

Indexing Textual Information

Ioannis N. Kouris

University of Patras, Greece

Christos Makris

University of Patras, Greece

Evangelos Theodoridis

University of Patras, Greece

Athanasios Tsakalidis

University of Patras, Greece

INTRODUCTION

Information retrieval is the computational discipline that deals with the efficient representation, organization, and access to information objects that represent natural language texts (Baeza-Yates, & Ribeiro-Neto, 1999; Salton & McGill, 1983; Witten, Mofiat, & Bell, 1999). A crucial subproblem in the information retrieval area is the design and implementation of efficient data structures and algorithms for indexing and searching information objects that are vaguely described. In this article, we are going to present the latest developments in the indexing area by giving special emphasis to: data structures and algorithmic techniques for string manipulation, space efficient implementations, and compression techniques for efficient storage of information objects.

The aforementioned problems appear in a series of applications as digital libraries, molecular sequence databases (DNA sequences, protein databases [Gusfield, 1997]), implementation of Web search engines, web mining and information filtering.

BACKGROUND

Dictionary Data Structures

The dictionary data structure stores a set S of n elements in order to support the operations of insertion, deletion, and the test of membership. A basic criterion for categorizing dictionary data structures is whether only comparisons are used, or the representation of elements for guiding the search is also employed. Typical representatives of the former group are *search trees* and of the latter *tries* and *hashing*. Search trees need $O(\log n)$ update/search time and $O(n)$ space and the most prominent examples of them are: AVL-trees, red-black trees, (α, b) -trees, $BB[\alpha]$ -trees and Weight Balanced B-trees (Arge, & Vitter, 1996; Cormen, Leiserson,

& Rivest, 1990; Mehlhorn, 1984). On the other hand, *tries* and *hashing* structures (Cormen, Leiserson, & Rivest, 1990; Czegh, Havas, & Majewski, 1997; Pagh, 2002) try to use the representation (for example, the value of the element written as a string of digits or the value itself), to compute directly the element's position in system's memory. The time and space complexities of these structures generally vary; however, it should be mentioned that a lately developed structure (Anderson & Thorup, 2001) answers both search and update operations in $O(\sqrt{\log n / \log \log n})$ time. This structure is also able to retrieve the largest element in the stored set smaller than a *query* element (*predecessor query*).

Finding Occurrences of Patterns

The string matching (or pattern matching) problem is one of the most frequently encountered and studied problems in the area of system/algorithm design. In this problem, we are searching for the occurrences of a pattern P in a sequence of symbols T . The naive $O(|P||T|)$ algorithm aligns the pattern at each one of the $O(|T|)$ possible positions of the sequence and executes $O(|P|)$ comparisons; however, there exist elegant, though complex, algorithms whose overall time complexity is linear, that is, $O(|P| + |T|)$. The most known linear time algorithms that achieve that are Knuth-Morris-Pratt (Knuth, Morris, & Pratt, 1977) and variants of the Boyer-Moore (Boyer & Moore, 1977) algorithm. For the case that we are searching for a set of patterns in the sequence, the Aho-Corasick automaton (Aho & Corasick, 1975) can be used. This automaton accepts all the patterns of the set and can be constructed in time linear to the sum of the lengths of them. Running the automaton with the characters of T , all the occurrences of the patterns are reported in $O(|T|)$ time.

On most of the modern applications, the patterns arrive in an on line manner and the $O(|T| + |P|)$ computational

time is prohibitive; there is need for indexing structures (indices) that can perform the queries as close as possible to $O(|P|)$ computational time, assuming that the text has been preprocessed *once*. The indices that try to satisfy this demand are divided in two categories: the *word-based* (or *keyword-based*) indices, which have been designed for sequences of symbols that can be divided in tokens/words, and the *full-text* (or *sequential scan*) indices, where the previous feature does not hold and the strings involved are non-tokenizable.

TEXT AND STRING DATA STRUCTURES

Word Based Indices

The most commonly used indexing structures in this category are *inverted files*, *signature files* and *bitmaps*. An inverted file consists of two parts: a structure for storing the set of all different words in the text and, for each such word, a list of the text positions where the word appearances are stored. Signature files are term-oriented structured based on hashing while bitmaps represent each document as a bit vector having length equal to the size of the lexicon. In typical applications compressed inverted files are considered to be superior to both signature files and bitmaps (Faloutsos, 1985; Zobel, Moffat, & Ramamohanarao, 1998).

More analytically, consider a document collection and a lexicon containing the terms that appear in the documents of the collection. An inverted file consists of a search structure containing all the distinct terms that appear in the lexicon and, for each distinct term, a list termed inverted (or postings) list storing the identifiers (usually integer numbers starting from 1) of the documents containing the term. The search structure can be implemented as a search tree or a hashing structure, and queries are evaluated by using the search structure to find the relevant terms, fetching the inverted lists for the corresponding terms, and then intersecting them for conjunctive queries or merging them for disjunctive queries. There are many algorithms for constructing inverted files while their performance can be significantly improved by employing compression techniques for representing the postings lists. The interested reader should find relevant material in Witten, Moffat, and Bell (1999) and Zobel, Moffat, and Ramamohanarao (1998).

On the other hand, in the signature file method, all the documents in the collection are stored sequentially; each document is hashed (mapped to) a distinct signature and all the signatures are stored in a signature file.

Finally, in bitmaps for every term in the lexicon a bit-vector is stored, each bit corresponding to a document, a bit

with value 1 means that the term appears somewhere in the document; otherwise it has value 0.

Full Text Indices

In general, these indices are more powerful from word-based indices since they answer a wider range of queries like arbitrary substring queries, motifs, and repetitions selection queries, and they have the ability to index non-tokenizable strings like biological sequences. The price for these extra capabilities is larger space consumption. The common feature of full-text indices is that they organize in a proper order (frequently lexicographic), all the suffixes of the sequence.

Static Indices: For the setting that texts do not underlie on updates (insertions/deletions of characters), there are many classical confrontations coming quite a few years ago. The most famous data structure of this type is the *suffix tree*. A suffix tree of a string T is a compacted trie of all suffixes of T . The suffix tree has linear, to the length $|T|$ of the string, construction time. For bounded alphabets, there exist the algorithms of Weiner (1973), McCreight (1976) and Ukkonen (1995) and, for arbitrary large alphabets, the algorithm of Farach (1997). The main counterparts of suffix trees are *suffix arrays*. A suffix array of a string T is an array that indicates the lexicographic order of all suffixes of T . Suffix arrays were proposed in Manber, and Myers (1993), where they described an $O(|T|\log|T|)$ construction time algorithm and how pattern matching queries can be performed in $O(|P|+\log|T|)$ time, assuming that the longest common prefixes among each pair of adjacent suffixes in the array have been precomputed. Later on, there were proposed three linear time construction algorithms in Karkkainen and Sanders (2002); Ko, and Aluru (2003); Kim, Sim, Park, and Park (2003). Suffix arrays, though having a slowdown term, are usually preferable from suffix trees due to their smaller space requirements. The interested reader should refer to Grossi, and Italiano (1993) and Gusfield (1997) to familiarize with the numerous applications of the aforementioned structures.

Dynamic Indices: In many applications, for example, text editors, log files and so forth, the sequences that underlie pattern matching queries change dynamically by insertions or deletions of symbols in arbitrary places. For this setting, we would like to answer the queries as close as possible to the desirable $O(|P|)$ time without preprocessing the whole sequence from the start and paying $O(|T|)$ computational time. The static structures from the previous subsection cannot achieve that since even a small change can modify all the suffixes and enforce a full reorganization of them. In Gu, Farach, and Pagli (1994), Ferragina (1994), Ferragina, and Grossi (1995a), Sahinalp, and Vishkin (1996), Ferragina and Grossi (1995b), Alstrup, Brodal, and Rauhe (2000), the researchers developed various indices, trying to keep the

4 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/indexing-textual-information/13840

Related Content

Understanding User Social Commerce Usage Intention: A Stimulus-Organism-Response Perspective

Tao Zhou (2019). *Information Resources Management Journal* (pp. 56-71).

www.irma-international.org/article/understanding-user-social-commerce-usage-intention/234443

How Do Actuaries Use Data Containing Errors?: Models of Error Detection and Error Correction

Barbara D. Klein (1997). *Information Resources Management Journal* (pp. 27-36).

www.irma-international.org/article/actuaries-use-data-containing-errors/51041

Modeling Customer-Related IT Diffusion

Shana L. Dardan, Ram L. Kumar and Antonis C. Stylianou (2009). *Best Practices and Conceptual Innovations in Information Resources Management: Utilizing Technologies to Enable Global Progressions* (pp. 251-263).

www.irma-international.org/chapter/modeling-customer-related-diffusion/5521

Digital Government and Individual Privacy

Patrick R. Mullen (2005). *Encyclopedia of Information Science and Technology, First Edition* (pp. 870-874).

www.irma-international.org/chapter/digital-government-individual-privacy/14351

Nazar Foods Company: Business Process Redesign Under Supply Chain Management Context

Vichuda Nui Polatoglu (2006). *Journal of Cases on Information Technology* (pp. 49-62).

www.irma-international.org/article/nazar-foods-company/3170