

RE-tree: an efficient index structure for regular expressions

Chee-Yong Chan, Minos Garofalakis, Rajeev Rastogi

Bell Laboratories, Lucent Technologies, 600 Mountain Ave., Murray Hill, NJ 07974, USA;
e-mail: {cychan, minos, rastogi}@research.bell-labs.com

Edited by R. Ramakrishnan. Received: September 16, 2002
Published online: July 8, 2003 – © Springer-Verlag 2003

Abstract. Due to their expressive power, regular expressions (REs) are quickly becoming an integral part of language specifications for several important application scenarios. Many of these applications have to manage huge databases of RE specifications and need to provide an effective matching mechanism that, given an input string, quickly identifies the REs in the database that match it. In this paper, we propose the RE-tree, a novel index structure for large databases of RE specifications. Given an input query string, the RE-tree speeds up the retrieval of matching REs by focusing the search and comparing the input string with only a small fraction of REs in the database. Even though the RE-tree is similar in spirit to other tree-based structures that have been proposed for indexing multidimensional data, RE indexing is significantly more challenging since REs typically represent infinite sets of strings with no well-defined notion of spatial locality. To address these new challenges, our RE-tree index structure relies on novel measures for comparing the relative sizes of infinite regular languages. We also propose innovative solutions for the various RE-tree operations including the effective splitting of RE-tree nodes and computing a “tight” bounding RE for a collection of REs. Finally, we demonstrate how sampling-based approximation algorithms can be used to significantly speed up the performance of RE-tree operations. Preliminary experimental results with moderately large synthetic data sets indicate that the RE-tree is effective in pruning the search space and easily outperforms naive sequential search approaches.

Keywords: Regular expressions – Index structure – Size measures – Sampling-based approximations

1 Introduction

Regular expressions (REs) provide an expressive and powerful formalism for capturing the structure of messages, events, and documents. Consequently, they have been used extensively in the specification of a number of languages for important application domains including the XPath pattern language for XML documents [7], the policy language of the border gateway protocol (BGP) for propagating routing information between au-

tonomous systems on the Internet [22], and the UNIX shell’s `grep` utility. Many of these applications have to manage large databases of RE specifications and need to provide an effective matching mechanism that, given an input string, quickly identifies all the REs in the database that match it. This “RE retrieval” problem is important for a variety of software components in the middleware and networking infrastructure of the Internet. We list some of these application domains below.

- **XML filtering.** Information dissemination applications extensively use document-filtering techniques to avoid flooding users with unnecessary information. The key idea is to maintain, for each user, a profile that captures the user’s interests/preferences and transmit to the user only documents that match the profile. While most existing systems for selective dissemination of information typically use simple keywords to specify profiles, the growing momentum of XML as the standard for information exchange on the Internet has recently prompted proposals that use more powerful RE-based languages for expressing user profiles [1]. The primary reason for this is that XML documents are structured, and thus REs provide a succinct syntax for creating more accurate and focused profiles. More specifically, REs can be used to specify matching constraints on the sequence of elements along specific paths in the XML document tree. Systems for filtering XML documents [1,5,11] Represent user interests using XPath expressions [7] that incorporate a restricted form of REs.
- **XML routing.** XML routers route XML documents based on their content. By directing XML traffic to the least loaded application server that is capable of processing a request, an XML router can balance load and boost Web site performance. Again, an RE-based language often provides a succinct and convenient tool for expressing the set of rules for routing incoming XML traffic. As an example, the Intel NetStructure XML Accelerator product [9] uses XPath 1.0 for specifying the rules for directing XML transactions to the appropriate application servers.
- **XML classification.** As XML becomes a popular standard for exchanging data on the Web, the volume of XML-encoded documents on the Web is expected to grow rapidly in coming years. Document type descriptors (DTDs) are RE-based expressions that serve the role of schemas spec-

ifying the internal structure of XML documents. However, DTDs are *not* mandatory, and an XML document may not always have an accompanying DTD. Thus, in many cases, it becomes necessary to identify the subset of a “universe” of known DTDs that match a new XML document. Identifying matching DTDs for XML documents, besides enabling document typing and classification, is also crucial for the efficient storage and effective querying of XML data [10].

- **BGP routing.** In the BGP4 Internet routing protocol [22], routers transmit advertisements to neighboring routers; each advertisement contains a destination IP address reachable from the router as well as the sequence of routing systems on the path from the router to the destination. In case a router receives multiple advertisements (from different neighbors) for the same destination, the BGP4 policy language allows for *priorities* to be assigned to based on the corresponding routing-system sequences. (The advertisement with the highest priority eventually determines the route to the destination). The BGP4 policy language essentially allows for network administrators to specify a set of REs (on routing-system sequences) and associate priorities with each RE. Advertisements containing routing-system sequences that match a specific RE are then assigned the corresponding priority.

Another application involves finding the root-cause of a fault in a network by matching the sequence of events (generated when a network fault occurs) to a database of RE patterns that capture the association between high-level fault patterns and their root causes. Clearly the above-mentioned applications have a critical need for a scalable and efficient structure that can index large numbers of REs and that can quickly retrieve all REs matching a given input string. For the XML-filtering application, the input string is a path in the XML document and the REs are user profiles, while for BGP routing the input string is the routing system sequence in an advertisement and the REs are the BGP policies.

In this paper, we propose the *RE-tree*, a novel index structure for performing fast retrievals of REs that match a given input string. The RE-tree, to the best of our knowledge, is the *first* truly scalable index structure that can handle the storage and retrieval of REs in their full generality. The only prior work along these lines that we are aware of are indexing schemes for filtering XML documents based on XPath expressions [1, 5, 11, 19]. However, while the XPath language allows rich patterns with tree structure to be specified, it lacks the full expressive power of REs (e.g., XPath does not permit the RE operators $*$, $|$ and \cdot to be arbitrarily nested), and thus extending these XML-filtering techniques to handle general REs may not be straightforward. Further, all of the XPath-based methods are designed for indexing main-memory resident data; in contrast, REs in the RE-tree are organized hierarchically (in a manner similar to the R-tree [14]), and thus the RE-tree is well suited for disk-resident data. Another possible main-memory-based approach would be to coalesce the automata for all the REs into a single *nondeterministic finite automaton* and then use this structure to determine the collection of matching REs. It is unclear, however, if the performance of such an approach would be superior to a simple sequential scan over the database

of REs; furthermore, it is not easy to see how such a scheme could be adapted for disk-resident RE data sets.

The task of indexing REs is challenging because REs typically represent *infinite* sets with no well-defined notion of spatial locality. While indexing of documents (which are essentially a finite set/bag of elements) has been extensively studied by the IR community [3], we are not aware of any indexing techniques for infinite sets. The RE-tree employs a number of novel and sophisticated techniques for indexing infinite regular languages, which we list below.

- **Novel measures for comparing the relative sizes of infinite regular languages.** We develop three novel measures for the size of an infinite regular language – the first counts the number of strings in the language that are less than or equal to a certain size, the second computes the “rate-of-growth” of the language between two consecutive windows, and the third *information-theoretic* measure is based on the cost (in bits) of encoding random samples in the language. The latter measure is inspired by a general observation in information theory that the cost of encoding a random element of a set is proportional to the set’s size.
- **Novel algorithms for splitting and generalizing a set of REs.** Like the R-tree, when an RE-tree index node overflows, the set of REs in the node are split and two new compact parent REs that are more general than the two subsets of REs (due to the split) are computed. We show that both splitting and generalizing a set of REs are NP-hard problems and present heuristics for these operations in the context of the RE-tree.
- **Novel sampling-based approximation algorithms for speeding RE-tree operations.** Specifically, we show how random samples of RE languages can be used to efficiently compute *approximate* counts (of the number of strings with a fixed length in the language) and samples for unions/intersections of regular languages. These counts and samples are important for the RE-tree split and generalize operations. Also, we devise a novel practical algorithm that employs a combination of dynamic programming and sampling to compute counts/samples for an RE using its nondeterministic finite automaton representation. Previous algorithms computed these after converting the RE to a deterministic finite automaton (see [17]), which can be fairly expensive.

To measure the effectiveness of the RE-tree in retrieving the set of REs that match an input query string, we conducted an extensive experimental study with synthetic data sets. Preliminary experimental results with moderately large synthetic data sets indicate that the RE-tree can drastically reduce the number of RE-comparison operations and easily outperforms naive sequential-search approaches, improving the overall search performance by up to an order of magnitude.

The remainder of this paper is organized as follows. We first present an overview of the RE-tree in Sect. 2 and identify the key design issues. Section 3 describes some fundamental algorithms for counting and sampling regular languages. Building on these algorithms, we propose two different measures for comparing the relative sizes of infinite regular languages in Sect. 4. Section 5 then describes in detail our algorithms for the various RE-tree operations. In Sect. 7, we present the results of our experimental study comparing the

Table 1. Notation

| Symbol | Description |
|----------------|---|
| Σ | Alphabet over which REs are defined |
| $\delta()$ | Automaton transition function |
| $ M $ | Number of states in automaton M |
| $L(M)$ | Language of automaton M (i.e., strings accepted by M) |
| $L_i(M)$ | Set of length- i strings accepted by automaton M |
| $ L(M) $ | Measure for size of $L(M)$ |
| $M_1 \cup M_2$ | Automaton that accepts strings in $L(M_1) \cup L(M_2)$ |
| $M_1 \cap M_2$ | Automaton that accepts strings in $L(M_1) \cap L(M_2)$ |
| m | Minimum occupancy for each RE-tree node |
| α | Maximum number of states for a bounding automaton |

RE-tree against sequential-search approaches. Finally, Sect. 9 concludes the paper with some directions for future research. The proofs of all theoretical results in this paper can be found in Appendix A.

2 Overview of RE-trees

An RE-tree indexes a large collection of REs such that, given an arbitrary input string w , REs in the collection that match w can be retrieved quickly and efficiently. In this section, we first present an overview of the RE-tree index structure. Although the design principles behind the RE-tree are similar in spirit to those in the R-tree spatial index [14], the application of these principles to indexing REs actually reveals a number of interesting algorithmic issues and trade-offs that require novel techniques. We identify these new design issues in Sect. 2.3, deferring the details of our solutions to Sect. 5. Table 1 summarizes some of the key notation used in this paper.

2.1 Index structure

An RE-tree is a dynamic, height-balanced, hierarchical index structure where leaf nodes contain data entries corresponding to the indexed REs and internal nodes contain “directory” entries that point to nodes at the next level of the index. Specifically, each leaf node entry is of the form (id, M) , where id is the unique identifier of an RE R and M is a *finite automaton* representing R [15]. Each internal node stores a collection of finite automata. And each node entry is of the form (M, ptr) , where M is a finite automaton and ptr is a pointer to some node N (at the next level) such that the following *containment property* is satisfied: if for a set of automata \mathcal{M} , $L(\mathcal{M}) = \bigcup_{M_i \in \mathcal{M}} L(M_i)$, and $\mathcal{M}(N)$ denotes the collection of automata contained in node N , then $L(\mathcal{M}(N)) \subseteq L(M)$. We refer to the automaton M as the *bounding automaton* for $\mathcal{M}(N)$. The containment property is key to improving the search performance of hierarchical index structures like RE-trees: if a query string w is not contained in $L(M)$, then it follows that $w \notin L(M_i)$ for all $M_i \in \mathcal{M}(N)$. As a result, the entire subtree rooted at N can be pruned from the search space. Clearly, the closer $L(M)$ is to $L(\mathcal{M}(N))$, the more effective this search-space pruning will be.

In general, there are an infinite number of bounding automata for $\mathcal{M}(N)$ with different degrees of precision from the least precise bounding automaton with $L(M) = \Sigma^*$ to the most precise bounding automaton, called the *minimal bounding automaton*, with $L(M) = L(\mathcal{M}(N))$.

Since the storage space for an automaton is dependent on its complexity (in terms of the number of its states and transitions), there is a space-precision trade-off involved in the choice of a bounding automaton for each internal node entry.¹ Thus, even though minimal bounding automata result in the best pruning due to their tightness, it may not be desirable (or even feasible) to always store minimal bounding automata in RE-trees since their space requirement can be too large (possibly exceeding the size of an index node), thus resulting in an index structure with a low fan-out. Therefore, to maintain a reasonable fan-out for RE-trees, we impose a space constraint on the maximum number of states (denoted by α) permitted for each bounding automaton in internal RE-tree nodes.

The automata stored in RE-tree nodes are, in general, non-deterministic finite automata (NFAs) with a minimum number of states [15]. Also, for space efficiency, we require that each individual RE-tree node contain at least m entries. An example RE-tree is illustrated in Fig. 1, where only the top three levels of internal nodes are shown. Each internal node has between two and three entries, with each M_i representing a bounding automaton; the details of some of these automata are shown in the table on the right in the form of REs (for conciseness). Note that each pair of parent-child node entries satisfies the containment property.

2.2 Search and maintenance algorithms

Search. The input query to the search algorithm is a string $w \in \Sigma^*$, and the outcome of the search is the set of all indexed REs whose languages contain w . As in an R-tree, searching in an RE-tree proceeds in a top-down manner starting with the root node and might require traversing multiple paths of the index structure. When searching an internal node N , the search string w is compared against each index entry (M, ptr) in N such that if w is accepted by M , then the search is continued at the node pointed by the pointer ptr ; otherwise, the search along that path is terminated. When the search reaches a leaf node N , each automaton M in N is compared against w , and if M accepts w , then M is returned as part of the result.

Insert. The algorithm to insert an automaton M (obtained as a result of converting the input RE R) into an RE-tree is comprised of two main phases. In the first phase, an optimal path of nodes (starting from the root node) to some leaf node N to insert the new automaton M is determined. The goal of performing the insertion along an optimal path is to try and ensure that automata along the path expand as little as possible (when M is inserted into N) since this would ensure better search performance. An optimal path of nodes is determined by choosing the best insertion node at each level (beginning from the root level) using algorithm CHOOSEBESTFA (to be discussed later).

¹ This is in contrast to R-trees, where the space of a bounding rectangle is independent of its precision, and therefore the bounding rectangles in R-trees are all minimal bounding rectangles.

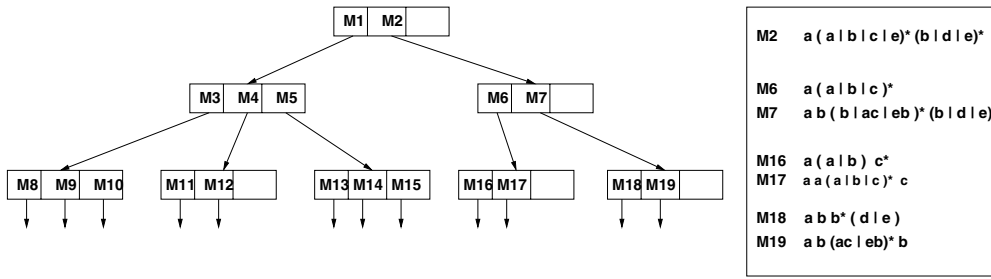


Fig. 1. RE-tree Example

Algorithm Search (N, w)**Input:** N is an index node being searched. w is a query string.**Output:** R is the set of qualified REs (i.e., $r \in R$ iff $w \in L(r)$) contained in the leaf nodes of the subtree rooted at node N .

- 1) $R = \emptyset$;
- 2) **if** (N is an internal node) **then**
- 3) **for each** node entry (M, ptr) **do**
- 4) **if** ($w \in L(M)$) **then**
- 5) $R = R \cup$ Search (N', w), where N' is the node pointed to by ptr ;
- 6) **else**
- 7) **for each** node entry (id, M) **do**
- 8) **if** ($w \in L(M)$) **then**
- 9) $R = R \cup \{M\}$;
- 10) **return** R ;

Fig. 2. Algorithm to search RE-tree

In the second phase, the new automaton M is inserted into the selected leaf node N . This is achieved by the algorithm in Fig. 3 (the input parameter ptr is set to R 's id). If N does not have sufficient space to accommodate M , then in step 6, N is split into another leaf node N' by distributing the collection of automata $\mathcal{M}(N) \cup \{M\}$ between the nodes N and N' using algorithm SPLITFA (to be discussed later).

In the recursive call to INSERT in step 15, two types of insertion updates are propagated along the selected path of nodes in a bottom-up manner from the leaf node N up to the root node ($\&N'$ denotes the address of N' in the step). The first type of insertion updates concerns the maintenance of the containment property for the chain of internal node entries along the selected path of nodes that led to N . Specifically, since the insertion of M has modified the collection of automata in N , we need to update N 's entry in its parent node N_p with a new bounding automaton for $\mathcal{M}(N)$, which is typically a more general (i.e., less precise) automaton than $\bigcup_{M_i \in \mathcal{M}(N)} M_i$, computed using algorithm GENERALIZEFA (to be discussed later); this update in turn modifies $\mathcal{M}(N_p)$, and therefore the update is propagated up to the root node.

The second type of insertion updates deals with node splits. If a node split has occurred at node N , a new internal node entry needs to be inserted into N 's parent node N_p to point to the newly created split node. The insertion of an internal node entry into N_p modifies $\mathcal{M}(N_p)$, which therefore necessitates an update propagation to maintain the containment property. Furthermore, the insertion into N_p might in turn result in further node splits; in the event that the root node itself is split, a new root node is created.

Delete. Deletion of an RE from the RE-tree is carried out in a manner similar to R-trees except that algorithm GENERALIZEFA is used to compute new bounding automata for nodes.

2.3 Design issues

At a high level, RE-trees are conceptually similar to other hierarchical, spatial index structures, like the R-tree [14]; this is also true for RE-tree search and maintenance algorithms. RE-tree search simply proceeds top-down along (possibly) multiple paths whose bounding automaton accepts the input string; RE-tree updates try to identify a “good” leaf node for insertion and can lead to node splits (or node merges for deletions) that can propagate all the way up to the root. There is, however, a fundamental difference between the RE-tree and the R-tree in the indexed data types: regular languages are much more complex objects than multidimensional rectangles. This difference mandates the development of novel algorithmic solutions for the core RE-tree operations. Our main goal for the RE-tree is to optimize search performance, and the key guiding principle for achieving this goal is to keep each bounding automaton M in every internal node as “tight” as possible. Thus, if M is the bounding automaton for $\mathcal{M}(N)$, then $L(M)$ should be as close to $L(\mathcal{M}(N))$ as possible. More specifically, the three core new problems that we need to address in the RE-tree context can be defined as follows (the algorithms shown in the parentheses correspond to our solutions and are described in detail later in the paper).

(P1) Selection of an optimal insertion node (algorithm CHOOSEBESTFA). The goal here is to choose an insertion path for a new RE that leads to “minimal expansion” in the bound-

Algorithm INSERT ($N, (M, ptr)$)**Input:** N is node to be updated. (M, ptr) is entry to be inserted into node N .

- 1) $M' := N' := \emptyset$;
- 2) **if** ($M \neq \emptyset$) **then**
- 3) **if** (N has room for M) **then**
- 4) Insert (M, ptr) into N ;
- 5) **else**
- 6) $(N, N') := \text{SPLITFA}(\mathcal{M}(N) \cup \{M\})$;
- 7) Let $M' := \text{GENERALIZEFA}(\mathcal{M}(N'))$;
- 8) Let $M := \text{GENERALIZEFA}(\mathcal{M}(N))$;
- 9) **if** (N is the root node) **then**
- 10) **if** ($N' \neq \emptyset$) **then**
- 11) Create a new root node N_r with two entries $(M, \&N)$ and $(M', \&N')$;
- 12) **else**
- 13) Let N_p be the parent node of N ;
- 14) Replace N 's old entry in N_p with $(M, \&N)$;
- 15) INSERT ($N_p, (M', \&N')$);

Fig. 3. Algorithm to propagate an insertion update in RE-tree

ing automaton in each internal node. Thus, given the collection of automata $\mathcal{M}(N)$ in an internal index node N and a new automaton M , we need to find the optimal $M_i \in \mathcal{M}(N)$ to insert M such that $|L(M_i) \cup L(M)| - |L(M_i)|$ is minimum (or, equivalently, $|L(M_i) \cap L(M)|$ is maximum).

(P2) Computing an optimal node split (algorithm SPLITFA). When splitting a set of REs during an RE-tree node-split, we seek to identify a partitioning that results in the minimal amount of “covered area” in terms of the languages of the resulting partitions. More formally, given the collection of automata $\mathcal{M} = \{M_1, M_2, \dots, M_k\}$ in an overflowed index node, we want to find the optimal partition of \mathcal{M} into two disjoint subsets \mathcal{M}_1 and \mathcal{M}_2 such that $|\mathcal{M}_1| \geq m$, $|\mathcal{M}_2| \geq m$ and $|L(\mathcal{M}_1)| + |L(\mathcal{M}_2)|$ is minimum.

(P3) Computing an optimal generalized automaton (algorithm GENERALIZEFA). During insertions, node-splits, or node-merges, we need to be able to identify a bounding automaton for a set of REs that does not cover too much “dead space.” Thus, given a collection of automata \mathcal{M} , we seek to find the optimal generalized automaton M such that $|M| \leq \alpha$, $L(\mathcal{M}) \subseteq L(M)$ and $|L(M)|$ is minimum.

The goal of the above operation specifications is to maximize the pruning during search by keeping bounding automata tight. In (P1), the automaton M_i for which $L(M_i)$ expands the least due to the insertion of M is chosen to insert M . The set of automata \mathcal{M} are split into two tight clusters in (P2), while in (P3) we are interested in finding the most precise automaton covering the set of automata in \mathcal{M} and with no more than α states. As discussed earlier, the restriction on the number of states is imposed to keep the fan-out of each node high since this is critical to improving search performance. Essentially (P3) represents the space-precision trade-off for bounding automata in RE-trees.

Note that while (P3) is unique to RE-trees, both (P1) and (P2) have their equivalents in R-trees. Examples of heuristics that have been proposed for (P1) in R-trees include minimizing the increase in the area of the minimum bounding rectangle (MBR) [14] and minimizing the overlap among the MBRs within the node [4]. The heuristics for (P2) in R-trees are

similar to those for (P1); examples include minimizing the total area of the two MBRs [14] and minimizing the area of the intersection between the two MBRs [4]. The aim of all these heuristics is to minimize the number of visits to nodes that do not lead to any qualifying data entries.

Although the “MBRs” in RE-trees (which correspond to regular languages) are very different from the MBRs in R-trees, the intuition behind minimizing the area of MBRs (total area or overlapping area) in R-trees should be effective for RE-trees as well. The counterpart for area in an RE-tree is $|L(M)|$, the size of the regular language for M . However, since a regular language is generally an infinite set, we need to develop new measures for the size of a regular language or for comparing the sizes of two regular languages.

3 Counting and sampling $L_n(M)$

Before defining size measures for infinite regular languages in the next section, we first describe two fundamental algorithms for counting and sampling that form the basis of our definitions. The first algorithm counts the number of length- n strings accepted by an automaton M (i.e., $|L_n(M)|$), while the second algorithm generates a random sample of $L_n(M)$. We begin by presenting the algorithms, originally proposed in [17], for the simpler case when M is a DFA. We then present novel practical counting and sampling algorithms for the case when M is an NFA.

Let $\beta(s, n)$ denote the number of distinct length- n paths that can be generated using an automaton M (with start state s_0) from state s to any accepting state.

3.1 Algorithms for DFAs

For the case where M is a DFA, $|L_n(M)| = \beta(s_0, n)$. Clearly, for $n = 0$, we have $\beta(s, 0) = 1$ if s is an accepting state, and 0 otherwise. For $n \geq 1$, $\beta(s, n)$ can be computed recursively as follows [17]: $\beta(s, n) = \sum_{x \in \Sigma, t = \delta(s, x)} \beta(t, n - 1)$. Thus, dynamic programming can be used to compute $\beta(s, i)$

Algorithm GENRANDOMSTRING (M, s, n)
Input: M is a DFA, s is a state in M .
 n is a length parameter, $n \geq 1$.
Output: A random string in $L_n(M)$.
1) **for** $i := n$ **downto** 1 **do**
2) Select transition x out of s from Σ with probability $\frac{\beta(t, i-1)}{\beta(s, i)}$, where $t = \delta(s, x)$;
3) Let y_i be the selected transition;
4) $s := \delta(s, y_i)$;
5) **return** $y_n \cdot y_{n-1} \cdots y_1$;

Fig. 4. Algorithm for generating a random string in $L_n(M)$

for all states s , first for $i = 0$ and then for successively increasing values of i (until $i = n$) using the computed results for $i - 1$. Since each $\beta(s, i)$ is computed by considering all states reachable from state s with a single transition, each $\beta(\cdot)$ value can be computed in $O(\min\{|\Sigma|, |M|\})$ time. Thus, since there are $|M|n$ values in $\beta(\cdot)$, it follows that there is an $O(n|M|\min\{|\Sigma|, |M|\})$ algorithm to compute $|L_i(M)|$ for all $1 \leq i \leq n$.

We now explain how to generate a random string from $L_n(M)$ based on the computed values of $\beta(\cdot)$. Consider the collection of all length- n strings that can be generated from some state s to any accepting state. Then the probability that a randomly selected string from this collection has $x \in \Sigma$ as the first symbol is given by $\frac{\beta(t, n-1)}{\beta(s, n)}$, where $t = \delta(s, x)$. Thus, a uniformly random string from $L_n(M)$ can be generated iteratively by randomly choosing a transition at each state beginning from the start state, as shown in Fig. 4. Algorithm GENRANDOMSTRING in the figure when invoked with input parameter $s = s_0$ returns a random string from $L_n(M)$. It is straightforward to show that if $s_0, s_n, s_{n-1}, \dots, s_1$ denote the sequence of states visited by the algorithm, then the probability that a string $w \in L_n(M)$ is returned by algorithm GENRANDOMSTRING is $\frac{\beta(s_n, n-1)}{\beta(s_0, n)} \times \frac{\beta(s_{n-1}, n-2)}{\beta(s_{n-1}, n-1)} \times \dots \times \frac{\beta(s_1, 0)}{\beta(s_1, 1)} = \frac{1}{\beta(s_0, n)} = \frac{1}{|L_n(M)|}$. Note that the time complexity of the algorithm is $O(n)$ if the function $\beta(\cdot)$ has been precomputed. Also, a uniform random sample of $L_n(M)$ of size k can be generated by repeatedly invoking GENRANDOMSTRING until it has returned k distinct strings.

3.2 Algorithms for NFAs

For the case where M is an NFA², $\beta(s_0, n) \geq |L_n(M)|$ since there can be multiple accepting paths in an NFA for a given string. However, the problem of computing $|L_n(M)|$ for an NFA M is #P-complete³ [17]. An unbiased estimator for $|L_n(M)|$ can be computed as follows. Let p be a uniformly

² The ϵ -transitions in M are assumed to be eliminated.

³ An enumeration problem belongs to #P if there is a nondeterministic algorithm such that for each instance I of the problem, the number of distinct “guesses” that lead to acceptance of I is exactly $|S(I)|$ (here, $S(I)$ is the number of candidate solutions for I) and such that the length of the longest accepting computation is bounded by a polynomial in size of I [12]. The class of #P-complete enumeration problems are believed to be even harder than NP-hard problems.

generated accepting path of length n in M , and let w be the string labeling p . Then $|L_n(M)| \approx \frac{\beta(s_0, n)}{q}$, where q is the number of accepting paths of w in M . Essentially, the unbiased estimator is computed by scaling down the total number of length- n accepting paths in M with the count of the number of accepting paths for a randomly generated length- n string. Note that since there is a one-to-one correspondence between an accepting string and an accepting path in a DFA, algorithm GENRANDOMSTRING (in Fig. 4) can be used to generate the path p . The number of accepting paths for w can be derived by traversing M with w .

However, as pointed out in [17], the above estimator can have a very large standard deviation. Although [17] proposes a more accurate, randomized algorithm for approximating $|L_n(M)|$, it is not very useful in practice due to its high time complexity of $O(n^{\log(n)})$.

In Fig. 5, we present a novel and practical algorithm for approximately counting strings in $L_n(M)$ for an NFA M that has a lower time complexity of $O(n^2|M|^2 \min\{|\Sigma|, |M|\})$. Like earlier approaches, the algorithm uses dynamic programming to compute counts $\beta(s, i)$. However, instead of computing in each $\beta(s, i)$, the number of length- i accepting paths from state s , the algorithm adjusts this total count of accepting paths by eliminating duplicate paths (recall that in an NFA, there can be multiple accepting paths that correspond to the same string). Thus, in our algorithm, each $\beta(s, i)$ captures $|L_i(M, s)|$, the number of length- i accepting strings in M from state s , more accurately.

When computing $\beta(s, i)$ in steps 5–11, algorithm COUNTSTRINGS subtracts from the total count of accepting paths $\sum_{x \in \Sigma, t \in \delta(s, x)} \beta(t, i-1)$ (due to the dynamic programming relationship)⁴, the number of paths for the same string that are counted multiple times. Thus, the key problem is estimating the number of these duplicate length- i accepting paths from s , which we solve as follows. First, observe that it is not possible for two paths to be identical if they are such due to transitions associated with different symbols; thus, duplicate paths can only result due to transitions out of s on the same symbol. Consequently, we only need to perform duplicate elimination from among the set of paths whose first transitions have the same symbol.

Suppose that t_1, t_2, \dots, t_l are the states due to transitions out of s on symbol x . Note that each $\beta(t_j, i-1)$ is an estimate of $|L_{i-1}(M, t_j)|$, which is the number of distinct accepting paths of length $i-1$ from state t_j . However, for a pair of states t_j, t_k , it is possible that $L_{i-1}(M, t_j)$ and $L_{i-1}(M, t_k)$ have strings in common. The two paths from t_j and t_k for each common string are identical and, with the two transitions from s to t_j and t_k (on symbol x), result in duplicate paths from s . Our strategy for eliminating such duplicates is, for each t_j , to subtract from $\beta(t_j, i-1)$ the number of strings in $L_{i-1}(M, t_j)$ that have already been counted earlier in $L_{i-1}(M, t_1), L_{i-1}(M, t_2), \dots, L_{i-1}(M, t_{j-1})$. In the algorithm, we estimate the number of these strings by generating a random sample r of $L_{i-1}(M, t_j)$ and scaling $\beta(t_j, i-1)$ by the fraction of r not contained in $\cup_{k=1}^{j-1} L_{i-1}(M, t_k)$. Thus, $\beta(t_j, i-1) * \frac{|r - (\cup_{k=1}^{j-1} L_{i-1}(M, t_k))|}{|r|}$ is a good estimate of the

⁴ Note that for an NFA M , $\delta(s, x)$ is a set of states.

Algorithm COUNTSTRINGS (M, n)
Input: M is an NFA, n is a length parameter, $n \geq 1$.
Output: Counts β of accepting strings of length less than or equal to n for each state in M .

- 1) **for each** state $s \in M$ **do** $\beta(s, 0) := 0$;
- 2) **for each** accepting state $s \in M$ **do** $\beta(s, 0) := 1$;
- 3) **for** $i := 1$ **to** n **do**
- 4) **for each** state $s \in M$ **do**
- 5) $\beta(s, i) := 0$;
- 6) **for each** symbol x such that there is a transition out of s on x **do**
- 7) Let t_1, \dots, t_l be the resulting states for transitions out of s on symbol x ;
- 8) $\beta(s, i) := \beta(s, i) + \beta(t_1, i - 1)$;
- 9) **for** $j := 2$ **to** l **do**
- 10) Let r be a random sample of $L_{i-1}(M, t_j)$ (generated by repeatedly invoking GENRANDOMSTRING ($M, t_j, i - 1$) a fixed number of times);
- 11) $\beta(s, i) := \beta(s, i) + \beta(t_j, i - 1) \times \frac{|r - (\cup_{k=1}^{j-1} L_{i-1}(M, t_k))|}{|r|}$;
- 12) **return** $\beta()$;

Fig. 5. Algorithm for computing approximate count of $L_n(M)$

count of strings in $L_{i-1}(M, t_j)$ not previously counted in $\cup_{k=1}^{j-1} L_{i-1}(M, t_k)$, and is added to $\beta(s, i)$ in step 11.

We must point out that each string returned by GENRANDOMSTRING (in step 10) is based on the previously computed counts $\beta()$, which only estimate the number of accepting strings. Thus, the sample r may not be a uniform random sample. Further, GENRANDOMSTRING assumes that M is a DFA and so it needs to be modified when M is an NFA. Details of the required modifications for NFAs can be found in Appendix B.

The overhead of generating the random sample r of length- $(i - 1)$ strings for states t_1, t_2, \dots, t_l increases the time complexity of the dynamic programming algorithm by a factor of $n|M|$ in the worst case, resulting in a worst-case time complexity of $O(n^2|M|^2 \min\{|\Sigma|, |M|\})$ for the algorithm. This is because $i \leq n$ and $l \leq |M|$. Note that the algorithm can also handle DFAs very efficiently and its time complexity reduces to $O(n|M| \min\{|\Sigma|, |M|\})$ if M is a DFA. The reason for this is that for a DFA, l is always 1, and as a result random samples r do not need to be generated.

4 Size measures for infinite regular languages

Estimating the size of a regular language $L(M)$ is much more difficult than computing $|L_n(M)|$ since, unlike $L_n(M)$, $L(M)$ can be an infinite set. In this section, we define three different measures that attempt to capture the size of infinite regular languages. We note that the techniques presented here have applications beyond indexing REs and can also be used for estimating the selectivities of REs, clustering REs, etc.

4.1 Definitions

Before we proceed to define our measures for the size of $L(M)$, we identify certain desirable properties for such measures. First, we note that it does not make sense to actually count the number of strings in an infinite regular language.

However, we do know that while it is impossible to assign a precise integer size to an infinite language, not all infinite languages are equal with respect to size. For example, the language for RE $(a|b)^*$ is *larger than* the language for $a(a|b)^*$ (since the latter is a proper subset of the former). Thus, we would like to define a measure for the size of a language that reflects our intuition of the “larger than” relationship between languages. We denote this measure of the size of $L(M)$ by $|L(M)|$.

We can formalize our intuition of the “larger than” relationship between regular languages as follows.

Definition 4.1. For a pair of automata M_i, M_j , we say that $L(M_i)$ is larger than $L(M_j)$ if there exists a positive integer N such that for all $k \geq N$, $\sum_{l=1}^k |L_l(M_i)| > \sum_{l=1}^k |L_l(M_j)|$.

The various definitions for measure $|L(M)|$ of the size of $L(M)$ that we present in the following subsections attempt to capture this “larger than” relationship. Specifically, for a pair of automata M_i, M_j , if $L(M_i)$ is larger than $L(M_j)$, then $|L(M_i)| > |L(M_j)|$.

4.2 Max-count measure

An obvious measure for the size of a regular language $L(M)$ is to count $L(M)$ up to some maximum length λ , i.e., $|L(M)| = |L_1(M)| + |L_2(M)| + \dots + |L_\lambda(M)|$. Clearly, the larger we choose λ to be, the more effective we can expect the max-count measure to be at reflecting the larger than relationship. However, for very large values of λ , the max-count measure may not be practical. Thus, a good compromise is to set λ to be equal to or slightly greater than $|M|$. This ensures that strings due to all accepting states and traversing a fair proportion of paths/cycles in M are counted in $|L(M)|$. This measure is particularly useful for applications where the maximum length of the query strings is known and its value is not too large. For such cases, λ can be set to the maximum query length.

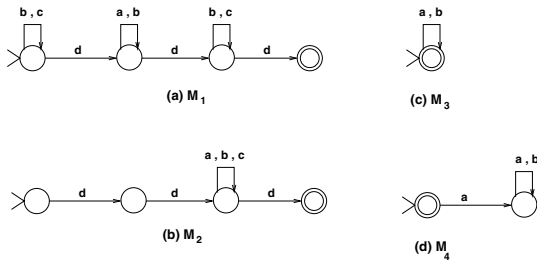


Fig. 6. Examples of automata

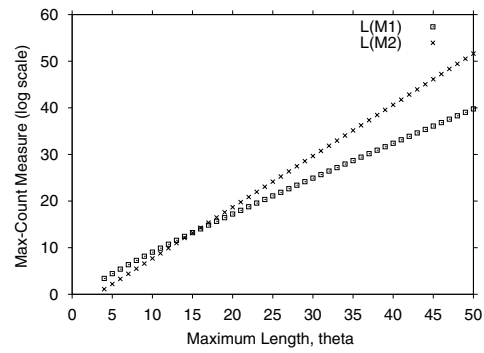
4.3 Rate-of-growth measure

For applications where information about the maximum length of the query strings is not known a priori, using the max-count measure can be problematic due to its sensitivity to the maximum length parameter value λ .

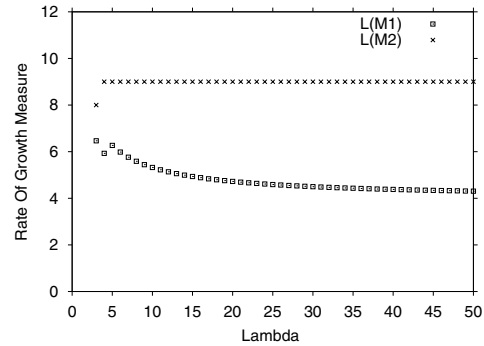
Example 4.2. Consider the two automata M_1 and M_2 in Fig. 6. Figure 7a compares $|L(M_1)|$ and $|L(M_2)|$ using the max-count measure for $3 \leq \lambda \leq 50$. Clearly, $L(M_2)$ is larger than $L(M_1)$ since $|L_n(M_2)|$ grows much faster than $|L_n(M_1)|$ as n increases. In fact, one can show that $|L_{n+3}(M_1)| = (n+1)(n+2)2^{(n-1)}$ and $|L_{n+3}(M_2)| = 3^n$ for $n > 0$. The crossover point for $|L(M_1)|$ and $|L(M_2)|$ using the max-count measure occurs for a maximum length of $\lambda = 16$, that is, $\sum_{l=1}^{\lambda} |L_l(M_2)| > \sum_{l=1}^{\lambda} |L_l(M_1)|$ for $\lambda \geq 16$. Thus, if the value of λ for the max-count measure is set at less than 16, then $L(M_1)$ would be considered, incorrectly, to be larger than $L(M_2)$.

To obtain a more robust measure of the size of infinite languages, we propose a second metric that attempts to capture the “rate-of-growth” of infinite regular languages. Intuitively, for a pair of automata M_i, M_j , if $|L_n(M_i)|$ grows at a faster rate than $|L_n(M_j)|$, then $L(M_i)$ must be larger than $L(M_j)$ since there would exist an N such that for all $k \geq N$, $\sum_{l=1}^k |L_l(M_i)| > \sum_{l=1}^k |L_l(M_j)|$. Thus, the rate of growth of a language can be a good measure of its size. In general, given an unknown, nonnegative function $f(\cdot)$, we can determine how fast $f(\cdot)$ is growing by computing the derivative of its cumulative function; i.e., by computing $F'(\lambda, \theta) = \frac{F(\lambda+\theta) - F(\lambda)}{\theta}$, where $F(k) = \sum_{l=1}^k f(l)$ is an increasing function. Applying this idea to compute the Rate-Of-growth measure of a regular language $L(M)$, we have $|L(M)| = F'(\lambda, \theta)$, where $f(l) = |L_l(M)|$; i.e., $|L(M)| = \frac{\sum_{l=\lambda+1}^{\lambda+\theta} |L_l(M)|}{\theta}$. However, like the max-count case, the sum $\sum_{l=\lambda+1}^{\lambda+\theta} |L_l(M)|$ may be sensitive to the values of λ and θ as well. For instance, for the two automata M_1 and M_2 in Example 4.2, if $\theta = 2$, then we have (for $\lambda = 13$) $|L_{13}(M_1)| + |L_{14}(M_1)| = 199168 > |L_{13}(M_2)| + |L_{14}(M_2)| = 157464$, but (for $\lambda = 14$) $|L_{14}(M_1)| + |L_{15}(M_1)| = 449536 < |L_{14}(M_2)| + |L_{15}(M_2)| = 472392$.

The derivative approach based on function F' works well for functions f with linear growth rates. However, since the rate of growth of an infinite regular language is generally non-linear, we propose an approximate measure of the size of a language as the rate of change of its size from one “window”



a



b

Fig. 7a,b. Comparison of size measures for $L(M_1)$ and $L(M_2)$. a Max-count measure. b Rate-of-growth measure

of lengths to the next consecutive “window” of lengths; that is, our rate-of-growth measure for the size of a regular language $L(M)$ is as follows:

$$|L(M)| = \frac{\sum_{l=\lambda+\theta}^{\lambda+2\theta-1} |L_l(M)|}{\sum_{l=\lambda}^{\lambda+\theta-1} |L_l(M)|} \quad (1)$$

where λ is some length parameter that denotes the start of the first window and θ is a window-size parameter. Essentially, Eq. 1 gives an indication of the rate of growth of $L(M)$, which alleviates some of the drawbacks of the max-count measure as illustrated in the following example.

Example 4.3. Consider again the two automata M_1 and M_2 in Fig. 6. Figure 7b compares $|L(M_1)|$ and $|L(M_2)|$ using the rate-of-growth measure, for $3 \leq \lambda \leq 50$ and $\theta = 2$ (note that similar trends are observed for larger values of θ). The graph correctly indicates that $L(M_2)$ is larger than $L(M_1)$ irrespective of the value for λ .

Like the max-count case, in order to ensure that strings involving a substantial portion of paths, cycles, and accepting states are counted in each window, parameters λ and θ should be chosen to be slightly greater than $|M|$.

4.4 Minimum description length (MDL)-based measure

Although the rate-of-growth measure appears to be a more robust metric than the max-count measure, there are cases where the rate-of-growth measure fails to capture the “larger than” relationship between regular languages, as illustrated by the following example.

Example 4.4. Consider the two automata M_3 and M_4 in Fig. 6. It is straightforward to observe that $|L_n(M_3)| = 2^n$ and $|L_n(M_4)| = 2^{n-1}$ since $L(M_4)$ only consists of strings that begin with a while $L(M_3)$ consists of all strings. Thus, for all $k \geq 1$, $\sum_{l=1}^k |L_l(M_3)| > \sum_{l=1}^k |L_l(M_4)|$ and $L(M_3)$ is larger than $L(M_4)$. However, using the rate-of-growth measure, it is not possible to determine which of $L(M_3)$ or $L(M_4)$ is larger. This is because for any value of λ and θ , $|L(M_3)|$ and $|L(M_4)|$ are both equal to 2^θ with the rate-of-growth measure.

To obtain a more robust measure of the size of infinite languages, we present an alternative metric that attempts to address some of the shortcomings of the max-count measure and is based on Rissanen’s *Minimum description length* (MDL) principle [21]. The MDL principle has been applied to a variety of problems (e.g., constructing decision trees [20], learning common patterns in a set of strings [18], inferring DTDs from a collection of XML data [13]). An important observation made in [13] is that for two given REs R_i and R_j , if R_i defines a larger language than R_j (i.e., R_i is less precise than R_j with respect to the same collection of input data sequences), then the cost of encoding an input data sequence using R_i is likely to be higher than using R_j . This observation is consistent with information theory since, in general, more bits are needed to specify an item that comes from a larger collection of items.

Our use of the MDL principle to define a measure of the size of a regular language is also inspired by a similar observation and is based on the following intuition: given two DFAs M_i and M_j , if $L(M_i)$ is larger than $L(M_j)$, then the *per-symbol-cost* of an MDL-based encoding of a random string in $L(M_i)$ using M_i is very likely to be higher than that of a string in $L(M_j)$ using M_j . The per-symbol-cost of encoding a string $w \in L(M)$ is essentially the ratio of the cost of an MDL-based encoding of w using M (defined below) to the length of w . Therefore, a reasonable measure for the size of a regular language $L(M)$ is the expected per-symbol-cost of an MDL-based encoding for a random sample of strings in $L(M)$. Let $\text{MDL}(M, w)$ denote the cost of an MDL-based encoding of a string $w \in L(M)$ using M and let r be a random sample of $L(M)$. The MDL-based measure for the size of $L(M)$ is then as Follows.⁵

$$|L(M)| = \frac{1}{|r|} \sum_{w \in r} \frac{\text{MDL}(M, w)}{|w|} \quad (2)$$

We next define the cost of an MDL-encoding of w using M . Suppose that $w = w_1.w_2.\dots.w_n \in L(M)$ and s_0, s_1, \dots, s_n is the unique sequence of states visited by w in M . Then,

$$\text{MDL}(M, w) = \sum_{i=0}^{n-1} \log_2(n_i) \quad (3)$$

where each n_i denotes the number of transitions out of state s_i , and $\log_2(n_i)$ is the number of bits required to specify the transition out of state s_i . Since the above formulation is based on counting the number of transitions, to obtain an accurate measure of the size of an infinite language, it is important that M does not contain any “nonessential” transitions (i.e., transitions that are not part of an accepting path that involves at least one cycle).⁶ Thus, M should be a minimal-state DFA

⁵ $|w|$ denotes the length string w .

⁶ For example, in Fig. 10a, if the state labeled 5 does not have an a self-loop transition, then both the a -transition from state 2 to state 5

without any nonessential transitions. The following example illustrates how the MDL-based encoding costs are computed.

Example 4.5. Consider the two automata M_1 and M_2 in Fig. 6 and two length-10 strings $w_1 = bbbdbbddd \in L(M_1)$ and $w_2 = ddbbbbbb \in L(M_2)$. The per-symbol-costs of encoding w_1 and w_2 using M_1 and M_2 , respectively, are as follows:

$$\frac{\text{MDL}(M_1, w_1)}{|w_1|} = \frac{10 \times \log_2(3)}{10} \approx 1.5850 \quad \text{and}$$

$$\frac{\text{MDL}(M_2, w_2)}{|w_2|} = \frac{(2 \times 0) + (8 \times \log_2(4))}{10} \approx 1.6000$$

The computed per-symbol-costs correctly indicate that $L(M_2)$ is larger than $L(M_1)$.

We still need to address the issue of generating the random sample r of $L(M)$ to encode. This can be carried out by selecting two values λ and θ , both slightly greater than $|M|$ and for a desired random sample of size k , generating $\frac{k|L_i(M)|}{\sum_{l=\lambda}^{\lambda+\theta-1} |L_l(M)|}$ random strings from each $L_i(M)$, $\lambda \leq i \leq \lambda + \theta - 1$ (using algorithm GENRANDOMSTRING). By sampling from each $L_i(M)$ in proportion to its count, we ensure that paths in the automaton that are traversed more frequently by accepting strings have a greater likelihood of being included in the sample r .

5 Algorithms for RE-tree operations

Now that we have robust measures for comparing the relative sizes of infinite regular languages, we are in a position to present details of the three RE-tree operations for choosing the best automaton, splitting a set of automata and computing a generalized automaton [problems (P1), (P2), and (P3) from Sect. 2.3]. The algorithms for the RE-tree operations presented in this section perform a number of standard operations on automata like union, intersection, and conversion of an NFA to a DFA. Efficient algorithms for these automata operations can be found in [15]. Note that while it is possible to compute the union/intersection of a pair of automata M_i, M_j in time proportional to $|M_i||M_j|$, the worst-case time complexity of converting an NFA M to a DFA can be exponential in $|M|$.

5.1 Algorithm CHOOSEBESTFA

From among the automata in $\mathcal{M}(N)$ contained in a node N , algorithm CHOOSEBESTFA returns the automaton M_i for which $|L(M_i) \cap L(M)|$ is maximum. For each $M_i \in \mathcal{M}(N)$, the algorithm constructs the DFA $M \cap M_i$ and computes $|L(M \cap M_i)|$ using either the max-count or the MDL metric from the previous section. The automaton M_i for which the language $L(M \cap M_i)$ is the largest is chosen for inserting M into the RE-tree.

as well as the b -transition from state 5 to state 6 are nonessential. The removal of such nonessential transitions from an automaton does not affect the infiniteness of its language.

Algorithm SPLITFA (\mathcal{M})**Input:** \mathcal{M} is a set of automata to be split.**Output:** Two sets of automata \mathcal{M}_1 and \mathcal{M}_2 resulting from the split.

- 1) Let $M_i, M_j \in \mathcal{M}$ be the pair of automata such that $|L(M_i \cup M_j)| - |L(M_i \cap M_j)|$ is maximum;
- 2) $\mathcal{M}_1 := \{M_i\}, \mathcal{M}_2 := \{M_j\}$;
- 3) $\mathcal{M} := \mathcal{M} - \{M_i, M_j\}$;
- 4) **while** ($\mathcal{M} \neq \emptyset$) **do**
- 5) **if** ($|\mathcal{M}| = m - |\mathcal{M}_1|$) **then**
- 6) $\mathcal{M}_1 := \mathcal{M}_1 \cup \mathcal{M}$;
- 7) **break**;
- 8) **if** ($|\mathcal{M}| = m - |\mathcal{M}_2|$) **then**
- 9) $\mathcal{M}_2 := \mathcal{M}_2 \cup \mathcal{M}$;
- 10) **break**;
- 11) Let $M_i \in \mathcal{M}$ be the automaton for which $\min\{|L(\mathcal{M}_1 \cup \{M_i\})| - |L(\mathcal{M}_1)|, |L(\mathcal{M}_2 \cup \{M_i\})| - |L(\mathcal{M}_2)|\}$ is minimum;
- 12) **if** ($|L(\mathcal{M}_1 \cup \{M_i\})| - |L(\mathcal{M}_1)| \leq |L(\mathcal{M}_2 \cup \{M_i\})| - |L(\mathcal{M}_2)|$) **then**
- 13) $\mathcal{M}_1 := \mathcal{M}_1 \cup \{M_i\}$;
- 14) **else**
- 15) $\mathcal{M}_2 := \mathcal{M}_2 \cup \{M_i\}$;
- 16) $\mathcal{M} := \mathcal{M} - \{M_i\}$;
- 17) **return** ($\mathcal{M}_1, \mathcal{M}_2$);

Fig. 8. Algorithm for splitting a set of automata**Algorithm GENERALIZEFA** (\mathcal{M}, α)**Input:** \mathcal{M} is a set of automata to be generalized. α is maximum number of states in generalized automaton.**Output:** The generalized automaton for \mathcal{M} .

- 1) Compute DFA M for $\cup_{M_i \in \mathcal{M}} M_i$
- 2) **while** ($|M| > \alpha$) **do**
 /* $M_{(s_i, s_j)}$ is the resulting automaton when states s_i, s_j in M are merged */
- 3) Let s_i, s_j be the pair of states in M for which $|L(M_{(s_i, s_j)})|$ is minimum among all pairs of states in M ;
- 4) Set M to be equal to the DFA for $M_{(s_i, s_j)}$;
- 5) **return** M ;

Fig. 9. Algorithm for generalizing a set of automata

5.2 Algorithm SPLITFA

For a set of automata $\mathcal{M} = \{M_1, \dots, M_k\}$, algorithm SPLITFA partitions \mathcal{M} into \mathcal{M}_1 and \mathcal{M}_2 such that $|\mathcal{M}_1| \geq m$, $|\mathcal{M}_2| \geq m$ and $|L(\mathcal{M}_1)| + |L(\mathcal{M}_2)|$ is minimum. Unfortunately, it can be shown that even if $L(M_i)$ for every $M_i \in \mathcal{M}$ is finite, the problem of optimally splitting \mathcal{M} is NP-hard (reduction from clique).⁷ Thus, one can expect that for infinite languages the problem is even more difficult.

Theorem 5.1. *Given finite sets of elements S_1, \dots, S_n , the problem of partitioning them into two sets \mathcal{S}_1 and \mathcal{S}_2 (each containing at least m elements) such that $|\cup_{S_i \in \mathcal{S}_1} S_i| + |\cup_{S_i \in \mathcal{S}_2} S_i|$ is minimum, is NP-hard.*

Since the problem of computing the optimal partitioning of \mathcal{M} is NP-hard, we resort to heuristics to generate the two sets \mathcal{M}_1 and \mathcal{M}_2 . Figure 8 contains the steps of the SPLITFA algorithm for splitting \mathcal{M} into two compact well-separated subsets.

⁷ Given a graph and a constant k , an instance of the clique problem seeks to determine whether the graph contains a clique of size at least k . (A clique is a set of nodes such that there is an edge between every pair of nodes.)

Algorithm SPLITFA is similar in spirit to the Quadratic Split algorithm from [14]; however, instead of trying to minimize the area of MBRs, our splitting algorithm attempts to reduce the size of regular languages. The algorithm begins by picking as seeds (in step 1) the two automata from \mathcal{M} whose languages have large nonoverlapping portions and then greedily assigns each remaining automaton M_i in \mathcal{M} to the set whose language increases the least due to the addition of M_i .

Note that step 1 requires $O(|\mathcal{M}|^2)$ automata intersection, union, and size measure computations to be performed, one for each pair of automata in \mathcal{M} . For efficiency purposes, the algorithm caches DFAs for \mathcal{M}_1 and \mathcal{M}_2 as well as $|L(\mathcal{M}_1)|$ and $|L(\mathcal{M}_2)|$ between successive iterations of the while loop. Thus, in step 11, the algorithm performs $2(|\mathcal{M}| - i)$ automata union and size measure computations in the i^{th} iteration of the while loop, two for each of the remaining automata in \mathcal{M} .

5.3 Algorithm GENERALIZEFA

Recall from Sect. 2.3 that the objective of the GENERALIZEFA algorithm is to generate for a set of automata \mathcal{M} an automaton M with no more than α states such that $L(\mathcal{M}) \subseteq L(M)$ and

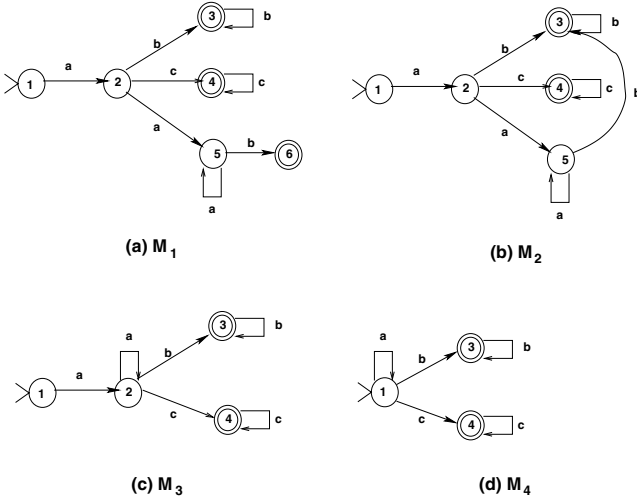


Fig. 10. Automaton generalization example

$|L(M)|$ are minimum. This problem can be shown to be NP-hard by reducing the partition problem [12] to it.

Theorem 5.2. *Given a set of automata \mathcal{M} , the problem of computing a DFA M with no more than α states and for which $L(\mathcal{M}) \subseteq L(M)$ and $|L(M)|$ are minimum is NP-hard.*

The algorithm for generalizing a set of automata \mathcal{M} is shown in Fig. 9. The algorithm greedily merges pairs of states (that result in the smallest regular languages) in the union automaton M until $|M| \leq \alpha$. Note that the automaton $M_{(s_i, s_j)}$ resulting from merging states s_i and s_j in M may not be a DFA. Thus, if we want to compute accurate counts and samples for $L_l(M_{(s_i, s_j)})$, we may need to convert $M_{(s_i, s_j)}$ into a DFA first. However, an option with lower overhead, if approximate counts and samples are acceptable, is to directly apply to $M_{(s_i, s_j)}$ the counting and sampling algorithms for NFAs described in Sect. 3. This would avoid converting NFA $M_{(s_i, s_j)}$ to a DFA, an operation which could be potentially expensive and also cause the number of states in the automaton to increase.

A further optimization is to not consider all pairs of states s_i, s_j in M as candidates for merging in step 3. One heuristic is to consider as candidates only those state pairs s_i, s_j that have similar incoming and outgoing transitions. The rationale here is that merging such similar states is more likely to result in automata with compact languages.

Example 5.3. Suppose we want to generalize the three automata for the REs abb^* , acc^* , and aa^*b . The DFA for the union of the three automata is depicted in Fig. 10a. Suppose $\alpha = 3$, that is, we would like the final generalized automaton to contain at most three states. We trace the steps of algorithm GENERALIZEFA on the union automaton M_1 . The first pair of states merged by the algorithm are states 3 and 6 since the language for the resulting automaton M_2 in Fig. 10b is $aa^*bb^*|acc^*$, which is smaller than the resulting languages when other pairs of states are merged. For instance, merging states 3 and 4 results in the language $(a(b|c)(b|c)^*)|aa^*b$, while $aa^*(bb^*|cc^*)$ is the language when states 2 and 5 are merged. In the next iteration, states 2 and 5 are merged to yield automaton M_3 in Fig. 10c whose language is $aa^*(bb^*$

$|cc^*$). And in the final iteration, states 1 and 2 are merged, and the final automaton M_4 in Fig. 10d containing three states is returned by the algorithm (i.e., $L(M_4)$ is $aa^*(bb^*|cc^*)$).

6 Optimizations

The RE-tree operations CHOOSEBESTFA and SPLITFA described in previous subsections required frequent computations of $|L(M_i \cap M_j)|$ and $|L(M_i \cup M_j)|$ to be performed for pairs of automata M_i, M_j . These computations can adversely affect RE-tree performance since construction of the intersection and union automaton M can be expensive. Further, the final automaton M may have many more states than the two initial automata M_i and M_j . This could be a problem since computing $|L(M)|$ using any of the size measures from Sect. 4 requires $|L_l(M)|$ for a range of l values to be calculated using algorithm COUNTSTRINGS. [The MDL metric also requires random samples for each $L_l(M)$ to be generated].

In this section, we show how sampling can be used to speed up the performance of the RE-tree operations CHOOSEBESTFA and SPLITFA. Specifically, we show that if we have counts and samples for each $L(M_i)$, then we can utilize this information to derive *approximate* counts and samples for $L(M_i \cap M_j)$ and $L(M_i \cup M_j)$ without incurring the overhead of constructing the automata $M_i \cap M_j$ and $M_i \cup M_j$ and without reinvoking algorithm COUNTSTRINGS on these new automata. Thus, by eliminating the building of union and intersection automata, the sampling-based approximation algorithms have the potential to considerably improve the overall performance of RE-trees.

6.1 Approximate counting and sampling techniques

Before we show how sampling can be used to speed up the computation of $|L(M_i \cap M_j)|$ and $|L(M_i \cup M_j)|$, we present results for approximating the sizes and generating uniform samples of unions and intersections of arbitrary sets. These results for finite sets constitute the building block for our sampling-based schemes for approximating counts and samples of the finite sets $L_l(M_i \cap M_j)$ and $L_l(M_i \cup M_j)$, which in turn form the basis of our three size measures for regular languages.

Counts and samples for set intersections. The following two theorems (see [8]) suggest how counts and samples for the intersection of two sets S_1 and S_2 can be generated using the count and sample information for one of them.

Theorem 6.1. *Suppose r_1 is a uniform random sample of set S_1 . Then $\frac{|r_1 \cap S_2| |S_1|}{|r_1|}$ is an unbiased estimator of the size of $S_1 \cap S_2$.*

Theorem 6.2. *Suppose r_1 is a uniform random sample of set S_1 . Then $r_1 \cap S_2$ is a uniform random sample of $S_1 \cap S_2$ with size $|r_1 \cap S_2|$.*

6.2 Counts and samples for set unions

Algorithm SAMPLEUNION in Fig. 11 returns a uniform random sample r of $S_1 \cup S_2$ under the assumption that S_1 and

S_2 are disjoint sets. Further, the size of sample r is at least $\min\{|r_1|, |r_2|\}$, where r_i is a uniform random sample of S_i . The input parameters to the algorithm include samples r_1 and r_2 of sets S_1 and S_2 , respectively, and the sizes of the two sets $|S_1|$ and $|S_2|$. Note that the sets themselves are not required by the algorithm – only their samples and their sizes. In the algorithm, variables k_1 and k_2 keep track of the number of elements in r that are chosen so far from samples r_1 and r_2 , respectively. During each iteration, in steps 4 and 5, sample r_1 is chosen as the set for selecting the next element from, with probability $\frac{|S_1| - k_1}{|S_1| + |S_2| - k_1 - k_2}$.

Theorem 6.3. *Let r_1 and r_2 be uniform random samples of sets S_1 and S_2 , respectively. If sets S_1 and S_2 are disjoint, then algorithm SAMPLEUNION returns a uniform random sample of $S_1 \cup S_2$.*

Note that the computed sample r contains at least $\min\{|r_1|, |r_2|\}$ elements. Also, from Theorem 6.3, it follows that algorithm SAMPLEUNION returns a uniform random sample of $S_1 \cup S_2$ only if sets S_1 and S_2 are disjoint. A possible alternative for obtaining an approximately uniform random sample of $S_1 \cup S_2$ if sets S_1 and S_2 are not disjoint is as follows. Consider the set $S_3 = S_2 - S_1$. Sets S_1 and S_3 are disjoint and can be passed as arguments to algorithm SAMPLEUNION. We already have r_1 and $|S_1|$. Further, due to Theorem 6.2, $r_2 - S_1$ is a uniform random sample of $S_3 = S_2 - S_1$ (since $r_2 - S_1$ is basically the same as $r_2 \cap (S_2 - S_1)$), which, due to the theorem, is a uniform random sample of $S_2 \cap (S_2 - S_1) = S_2 - S_1$. The problem, however, is that it is not possible to accurately determine $|S_3|$ from $r_1, r_2, |S_1|$ and $|S_2|$. Thus, the only available option is to use an estimate of $|S_3|$ instead of the accurate value. We use the value $\frac{|r_2 - S_1| |S_2|}{|r_2|}$, which can be shown to be an unbiased estimator of $|S_3|$ using Theorem 6.1, as an estimate of $|S_3|$. Thus, by passing the samples and (approximate) sizes for disjoint sets S_1 and S_3 to algorithm SAMPLEUNION, we can obtain an approximately uniform random sample of $S_1 \cup S_2$. Note also that $|S_1 \cup S_2| = |S_1| + |S_3|$, and thus, due to Theorem 6.1, $|S_1| + \frac{|r_2 - S_1| |S_2|}{|r_2|}$ is an unbiased estimator of $|S_1 \cup S_2|$.

6.3 Algorithm CHOOSEBESTFA

The CHOOSEBESTFA algorithm requires that $|L(M \cap M_i)|$ be estimated, where M is the automaton being inserted into the tree and M_i is an automaton in the current RE-tree node. In the following, we show how a random sample of $L_l(M)$ and knowledge of $|L_l(M)|$ for a range of l values can be used to compute $|L(M \cap M_i)|$ very fast for all three size measures proposed earlier. Thus, for an automaton M being inserted into the RE-tree index, our algorithm only requires us to generate counts and samples for $L_l(M)$ for specific values of l .

Max-count. For the max-count measure, we need to compute $|L_l(M \cap M_i)|$ for l values in the range $[1, \lambda]$. We assume that for these l values, we have already computed $|L_l(M)|$ and $\hat{L}_l(M)$, a uniform random sample of $L_l(M)$. From Theorem 6.1 it follows that $\frac{|\hat{L}_l(M) \cap L(M_i)| |L_l(M)|}{|\hat{L}_l(M)|}$ is an

unbiased estimate of $|L_l(M \cap M_i)|$. Here $\hat{L}_l(M) \cap L(M_i)$ can be computed efficiently by simply checking for each string in sample $\hat{L}_l(M)$ if it is accepted by M_i . Thus, with our sampling approach we only need to compute counts and samples for $L_l(M)$ (the automaton being inserted) once at the beginning. There is no need to either construct the automaton for each $M \cap M_i$ or count $L_l(M \cap M_i)$ using the dynamic programming algorithm COUNTSTRINGS.

MDL. For computing $|L(M \cap M_i)|$ using the MDL measure, we need to generate random samples of $L_l(M \cap M_i)$ for $\lambda \leq l \leq \lambda + \theta - 1$. Suppose that $\hat{L}_l(M)$ denotes the random sample of $L_l(M)$. Then, due to Theorem 6.2, $\hat{L}_l(M) \cap L(M_i)$ is a uniform random sample of $L_l(M \cap M_i)$. Thus, using our sampling-based approach, sample $\hat{L}_l(M)$ needs to be computed only once in the beginning. In addition, we also need to construct the intersection automaton for each $M \cap M_i$ to encode the strings in the sample. However, we completely eliminate the need to count $L_l(M \cap M_i)$ using algorithm COUNTSTRINGS or generate new samples for $L_l(M \cap M_i)$ using algorithm GENRANDOMSTRING.

6.4 Algorithm SPLITFA

In step 1 of the SPLITFA algorithm (see Fig. 8), we need to estimate $|L(M_i \cup M_j)|$ and $|L(M_i \cap M_j)|$ for pairs of automata $M_i, M_j \in \mathcal{M}$. Similarly, in step 11, for each remaining automaton $M_i \in \mathcal{M}$ we need to estimate $|L(\mathcal{M}_1 \cup \{M_i\})|$ and $|L(\mathcal{M}_2 \cup \{M_i\})|$. Below, we show how sampling-based techniques can reduce the number of times algorithm COUNTSTRINGS is invoked from $O(|\mathcal{M}|^2)$ down to $O(|\mathcal{M}|)$. Further, with sampling, the DFA's input to the counting algorithm are the much smaller initial automata in \mathcal{M} rather than the union or intersection automata. In the following, we only show how step 1 can be made more efficient; however, the same techniques can also be used to speed up step 11.

Max-count. For this measure we need to compute $|L_l(M_i \cap M_j)|$ and $|L_l(M_i \cup M_j)|$ for a range of l values. Suppose for each $M_i \in \mathcal{M}$ we computed $|L_l(M_i)|$ and a random sample $\hat{L}_l(M_i)$ of $L_l(M_i)$. Then, from Theorem 6.1 it follows that $\frac{|\hat{L}_l(M_i) \cap L(M_j)| |L_l(M_i)|}{|\hat{L}_l(M_i)|}$ is an unbiased estimate of $|L_l(M_i \cap M_j)|$. Also, as described in Sect. 6.1, $|L_l(M_i \cup M_j)|$ can be estimated as $|L_l(M_i)| + \frac{|\hat{L}_l(M_j) - L(M_i)| |L_l(M_j)|}{|\hat{L}_l(M_j)|}$. Thus, our sampling-based approach only requires counts and samples to be generated for $O(|\mathcal{M}|)$ automata as opposed to $O(|\mathcal{M}|^2)$ automata without sampling (see Sect. 5.2). Further, it does not require that a single union and intersection automaton be constructed, while the original approach without sampling required a union and intersection automaton to be constructed for each pair $M_i, M_j \in \mathcal{M}$ and counts to be computed subsequently for these much larger automata.

MDL. For computing $|L(M_i \cap M_j)|$ and $|L(M_i \cup M_j)|$ using the MDL measure we need to generate random samples of $L_l(M \cap M_i)$ for $\lambda \leq l \leq \lambda + \theta - 1$. Suppose for each $M_i \in \mathcal{M}$ we computed $|L_l(M_i)|$ and a random sample $\hat{L}_l(M_i)$ of $L_l(M_i)$. Then, from Theorem 6.2 it fol-

Algorithm SAMPLEUNION ($r_1, r_2, |S_1|, |S_2|$)
Input: r_1, r_2 are random samples of sets S_1 and S_2 .
 $|S_1|, |S_2|$ are sizes of sets S_1 and S_2 .
Output: A random sample of $S_1 \cup S_2$.

- 1) $r := \emptyset$;
- 2) $k_1 := k_2 := 0$;
- 3) **while** ($r_1 \neq \emptyset$ **and** $r_2 \neq \emptyset$) **do**
- 4) Let *rand* be a uniform random real value between 0 and 1;
- 5) **if** ($\text{rand} \leq \frac{|S_1| - k_1}{|S_1| + |S_2| - k_1 - k_2}$) **then**
- 6) $i := 1$;
- 7) **else**
- 8) $i := 2$;
- 9) Choose an element e uniformly and randomly from r_i ;
- 10) $r := r \cup \{e\}$;
- 11) $r_i := r_i - \{e\}$;
- 12) $k_i := k_i + 1$;
- 13) **return** r ;

Fig. 11. Algorithm for computing random sample of the union of two disjoint sets

lows that $\hat{L}_i(M_i) \cap L(M_j)$ is a uniform random sample of $L_i(M \cap M_i)$. Also, an approximately uniform random sample of $L_i(M \cup M_i)$ can be computed using algorithm SAMPLEUNION described in Fig. 11. Thus, our sampling-based approach only requires counts and samples to be generated for $O(|\mathcal{M}|)$ automata as opposed to $O(|\mathcal{M}|^2)$ automata without sampling (see Sect. 5.2). Further, counting and sampling using our approach are performed on the individual automata in \mathcal{M} , which are much smaller than the union and intersection automata on which the original approach performs counting and sampling. Note, however, that our sampling methods do require union and intersection automata to be constructed in order to encode the samples.

7 Experimental evaluation

To determine the effectiveness of the RE-tree, we compare its performance against the sequential file approach, which stores the REs in a flat file and searches the entire file sequentially for each search query. As we noted in Sect. 1, we are not aware of any disk-based indexing method for indexing REs (in their full generality). Our experimental results (based on synthetic data sets) indicate that the RE-tree approach offers a significant performance improvement (up to an order of magnitude) over the sequential search approach.

Data sets: each synthetic data set is comprised of clusters of similar REs that are generated using a synthetic RE generator, where the number of clusters and the size of each cluster are controlled by the input parameters c_{num} and c_{size} , respectively. The REs in each cluster are similar in terms of both their content (i.e., symbol distribution) and their structure (i.e., parse tree). Content-wise, each cluster is associated with some subset of the alphabet $\Sigma' \subseteq \Sigma$, referred to as its *hot alphabet*, such that each RE symbol is drawn from Σ' with a probability of ρ and drawn from $(\Sigma - \Sigma')$ with a probability of $(1 - \rho)$. The parameter ρ applies to all the clusters and is referred to as the *hot probability*. Each cluster is mapped to its hot alphabet as follows: first, the alphabet Σ is arbitrarily partitioned into n_σ disjoint subsets $\Sigma = \bigcup_{i=1}^{n_\sigma} \Sigma_i$ such that each

subset Σ_i consists of $\frac{|\Sigma|}{n_\sigma}$ symbols, where n_σ is another input parameter.⁸ Each cluster is then randomly mapped to one of the n_σ subsets as its hot alphabet.

In terms of structural similarity, each cluster of similar REs is generated by first randomly generating a *seed RE* R_i and then using R_i to derive the other REs in the cluster by making a random modification to the parse tree of R_i . Note that in a parse tree for an RE, the leaf nodes correspond to symbols in Σ while the internal nodes correspond to one of the three operators: union, concatenation, or kleene star. The structure of each seed RE is controlled by two parameters θ_E and θ_e such that each seed RE is of the form $R_i = E_1.E_2.\dots.E_n$, where n is a random integer between 1 and θ_E and each E_j has one of the following four forms: (1) $(e_1|e_2|\dots|e_{n_j})$, (2) $(e_1.e_2.\dots.e_{n_j})$, (3) $(e_1|e_2|\dots|e_{n_j})^*$, or (4) $(e_1.e_2.\dots.e_{n_j})^*$, where n_j is a random integer between 1 and θ_e . A similar RE is derived from a seed RE R_i by adding a new internal node to its parse tree (corresponding to either the kleene star operator or union operator) at a randomly selected location in the parse tree. For instance, when adding a new union operator, one of the operands of the union operator is one of the existing nodes in the parse tree while the other operand is a new leaf node corresponding to a randomly generated symbol. Examples of random regular expressions generated by our synthetic RE generator include “ $g.(d|i) * .(e|a|m)$ ” and “ $(q|b).c.(f|k)^*$ ”.

Queries: For each data set, we generate a set of 1000 random queries, where each query $w \in \Sigma^*$ is generated as follows. First, randomly select an RE R from the data set and randomly generate an integer n between 5 and 10. Next, generate a random length- n string w that matches R ; if no such string exists, we repeat the generation process with another randomly chosen pair of values for R and n .

Algorithms: For the sequential file approach, a file of automata is created for each data set by packing the collection of automata into as few pages as possible. To find matching REs for a query string w requires sequentially checking against each automaton in the file. For the RE-tree approach, our implementation is based on the algorithms presented in

⁸ For simplicity, assume that $|\Sigma|$ is a multiple of n_σ .

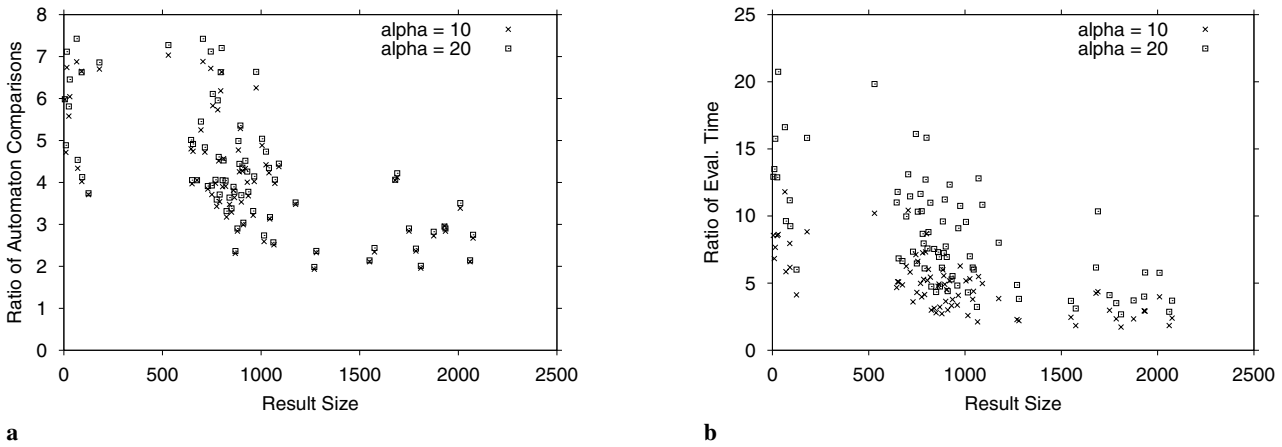


Fig. 12. Varying α , $D = 50000$, $\rho = 0.75$. **a** Ratio of automaton comparisons, $\frac{N_{seq}(n)}{N_{rt}(n)}$. **b** Ratio of evaluation time, $\frac{T_{seq}(n)}{T_{rt}(n)}$

Sects. 2–5. We use the MDL measure for comparing sizes of regular languages since in our experiments we found it to be more effective than the max-count and rate-of-growth metrics. The RE-tree operations are implemented as described in Sect. 5, with the sampling-based optimizations discussed in Sect. 6.

The experiments were conducted on a 2.0-GHz Intel Pentium IV machine with 1 GB memory running Red Hat Linux 7.2. Both the sequential file and RE-tree methods were implemented as Unix files on a local disk. For each query and each method we measured both the evaluation time (including both CPU and I/O times) as well as the number of automaton comparisons (i.e., the number of automata that were checked against the query). Each query was run 10 times (for each method) to compute its average measurement values.

7.1 Experimental results

This section presents experimental results that compare the performance of the two approaches by varying three main parameters: (1) α , the maximum size of a bounding automaton; (2) ρ , the hot probability; and (3) the size of the data set (number of REs) given by $D = c_{num} \times c_{size}$. Table 2 summarizes the values of the parameters used in our experiments.

The performance results are presented in terms of two ratios, $\frac{N_{seq}(n)}{N_{rt}(n)}$ and $\frac{T_{seq}(n)}{T_{rt}(n)}$, where $N_{seq}(n)$ ($N_{rt}(n)$) is the average number of automaton comparisons incurred by the sequential (RE-tree) approach for queries with a result size (that is, number of matching REs) of n , and $T_{seq}(n)$ ($T_{rt}(n)$) is the average evaluation time incurred by the sequential (RE-tree) approach for queries with a result size of n . Note that the number of automaton comparisons incurred by the sequential file approach $N_{seq}(n)$ is always equal to D , the size of the data set.

Figure 12 shows the performance results as α is varied. The graphs indicate that RE-trees outperform the sequential file approach by up to a factor of 7 and 20, respectively, for the number of automaton comparisons and search time. Our results indicate that both the number of automaton comparisons (Fig. 12a) and running time (Fig. 12b) improve with larger values of α . This is because as α increases, the precision of the

Table 2. Experimental parameters and values

| Param | Meaning | Value |
|------------|---|-----------------|
| P | Page Size (in KB) | 4 |
| $ \Sigma $ | Size of alphabet | 20 |
| θ_E | Max. number of first-level expressions (per RE) | 3 |
| θ_e | Max number of second-level symbols (per first-level expression) | 3 |
| n_σ | Number of alphabet subsets | 5 |
| α | Maximum number of states for bounding automata | 10, 20 |
| c_{num} | Number of RE clusters | 250, 500 |
| c_{size} | Size of RE cluster | 200 |
| D | Size of data set given by $c_{num} \times c_{size}$ | 50000, 100000 |
| ρ | Probability that a symbol belongs to “hot” alphabet subset | 0.25, 0.5, 0.75 |

bounding automata generally becomes higher, thereby resulting in better pruning of “false-drops” (i.e., path traversals that do not lead to any qualifying REs) and hence fewer number of automaton comparisons and index page accesses.

Figure 13 depicts the performance results as ρ is varied. Since the similarity of the REs in each cluster becomes higher with a larger value of ρ , the RE-tree approach is able to improve its filtering (with tighter bounding automata) thereby resulting in less false drops and hence improved query evaluation. Therefore, the performance of RE-trees improves as ρ increases.

Figure 14 illustrates the performance results as D is varied. We note that the performance gain of RE-trees generally increases with a larger value of D , indicating that the effective pruning of search space becomes even more crucial for search performance as the data size increases.

Our experimental results indicate that the performance of both the sequential file approach and RE-tree approach are dominated by CPU time, with at least 75% of the total evaluation time being consumed for CPU processing. Thus, these results demonstrate that the RE-tree approach is very effective in pruning the search space, thereby leading to a significant

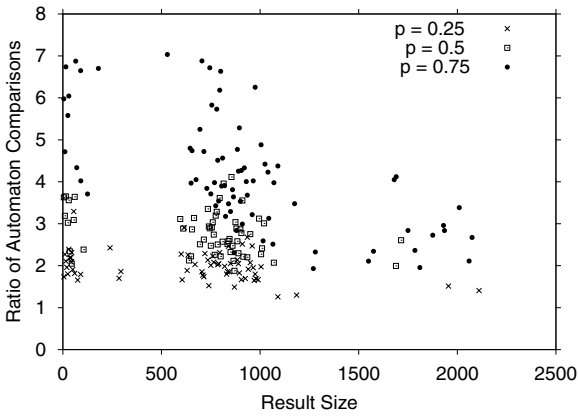
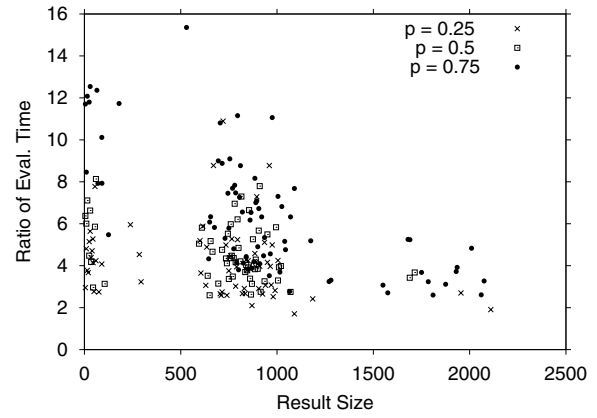
**a****b**

Fig. 13a,b. Varying ρ , $D = 50000$, $\alpha = 10$. **a** Ratio of automaton comparisons, $\frac{N_{seq}(n)}{N_{rt}(n)}$. **b** Ratio of evaluation time, $\frac{T_{seq}(n)}{T_{rt}(n)}$

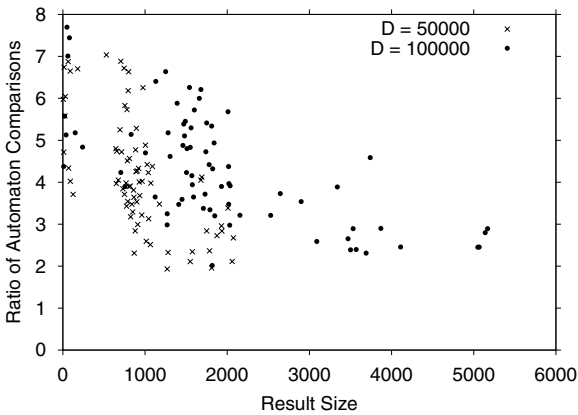
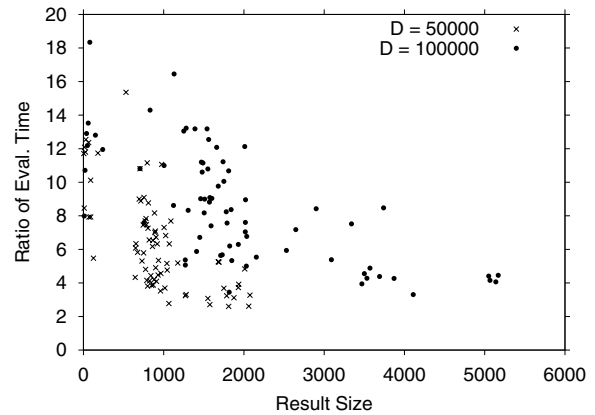
**a****b**

Fig. 14a,b. Varying D , $\alpha = 10$, $\rho = 0.75$. **a** Ratio of automaton comparisons, $\frac{N_{seq}(n)}{N_{rt}(n)}$. **b** Ratio of evaluation time, $\frac{T_{seq}(n)}{T_{rt}(n)}$

reduction in the number of automata comparisons compared to the sequential approach.

8 Related work

Although several indexing methods have been proposed to speed up the search of textual data with REs (e.g., bit-parallel implementation of NFA [23] and suffix trees [2]), there is very little work on the reverse indexing problem involving the retrieval of REs that match an input string. Besides the recent main-memory indexes proposed for filtering XML documents using XPath expressions [1,5,11,19], which is a specialized class of REs, we are not aware of any work on (disk-based) indexes for general REs.

Recently, several indexing methods [1,5,11,19] have been proposed for filtering incoming XML documents against a collection of user profiles expressed as XPath queries, where each XPath query specifies patterns pertaining to paths in the document graph. Our work in this paper (which is a full version of [6]) differs from these indexing schemes in two important aspects. First, the class of REs supported by XPath queries is more specialized as both the union and kleene operators are to be used only with the entire alphabet (i.e., as Σ and Σ^* ,

respectively). Consequently, the associated finite automata are simpler as there are no cycles (besides self-loops) and each state has at most two outgoing transitions: a self-loop transition to itself and a transition to one other state (labeled with either a single letter of the alphabet or the entire alphabet). This is in contrast to our work that deals with REs in their full generality. Second, as opposed to our hierarchical disk-based index approach, the existing XPath-based indexing schemes are designed for main-memory resident data.

A related problem is that of indexing sets of objects, which has been addressed in various contexts, including document retrieval (where each document is characterized by a set of keywords) [3] and evaluation of set predicates (involving set-valued attributes) [16]. Inverted lists and signature files are two well-known techniques that have been developed for this problem [3]. However, since these methods are targeted for indexing finite sets of objects, they are not appropriate for indexing infinite regular languages.

9 Conclusions

In this paper, we presented the RE-tree, which is a novel index structure for performing fast retrievals of REs that match

a given input string. In order to overcome the challenge of indexing the infinite regular sets (corresponding to REs) with no well-defined notion of spatial locality, we developed novel measures for comparing the relative sizes of infinite regular languages; these range from “rate-of-growth” estimates to encoding costs of sample strings using the corresponding REs. We also proposed innovative solutions for the various RE-tree operations, including the effective splitting of RE-tree nodes and computing a “tight” bounding RE (under a space constraint) for a collection of REs. Finally, we showed how sampling-based approximation algorithms can be used to significantly speed up the performance of RE-tree operations.

Our experimental results with synthetic data sets clearly demonstrate that the RE-tree index is significantly more effective than performing a sequential search for matching REs and in a number of cases outperforms sequential search by up to an order of magnitude. An interesting direction for future work is the application of the counting, sampling, and size estimation techniques for regular languages that we developed in the paper to other problem domains like selectivity estimation of REs, clustering REs, etc.

References

1. Altinel M, Franklin MJ (2000) Efficient filtering of XML documents for selective dissemination of information. In: Proceedings of VLDB, Cairo, September 2000, pp 53–64
2. Baeza-Yates R, Gonnet G (1996) Fast text searching for regular expressions or automaton searching on a trie. *J ACM* 43(6):915–936
3. Baeza-Yates R, Ribeiro-Neto B, Navarro G (1999) Indexing and searching. In: Modern information retrieval. ACM Press, New York, pp 191–228
4. Beckmann N, Kriegel H-P, Schneider R, Seeger B (1990) The R*-Tree: an efficient and robust access method for points and rectangles. In: Proceedings of SIGMOD, Atlantic City, NJ, May 1990, pp 322–331
5. Chan C-Y, Felber P, Garofalakis M, Rastogi R (2002) Efficient filtering of XML documents with XPath expressions. In: Proceedings of the 18th international conference on data engineering, San Jose, February 2002, pp 235–244
6. Chan CY, Garofalakis M, Rastogi R (2002) RE-Tree: an efficient index structure for regular expressions. In: Proceedings of VLDB, Hong Kong, August 2002, pp 251–262
7. Clark J, DeRose S (1999) XML path language (XPath) v. 1.0. W3C Recommendation, <http://www.w3.org/TR/xpath>
8. Cochran WG (1977) Sampling techniques. Wiley, New York
9. Intel Corporation (2000) Intel NetStructure XML accelerators, <http://www.intel.com/netstructure/products/xml.accelerators.htm>
10. Deutsch A, Fernandez M, Suci D (1999) Storing semistructured data with STORED. In: Proceedings of the ACM SIGMOD conference on management of data, Philadelphia, June 1999, pp 431–442
11. Diao Y, Fischer P, Franklin MJ, To R (2002) YFilter: efficient and scalable filtering of XML documents. In: Proceedings of the 18th international conference on data engineering, San Jose, February 2002, pp 341–342
12. Garey MR, Johnson DS (1979) Computers and intractability: a guide to the theory of NP-Completeness. Freeman, San Francisco
13. Garofalakis M, Gionis A, Rastogi R, Seshadri S, Shim K (2000) XTRACT: a system for extracting document type descriptors from XML documents. In: Proceedings of SIGMOD, Dallas, May 2000, pp 165–176
14. Guttman A (1984) R-Trees: a dynamic index structure for spatial searching. In: Proceedings of SIGMOD, Boston, June 1984, pp 47–57
15. Hopcroft JE, Ullman JD (1979) Introduction to automata theory, languages, and computation. Addison-Wesley, Reading, MA
16. Ishikawa Y, Kitagawa H, Ohbo N (1993) Evaluation of signature files as set access facilities in OODBs. In: Proceedings of SIGMOD, Washington, DC, May 1993, pp 247–256
17. Kannan S, Sweedyk Z, Mahaney S (1995) Counting and random generation of strings in regular languages. In: Proceedings of SODA, San Francisco, January 1995, pp 551–557
18. Kilpelainen P, Mannila H, Ukkonen E (1995) MDL learning of unions of simple pattern languages from positive examples. In: Proceedings of EuroCOLT, Barcelona, March 1995, pp 252–260
19. Lakshmanan LVS, Parthasarathy S (2002) On efficient matching of streaming XML documents and queries. In: Proceedings of EDBT, Prague, March 2002, pp 142–160
20. Quinlan JR, Rivest RL (1989) Inferring decision trees using the minimum description length principle. *Information and Computation* 80:227–248
21. Rissanen J (1978) Modeling by shortest data description. *Automatica* 14:465–471
22. Stewart JW (1998) BGP4 inter-domain routing in the Internet. Addison-Wesley, Reading, MA
23. Wu S, Manber U (1992) Fast text searching allowing errors. *Commun ACM* 35(10):83–91

A Proofs of theorems

Proof of Theorem 5.1: We show that the set-partitioning problem is NP-hard by reducing the clique problem associated with it. The clique problem is as follows: given an undirected graph $G = (V, E)$ and a constant K , does G contain a clique of size K or more? Let $|V| = n$, $deg(v_i)$ be the degree of vertex $v_i \in V$ and D the maximum degree of any vertex in V . We construct the sets S_1, \dots, S_{n+1} from G as follows. There is a distinct set S_i corresponding to each vertex v_i in V . In addition, there is a distinct element e for each undirected edge (v_i, v_j) in E . Set S_i corresponding to vertex v_i contains all elements for edges (v_i, v_j) incident on vertex v_i . In addition, S_i contains $D - deg(v_i)$ new elements. Finally, S_{n+1} contains all of the elements in sets S_1, \dots, S_n corresponding to vertices in V . Note that the new elements added to each S_i (that do not correspond to edges in E) only appear in S_i and S_{n+1} . Also, the total number of elements is $|S_{n+1}|$.

We show that G contains a clique of at least of size K if and only if there exists a partitioning into S_1 and S_2 , each set containing at least K elements (thus $m = K$) and such that $|\cup_{S_i \in S_1} S_i| + |\cup_{S_i \in S_2} S_i|$ is less than or equal to $|S_{n+1}| + DK - \binom{K}{2}$. Observe that in the optimal partitioning of sets S_1, \dots, S_{n+1} , S_{n+1} is contained in one of the sets, say, S_1 , and the other set S_2 contains K sets from S_1, \dots, S_n (since S_{n+1} contains all the elements). Further, the cost of the partitioning $|\cup_{S_i \in S_1} S_i| + |\cup_{S_i \in S_2} S_i|$ is equal to $|S_{n+1}| + |\cup_{S_i \in S_2} S_i|$.

Suppose G contains a clique of size at least K . We construct the sets S_1 and S_2 as follows. For every vertex v_i in the clique, S_2 contains the set S_i . S_1 contains the remaining sets not contained in S_2 (including S_{n+1}). Thus,

$|\cup_{S_i \in \mathcal{S}_1} S_i| = |S_{n+1}|$. Further, each set S_i in \mathcal{S}_2 contains D elements. However, $K - 1$ of these are due to edges between v_i and other vertices in the clique. Since there are $\binom{K}{2}$ of these edges between vertices of the clique and the $\binom{K}{2}$ elements for these edges occur in two sets in \mathcal{S}_2 , it follows that the total number of elements in $\cup_{S_i \in \mathcal{S}_2} S_i$ is $DK - \binom{K}{2}$.

On the other hand, suppose that there is a partitioning of sets into \mathcal{S}_1 and \mathcal{S}_2 whose cost is less than or equal to $|S_{n+1}| + DK - \binom{K}{2}$. Without loss of generality, let \mathcal{S}_1 contain the set S_{n+1} and \mathcal{S}_2 contain K other sets from S_1, \dots, S_n . Thus, \mathcal{S}_2 must contain less than or equal to $DK - \binom{K}{2}$ distinct elements. Of the total of DK elements in the K sets in \mathcal{S}_2 , only elements corresponding to edges connecting vertices for sets in \mathcal{S}_2 are counted twice. Thus, if \mathcal{S}_2 contains less than or equal to $DK - \binom{K}{2}$ distinct elements, then there must be at least $\binom{K}{2}$ edges between the K vertices for \mathcal{S}_2 . This is possible only if the K vertices form a clique in which an edge exists between every pair of the K vertices. \square

Proof of Theorem 5.2: We reduce the partition problem to the problem of computing the desired DFA M . In the partition problem, there is a set of elements S with each element $e \in S$ associated with a weight $w(e)$. The idea is to partition S into two subsets S_1 and S_2 such that the sum of weights of elements in each subset is equal.

Given an instance of the partition problem, we map it to the automaton computation problem as follows. For each $e_i \in S$ we add to \mathcal{M} the automaton that accepts the language for regular expression $R_i = (e_{i1} | \dots | e_{iw(e_i)})^*$ (here each symbol e_{ij} is distinct). The DFA for the union of the set of regular expressions in \mathcal{M} essentially has a start state s_0 and a separate state s_i for each automaton M_i . Further, there are transitions from a to s_i for every symbol e_{ij} ($j \leq w(e_i)$). Also, there are transitions from s_i to itself for every symbol e_{ij} ($j \leq w(e_i)$). Finally, every state in the union automaton is an accept state. Suppose $B = \sum_{e_i \in S} w(e_i)/2$. We show that there exists an automaton M more general than the union automaton with at most 2 states and such that $|L_l(M)| \leq (l+1)B^l$ if and only if there exists a partitioning of S .

Suppose there exists a partitioning of S into sets S_1 and S_2 . Then, consider the DFA (derived from the union automaton) in which all the states s_i corresponding to $e_i \in S_1$ are merged with the start state s_0 and all the remaining states are merged among themselves. One can easily show that $|L_l(M)| \leq (l+1)B^l$ for this automaton with two states (since each state has exactly B self-transitions). On the other hand, suppose that there exists an automaton with two states such that $|L_l(M)| \leq (l+1)B^l$. Then we show that for all the states s_i corresponding to elements e_i that are merged with start state s_0 , the sum of $w(e_i) = B$. Suppose that this is not the case. Then, of the two states, one must have transitions to itself for at least $B+1$ symbols. Thus, $|L_l(M)| \geq (B+1)^l$, which contradicts the fact that $|L_l(M)| \leq (l+1)B^l$ for all l . Consequently, for all the states s_i corresponding to elements e_i that are merged with start state s_0 , the sum of $w(e_i) = B$. \square

Proof of Theorem 6.1: Let $N_1 = |S_1|$ and $n_1 = |r_1|$. Each random sample r_1 of S_1 with size n_1 is selected with a uniform probability of $1/\binom{N_1}{n_1}$. Thus, the expected value of the estimator is given by $(N_1/n_1) \times (1/\binom{N_1}{n_1}) \times \sum_{r_1} |(r_1 \cap S_2)|$.

In order to estimate $\sum_{r_1} |(r_1 \cap S_2)|$, we note that each element of $S_1 \cap S_2$ can occur in $\binom{N_1-1}{n_1-1}$ samples of size n_1 . Thus, $\sum_{r_1} |(r_1 \cap S_2)| = \binom{N_1-1}{n_1-1} \times |S_1 \cap S_2|$. Thus, the expected value of the estimator is $|S_1 \cap S_2|$. \square

Proof of Theorem 6.2: Let $N_1 = |S_1|$, $k = |r_1 \cap S_2|$ and $n_1 = |r_1|$, and $N = |S_1 \cap S_2|$. We need to show that the probability of a set of size k being chosen is $1/\binom{N}{k}$. This is essentially the probability of a random sample r_1 (of size n_1) containing the specific k elements of $S_1 \cap S_2$. The number of subsets of S_1 of size n_1 containing the specific k elements is $\binom{N_1-N}{n_1-k}$. The total number of distinct subsets of S_1 of size n_1 containing any k elements from $S_1 \cap S_2$ is $\binom{N}{k} \times \binom{N_1-N}{n_1-k}$. Thus, the probability of $r_1 \cap S_2$ containing the specific k elements is equal to $\binom{N_1-N}{n_1-k} / (\binom{N}{k} \times \binom{N_1-N}{n_1-k})$, which is equal to $1/\binom{N}{k}$. \square

Proof of Theorem 6.3: Let $N_1 = |S_1|$, $N_2 = |S_2|$ and $n = |r|$. Consider any set r returned by algorithm SAMPLEUNION. We show that this is a uniform random sample of size n by showing that the probability of selecting r is $1/\binom{N_1+N_2}{n}$. Consider an arbitrary permutation of elements in r . We show that the i^{th} element of the permutation is chosen with probability $1/(N_1 + N_2 - i + 1)$. Suppose that, of the previous $i-1$ elements, i_1 are chosen from r_1 and i_2 are chosen from r_2 (note that $i_1 + i_2 = i-1$). Further, without loss of generality, the i^{th} element belongs to r_1 . Then the probability of choosing the i^{th} element is $(N_1 - i_1)/(N_1 + N_2 - i_1 - i_2) \times (|r_1| - i_1)/(N_1 - i_1) \times 1/(|r_1| - i_1)$, which is equal to $1/(N_1 + N_2 - i + 1)$. Of the three terms in the probability computation, the first is the probability of selecting an element from S_1 , the second is the probability that the i^{th} element belongs to r_1 after the previous i_1 elements have been deleted from it, and the final term is the probability that the i^{th} element is chosen from r_1 (without the previous i_1 elements). We elaborate further on the second term – once i_1 elements are deleted from r_1 , an argument similar to that used in Theorem 6.2 can be used to show that r_1 without the i_1 elements is a uniform random sample of S_1 without the i_1 elements.

Thus, from the above discussion it follows that the probability of choosing the i^{th} element in any given permutation of r is $1/(N_1 + N_2 - i + 1)$. Thus, the probability of choosing all the elements in any permutation of r is $(N_1 + N_2 - n)! / (N_1 + N_2)!$. Since there are $n!$ different permutations of r , the probability that r is selected is $((N_1 + N_2 - n)! \times n!) / (N_1 + N_2)!$, which is essentially $1/\binom{N_1+N_2}{n}$. This proves the theorem. \square

B Generating random strings for NFAs

In this section, we explain how algorithm GENRANDOMSTRING can be modified to generate random strings for NFAs.

Essentially, for an NFA, since there can be multiple transitions for the same symbol x out of a state s of M , choosing the (symbol, state) pair for the next transition needs to be carried out in two phases (step 2 of algorithm GENRANDOMSTRING). Suppose that for each value of i from 1 to n , for each state s and each symbol x , during the execution of steps 7–11 of algorithm COUNTSTRINGS, we store in variable $\gamma_{xt_j}(s, i)$ the contribution of symbol x and state t_j ($1 \leq j \leq l$) to

$\beta(s, i)$. Thus, $\gamma_{xt_1}(s, i) = \beta(t_1, i - 1)$ (step 8) and for $j \geq 2$, $\gamma_{xt_j}(s, i) = \beta(t_j, i - 1) * \frac{|r - (\cup_{k=1}^{j-1} L_{i-1}(M, t_k))|}{|r|}$ (step 11). Essentially, $\gamma_{xt_j}(s, i)$ stores the number of distinct length- i accepting paths from s whose first transition is on x and that pass through t_j but have not been counted earlier in $\gamma_{xt_1}(s, i), \dots, \gamma_{xt_{j-1}}(s, i)$. Also, let $\gamma_x(s, i) = \sum_{t_j \in \delta(s, x)} \gamma_{xt_j}(s, i)$.

We now describe the two phases when choosing a (symbol, state) pair for the next transition when generating a length- i random string from state s . In the first phase, symbol x for the transition is chosen with probability $\frac{\gamma_x(s, i)}{\beta(s, i)}$ (this is

approximately the fraction of distinct length- i accepting paths from s beginning with symbol x). Once a symbol x has been chosen, in the second phase, if $t_j \in \delta(s, x)$, then state t_j for the transition on x is chosen with probability $\frac{\gamma_{xt_j}(s, i)}{\gamma_x(s, i)}$ (this is approximately equal to the fraction of distinct length- i accepting paths from s whose first transition is on x and that pass through t_j). Note that during the i^{th} iteration of the “for” loop in step 3 of algorithm COUNTSTRINGS, random strings of length $i - 1$ need to be generated in step 10. This only requires γ for values less than i for each state, which should have been computed in previous iterations of the algorithm.