

文章解析アクセラレータ(1) -形態素抽出マシンの試作-

福島 傑一 大山 裕 宮井 均

日本電気株式会社 C & C システム研究所

本稿では、文章解析処理の高速化への新しいアプローチとして、文章解析アクセラレータハードウェアを提案する。そして、その第1ステップとして試作した形態素抽出マシンについて報告する。形態素抽出処理は、単語辞書を検索することによって、入力されたテキストに出現した可能性のある全単語（全形態素）を抽出する処理である。文章解析の最初に不可欠な処理である上に、処理時間に占める比率が大きいことから、高速化が強く望まれていた。試作した形態素抽出マシンでは、8万見出しの単語辞書を用いた、1万文字のテキストに対する形態素抽出処理が、1秒程度で実行できる。この処理速度は、従来のPC9800上のソフトウェアと比べて、100～1000倍という高速なものである。本マシンは、①単語を構成する全文字の同時照合、②先頭文字による辞書範囲の絞り込み、③シフトレジスタによるテキストの順送り、などの手法を組み合わせた形態素抽出の新しいハードウェアアルゴリズムを実現している。各文字が複数の候補をもったテキストについても、候補数に対して線形オーダーの時間で、形態素抽出処理が行なえる。

Natural Language Parsing Accelerators (1)
- An Experimental Machine for Morpheme Extraction -

Toshikazu FUKUSHIMA, Yutaka OHYAMA and Hitoshi MIYAI

C&C Systems Research Laboratories, NEC Corporation

1-1, Miyazaki 4-chome, Miyamae-ku, Kawasaki City, Kanagawa 213, Japan

(fuku@tsl.cl.nec.co.jp, ohyama@tsl.cl.nec.co.jp, miya@tsl.cl.nec.co.jp)

This paper proposes natural language parsing accelerators, as a new approach to parsing speed-up. It also describes an experimental machine for morpheme extraction, which is designed as the first step toward achieving the natural language parsing accelerators. In the morpheme extraction process, all morphemes, which are considered probably used to construct input text, are extracted by searching a morpheme dictionary. Speed-up requirement for this process is very strong, because this is indispensable process for the parsing, and spends a large percentage of parsing time. The experimental machine can extract morphemes from 10000 character Japanese text by searching an 80000 morpheme dictionary in one second. It runs 100-1000 times as fast as morpheme extraction programs on personal computers. It realizes a new hardware algorithm, which uses parallel comparators, an index and shift registers. It can treat multiple text streams, which are composed of character candidates, as well as one stream of text. The algorithm is implemented on it in linear time for the number of candidates.

1 はじめに

1-1 文章解析処理の高速化の効果

近年、計算機による文章解析技術は、ワードプロセッサや機械翻訳システムなどのかたちで実用化され始めている。この実用化をさらに進めてゆくためには、より深い解析の追究による高精度化と並行して、処理の高速化を行なってゆく必要がある。文章解析処理を高速化することは、次の(1)～(3)のようなかたちで文章解析応用の拡大・実用化につながる。

- (1) レスポンスの高速化
 - (2) 处理対象ドキュメントの多量化
 - (3) 处理の高精度化

(1)は最も直接的な効果である。自動通訳・音声認識における言語処理、ワードプロセッサのかな漢字変換、データベース検索の自然言語インターフェースなどのように、ユーザの入力に対する実時間の応答が要求される応用で望まれている。

一方、機械翻訳、ドキュメント検査（自動校正）、文字認識における言語処理（大量ドキュメント入力）、データベースへ登録するための自動キーワード付けなどの応用では、非常に多量のドキュメントの一括処理を行なう。(2)は、このような応用のメリットを高め（人手と比較した場合の所要時間を短縮し）、その実用化を促進する。

(3)は、間接的な効果である。処理が高速化できると、同一の時間内に、これまでより豊富な処理が行なえる。そこで、より大規模の辞書・知識ベースを参照したり、意味・文脈処理により多くの時間を割り当てたりすれば、一定の処理時間・レスポンス時間を保ったままで、解析精度を上げることが可能となる。

1-2 文章解析アクセラレータの提案

文章解析処理の高速化は、これまで、おもに、逐次型の汎用計算機の発展と、その上のソフトウェア／アルゴリズムの改良に支えられてきた。しかし、上記のような高速化の要求は、それ以上に強い。

このさらなる高速化の要求に対して、近年、汎用並列計算機が期待されている。並列論理型プログラミング言語で記述された構文解析プログラム[5]や、P-RAM (Parallel Random Access Machine) を想定した並列構文解析アルゴリズム[6][7]などが考案されている。しかし、汎用並列計算機の実用化にはもう少し時

間を要すると思われ、また、実際に実行した場合には、共有メモリに対するアクセス競合や、プロセス間通信のボトルネック問題などが発生して、十分な性能が得られない可能性がある。

一方、パターン認識の分野では、様々な専用ハードウェアが開発され、高速化が実現されてきた。そのなかには、文章解析処理へ適用できる可能性をもつものもある。例えば、音節の種類だけのアドレス空間をもつRAMを用いた音声認識用の辞書照合方式[3]、マルチプロセッサにより複数テキスト位置で並列に辞書検索を行なう形態素抽出方式[4]、シストリックアレイを用いてCYK法やEarley法を並列化した文脈自由言語認識方式[8][9]などである。

筆者らは、今後、文章解析処理についても、ハードウェアアプローチが必要と考える。そこで、文章解析処理を高速化する専用ハードウェア（文章解析アクセラレータ：Natural Language Parsing Accelerators）を提案する。

2 形態素抽出処理のハードウェア化

2-1 形態素抽出処理の選択

形態素抽出処理は、単語辞書を検索することによって、入力されたテキストに出現した可能性のある全単語（全形態素）を抽出する処理である（図1参照）。通常、形態素抽出処理に引き続いて、文法的な情報をもとに単語間の接続可否を調べて文節を形成する接続検定処理が行なわれる。入力された文章（日本語ではべた書きされた文章）を単語列（形態素列）に分割し、

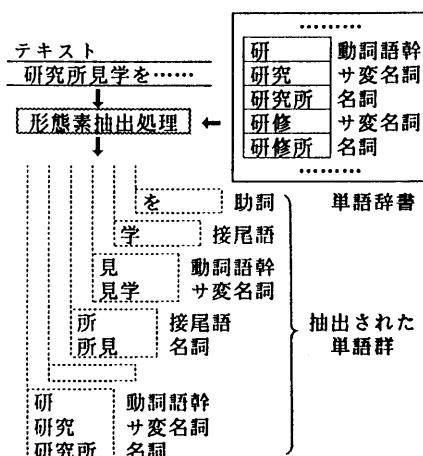


図1 形能素抽出処理

各単語の品詞を認定する形態素解析処理は、それら形態素抽出処理と接続検定処理をあわせたものを指す。文章解析処理は、さらに、構文解析処理・意味処理・文脈処理などを含んでいる。

形態素抽出処理は、文章解析の最初に不可欠な処理である上に、処理時間に占める比率が大きい（形態素解析では処理時間の7割を越える）ことから、高速化が強く望まれる。一方、ソフトウェア的には手法が確立されており、処理の大半が単純な照合処理の繰り返しだることから、ハードウェア化に適している。

また、近年、並列構文解析方式が研究され始めているが[5][6][7]、実際にインプリメントした場合には、共有メモリに置かれる辞書のアクセス時間が性能のボトルネックになることが予想される。この問題を軽減するためにも、形態素抽出処理のような辞書周辺処理の高速化は重要である。

以上の点から、高速化が望まれ、かつ、ハードウェア化に適した処理単位として、まず形態素抽出処理を取り上げ、そのハードウェア化を行なった[1]。

2-2 ハードウェア化の方針

ハードウェアを想定して形態素抽出処理を並列化する1つの方法として、テキストと単語辞書とを照合する際に、次のような多重化を行なうことが考えられる。

Ⓐ 照合対象単語の多重化（図2参照）

Ⓑ 単語抽出位置の多重化（図3参照）

従来のソフトウェアによる形態素抽出処理では、テキストのある位置から部分文字列を切り出して、その部分文字列に一致する単語を、単語辞書から検索する処理が繰り返される。その単語辞書検索の際には、部分文字列と単語との照合が、辞書内の単語ごとに逐次行なわれる（一対一の照合）。これに対して、Ⓐでは、テキストのある位置から切り出した部分文字列と、単語辞書内の複数の単語とを同時に照合する（一対多の照合）。一方、Ⓑは、テキストにおける形態素抽出を行なう位置を多重化するものである（多対一の照合）。すなわち、テキストの1文字目・2文字目・3文字目・……から同時に抽出処理を行なう。

Ⓐの方式は、従来、連想メモリ[10]を用いて実現されている。ただし、現状の連想メモリLSIは1個の容量が数K～数十Kビットであるから、数万見出しの単語辞書を一度に登録するには、数千個を越える多数の連想メモリLSIが必要になる。数個～十数個とい

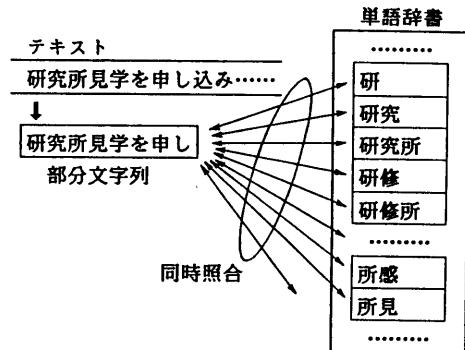


図2 照合対象単語の多重化

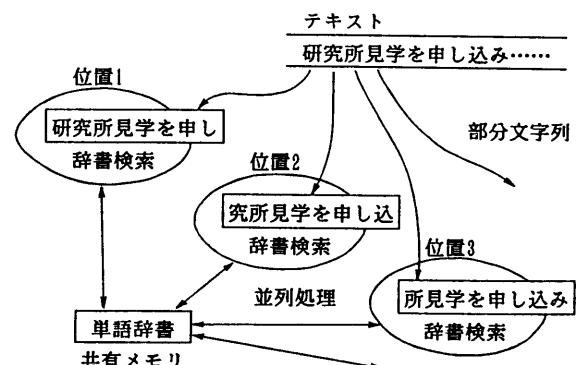


図3 単語抽出位置の多重化

う現実的な数の連想メモリLSIを想定し、単語辞書を1回で登録できるサイズの複数個に分割して、その登録・検索を繰り返すようにすると、今度は、繰り返し登録に要する時間が障壁となり、極めて多量のテキストを一括処理するときでないとメリットが出なくなる[2]。

Ⓑの方式は、従来、テキストの各文字位置に処理ユニットを対応付けたマルチプロセッサ形式で実現されている[4]。この場合、単語辞書は共有メモリに置かれ、複数処理ユニットからのアクセス競合の調整機構が設けられる。処理ユニット数を増加すると、処理速度は向上するが、それにつれてアクセス競合の頻度が増すので、処理速度の向上は頭打ちになる。文字列検索LSI[11]に1文字ずつ位置をずらしてテキストを登録し、検索時に単語辞書を入力する方法で、Ⓑを実現することもできるが、その場合は、単語の抽出される順が位置的にランダムとなり、以降の処理（接続検定処理など）との整合が悪いという欠点がある[2]。いずれの実現方法とも、ハードウェアの構造が複雑になり、価格性能比の面では現実的でない。

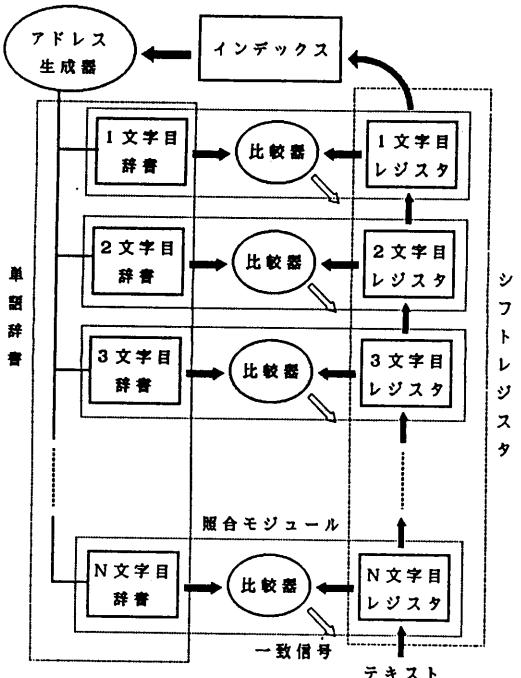


図4 基本的な構成

以上のような検討にもとづき、筆者らは、ⒶⒷのよ
うな多重化の方針はとらず、基本的な1回の照合処理
のサイクルを可能な限り短縮するアプローチをとった。

3 形態素抽出ハードウェアアルゴリズム

3-1 基本的な構成・動作

考案した形態素抽出ハードウェアアルゴリズムを実行するための基本的な構成は、図4の通りである。おもな構成要素は、単語辞書、シフトレジスタ、インデックス、アドレス生成器、比較器である。

単語辞書は、図5に示すような形式で、1文字目辞書～N文字目辞書のN個に分割されている。長さがNより短い単語については、Nに満たない部分に、あらかじめ定められた残余記号が格納されている（図5では「*」が残余記号を表わしている）。

シフトレジスタは、1文字目レジスタ～N文字目レジスタで構成され、そのなかにはテキストのN文字分が格納される。シフトレジスタでは、図6に示すような文字列の順送りが行なえる。

インデックス（先頭文字表）からは、ある文字で始まる単語群の、単語辞書における存在範囲が得られる（図5参照）。

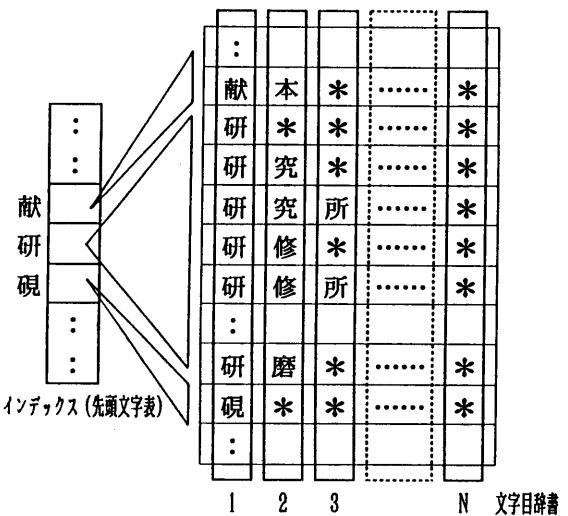


図5 単語辞書の構成とインデックスとの関係

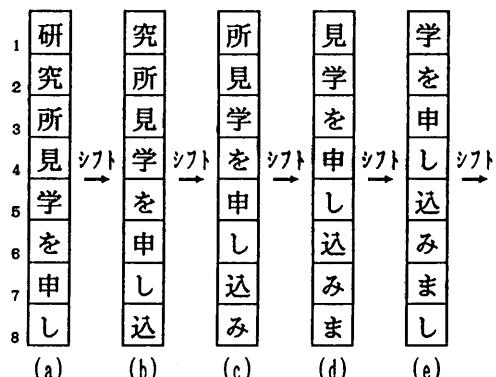
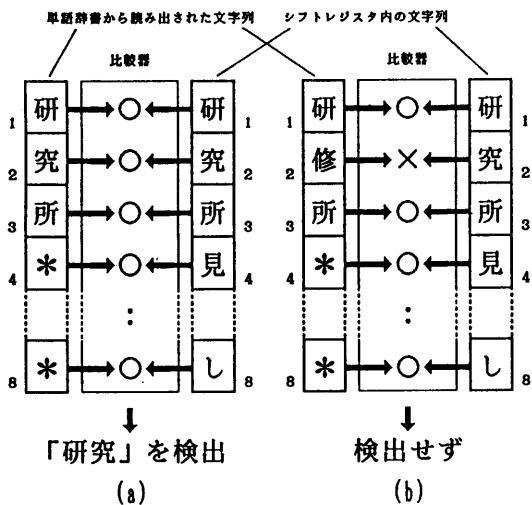


図6 シフトレジスタの動作 (N=8のときの例)

アドレス生成器は、1文字目辞書～N文字目辞書に
対して、同一のアドレスを与える、かつ、そのアドレス
を同時に変化させる。変化させるアドレスの範囲は、
シフトレジスタの1文字目と同じ文字で始まる単語群
の存在範囲であり、インデックスから得る。

比較器は、単語辞書とシフトレジスタの各文字位置
を対応付けるように、全部でN個設けられている。n
文字目 ($1 \leq n \leq N$) に対応する比較器は、n文字目
辞書から読み出された文字が、シフトレジスタのn文字
目が残余記号のどちらかに一致したときに、一致信
号を出力する（残余記号はワイルドカードの働きをす
る）。単語辞書内の1単語に対応するN文字、および
シフトレジスタ内のN文字は、同時に読み出せるよう
になっており、上記の照合処理は、N個の比較器で同
時に行なわれる。全比較器から一致信号が出力された



ときに、単語（形態素）が1個検出されたことになる（図7参照）。

さて、上記のような構成にもとづく、形態素抽出のハードウェアアルゴリズムは、以下の通りである。

■メインプロシジャー：

【ステップ1】シフトレジスタに入力テキストの先頭のN文字をロードする。

【ステップ2】1文字目レジスタにテキスト末尾マークが到達するまで、プロシジャー1を繰り返す。

■プロシジャー1：

【ステップ1】1文字目レジスタと同じ文字で始まる単語群の単語辞書内アドレス範囲を、インデックスより得て、その先頭アドレスを単語辞書のカレントアドレスとする。

【ステップ2】カレントアドレスが上記のアドレス範囲内にある間は、プロシジャー2を繰り返す。

【ステップ3】シフトレジスタ内の文字列を、1文字分順送りする。

■プロシジャー2：

【ステップ1】カレントアドレスにおいて単語辞書から読み出されたN文字と、シフトレジスタ内のN文字とを、N個の比較器で同時に照合し、単語の検出を判定する。

【ステップ2】カレントアドレスをインクリメントする。

以下では、動作例を簡単に説明する。

「研究所見学を申し込みました……」という入力

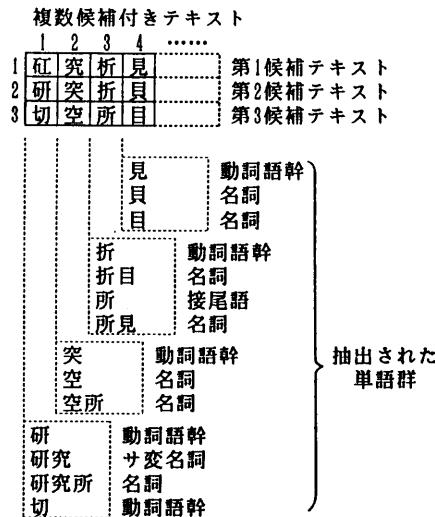


図8 複数候補付きテキストと形態素抽出結果 (M = 3 のときの例)

テキストに対して、メインプロシジャーのステップ1の直後に、シフトレジスタの内容は、図6の(a)のようになる。プロシジャー1のステップ3ごとに、シフトレジスタ内の文字列は、(a)→(b)、(b)→(c)、(c)→(d)、(d)→(e)、………と移動する。プロシジャー1は、図6の(a)(b)(c)(d)(e)………の各々に対して実行される。

図6の(a)の状態では、シフトレジスタの1文字目が「研」である。そこで、プロシジャー1のステップ1では、単語辞書における「研」で始まる単語群のアドレス範囲が、インデックスより得られる。図5の単語辞書では、「研」「研究」「研究所」「研修」「研修所」………「研磨」が、そのアドレス範囲に該当する。そして、そのアドレス範囲内の各単語について、プロシジャー2のステップ1の照合処理が順次実行される。

以上を簡単にまとめると、本アルゴリズムは、次のような3つの手法を組み合わせたものである。

- ①単語を構成する全文字の同時照合
- ②先頭文字による辞書範囲の絞り込み
- ③シフトレジスタによるテキストの順送り

非常にシンプルな手法の組み合わせであるため、構成が簡単になり、2-2で述べた方針通り、照合サイクルの短縮による高速化が可能となる。

3-2 複数候補付きテキストへの拡張

本節では、3-1に示したアルゴリズムを、テキストの各文字が複数個 (=M個とする) の候補をもつ場

合に拡張する。複数候補付きテキストは、文字認識や音声認識の結果などにおいて発生する。図8に、複数候補付きテキストと、それに対する形態素抽出結果の例を示した。

図9は、拡張アルゴリズムを実行するための構成である。図9の構成は、図4の構成に対して、次の3点の修正が加えられている。

まず、シフトレジスタをM個設ける。M個のシフトレジスタは同期して、テキストの同一位置のM個の候補を同時に順送りする。

次に、n文字目に対応する比較器は、n文字目辞書から読み出された文字と、M個のシフトレジスタのM種類のn文字目とを照合するようとする。その比較器は、n文字目辞書から読み出された文字が、残余記号、または、M個のシフトレジスタのいずれかのn文字目に一致したときに、一致信号を出力する。

第三に、切り換え器を設ける。切り換え器は、M個のシフトレジスタの1文字目を、順番に切り換えてインデックスへ渡す。

このような構成を用いて、複数候補付きテキストに対する形態素抽出を行なうアルゴリズムは、3-1に示したアルゴリズムにおいて、プロシジャー1を次のようなプロシジャーに置き換えればよい。

■ プロシジャー1：

【ステップ1】カレントレベルを第1候補とする。

【ステップ2】カレントレベルがM以下である間、プロシジャー1.5を繰り返す。

【ステップ3】シフトレジスタ内の文字列を、1文字分順送りする。

■ プロシジャー1.5：

【ステップ1】カレントレベルのシフトレジスタの1文字目と同じ文字で始まる単語群の単語辞書内アドレス範囲を、インデックスより得て、その先頭アドレスを単語辞書のカレントアドレスとする。

【ステップ2】カレントアドレスが上記のアドレス範囲内にある間は、プロシジャー2を繰り返す。

【ステップ3】カレントレベルをインクリメントする。

図10は、このような拡張を行なった場合の、単語の照合（プロシジャー2のステップ1）の例である。単語辞書から読み出された1個の単語の照合に要する時間は、候補数Mによらない。ただし、シフトレジスタの1文字目がM通りになるので、テキストの各位置で

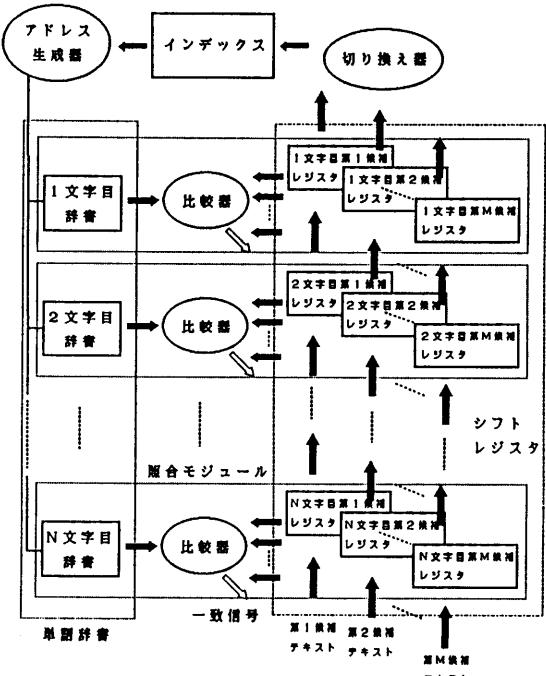


図9 複数候補付きテキストへ拡張した構成

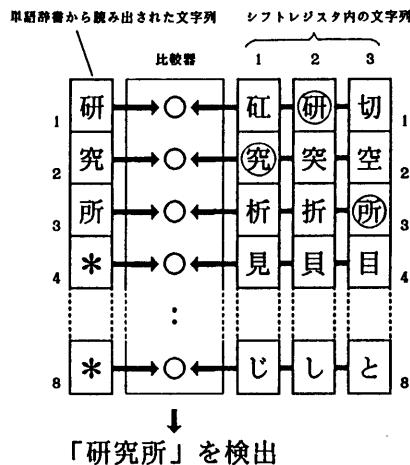


図10 複数候補付きテキストでの同時照合
(N=8、M=3のときの例)

照合する単語数（アドレス範囲）は、平均してM倍になる。

従来のソフトウェアでは、候補数がM個になったときの形態素抽出の所要時間は、候補が1個の場合に比較して、Mの累乗オーダー（最悪で M^N 倍）になる。それに対して、上記の方式では、所要時間がM倍（線形オーダー）に抑えられる。

本アルゴリズムは、音声認識を対象とした浜口らのアルゴリズム[3]と基本的な考え方は共通している。浜口らのアルゴリズムでは、比較器を用いず、文字（音節）の種類だけのアドレス空間をもつRAMへ候補をマッピングして一致を判定する。したがって、音声のように種類が少なく候補数が多い対象には浜口らのアルゴリズムが適しているが、漢字を含み字種の多いテキストには筆者らのアルゴリズムが適している。さらに、シフトレジスタによるテキストの順送り機構を備えている点も、筆者らのアルゴリズムの特長である。

4 形態素抽出マシン

4-1 マシンの概要・諸元

試作した形態素抽出マシンは、3-2に示した形態素抽出アルゴリズムを実行する専用ハードウェアである。パソコン（PC9800シリーズ）のバックエンダーマシンとして動作する。約80個のメモリIC（総容量：約2MB）と約500個のロジックICを載せた、合計12枚のボードから構成されている。図11にその外観を示した。

アルゴリズムのパラメータは、 $N = 8$ （16まで拡張可能）、 $M = 1 \sim 3$ とした。漢字かな混じりテキストを対象としており、扱う文字コードは16ビットである（各メモリやレジスタのデータ幅も同様）。単語辞書は約8万見出しのものを用いている（自立語だけでなく付属語・接辞なども含む）。

4-2 評価

本マシンは10MHzのクロック（100n秒サイクル）で動作している。1単語の照合処理（3-1に示したプロシジャー2）は、3クロック（=300n秒）で実行できる。このとき、形態素抽出の所要時間は、およそ次の式のようになる。

$$A \times D \times L \times M \times 300 \text{ [n秒]}$$

ここで、Dは単語辞書の見出し数、Lはテキスト長、Mは文字位置当りの文字候補数、Aはインデクシング係数（Dに対するインデックスによって絞り込まれた見出しの比率）である。

種々のテキストに関して、形態素抽出所要時間を測定した結果を、図12に示す。単語辞書の見出し数Dは前述の通り約8万件である。この結果は、形態素抽

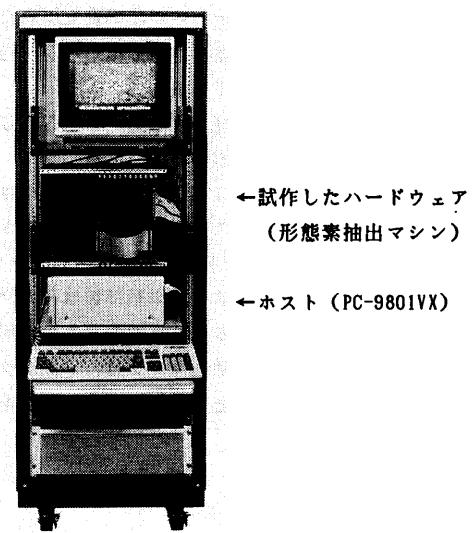


図11 システムの外観

形態素抽出所要時間 [秒]

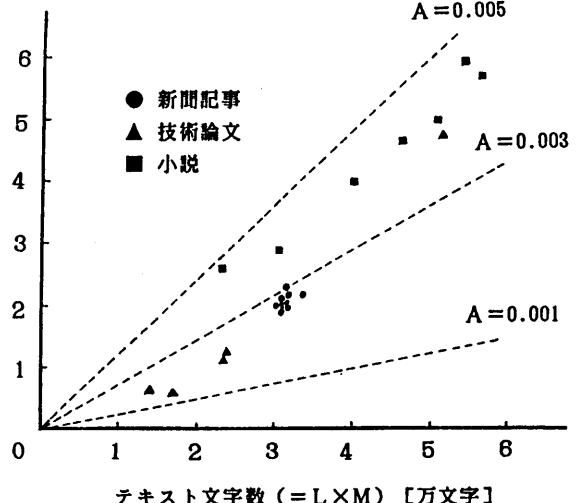


図12 形態素抽出マシンの形態素抽出所要時間

形態素抽出方式		テキスト1	テキスト2	テキスト3	テキスト4
ソ	二分探索法[12]	564	642	615	673
ウフ	先頭文字表利用二分探索法	133	153	147	155
エト	順序ハッシュ法[13][14]	406	440	435	416
ア	木構造見出しの桁探索法[12]	52	58	54	54
形態素抽出マシン		0.56	0.50	0.51	0.44
(単位:秒)					

図13 5千文字テキストに対する処理時間の比較

出マシンによれば、8万見出しの単語辞書を用いた、1万文字のテキストに対する形態素抽出処理が、約1秒で実行できることを示している。

また、4種類のテキスト（いずれも、 $L = 5000$ 、 $M = 1$ ）に対する形態素抽出の所要時間を、ソフトウェアと比較した結果を図13に示す。単語辞書の見出し数Dは、これまでと同様に8万件である。ソフトウェアはPC-98XL²（CPU: 80386、クロック：16MHz）上で動作させ、単語辞書はRAMに置いた。パソコンとしてはかなり高性能なものを用いたにもかかわらず、形態素抽出マシンは、ソフトウェアの100倍～1000倍という高速な処理速度を達成している。

5 おわりに

文章解析処理の高速化への新しいアプローチとして、文章解析アクセラレータを提案し、その第1ステップとして、形態素抽出マシンを試作した。本マシンでは、8万見出しの単語辞書を用いた、テキスト1万文字に対する形態素抽出処理が、1秒程度で実行できる。この処理速度は、パソコン上のソフトウェアの100倍～1000倍という高速なものである。本マシンで実現した形態素抽出のための新しいハードウェアアルゴリズムは、各文字が複数の候補をもったテキストも対象とことができ、かつ、候補数に対して線形オーダーの処理時間を達成している。

上記の高速性は、シンプルな構成により得られている。このシンプルな構成は、VLSI化に適しており、現在、その検討を進めている。形態素抽出LSIは、既に商用化されている機械翻訳システムなどの処理速度の向上をもたらすであろうし、多量のドキュメントを対象とした自動校正や自動キーワード付けなどの実用化を促進することも期待できる。また、人間の入力速度に追従するように各種ソフトウェア手法を導入しているかな漢字変換についても、ソフトウェアの負担が軽減され、さらに、従来よりもはるかに大語彙の辞書を用いることも可能になる。

今後は、文章解析アクセラレータのアプローチにしたがって、形態素抽出処理以外についても、ハードウェア化を検討してゆく予定である。文章解析の枠組みは長期的な研究課題であり、その途上においては、ソフトウェアのもつ柔軟性は必要である。しかし、今日、手法がある程度確立された部分もあり、それらについ

ては、積極的なハードウェア化を考えてゆく。さらに、これまでのソフトウェアの枠組みのハードウェア化にとどまらず、専用ハードウェアの可能性を、文章解析の新しい枠組みに結び付けることも考えてゆきたい。

参考文献

- [1] 福島・他、文章解析アクセラレータ(1)－形態素抽出マシンの試作－、情処39全大、2F-5、1989年。
- [2] 福島・他、多重照合型形態素抽出方式に関する検討、情処39全大、1F-5、1989年。
- [3] 浜口・他、音声日本語入力システムにおける高速な言語処理のための辞書照合アルゴリズム、信学論、J70-D(8)、1987年。
- [4] 中村・他、形態素抽出アルゴリズムの高速処理方式、情処37全大、5B-9、1988年。
- [5] 松本、論理文法の並列構文解析、情処論、29(4)、1988年。
- [6] W.Rytter, Parallel Time $O(\log n)$ Recognition of Unambiguous Context-free Languages, Info. and Comput., 73, 1987.
- [7] 峰・他、文脈自由文法の並列構文解析、情処研報、NL-73-1、1989年。
- [8] K.H.Chu, et al., VLSI Architectures for High Speed Recognition of Context-free Languages, 9th Annu. Int. Sympo. Comput. Arch., 1982.
- [9] Y.T.Chiang, et al., Parallel Parsing Algorithms and VLSI Implementation for Syntactic Pattern Recognition, IEEE Trans., PAMI-6(3), 1984.
- [10] 小倉・他、連想メモリLSIの現状と今後、信学誌、69(7)、1986年。
- [11] 山田・他、文字列検索LSI、信学技報、CAS87-25、1987年。
- [12] D.E.Knuth, The Art of Computer Programming, Vol.3: Sorting and Searching, Addison-Wesley, 1973.
- [13] O.Amble, et al., Ordered Hash Table, Comput.J., 17(2), 1974.
- [14] 横山・他、二次記憶上的大規模語彙を用いる自然言語処理システム、情処論、29(6)、1988年。