

A Fast and Space-Economical Algorithm for the Tree Inclusion Problem

Yangjun Chen

University of Winnipeg, Canada

Yibin Chen

University of Winnipeg, Canada

INTRODUCTION

Let T be a rooted tree. We say that T is *ordered* and *labeled* if each node is assigned a symbol from an alphabet Σ and a left-to-right order among siblings in T is specified. Let v be a node different of the root in T with parent node u . Denote by $delete(T, v)$ the tree obtained by removing the node v from T , by which the children of v become part of the children of u as illustrated in Figure 1.

Given two ordered labeled trees P and T , called the pattern and the target, respectively. We may ask: Can we obtain pattern P by deleting some nodes from target T ? That is, is there a sequence v_1, \dots, v_k of nodes such that for

$$T_0 = T \text{ and}$$

$$T_{i+1} = delete(T_i, v_{i+1}) \text{ for } i = 0, \dots, k-1,$$

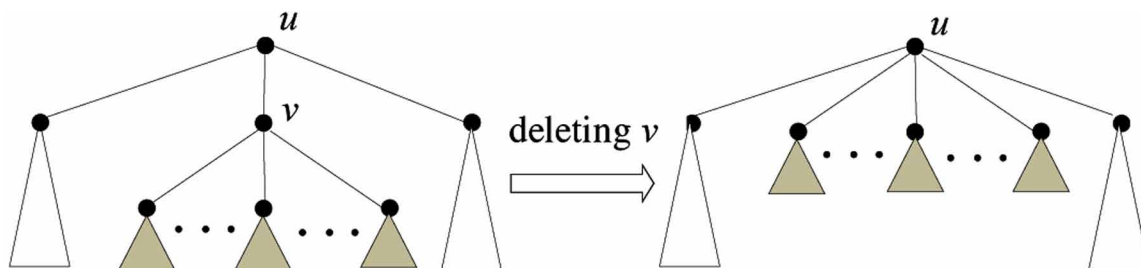
we have $T_k = P$? If this is the case, we say, P is included in T (Kilpeläinen and Mannila, 1995).

Such a problem is called the *tree inclusion problem*. It has many applications in the computer engineering as described below.

BACKGROUND

The first interesting application of the tree inclusion is used as an important query primitive for XML data (Mannila and R  iha, 1990), where a structured document database is considered as a collection of parse trees that represent the structure of the stored texts and the tree inclusion is used as a means of retrieving information from them. As an example, consider the tree shown in Figure 2, representing an XML document for the book *Arts of Programming* authored by (Knuth, 1969). One might want to find this book in an XML database by forming a pattern tree as shown in Figure 3 as a query, which can be obtained by deleting some nodes from the tree shown in Figure 2. Thus, a tree inclusion checking needs to be conducted to evaluate this query.

Figure 1. Illustration of node deletion



DOI: 10.4018/978-1-5225-2255-3.ch391

Figure 2. A XML document (target) tree

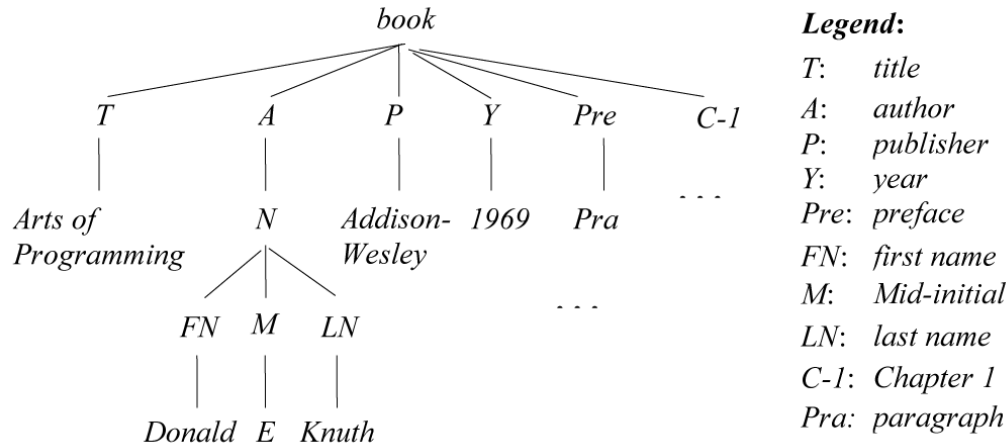
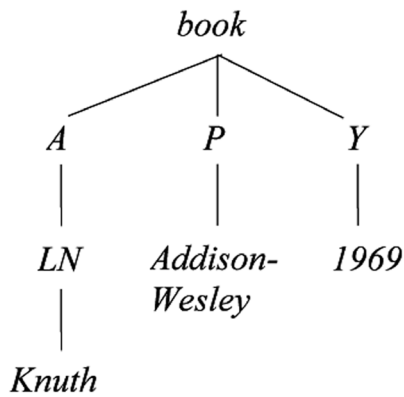


Figure 3. A pattern tree



Another application of this problem is to query the grammatical structures of *English* sentences, which can also be represented as an ordered tree since a sentence can always be divided into several ordered components such as noun phrases, verb phrases, and adverbs; and a noun phrase itself normally contains an article and a noun while a verb phrase may contain a verb, a noun phrase, an adverb, and so on. To check whether a concrete sentence is grammatically correct, we will represent it as a pattern tree and make a tree inclusion checking against some target grammatical tree structures.

A third application of the ordered tree inclusion is the video content-based retrieval. According to (Rui and Huang, 1999), a video can be successfully

decomposed into a hierarchical tree structure, in which each node represents a scene, a group, a shot, a frame, a feature, and so on. Especially, such a tree is an ordered one since the temporal order is very important for video. Some other areas, in which the ordered tree inclusion finds its applications, are the scene analysis, the computational biology, such as *RNA* structure matching (Lyngs, Zuker and Pedersen, 1999), and the data mining, such as tree mining discussed in (Zaki, 2002), just to name a few.

Up to now, the best algorithm for this problem requires $O(|T| + |P|)$ space and $\Theta(|T| \cdot |\text{leaves}(P)|)$ time (Alonso and Schott, 1993; Chen, 1998; Chen and Chen, 2006; Bille and Gørtz, 2011), where $\text{leaves}(P)$ stands for the set of the leaves of P .

In this chapter, we propose an efficient algorithm for this problem. Its time and space complexities are bounded by $O(|T| \cdot \min\{h_p, |\text{leaves}(P)|\})$, and $O(|T| + |P|)$, respectively, where h_p is the height of P , defined to be the number of edges on the longest downward path from the root to a leaf node.

BASIC DEFINITIONS

In this section, we mainly define the notations that will be used throughout the paper. Let T be a

11 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/a-fast-and-space-economical-algorithm-for-the-tree-inclusion-problem/184158

Related Content

An Overview for Non-Negative Matrix Factorization

Yu-Jin Zhang (2015). *Encyclopedia of Information Science and Technology, Third Edition* (pp. 1631-1641).
www.irma-international.org/chapter/an-overview-for-non-negative-matrix-factorization/112568

Should Innovation Knowledge be Assessed?

Fawzy Soliman (2015). *Encyclopedia of Information Science and Technology, Third Edition* (pp. 4699-4708).
www.irma-international.org/chapter/should-innovation-knowledge-be-assessed/112912

Detection of Automobile Insurance Fraud Using Feature Selection and Data Mining Techniques

Sharmila Subudhi and Suvasini Panigrahi (2018). *International Journal of Rough Sets and Data Analysis* (pp. 1-20).
www.irma-international.org/article/detection-of-automobile-insurance-fraud-using-feature-selection-and-data-mining-techniques/206874

Lean Logistics of the Transportation of Fresh Fruit Bunches (FFB) in the Palm Oil Industry

Cheah Cheng Teik and Yudi Fernando (2018). *Encyclopedia of Information Science and Technology, Fourth Edition* (pp. 5422-5432).
www.irma-international.org/chapter/lean-logistics-of-the-transportation-of-fresh-fruit-bunches-ffb-in-the-palm-oil-industry/184245

Mobilization: Decision Theory

Idongesit Williams (2021). *Encyclopedia of Information Science and Technology, Fifth Edition* (pp. 1436-1450).
www.irma-international.org/chapter/mobilization/260278