

# Nearest Common Ancestors : A Survey and a new Distributed Algorithm

Stephen Alstrup  
Dept. of Theory  
IT University of Copenhagen  
Copenhagen, DK-2400  
Denmark  
stephen@it-c.dk

Haim Kaplan  
School of Computer Science  
Tel Aviv University  
Tel Aviv  
Israel  
haimk@math.tau.ac.il

Cyril Gavoille  
LaBRI  
Université de Bordeaux  
33405 Talence cedex France  
gavoille@labri.fr

Theis Rauhe  
Dept. of Theory  
IT University of Copenhagen  
Copenhagen, DK-2400  
Denmark  
theis@it-c.dk

## ABSTRACT

Several papers describe linear time algorithms to preprocess a tree, such that one can answer subsequent nearest common ancestor queries in constant time. Here, we survey these algorithms and related results. A common idea used by all the algorithms for the problem is that a solution for complete binary trees is straightforward. Furthermore, for complete binary trees we can easily solve the problem in a distributed way by labeling the nodes of the tree such that from the labels of two nodes alone one can compute the label of their nearest common ancestor. Whether it is possible to distribute the data structure into short labels associated with the nodes is important for several applications such as routing. Therefore, related labeling problems have received a lot of attention recently.

Previous optimal algorithms for nearest common ancestor queries work using some mapping from a general tree to a complete binary tree. However, it is not clear how to distribute the data structures obtained using these mappings. We conclude our survey with a new simple algorithm that labels the nodes of a rooted tree such that from the labels of two nodes alone one can compute in constant time the label of their nearest common ancestor. The labels assigned by our algorithm are of size  $O(\log n)$  bits where  $n$  is the number of nodes in the tree. The algorithm runs in  $O(n)$  time.

## Categories and Subject Descriptors

E.1 [Data Structures]: Distributed data structures, Trees;  
E.4 [Coding and Information Theory]: Data compaction

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SPAA'02, August 10-13, 2002, Winnipeg, Manitoba, Canada.  
Copyright 2002 ACM 1-58113-529-7/02/0008 ...\$5.00.

and compression; F.2 [Analysis of Algorithms and Problem Complexity]: General

## General Terms

Algorithms, Theory

## Keywords

Ancestor, Distributed, Labeling

## 1. INTRODUCTION

Let  $T$  be a rooted tree. A node  $x \in T$  is an *ancestor* of a node  $y \in T$  if the path from the root of  $T$  to  $y$  goes through  $x$ . A node  $v \in T$  is a *common ancestor* of  $x$  and  $y$  if it is an ancestor of both  $x$  and  $y$ . The *nearest common ancestor*, NCA, of two nodes  $x, y$  is the common ancestor of  $x$  and  $y$  whose distance to  $x$  (and to  $y$ ) is smaller than the distance to  $x$  of any other common ancestor of  $x$  and  $y$ . We denote the NCA of  $x$  and  $y$  by  $nca(x, y)$ . The problem of efficiently computing NCAs has been studied extensively over the last three decades in an online and offline setting, and in various models of computation (see e.g. [34, 2, 3, 27, 41, 11, 7]).

The problem of finding NCAs has numerous applications. A procedure solving it is used by algorithms for finding a maximum weighted matching in a graph [22], a minimum spanning tree in a graph [30, 14], and a dominator tree in a directed flow-graph [3]. It is also proved useful in several string algorithms [25, 16], for dynamic planarity testing [45], in network routing [43], for solving various geometric problems [9] including range searching [23], for finding evolutionary trees [17], and in bounded tree-width algorithms [10].

One of the most fundamental results on computing NCAs is that of Harel and Tarjan [27, 26]. Harel and Tarjan describe a linear time algorithm to preprocess a tree and build a data structure that allows subsequent NCA queries to be answered in constant time. Subsequent to this result of Harel and Tarjan several simpler algorithms with essentially the same properties but better constant factors have been proposed [39, 33, 23, 41, 8] (see also Section 2.2). These

algorithms, including the one of Harel and Tarjan, use the observation that it is rather easy to solve the problem when the input tree is a complete binary tree.

To solve the problem when the input is a complete binary tree one has to label the nodes by their index in an inorder traversal of the tree. If the tree has  $n$  nodes each such number occupies  $\ell = \lceil \log n \rceil$  bits<sup>1</sup>. We assume that the least significant bit is the rightmost and its index is zero. Let  $\text{inorder}(x)$  and  $\text{inorder}(y)$  be the inorder indexes of  $x$  and  $y$ , respectively. Let  $i$  be the maximum among 1) the index of the leftmost bit in which  $\text{inorder}(x)$  and  $\text{inorder}(y)$  differ, 2) the index of the rightmost 1 in  $\text{inorder}(x)$ , 3) the index of the rightmost 1 in  $\text{inorder}(y)$ . It is easy to prove by induction that the inorder number of  $\text{nca}(x, y)$  consists of leftmost  $\ell - i$  bits of  $\text{inorder}(x)$  (or  $\text{inorder}(y)$  if the maximum above was the index of the rightmost 1 in  $\text{inorder}(y)$ ) followed by a 1 and  $i - 1$  zeros.

Note that this algorithm for complete binary trees is distributed in the sense that it constructs the inorder index of  $\text{nca}(x, y)$  from the inorder indices of  $x$  and  $y$  alone without accessing the original tree or any other global data structure. In case the application does not identify a node by its inorder index then tables converting identifiers to inorder indices and vice versa have to be constructed. In the common case where a node is identified by a pointer to a structure representing it, we can simply store the inorder index as an extra field in that structure. This will make it possible to get the inorder index of a node giving a pointer to its corresponding structure in constant time. In order to map the inorder index of the NCA back to a pointer to the corresponding node, a table of  $O(2^\ell) = O(n)$  entries can easily be constructed.

The algorithms in [27, 39, 33, 23, 41, 8] for general trees all work using some mapping of the tree to a completely balanced binary tree, thereby exploiting the fact that for completely balanced binary trees the problem is easier. Different algorithms differ by the way they do the mapping. Unfortunately, all algorithms, as a result of doing the mapping, have to use some precomputed auxiliary data structures in addition to the labels of the nodes, in order to compute NCAs. Thus, in contrast with the algorithm for completely balanced binary trees, the algorithms for general trees do not allow to compute a unique identifier of  $\text{nca}(x, y)$  from short labels associated with  $x$  and  $y$  alone.

In this paper we show how to label the nodes of an arbitrary tree so that from the labels of  $x$  and  $y$  alone one can determine the label of  $\text{nca}(x, y)$ . In particular we give an algorithm that proves the following theorem.

**THEOREM 1.** *There is a linear time algorithm that labels the  $n$  nodes of a rooted tree  $T$  with labels of length  $O(\log n)$  bits such that from the labels of nodes  $x, y \in T$  alone, one can compute the label of  $\text{nca}(x, y)$  in constant time.*

Our algorithm is as simple as the simplest among the non distributed algorithms mentioned above. In a scenario when nodes have to be identified by some predetermined identifiers one can use our algorithm together with a table converting labels to those predetermined identifiers and vice versa. Thus our algorithm provides an alternative to any of the algorithms mentioned above also in a non distributed settings. Theorem 1 should be put in contrast with a recent lower

bound of Peleg [38]. Peleg shows that labels which allow to compute directly a  $\Theta(\log n)$  bit predetermined identifier of  $\text{nca}(x, y)$  from the labels of  $x$  and  $y$  alone must be of length  $\Omega(\log^2 n)$  bits. So to obtain a theorem like Theorem 1 one has to exploit the freedom to choose the identifiers of the nodes. A distinguished feature in our algorithm is its use of alphabetic codes [24] to generalize the inorder approach for complete binary trees.

A labeling scheme to identify NCAs may be useful in routing messages on tree networks. When a message has to be sent from a node  $x$  to a node  $y$  in a tree it has to go through  $\text{nca}(x, y)$ . Therefore the ability to compute the identifier of  $\text{nca}(x, y)$  from the identifiers of  $x$  and  $y$  may prove useful. Our particular NCA labeling scheme also allows to identify the first edge on the shortest path from  $x$  to  $y$  from the labels of  $x$  and  $y$  alone. In that sense it generalizes recent labeling schemes for routing on trees [18, 43] (see also Section 3). Another possible application of our algorithm arises in XML search engines. Such search engines typically maintain a reverse index. This index is a hash table mapping each word or a name of a tag to all XML-documents containing it. The engine can exploit the fact that an XML-document is essentially a tree and label each such tree using our algorithm. Then it can attach to each occurrence of a word in a document the label of the corresponding node. By doing that the engine can process sophisticated queries by accessing only the hash table rather than the documents themselves.

**Outline of the paper:** In Section 2 we survey previous algorithms for computing NCAs, including related results and applications. One of the most important applications of NCAs is range searching. Up to now the connection between range searching and NCA has been given through several reductions. In the survey we make observations which give a more direct relation between these two problems. In Section 4 we prove Theorem 1. Since our new algorithm makes it possible to identify the NCA from the two input nodes alone, its potential applications may avoid several lookups to perhaps expensive external memory. Further motivation for our new approach is given in Section 3.

## 2. NCA AND RELATED RESULTS

First we describe the relationship between NCAs and range searching, since most of the simple NCA algorithms are constructed using this relationship. Next we survey previous NCA algorithms for the RAM model of computation where the tree is static. Finally, we list results for finding NCAs in other models of computation and in cases where the topology of the tree can change. We also discuss few closely related problems.

### 2.1 NCA and discrete range searching

Gabow, Bentley, and Tarjan [23] observed that the one-dimensional Discrete Range Searching (DRS) problem is equivalent to the NCA problem. The DRS problem is defined as follows. Given a sequence of real numbers  $x_1, \dots, x_n$ , preprocess the sequence so that one can answer efficiently subsequent queries of the form “given a pair of indices  $(i, j)$ , what is the index of a maximum element among  $x_i, \dots, x_j$ ?”. We denote such query by  $\text{max}(i, j)$ . The DRS problem is a fundamental geometric searching problem and orthogonal range searching problems in two and higher dimensions are often

<sup>1</sup>All the logarithms are in base two.

reduced to DRS [23]. Gabow et al. (see also [32]) reduce the DRS problem to the NCA problem by constructing a Cartesian tree for the sequence  $x_1, \dots, x_n$ , as defined by Vuillemin [44]. The *Cartesian tree* of the sequence  $x_1, \dots, x_n$  is a binary tree with  $n$  nodes each containing a number  $x_i$ . Let  $x_j = \max\{x_1, \dots, x_n\}$ . The root of the Cartesian tree for  $x_1, \dots, x_n$  contains  $x_j$  (and possibly also the index  $j$ ). The left subtree of the root is a Cartesian tree for  $x_1, \dots, x_{j-1}$ , and the right subtree of the root is a Cartesian tree for  $x_{j+1}, \dots, x_n$ . Vuillemin shows how to construct the Cartesian tree of  $x_1, \dots, x_n$  in  $O(n)$  time. It is easy to see that the maximum among  $x_i, \dots, x_j$  corresponds to the NCA of the node containing  $x_i$  and the node containing  $x_j$ .

Gabow et al. also show how to reduce the NCA problem to the DRS problem. Given a tree we first construct a sequence of its nodes by doing a depth first traversal [42]. Each time we visit a node we add it to the end of the sequence so that each node appears in the sequence as many times as its degree. Note that this sequence is a prefix of an Euler tour of the tree. Let  $\text{depth}(x)$  be the depth of a node  $x$ . From the sequence of nodes we obtain a sequence of integers by replacing each occurrence of a node  $x$  by  $-\text{depth}(x)$ . To compute  $\text{nca}(x, y)$  we pick arbitrary two elements  $x_i, x_j$  in the sequence representing  $x$  and  $y$ , respectively, and compute the index of the maximum among  $x_i, \dots, x_j$ . It is easy to see that the node corresponding to this maximum element is  $\text{nca}(x, y)$ .

Combining the equivalence between the NCA problem and the DRS problem with Theorem 1 we obtain the following corollary.

**COROLLARY 2.** *Let  $x_1, \dots, x_n$  be a sequence of  $n$  real numbers. We can assign in linear time a label of length  $O(\log n)$  bits to each element, such that, given the labels of  $x_i, x_j$ , the label of a maximum among  $x_i, \dots, x_j$  can be computed in constant time from the labels of  $x_i$  and  $x_j$  alone.*

If the value  $x_i$  is required rather than its label we need one lookup in a global table to map a label to a real number. We can implement this mapping either using sufficient space, using linear space and expected linear preprocessing time [21], or using linear space with  $O(n \log n)$  deterministic preprocessing time [36].

## 2.2 Previous static algorithms for computing NCA

As mentioned in Section 1, Harel and Tarjan [27] were the first who described how to preprocess a tree in linear time such that one can answer NCA queries in constant time per query. However they already point out in [27] that one could simplify some parts of their algorithm. Subsequently, Schieber and Vishkin [41] contracted some of the preprocessing steps of Harel and Tarjan into a single step and obtained a simpler algorithm which they could also parallelize. Powell [39] describes a simplification of the algorithm of Schieber and Vishkin.

Berkman and Vishkin [8] describe a simpler algorithm which is based on the equivalence between the NCA problem and the DRS problem. Let  $x_1, \dots, x_n$  be the input to the DRS problem. We build a complete binary tree with leaves  $x_1, \dots, x_n$  from left to right. With each internal node we associate a *left value* which is the maximum among the values in its left subtree and a *right value* which is the maximum among the values in its right subtree. We also associate

with each leaf  $l$  two tables each with  $\log n$  entries, denoted by  $\text{right}_l$  and  $\text{left}_l$ . For a leaf  $l$  the value  $\text{right}_l(i)$  is the maximum among the right values of the  $i$  closest ancestors of  $l$ , for which  $l$  belongs to the left subtree of the ancestor. The value of  $\text{left}_l(i)$  is defined similarly. To answer a query of the form  $\max(i, j)$  we find the depth  $d$  of  $\text{nca}(x_i, x_j)$  using a simple NCA algorithm for complete binary trees. Then we return the maximum among  $\text{right}_{x_i}(d-1)$ , and  $\text{left}_{x_j}(d-1)$ . Obviously, the drawback of this simple algorithm is that it requires  $O(n \log n)$  preprocessing time and space.

To overcome this difficulty Berkman and Vishkin observed that the sequence of depths of nodes on an Euler tour of the tree (obtained by reducing the NCA problem to the DRS problem) consists of integers such that the absolute difference between consecutive integers is at most 1. Based on this observation they suggested a three level algorithm for the DRS (and the NCA) problem which runs in linear time. In their scheme the sequence is partitioned into *blocks* each of size  $O(\log n)$ , and each block is further partitioned into *microblocks* each of size  $O(\log \log n)$ . They use the observation that  $|x_i - x_{i+1}| = 1$ , for  $1 \leq i \leq n-1$ , to precompute the answers to all possible queries inside all possible microblocks in one large table. Then they apply the nonlinear algorithm described above to each sequence of microblocks in a block, and to the sequence of blocks, representing each block or microblock by the maximum element in it. Since the number of microblocks in a block is  $\log n / \log \log n$  and the number of blocks is  $n / \log n$  the resulting structures are of linear size. We answer a maximum query in constant time using the algorithm above on the relevant blocks and two lookups in the table of precomputed answers to queries in microblocks. Farach and Bender [33] show that in fact one can in linear time precompute all answers to queries inside blocks if the size of the block is  $\frac{1}{2} \log n$ . Therefore the microblock level of [8] is not needed and a simpler two level algorithm based using the same ideas exists.

The algorithm of Berkman and Vishkin implicitly reduces the NCA problem for general trees to the NCA problem for completely balanced binary trees. By solving the NCA problem for completely balanced binary trees they obtain an algorithm for DRS on inputs where  $|x_i - x_{i+1}| = 1$ , for  $1 \leq i \leq n-1$ . Then they reduce the general NCA problem to this special case of DRS. Using the reduction from DRS to NCA described in Section 1 this also implies a solution to DRS for general sequences.

Last, we observe that in fact one can solve the DRS problem for general sequences directly, without going through the long reduction sequence to the special case of DRS just described. The solution to DRS described above exploited the restriction that  $|x_i - x_{i+1}| = 1$  only to solve small subproblems of size  $\frac{1}{2} \log n$ . Alternatively, we can solve a DRS problem  $x_1, \dots, x_m$  of size  $m = \lfloor \log n \rfloor$  as follows. For each  $x_i$  we compute  $g(x_i) = \max\{k \mid k < i \text{ and } x_k > x_i\} \cup \{-1\}$ . It is not hard to see that one can compute  $g(x_i)$  for every  $1 \leq i \leq m$ , in  $O(m)$  time by maintaining the largest elements in all the suffixes of the prefix of the sequence so far traversed in a stack. Based on the values of  $g(x_i)$  we associate a label,  $l(x_i)$ , of size  $m$  bits with  $x_i$ , for every  $1 \leq i \leq m$ . The  $j$ th bit of  $l(x_i)$  corresponds to element  $x_j$ ,  $1 \leq j \leq m$ , and it is set if and only if  $j < i$ ,  $x_j > x_i$ , and for every  $x_k$ ,  $j < k < i$ ,  $x_k < x_j$ . We compute these labels recursively as follows. The label  $l(x_1)$  is 0. For  $i > 1$ ,  $l(x_i)$  is the same as  $l(x_{g(x_i)})$  but with bit  $g(x_i)$  also set. Now

for  $i < j$ , we find  $\max(i, j)$  from  $l(x_j)$  as follows. We clear all bits with index smaller than  $i$  in  $l(x_j)$ , getting a word  $w$ . Then we return  $x_j$  if  $w = 0$ , and otherwise we return  $x_{\text{lsb}(w)}$ , where  $\text{lsb}(w)$  is the index of the least significant bit in  $w^2$ . Using this technique for the blocks we obtain simpler DRS and NCA algorithms.

None of these algorithms however implies Theorem 1 or Corollary 2. If we try to distribute the data structures of any of the algorithms described so far so that it is possible to answer queries from the labels of the corresponding nodes alone, then the best labeling schemes we can get may contain labels of length  $\Omega(\log^2 n)$  bits.

A somewhat different and interesting approach for computing NCAs on trees of depth  $O(\log n)$  has been suggested by King [31]. King's algorithm labels each edge of the tree by either 0 or 1 randomly by flipping a fair coin. Then we label each node  $v$  with the concatenation of the labels of the edges on the path from the root to  $v$ . To find  $\text{nca}(x, y)$  we first find the length,  $d$ , of the longest common prefix of the label of  $x$  and the label of  $y$ . Let  $\text{ancestor}(x, d)$  be the ancestor of depth  $d$  of  $x$ . If  $\text{ancestor}(x, d)$  is also an ancestor of  $y$  then it is the NCA of  $x$  and  $y$ . Otherwise for  $i = 1, 2, \dots$  we check whether  $\text{ancestor}(x, d - i)$  is an ancestor of  $y$  until we find such a node which must also be the NCA of  $x$  and  $y$ . It is easy to see that we would have to perform  $k$  ancestor queries with probability at most  $1/2^k$ . Therefore, on average, we need to perform a constant number of ancestor queries that involve ancestors of  $x$  at a particular level. The drawback of this scheme is that it requires external data structures to compute the ancestor of a node with a specified depth and to answer ancestor queries (See e.g. [40] and [13, 4]).

### 3. LABELING SCHEMES FOR DISTRIBUTED ENVIRONMENT

Motivated by applications in the construction of XML search engines and network routing [1, 29, 37] labeling schemes that allow ancestor queries have recently been developed. Santoro and Khatib [40] suggested to label the leaves of the tree from left to right by consecutive integers, and then to label each internal node by the pair of the labels of its leftmost and rightmost leaf descendants. One can then answer ancestor queries by checking if the corresponding intervals are nested. Clearly the maximum length of a label according to this scheme is  $2 \log n$  bits. Recently, Alstrup and Rauhe [6] building upon the work of Abiteboul, Kaplan, and Milo [1] gave a more complicated recursive labeling scheme for ancestor queries that generates labels of length at most  $\log n + O(\sqrt{\log n})$  bits. A lower bound showing that  $\log n + \Omega(\log \log n)$  bits are needed is given in [5]. This lower bound also holds for NCA queries which are more general.

Labeling schemes for parent queries, and other functions have also been studied [28, 29, 38]. Unfortunately neither of these labeling schemes allows to identify the nearest common ancestor of  $x$  and  $y$  when  $x$  and  $y$  are unrelated. The labeling scheme that we describe in this paper do allow this extra functionality.

In [43] and [18] it is demonstrated how compact labelings

<sup>2</sup>As previous NCA algorithms (e.g. [27, 8]) we assume that bit operations, on words of size  $\log n$ , can be precomputed in  $O(n)$  time if they are not supported to begin with.

of a tree can improve routing in trees and graphs. In the routing problem one assigns two labels to every node of the graph. The first label is the *address* of the node whereas the second label is a data structure called *local routing table*. The labels are assigned such that at every source node  $x$ , and for every destination node  $y$  one can find the first edge (or an identifier of that edge) on a shortest path from  $x$  to  $y$  from the local routing table of  $x$  and the address of  $y$ . So the path from  $x$  to  $y$  is built in a distributed way by all intermediate nodes encountered along the way. The goal is to obtain such labeling with labels as small as possible.

Cowen in [12] shows how to construct a labeling for trees that uses  $3 \log n$  bits for the addresses and  $O(\min\{d \log n, \sqrt{n} \log n\})$  bits for the local routing table, where  $d$  is the degree of the node. In [15], Eilam et al. show that any labeling in which the length of the addresses is at most  $\log n$  bits will require  $\Omega(\sqrt{n})$  bits for the routing table in some trees. However, Gavioille and Fraigniaud [18] showed that  $c \log n$  bits, for a small constant  $c$ , are enough to encode both the address and the local routing table of each node of the tree. Thorup and Zwick [43] even show that both the address and the local routing table of each node can be encoded together using  $\log n + O(\log n / \log \log n)$  bits (so  $c = 1 + o(1)$  suffices). Thorup and Zwick also show how to use their tree labeling to construct labeling for general graphs based on tree covering.

The labeling schemes for routing mentioned above do not allow to determine NCAs without some extensions. The labeling scheme that we suggest in this paper, however, allows not only to identify  $\text{nca}(x, y)$  from the labels of  $x$  and  $y$  but also to identify the first edge on the shortest path from  $x$  to  $y$ . Therefore it is more general than the routing labeling schemes.

A particular routing context where a labeling scheme for NCA may prove useful is the design of routing labeling schemes for graphs that are close to a tree, say  $c$ -decomposable graphs [19, 20] for  $c = O(1)$ . These graphs admit a tree-decomposition  $T$  where each node of  $T$  represents a separator of the graph of size at most  $c$ . So, one can use a compact labeling of  $T$  in order to determine the nearest separator  $S$  between the source  $x$  and the destination  $y$ . If  $x$  is contained in the component  $X \in T$ , and  $y$  is contained in the component  $Y \in T$ , then  $S$  is precisely the nearest common ancestor of  $X$  and  $Y$ . Since, any path from  $x$  to  $y$  has to cross some node of  $S$ , computing  $\text{nca}(X, Y)$  from the labels of  $x$  and  $y$  may simplify the local routing tables.

## 4. AN NCA ALGORITHM FOR DISTRIBUTED ENVIRONMENT

### 4.1 Preliminaries

We denote by  $\langle y \rangle_k$  a sequence of objects  $y_1, y_2, \dots, y_k$  (such as integers or binary strings). For binary strings  $a, b \in \{0, 1\}^*$ ,  $a <_{\text{lex}} b$  if and only if  $a$  precedes  $b$  in the lexicographic order on binary strings. I.e.,  $a$  is prefix of  $b$  or the first bit in which  $a$  and  $b$  differ is 0 in  $a$  and 1 in  $b$ . An *alphabetic* sequence  $\langle y \rangle_k$  is a sequence of binary strings  $\langle b \rangle_k$ ,  $b_i \in \{0, 1\}^*$ , where  $b_i <_{\text{lex}} b_j$ , for all  $1 \leq i < j \leq k$ . Let  $|s|$  denote the length of a binary string  $s \in \{0, 1\}^*$ . Observe that given machine words that contain  $a$ ,  $b$ ,  $|a|$ , and  $|b|$ , respectively, in their least significant bits, it is possible

to determine whether  $a <_{\text{lex}} b$  in a constant number of operations. (For instance, first align  $a$  and  $b$  by shifting to the left the smallest string, and then use standard integer comparison operators on the resulting words. If the two words are then equal then we break the tie according to the length of  $a$  and  $b$ ). When the strings  $\langle b \rangle_k$  are also prefix-free (no string is a prefix of another) we call  $\langle b \rangle_k$  an *alphabetic code*. The following lemma due to Gilbert and Moore [24] states the result which we need for alphabetic codes.

**LEMMA 3** (GILBERT AND MOORE [24]). *A sequence  $\langle y \rangle_k$  of positive numbers with  $n = \sum_{i=1}^k y_i$  has an alphabetic code  $\langle b \rangle_k$  where  $|b_i| \leq \log n - \log y_i + O(1)$  for all  $i$ .*

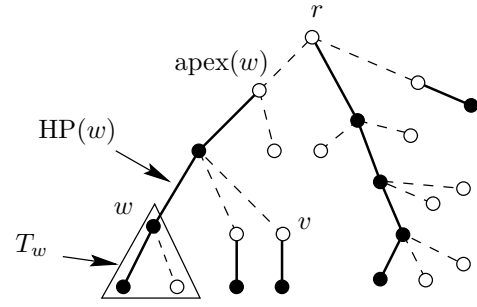
For our purposes it would suffice that  $\langle b \rangle_k$  is alphabetic (strings can be prefixes of one another). We can construct an alphabetic sequence satisfying the length bounds in Lemma 3 for an integer sequence  $\langle y \rangle_k$  in  $O(k)$  time as follows. Let  $s_i = \sum_{j=1}^i y_j$ ,  $i = 1, \dots, k-1$ , and let  $s_0 = 0$ . Also let  $I_i = [s_{i-1} + 1, s_i]$ ,  $i = 1, \dots, k$ , and  $f_i = \lfloor \log y_i \rfloor$ . In the interval  $I_i$  there must be a number  $z_i$  such that  $z_i \bmod 2^{f_i} = 0$ . If  $s_{i-1} + 1 \bmod 2^{f_i} = 0$ , then  $z_i = s_{i-1} + 1$ . Otherwise  $z_i = s_{i-1} + 1 - ((s_{i-1} + 1) \bmod 2^{f_i}) + 2^{f_i}$ . Hence,  $z_i$  can be represented in a word with  $w = \lceil \log n \rceil$  bits, having the  $f_i$  less significant bits set to 0. Then we can let  $b_i$  be the bit string consisting of the  $w - f_i$  most significant bit from  $z_i$ . Thus we get an alphabetic sequence where  $1 \leq |b_i| \leq \log n - \log y_i + 2$ .

The algorithm runs in  $O(k)$  time if machine operations to compute shifting, modulus, and discrete logarithm on  $O(\log n)$  bit words are supported. In a machine that does not support such operations we can use  $O(n)$  preprocessing time and space construct a table representing these functions. This will only increase the preprocessing time of our labeling algorithm which is described in Section 4.2 by a constant factor. Note that Mehlhorn [35] gives a somewhat more complicated algorithm which produces an *alphabetic code*,  $\langle b \rangle_k$ , for an arbitrary sequence  $\langle y \rangle_k$  of positive real numbers with the same bound on the lengths, i.e.,  $|b_i| \leq \log n - \log y_i + 2$ . Mehlhorn's algorithm can be implemented to run in  $O(k)$  time.

## 4.2 Labeling for NCA queries

As in [27] we divide the tree into disjoint paths. For a tree  $T$ , let  $|T|$  denote the number of nodes of  $T$ . Let  $T_v$  be the subtree rooted by  $v$ , and let  $\text{size}(v) = |T_v|$ . Let  $\text{parent}(v)$  be the parent of  $v$ , and  $\text{children}(v)$  be the set of children of  $v$ . We classify each node of  $T$  as either *heavy* or *light* as follows. The root is light. For each internal node  $v$ , we pick a child  $w$  of  $v$ , where  $\text{size}(w) = \max \{ \text{size}(z) \mid z \in \text{children}(v) \}$  and classify  $w$  as *heavy*. We classify each of the remaining children of  $v$  as *light*. We call an edge to a light child a *light edge*, and an edge to a heavy child a *heavy edge*. For a node  $v$  with a heavy child  $w$ , let  $\text{lsize}(v) = \text{size}(v) - \text{size}(w)$ . The nearest ancestor of  $v$  which is light (possibly  $v$  itself if  $v$  is light) is denoted by  $\text{apex}(v)$ . By removing the light edges  $T$  is partitioned into paths, which we call *heavy paths*. A node  $w$  belongs to the same heavy path as the nodes of the set  $\text{HP}(w) = \{ v \mid v \in T, \text{apex}(v) = \text{apex}(w) \}$ . See Figure 1.

First we compute  $\text{size}(v)$ ,  $\text{lsize}(v)$ ,  $\text{apex}(v)$  and the partition of  $T$  into heavy paths. It is easy to see that one can compute this information in linear time. Next we assign a heavy label,  $\text{hlabel}(v)$ , to each node  $v \in T$ , and a light label



**Figure 1:** Heavy and light nodes are black and white points respectively. Heavy paths are subtrees composed of heavy edges only (solid edge). In this example,  $\text{apex}(v) = v$ .

$\text{label}(v)$  to each light node  $v$ . These labels are defined as follows.

For the root  $r$ ,  $\text{label}(r)$  is the empty string. Then, for each light node  $w \neq r$ ,  $\text{label}(w)$  is a binary string such that

- Light Label:  
 $\text{label}(w) \notin \{ \text{label}(z) \mid z \neq w, z \in \text{children}(p) \}$ , where  $p = \text{parent}(w)$ .

Let  $w$  be any node. Then  $\text{hlabel}(w)$  is a binary string such that

- Heavy label:  
 $\text{hlabel}(w) <_{\text{lex}} \min_{\text{lex}} \{ \text{hlabel}(z) \mid z \neq w, z \in T_w \cap \text{HP}(w) \}$ , where  $\text{lex}$  is the lexicographic order of two strings.

We assign heavy and light labels using the algorithm to assign alphabetic strings to a sequence of positive numbers described in Section 4.1. We associate heavy labels with the nodes on each heavy path separately. We associate the weight  $\text{lsize}(v)$  with each node  $v$  on a heavy path and construct the alphabetic strings for the nodes on the path using the algorithm in Section 4.1. We assign light labels similarly. Here we apply the algorithm of Section 4.1 to the light children of each node  $u$  separately. When applying the algorithm we associate the weight  $\text{size}(v)$  with each light child  $v$  of  $u$ . We remark that light labels have no alphabetic restriction imposed on them so in fact one can obtain even shorter light labels using a simpler algorithm.

Next we assign a label  $l(v)$  to each node  $v \in T$  topdown as follows. We define  $l(\text{parent}(r))$ , and  $\text{label}(r)$ , to be the empty string. Then for every node  $v$ ,  $l(v) = l(\text{parent}(\text{apex}(v))) \cdot \text{label}(\text{apex}(v)) \cdot \text{hlabel}(v)$  (We use the  $\cdot$  operator for concatenation of strings).

It follows from our definition that a label  $l(v)$  consists of the concatenation of alternating heavy and light labels, thus  $l(v) = h_1 \cdot l_1 \cdot h_2 \cdot l_2 \cdot \dots$ . For a bit string  $s$ , we let  $s[i]$  be the  $i^{\text{th}}$  bit of  $s$ . In addition to  $l(v)$  we need a label  $k(v)$  of the same length, where  $k(v)[i] = 1$  if and only if  $l(v)[i]$  is the beginning of either a light or heavy label. The label  $k(v)$  is well defined because  $|\text{label}(v)|, |\text{hlabel}(v)| \geq 1$ . The label,  $\text{label}(v)$ , assigned to a node consists of the concatenation of  $l(v)$  and  $k(v)$ . Once heavy and light labels of all the nodes have been computed, one can perform a depth first traversal [42] of the tree and compute  $l(v)$  and  $k(v)$  in linear time.

Before showing how to compute the NCA given two labels we bound the length of the labels. Using the alphabetic strings from the preliminary section, we have

- $|\text{label}(w)| \leq \log \text{lsize}(\text{parent}(w)) - \log \text{size}(w) + O(1)$ , and
- $|\text{hlabel}(w)| \leq \log \sum_{v \in \text{HP}(w)} \text{lsize}(v) - \log \text{lsize}(w) + O(1)$

LEMMA 4. For any node  $v \in T$ ,  $|\text{label}(v)| = O(\log n)$ .

PROOF. First, note that  $|\text{label}(w)| = |l(w)| + |k(w)| = 2|l(w)|$ . The number of heavy and light labels in  $l(w)$  is at most  $\log n$ , since there is at most  $\log n$  light edges on a path from a leaf to the root. Hence, the additive  $O(1)$  terms in the light and heavy labels of  $l(w)$  sum to  $O(\log n)$ . So, in the remaining analysis we disregard the additive constant term. We prove by induction on the depth of  $w$  that  $|l(w)| \leq \log n - \log \text{lsize}(w)$ . For a node  $w$ , where  $\text{apex}(w) = r$ , in a tree with root  $r$ ,  $|l(w)| = |\text{hlabel}(w)| \leq \log n - \log \text{lsize}(w)$  so the statement holds. For another node  $w$ , where  $\text{apex}(w) \neq r$ ,  $l(w) = l(\text{parent}(\text{apex}(w))) \cdot \text{hlabel}(\text{apex}(w)) \cdot \text{hlabel}(w)$ . Now, by induction,  $|l(\text{parent}(\text{apex}(w)))| \leq \log n - \log \text{lsize}(\text{parent}(\text{apex}(w)))$ , and also by the definitions of the heavy labels and the light labels we have  $|\text{hlabel}(\text{apex}(w))| \leq \log \text{lsize}(\text{parent}(\text{apex}(w))) - \log \text{size}(\text{apex}(w))$ , and  $|\text{hlabel}(w)| \leq \log \sum_{v \in \text{HP}(w)} \text{lsize}(v) - \log \text{lsize}(w)$ . By summing up the last three bounds and using the fact that  $\sum_{v \in \text{HP}(w)} \text{lsize}(v) \leq \text{size}(\text{apex}(w))$  we obtain the lemma.  $\square$

The next lemma indicates how to answer NCA queries.

LEMMA 5. Let  $x$  and  $y$  be two vertices of  $T$ ,

- 1 If  $l(x) = h_1 \cdot l_1 \cdots h_i \cdot l_i \cdot t$  and  $l(y) = h_1 \cdot l_1 \cdots h_i \cdot l'_i \cdot t'$ , where  $l_i \neq l'_i$  or  $l_i \cdot t$  is empty or  $l'_i \cdot t'$  is empty, then  $l(\text{nca}(x, y)) = h_1 \cdot l_1 \cdots h_{i-1} \cdot l_{i-1} \cdot h_i$ .
- 2 If  $l(x) = h_1 \cdot l_1 \cdots h_i \cdots$  and  $l(y) = h_1 \cdot l_1 \cdots h'_i \cdots$ , where  $h'_i \neq h_i$ , then  $l(\text{nca}(x, y)) = h_1 \cdot l_1 \cdots h_{i-1} \cdot l_{i-1} \cdot \min_{\text{lex}} \{h_i, h'_i\}$ .

PROOF. Let  $z = \text{nca}(x, y)$ . By definition  $l(\text{parent}(\text{apex}(z)))$  is a prefix of both  $l(x)$  and  $l(y)$ . Let  $w$  be the heavy child of  $z$ . If  $x \in T_a$ ,  $y \in T_b$ ,  $a, b \in \text{children}(z) \setminus \{w\}$  or  $x$  is an ancestor to  $y$ , then Case 1 occurs; otherwise Case 2 occurs.  $\square$

To implement the calculation of  $\text{nca}(x, y)$  described in Lemma 5 we need to be able to identify the maximum  $i$  such that both  $l(x)$  and  $l(y)$  have  $h_1 \cdot l_1 \cdots h_i \cdot l_i$  as a prefix. To that end we calculate the length  $j_1$  of the maximum common prefix of  $l(x)$  and  $l(y)$ , and the length  $j_2$  of the maximum common prefix of  $k(x)$  and  $k(y)$ . Then we take  $j = \min \{j_1, j_2\}$  and calculate whether this bit ( $j$ th from the left) is part of a light label or a heavy label. To easily determine whether bit  $j$  occurs within a light label or a heavy label we can for example add to the label of  $v$  a mask containing a 1 in every bit that belongs to a light label or alternatively use a table that counts the parity of the number of bits set to one in a word. Once we identify whether  $j$  occurs within a light or heavy label we can extract the label of  $\text{nca}(x, y)$  using straightforward bit manipulations.

## Remarks:

1) In case we want to solve the unrestricted DRS problem using a Cartesian tree, we can let all the light labels be an empty string, since a Cartesian tree is binary. This will improve the constant factor for the length of the labels.

2) If we use an alphabetic code for the heavy labels (no heavy label on a path is a prefix of another) and any prefix code for the light labels then  $l(v)$  uniquely identifies  $v$ , and there is no need to look at  $k(v)$  in order to figure out the maximum  $i$  such that  $h_1 \cdot l_1 \cdots h_i \cdot l_i$  is a prefix of both  $l(x)$  and  $l(y)$ .

## 5. REFERENCES

- [1] S. Abiteboul, H. Kaplan, and T. Milo. Compact labeling schemes for ancestor queries. In *12<sup>th</sup> Annual ACM Symposium on Discrete Algorithms (SODA)*, pages 547–556, 2001.
- [2] A. V. Aho, Y. Sagiv, T. G. Szymanski, and J. D. Ullman. Inferring a tree from lowest common ancestors with an application to the optimization of relational expressions. *SIAM Journal on Computing*, 10(3):405–421, 1981.
- [3] A.V. Aho, J.E. Hopcroft, and J.D. Ullman. On finding lowest common ancestor in trees. *SIAM Journal on Computing*, 5(1):115–132, 1976. See also STOC 1973.
- [4] S. Alstrup and J. Holm. Improved algorithms for finding level ancestors in dynamic trees. In *27<sup>th</sup> International Colloquium on Automata, Languages and Programming (ICALP)*, volume 1853 of LNCS, pages 73–84. Springer-Verlag, 2000.
- [5] S. Alstrup and T. Rauhe. Labeling scheme lower bounds for ancestor, sibling and connectivity. Technical Report TR-2001-10, IT University of Copenhagen, 2001.
- [6] S. Alstrup and T. Rauhe. Improved labeling scheme for ancestor queries. In *13<sup>th</sup> Annual ACM Symposium on Discrete Algorithms (SODA)*, 2002.
- [7] S. Alstrup and M. Thorup. Optimal pointer algorithms for finding nearest common ancestors in dynamic trees. *Journal of Algorithms*, 35(2):169–188, 2000.
- [8] O. Berkman and U. Vishkin. Recursive star-tree parallel data structure. *SIAM Journal on Computing*, 22(2):221–242, 1993.
- [9] S. Carlsson and B. J. Nilsson. Computing vision points in polygons. *Algorithmica*, 24(1):50–75, 1999.
- [10] S. Chaudhuri and C. D. Zaroliagis. Shortest paths in digraphs of small treewidth. Part II: Optimal parallel algorithms. *Theoretical Computer Science*, 203(2):205–223, August 1998.
- [11] R. Cole and R. Hariharan. Dynamic lca queries on trees. In *10<sup>th</sup> Annual ACM Symposium on Discrete Algorithms (SODA)*, 1999.
- [12] L. J. Cowen. Compact routing with minimum stretch. *Journal of Algorithms*, 38:170–183, 2001.
- [13] P. F. Dietz. Finding level-ancestors in dynamic trees. In *2nd Work. on Algo. and Data Struct.*, volume 1097 of LNCS, pages 32–40, 1991.
- [14] B. Dixon, M. Rauch, and R. E. Tarjan. Verification and sensitivity analysis of minimum spanning trees in

- linear time. *SIAM Journal on Computing*, 21(6):1184–1192, 1992.
- [15] T. Eilam, C. Gavoille, and D. Peleg. Compact routing schemes with low stretch factor. In *17<sup>th</sup> Annual ACM Symposium on Principles of Distributed Computing (PODC)*, pages 11–20, August 1998.
  - [16] M. Farach-Colton. Optimal suffix tree construction with large alphabets. In *38<sup>th</sup> Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 137–143, 1997.
  - [17] M. Farach-Colton, S. Kannan, and T. Warnow. A robust model for finding optimal evolutionary trees. *Algorithmica*, 13(1/2):155–179, 1995.
  - [18] P. Fraigniaud and C. Gavoille. Routing in trees. In *28<sup>th</sup> International Colloquium on Automata, Languages and Programming (ICALP)*, volume 2076 of LNCS, pages 757–772, 2001.
  - [19] G. N. Frederickson and R. Janardan. Separator-based strategies for efficient message routing. In *27<sup>th</sup> Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 428–437, October 1986.
  - [20] G. N. Frederickson and R. Janardan. Space-efficient message routing in  $c$ -decomposable networks. *SIAM Journal on Computing*, 19(1):164–181, February 1990.
  - [21] M. L. Fredman, J. Komlós, and E. Szemerédi. Storing a sparse table with  $O(1)$  worst case access time. *Journal of the ACM*, 31(3):538–544, July 1984.
  - [22] H. N. Gabow. Data structure for weighted matching and nearest common ancestors with linking. In *1<sup>st</sup> Annual ACM Symposium on Discrete Algorithms (SODA)*, pages 434–443, 1990.
  - [23] H. N. Gabow, J. L. Bentley, and R. E. Tarjan. Scaling and related techniques for geometry problems. In *16<sup>th</sup> Annual ACM Symposium on Theory of Computing (STOC)*, pages 135–143, 1984.
  - [24] E. N. Gilbert and E. F. Moore. Variable-length binary encodings. Technical report, Bell System Technical Journal, 1959.
  - [25] D. Gusfield. *Algorithms on Strings, Trees, and Sequences*. Cambridge University Press, 1997.
  - [26] D. Harel. A linear time algorithm for the lowest common ancestors problem (extended abstract). In *21<sup>st</sup> Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 308–319, 1980.
  - [27] D. Harel and R. E. Tarjan. Fast algorithms for finding nearest common ancestors. *SIAM Journal Computing*, 13(2):338–355, 1984.
  - [28] S. Kannan, M. Naor, and S. Rudich. Implicit representation of graphs. *SIAM Journal on Discrete Mathematics*, 5(4):596–603, 1992. Preliminary version appeared in STOC 1988.
  - [29] H. Kaplan and T. Milo. Short and simple labels for small distances and other functions. In *7<sup>th</sup> International Workshop on Algorithms and Data Structures (WADS)*, volume 2125 of LNCS, pages 32–40. Springer, August 2001.
  - [30] D. R. Karger, P. N. Klein, and R. E. Tarjan. A randomized linear-time algorithm to find minimum spanning trees. *Journal of the ACM*, 42(2):321–328, 1995.
  - [31] V. King. private communication, July 1999.
  - [32] Janos Komlos. Linear verification for spanning trees. *Combinatorica*, pages 57–65, 1985.
  - [33] M. Farach-Colton M. A. Bender. The LCA problem revisited. In *4<sup>th</sup> LATIN*, pages 88–94, 2000.
  - [34] D. Maier. A space efficient method for the lowest common ancestor problem and an application to finding negative cycles. In *18<sup>th</sup> Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 132–141, 1977.
  - [35] K. Mehlhorn. A best possible bound for the weighted path length of binary search trees. *SIAM Journal on Computing*, 6(2):235–239, June 1977.
  - [36] R. Pagh. Faster deterministic dictionaries. In *11<sup>th</sup> Annual ACM Symposium on Discrete Algorithms (SODA)*, pages 487–49, 2000.
  - [37] D. Peleg. Proximity-preserving labeling schemes and their applications. In *25<sup>th</sup> International Workshop, Graph-Theoretic Concepts in Computer Science (WG)*, volume 1665 of LNCS, pages 30–41, 1999.
  - [38] D. Peleg. Informative labeling schemes for graphs. In *25<sup>th</sup> International Symposium on Mathematical Foundations of Computer Science (MFCS)*, volume 1893 of LNCS, pages 579–588. Springer, August 2000.
  - [39] P. Powell. A further improved LCA algorithm. Technical Report TR90-01, University of Minneapolis, 1990.
  - [40] N. Santoro and R. Khatib. Labeling and implicit routing in networks. *The Computer Journal*, 28:5–8, 1985.
  - [41] B. Schieber and U. Vishkin. On finding lowest common ancestors: Simplification and parallelization. *SIAM Journal of Computing*, 17:1253–1262, 1988.
  - [42] R. E. Tarjan. Depth-first search and linear graph algorithms. *SIAM Journal on Computing*, 1(2):146–160, June 1972.
  - [43] M. Thorup and U. Zwick. Compact routing schemes. In *13<sup>th</sup> ACM Symposium on Parallel Algorithms and Architectures (SPAA)*, 2001.
  - [44] J. Vuillemin. A unifying look at data structures. *Communications of the ACM*, 23(4):229–239, 1980.
  - [45] J. Westbrook. Fast incremental planarity testing. In Werner Kuich, editor, *19<sup>th</sup> International Colloquium on Automata, Languages and Programming (ICALP)*, volume 623 of LNCS, pages 342–353. Springer-Verlag, 1992.