

文字連接を用いたフルテキスト検索の高速化 7H-5

宮原末治 小橋史彦
(NTT ヒューマンインターフェース研究所)

1. まえがき 多量に発生する情報を、事前加工なしにフルテキストの形で蓄積し、情報が必要になった時点で、自然言語を用いて検索するフルテキスト形データベース検索システムについて検討している¹⁾。この中で、検索は「問い合わせ」と「検索結果からの適合文の選択」とを繰り返す仕事となり、これを効率よく行うために処理の高速化が要求され、ハードウェアで実現する試みがある²⁾。

本報告ではソフトウェアによってテキスト検索を高速に行う方法として、文字連接情報を索引に用いた検索法の検討結果について述べる。

2. 外部条件と処理内容 テキスト検索を行うためのハードウェアにはELISを用いた³⁾。ELISはマイクロコード長が64ビットのマイクロプログラム制御で動作する装置であり、マシンサイクルは120nsecである。

この装置にテキストとして新聞記事を登載し、文を単位とした検索を行うことにした。検索の処理は、フルテキスト検索と同等の処理を実現することを狙いに、下記のような7つの処理ステップで構成した。

- ① 前処理（索引作成とテキストのリスト形式への変換）
 - ② オペレータによる問い合わせ（検索文）の入力
 - ③ 入力文からの単語（キー単語）抽出とシソーラスによる類語展開（これらの単語を総称して検索単語と呼ぶ）
 - ④ 索引を用いた候補文の抽出
 - ⑤ 候補文サーチによる検索単語の種類数と出現数との検出
 - ⑥ 候補文の適合性の計算（キー単語ごとの種類数の総和に検索単語の出現数を加えた値）とその順位付け¹⁾
 - ⑦ 優先付けされた検索結果の表示・出力
- この処理の中で①の索引付けの方法を変更した場合に必要なメモリ容量と④, ⑤, ⑥の処理を合わせた処理速度との関係、およびその時のフルテキストサーチの検索結果との差分を調査比較することとした。

3. 実験と評価 高速化の方法としては、語句が出現する文番号を索引として予め用意しておき、検索の際には検索文の検索単語に出現する語句から候補文を求め、その候補文に対してのみテキストサーチの処理を施し、検索単語と同じ単語が存在するか否かを検査する方法を採用した。

表1. 索引の付与方法とメモリ容量・候補文数の関係

方法	索引の付与方法	平均候補文数	メモリ容量 [*] (Kバイト)	索引の種類数
イ	テキストからの文字索引	62	655	2,084
ロ	テキストからの連接索引	4	672	22,025
ハ	解析単語からの連接索引	5	487	13,953

* テキスト(被検索文)の大きさは約605Kバイト(約5,000文)

* A Study of Hight-speed Full-Text Retrieval using Character Transition Probability.

* Sueharu MIYAHARA, Fumihiko OBASI

* NTT Human Interface Laboratories

3.1 索引の作成： 索引は文字単位にどの文にその文字が出現したかを記憶する方法(イ)、テキスト中の2文字の連接を索引にする方法(ロ)、テキストを言語解析して単語を抽出し、その単語から連接を調べて索引にする方法(ハ)を採用した。なお、平仮名は索引から除外し、その対策として平仮名以外の索引から得られた候補文に対して平仮名も含めてサーチする方法を探った。これら索引付けの方法とその時の索引の状態を表1に示す。

3.2 検索方式と検索結果： 検索方式と検索時間との関係を表2に示す。これは検索文5例の平均を表わしたもので、1つの検索文あたりキー単語は平均5個、検索単語は平均12単語である。方式Aの検索では検索単語内に索引が複数個連続する場合、共通して出現する文番号を候補文と見なしてサーチ処理を行い高速化した。測定の結果、全文をサーチして検索を行った場合に約19.5秒を要したテキストが、言語解析した単語を用いて2文字の連接索引を作成した場合、約7倍の速度で検索でき、上位10個の検索結果に変化がなかった。これは検索文あたり候補文が平均550個と多かったためと考える。なおこのとき、テキスト容量の約80%に相当するメモリが索引用として必要になった。

3.3 考察： ① 索引をより細かな単位にすると、候補文が増加してテキストサーチの時間は増加するが、検索漏れを少なくできる利点がある。しかし、(イ)や(ロ)の方法では索引のデータ量の増加から平仮名を索引に持つことは困難と思われるが、(ハ)の方法は可能と考える。
 ② 上記以外にもテキストを言語解析して単語を索引にする方法が在るが、この方法は検索速度がやや速くなるが、部分一致の文字列を検索でき出来ないと云う問題が残る。

4. むすび フルテキスト検索を高速に行うため、検索対象テキストを言語解析して得られた単語から連接索引を作成して検索する方法を検討した。その結果、テキスト容量の約80%のメモリが索引用として必要になるが、検索精度をほとんど低下させることなく、検索速度を全文サーチに較べ、約7倍に向上できることが分かった。今後は索引作成の高速化やシステム操作や検索・確認の容易化について検討していく。

謝辞 本研究の機会を与えて戴いた当研究所・川嶋言語メディア研究部長、協力戴いた加納英文氏に深謝する。

表2. 検索方式と処理時間の関係

方式	処理(検索と順位付け)	索引	処理時間 [*] (秒)
A	文字の指示する文をサーチ	イ	3.06
B	連接の指示する文をサーチ	ロ	2.94
C	連接の指示する文をサーチ	ハ	2.73
D	全文サーチ	無	19.50

* 検索文5例の平均、 * 被検索文の長さは平均55文字

文献 1) 宮原他 “文書情報蓄積検索…”, 平1情處全大, 2N-6 (1989-10).

2) 加藤他 “大規模文書情報システム用キヤ…”, 情報学基礎, 14-6 (1989-7).

3) 日比野他 “LISPマシンELISの基本…”, 情報処理研究, 12-15 (1980).