

Indexing and Compressing Text

Ioannis 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 (Salton & McGill, 1983; Witten, Moffat, & Bell, 1999; Baeza-Yates & Ribeiro-Neto, 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

DOI: 10.4018/978-1-4666-5888-2.ch173

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 (Mehlhorn, 1984; Cormen, Leiserson, & Rivest, 1990; Arge & Vitter, 1996). On the other hand, *tries* and *hashing* structures (Cormen, Leiserson, & Rivest, 1990; Pagh, 2002; Czegh, Havas, & Majewski, 1997) 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 (Andersson & 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, i.e. $O(|P| + |T|)$. The most known linear time algorithms that achieve that are Knuth, Morris, and 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 closer 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, & Bell, 1999; Zobel, Moffat, & Ramamohanarao, 1998).

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 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

7 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-and-compressing-text/112585

Related Content

Hybrid TRS-FA Clustering Approach for Web2.0 Social Tagging System

Hannah Inbarani H and Selva Kumar S (2015). *International Journal of Rough Sets and Data Analysis* (pp. 70-87).

www.irma-international.org/article/hybrid-trs-fa-clustering-approach-for-web20-social-tagging-system/122780

A Rough Set Theory Approach for Rule Generation and Validation Using RSES

Hemant Rana and Manohar Lal (2016). *International Journal of Rough Sets and Data Analysis* (pp. 55-70).

www.irma-international.org/article/a-rough-set-theory-approach-for-rule-generation-and-validation-using-rses/144706

Identification of Heart Valve Disease using Bijective Soft Sets Theory

S. Udhaya Kumar, H. Hannah Inbarani, Ahmad Taher Azar and Aboul Ella Hassanien (2014). *International Journal of Rough Sets and Data Analysis* (pp. 1-14).

www.irma-international.org/article/identification-of-heart-valve-disease-using-bijective-soft-sets-theory/116043

Illness Narrative Complexity in Right and Left-Hemisphere Lesions

Umberto Giani, Carmine Garzillo, Brankica Pavic and Maria Piscitelli (2016). *International Journal of Rough Sets and Data Analysis* (pp. 36-54).

www.irma-international.org/article/illness-narrative-complexity-in-right-and-left-hemisphere-lesions/144705

The Past, Present, and Future of UML

Rebecca Platt and Nik Thompson (2018). *Encyclopedia of Information Science and Technology, Fourth Edition* (pp. 7481-7487).

www.irma-international.org/chapter/the-past-present-and-future-of-uml/184445