

フルテキストサーチシステム Bibliotheca/TS の開発(1)
—システムの概要—

2C-3

浅川悟志 川下靖司 坂田淳 畠山敦
(株) 日立製作所

1.はじめに

近年、ワードプロセッサを始めOA機器の普及拡大に伴い、作成される文書情報も急速に増加してきており、近い将来膨大な量に達するものと予想されている。このため、大量の文書情報を簡単に蓄積検索できる文書情報検索システムに対する要求が高まりつつある。また、既存システムにおいても、文書データベースの大規模化に伴う絞り込み率の低下や、技術用語の目まぐるしい変遷に起因する検索精度の低下などが大きな問題となってきた。

こうした要求や問題に応えるため、インデックスを用いない自由な言葉による検索を目的としたフルテキストサーチ技術の研究を行ってきている。本研究の中で、その実用性を確かめるため、既にテキストサーチマシン(TSM-I)を試作し、発表してきた。^[1]

今回は、そのプロトタイプの技術を活かしソフトウェアのみで動作する、ワークステーション3050ベースのフルテキストサーチシステム"Bibliotheca/TS"を開発したので、報告する。

2.システムの概要2.1構成

Bibliotheca/TSは、TCP/IPプロトコルを用いたクライアント・サーバシステムである。図1に構成を示す。

(1) BIBLIOTHECA/TS-CL(Text Search - Client)

サーバに対して検索要求を出すためのクライアントソフトウェアである。GUIにはMotifを採用している。検索条件入力、結果一覧表示、テキストハイライト表示、履歴保持などの機能を有する。また、APIとして、Cのライブラリ関数群を提供するようにした。

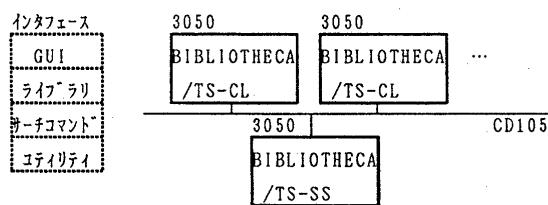


図1 Bibliotheca/TS構成概要

(2) BIBLIOTHECA/TS-SS(TS - Software Server)

クライアントからの検索要求(サーチコマンド)を処理するためのサーバソフトウェアである。同時に複数台のクライアントと通信セッションを確立できるようにした。また、DB構築やメンテナンスを行うためのユーティリティ機能も持たせるようにした。

2.2特長

Bibliotheca/TSは、以下の特長を持つ。

(1)階層プリサーチ方式

階層プリサーチとは、検索タームの存在する可能性がない文書を読み飛ばすことで、フルテキストサーチの検索速度を加速する方式である。スキャンの必要性を判定するために、テキストデータを単語レベルで圧縮した凝縮テキストと、文字レベルで圧縮した文字成分表を用いた。特に今回は、文字成分表として2文字単位で文書中の使用文字を登録する連接文字成分表を採用した。

この階層プリサーチ方式を用いた検索処理の流れは、図2のようになる。まず、与えられた検索タームを2文字単位に分割し、文字成分表を用いて検索ターム中の全ての連接文字を含んでいる文書を候補として抽出する。次に、得られた候補文書について、その凝縮テキストをスキャンし、単語レベルでの照合を行う。そして最後に、必要に応じてテキストデータをスキャンし、最終の検索結果を得る。

この方式により、高速かつ検索もれのないフルテキストサーチを実現することができる。

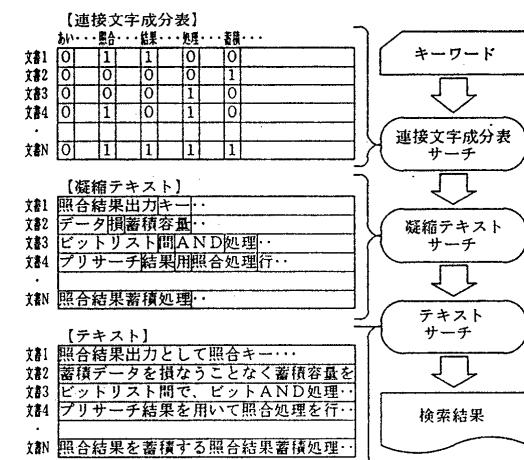


図2 階層プリサーチ方式

(2) サーチエンジンエミュレータ

TSM-Iにおいてハードウェアで実現していた文字列照合処理を、ソフトウェアでエミュレートし、コストの低減を図った。複数個の検索タームを一度のテキストスキャンで探索するアルゴリズムを用いた。

(3) 同義語・異表記展開処理

ユーザの指定した検索タームと同義なタームを、同義語辞書を用いて展開するようにした。また、外来語などのカタカナ表記のゆれは、部分文字列の置換規則により、異表記を生成するようにした。これらを一括して探索することにより、フルテキストサーチの課題である検索もれの解消を図った。

(4) 複合条件判定処理

指定された検索タームが存在する文書を抽出する機能の他に、検索ターム間の論理条件や、テキストデータ内での検索タームの出現位置関係を条件として指定し検索する機能を持たせた。この機能により、きめの細かな検索が行えるため、的確な目的の文書を検索することが可能である。

3. 性能評価

階層プリサーチの傾向として、文字成分表での絞り込み率や、検索条件の特性により処理速度が大きく左右される。例えば、文字成分表サーチの結果、多くの文書がヒットすると、それだけ多くの凝縮テキストやテキストデータをスキャンしなければならない。また、検索タームが单一文字種であれば、凝縮テキストまで検索処理が終了するのに対し、複合文字種の場合にはテキストデータのサーチを行う必要が生じる。

これらの要因から、以下4種類の条件下で検索処理時間を測定した。

(1) 単純条件

漢字あるいはカタカナによる単一文字種で検索タームを指定する場合は、凝縮テキストで検索処理が終了する。

(2) 論理条件

漢字あるいはカタカナによる単一文字種で検索タームを指定し、かつ論理条件を指定する場合は、凝縮テキストで検索処理が終了する。ただし、検索ターム間の論理条件判定を行うため、複合条件判定プログラムにかかる負荷が大きい。

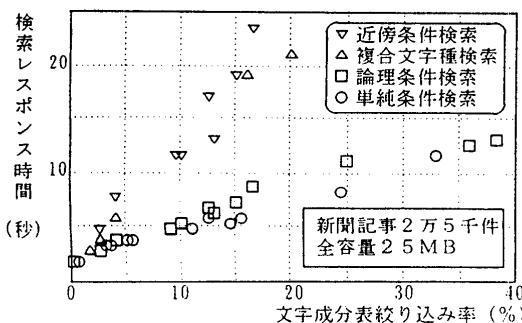


図3 システム検索性能

(3) 複合文字種

漢字とカタカナの組合せなど、異なる文字種での単純条件検索の場合は、条件式は(1)と同じだが、内部的には凝縮テキストを検索タームの部分文字列で論理条件検索し、テキストを単純条件検索することになる。

(4) 近傍条件

検索ターム間のテキスト中の位置関係を指定する場合は、凝縮テキストを検索ターム間の論理条件で、テキストを近傍条件で検索することになる。

以上の条件での測定結果を図3に示す。検索時間は、クライアント側での待ち時間、すなわち検索レスポンス時間を示すものである。データベースには新聞記事2万5千件を用いた。テキストデータの容量は約25MB、凝縮テキストの圧縮率は50%である。図から、検索処理のオーバヘッドは約2秒であることが分かる。これは、通信と条件式の解析およびサーバ内部でのコマンド実行に伴うオーバヘッドである。

単純条件と論理条件は凝縮テキストで処理が終了するため、複合文字種検索や近傍条件と比べ半分以下の時間となっている。複合文字種と近傍条件及び単純条件と論理条件の間で顕著な差がみられるが、これは複合条件判定処理にかかる負荷に比べて文字列探索処理のサーチエンジンエミュレータにかかる負荷の方が、はるかに大きいことを表している。

以上の結果から、文字成分表の絞り込み率が10%程度のときは、単純条件あるいは論理条件検索時に、5~6秒で検索処理が終了することが分かる。また、近傍条件あるいは複合文字種検索時には、10~15秒で検索が終了する。最も早く検索結果が得られるのは、文字成分表の結果が0件に近いときで、2秒で検索ができる。このことから、DBの検索処理速度は、1.5MB/s ~ 1.2MB/sであるといえる。

4. おわりに

今回開発したBibliotheca/TSは、以下の特長を持つ。

(1) テキストをスキャンする前に検索対象を高精度で絞り込む連接文字成分表方式

(2) 複数個の検索タームの照合処理を一度のテキストスキャンで行うサーチエンジンエミュレータ

(3) フルテキストサーチでの検索もれを解消する同義語・異表記展開処理方式

(4) 検索ターム間の論理条件や近傍条件を判定する複合条件判定処理方式

(5) ネットワーク上の他のワークステーションから検索要求が出せるクライアント・サーバシステム

また、検索処理速度は、典型例で1.5MB/s ~ 1.2MB/sという性能が得られた。

参考文献

[1] 加藤, 他, 「大規模文書情報システム用テキストサーチマシンの研究」, 情報学基礎14-6 (89.7)

[2] 有川, 他, 「テキストデータベース管理システムSIGMAとその利用」, 情報学基礎14-6 (89.7)