

# EFFICIENT STRING MATCHING

---

Science fair

Fredrik Koppelow Eliassen

# THE PAPER

---

Programming  
Techniques

Glenn Manacher  
Editor

---

## Efficient String Matching: An Aid to Bibliographic Search

Alfred V. Aho and Margaret J. Corasick  
Bell Laboratories

---

**This paper describes a simple, efficient algorithm to locate all occurrences of any of a finite number of keywords in a string of text. The algorithm consists of constructing a finite state pattern matching machine from the keywords and then using the pattern matching machine to process the text string in a single pass. Construction of the pattern matching machine takes time proportional to the sum of the lengths of the keywords. The number of state transitions made by the pattern matching machine in processing the text string is independent of the number of keywords. The algorithm has been used to improve the speed of a library bibliographic search program by a factor of 5 to 10.**

### 1. Introduction

In many information retrieval applications it is necessary to be able to locate all occurrences of user-specified keywords or phrases in text. This paper describes a simple algorithm to locate all occurrences of a finite number of keywords and phrases in a text string.

The approach should be familiar to those familiar with finite automata. The algorithm is described in two parts. In the first part we construct from a set of keywords a finite state pattern matching machine. In the second part we apply the text string as input to the machine. The machine signals a match for a keyword.

Using finite state machines in information retrieval applications is not new [4, 8, 17], but has been frequently shunned by programmers. The reason for this reluctance on the part of programmers is due to the complexity of programming the algorithms for constructing finite state machines from expressions [3, 10, 15], particularly when many techniques are needed [2, 14]. The efficient finite state pattern matching machine we constructed quickly and simply from regular expressions, namely those that can be expressed as a set of keywords. Our approach compares favorably with the Knuth-Morris-Pratt algorithm [11] for constructing finite state machines.

Perhaps the most interesting feature of the algorithm is the amount of improvement that it gives over more conventional approaches to finite state pattern matching algorithms for bibliographic search programs. The paper

# THE AHO-CORASICK ALGORITHM

---

- Quickly locate occurrences of user-specified keywords from an arbitrary text string.
- 2 parts:
  - Construction of the finite state pattern machine
  - Locate keywords from string input

# THE RELEVANCE OF THE ALGORITHM TO THE COURSE

---

- Search