the alphabet size $|\mathcal{B}|$ becomes large and the space usage can increase significantly. Therefore, it is necessary to determine a suitable λ . In Section 6, we empirically determine 32 and 64 to be favorable values for λ .

3.4 Representation Scheme

To use standard trie techniques, we split up \mathcal{T}_S^c into two parts:

- (1) a (standard) trie structure $\mathcal{T}_{\mathcal{D}}$ for a set of strings $\mathcal{D} \subset \mathcal{B}^*$ to represent \mathcal{T}_S^c with the difference that it assigns a node to a unique id instead of its node label, and
- (2) an associative array that maps the ids of the nodes of $\mathcal{T}_{\mathcal{D}}$ to their corresponding node labels, called *node label map (NLM)*.

For example, in Figure 4, the trie $\mathcal{T}_{\mathcal{D}}$ built on the string set $\mathcal{D} = \{(i, 5)(q, 0), (i, 5)(a, 1), \phi(i, 1)\}$ and the NLM stores node labels $L_{u_1}, L_{u_2}, \dots, L_{u_6}$ to be accessed by the respective node ids u_1, u_2, \dots, u_6 .

Node-Label-Map. NLM dynamically manages node labels depending on the node ids assigned. As explained in Section 1, we use the m-Bonsai [49] and FK-hash [22] representations for $\mathcal{T}_{\mathcal{D}}$. Moreover, we design the NLM data structures for m-Bonsai and FK-hash individually, which we respectively present in Sections 4 and 5.

Trie Representation $\mathcal{T}_{\mathcal{D}}$. To discuss the representation approaches in the next sections, we define $\mathcal{T}_{\mathcal{D}}$ to be a dynamic trie with n nodes whose edge labels are characters drawn from the alphabet \mathcal{B} of size $\sigma = |\mathcal{A}| \cdot \lambda$. Although the number of nodes n depends on λ , we write $n := n(\lambda)$ for simplicity. $\mathcal{T}_{\mathcal{D}}$ supports the following operations:

- $\text{addchild}(u, c)$ adds a new child of u with branching character $c \in \mathcal{B}$ and returns its id.
- $\text{getchild}(u, c)$ returns the id of the child v of u with branching character $c \in \mathcal{B}$ if v exists, or returns \perp otherwise.

Motivation for m-Bonsai and FK-hash. We briefly review some common trie representations and point out their suitability for $\mathcal{T}_{\mathcal{D}}$. The simplest representation is a list trie [6, Chapter 2.3.2], which transforms an arbitrary trie to its first-child next-sibling representation. In this representation, each node of the list trie stores its branching character, a pointer to its first child, and a pointer to its next sibling. The list trie represents $\mathcal{T}_{\mathcal{D}}$ in $2n \log n + n \log \sigma$ bits and supports addchild and getchild in $O(\sigma)$ time; however, the operation time becomes problematic if $\sigma = |\mathcal{A}| \cdot \lambda$ is large. Another representation is a ternary search trie (TST) [13] that reduces the time complexity of the list trie to $O(\log \sigma)$; however, the space usage grows to $3n \log n + n \log \sigma$ bits. A well-known time- and space-efficient representation is the double array [2]. Its space usage is $2n \log n$ bits in the best case, while supporting getchild in $O(1)$ time; however, a double array for a large alphabet tends to be sparse in practice. Actually, we are only aware of dynamic double-array implementations handling byte characters (e.g., [38, 64]). Judy [11] and ART (adaptive radix tree) [43] are trie representations that dynamically choose suitable data structures for the trie topology; however, both are also designed for byte characters. As each trie node is associated with an id, compact tries like the z-fast trie [12] representing only $O(|\mathcal{D}|)$ nodes explicitly become inefficient with this requirement.

Compared to these trie representations, m-Bonsai and FK-hash have better complexities. m-Bonsai can represent $\mathcal{T}_{\mathcal{D}}$ in $cn(\log \sigma + O(1))$ bits of expected space for a constant $c > 1$, while supporting getchild and addchild in $O(1)$ expected time [49]. Compared to that, FK-hash needs $cn \log n$ additional bits of expected space, but supports faster insertions in practice.

A straightforward solution to provide the NLM for m-Bonsai and FK-hash is to store the node labels as satellite data in the respective hash table. However, by doing so, we would waste space for each unoccupied entry in the hash table. In the following, we present efficient solutions for the NLM tailored to m-Bonsai and FK-hash.

4 REPRESENTATION BASED ON M-BONSAI

This section presents our approach based on m-Bonsai [49]. m-Bonsai represents trie nodes as entries in a closed hash table that, spoken informally, compactify the stored keys with compact hashing [41].

Outline. We present a plain and a compact form of the $\mathcal{T}_{\mathcal{D}}$ representation based on m-Bonsai. We refer to the former as PBT (Plain m-Bonsai Trie), which is a non-compact variant of m-Bonsai. PBT can be useful for fast implementation although it has not been considered in any applications yet. We refer to the latter as CBT (Compact m-Bonsai Trie) as it uses the original m-Bonsai implementation. We describe PBT and CBT in Sections 4.1 and 4.2, respectively. In both variants, we maintain a hash table H of size m with the *load factor* $\alpha = n/m \leq 1$ to store n nodes. In Section 4.3, we propose a linear-time growing algorithm based on the approach of Arroyuelo et al. [3]. Finally, in Section 4.4, we propose NLM data structures designed for PBT and CBT.

4.1 Plain Trie Representation

PBT uses a hash function $h : \mathbb{N} \rightarrow \mathbb{N}$. Trie nodes are elements in the hash table. As their locations in the hash table are fixed unless the hash table is rebuilt, we use these locations as *node ids*. In other words, the id of a node located at $H[u]$ is u . $\text{addchild}(u, c)$ is performed as follows. We first compose the hash key $k = (u, c) \in \{0, 1, \dots, m-1\} \times \mathcal{B}$ and then compute its *initial address* $i = h(k) \bmod m$.⁷

⁷This paper defines $a \bmod b$ as $a - b \cdot \lfloor a/b \rfloor$.

Let i' be the first vacant address from i determined by linear probing. We create the new child by $H[i'] \leftarrow k$. That is, the id of the new child becomes i' . `getchild` can be also computed in the same manner. If h is fully independent and uniformly random, the operations can be performed in $O(1)$ expected time. PBT uses $m \lceil \log(m\sigma) \rceil$ bits of space.

Practical Implementation. The table size m is a power of two in order to quickly compute the modulo operation of $h(k) \bmod m$ by using the bitwise AND operation $h(k) \& (m - 1)$ [47, Section 4.4]. We set the maximum load factor to $\hat{\alpha} := 0.9$. If α reaches $\hat{\alpha}$ during an update, we double the size of the hash table by the growing algorithm described in Section 4.3. We set the initial capacity of the hash table to $m = 2^{16}$. Our hash function h is a XorShift hash function⁸ derived from [51].

4.2 Compact Trie Representation

CBT reduces the space usage of PBT with the compact hashing technique [41]. Locating nodes on a compact hash table is identical to PBT with the difference that CBT uses a bijective transform $h : \{0, 1, \dots, m\sigma - 1\} \rightarrow \{0, 1, \dots, m\sigma - 1\}$ that maps a key k to its hash value $h(k) \bmod m$ and its quotient $\lfloor h(k)/m \rfloor$. Instead of k , the compact hash table stores only its quotient $\lfloor h(k)/m \rfloor$ in $H[i']$. The hash value $h(k)$ can be restored from the initial address $i = h(k) \bmod m$ and the quotient $H[i'] = \lfloor h(k)/m \rfloor$, where i' is the first empty slot at or after the initial address i . The original key k can also be restored from the hash value $h(k)$ since h is bijective. Therefore, `addchild` and `getchild` can be performed in the same manner as PBT if the corresponding initial address i can be identified from the location i' .

The remaining problem is how to identify the corresponding initial address i from i' . Poyias et al. [49] solved this problem by introducing a *displacement array* D such that $D[i']$ keeps the number of probes from i to i' , that is, $D[i'] = (i' - i) \bmod m$. Given a location i' , one can compute the corresponding initial address i with $(i' - D[i']) \bmod m$. Although a value in D is at most $m - 1$, the average value becomes small if h is fully independent and uniformly random and the load factor α is small. Poyias et al. [49] demonstrated that D can be represented in $O(m)$ bits using CDRW (Compact Dynamic ReWritable) arrays. As H takes $m \lceil \log \sigma \rceil$ bits for the quotients, CBT can represent \mathcal{T}_D in $m \log \sigma + O(m)$ expected bits of space.

Practical Representation of the Displacement Array. The representation of D with the CDRW array seems impractical. Poyias et al. [49] gave an alternative practical representation, where D is represented by three data structures D_1 , D_2 and D_3 as follows.

- (1) D_1 is a simple array of length m in which each element uses Δ_1 bits for a constant $\Delta_1 > 1$.
- (2) D_2 is a *compact hash table* (CHT) described by Cleary [17], which stores keys from $\mathcal{U} = \{0, 1, \dots, m - 1\}$ and values from $\{0, 1, \dots, 2^{\Delta_2} - 1\}$ for a constant $\Delta_2 > 1$. The keys are stored in a closed hash table of length $m' < m$ through the compact hashing technique [41], where m' is a power of two (a property that is in common with m). In detail, the hash table consists of
 - a bijective transform $h : \mathcal{U} \rightarrow \mathcal{U}$,
 - an integer array Q of length m' to store the quotients of the keys (i.e., entry indices of D) representable in $\log(m/m')$ bits,
 - an integer array F of length m' to store displacement values of D representable in Δ_2 bits, and
 - two bit arrays each of length m' storing the displacement values of the quotients in Q (not to be confused with the displacement values stored in F).

⁸<http://xorshift.di.unimi.it/splitmix64.c>.

On inserting a key $k \in \mathcal{U}$, we store its quotient $\lfloor h(k)/m' \rfloor$ in the first vacant slot in Q starting at the initial address $h(k) \bmod m'$. The collisions in Q are therefore resolved with linear probing. However, this collision resolution poses the same problem as in CBT, as additional displacement information is required to restore the initial address of a stored quotient in Q . Cleary solves this problem by using two bit arrays (see [17]). Finally, $F[i]$ stores the value associated with the key whose quotient is stored in $Q[i]$. Since F uses $m'\Delta_2$ bits of space, D_2 uses $m' \log(m/m') + m'\Delta_2 + 2m'$ bits of space in total.

- (3) D_3 is a standard associative array that maps keys from \mathcal{U} to values from \mathcal{U} . In our implementation, D_3 is a closed hash table with linear probing. Given m'' is the capacity of D_3 , D_3 takes $2m'' \log m$ bits.

The representation of the entry $D[i]$ for an integer i depends on its actual value:

- (1) If $D[i] < 2^{\Delta_1} - 1$, then we store $D[i]$ in the Δ_1 bits of $D_1[i]$.
- (2) If $2^{\Delta_1} - 1 \leq D[i] < 2^{\Delta_1} + 2^{\Delta_2}$, we represent $D[i]$ by the key-value pair $(i, D[i] - 2^{\Delta_1})$ stored in D_2 .
- (3) Finally, if $D[i] \geq 2^{\Delta_1} + 2^{\Delta_2}$, we represent $D[i]$ by the key-value pair $(i, D[i])$ stored in D_3 .

In the experiments, we set $\Delta_1 = 4$ and $\Delta_2 = 7$. We set the initial capacities of D_2 and D_3 to $m' = 2^{12}$ and $m'' = 2^6$, respectively. We set the maximum load factor of D_2 and D_3 to 0.9. If the actual load factor of D_2 (resp. D_3) reaches the maximum load factor 0.9, we double the size of D_2 (resp. of D_3).

Design of the Bijective Transform. Since we assume that m , m' , and σ are powers of two, the bijective transform is $h : \{0, 1, \dots, 2^z - 1\} \rightarrow \{0, 1, \dots, 2^z - 1\}$ for some z . We design this function as the concatenation of two bijective functions $h = h_1 \circ h_2$, where $h_1(x) = x \oplus \lfloor x/2^a \rfloor$ for an integer a larger than $\lfloor z/2 \rfloor$ and $h_2(x) = xp \bmod 2^z$ for a large prime p smaller than 2^z . h_1 is based on the XorShift random number generators [44], where the inverse function h_1^{-1} is given by $h_1^{-1}(x) = h_1(x)$. The inverse function h_2^{-1} of h_2 is given by $h_2^{-1}(x) = xp^{-1} \bmod 2^z$, where $p^{-1} \in \{1, 2, \dots, 2^z - 1\}$ is the multiplicative inverse of p such that $pp^{-1} \bmod 2^z = 1$ (see [42] for details). By construction, the inverse function h^{-1} of h is $h^{-1} = h_2^{-1} \circ h_1^{-1}$. Our hash function is inspired by the SplitMix algorithm [51].

4.3 Linear-Time Growing Algorithm

If the load factor α of hash table H of length m reaches the maximum load factor $\hat{\alpha}$, we create a new hash table H' (and a new displacement array D' for CBT) of length $2m$ and relocate all nodes to H' . Since a node depends on the position of its parent in H , we can relocate a node only after having relocated all its ancestors. This can be done in a top-down traversal (e.g., in BFS or DFS order) of the tree during which all children of a node are successively selected. However, because selecting all children of a node is performed by checking `getchild` for *all* possible characters in \mathcal{B} , the relocation based on a top-down traversal needs $O(n\sigma)$ expected time and is therefore only for tiny alphabets practical. Here we describe a bottom-up approach that is based on the approach by Arroyuelo et al. [3]. This approach, called *growing algorithm*, runs in $O(n)$ expected time. A pseudo code of it is shown in Algorithm 1.

Given a trie $\mathcal{T}_{\mathcal{D}}$ with a hash table H of length m , the algorithm constructs an equivalent trie $\mathcal{T}'_{\mathcal{D}}$ with a hash table H' of length $2m$. To explain the algorithm, we define two operations `getedge(u)` returning the branching character of node u and `getparent(u)` returning the parent id of node u . They can be computed in constant time because $H[u]$ explicitly stores the branching character and the parent id as the hash key in PBT. CBT can also restore the hash key from $H[u]$ and $D[u]$.

Algorithm 1 Linear-time growing algorithm of PBT and CBT**Input:** Trie $\mathcal{T}_{\mathcal{D}}$ with hash table H of size m **Output:** Equivalent trie $\mathcal{T}'_{\mathcal{D}}$ with hash table H' of size $2m$

```

1: Create an empty trie  $\mathcal{T}'_{\mathcal{D}}$  with hash table  $H'$  of size  $2m$  and create its root
2: Create an integer array Map and a bit array Done, each of length  $m$ 
3: Initialize Done[ $i$ ]  $\leftarrow \emptyset$  for all  $i$ 
4: Done[ $u_1$ ]  $\leftarrow 1$  and Map[ $u_1$ ]  $\leftarrow u'_1$ , where  $u_1$  and  $u'_1$  are the root ids of  $\mathcal{T}_{\mathcal{D}}$  and  $\mathcal{T}'_{\mathcal{D}}$ , respectively
5: for  $i = 0, \dots, m - 1$  do
6:   if  $H[i]$  is empty then continue
7:    $u \leftarrow i$  and  $\pi \leftarrow$  empty string
8:   while Done[ $u$ ]  $\neq 1$  do ▷ Climb up  $\mathcal{T}_{\mathcal{D}}$ 
9:      $\pi \leftarrow \mathcal{T}_{\mathcal{D}}.\text{getedge}(u) + \pi$  ▷ Prepend ancestor to  $\pi$ 
10:     $u \leftarrow \mathcal{T}_{\mathcal{D}}.\text{getparent}(u)$ 
11:   end while
12:    $u' \leftarrow$  Map[ $u$ ]
13:   for  $c \in \pi$  do ▷ Walk down the computed path
14:      $u \leftarrow \mathcal{T}_{\mathcal{D}}.\text{getchild}(u, c)$  and  $u' \leftarrow \mathcal{T}'_{\mathcal{D}}.\text{addchild}(u', c)$ 
15:     Map[ $u$ ]  $\leftarrow u'$  and Done[ $u$ ]  $\leftarrow 1$ 
16:   end for
17: end for
18: output  $\mathcal{T}'_{\mathcal{D}}$ 

```

In the growing algorithm, we initially define two auxiliary arrays Map and Done: Map is an integer array and Done is a bit array, each of length m . We store in Done[u] a 1 after relocating the node stored in $H[u]$. We keep the invariant that whenever Done[u] = 1, then Map[u] stores the position in H' of the node stored in $H[u]$. All bits in Done are initialized by \emptyset except for the root. We scan H from left to right and perform the following steps for each non-vacant slot i . We first set u to i and π to an empty string, and then climb up the path from the node u to the root. We prematurely stop when encountering a node v with Done[v] = 1. In this case, all ancestors of v have already been relocated such that there is no need to visit them again. Subsequently, we walk down the computed path π while relocating the visited nodes. Since we do not reprocess already visited nodes, we can perform the node relocation in $O(m) + O(n) = O(n)$ expected time, with $n = \hat{\alpha} \cdot m$ for a constant loaf factor $\hat{\alpha}$.

Extra Working Space. Algorithm 1 maintains the auxiliary arrays Map of $m \lceil \log(2m) \rceil$ bits, Done of m bits and π of $h \lceil \log \sigma \rceil$ bits, where h is the height of $\mathcal{T}_{\mathcal{D}}$. Thus, the extra working space is $m \lceil \log m \rceil + 2m + h \lceil \log \sigma \rceil$ bits if we create the auxiliary arrays naively. However, the working space of Map can be shared with H because $H[i]$ for Done[i] = 1 is no longer needed. In PBT, the working space of Map can be fully placed in H because the space of H is $m \lceil \log(m\sigma) \rceil$ bits and σ is at least 2 in practice.⁹ Based on this in-place approach, the extra working space of Algorithm 1 is only $m + h \lceil \log \sigma \rceil$ bits, taking account for Done and π in PBT. In practice, the space of π is negligible because h is bounded by the maximum length of keywords in \mathcal{S} and $h \ll m$.

In CBT, H uses only $m \lceil \log \sigma \rceil$ bits. As $\sigma \ll m$ in most scenarios, it is difficult to completely store Map in H ; however, we can also use the space of D_1 , which is $m\Delta_1$ bits. If $\lceil \log(2m) \rceil \leq \lceil \log \sigma \rceil + \Delta_1$, Map can be fully placed in H and D ; otherwise, the extra working space of $m(\lceil \log(2m) \rceil - \lceil \log \sigma \rceil - \Delta_1)$ bits for Map is needed in addition to that of Done and π .

⁹Even for $\sigma = 1$, a simple bit array suffices.

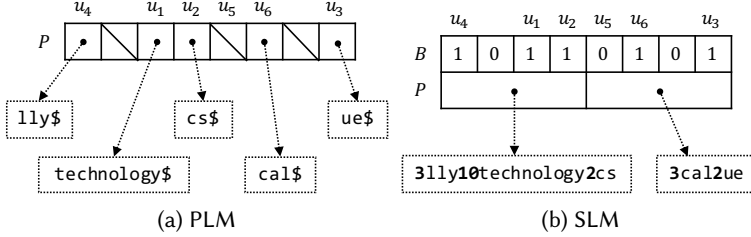


Fig. 5. Examples of NLM in m-Bonsai for the DynPDT in Figure 4.

4.4 NLM Data Structures

In m-Bonsai, the node ids are values drawn from the universe $[0, m)$ whose randomness depend on the used hash function. As the task of an NLM data structure is to map node ids to their respective node labels, an appropriate NLM data structure for m-Bonsai is a dynamic associative array that stores node label strings L_i for arbitrary integer keys $i \in [0, m)$. In what follows, we first present a plain approach and then show how to compactify it.

Plain NLM. The simplest approach is to use a pointer array P of length m such that $P[i]$ stores the pointer to L_i or \perp if no node with id i exists. We refer to the approach as PLM (Plain Label Map). Figure 5a shows an example of PLM. Given a node of id i , PLM can obtain L_i through $P[i]$ in $O(1)$ time. However, P takes $mw = O(m \log m)$ bits, where the word size is $w = \Theta(\log m)$. This space consumption is obviously large.

Sparse NLM. We present an alternative compact approach that reduces the pointer overhead of PLM in a manner similar to Google’s sparse hash table [26]. In this approach, we divide the node labels into *groups* of $\ell = \Theta(w)$ labels over the ids. That is, the first group consists of $L_0, L_1, \dots, L_{\ell-1}$, the second group consists of $L_\ell, L_{\ell+1}, \dots, L_{2\ell-1}$, and so on. Moreover, we introduce a bitmap B such that $B[i] = 1$ iff L_i exists. We concatenate all node labels L_i with $B[i] = 1$ of the same group together, sorted in the id order. The length of P becomes $\lceil m/\ell \rceil$ by maintaining, for each group, a pointer to its concatenated label string. We refer to the approach as SLM (Sparse Label Map).

With the array P and the bitmap B , we can access L_i as follows: If $B[i] = 0$, we are done since L_i does not exist in this case; otherwise, we obtain the concatenated label string storing L_i from $P[g]$, where $g = \lfloor i/\ell \rfloor$. Given $j = \sum_{k=0}^{i \bmod \ell} B_g[k]$ for the bit chunk $B_g := B[g\ell, (g+1)\ell)$, L_i is the j -th node label of the concatenated label string. As $\ell = \Theta(w)$, counting the occurrences of 1s in chunk B_g is supported in constant time using the *popcount* operation [25]. It is left to explain how to search L_i in the respective concatenated label string. For that we present two representations of the concatenated label strings:

- (1) If the node labels are straightforwardly concatenated (e.g., the second group in Figure 5a is `calue` in $\ell = 4$), we can sequentially count the `$` delimiters to find the $(j-1)$ -th delimiter marking the ending of the $(j-1)$ -th stored string, after which L_i starts. We can therefore extract L_i in $O(\ell\Lambda)$ time, where Λ again denotes the maximum length of all node labels.
- (2) We can shorten the scan time with the *skipping* technique used in array hashing [8]. This technique puts its length in front of each node label via some prefix encoding such as `VByte` [61]. Note that we can omit the terminators of each node label. The skipping technique allows us to jump ahead to the start of the next node label; therefore, the scan is supported in $O(\ell)$ time. Figure 5b shows an example of SLM with the skipping technique.

Regarding the space usage of SLM, P and B use $w\lceil m/\ell \rceil$ and m bits, respectively. For $\ell = \Theta(w)$, the total space usage becomes $O(m)$ bits, which is smaller than mw bits in PLM; however, the access time is $O(w) = O(\log m)$.

5 REPRESENTATION BASED ON FK-HASH

This section presents our DynPDT representation approaches based on FK-hash [22]. The basic idea of FK-hash is the same as that of m-Bonsai. The difference is that FK-hash incrementally assigns node ids and explicitly stores them as values in the hash table, while m-Bonsai uses the locations of the stored elements of the hash table as node ids. Although FK-hash uses more space than m-Bonsai, the assignment of node ids simplifies the growing algorithm.

Outline. In the same manner as m-Bonsai, we consider a plain and a compact representation based on FK-hash. In Section 5.1 we present both representations. In Section 5.2 we propose NLM data structures designed for FK-hash.

5.1 Trie Representations

Like m-Bonsai, FK-hash locates nodes on a closed hash table H of length m , but does not use the addresses of H as node ids. FK-hash incrementally assigns node ids from zero and explicitly stores them in an integer array M of length m . In other words, when creating the u -th node by storing it in $H[i]$, its node id is u , which is stored in $M[i]$. In a way similar to m-Bonsai, $\text{addchild}(u, c)$ is performed as follows: We compose the key $k = (u, c)$, hash it with h , and then search the first vacant slot $H[i']$ from $i = h(k) \bmod m$ by linear probing. Given u_{\max} is the currently largest node id, we assign the id $v = u_{\max} + 1$ to the new child, and set $H[i'] = k$ and $M[i'] = v$. The displacement information $i' - i$ is maintained analogously to m-Bonsai.

In the same manner as m-Bonsai, we can think of two representations depending on whether H is compactified or not. The non-compact one is referred to as PFKT (Plain FK-hash Trie). The compact one is referred to as CFKT (Compact FK-hash Trie). Compared to PBT and CBT, PFKT and CFKT keep an additional integer array M and require $m\lceil \log n \rceil$ additional bits of space.

Table Growing. An advantage of FK-hash is that growing the hash table is done in the same manner as in standard closed hash tables. In detail, H can be enlarged by scanning nodes on H from left to right and relocating the nodes in a new hash table H' of length $2m$. The growing algorithm takes $O(m)$ expected time. This time complexity is identical to that of Algorithm 1; however, the growing algorithm of FK-hash is faster in practice because of its simplicity. In addition, no auxiliary data structure is needed like Map and Done used by Algorithm 1.

5.2 NLM Data Structures

Like in Section 4.4, we introduce PLM and SLM adapted to FK-hash. Figure 6 shows an example for each of them. Although PLM in FK-hash is basically identical to that in m-Bonsai, SLM can be simplified as follows.

In m-Bonsai, it is necessary to identify whether L_i exists and the rank of L_i in the group because node ids are randomly assigned; therefore, we introduced a bitmap B of length m and utilized the popcount operation. In FK-hash, however, such a bitmap is not needed because node ids are incrementally assigned. Put simply, a node label L_i is stored in the group of id $g = \lfloor i/\ell \rfloor$ and located at the $(i \bmod \ell)$ -th position in the group. When using the skipping technique, care has to be taken for the step nodes whose node labels are empty. For each of them, we put the length 0 in its corresponding concatenated label string. For example, we put a '0' in the second concatenated label string for the step node u_5 in Figure 6b. Finally, we can insert a new node label by appending it to the last concatenated label string.

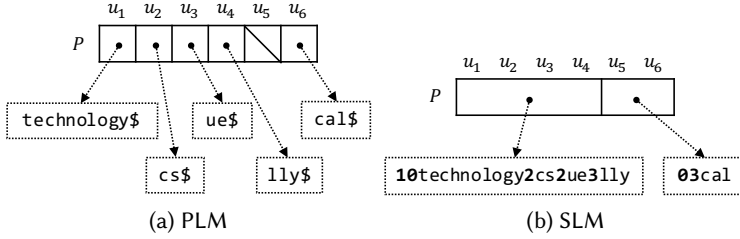


Fig. 6. Examples of NLM in FK-hash for the DynPDT in Figure 4.

6 EXPERIMENTS

In this section we evaluate the practical performance of DynPDT. The source code for our experiments are available at https://github.com/kampersanda/dictionary_bench.

6.1 Setup

We conducted all experiments on one core of a quad-core Intel Xeon CPU E5-2680 v2 clocked at 2.80 Ghz in a machine with 256 GB of RAM, running the 64-bit version of CentOS 6.10 based on Linux 2.6. We implemented our data structures in C++17. We compiled the source code with g++ (version 7.3.0) in optimization mode `-O3`. We used 4-byte integers for the values associated with the keywords.

Datasets. Our benchmarks are based on the following eight real-world datasets:

- GeoNames consists of 7 million different names for the geographic points provided by the GeoNames database.¹⁰ Managing such geographic identifiers within a limited resource is essential in modern geographic information systems as described in [45]. We obtained the geographic names by extracting the *asciiname* column of the GeoNames dump in the same manner as [45].
- AOL consists of 10 million different search queries in the AOL database, which is a huge collection of 20 million search queries from 650,000 users sampled over three months.¹¹ The dataset contains keywords written in natural English, which has been often used to benchmark search algorithms such as [30].
- Wiki consists of 14 million different page titles from the English Wikipedia dump at September 2018.¹² As the dataset contains various special characters encoded in UTF-8, the alphabet size is larger than that of AOL. It is also a well-used dataset to benchmark search algorithms such as [5, 30, 39].
- DNA consists of all 12-mers (i.e., substrings of length 12) found in the DNA dataset from the Pizza&Chili corpus.¹³ Among the used datasets, it has the smallest alphabet and the shortest keywords. The number of keywords is 15 million. In bioinformatics, popular alignment software need to manage such keywords within limited space as described in [45].
- LUBMS consists of 53 million different URIs extracted from the RDF dataset generated by the Lehigh University Benchmark [31] for 1,600 universities.¹⁴ Modern RDF systems [62, 63]

¹⁰<http://download.geonames.org/export/dump/>

¹¹<http://www.cim.mcgill.ca/~dudek/206/Logs/AOL-user-ct-collection/>

¹²<https://dumps.wikimedia.org/enwiki/>

¹³<http://pizzachili.dcc.uchile.cl/texts/dna/>

¹⁴The dataset is distributed under the name 'DS5' at <https://exascale.info/projects/web-of-data-uri/>.

Table 1. Statistics for the datasets used. Size is the total length of the keywords, n is the number of all distinct keywords in millions (M), MinLen (resp. MaxLen and AveLen) is the maximum (resp. minimum and average) length of the keywords, and $|\mathcal{A}|$ is the actual alphabet size of the keywords.

	Size	n	MinLen	MaxLen	AveLen	$ \mathcal{A} $
GeoNames	109 MiB	7.3 M	2	152	15.7	99
AOL	224 MiB	10.2 M	2	523	23.2	85
Wiki	286 MiB	14.1 M	2	252	21.2	200
DNA	189 MiB	15.3 M	13	13	13.0	16
LUBMS	3.1 GiB	52.6 M	10	80	63.7	57
LUBML	13.8 GiB	230.1 M	10	80	64.2	57
UK	2.7 GiB	39.5 M	17	2,030	72.4	103
WebBase	6.6 GiB	118.2 M	10	10212	60.2	223

encode URIs in a huge set into unique integers by using a dynamic keyword dictionary. The dataset is evaluated in [46] to analyze the performances of RDF systems.

- LUBML consists of 230 million different URIs extracted from the RDF dataset generated by the Lehigh University Benchmark [31] for 7,000 universities.¹⁵ The dataset is a larger version of LUBMS. It is also evaluated in [46].
- UK consists of 40 million different URLs obtained from a 2005 crawl of the .uk domain performed by UbiCrawler [14].¹⁶ URLs are traditionally used to benchmark search algorithms for long strings such as [5, 7, 30, 39]. Also, the modern Web crawler [59] manages a huge set of URLs by using a dynamic keyword dictionary.
- WebBase consists of 118 million different URLs of a 2001 crawl performed by the WebBase crawler [34].¹⁷ The dataset is larger than UK and also used in previous experiments of keyword dictionaries such as [30].

Table 1 summarizes relevant statistics for each dataset.

6.2 Average Height

We evaluate the average height of the DynPDT \mathcal{T}_S^c built on our datasets. The average height of \mathcal{T}_S^c is the arithmetic mean of the heights of all nodes over the number of nodes, omitting step nodes in the calculation. Although the average height is an important measure related to the average number of random accesses, we cannot *a priori* predict the average height of DynPDT because this number depends on the insertion order of the keywords. To reason about the quality of the average height, we study it in relation to the following known lower and upper bounds on it: The lower bound is the average height of the path-decomposed trie created by the centroid path decomposition [1, Corollary 3]. The upper bound is the average height of the path-decomposed trie created by always choosing the child whose subtree has the *fewest* number of leaves.

Table 2 shows the experimental results of the average heights of \mathcal{T}_S^c and \mathcal{T}_S for all the datasets. To analyze the performance of DynPDT in our experiments, we constructed DynPDT dictionaries by inserting keywords in random order. For that, we shuffled the dataset with the Fisher–Yates shuffle algorithm [20]. Naturally, the actual average heights of \mathcal{T}_S^c are between their lower and upper bounds, and those of \mathcal{T}_S are the same as AveLen. The upper bounds are more than twice as large as

¹⁵Although this dataset is not distributed, one can obtain the identical dataset through the LUBM data generator (called UBA) at <http://swat.cse.lehigh.edu/projects/lubm/>.

¹⁶<http://law.di.unimi.it/webdata/uk-2005/>

¹⁷<http://law.di.unimi.it/webdata/webbase-2001/>

Table 2. Experimental results of the average heights of \mathcal{T}_S^c and \mathcal{T}_S denoted by AveHeight. Also, AveHeightLB and AveHeightUB are the lower bound and the upper bound of the average height of \mathcal{T}_S^c , respectively (defined in Section 6.2). AveHeightLB is the average height of the path-decomposed trie obtained by the centroid path decomposition. AveHeightUB is the average height of the path-decomposed trie obtained by the path decomposition selecting children with the fewest leaves.

	GeoNames	AOL	Wiki	DNA	LUBMS	LUBML	UK	WebBase
AveHeight of \mathcal{T}_S^c	6.0	6.2	6.3	9.0	7.5	7.9	7.8	7.3
AveHeightLB of \mathcal{T}_S^c	5.2	5.2	5.3	8.9	6.6	7.4	6.0	6.2
AveHeightUB of \mathcal{T}_S^c	8.5	10.5	9.7	10.7	11.8	12.4	14.7	15.4
AveHeight of \mathcal{T}_S	15.7	23.2	21.2	13.0	63.7	64.2	72.4	60.2

Table 3. Experimental results of PDT-CFK for various values of the parameter λ . Steps is the proportion of the number of step nodes among all nodes in DynPDT, Space is the working space in GiB, and Time is the elapsed time for the construction in seconds.

λ	Wiki			LUBMS			UK		
	Steps	Space	Time	Steps	Space	Time	Steps	Space	Time
4	19.55%	0.36	20.9	7.75%	0.78	116.1	37.60%	1.22	116.5
8	6.12%	0.27	18.2	2.83%	0.78	96.8	15.51%	1.20	96.2
16	1.32%	0.27	17.1	0.31%	0.78	84.9	5.22%	1.20	87.6
32	0.12%	0.27	17.3	0.02%	0.79	83.3	1.31%	1.20	86.0
64	0.00%	0.27	17.2	0.00%	0.80	82.9	0.23%	1.21	85.4
128	0.00%	0.27	17.4	0.00%	0.80	82.9	0.04%	1.22	85.3
256	0.00%	0.28	17.2	0.00%	0.81	83.6	0.01%	1.22	85.2
512	0.00%	0.28	17.2	0.00%	0.82	83.3	0.00%	1.23	85.4
1024	0.00%	0.28	17.3	0.00%	0.83	83.1	0.00%	1.24	85.4

the lower bounds for AOL, UK, and WebBase; however, the upper bounds were up to 5.4x smaller than the average heights of \mathcal{T}_S due to the path decomposition, especially for long keywords such as URIs. Therefore, the incremental path decomposition can make dynamic keyword dictionaries more cache-friendly, especially for long keywords even if the insertion order is inconvenient and the average height is close to the upper bound.

6.3 Parameter for Step Nodes

The parameter λ influences the number of step nodes. We analyze the space and time performance of DynPDT when varying the parameter λ . In this experiment, we constructed DynPDT dictionaries for each parameter $\lambda \in \{4, 8, 16, \dots, 1024\}$ on the datasets Wiki, LUBMS and UK, and observed the working space and the construction time. For the DynPDT representation, we tested the combination of CFKT and SLM with $\ell = 16$, referred to as PDT-CFK in the following. As described in Section 6.2, the dictionary was constructed by inserting keywords in random order. The working space was measured by checking the maximum resident set size (RSS) required during the online construction.

Table 3 shows the experimental results for construction. Since λ has a direct impact on σ , which influences the space usage of H , the working space depends on the value of λ . Although this dependency looks like λ and the taken space are in direct correlation, for Wiki and UK, the working

Table 4. Proportion of the number of step nodes to the total number of nodes in DynPDT. Bold font indicates the results with the smallest λ such that Steps is less than 1%. AveNLL is the average length of the node labels.

	Steps						AveNLL
	$\lambda = 4$	$\lambda = 8$	$\lambda = 16$	$\lambda = 32$	$\lambda = 64$	$\lambda = 128$	
GeoNames	6.34%	1.44%	0.28%	0.04%	0.00%	0.00%	6.1
AOL	22.16%	6.83%	1.26%	0.11%	0.01%	0.00%	10.6
Wiki	19.55%	6.12%	1.32%	0.12%	0.00%	0.00%	8.7
DNA	0.11%	0.00%	0.00%	0.00%	0.00%	0.00%	1.4
LUBMS	7.75%	2.83%	0.31%	0.02%	0.00%	0.00%	3.7
LUBML	7.66%	2.80%	0.31%	0.02%	0.00%	0.00%	3.7
UK	37.60%	15.51%	5.22%	1.31%	0.23%	0.04%	18.0
WebBase	24.15%	8.94%	2.46%	0.50%	0.08%	0.02%	11.1

spaces for $\lambda = 4$ (i.e., 0.36 GiB and 1.22 GiB respectively) were not the smallest. For Wiki, the reason for this is that many step nodes raised the load factor α and involved an additional enlargement of the hash table. Specifically, the enlargements were conducted nine times with $\lambda = 4$, although they were conducted eight times with $\lambda \geq 8$. For UK, this reason is that the high load factor α caused by a huge number of step nodes raised the average displacement value stored in D and involved the use of D_2 and D_3 , although no additional enlargement was conducted. Regarding the time performance, this huge number of step nodes slowed down the construction. Therefore, a too small parameter λ can involve large space requirements and long construction times. On the other hand, when $16 \leq \lambda$, the working space and construction time do not significantly vary.

From this observation, we derive two facts for λ : On the one hand, the most important recommendation is not to choose a parameter λ that is too small. On the other hand, choosing a large parameter λ is not a significant problem because the space and time performance do not significantly decrease as λ grows. For example, when $\lambda = 32$ on Wiki, the proportion of step nodes is 0.12%; however, even with a larger parameter λ such as 512 or 1024, the working space and construction time are almost the same. Table 4 shows Steps for each parameter λ and the average length of the node labels (denoted by AveNLL) for all the datasets. Even for long keywords like URLs (i.e., UK), AveNLL is bounded by 18.0 and Steps is within 1% of all nodes when $\lambda = 64$. Among the tested values for λ , we suggest setting λ to 32 or 64 for keywords whose length is not much longer than that of the URL datasets.

6.4 Comparison among DynPDT Representations

We compared the performance of our DynPDT representations, for which we benchmarked the following six combinations:

- PDT-PB is the combination of PBT and PLM,
- PDT-SB is the combination of PBT and SLM,
- PDT-CB is the combination of CBT and SLM,
- PDT-PFK is the combination of PFKT and PLM,
- PDT-SFK is the combination of PFKT and SLM, and
- PDT-CFK is the combination of CFKT and SLM.

We evaluated the working space during the construction and the running times of insert and lookup. Like in Section 6.3, we constructed each dictionary and measured its working space. To measure the lookup time, we chose 1 million random keywords from each dataset. The running

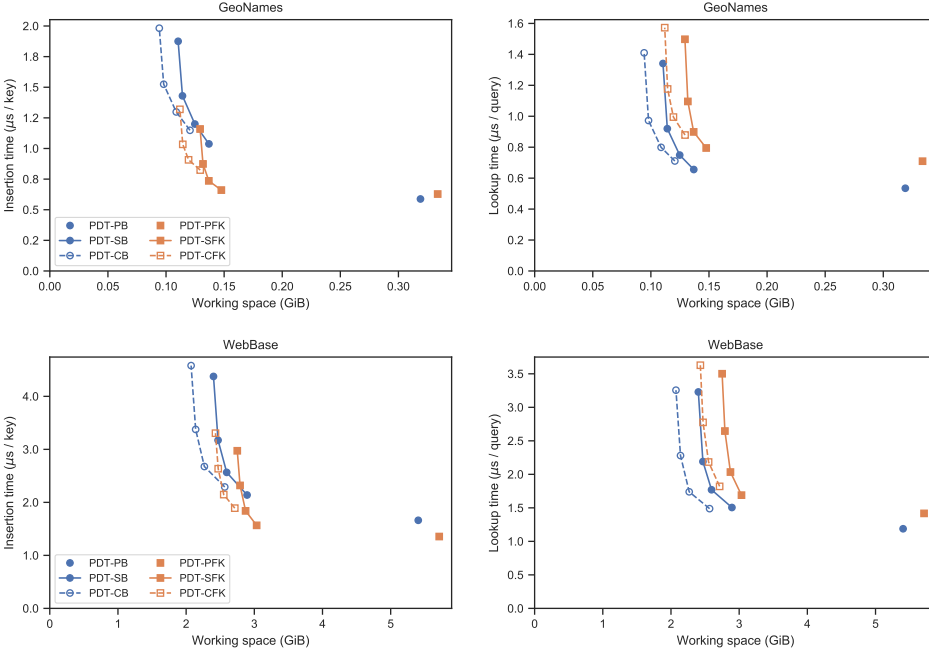


Fig. 7. Experimental results for combinations of DynPDT representations.

times are the average of 10 runs. For SLM, we tested $\ell \in \{8, 16, 32, 64\}$. For λ , we chose the smallest value among those from Table 4 where Steps is less than 1%.

Figure 7 shows the experimental results for GeoNames and WebBase. Regarding the representations using SLM, the working space is the largest but the running times are the shortest with $\ell = 8$, and vice versa with $\lambda = 64$. In other words, for each representation in the plots, the rightmost and lowest result is the one with $\ell = 8$, and the leftmost and highest result is the one with $\ell = 64$.

We observe that

- SLM significantly reduces the working space of PLM. Compared to PDT-PB, PDT-SB is 57–65% smaller for GeoNames and 46–56% smaller for WebBase. Compared to PDT-PFK, PDT-SFK is 56–61% smaller for GeoNames and 47–52% smaller for WebBase.
- Regarding the representations based on m-Bonsai, the insert time of SLM is slower than that of PLM because inserting a new node label into the group is costly. When $\ell = 8$, the insertion of PDT-SB is 29-163% slower than that of PDT-PB; however, the lookup times are competitive.
- Regarding the representations based on FK-hash, SLM with $\ell = 8$ is competitive to PLM with respect to the insert time because the update algorithm is simple. Also, the lookup times are competitive.
- The time performance of SLM with large group sizes ($\ell = 32$ or 64) is worse than that of SLM with small group sizes ($\ell = 8$ or 16). For example, for GeoNames, PDT-SB with $\ell = 64$ is 19% smaller but 81–105% slower than PDT-SB with $\ell = 8$.
- The compact trie representations CBT and CFKT are more lightweight but slower than the plain representations PBT and PFKT; however, the differences are small. For example, PDT-SB is 12% smaller but 8–11% slower than PDT-CB for GeoNames.

- The representations based on m-Bonsai are smaller than those based on FK-hash. Also regarding the lookup time, the m-Bonsai representations are faster. However, regarding the insert time, the FK-hash representations are faster because the growing algorithm is simple.

6.5 Comparison with Existing Data Structures

We compare the performance of DynPDT with existing data structures. We exhaustively tested existing implementations of dynamic keyword dictionaries such as open-source dynamic hash containers [28, 53, 56] and recent dynamic trie indexes [52, 58]. However, compared to DynPDT, most of them consumed significantly more space. For our benchmarks, we selected the following four space-efficient implementations:¹⁸

- ArrayHash is a cache-conscious hash table with string keys [8].
- HAT is a hybrid data structure of the burst trie [33] and ArrayHash [7].
- Judy is a trie-based dictionary implementation developed at Hewlett-Packard Research Labs [11].
- Cedar developed by Yoshinaga [64] is an efficient dictionary implementation based on dynamic double-array tries [2].

For ArrayHash and HAT, we used Tessil’s implementations [54, 55]. From the three implementation variations of Cedar, we took one based on a reduced trie [64] and one based on prefix trie [2], and denote them by Cedar-R and Cedar-P, respectively. Cedar-R is suitable for short keywords¹⁹, whereas Cedar-P is suitable for the general case.

We evaluated the working space and the running times in the same manner as Section 6.4. Figure 8 shows the experimental results for the four datasets GeoNames, AOL, Wiki, and DNA consisting of short keywords. Figure 9 shows the experimental results for the four datasets LUBMS, LUBML, UK, and WebBase consisting of long keywords. For our methods, we only plot the results of PDT-SB, PDT-CB, PDT-SFK and PDT-CFK, setting ℓ to 8, 16, or 32. To keep focus on the competitive contestants in the plots, we omitted some weaker instances, namely the DynPDT dictionaries with $\ell = 64$ and the dictionaries with PLM. The former are too slow, while the latter take too much working space. Only for DNA, we plotted the results of Cedar-R instead of Cedar-P because Cedar-R is superior on that instance. For LUBML and WebBase, we were not able to run our experiments with Cedar because the resulting number of trie nodes becomes too large to be representable in Cedar based on 32-bit pointers. For the long keywords (Figure 9), we omitted the results of ArrayHash because its working space is too large. For example, ArrayHash is 143% larger than HAT for LUBMS.

Based on Figure 8 showing the evaluation for short keywords, we can state the following observations:

- The DynPDT dictionaries are the smallest. PDT-CB for $\ell = 32$ is 25–48% smaller than the existing smallest data structures (Cedar-R for DNA and HAT for the others). PDT-CFK with $\ell = 32$ is 29–39% smaller than HAT for the datasets except DNA.
- Regarding the insert time, HAT is the fastest. Except for DNA, the DynPDT dictionaries based on FK-hash, PDT-SFK and PDT-CFK, are competitive to the other data structures.
- Regarding the lookup time, ArrayHash is the fastest. Except for DNA, the DynPDT dictionaries based on m-Bonsai, PDT-SB and PDT-CB, are competitive to Judy.
- For DNA consisting of short keywords, the DynPDT dictionaries are not efficient because the merits of the path decomposition applied to a trie with only short paths become negligible to

¹⁸All the experimental results are shown in Appendix A.

¹⁹We cannot be more concrete here since the efficiency of the heuristics of these data structures do not merely depend on the keyword lengths.

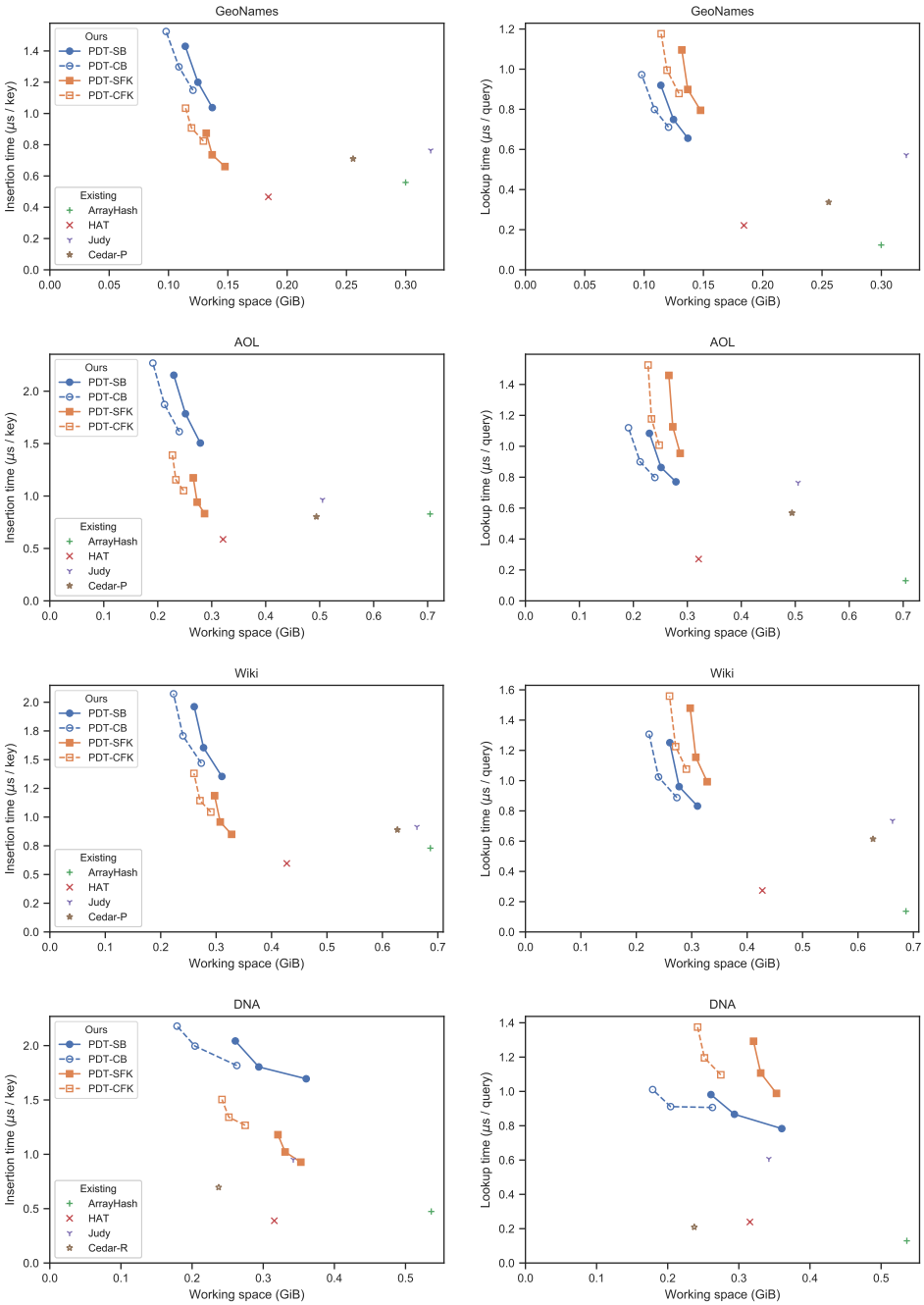


Fig. 8. Experimental results for short keywords with $\ell = 8, 16, \text{ and } 32$.

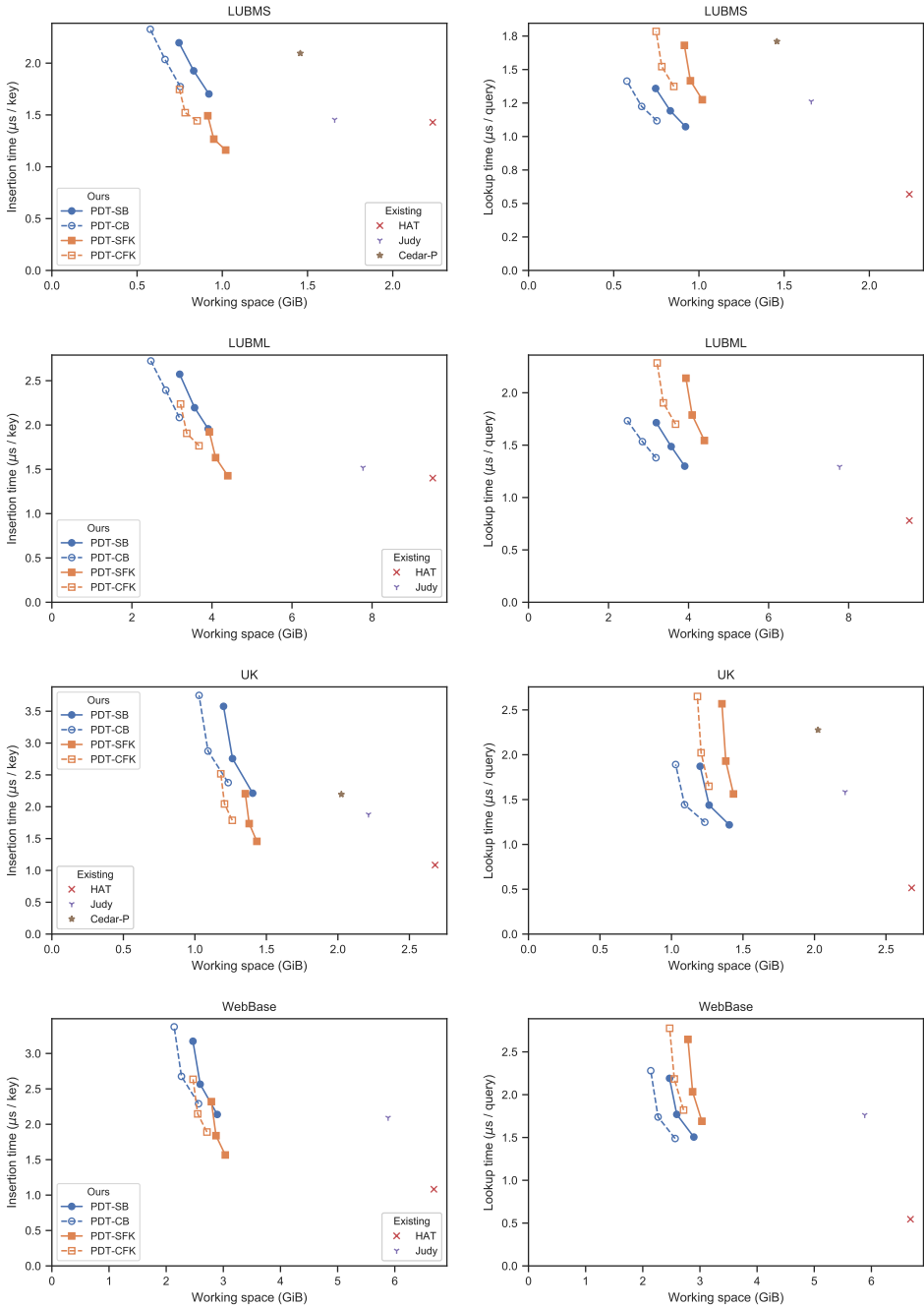


Fig. 9. Experimental results for long keywords with $\ell = 8, 16, \text{ and } 32$.

the additional burden of representing the trie with two separate data structures, one for its path-decomposed trie topology and one for its node labels.

Based on Figure 9 showing the evaluation for long keywords, we can state the following observations:

- The DynPDT dictionaries are the smallest for all the datasets. When $\ell = 32$, PDT-CB is 49–60% smaller than Cedar-P for LUBMS and UK, and is 64–68% smaller than Judy for LUBML and WebBase. When $\ell = 32$, PDT-CFK is 42–49% smaller than Cedar-P for LUBMS and UK, and is 58–59% smaller than Judy for LUBML and WebBase.
- Regarding the insert time, PDT-SFK is competitive to the other data structures.
- Regarding the lookup time, HAT is the fastest although its working space is large. Compared to PDT-SB with $\ell = 8$, HAT is 40–78% faster but 48–61% larger.
- In many cases, the DynPDT dictionaries outperform Judy and Cedar-P. For example, PDT-SFK with $\ell = 8$ is 48% smaller and 4–25% faster than Judy for WebBase. PDT-CB with $\ell = 8$ is 48% smaller and 15–35% faster than Cedar-P for LUBMS.

Summary. Throughout all dataset instances, DynPDT is the smallest data structure. Especially for long keywords such as URIs, our dictionaries are space-efficient and fast thanks to the path decomposition; however, they are not efficient for extremely short keywords because the path decomposition does not work well on such instances. In summary, DynPDT is useful for in-memory applications handling massive datasets consisting of long keywords.

For example, the RDF database system Diplodocus [62, 63] encodes every URI as an integer number through a dynamic keyword dictionary because the fixed-size integers can be handled more efficiently than the original strings having variable lengths. Since the encoding time is a significant part of the query execution time on the Diplodocus system, Mavlyutov et al. [46] experimentally compared a series of dynamic keyword dictionaries. Actually, LUBMS and LUBML of our datasets are exactly those evaluated in [46]. They concluded that HAT is a good data structure taking aspects like working space and time performance into account.²⁰ However, as demonstrated in our experiments, our DynPDT dictionaries can maintain the URI datasets in space up to 74% smaller than HAT, while keeping competitive insertion times. Although DynPDT’s slow lookup time is a drawback compared to HAT, maintaining massive RDF database systems in main-memory is essential, and we believe that DynPDT’s high memory efficiency will contribute to the future of Semantic Web applications.

7 CONCLUSION

We presented a novel data structure for dynamic keyword dictionaries — called DynPDT — which is applicable to scalable string data processing. For that, we applied path decomposition and utilized the recent hash-based trie representations m-Bonsai and FK-hash. We demonstrated with experiments on real-world massive datasets that the memory footprint of DynPDT is the smallest within a careful selection of efficient dynamic keyword dictionaries. It is especially efficient for long keywords due to the path decomposition approach.

Our results pave new ways for major improvements in various existing systems because the dynamic keyword dictionary problem is a common task in applications such as vocabulary accumulation for inverted-index construction [33], RDF database systems [62, 63], in-memory OLTP (online transaction processing) database systems [43], Web crawlers [59], and search engines [15, 16]. DynPDT can contribute to those systems especially by reducing their memory requirements. Although we have put the focus on the keyword dictionary problem in this paper, DynPDT as a general

²⁰Judy and Cedar were not evaluated in [46].

data structure is of independent interest, being useful for applications handling dynamic tries. An interesting application is the LZD compression [9, 27], a variation of the LZ78 compression [65]. Since the LZD algorithm maintains long factors (or strings) in a dynamic trie, we are confident that the incremental path decomposition on such a trie will have performance benefits.

Our future plans for DynPDT are as follows.

- The *burst trie* developed by Heinz et al. [33] maintains sparse subtrees in a trie in dynamic containers of strings by collapsing the subtrees. DynPDT would be suited as an alternative container representation to enhance the memory efficiency of the burst trie.
- In our experiments, we implemented the second data structure of the displacement array D_2 through the CHT by Cleary [17], following the original m-Bonsai approach [49]. Recently, Köppl et al. [42] developed space-efficient hash tables with separate chaining and compact hashing. Although the CHT needs additional displacement information (i.e., two bit arrays), his hash tables do not need such additional information. We expect that his hash tables are suitable representations of D_2 .

ACKNOWLEDGMENTS

We thank Kazuya Tsuruta for kindly providing us the implementations used in [58]. We thank the anonymous reviewers for their helpful comments. A part of this work was supported by JSPS KAKENHI Grant Numbers 17J07555 and JP18F18120.

REFERENCES

- [1] Daigle Alexandre. 2016. *Optimal path-decomposition of tries*. Ph.D. Dissertation. University of Waterloo.
- [2] Jun'ichi Aoe. 1989. An efficient digital search algorithm by using a double-array structure. *IEEE Transactions on Software Engineering* 15, 9 (1989), 1066–1077. <https://doi.org/10.1109/32.31365>
- [3] Diego Arroyuelo, Rodrigo Cánovas, Gonzalo Navarro, and Rajeev Raman. 2017. LZ78 compression in low main memory space. In *Proceedings of the 24th International Symposium on String Processing and Information Retrieval (SPIRE)*. 38–50. https://doi.org/10.1007/978-3-319-67428-5_4
- [4] Diego Arroyuelo, Pooya Davoodi, and Srinivasa Rao Satti. 2016. Succinct dynamic cardinal trees. *Algorithmica* 74, 2 (2016), 742–777. <https://doi.org/10.1007/s00453-015-9969-x>
- [5] Julian Arz and Johannes Fischer. 2018. Lempel–Ziv-78 compressed string dictionaries. *Algorithmica* 80, 7 (2018), 2012–2047. <https://doi.org/10.1007/s00453-017-0348-7>
- [6] Nikolas Askitis. 2007. *Efficient data structures for cache architectures*. Ph.D. Dissertation. RMIT University.
- [7] Nikolas Askitis and Ranjan Sinha. 2010. Engineering scalable, cache and space efficient tries for strings. *The VLDB Journal* 19, 5 (2010), 633–660. <https://doi.org/10.1007/s00778-010-0183-9>
- [8] Nikolas Askitis and Justin Zobel. 2005. Cache-conscious collision resolution in string hash tables. In *Proceedings of the 12th International Symposium on String Processing and Information Retrieval (SPIRE)*. 91–102. https://doi.org/10.1007/11575832_11
- [9] Golnaz Badkobeh, Travis Gagie, Shunsuke Inenaga, Tomasz Kociumaka, Dmitry Kosolobov, and Simon J Puglisi. 2017. On two LZ78-style grammars: Compression bounds and compressed-space computation. In *Proceedings of the 24th International Symposium on String Processing and Information Retrieval (SPIRE)*. 51–67.
- [10] Ricardo Baeza-Yates and Berthier Ribeiro-Neto. 2011. *Modern information retrieval* (2nd ed.). Vol. 463. Addison Wesley, Boston, MA, USA.
- [11] Doug Baskins. 2002. *A 10-minute description of how Judy arrays work and why they are so fast*. <http://judy.sourceforge.net/doc/10minutes.htm>
- [12] Djamel Belazzougui, Paolo Boldi, and Sebastiano Vigna. 2010. Dynamic z-fast tries. In *Proceedings of the 17th International Symposium on String Processing and Information Retrieval (SPIRE)*. 159–172. https://doi.org/10.1007/978-3-642-16321-0_15
- [13] Jon L. Bentley and Robert Sedgewick. 1997. Fast algorithms for sorting and searching strings. In *Proceedings of the 8th ACM-SIAM Symposium on Discrete Algorithms (SODA)*, Vol. 97. 360–369.
- [14] Paolo Boldi, Bruno Codenotti, Massimo Santini, and Sebastiano Vigna. 2004. Ubcrawler: A scalable fully distributed Web crawler. *Software: Practice and Experience* 34, 8 (2004), 711–726. <https://doi.org/10.1002/spe.587>
- [15] Brazil Inc. 2019. Groonga: An open-source fulltext search engine and column store. <http://groonga.org/>

- [16] Michael Busch, Krishna Gade, Brian Larson, Patrick Lok, Samuel Luckenbill, and Jimmy Lin. 2012. Earlybird: Real-time search at twitter. In *Proceedings of the 28th international conference on data engineering (ICDE)*. 1360–1369. <https://doi.org/10.1109/ICDE.2012.149>
- [17] John G. Cleary. 1984. Compact hash tables using bidirectional linear probing. *IEEE Trans. Comput.* 33, 9 (1984), 828–834. <https://doi.org/10.1109/TC.1984.1676499>
- [18] Armon Dadgar. 2012. Libart: Adaptive radix trees implemented in C. <https://github.com/armon/libart>
- [19] John J. Darragh, John G. Cleary, and Ian H. Witten. 1993. Bonsai: A compact representation of trees. *Software: Practice and Experience* 23, 3 (1993), 277–291. <https://doi.org/10.1002/spe.4380230305>
- [20] Richard Durstenfeld. 1964. Algorithm 235: random permutation. *Commun. ACM* 7, 7 (1964), 420. <https://doi.org/10.1145/364520.364540>
- [21] Paolo Ferragina, Roberto Grossi, Ankur Gupta, Rahul Shah, and Jeffrey Scott Vitter. 2008. On searching compressed string collections cache-obliviously. In *Proceedings of the 27th Symposium on Principles of Database Systems (PODS)*. 181–190. <https://doi.org/10.1145/1376916.1376943>
- [22] Johannes Fischer and Dominik Köppl. 2017. Practical evaluation of Lempel-Ziv-78 and Lempel-Ziv-Welch tries. In *Proceedings of the 24th International Symposium on String Processing and Information Retrieval (SPIRE)*. 191–207. https://doi.org/10.1007/978-3-319-67428-5_16
- [23] Edward Fredkin. 1960. Trie memory. *Commun. ACM* 3, 9 (1960), 490–499. <https://doi.org/10.1145/367390.367400>
- [24] Simon Gog, Timo Beller, Alistair Moffat, and Matthias Petri. 2014. From theory to practice: Plug and play with succinct data structures. In *Proceedings of the 13th International Symposium on Experimental Algorithms (SEA)*. 326–337. https://doi.org/10.1007/978-3-319-07959-2_28
- [25] Rodrigo González, Szymon Grabowski, Veli Mäkinen, and Gonzalo Navarro. 2005. Practical implementation of rank and select queries. In *Poster Proceedings of the 4th Workshop on Experimental and Efficient Algorithms (WEA)*. 27–38.
- [26] Google Inc. 2005. Sparsehash: C++ associative containers. <https://github.com/sparsehash/sparsehash>
- [27] Keisuke Goto, Hideo Bannai, Shunsuke Inenaga, and Masayuki Takeda. 2015. LZD factorization: Simple and practical online grammar compression with variable-to-fixed encoding. In *Proceedings of the 26th Annual Symposium on Combinatorial Pattern Matching (CPM)*. 219–230. https://doi.org/10.1007/978-3-319-19929-0_19
- [28] Popovitch Gregory. 2016. Sparsepp: A fast, memory efficient hash map for C++. <https://github.com/greg7mdp/sparsepp>
- [29] Roberto Grossi and Giuseppe Ottaviano. 2013. Design of practical succinct data structures for large data collections. In *Proceedings of the 12th International Symposium on Experimental Algorithms (SEA)*. 5–17. https://doi.org/10.1007/978-3-642-38527-8_3
- [30] Roberto Grossi and Giuseppe Ottaviano. 2014. Fast compressed tries through path decompositions. *ACM Journal of Experimental Algorithmics* 19, 1 (2014), Article 1.8. <https://doi.org/10.1145/2656332>
- [31] Yuanbo Guo, Zhengxiang Pan, and Jeff Heflin. 2005. LUBM: A benchmark for OWL knowledge base systems. *Web Semantics: Science, Services and Agents on the World Wide Web* 3, 2 (2005), 158–182. <https://doi.org/10.1016/j.websem.2005.06.005>
- [32] Harold Stanley Heaps. 1978. *Information retrieval: Computational and theoretical aspects*. Academic Press, Inc., Orlando, FL, USA.
- [33] Steffen Heinz, Justin Zobel, and Hugh E. Williams. 2002. Burst tries: A fast, efficient data structure for string keys. *ACM Transactions on Information Systems* 20, 2 (2002), 192–223. <https://doi.org/10.1145/506309.506312>
- [34] Jun Hirai, Sriram Raghavan, Hector Garcia-Molina, and Andreas Paepcke. 2000. WebBase: A repository of Web pages. *Computer Networks* 33, 1 (2000), 277–293. [https://doi.org/10.1016/S1389-1286\(00\)00063-3](https://doi.org/10.1016/S1389-1286(00)00063-3)
- [35] Bo-June Paul Hsu and Giuseppe Ottaviano. 2013. Space-efficient data structures for top-k completion. In *Proceedings of the 22nd International Conference on World Wide Web (WWW)*. 583–594. <https://doi.org/10.1145/2488388.2488440>
- [36] Jesper Jansson, Kunihiko Sadakane, and Wing-Kin Sung. 2015. Linked dynamic tries with applications to LZ-compression in sublinear time and space. *Algorithmica* 71, 4 (2015), 969–988. <https://doi.org/10.1007/s00453-013-9836-6>
- [37] Shunsuke Kanda. 2018. *Space- and time-efficient string dictionaries*. Ph.D. Dissertation. Tokushima University.
- [38] Shunsuke Kanda, Yuma Fujita, Kazuhiro Morita, and Masao Fuketa. 2018. Practical rearrangement methods for dynamic double-array dictionaries. *Software: Practice and Experience* 48, 1 (2018), 65–83. <https://doi.org/10.1002/spe.2516>
- [39] Shunsuke Kanda, Kazuhiro Morita, and Masao Fuketa. 2017. Compressed double-array tries for string dictionaries supporting fast lookup. *Knowledge and Information Systems* 51, 3 (2017), 1023–1042. <https://doi.org/10.1007/s10115-016-0999-8>
- [40] Shunsuke Kanda, Kazuhiro Morita, and Masao Fuketa. 2017. Practical implementation of space-efficient dynamic keyword dictionaries. In *Proceedings of the 24th International Symposium on String Processing and Information Retrieval (SPIRE)*. 221–233. https://doi.org/10.1007/978-3-319-67428-5_19
- [41] Donald E. Knuth. 1998. *The art of computer programming, 3: sorting and searching* (2nd ed.). Addison Wesley, Redwood City, CA, USA.

- [42] Dominik Köppl, Simon J Puglisi, and Rajeev Raman. 2020. Fast and simple compact hashing via bucketing. In *Proceedings of the 18th International Symposium on Experimental Algorithms (SEA)*. in press.
- [43] Viktor Leis, Alfons Kemper, and Thomas Neumann. 2013. The adaptive radix tree: ARTful indexing for main-memory databases. In *Proceedings of the IEEE 29th International Conference on Data Engineering (ICDE)*. 38–49. <https://doi.org/10.1109/ICDE.2013.6544812>
- [44] George Marsaglia. 2003. Xorshift RNGs. *Journal of Statistical Software* 8, 14 (2003), 1–6. <https://doi.org/10.18637/jss.v008.i14>
- [45] Miguel A. Martínez-Prieto, Nieves R. Brisaboa, Rodrigo Cánovas, Francisco Claude, and Gonzalo Navarro. 2016. Practical compressed string dictionaries. *Information Systems* 56 (2016), 73–108. <https://doi.org/10.1016/j.is.2015.08.008>
- [46] Ruslan Mavlyutov, Marcin Wylot, and Philippe Cudre-Mauroux. 2015. A comparison of data structures to manage URIs on the Web of data. In *Proceedings of the 12th European Semantic Web Conference (ESWC)*. 137–151. https://doi.org/10.1007/978-3-319-18818-8_9
- [47] Vincent Migliore, Benoît Gérard, Mehdi Tibouchi, and Pierre-Alain Fouque. 2019. Masking dilithium. In *Proceedings of the 17th International Conference on Applied Cryptography and Network Security (ACNS)*. 344–362. https://doi.org/10.1007/978-3-030-21568-2_17
- [48] Andreas Poyias, Simon J. Puglisi, and Rajeev Raman. 2017. Compact dynamic rewritable (CDRW) arrays. In *Proceedings of the 19th Workshop on Algorithm Engineering and Experiments (ALENEX)*. 109–119. <https://doi.org/10.1137/1.9781611974768.9>
- [49] Andreas Poyias, Simon J Puglisi, and Rajeev Raman. 2018. m-Bonsai: A practical compact dynamic trie. *International Journal of Foundations of Computer Science* 29, 08 (2018), 1257–1278. <https://doi.org/10.1142/S0129054118430025>
- [50] Nicola Prezza. 2017. A framework of dynamic data structures for string processing. In *Proceedings of the 16th International Symposium on Experimental Algorithms (SEA)*, Vol. 75. 11:1–11:15. <https://doi.org/10.4230/LIPIcs.SEA.2017.11>
- [51] Guy L. Steele Jr, Doug Lea, and Christine H. Flood. 2014. Fast splittable pseudorandom number generators. In *Proceedings of the 14th ACM International Conference on Object Oriented Programming Systems Languages & Applications (OOPSLA)*. 453–472. <https://doi.org/10.1145/2714064.2660195>
- [52] Takuya Takagi, Shunsuke Inenaga, Kunihiko Sadakane, and Hiroki Arimura. 2016. Packed compact tries: A fast and efficient data structure for online string processing. In *Proceedings of the 27th International Workshop on Combinatorial Algorithms (IWOCA)*. 213–225. https://doi.org/10.1007/978-3-319-44543-4_17
- [53] Tessil. 2016. Hopscotch-map: C++ implementation of a fast hash map and hash set using hopscotch hashing. <https://github.com/Tessil/hopscotch-map>
- [54] Tessil. 2017. Array-hash: C++ implementation of a fast and memory efficient hash map and hash set specialized for strings. <https://github.com/Tessil/array-hash>
- [55] Tessil. 2017. Hat-trie: C++ implementation of a fast and memory efficient HAT-trie. <https://github.com/Tessil/hat-trie>
- [56] Tessil. 2017. Robin-map: C++ implementation of a fast hash map and hash set using robin hood hashing. <https://github.com/Tessil/robin-map>
- [57] Kazuya Tsuruta, Dominik Köppl, Shunsuke Kanda, Yuto Nakashima, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. 2019. Dynamic packed compact tries revisited. *CoRR* (2019). arXiv:1904.07467
- [58] Kazuya Tsuruta, Dominik Köppl, Shunsuke Kanda, Yuto Nakashima, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. 2020. c-trie++: A dynamic trie tailored for fast prefix searches. In *Proceedings of the Data Compression Conference (DCC)*. in press.
- [59] Takanori Ueda, Koh Satoh, Daichi Suzuki, Kenji Uchida, Kousuke Morimoto, Sayaka Akioka, and Hayato Yamana. 2013. A parallel distributed Web crawler consisting of producer-consumer modules. *IPSJ Transactions on Database* 6, 2 (2013), 85–97.
- [60] Terry A. Welch. 1984. A technique for high-performance data compression. *IEEE Computer* 52 (1984). <https://doi.org/10.1109/MC.1984.1659158>
- [61] Hugh E. Williams and Justin Zobel. 1999. Compressing integers for fast file access. *Computer Journal* 42, 3 (1999), 193–201. <https://doi.org/10.1093/comjnl/42.3.193>
- [62] Marcin Wylot, Philippe Cudre-Mauroux, and Paul Groth. 2014. TripleProv: Efficient processing of lineage queries in a native RDF store. In *Proceedings of the 23rd International Conference on World Wide Web (WWW)*. 455–466. <https://doi.org/10.1145/2566486.2568014>
- [63] Marcin Wylot, Jigé Pont, Mariusz Wisniewski, and Philippe Cudré-Mauroux. 2011. dipLODocus[RDF] – short and long-tail RDF analytics for massive Webs of data. In *Proceedings of the 10th International Semantic Web Conference (ISWC)*. 778–793. https://doi.org/10.1007/978-3-642-25073-6_49
- [64] Naoki Yoshinaga and Masaru Kitsuregawa. 2014. A self-adaptive classifier for efficient text-stream processing. In *Proceedings of the 24th International Conference on Computational Linguistics (COLING)*. 1091–1102.

[65] Jacob Ziv and Abraham Lempel. 1978. Compression of individual sequences via variable-rate coding. *IEEE Transactions on Information Theory* 24, 5 (1978), 530–536. <https://doi.org/10.1109/TIT.1978.1055934>

A EXPERIMENTAL RESULTS

Within the same setting as in Section 6.5, we present an extended evaluation including the following contestants:

- STLHash is the hash table `std::unordered_map` of the C++ standard library.
- GoogleDense is the hash table implementation `google::dense_hash_map` of Google [26].
- Sparsepp is Gregory Popovitch’s space-efficient hash container implementation derived from Google’s sparse hash table [28].
- Hopscotch is Tessil’s hash table implementation using hopscotch hashing [53].
- Robin is Tessil’s hash table implementation using robin hood hashing [56].
- ART is Armon Dadgar’s implementation [18] of the adaptive radix tree [43].

Further, we include the following implementations, which are also used and studied in the experimental section of [57]:

- PCT-Bit is a packed compact trie using bit parallelism [52].
- PCT-Hash is a packed compact trie using additionally STLHash as a dictionary in each micro trie [52].
- ZFT is Tsuruta’s C++ implementation of the z-fast trie [12].
- CTrie++ is a trie [57] combining aspects of the z-fast trie with the packed compact trie.

Table 5 shows the results for the datasets consisting of short keywords (i.e., GeoNames, AOL, Wiki and DNA). Table 6 shows the results for the datasets consisting of long keywords (i.e., LUBMS, LUBML, UK and WebBase). In these tables, Space is the working space in GiB, Insert is the average insertion time in microseconds, and Lookup is the average lookup time in microseconds. For SLM of DynPDT, the results with $\ell = 16$ are shown. Concerning PCT-Bit, PCT-Hash, ZFT and CTrie++, we could not obtain some results for large datasets because the resulting trie was too large to fit into RAM.

Table 5. Experimental results for short keywords.

(a) GeoNames				(b) AOL			
	Space	Insert	Lookup		Space	Insert	Lookup
PDT-PB	0.32	0.59	0.53	PDT-PB	0.52	1.01	0.65
PDT-SB	0.12	1.20	0.75	PDT-SB	0.25	1.78	0.86
PDT-CB	0.11	1.30	0.80	PDT-CB	0.21	1.87	0.90
PDT-PFK	0.33	0.63	0.71	PDT-PFK	0.52	0.80	0.80
PDT-SFK	0.14	0.74	0.90	PDT-SFK	0.27	0.94	1.13
PDT-CFK	0.12	0.91	0.99	PDT-CFK	0.23	1.16	1.18
STLHash	0.58	0.44	0.24	STLHash	1.01	0.52	0.26
GoogleDense	0.73	0.37	0.12	GoogleDense	1.72	0.67	0.15
Sparsepp	0.42	0.58	0.15	Sparsepp	0.77	0.76	0.18
Hopscotch	0.84	0.40	0.10	Hopscotch	1.04	0.51	0.12
Robin	0.84	0.31	0.10	Robin	1.79	0.47	0.12
ArrayHash	0.30	0.56	0.12	ArrayHash	0.70	0.83	0.13
HAT	0.18	0.47	0.22	HAT	0.32	0.59	0.27
Judy	0.32	0.76	0.57	Judy	0.51	0.97	0.76
ART	0.59	0.85	0.56	ART	0.91	0.96	0.77
Cedar-R	0.47	0.68	0.39	Cedar-R	1.07	0.87	0.61
Cedar-P	0.26	0.71	0.34	Cedar-P	0.49	0.80	0.57
PCT-Bit	2.96	9.12	10.43	PCT-Bit	4.07	12.42	14.34
PCT-Hash	5.13	7.74	5.49	PCT-Hash	7.66	10.26	6.99
ZFT	1.13	3.25	2.42	ZFT	1.82	3.84	2.57
CTrie++	1.34	2.58	0.79	CTrie++	2.12	3.01	1.10
(c) Wiki				(d) DNA			
	Space	Insert	Lookup		Space	Insert	Lookup
PDT-PB	0.64	0.98	0.68	PDT-PB	0.84	1.18	0.65
PDT-SB	0.28	1.60	0.96	PDT-SB	0.29	1.80	0.87
PDT-CB	0.24	1.71	1.02	PDT-CB	0.20	2.00	0.91
PDT-PFK	0.67	0.79	0.86	PDT-PFK	0.80	0.85	0.88
PDT-SFK	0.31	0.96	1.15	PDT-SFK	0.33	1.02	1.11
PDT-CFK	0.27	1.14	1.22	PDT-CFK	0.25	1.34	1.20
STLHash	1.29	0.50	0.27	STLHash	1.02	0.91	0.34
GoogleDense	1.64	0.54	0.14	GoogleDense	1.25	0.24	0.09
Sparsepp	0.97	0.69	0.18	Sparsepp	0.67	0.50	0.13
Hopscotch	1.08	0.42	0.13	Hopscotch	1.50	0.27	0.07
Robin	1.83	0.41	0.12	Robin	1.50	0.26	0.08
ArrayHash	0.69	0.73	0.14	ArrayHash	0.54	0.47	0.13
HAT	0.43	0.60	0.27	HAT	0.32	0.39	0.24
Judy	0.66	0.92	0.74	Judy	0.34	0.95	0.61
ART	1.23	1.00	0.73	ART	1.01	0.65	0.63
Cedar-R	1.19	0.89	0.59	Cedar-R	0.24	0.70	0.21
Cedar-P	0.63	0.89	0.61	Cedar-P	0.31	0.67	0.24
PCT-Bit	5.67	11.85	13.79	PCT-Bit	7.05	8.44	9.91
PCT-Hash	10.11	9.48	7.09	PCT-Hash	8.45	5.44	6.86
ZFT	2.24	3.56	2.64	ZFT	2.57	3.30	2.88
CTrie++	2.92	3.01	1.09	CTrie++	2.50	2.55	0.78

Table 6. Experimental results for long keywords.

(a) LUBMS				(b) LUBML			
	Space	Insert	Lookup		Space	Insert	Lookup
PDT-PB	2.37	1.62	1.10	PDT-PB	10.1	1.43	1.00
PDT-SB	0.83	1.93	1.19	PDT-SB	3.6	2.20	1.49
PDT-CB	0.66	2.04	1.22	PDT-CB	2.8	2.39	1.53
PDT-PFK	2.46	1.09	1.14	PDT-PFK	10.7	1.32	1.39
PDT-SFK	0.95	1.27	1.42	PDT-SFK	4.1	1.63	1.79
PDT-CFK	0.78	1.52	1.52	PDT-CFK	3.4	1.91	1.90
STLHash	7.47	0.61	0.51	STLHash	32.9	0.67	0.59
GoogleDense	9.93	0.89	0.28	GoogleDense	40.1	0.99	0.32
Sparsepp	6.22	0.83	0.39	Sparsepp	27.3	0.89	0.44
Hopscotch	6.87	0.70	0.27	Hopscotch	41.2	1.01	0.27
Robin	9.87	0.61	0.26	Robin	41.2	0.69	0.28
ArrayHash	5.44	0.98	0.30	ArrayHash	21.9	1.06	0.32
HAT	2.23	1.43	0.57	HAT	9.5	1.40	0.78
Judy	1.66	1.45	1.26	Judy	7.8	1.52	1.29
ART	5.83	0.91	0.77	ART	25.8	1.05	0.93
Cedar-R	1.97	1.79	1.66	Cedar-R	n/a	n/a	n/a
Cedar-P	1.46	2.10	1.71	Cedar-P	n/a	n/a	n/a
PCT-Bit	n/a	n/a	n/a	PCT-Bit	n/a	n/a	n/a
PCT-Hash	n/a	n/a	n/a	PCT-Hash	n/a	n/a	n/a
ZFT	9.27	6.33	5.65	ZFT	n/a	n/a	n/a
CTrie++	8.13	4.25	2.43	CTrie++	n/a	n/a	n/a
(c) UK				(d) WebBase			
	Space	Insert	Lookup		Space	Insert	Lookup
PDT-PB	2.32	1.45	0.94	PDT-PB	5.4	1.66	1.19
PDT-SB	1.26	2.76	1.44	PDT-SB	2.6	2.57	1.77
PDT-CB	1.09	2.87	1.44	PDT-CB	2.3	2.68	1.74
PDT-PFK	2.32	1.27	1.24	PDT-PFK	5.7	1.36	1.42
PDT-SFK	1.38	1.74	1.93	PDT-SFK	2.9	1.84	2.03
PDT-CFK	1.21	2.04	2.02	PDT-CFK	2.6	2.15	2.18
STLHash	6.05	0.67	0.50	STLHash	16.3	0.64	0.55
GoogleDense	10.50	1.09	0.27	GoogleDense	19.5	0.82	0.28
Sparsepp	5.06	0.96	0.37	Sparsepp	13.5	0.83	0.43
Hopscotch	6.23	0.75	0.25	Hopscotch	20.3	0.93	0.24
Robin	9.23	0.63	0.25	Robin	20.3	0.64	0.26
ArrayHash	5.91	1.16	0.28	ArrayHash	10.4	0.98	0.29
HAT	2.68	1.08	0.51	HAT	6.7	1.08	0.55
Judy	2.21	1.88	1.59	Judy	5.9	2.09	1.76
ART	5.17	1.64	1.19	ART	14.0	1.76	1.45
Cedar-R	7.37	2.24	2.30	Cedar-R	n/a	n/a	n/a
Cedar-P	2.02	2.20	2.28	Cedar-P	n/a	n/a	n/a
PCT-Bit	18.05	25.92	33.49	PCT-Bit	n/a	n/a	n/a
PCT-Hash	n/a	n/a	n/a	PCT-Hash	n/a	n/a	n/a
ZFT	7.53	6.20	5.03	ZFT	19.6	5.77	5.06
CTrie++	8.17	4.75	2.86	CTrie++	23.5	5.12	3.12