

自由語による全文検索のための 2N-6 テキストサーチマシン TSM-I

加藤 寛次, 藤澤 浩道, 川口 久光, 畠山 敦, 大山 光男, 兼岡 則幸
(株)日立製作所 中央研究所

1. はじめに

近年, ワードプロセッサやパーソナルコンピュータ, ワークステーションなど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, 大量の文書情報を一般のユーザが簡単に蓄積・検索できる文書情報検索装置に対する要求が高まりつつある。

こうした要求に応えるものとして, インデックス情報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(自由語全文検索)技術の研究を行なってきた。ここでは, そのためのテキストサーチマシン TSM のアーキテクチャと試作したプロトタイプについて報告する。

2. フルテキストサーチの課題

自由語によるインデックスを用いないフルテキストサーチを実現するためには, 以下に示す課題を解決する必要がある。

- (1) ユーザが指定したキーワードとテキスト中に記述されたキーワードの記述・表記形態の違いによる検索漏れの防止。
- (2) 実用上許容し得る時間内で検索を実行できる高速なスキャン型検索処理の実現。
- (3) 小規模システムから大規模システムまで自由に構成でき, かつデータベースの増加にも容易に対応できるシステムアーキテクチャの実現。

3. フルテキストサーチ方式

上記の課題に対して, 本テキストサーチマシン TSM では以下のようなフルテキ

ストサーチ方式を採った。

- (1) 記述・表記違いの検索漏れの防止
ユーザが指定したキーワードの同義語と異表記語をシステム内部で自動的に生成して, これらをまとめてキーワードとして検索する同義語・異表記検索方式の採用。

- (2) 高速なスキャン型検索処理の実現

① テキストデータの高速な読出し

複数台の磁気ディスク装置を並列に並べ, 同時にテキストデータを読み出す集合型磁気ディスク(マルチディスクユニット MDU)制御方式の採用。

② 高速な多重文字列サーチ

同義語や異表記語を含めて複数のキーワードを, テキストのただ1回の走査で一括

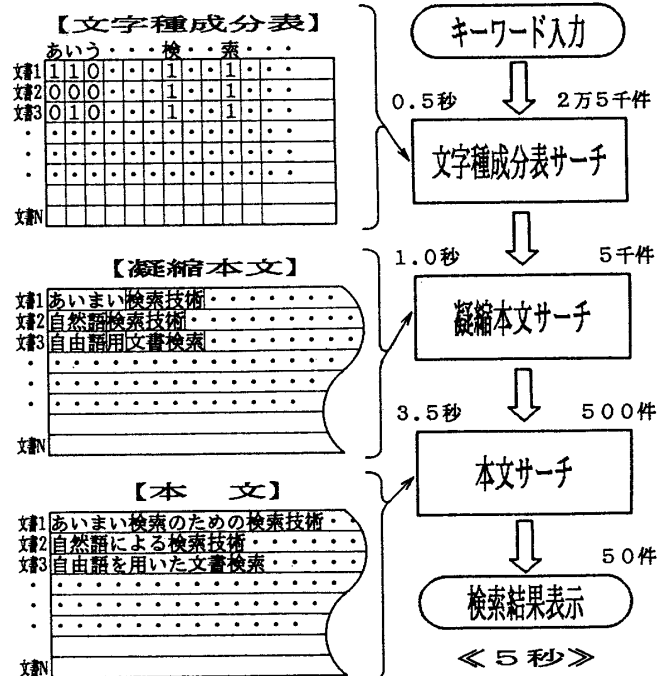


図1. サロゲート型検索加速方式

してサーチできる有限オートマトン型検索プロセッサ(サーチエンジンSE)の採用。

③フルテキストサーチの加速

あらかじめテキスト本文から自動的に作成、情報圧縮しておいた「文字種成分表」と「凝縮本文」を用いた階層的検索によって、検索速度を等価的に高めるサロゲート型検索加速方式の採用。(第1図参照)

④複合条件の高速判別処理

キーワード間の近傍、文脈及び論理条件に関する判別を一貫処理するパイプライン型複合条件判別方式の採用。

(3) システムの拡張性確保

それぞれ1台のサーチエンジンSEとマルチディスクユニットMDUで検索ユニットSUを構成し、これを複数ユニットまとめて1台のサーチマシンTSMとし、更にこのサーチマシンTSMを複数台並列にLANに接続できるようにしたビルディングブロック型アーキテクチャの採用。

4. プロトタイプを試作

今回試作したプロトタイプは、12台の小型磁気ディスクと1台のサーチエンジンSEで構成された1サーチユニットSUで

あり、これは2050/32を介してLAN(CSMA/CD)へ接続される構造になっている。マルチディスクユニットMDUの最大読出し速度は10MB/s、サーチエンジンSEの最大照合速度は20MB/sである。(第2図参照)

このサーチユニットSUには平均20MBの容量の文書が2万5千件格納でき、約1千語の同義語・異表記語を含めて約5秒で検索できる。等価的には100MB/sの検索速度が得られる。

これを用いれば、約2年半分の新聞記事が5秒で検索できることになる。

5. おわりに

サロゲート型検索加速方式を用いることによりスキャン型のフルテキストサーチを実用的な速度で実現できることを確認した。

今後は、実データに基づきフルテキストサーチの有効性を検証して行きたい。

参考文献

- (1) 畠山 他, 自由語検索のための同義語・異表記展開方式, 第39回情処全大
- (2) 川口 他, 自由語検索のための高速文字列検索方式, 第39回情処全大

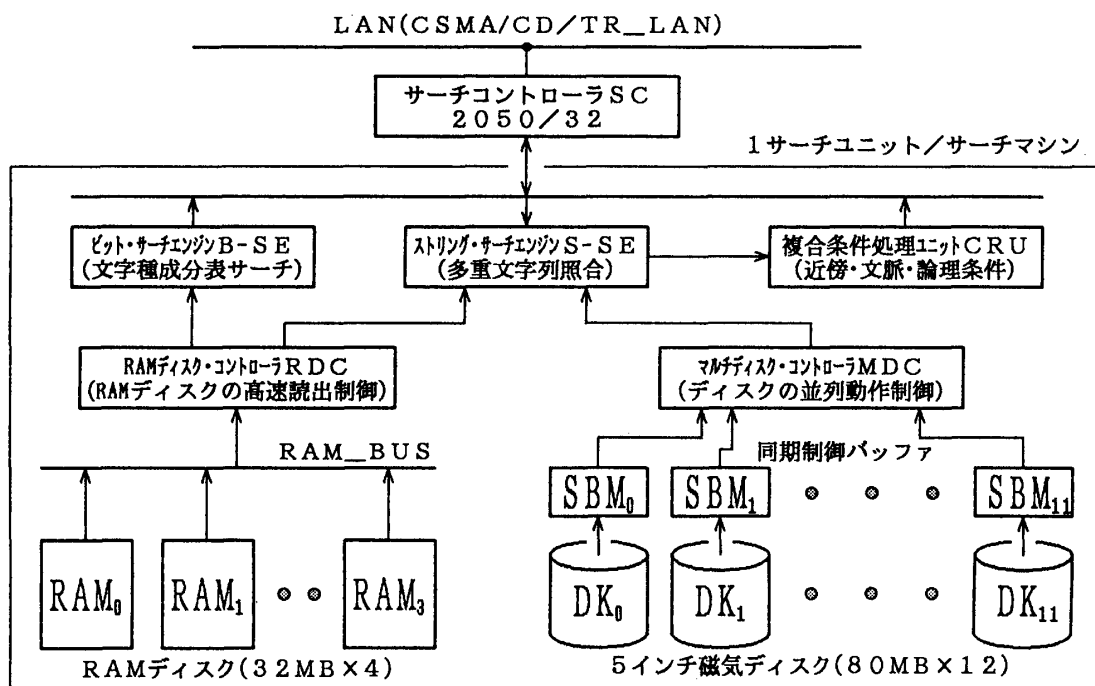


図2. 試作テキストサーチマシンの構成