

テキストDBマシンと検索プロセッサのアーキテクチャ

高橋恒介 山崎哲矢、 西塚博文 本村真人

日本電気(株)、C&Cシステム研究所、研究開発技術本部、マイクロエレクトロニクス研究所

ここではマルチプロトコル文書データの検索を高速化するデータベース(DB)マシンアーキテクチャと入れ子構造で階層的なデータの検索プロセッサのアーキテクチャを提案する。文書データの多様化、マルチプロトコル化に対応して、シソーラス、文書データをOSIメッセージのプロトコルASN.1で記述し、シリコンデバイスを使う高速補助記憶(FAM)に格納する。FAMに直結されるOSIメッセージ検索プロセッサMSPと文字列検索プロセッサDISPによりテキストデータ部分の選択的検索を高速化する。実現性の確認のために開発されたMSPとDISPのLSIを使用するDBマシンのシステム構成と動作が検討される。

Architecture for Textual DB Machines and Search Processor LSIs

Kousuke Takahashi, Testuya Yamazaki, Hirofumi Nishizuka, and Masato Motomura

C&C System Research Laboratories, R&D Planning & Technical Service Division
and Microelectronics Research Laboratories of NEC Corp.

4-1-1, Miyazaki, Miyamae-ku, Kawasaki 216, Japan

This paper proposes textual DB machine architecture and nested data search processor architectures to accelerate the multi-protocol document retrieval. Document and thesaurus data are described by the rule ASN.1 for the OSI messages, stored in the fast auxiliary memory (FAM) using silicone file devices, and searched by a combination of an OSI message search processor MSP and a string search processor (DISP), which have been realized as LSI chips. The DB machine system using LSI chips is also discussed.

1. まえがき

テキスト検索の対象を文書一般に広げると、それは文字列以外に、文字でない部分（写真、グラフ、表、図形）を含むマルチメディア情報であるが、メディア毎でのデータの作成され方の不統一を考えると、文書はマルチプロトコル情報となる。そのような情報の検索をここで考える。

図面や表を含む文書をイメージデータとして記憶している図1の電子ファイリル検索システムでは、従来、データ検索が文書の書誌データを貯えるインデックスファイルやテキストファイルの検索で高速化されてきた。ソフトウェアによる高速化に限界があつて、専用ハードウェアを使うDBマシンが考え出された。図1 (b) のようにインデックスファイルの検索を高速化するDBマシン、図1 (c) のようにテキストファイルの全文検索を高速化するDBマシン、図1 (d) のようにシグナチャファイルのような圧縮テキストを作ってその検索を高速化するDBマシンも開発された。

情報検索を問題解決に必要な知識獲得行動と捉えたと、書誌データを含むインデックスファイルとテキストファイルとシソーラスは一体にして検索できるようにすべきである。その時の問題は、ファイルデータのプロトコル（日本語／英語／外国語、ビット列／コード列、構造／非構造、圧縮／暗号符号、図／表など）の違いへの対応にある。通信ネットワークを介してファイルデータが交換されるような情報社会が到来するとすれば、文書DB設計にもOSIの階層アーキテクチャを念頭におくべきとの指摘がある[2,3]。

異質な文書データがメッセージ列に統合されることによつて文書検索が広く役立つものと考えられる。そこで、これまでの様なテキストデータの検索から、マルチプロトコルの文書データ検索への拡張のアプローチを提案する。

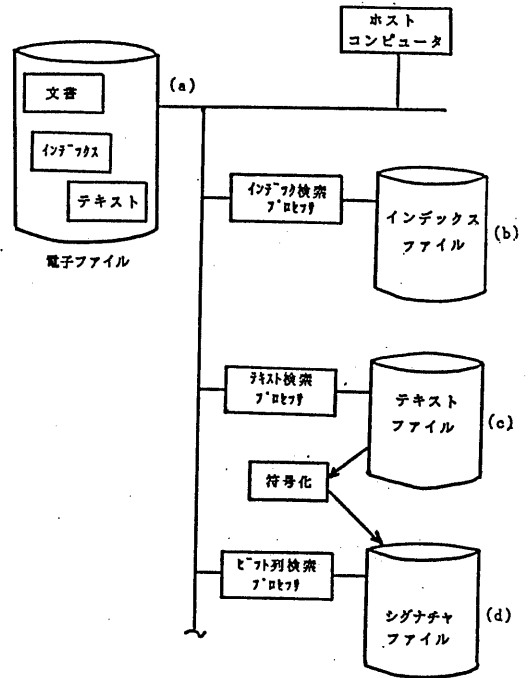


図1、電子ファイルの検索システム例

2. マシンアーキテクチャ

ここで提案するテキストDBマシンは、図2に示すように、高速補助記憶（FAM）と2段構成の検索プロセッサの組み合わせで、書誌データ、シソーラスを含めた文書データのテキスト部分の高速検索を可能にする。図1 (b, c, d) に代るものとなる。構成を簡単にする代りに文書データは完全な非構造から僅かだけ構造化されたものになる。

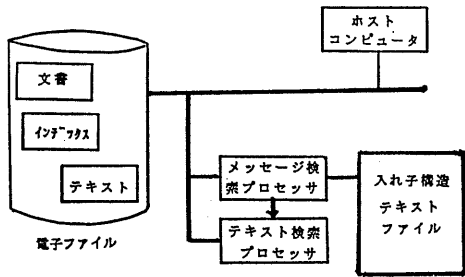


図2、テキストDBマシン

このようなDBマシンはコンピュータネットワークの端末側で使われ、電子ファイルはネットワーク・サーバ側に用意されるとする。ユーザは図2のDBマシンに格納されたマルチプロトコル・マルチメディア文書情報を利用し、問題解決に必要な知識の獲得を行なう。

このDBマシンの設計概念の概略は、

- 1、構造化されたインデックスファイルやシソーラスと非構造のテキストファイルを統合する事、
- 2、ファイルデータをOS Iプロトコル (ASN.1) で記述し、マルチメディア文書の出版・交換への対応の可能性を残す事、
- 3、ファイルデータを入力子で階層構造のOS Iメッセージで合成し、メッセージ毎のプロトコルをメッセージ検索プロセッサで識別する事、
- 4、メッセージのフィールド別に検索条件を用意し、マルチフィールドの検索条件でファイルデータを一括検索する事。

ここに、フィールドは技術分野や文書の内容の分野、メディア、インデックスの属性、プロトコル別処理機能、オブジェクトのクラスなどである。

図3は、従来の一様なテキストデータの検索(a)と、上記設計概念により入力子で階層構造を含むマルチプロトコルファイルデータの一括検索(b)との違いを示す。

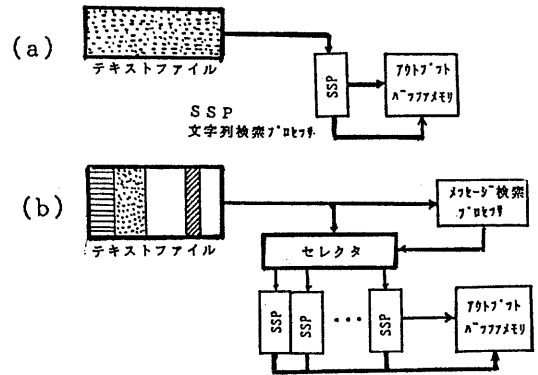


図3、テキスト検索の旧方式(a)と新方式(b)

図3 (b) は指定フィールド別に文字列検索プロセッサSSPを用意し、文字列検索動作の前にフィールド識別を行ない、その結果で指定フィールドのテキストデータを該当するSSPへ導くことを示す。これを行なうために、ファイルデータの内容(OS Iメッセージ)を解釈し、メッセージの先頭でそのフィールドコードを検出するプロセッサを用意する。

図3 (b) を実現する時の主な技術課題は、

- (イ) テキストデータの入れ子構造OS Iメッセージへの変換ソフトと記憶管理ソフト、
- (ロ) 入れ子構造メッセージの選択的読み出しと、フィールドコードのリアルタイム出力を行なうメッセージ検索プロセッサ、
- (ハ) フィールド別検索条件の設定可能でフィールド照合と文字列照合を並行して行なう文字列検索プロセッサの開発
- (ニ) フィールド別シソーラス(用語辞書)の一般用と個人用の準備、である。

シソーラスは一般用、個人用で異なる知識情報の一種であり、情報検索は検索者の問題解決に必要な知識情報を獲得する1プロセスである。質問文を与えて、知識を個人用シソーラスに書き加えることである。

検索者は質問文から一般用シソーラスを使って検索条件を抽出し、それを検索プロセッサに登録して、辞書検索で意味や同義語や関連語を知り、文書のテキスト検索で使用例を知る。検索結果を検索文字列と対応づけて個人用シソーラスに組み込む。

個人別シソーラスは一般用シソーラスが階層構造化されていれば、同じような階層構造になる。文書検索回数を重ねる事を通して、個人シソーラスを充実させる事が知識獲得に相当し、新旧検索条件を個人用シソーラスをベースに比較照合し、全文走査無しで検索結果を得る確率を高めることが性能を改善する学習効果に対応する[7]。

なお、DBマシンはイメージプロセッサの付加によってテキストデータ以外のデータも処理可能とする。MSPがテキストでないデータを区別し、専用メモリを介しイメージプロセッサにデータ振り分けるからである。この点はここでは議論しない。

3、ファイルデータのデータ構造

ここでは、異機種コンピュータ間での文書データ交換に使われる抽象構文規則ASN.1を利用し、文書データをOSIメッセージの集合として記述する。マルチメディア・マルチプロトコルの文書要素データをプロトコルが個別に指定されるメッセージで表し、階層的に関連する文書要素データを入れ子（階層）構造で連結する。結果はバイトシリアルなメッセージ列となる。

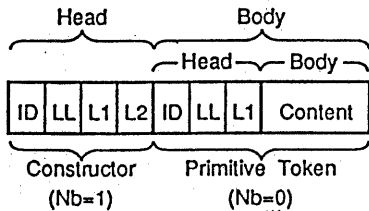


図4、入れ子構造OSIメッセージ例

OSIのメッセージは、図4のように文書要素データをボディとし、その前にボディへの処理機能を示す識別子（ID）と長さコード（LL、L）から成るヘッダを付けたものである。識別子の上位3ビット目のネストビットが'1'か'0'かによってボディ部に入れ子を持つか否かを示す[4]。ヘッダの付加は、情報通信ネットワークを使って文書交換を行なう時のフォーマットに準じており、かなり緩やかな構造化と言える。

OSIメッセージの合成は図5のように行なえる。(a)は文書データが一般的なテキストデータだけであるときのメッセージであって、先頭に1つだけヘッダを持つことを示す。(b)は書誌データと本文から成る文書データが2つのメッセージの先頭にコンストラクタと呼ばれるヘッダを持つことを示す。(c)は書誌データに対応するメッセージが入れ子で属性別のデータを含むことを示す。(d)は本文に対応するメッセージも入れ子で色々な文書要素を含むことを示す。文書要素が写真図面であってもよい。

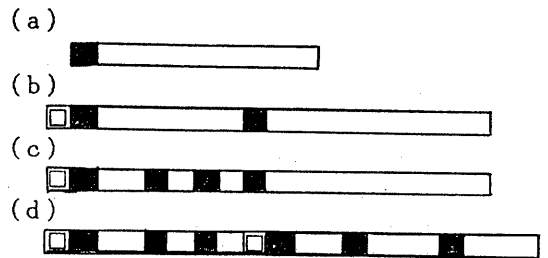


図5、文書データの入れ子による階層構造化

文書要素データを入れ子構造のメッセージ列に変換する手順は、文書要素データをカッコ（または）とスラッシュ/でクラス分けし、（の数でネストレベルを指示するか、文書要素の境界を改行で、階層を改行後の文字列印字開始位置で指示するか、である。長さコードは文字数の計算によって決る。

図5 (a)のように、ヘッダが1つ付加されただけでも、ヘッダ部のIDコードの下位5ビットでボディ部のクラスを任意に分類できる。1ビットをイメージと文字コードの区別に使い、残り4ビットを情報の分野の識別に使うとしても、16フィールドの文書要素をカバー出来る。

図5 (c, d)のように、構造化されたDBのインデックスファイルを文書データと別に含む場合には、入れ子を2層にする。ここでは、1階層目で16種のデータが記憶され、2階層目の中に16フィールドに分類されたデータが含まれる。

入れ子の深さが増すと、メッセージ列は図6のようなシンタックスツリーと対応づけて読まれる。このときの階層の高さをネストレベルNLと呼ぶ。各文書または文書集合の構造が複雑であっても、図6 (b)のように分類できれば、(a)のようなメッセージ列に変換できる。このことは、文書データだけでなく、書誌データ、シソーラス、検索結果などのデータ記述についても同じである。

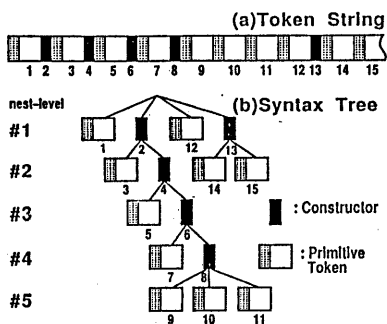


図6、メッセージ列と階層トリーとの対応付け

文書データや書誌データやシソーラスや検索結果をメッセージ列に変換したあとは、バイトシリアルなファイルデータとしてFAMに記憶する。以下はこのようなメッセージ列の検索システムを考える。

4、OSIメッセージ検索プロセッサ

OSIメッセージ列の検索には図7に示すようなメッセージ検索プロセッサMSPが必要となる。入力はメッセージ列であり、出力はメッセージ列をFAMから呼び出すアドレス、メッセージ列のヘッダ部の解析結果 (ID、NLなど) である。

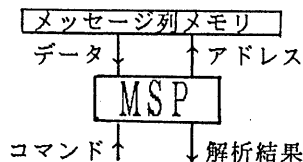


図7、MSPのLSIチップのイメージ

メッセージのヘッダ部の解析は各メッセージの先頭がヘッダ部のIDコードであること (図4) をベースに、図8 (a) の順序論理に従って行なわれる。1バイト目でIDとネストビットNbを、2バイト目で長さ形式LbとLコード数を、3バイト目以降でLコードを読み取り、ボディ部長さBLや次のメッセージの開始位置 (HA) を計算する。その間にID、NL、Nbの照合を行ない、照合結果Zbで次に読み出すアドレスを計算する。

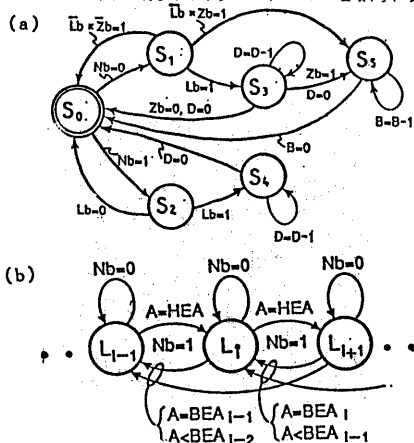


図8、ヘッダ部解析の状態遷移図 (a) とネストレベル判別のための状態遷移図 (b)

ネストレベル判別には図8 (b) の状態遷移図が使われる。Nb が1のIDコードを持つヘッダ(コンストラクタ)がきた時に入れ子のメッセージが続き、そのメッセージのNLが1つ高くなる。入れ子のメッセージを読み終わると、NLは元の高さにもどる。入れ子メッセージの中に入れ子メッセージが続くとNLは1段づつ高くなるが、深い方のメッセージが読み終わり次第、NLは下がる。何段もまとめて下がる場合もあるため並列処理が必要となる[5]。

メッセージ検索プロセッサMSPは図8 (a) の状態遷移図に従った順序論理回路、ミスマッチメッセージのボディ部を読み飛ばすアドレス計算回路、図8 (b) に従って次のメッセージのNLを決定するネストレベルモニタ回路、検索モード選択回路から成る。図9は72ピンPGAに実装されるLSIチップを示し、表1はMSPのLSIチップの諸特性を示す[5]。

検索モードや検索範囲などのコマンドは始めに設定される。MSPの検索動作モードは次のとおりである。

- 1, 実時間検索モード：ヘッダ部を除いた全文検索
- 2, ヘッダ検索モード：ヘッダ部だけの検索
- 3, 断続的検索モード：メッセージ単位で検索停止
- 4, スキップ検索モード：マッチメッセージの検索

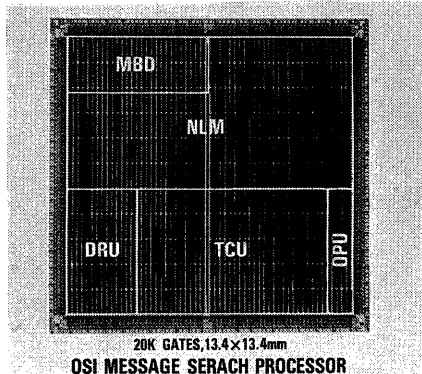


図9、メッセージ検索プロセッサMSP

表1、MSPのLSIチップの諸特性

最大検索範囲	32ビットアドレス空間
検索処理速度	280万オクテット/秒 (5.6MHz)
入力データ	OSIプロトコルX.409 又はASN, 1で記述のもの
出力データ	アドレス/マッチトークン HA, BA, ID/NL, BL
使用電源単一	5V+10%、約75mA (5MHz)
LSIパッケージ	72ピンPGA、TTL

5、高速補助メモリFAM

MSPが発生するアドレスに従ってメッセージ列のファイルデータを読み出せるFAMとしては大容量・低コストで、低電力のRAMが望ましい。シリコンファイルはDRAMをポータブルコンピュータ用に改良したものでありDRAM並みの価格で消費電流が小さい。非動作時をオートリフレッシュモードにし、データ保持電流が30uAにできる。300チップなら約10mAになる。この消費電流は1AHの電池で100時間のデータ保持を許す。停電や電源故障の期間のデータ保持には十分である。なお、300チップは4Mb RAMの場合、1Gビットに相当する。

シリコンファイルでは多数のチップを使っても同時に動作するチップが少ないので動作時の消費電流も小さい。1Gb当たりで1A以下となる。データ転送レートは5~20MB/s以上であり、ページモードでは更に高い。フラッシュメモリに較べると書き換え動作が読み出しと同じ速度で実行されるメリットがある。磁気ディスクに較べると、MBの価格が1桁以上高いが、性能価格比で見ると、すでに、シリコンファイルが有利になる。そこで、以下ではFAMにシリコンファイルメモリSFMを採用することとする。

6、文字列検索プロセッサ

文字列検索方式としてAM法、CA法、FSA法、DP法とPSL法があるが、LSI化にはPSL(プログラマブル順序論理)法が良い。連想メモリCAMに順序論理回路SLCを結合させたPSL法の文字列検索プロセッサSSPをLSI化すると、可変長の文字列照合、アンカー/ノンアンカー文字列照合、固定長・可変長ドントケア文字列照合、ワイルドカード文字列照合、あいまい文字列照合の機能が可能になり、かつ、512文字の文字列登録と、64以下の検索文字列の並列照合と、20MB/sでの文字列検索も可能になる[6,7]。

しかし、このままの文字列検索プロセッサでは階層構造をもったファイルデータの検索に使えない。OSIメッセージ列のファイルデータから指定プロトコルのメッセージに含まれるテキストデータを選び出せば、SSPで検索が可能になる。

メッセージ検索プロセッサMSPはヘッダ部を検出し、特定フィールド(プロトコル)のテキストデータだけを選択してくれる。問題は、図3(b)に示したように、複数のフィールドにまたがる検索文字列により複数フィールドのテキストデータを1回のスキャンで検索させられないかという事である。

m個のSSPを1チップで実現するLSI設計に際しては、単にLSI規模を高めるだけでなく、SSPのCAM部のみをm個に増やしSLC部をそれらで共用させるVLSIアーキテクチャと、m個の中の1つのCAMの選択信号を、メッセージのフィールド(ID・NL)照合によって発生するマッチ回路の内蔵が望まれ、図10に示すような辞書検索プロセッサDISPが開発された[8,9]。文字記憶容量が16倍に増え、1チップが16個分のSSPに相当する。それにも係わらず、検索速度がSSPの場合以上と報告されている。

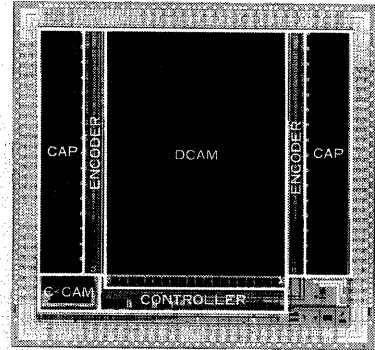


図10、DISPのLSIチップ写真[8,9]

DISPのアーキテクチャはCAM部をSRAMで実現し、SRAMの大容量化に合わせ、同じチップ面積でCAM部の記憶容量を増大可能にし、今後の大容量FAMや大規模シソーラスへの対応を容易にする。なお、DISPは複雑な階層構造の用語辞書を検索するが、辞書検索専用ではない。

6、DBマシンのシステム

したがって、マルチプロトコル文書データ中からテキストデータ部分を選び出してテキスト検索を可能にするDBマシンのシステムは、図11に示すようになる。SFMを用いたFAMとそれをアクセスするMSPとFAMから読み出されたデータを検索するDISPと、検索結果の出力バッファメモリOBMが主な構成要素である。

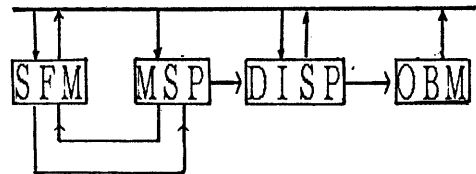


図11、テキストDBマシンのシステム構成

実際にはOBMの後に問い合わせ解読プロセスQRPが必要で、検索結果がOBMに集められた後に、DISPへ検索条件式を入力し、それに対応したOBMの出力をQRPで処理する。

このシステムで行なわれる検索処理は、質問文と一般シソーラスとの照合、シソーラス間の用語文字列照合、シソーラスと文書データの照合、文書集合検索、引用文献欄の検索、個人シソーラスの検索などである。

検索速度の高くなる理由は以下のようになる。

- (1) 指定フィールドに含まれる各文書データの指定部分を検索する場合に、指定されない部分をスキップしながら、FAMの中のバイトシリアルなメッセージ列を検索すること、
- (2) 複数のフィールドの文書を1度で検索出来ること、
- (3) 問題解決に必要な用語を個人用シソーラスの検索で確認し、無駄な走査を省けること、
- (4) FAMやDISPやMSPが高速に動作すること、

7. 結論

マルチプロトコルの文書データをOSIに準拠した抽象構文規則ASN.1で記述することを考え、そのような文書データの検索を可能にするテキストDBマシンと検索プロセッサのアーキテクチャを提案し、検索プロセッサMSP、DISPのLSI開発結果を示した。FAMには低電力のシリコンファイルメモリが適当である事を示した。

DBマシンは単純な構成で、多様な文書データ、シソーラスの高速検索の可能性を与える。今後の課題はSFM、MSP、DISPを組み合わせたシステム実験で具体的に情報検索を評価することである。

最後に、本研究の機会と御討議を頂いた山本昌弘所長並びに関係者各位に謝意を表す。

<参考文献>

- 1, C.W. Bachman; 'A Personal Chronicle: Creating Better Information Systems with some Guiding Principles', IEEE Trans. Vol. KE-1(1), Mar. 1989
- 2, 加藤・藤澤、その他「全文検索用テキストサーチマシンの開発」、電子情報通信学会、技術研究報告DE89-38, Dec. 1989
- 3, I.A. Macleod; 'A Query Language for Retrieving Information from Hierarchic Text Structure', The Computer Journal, Vol. 34(3), p254, 1991
- 4, 高橋、中川路「構文定義用言語ASN.1の特質と処理系の現状と動向」情報処理 Vol. 31(1), Jan. 1990
- 5, 高橋、高橋、山崎「マルチメディア文書の記憶・検索とプロトコル解析LSI」、電子情報通信学会、研究技報CPSY90-91, Jan. 1991
- 6, K. Takahashi, H. Yamada and M. Hirata; 'A String Search Processor LSI', JIP, J. of Information Processing, Vol. 13(2), 1990
- 7, 高橋、永井、山田、「テキストデータベースの学習型検索方式」電子情報通信学会、DE88-3, May 1988
- 8, M. Motomura, et al; 'A 1.2-Million transistor, 33MHz, 20-bit Dictionary Search Processor with a 160kb CAM', Digest of Technical Paper, TAM5.4, IEEE ISSCC'90, Feb. 1990
- 9, 本村、その他「120万トランジスタ辞書検索プロセッサ(DISP)ーその構成と機能ー」、1990電子情報通信学会春期全国大会C-662, p 5-2