

Fast Prefix Matching of Bounded Strings

ADAM L. BUCHSBAUM, GLENN S. FOWLER, BALACHANDER KRISHNAMURTHY,
KIEM-PHONG VO, and JIA WANG

AT&T Labs

Longest Prefix Matching (LPM) is the problem of finding which string from a given set is the longest prefix of another, given string. LPM is a core problem in many applications, including IP routing, network data clustering, and telephone network management. These applications typically require very fast matching of bounded strings, i.e., strings that are short and based on small alphabets. We note a simple correspondence between bounded strings and natural numbers that maps prefixes to nested intervals so that computing the longest prefix matching a string is equivalent to finding the shortest interval containing its corresponding integer value. We then present *retries*, a fast and compact data structure for LPM on general alphabets. Performance results show that retries often outperform previously published data structures for IP look-up. By extending LPM to general alphabets, retries admit new applications that could not exploit prior LPM solutions designed for IP look-ups.

Categories and Subject Descriptors: D.2.8 [Software Engineering]: Metrics—*Performance Measures*; E.1 [Data Structures]: Arrays; Records; E.2 [Data Storage Representations]: Linked Representations; F.2.2 [Analysis of Algorithms and Problem Complexity]: Nonnumerical Algorithms and Problems—*Pattern Matching*

General Terms: Algorithms, Experimentation, Measurement, Performance, Theory

Additional Key Words and Phrases: IP routing, prefix matching, table look-up, tries

1. INTRODUCTION

Longest Prefix Matching (LPM) is the problem of determining from a set of strings the longest one that is a prefix of some other, given string. LPM is at the heart of many important applications. Internet Protocol (IP) routers [Fuller et al. 1993] routinely forward packets by computing from their routing tables the longest bit string that forms a prefix of the destination address of each packet. Krishnamurthy and Wang [2000] describe a novel method to cluster Web clients by identifying a set of IP addresses that with high probability are under common administrative control and topologically close together. Clustering information like this has important applications ranging from network design and management to providing on-line quality-of-service differentiation based on the origin of a request. The proposed

An extended abstract appeared in *Proc. 5th Wks. on Algorithm Engineering and Experiments (ALENEX'03)*, Baltimore, MD, 2003, pp. 128–139.

Authors' address: AT&T Labs—Research, Shannon Laboratory, 180 Park Avenue, Florham Park, NJ 07932, USA. Emails: {alb,gsf,bala,kpv,jiawang}@research.att.com.

Permission to make digital/hard copy of all or part of this material without fee for personal or classroom use provided that the copies are not made or distributed for profit or commercial advantage, the ACM copyright/server notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

© 2004 ACM 0000-0000/2004/0000-0001 \$5.00

clustering approach is *network aware* in that addresses are grouped based on prefixes in snapshots of Border Gateway Protocol (BGP) routing tables.

Telephone network management and marketing applications often classify regions in the country by area codes or combinations of area codes and the first few digits of the local phone numbers. For example, the state of New Jersey is identified by area codes such as 201, 908, and 973. In turn, Morris County in New Jersey is identified by longer telephone prefixes like 908876 and 973360. These applications typically require computing in seconds or minutes summaries of calls originating and terminating at certain locations from daily streams of telephone calls, up to hundreds of millions records at a time. This requires very fast classification of telephone numbers by finding the longest matching telephone prefixes.

Similar to other string matching problems [Hirschberg 1977; Jacobson and Vo 1992; Sankoff and Kruskal 1983] with practical applications [Arnold 1977; Chen et al. 2000], LPM solutions must be considered in the context of the intended use to maximize performance. The LPM applications discussed above have some common characteristics:

- Look-ups overwhelmingly dominate updates of the prefix sets. A router may route millions of packets before its routing table changes. Similarly, telephone number classifications rarely change, but hundreds of millions of phone calls are made daily.

- The look-up rate is extremely demanding. IP routing and clustering typically require LPM performance of 200 nanoseconds per look-up or better. This severely limits the number of machine instructions and memory references allowed.

- Prefixes and strings are bounded in length and based on small alphabets. For example, current IP addresses are 32-bit strings, and U.S. telephone numbers are 10-digit strings.

The first two characteristics mean that certain theoretically appealing solutions based on, e.g., suffix trees [McCreight 1976], string prefix matching [Breslauer 1995; Breslauer et al. 1998], or dynamic string searching [Ferragina and Grossi 1995] are not applicable, as their performance would not scale. Fortunately, the third characteristic means that specialized data structures can be designed with the desired performance levels. There are many papers in the literature proposing schemes to solve the IP routing problem with various tradeoffs based on memory consumption or memory hierarchies [Crescenzi et al. 1999; Degermark et al. 1997; Doeringer et al. 1996; Ergun et al. 2001; Feldmann and Muthukrishnan 2000; Lampson et al. 1999; Nilsson and Karlsson 1999; Sklower 1991; Srinivasan and Varghese 1999; Waldvogel et al. 2001]. We are not aware of any published work that generalizes to bounded strings such as telephone numbers, however.

Work on routing Internet packets [Lampson et al. 1999] exploits a simple relationship between IP prefixes and nested intervals of natural numbers. We generalize this idea to a correspondence between bounded strings and natural numbers, which shows that solutions to one instance of LPM may be usable for other instances. We present *retries*, a novel, fast, and compact data structure for LPM on general alphabets; and we perform simulation experiments based on trace data from real applications. On many test sets, *retries* outperform other published data structures

00100000/3	a
00101000/5	b
11000000/2	c
11010000/4	d

Fig. 1. An example prefix set for matching 8-bit strings.

for IP routing, often by significant margins. By extending LPM to general alphabets, retries also admit new applications that could not exploit prior LPM solutions designed for IP look-ups.

2. PREFIXES AND INTERVALS

Let A be an alphabet of finite size $\alpha = \delta + 1$. Without loss of generality, assume that A is the set of natural numbers in the range $[0, \delta]$. Otherwise, map A 's elements to their ranks in any fixed, arbitrary order. Then, we can think of elements of A as digits in base α so that a string $s = s_1 s_2 \cdots s_k$ over A represents an integer $v = s_1 \alpha^{k-1} + s_2 \alpha^{k-2} + \cdots + s_k$. We denote $\iota(s) = v$ and $\sigma(v) = s$. When we work with fixed-length strings, we shall let $\sigma(v)$ have enough 0's padded on the left to gain this length. For example, when the string 1001 represents a number in base 2, $\iota(1001)$ is the decimal value 9. Conversely, in base 3 and with prescribed length 6, $\sigma(9)$ is the string 000100. Clearly, for any two strings s and t with equal lengths, $\iota(s) < \iota(t)$ if and only if s precedes t lexicographically.

2.1 Longest matching prefixes and shortest containing intervals

Let m be some fixed integer. Consider $A^{\leq m}$ and A^m , respectively the sets of strings over A with lengths no greater than m and lengths exactly equal to m . Let $P \subseteq A^{\leq m}$, and with each $p \in P$ let there be associated a data value; the data values need not be mutually distinct. We define an LPM instance (A, m) as the problem of finding the data value associated with the longest string in P that is a prefix of some string $s \in A^m$. P is commonly called the *prefix set*, and its elements are called *prefixes*. Following a convention in IP routing, we shall write s/k to indicate the length- k prefix of a string s ($k \leq \text{len}(s)$).

To show examples of results as they develop, we shall use the binary alphabet $A = \{0, 1\}$ and maximum string length $m = 8$. Figure 1 shows an example prefix set of four strings and associated values. For example, the first string in the set would best match the string 00100101, yielding result a. On the other hand, the second string would best match 00101101, yielding b. It is merely coincidental that all prefixes in this example end with a “1” bit; in practice prefixes can end in any alphabet symbol.

For any string s in $A^{\leq m}$, let $s^0 = s0 \cdots 0$ and $s^\delta = s\delta \cdots \delta$ be two strings in which enough 0's and δ 's are used to pad s to length m . Using the above correspondence between strings and integers, s can be associated with the closed interval of integers $[\iota(s^0), \iota(s^\delta)]$. This interval is denoted $I(s)$, and its *length* is $\iota(s^\delta) - \iota(s^0) + 1$.

Now let v be in $I(s)$, and consider the string $\sigma(v)$ with appropriate padding of 0's on the left to make a string of length m . By construction, $\sigma(v)$ must agree with s^0 and s^δ up to the length of s . On the other hand, if $v < \iota(s^0)$, then $\sigma(v)$ lexicographically precedes s^0 , so s cannot be a prefix of $\sigma(v)$. A similar argument

00100000/3	a	[32,63]
00101000/5	b	[40,47]
11000000/2	c	[192,255]
11010000/4	d	[208,223]

Fig. 2. Prefixes and corresponding intervals.

applies to the case when $v > \iota(s^\delta)$. Thus, we have:

OBSERVATION 2.1. *Let s be a string in $A^{\leq m}$ and $v < \alpha^m$. Then s is a prefix of $\sigma(v)$ if and only if v is in $I(s)$.*

For any prefix set P , we use $I(P)$ to denote the set of intervals associated with prefixes in P . Now consider two prefixes p_1 and p_2 and their corresponding intervals $I(p_1)$ and $I(p_2)$. Applying Observation 2.1 to the endpoints of these intervals shows that either the intervals are completely disjoint or one is contained in the other. Furthermore, $I(p_1)$ contains $I(p_2)$ if and only if p_1 is a prefix of p_2 . Next, when s has length m , $\iota(s^0) = \iota(s^\delta) = \iota(s)$. Observation 2.1 asserts that if p is a prefix of s then $I(p)$ must contain $\iota(s)$. The nested property of intervals in a prefix set P then gives:

THEOREM 2.2. *Let P be a prefix set and s a string in A^m . Then p is the longest prefix matching s if and only if $I(p)$ is the shortest interval in $I(P)$ containing $\iota(s)$.*

Figure 2 shows the correspondence between prefixes and intervals. For example, the string 00101101 with numerical value 45 would have [40,47] as the shortest containing interval, giving **b** as the matching result.

Two intervals that are disjoint or nested are called *nested intervals*. Theorem 2.2 enables treating the LPM problem as that of managing a collection of mutually nested intervals with the following basic operations.

Insert(a, b, v). Insert new interval $[a, b]$ with associated data value v . $[a, b]$ must contain or be contained in any interval it intersects.

Retract(a, b). Delete existing interval $[a, b]$.

Query(p). Determine the value associated with the shortest interval, if any, that contains integer p .

When m and α are small, standard computer integer types suffice to store the integers arising from strings and interval endpoints. As we shall see later, this allows construction of practical data structures for LPM based on integer arithmetic.

2.2 Equivalence among LPM instances and prefix sets

A data structure solving an (A, m) instance can sometimes be used for other instances, as follows. Let (B, n) be another instance of the LPM problem with β the size of B and n the maximal string length. Suppose that $\alpha^m \geq \beta^n$. Since the integers corresponding to strings in $B^{\leq n}$ are less than α^m , they can be represented as strings in $A^{\leq m}$. Furthermore, let $I(p)$ be the interval corresponding to a prefix p in $B^{\leq n}$. Each integer in $I(p)$ can be considered an interval of length 1, so it is representable as a prefix of length m in $A^{\leq m}$. Thus, each string and prefix set in (B, n) can be translated into some other string and prefix set in (A, m) . We have shown:

```

while (lo <= hi)
{
  for (i = 0; i < m; ++i)
    if ((lo % A[i+1]) != 0 ||
        (lo + A[i+1] - 1) > hi)
      break;
  itvl2pfx(lo, lo + A[i] - 1);
  lo += A[i];
}

```

Fig. 3. Constructing the prefixes covering an interval $[lo, hi]$.

00100000/3	a	[32,63]	[32,39]	a	00100000/5
00101000/5	b	[40,47]	[40,47]	b	00101000/5
			[48,63]	a	00110000/4
11000000/2	c	[192,255]	[192,207]	c	11000000/4
11010000/4	d	[208,223]	[208,223]	d	11010000/4
			[224,255]	c	11100000/3

Fig. 4. Left: A prefix set and the corresponding nested intervals. Right: The corresponding disjoint intervals and the equivalent set of disjoint prefixes.

THEOREM 2.3. *Let (A, m) and (B, n) be two instances of the LPM problem in which the sizes of A and B are α and β respectively. Then any data structure solving LPM on (A, m) can be used to solve (B, n) as long as $\alpha^m \geq \beta^n$.*

Using single values in an interval to generate prefixes is inefficient. Let $[lo, hi]$ be some interval where $lo < \alpha^m$ and $hi < \alpha^m$. Figure 3 shows an algorithm (in C) for constructing prefixes from $[lo, hi]$. Simple induction on the quantity $hi - lo + 1$ shows that the algorithm constructs the minimal set of subintervals covering $[lo, hi]$ such that each subinterval corresponds to a single prefix in $A^{\leq m}$. We assume an array $A[]$ such that $A[i]$ has the value α^i . The function `itvl2pfx()` converts an interval into a prefix by inverting the process described earlier of mapping a prefix into an interval. Note that such a prefix will have length $m - i$.

Given a nested set of intervals I , we can construct a minimal set of prefixes $P(I)$ such that (1) $I(p)$ and $I(q)$ are disjoint for $p \neq q$; and (2) finding the shortest interval in I containing some integer i is the same as finding the longest prefix in $P(I)$ matching $\sigma(i)$. This is done as follows.

(1) Sort the intervals in I by their low ends, breaking ties by taking longer intervals first. The nested property of the intervals means that $[i, j] < [k, l]$ in this ordering if $j < k$ or $[i, j]$ contains $[k, l]$.

(2) Build a new set of intervals by adding the sorted intervals in order. When an interval $[i, j]$ is added, if it is contained in some existing interval $[k, l]$, then in addition to adding $[i, j]$, replace $[k, l]$ with at most two new disjoint intervals, $[k, i - 1]$ and $[j + 1, l]$, whenever they are of positive length.

(3) Merge adjacent intervals that have the same data values.

(4) Apply the algorithm in Figure 3 to each of the resulting intervals to construct the new prefixes.

Figure 4 shows how the nested intervals are split into disjoint intervals. These intervals are then transformed into a new collection of prefixes. A property of the new prefixes is that none of them can be a prefix of another.

From now on, we assert that every considered prefix set P shall represent disjoint intervals. If not, we convert it into the equivalent set of prefixes $P(I(P))$ as discussed.

3. THE RETRIE DATA STRUCTURE

Theorem 2.3 asserts that any LPM data structure for one type of string can be used for other LPM instances as long as alphabet sizes and string lengths are within bounds. For example, 15-digit international telephone numbers fit in 64-bit integers, so data structures for fast look-ups of IPv4 32-bit addresses are potentially usable, with proper extensions, for telephone number matching. Unfortunately, many of these are too highly tuned for IP routing to be effective in the other applications that we consider, such as network address clustering and telephone number matching (Section 4). We next describe the *retrieve* data structure for fast LPM queries in $A^{\leq m}$. We compare it to prior art in Section 5.

3.1 The basic retrieve scheme

Let P be a prefix set over $A^{\leq m}$. Each prefix in P is associated with some data value, an integer in a given range $[0, D]$. We could build a table of size α^m that covers the contents of the intervals in P . Then the LPM of a string can be queried with a single index. Such a table would be impossibly large for interesting values of α and m , however. We therefore build a smaller, recursively structured table, which we call a *radix-encoded trie* (or *retrieve*). The basic structure is as follows.

- The top-level table is indexed by some number of left digits of a given string. Each entry in this table points to another table, indexed by some of the following digits, and so on.

- There are two types of tables: *internal* and *leaf*. An entry of an internal table has a pointer to the next-level table and indicates whether that table is an internal or leaf table. An entry of a leaf table contains the data associated with the prefix matched by that entry.

- All internal tables are kept in a single array *Node*, and all leaf tables are kept in a single array *Leaf*.

We discuss in Section 3.2 how to minimize the total space used by and bound the depth of the structure. The size of a leaf entry depends on the maximum data value associated with any given prefix. For example, if the maximum data value is $< 2^8$, then a leaf entry can be stored in a single byte, while a maximum data value between 2^8 and 2^{16} means that leaf data must be stored using 2 bytes each. For fast computation, the size of an internal table entry is chosen so that the entry would fit in some convenient integer type. We partition the bits of this type into three fields: *index*, *type*, and *offset*.

Index specifies the number of digits used to index the next-level table, which thus has α^{index} entries;

```

for (node = root, shift = m; ; sv %= A[shift])
{
  shift -= node >> (obits+1);
  if (node & (1 << obits))
    node = Node[(node & ((1 << obits) - 1)) + sv/A[shift]];
  else return Leaf[(node & ((1 << obits) - 1)) + sv/A[shift]];
}

```

Fig. 5. Searching a retrie for a longest prefix; *obits* is the number of bits reserved for *offset*.

type is a single bit, which if 1 indicates that the next level is internal, and if 0 indicates that the next level is leaf;

offset specifies the offset into the *Node* or *Leaf* array at which the next-level table begins.

Let w be the word size in bits of an internal-entry type. If b_o bits are reserved for *offset*, then the maximum size for the *Node* and *Leaf* arrays is $\alpha^{b_o/\lg \alpha}$, the largest power of α no greater than 2^{b_o} , which thus upper bounds the size of any table. Now if b_i bits are reserved for *index*, then the maximum table size is also $\alpha^{2^{b_i}-1}$. Therefore, b_i and b_o should be chosen so that these two values are about equal, i.e., so that b_i is about $\lg(1 + b_o/\lg \alpha)$, keeping in mind that $b_i + b_o + 1 = w$. In practice, we often use 32-bit integers for internal table entries. For $\alpha = 2$, we thus choose $b_o = 26$ and $b_i = 5$, since $\lg(1 + 26/\lg 2)$ is about 4.8. These choices allow the sizes of the *Node* and *Leaf* arrays to be up to 2^{26} , which is ample for our applications.

Given a retrie built for some prefix set $P \subseteq A^{\leq m}$, let *root* be a single internal table entry that describes how to index the top-level table. Let $A[i]$ be an integer array such that $A[i] = \alpha^i$. Let s be a string in A^m with integer value $sv = \iota(s)$. Figure 5 shows the algorithm to compute the datum associated with the LPM of s .

Figure 6 shows a 3-level retrie for the prefixes and associated data values shown in Figure 4. The *Node* array begins with the top-level internal table. Indices to the left of table entries are in binary and with respect to the head of the corresponding table within the *Node* or *Leaf* array. Each internal-table entry has three fields as discussed. All internal-table entries with *offset* = *nil* indicate some default value for strings without any matching prefix. For example, the string 00101101 is matched by first stripping off the starting 2 bits, 00, to index entry 0 of the top-level table. The *type* of this entry is 1, indicating that the next level is another internal table. The *offset* of the entry points to the base of this table. The *index* of the entry indicates that one bit should be used to index the next level. Then the indexed entry in the next-level table points to a leaf table. The entries of this table are properly defined so that the fourth and fifth bits of the string, 01, index the entry with the correct data: **b**.

A retrie with k levels enables matching with at most k indexing operations. This guarantee is important in applications such as IP forwarding. Smaller k 's mean larger look-up tables, so it is important to ensure that a retrie with k levels uses minimal space. We next discuss how to do this using dynamic programming.

3.2 Optimizing the basic retrie scheme

Given a prefix set P , define $len(P) = \max\{len(p) : p \in P\}$. For $1 \leq i \leq len(P)$, let $P^{\leq i}$ be the subset of prefixes with lengths $\leq i$. Then let $L(P, i)$ be the partition

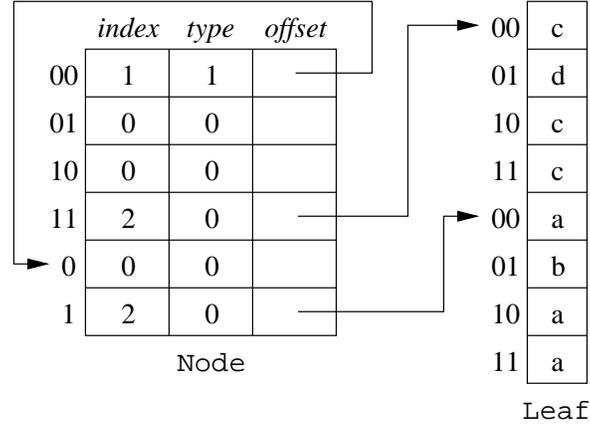


Fig. 6. A retrieve data structure.

of $P - P^{\leq i}$ into equivalence classes induced by the left i digits. That is, each part Q in $L(P, i)$ consists of all prefixes longer than i and having the same first i digits. Now, let $strip(Q, i)$ be the prefixes in Q with their left i digits stripped off. Such prefixes represent disjoint intervals in the LPM instance $(A, m - i)$. Finally, let c_d be the size of a data (leaf table) entry and c_t the size of an internal-table entry.

Following is the dynamic program to compute $S(P, k)$, the optimal size of a retrieve data structure for a prefix set P using at most k levels; the value of k depends on the application.

$$S(P, k) = \begin{cases} c_d \alpha^{\text{len}(P)} & \text{if } P = \emptyset \text{ or } k = 1, \\ \min \{c_d \alpha^{\text{len}(P)}, D(P, k)\} & \text{otherwise;} \end{cases}$$

$$\text{where } D(P, k) = \min_{1 \leq i < \text{len}(P)} \left\{ c_t \alpha^i + c_d |P^{\leq i}| + \sum_{Q \in L(P, i)} S(strip(Q, i), k - 1) \right\}.$$

The first case states that a leaf table is built whenever the prefix set is empty or only a single level is allowed. The second case recurses to compute the optimal retrieve and its size. The first part of the minimization considers the case when a single leaf table is constructed. The second part of the minimization—the computation of $D(P, k)$ —considers the cases when internal tables may be constructed. In these cases, the first term, $c_t \alpha^i$, expresses the size of the constructed internal table to be indexed by the first i digits of a string. The second term, $c_d |P^{\leq i}|$, expresses the fact that each prefix short enough to end at the given level requires a single leaf-table entry for its data value. The last term, $\sum_{Q \in L(P, i)} S(strip(Q, i), k - 1)$, recurses down each set $strip(Q, i)$ to compute a retrieve with optimal size for it.

Each set $strip(Q, i)$ is uniquely identified by the string formed from the digits stripped off from the top level of the recursion until $strip(Q, i)$ is formed. The number of such strings is bounded by $|P| \text{len}(P)$. Each such set contributes at most one term to any partition of the $\text{len}(P)$ bits in a prefix. For $k \leq \text{len}(P)$, the dynamic program examines all partitions of $[1, \text{len}(P)]$ with $\leq k$ parts. The number of such partitions is $O(\text{len}(P)^{k-1})$. Thus, we have:

THEOREM 3.1. $S(P, k)$ can be computed in $O(|P|len(P)^k)$ time.

In practice, $len(P)$ is bounded by a small constant, e.g., 32 for IP routing and 10 for U.S. phone numbers. Since the number of levels cannot exceed $len(P)$, the dynamic program essentially runs in time linear in the number of prefixes.

3.3 Superstring lay-out of leaf tables

The internal and leaf tables are sequences of elements. In the dynamic program, we consider each table to be instantiated in the *Node* or *Leaf* array as it is created. We can reduce memory usage by exploiting redundancies, especially in the leaf tables. For example, in IP routing, the leaf tables often contain many similar, long runs of relatively few distinct data values. Computing a short superstring of the tables reduces space very effectively. Since computing shortest common superstrings (SCS) is MAX-SNP hard [Blum et al. 1994], we experiment with three heuristics.

- (1) The trivial superstring is formed by concatenating the leaf tables.
- (2) The left and right ends of the leaf tables are merged in a best-fit order.
- (3) A superstring is computed using a standard greedy SCS approximation [Blum et al. 1994].

Both methods 2 and 3 effectively reduce space usage. (See Section 4.) In practice, however, method 2 gives the best trade-off between computation time and space.

Finally, we note that it is possible to incorporate superstring computation of leaf tables into the dynamic program to estimate more accurately the actual sizes of the leaf tables. This would better guide the dynamic program to select an optimal overall lay-out. The high cost of superstring computation makes this impractical, however. Thus, the superstring lay-out of leaf tables is done only after the dynamic program finishes.

4. APPLICATIONS AND PERFORMANCE

We consider three applications: IP routing, network clustering, and telephone service marketing. Each of these exhibits a different set of characteristics. In current IP routing, the strings are 32 bits long, and the number of distinct data values, i.e., next-hops, is small. For network clustering, we merge several BGP tables together and use either the prefixes or their lengths as data values so that after a look-up we can retrieve the matching prefix itself. In this case, either the number of data values is large or there are many runs of data values. For routing and clustering, we compared retries to data structures with publicly available code: LC-tries [Nilsson and Karlsson 1999], which are conceptually quite similar, and the compressed-table data structure of Crescenzi et al. [1999] (CDG), which is among the fastest IP look-up data structures reported in the literature. (See Section 5.) We used the authors' code for both benchmarks and included in our test suite the FUNET router table and traffic trace used in the LC-trie work [Nilsson and Karlsson 1999]. (See comment in the discussion regarding the FUNET data.) These data structures are designed specifically for IP prefix matching. The third application is telephone service marketing, in which strings are telephone numbers.

CDG and retries provide constant-time look-ups wherever they are applicable: CDG guarantees answers in three memory accesses, and we used depth-2 and -3

Table I. Number of entries, next-hops, and data structure sizes for tables used in IP routing experiment.

Routing table	Entries	Next-hops	Data struct. size (KB)				
			retrie			lctrie	CDG
			-FL	-LR	-GR		
AADS	32505	38	1069.49	866.68	835.61	763.52	4446.37
ATT	71483	45	2508.79	2231.89	2180.21	1659.52	15601.92
FUNET	41328	18	506.57	433.00	411.36	967.36	682.93
MAE-WEST	71319	38	1241.14	1040.37	1000.52	1654.06	5520.26
OREGON	118190	33	3828.93	3107.78	3035.73	2711.16	12955.85
PAIX	17766	28	912.94	741.92	723.85	417.74	3241.31
TELSTRA	104096	182	2355.03	2023.08	1971.78	2384.66	9863.96

retries in our experiments. In contrast, LC-tries can incur look-up time linear in the prefix length, although this would be atypical in practice.

4.1 IP routing

Table I and Figure 7 summarize the routing tables we used. They report how many distinct prefixes and next-hops each contained and the sizes of the data structures built on each. Retrieve-FL (rsp., -LR, -GR) is a depth-2 retrieve with catenated (rsp., left-right/best-fit merge, greedy) record layout. For routing, we limited depth to 2 to emphasize query time. We use deeper retrievals in Section 4.2. ATT is a routing table from an AT&T BGP router; FUNET is the routing table reported in the LC-trie work [Nilsson and Karlsson 1999]; TELSTRA comes from Telstra Internet [2001]; the other tables are described by Krishnamurthy and Wang [2000].

We timed the LPM queries engendered by router traffic. Since we lacked real traffic traces for the tables other than ATT and FUNET, we constructed random traces by choosing, for each table, 100,000 prefixes uniformly at random (with replacement) and extending them to full 32-bit addresses. We used each random trace in the order generated and also with the addresses sorted lexicographically to present locality that might be expected in a real router. We generated random traces for the ATT and FUNET tables as well, to juxtapose real and random data. We processed each trace through the respective table 100 times, measuring average LPM query time; each data structure was built from scratch for each unique prefix table. We also recorded data structure build times.

We performed this evaluation on two machines: an SGI Challenge (400 MHz R12000) with split 32 KB L1 data and instruction caches, 8 MB unified L2 cache, and 12 GB main memory, running IRIX 6.5; and a 1 GHz Pentium III with split 16 KB L1 data and instruction caches, 256 KB unified L2 cache, and 256 MB main memory, running Linux 2.4.6. Each time reported is the median of five runs. Table II reports the results, and Figures 8–10 plot the query times.

LC-tries were designed to fit in L2 cache and do so on all the tables on the SGI but none on the Pentium. Retries behave similarly, although they were not designed to be cache resident. CDG fits in the SGI cache on AADS, FUNET, MAE-WEST, and PAIX. Retries uniformly outperformed LC-tries, sometimes by an order of magnitude, always by a factor exceeding three. CDG significantly outperformed

Table II. Build and query times for routing. “T” designates times over real router traces; “SR” designates times over sorted random traces; “R” designates times over random traces.

Routing table	Data struct.	SGI				Pentium			
		Build (ms)	Query (ns)			Build (ms)	Query (ns)		
			T	SR	R		T	SR	R
AADS	retrie-FL	195		15	20	150		25	48
	retrie-LR	225		15	20	180		24	42
	retrie-GR	4157		15	20	3010		22	40
	lctrie	113		160	215	100		153	380
	CDG	214		17	144	220		25	69
ATT	retrie-FL	459	18	16	31	370	20	32	83
	retrie-LR	518	17	16	22	440	18	31	68
	retrie-GR	9506	17	16	22	7220	19	31	66
	lctrie	245	163	159	214	270	146	181	452
	CDG	1011	64	39	224	1010	42	43	85
FUNET	retrie-FL	166	14	14	18	120	14	20	28
	retrie-LR	190	14	14	17	130	14	21	24
	retrie-GR	1650	14	14	17	1020	14	19	24
	lctrie	136	134	149	199	140	111	153	381
	CDG	102	15	14	20	70	8	14	26
MAE-WEST	retrie-FL	340		14	21	270		28	65
	retrie-LR	388		15	20	310		26	52
	retrie-GR	8241		15	21	6950		26	49
	lctrie	233		155	210	250		176	454
	CDG	325		19	152	310		33	73
OREGON	retrie-FL	447		15	67	350		34	76
	retrie-LR	512		16	23	390		31	59
	retrie-GR	8252		16	23	7740		31	56
	lctrie	383		155	341	420		206	464
	CDG	949		26	142	1050		41	81
PAIX	retrie-FL	118		15	20	90		24	36
	retrie-LR	132		14	18	90		22	30
	retrie-GR	2100		14	18	1360		22	28
	lctrie	66		162	213	60		146	284
	CDG	136		15	123	150		22	58
TELSTRA	retrie-FL	508		16	27	420		37	86
	retrie-LR	572		16	21	480		31	62
	retrie-GR	10000		15	21	8510		31	56
	lctrie	343		155	311	390		201	458
	CDG	599		26	180	610		47	83

retries on the real and sorted random traces for FUNET on the Pentium, but this advantage disappeared for the random trace and also for all the FUNET traces on the SGI. This suggests the influence of caching effects. Also, the numbers of prefixes and next-hops for FUNET were relatively low, and CDG is sensitive to these sizes. On the larger tables (ATT, MAE-WEST, OREGON and TELSTRA), retries significantly outperformed CDG, even for the real and sorted random traces for ATT (on both machines). As routing tables are continually growing, with 250,000 entries expected by the year 2005, we expect that retries will outperform CDG on real data. Finally, the FUNET trace was filtered to zero out the low-order 8 bits of each address for privacy purposes [Nilsson 2001] and is likely not a true

trace for the prefix table, which contains some prefixes longer than 24 bits.

The data suggest that the non-trivial superstring retrie variations significantly reduce retrie space. As might be expected, the greedy superstring approximation is comparatively slow, but the best-fit merge runs with little time degradation over the trivial superstring method and still provides significant space savings. The FUNET results, in particular on the real and sorted random traces, suggest that CDG benefits from the ordering of memory accesses more than retries benefit from the superstring layouts. A first-fit superstring merging strategy might be useful in testing this hypothesis.

There is a nearly uniform improvement in look-up times from retrie-FL to retrie-LR to retrie-GR even though each address look-up performs exactly the same computation and memory accesses in all three cases. This suggests beneficial effects from the hardware caches. We believe that this is due to the overlapping of leaf tables in retrie-LR and retrie-GR, which both minimizes space usage and increases the hit rates for similar next-hop values.

The data also suggest that LPM queries on real traces run significantly faster than on random traces. Again this suggests beneficial cache effects, from the locality observed in our IP traffic traces. Real trace data is thus critical for accurate measurements, although random data seem to provide an upper bound to real-world performance.

Finally, while retries take longer to build than LC-tries (and sometimes CDG), build time (for -FL and -LR) is acceptable, and query time is more critical to routing and on-line clustering, which we assess next.

4.2 Network clustering

For clustering, we combined the routing tables used above. There were 168,161 unique prefixes in the tables. The goal of clustering is to recover the actual matching prefix for an IP address, thereby partitioning the addresses into equivalence classes [Krishnamurthy and Wang 2000]. PREF assigns each resulting prefix itself as the data value. LEN assigns the length of each prefix as its data value, which is sufficient to recover the prefix given the input address. PREF, therefore, has 168,161 distinct data values, whereas LEN has only 32.

We built depth-2 and -3 retries and LC-tries for PREF and LEN. Table III and Figure 11 detail the data structure sizes. Note the difference in retrie sizes for the two tables. The relative sparsity of data values in LEN produces a much smaller *Leaf* array, which can also be more effectively compressed by the superstring methods. Also note the space reduction achieved by depth-3 retries compared to depth-2 retries. Depth-3 retries are smaller than LC-tries for this application, yet, as we will see, outperform the latter. CDG could not be built on either PREF or LEN. CDG assumes a small number of next-hops and exceeded memory limits for PREF on both machines. CDG failed on LEN, because the number of runs of equal next-hops was too large. Here the difference between the IP routing and clustering applications of LPM becomes striking: retries work for both applications, but CDG cannot be applied to clustering.

We timed the clustering of Apache and EW3 server logs. Apache contains IP addresses recorded at www.apache.org. The 3,461,361 records gathered in late 1997 had 274,844 unique IP addresses. EW3 contains IP addresses of clients visiting

Table III. Data structure sizes for tables used in clustering experiment.

Routing table	Data struct. size (KB)						
	depth-2 retriee			depth-3 retriee			lctrie
	-FL	-LR	-GR	-FL	-LR	-GR	
PREF	13554.87	13068.88	13054.80	3181.63	2878.24	2889.90	3795.91
LEN	5704.37	4938.74	4785.66	1400.84	1045.53	990.33	3795.91

Table IV. Build and query times for clustering.

Machine	Table	Operation	Times (build: ms) (query: ns)						lctrie
			depth-2 retriee			depth-3 retriee			
			-FL	-LR	-GR	-FL	-LR	-GR	
SGI	PREF	build	1732	2028	25000	6919	7478	43000	801
		query (Apache)	20	19	19	36	35	35	136
		query (EW3)	20	19	19	36	36	36	139
	LEN	build	1299	1495	28000	4299	4476	25000	588
		query (Apache)	15	16	16	32	32	32	135
		query (EW3)	15	16	16	33	32	32	139
Pentium	PREF	build	1300	1570	19000	4670	5060	32000	750
		query (Apache)	26	26	26	41	40	40	121
		query (EW3)	27	27	27	43	42	42	129
	LEN	build	990	1170	25000	2860	3070	18000	640
		query (Apache)	21	21	21	35	34	34	117
		query (EW3)	23	21	21	37	35	35	124

three large corporations whose content were hosted by AT&T's Easy World Wide Web. Twenty-three million entries were gathered during February and March, 2000, representing 229,240 unique IP addresses. The experimental setup was as in the routing assessment. In an on-line clustering application, such as in a Web server, the data structures are built once (e.g., daily), and addresses are clustered as they arrive. Thus, query time is paramount. Retries significantly outperform LC-tries for this application, even at depth 3, as shown in Table IV and Figure 12.

4.3 Telephone service marketing

In our telephone service marketing application, the market is divided into regions, each of which is identified by a set of telephone prefixes. Given daily traces of telephone calls, the application classifies the callers by regions, updates usage statistics, and generates a report. Such reports may help in making decisions on altering the level of advertisement in certain regions. For example, the set of prefixes identifying Morris County, NJ, includes 908876 and 973360. Thus, a call originating from a telephone number of the form 973360XXXX would match Morris County, NJ.

Table V shows performance results (on the SGI) from using a retriee to summarize telephone usage in different regions of the country for the first half of 2001. The second column shows the number of call records per month used in the experiment. Since this application is off-line, we consider the total time required to classify all the numbers. The third column shows this time (in seconds) for retries, and the fourth column contrasts the time using a binary search approach for matching. This shows the benefit of using retries in this application. Previous IP look-up

Table V. Time to classify telephone numbers.

Month	Counts	Retrie (s)	Bsearch (s)
1	27,479,712	24.35	83.51
2	25,510,814	22.37	74.73
3	28,993,583	25.49	84.60
4	28,452,823	24.94	80.76
5	29,786,302	26.11	84.86
6	28,874,669	25.27	80.79

data structures, which we review next, do not currently extend to this alphabet, although they could be extended using the correspondence between bounded strings and natural numbers.

5. COMPARISON TO RELATED WORK

The popularity of the Internet has made IP routing an important area of research. Several LPM schemes for binary strings were invented in this context. To the best of our knowledge, Amir et al. [1995] were the first to exploit the equivalence of longest common prefixes and nested intervals, via an intermediate reduction to the nearest marked ancestor problem on trees. Lampson et al. [1999] applied this concept to perform IP look-ups; Gupta et al. [2000] later provided a worst-case analysis. Ergun et al. [2001] considered biased access to ranges. Feldmann and Muthukrishnan [2000] generalized the idea to packet classification. We generalized this idea to non-binary strings and showed that LPM techniques developed for strings based on one alphabet can also be used for strings based on another. Thus, under the right conditions, the data structures invented for IP routing can be used for general LPM. Encoding strings over arbitrary alphabets as reals and searching in that representation goes back at least to arithmetic coding; see, e.g., Cover and Thomas [1991].

Retries are in the general class of multi-level table look-up schemes used for both hardware [Gupta et al. 1998; Huang et al. 1999; Moestedt and Sjödin 1998] and software [Degermark et al. 1997; Doeringer et al. 1996; Nilsson and Karlsson 1999; Sklower 1991; Srinivasan and Varghese 1999] implementations for IP routing. Waldvogel et al. [2001] discuss a different type of software-based scheme, which uses binary search over multiple hash tables. Since modern machines use memory hierarchies with sometimes dramatically different performance levels, some of these works attempt to build data structures conforming to the memory hierarchies at hand. Both the LC-trie scheme of Nilsson and Karlsson [1999] and the multi-level table of Srinivasan and Varghese [1999] attempt to optimize for L2 caches by adjusting the number of levels to minimize space usage. Efficient implementations, however, exploit the binary alphabet of IP addresses and prefixes. Further distinguishing retries from LC-tries are that retries have a single-phase search procedure, whereas LC-tries require backtracking once a the appropriate leaf is discovered; and that space-efficient retries can be constructed to guarantee constant look-up time.

Cheung and McCanne [1999] took a more general approach to dealing with memory hierarchies that includes the use of prefix popularity. They consider a multi-level table scheme similar to retries and attempt to minimize the space usage of popu-

lar tables so that they fit into the fast caches. Since the cache sizes are limited, they must solve a complex constrained optimization problem to find the right data structure. L1 caches on most machines are very small, however, so much of the gain comes from fitting a data structure into L2 caches. In addition, the popularity of prefixes is a dynamic property and not easy to approximate statically. Thus, we focus on bounding the number of memory accesses and minimizing memory usage. We do this by (1) separating internal tables from leaf tables so that the latter can use small integer types to store data values; (2) using dynamic programming to optimize the lay-out of internal and leaf tables given some bound on the number of levels, which also bounds the number of memory accesses during queries; and (3) using a superstring approach to minimize space usage of the leaf tables. The results in Section 4 show that we achieve both good look-up performance and small memory usage.

Crescenzi et al. [1999] introduced a compressed-table data structure for IP look-up. The key idea is to identify runs induced by common next-hops among the 2^{32} implicit prefixes to compress the entire table with this information. This works well when the number of distinct next-hops is small and there are few runs, which is mostly the case in IP routing. The compressed-table data structure is fast, because it bounds the number of memory accesses per match. Unfortunately, in network clustering applications, both the number of distinct next-hop values and the number of runs can be quite large. Thus, this technique cannot be used in such applications. Section 4 shows that retries often outperform the compressed-table data structure for IP routing and use much less space.

Finally, there is related work on the more general problem of packet filtering [Eppstein and Muthukrishnan 2001; Srinivasan et al. 1998; Suri and Varghese 1999; Warkhede et al. 2001], which typically subsumes IP routing but also includes higher dimensional queries. As we compared algorithms specifically designed for IP routing, we did not consider algorithms for the more general problem.

6. CONCLUSIONS

We considered the problem of performing LPM on short strings with bounded alphabets. We showed how to map such strings into the integers so that small strings would map to values representable in machine words. This enabled the use of standard integer arithmetic for prefix matching. We then presented retries, a novel, multi-level table data structure for LPM. Experimental results were presented showing that retries outperform other comparable data structures.

A number of open problems remain. A dynamic LPM data structure that performs queries empirically fast remains elusive. Build times for static structures are acceptable for present applications, but the continual growth of routing tables will likely necessitate dynamic solutions in the future. As with general string matching solutions, theoretically appealing approaches, e.g., based on interval trees [Preparata and Shamos 1988], do not exploit some of the peculiar characteristics of these applications. Feldmann and Muthukrishnan [2000] report partial progress. We have a prototype based on nested-interval maintenance but have not yet assessed its performance.

Our results demonstrate that LPM data structures perform much better on real

trace data than on randomly generated data. Investigating the cache behavior of LPM data structures on real data thus seems important. Compiling benchmark suites of real data is problematic given the proprietary nature of such data, so work on modeling IP address traces for experimental purposes is also worthwhile. Such data would also allow other types of experimental studies of LPM data structures, in particular empirical validations of theoretical analyses, which are not currently possible without models upon which to base predictions.

ACKNOWLEDGMENTS

We thank John Linderman for describing the use of prefixes in telephone service marketing. We also thank the anonymous referees who made many helpful comments on both this manuscript and our original ALENEX extended abstract.

REFERENCES

- AMIR, A., FARACH, M., IDURY, R. M., POUTRÉ, J. A. L., AND SCHÄFFER, A. A. 1995. Improved dynamic dictionary matching. *Inf. Comp.* 119, 2, 258–282.
- ARNOLD, K. C. R. C. 1977. *Screen Updating and Cursor Movement Optimization: A Library Package*. 4.2BSD UNIX Programmer’s Manual.
- BLUM, A., LI, M., TROMP, J., AND YANNAKAKIS, M. 1994. Linear approximation of shortest superstrings. *J. ACM* 41, 4, 630–647.
- BRESLAUER, D. 1995. Fast parallel string prefix-matching. *Theor. Comput. Sci.* 137, 2, 268–278.
- BRESLAUER, D., COLUSSI, L., AND TONIOLO, L. 1998. On the comparison complexity of the string prefix-matching problem. *J. Alg.* 29, 1, 18–67.
- CHEN, Y.-F., DOUGLIS, F., HUANG, H., AND VO, K.-P. 2000. TopBlend: An efficient implementation of HtmlDiff in Java. In *Proc. 2000 World Conf. on the WWW and Internet*. AACE, Norfolk, VA.
- CHEUNG, G. AND MCCANNE, S. 1999. Optimal routing table design for IP address lookup under memory constraints. In *Proc. 18th Ann. Joint Conf. of the IEEE Computer and Communications Societies*. Vol. 3. IEEE, Los Alamitos, CA, 1437–1444.
- COVER, T. M. AND THOMAS, J. A. 1991. *Elements of Information Theory*. John Wiley & Sons, New York.
- CRESCENZI, P., DARDINI, L., AND GROSSI, R. 1999. IP address lookup made fast and simple. In *Proc. 7th European Symp. on Algorithms*. Lecture Notes in Computer Science, vol. 1643. Springer-Verlag, Berlin, 65–76.
- DEGERMARK, M., BRODNIK, A., CARLSSON, S., AND PINK, S. 1997. Small forwarding tables for fast routing lookups. In *Proc. ACM SIGCOMM ’97*. ACM, New York, 3–14.
- DOERINGER, W., KARJOTH, G., AND NASSEHI, M. 1996. Routing on longest-matching prefixes. *IEEE/ACM Trans. Network.* 4, 1, 86–97. *Erratum*, 1997, 5, 1, 600.
- EPPSTEIN, D. AND MUTHUKRISHNAN, S. 2001. Internet packet filter management and rectangle geometry. In *Proc. 12th ACM-SIAM Symp. on Discrete Algorithms*. SIAM, Philadelphia, 827–835.
- ERGUN, F., MITTRA, S., SAHINALP, S. C., SHARP, J., AND SINHA, R. K. 2001. A dynamic lookup scheme for bursty access patterns. In *Proc. 20th Ann. Joint Conf. of the IEEE Computer and Communications Societies*. Vol. 3. IEEE, Los Alamitos, CA, 1444–1453.
- FELDMANN, A. AND MUTHUKRISHNAN, S. 2000. Tradeoffs for packet classification. In *Proc. 19th Ann. Joint Conf. of the IEEE Computer and Communications Societies*. Vol. 3. IEEE, Los Alamitos, CA, 1193–1202.
- FERRAGINA, P. AND GROSSI, R. 1995. A fully-dynamic data structure for external substring search. In *Proc. 27th ACM Symp. on Theory of Computing*. ACM, New York, 693–702.

- FULLER, V., LI, T., YU, J., AND VARADHAN, K. 1993. *Classless Inter-Domain Routing (CIDR): An Address Assignment and Aggregation Strategy*. Internet Engineering Task Force (www.ietf.org). RFC 1519.
- GUPTA, P., LIN, S., AND MCKEOWN, M. 1998. Routing lookups in hardware and memory access speeds. In *Proc. 17th Ann. Joint Conf. of the IEEE Computer and Communications Societies*. Vol. 3. IEEE, Los Alamitos, CA, 1240–1247.
- GUPTA, P., PRABHAKAR, B., AND BOYD, S. 2000. Near-optimal routing lookups with bounded worst case performance. In *Proc. 19th Ann. Joint Conf. of the IEEE Computer and Communications Societies*. Vol. 3. IEEE, Los Alamitos, CA, 1184–1192.
- HIRSCHBERG, D. S. 1977. Algorithms for the longest common subsequence problem. *J. ACM* 24, 4, 664–675.
- HUANG, N.-F., ZHAO, S.-M., PAN, J.-Y., AND SU, C.-A. 1999. A fast IP routing lookup scheme for gigabit switching routers. In *Proc. 18th Ann. Joint Conf. of the IEEE Computer and Communications Societies*. Vol. 3. IEEE, Los Alamitos, CA, 1429–1436.
- JACOBSON, G. AND VO, K.-P. 1992. Heaviest increasing/common subsequence problems. In *Proc. 3rd Symp. on Combinatorial Pattern Matching*. Lecture Notes in Computer Science, vol. 644. Springer-Verlag, Berlin, 52–65.
- KRISHNAMURTHY, B. AND WANG, J. 2000. On network-aware clustering of Web clients. In *Proc. ACM SIGCOMM '00*. ACM, New York, 97–110.
- LAMPSON, B., SRINIVASAN, V., AND VARGHESE, G. 1999. IP lookups using multiway and multi-column search. *IEEE/ACM Trans. Network.* 7, 3, 324–334.
- MCCREIGHT, E. M. 1976. A space-economical suffix tree construction algorithm. *J. ACM* 23, 2, 262–272.
- MOESTEDT, A. AND SJÖDIN, P. 1998. IP address lookup in hardware for high-speed routing. In *Proc. Hot Interconnects VI*. IEEE, Los Alamitos, CA.
- NILSSON, S. 2001. Personal communication.
- NILSSON, S. AND KARLSSON, G. 1999. IP-address lookup using LC-tries. *IEEE J. Sel. Areas Comm.* 17, 6, 1083–1092.
- PREPARATA, F. P. AND SHAMOS, M. I. 1988. *Computational Geometry: An Introduction*. Texts and Monographs in Computer Science. Springer-Verlag, Berlin.
- SANKOFF, D. AND KRUSKAL, J. B. 1983. *Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparisons*. Addison Wesley, Reading, MA.
- SKLOWER, K. 1991. A tree-based routing table for Berkeley UNIX. In *Proc. USENIX Winter 1991 Technical Conf.* USENIX, Berkeley, CA, 93–104.
- SRINIVASAN, V. AND VARGHESE, G. 1999. Fast address lookup using controlled prefix expansion. *ACM Trans. Comput. Syst.* 17, 1, 1–40.
- SRINIVASAN, V., VARGHESE, G., SURI, S., AND WALDVOGEL, M. 1998. Fast and scalable layer four switching. In *Proc. ACM SIGCOMM '98*. ACM, New York, 191–202.
- SURI, S. AND VARGHESE, G. 1999. Packet filtering in high speed networks. In *Proc. 10th ACM-SIAM Symp. on Discrete Algorithms*. SIAM, Philadelphia, 969–970.
- Telstra Internet 2001. <http://www.telstra.net/ops/bgptab.txt>.
- WALDVOGEL, M., VARGHESE, G., TURNER, J., AND PLATTNER, B. 2001. Scalable high speed prefix matching. *ACM Trans. Comput. Syst.* 19, 4, 440–482.
- WARKHEDE, P., SURI, S., AND VARGHESE, G. 2001. Fast packet classification for two-dimensional conflict-free filters. In *Proc. 20th Ann. Joint Conf. of the IEEE Computer and Communications Societies*. Vol. 3. IEEE, Los Alamitos, CA, 1434–1443.

Received Month Year; revised Month Year; accepted Month Year

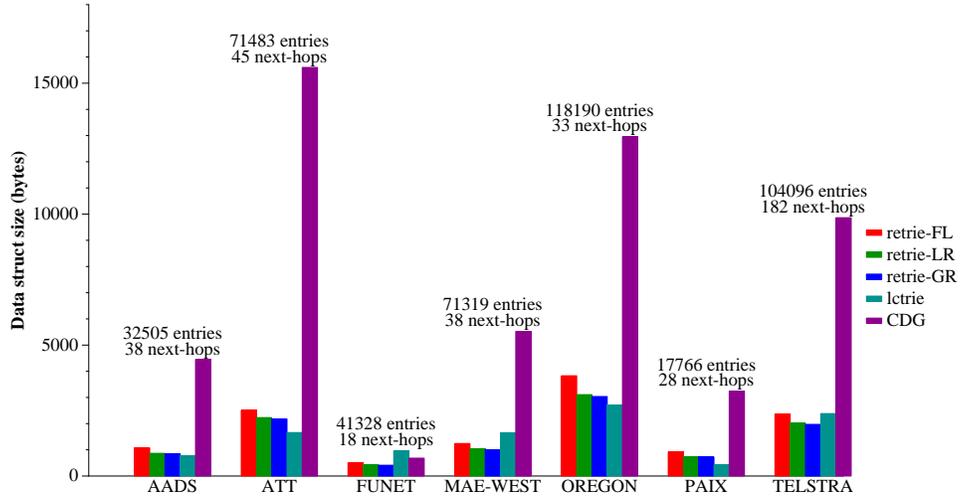


Fig. 7. Data structure sizes for tables used in IP routing experiment.

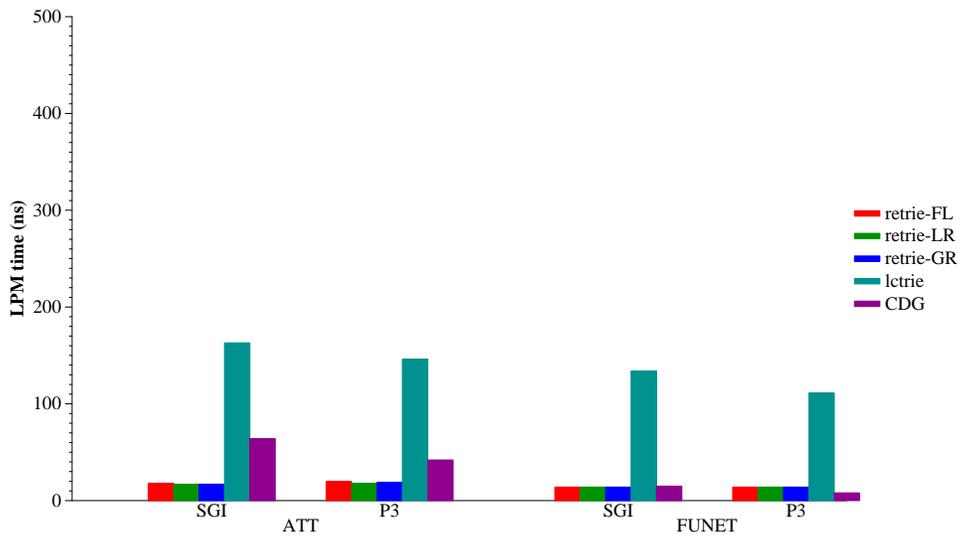


Fig. 8. Query times for routing using real traffic data.

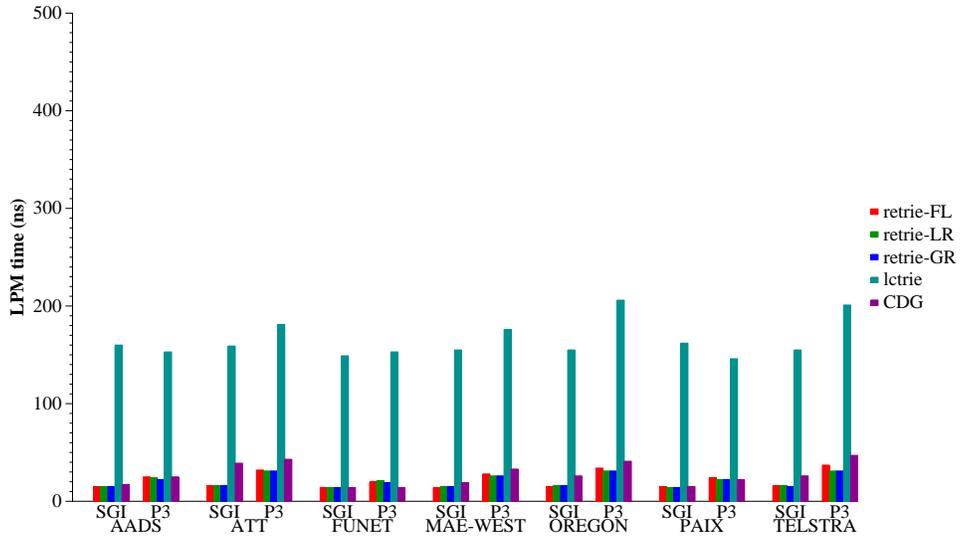


Fig. 9. Query times for routing using sorted random traffic data.

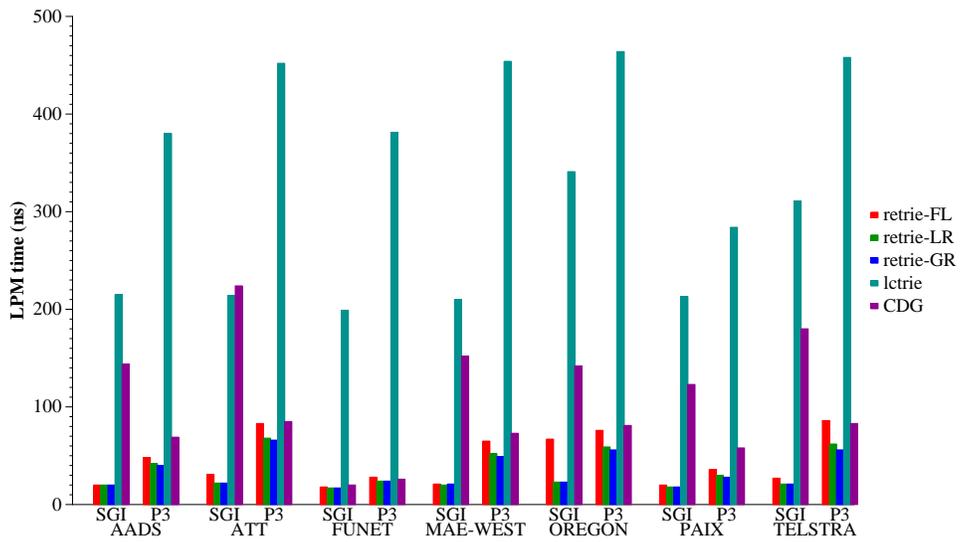


Fig. 10. Query times for routing using random traffic data.

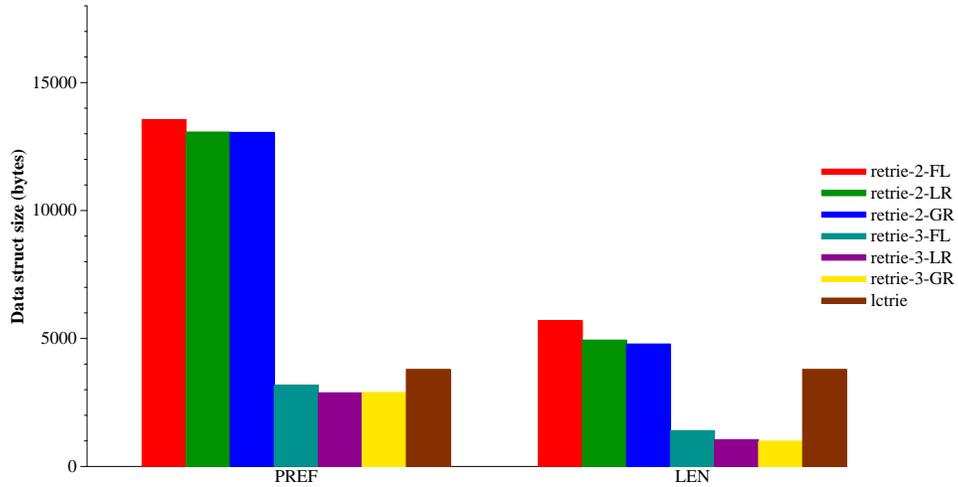


Fig. 11. Data structure sizes for tables used in clustering experiment.

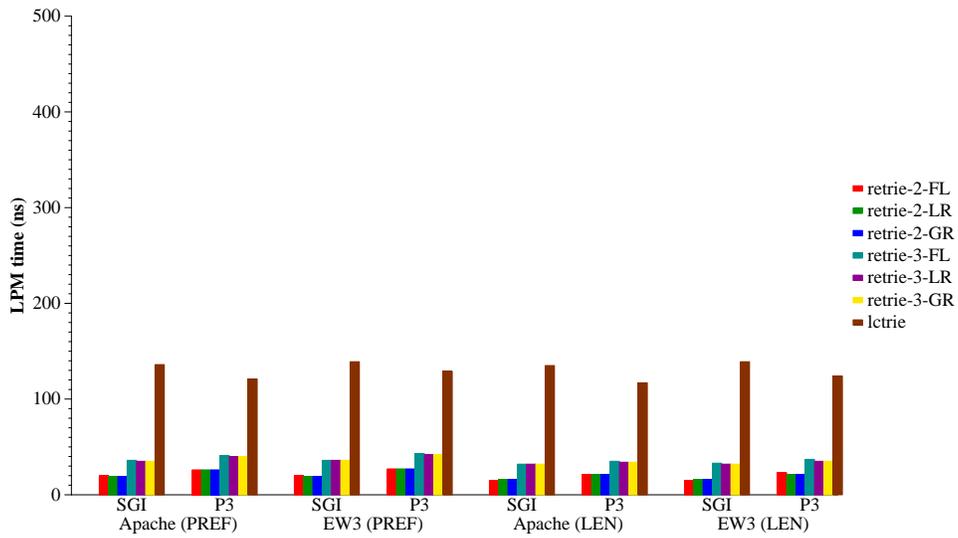


Fig. 12. Query times for clustering.