

# 音声認識とDP

迫江 博昭

## 1. まえがき

ダイナミックプログラミング(DP)が音声認識の分野で重用されるようになって15年になる。この間、多数の研究がなされ、幾多の論文が発表され、DPは音声認識の基本的な手法としての地位を確立した。今日では脱DPを模索するのが、一部研究者の命題となっているほど、この分野でのDPの存在は大きい。しかし、研究の中心となった人たちがORとは縁が遠かったため、DPそのものの研究グループには、音声認識とDPとのかかわりは、ほとんど知られていなかったようである。

著者自身も、DPの入門書を一読して、その考え方を自分にわかる範囲で把握したあとは、もっぱらその応用を考えるのみで、DPそのものの理解を深めるための努力はしなかった。元祖 R. Bellman 氏に対しても、伝説上の人物といった認識しかなく、発表した大半の論文でも氏の著作を引用するのを省略してしまっている。DPはそれほど有名であったというのが、故人となられたベルマン氏に対する言いわけである。

本稿では音声認識に対するDPの応用を解説する。筆者はこの研究の歴史のはじまりから深く関係し、現役研究者としての一生をこのテーマにさげたものである。このため解説に筆者の主観が

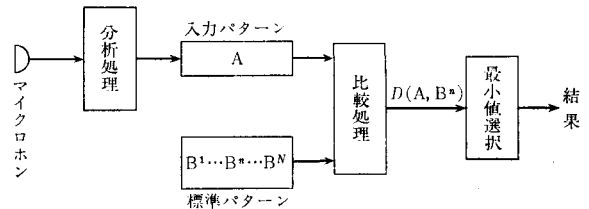


図1 パターンマッチングによる音声認識

相当程度入ることになると思う。また、いくつかの術語に、DPの教科書とは異なったものを用いる。これらは音声認識の分野の方言として定着してしまったものである。

## 2. 音声認識入門

音声認識のモデルとしては種々のものが提案されているが、最も基本的なものを図1に示す。パターンマッチング法にもとづくものである。入力音声波形は分析部で周波数分析され、短時間スペクトラムの時系列パターンに変換される。図2に数字“3”のパターンの例を示す。数式的には

$$A = a_1 a_2 \cdots a_i \cdots a_l \quad (1)$$

と、ベクトルの時系列として表現する。 $a_i$ は時刻*i*でのスペクトラム特徴を示すベクトルである。

単語名を番号*n*で示すこととし $\{n | n = 1, 2, \dots, N\}$ なる*N*個の単語セットを考える。単語*n*には標準パターン $\{B^1, B^2, \dots, B^n, \dots, B^N\}$ が用意されている。これらは(1)のAと同じくベクトルの時系列である。代表的なものを、添字を省略して

$$B = b_1 b_2 \cdots b_j \cdots b_l \quad (2)$$

さこえ ひろあき 日本電気㈱ C&Cシステム研究所

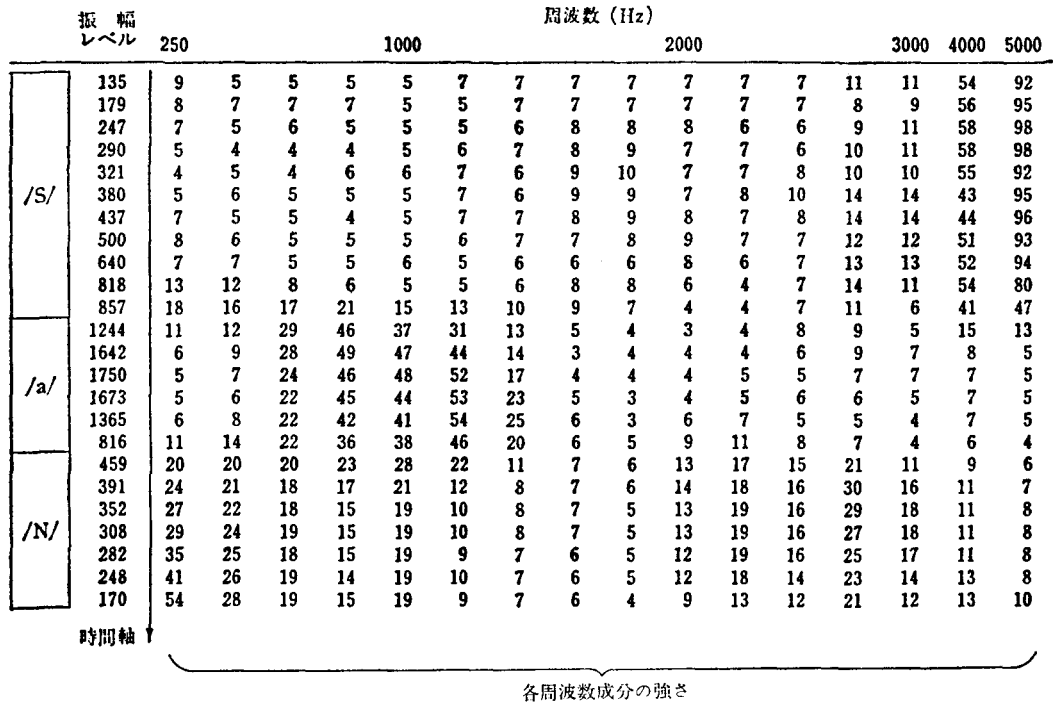


図 2 音声パターンの例(数字“3”)

と記す。

未知の入力パターンAが与えられると、これらの標準パターンとのあいだで比較処理(マッチング)が行なわれ、相異の度合いを示す尺度として、距離  $D(A, B^n)$  が算出される。これらの距離が最小となる単語名  $n=\hat{n}$  が認識結果として出力される。

以上は最も簡単な離散単語の認識の場合で、入力パターンも標準パターンも孤立発声した単語パターンとして扱われる。利用する各人の音声パターンを標準パターンとして登録する型のものを特定話者型と呼ぶ。この型のものでは、本人の標準パターンを用いることによって声質の個人差の影響を避けるのである。不特定の話者を対象とする型では、多数の人の音声パターンの中から、代表的なパターン(通常各単語ごとに複数個)をクラスタリング手法で選び、標準パターンとする。

2個以上の単語を連続発声したものの認識のためには、図1のモデルを若干修正する必要がある。

これについては、4章で説明する。

### 3. 離散単語認識における DP

#### 3.1 音声パターンの時間軸歪

パターンの比較をどう行なうかは、音声認識の性能を大きく左右する。音声パターンに生ずる各種の変動を十分考慮した方式を用いないと良好な精度は期待できない。主要な変動要因として時間軸歪がある。同一人が同じ単語を発声しても、まったく同一の速度で発声するとは限らない。このため、図3に示すように、入力パターンAと標準パターンBとのあいだに歪が発生する。このため時間的に対応するベクトル  $a_i$  と  $b_i$  との距離の総和を求めるといった比較方法は、きわめて根拠のない方法となる。以前は線形な伸縮によって長さを揃えてからマッチングするという方法(線形伸縮法)が試みられたが不十分であった。音声パターン中の各部は独立に伸縮するので全体としては複雑な非線形伸縮となる。線形伸縮は1次近似に

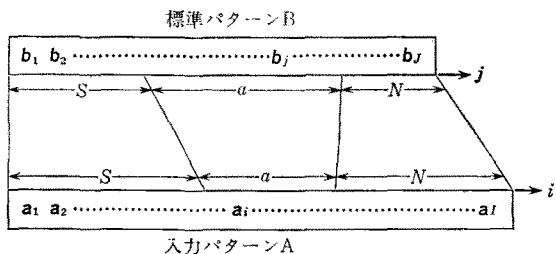


図3 音声パターンの時間軸歪

しかなかったのである。

### 3.2 DPマッチング

時間軸歪の問題をDPの導入によって解決できた。アルゴリズムには、いくつかの分化があるがDPマッチングと総称されている。

図4は入力パターンAを横軸に、標準パターンBを縦軸に配したものである。この図で、標準パターンの時間軸  $j$  を関数  $j(i)$  によって入力パターンAの時間軸  $i$  に対応づけ、時間軸歪を除くことを考える。

関数  $j(i)$  をどういう基準で定めるとよいのであろうか。そこで、時間軸歪がマッチングにおよぼす影響を定性的に考えてみる。AとBとがまったく同じパターンであるならばその間の距離は0である。時間軸歪が生じ、その度合いが大きくなると、両者のあいだの距離も大となる。したがってAとBとの距離が0となるように写像  $j(i)$  を定めればよい。一般には、AとBのあいだには時間軸歪以外の他の変動要因による歪も入っているので、距離0とはならないが、距離を最小にすると時間軸歪を除去できたと考えてよい。すなわち

$$D(A, B) = \min_{j=j(i)} \left[ \sum_{i=1}^I d(i, j) \right] \quad (3)$$

なる最小化問題として定式化できる。ここに  $d(i, j) = \|a_i - b_j\|$  で、ベクトル  $a_i$  と  $b_j$  との距離である。

図4と(3)式を並べると、おのずからDPの適用が連想される。すなわち図4の各  $(i, j)$  座標に距離  $d(i, j)$  が対応して

いるとき、それらの総和が最小になるようにするという最小コスト問題となっているのである。

具体的にDPを適用する前に、写像の関数  $j(i)$  の性質について考えてみる。ここでは一時、時間軸  $i, j$  は連続的と考えて説明する。

- 1)  $j(i)$ は連続関数。(時間から時間への写像だから当然)
- 2)  $j(i)$ は単調増加。(現象の前後関係を保つため)
- 3)  $j(1)=1, j(I)=J$ 。(始端や終端が切れてしまわないため)
- 4)  $j(i) \sim i$ 。(実際には時間軸歪があまり極端になることはないという事実に対応)

これらの性質が(3)式の最小化問題に対する制約条件となる。ただし、実際には離散系であるので、これらの条件は一部近似的に実現される。

次節に入る前に次の点を注意しておく。パターンマッチングの目的には(3)式の最小値  $D(A, B)$  が求められればよいので、最適な写像  $j(i)$  の決定は省略する。 $D(A, B)$ はAとBとのあいだの時間軸歪を除去しても、なお残る距離であって、時間正規化距離と呼ばれることもある。(一般的に言っ

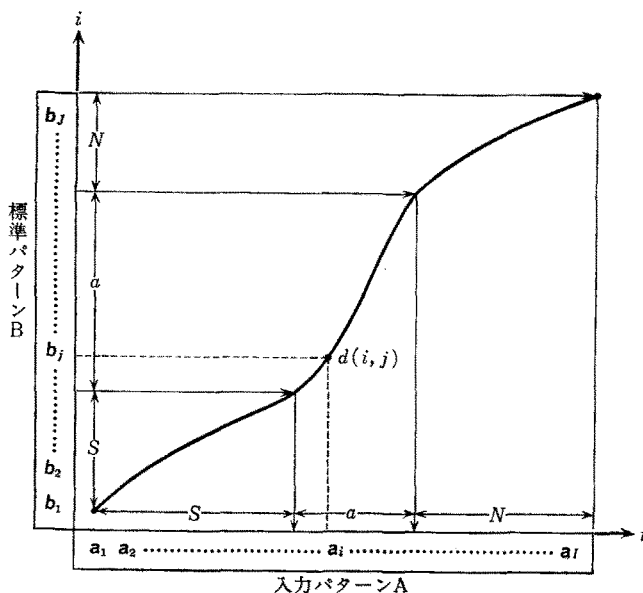


図4 最適経路問題への変換

て距離の公理を満足するものではないが、慣用的に距離と呼んでいる)

### 3.3 実行アルゴリズム

前節の(3)式の最小化問題を解くのであるが、条件1), 2), 3), 4)の実現の仕方によって種々の実行アルゴリズムが分化している。以下には最も簡単なものを示す。

図5の $(i, j)$ 格子点で次のDP計算を行なり。

○初期条件

$$g(1, 1) = d(1, 1) \quad (4)$$

○漸化式

$$g(i, j) = d(i, j) + \min \begin{bmatrix} g(i-1, j) \\ g(i-1, j-1) \\ g(i-1, j-2) \end{bmatrix} \quad (5)$$

○整合窓条件

$$j-r \leq i \leq j+r \quad (6)$$

これによって距離((3)式の最小値)は

○距離

$$D(A, B) = g(I, J) \quad (7)$$

と定まる。

漸化式(5)は時刻 $i$ における点 $(i, j)$ に至る経路として、図5中に示す3種を許すという形になっている。これは前節の1), 2)の性質を離散座標系で近似的に実現したものである。3)の性質は(4)式と(7)式で実現されている。4)は整合窓条件(6)式によって0次近似として実現されている。

(6)式の条件は、制約条件によって問題が簡化されるというDPの好ましい性質にマッチしたものになっている。漸化式(5)の計算を、図5の整合窓の範囲内に限定し、処理量を低減している。

表1 DPマッチングの効果

方法	誤認識
線形伸縮法	5.9%
DPマッチング	0.8%

認識対象は地名50語

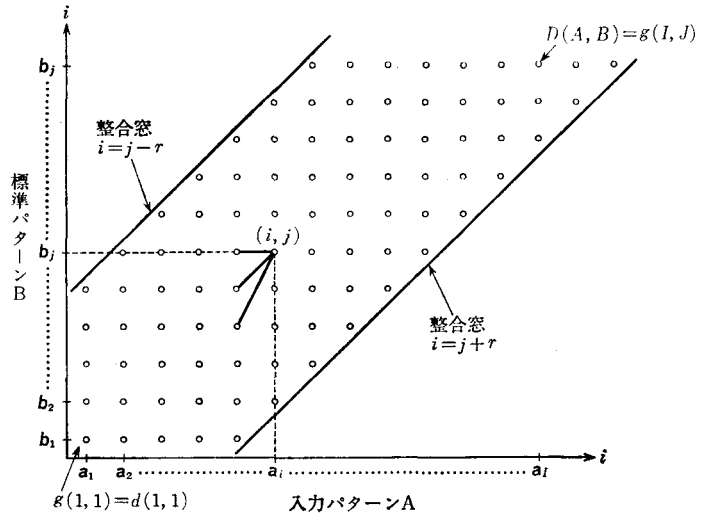


図5 DPの計算

また音声認識の分野では、音声パターンの終端 $i=I$ が不明な場合がある。この場合には、終端開放の条件で(5)式を計算すると

$$J-r \leq I \leq J+r$$

の範囲の各終端候補 $I$ に対して(7)式により距離 $D(A, B)$ が逆列に評価できる。これもDPならの性質であるが、後に述べる連続単語認識などで、有効なものである。

以上のようなDPの適用によって、非線形な時間軸歪み問題が実用レベルとして解決された(1970)。表1にDPマッチングと、代表的な従来手法である線形伸縮法との比較実験の結果を示す。誤認識が約1/7に低減されるという効果が得られている。

図6にDP漸化式の変形の例を示す。(a)が図5(すなわち(5)式)のものである。(c)は写像 $j'(i)$ の傾斜を近似的に $1/2 \leq j' \leq 2$ と制限した型であり、前節の条件4)を1次近似まで実現したものである。この型のものが、認識精度の点で最も良いことが実験的に確認されている。表1のデータは、この型を用いた場合のものである。

以上で離散単語認識のDPマッチングの説明を終る。不可能を可能にしたと言えるほどDPの効果は顕著であった。(3)式あるいは図4の最適経路問題に対するDPの効用は多くの読者に自明の

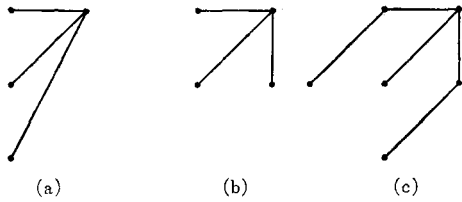


図 6 経路条件の変形

ことと思う。DPによって導びかれた(5)式の漸化式には距離  $d(i, j)$  の計算, 最小値選択, 加算, と3種の処理が含まれているが, 計算量的にはベクトル  $a_i$  と  $b_j$  の距離を計算するための処理が大半を占める。すなわち, DPの効用があまりにも大きいので, DPマッチングのための計算は, ベクトル間距離計算という, DPの責任外の処理のために, ほとんどが費やされるのである。

#### 4. 連続単語認識におけるDP

いくつかの単語を連続して発声したものを認識しようとする場合にはセグメンテーションの問題が発生する。連続した単語と単語の境界が不明なので, これを決定しないと前章のような単語単位でのパターンマッチングが適用できない。このためセグメンテーションに関する研究が数多くなされたが, 境界そのものがきわめて曖昧模範なものであるため, 信頼性のある手法は確立されていない。極端には“21”( /ni itʃi/ )のように同じ音  $i$  がつづいて, 境界決定が原理上不可能な場合もあるのである。

このため, 図7に示すように, 入力パターンを分解するかわりに, 単語標準パターンを接続して連続標準パターンを合成し, パターンマッチングを行なうという原理を採用した。4桁の連続数字(棒読み)の例で説明すると, まず“0”~“9”の標準パターン  $B^0 \sim B^9$  を用意する。入力音声パターンAとして, 連続数字“5892”(実際には未知)が入力されたとする。

4桁の数字“0000”から“9999”までの数字(一般に  $n(1), n(2), n(3), n(4)$  とする)の連続標準パターン  $B^{n(1)} \oplus B^{n(2)} \oplus B^{n(3)} \oplus B^{n(4)}$  を合成して入力パターンAとDPマッチングを行なう。距離  $D(A, B^{n(1)} \oplus B^{n(2)} \oplus B^{n(3)} \oplus B^{n(4)})$  をそれぞれ求め, 最小となる  $n(1), n(2), n(3), n(4)$  を決定し, これを認識結果とする。一般的には次の最小化問題となる。

$$\min_{\{n(x)\}} [D(A, B^{n(1)} \oplus \dots \oplus B^{n(x)} \oplus \dots \oplus B^{n(k)})] \quad (8)$$

この原理によると確かにセグメンテーションの問題を回避できる。しかし, 新たに順列  $\{n(x)\}$  に関する最小化という問題が生じた(距離  $D$  を算出するための写像関数  $j(i)$  に関する最小化が埋め込まれているが, これは前章で解決済みである)。この最小化問題を総当り法で解くと,  $k$  数字連続の場合,  $10^k$  種の連続標準パターンに関してDPマッチングのくりかえしが必要となる。そこで再びDPの適用を試みる。

入力パターンAの  $i=l+1$  から  $i=m$  までの部分を部分パターン  $A(l, m) = a_{l+1} \dots a_m$  と定義する。(8)式を解くための補助変数として, 入力パターンAの中に  $(k-1)$  個の区分点  $l(1) \dots l(x) \dots l(k-1)$  を仮定する。これによって入力パターンAは  $A = A(l(0), l(1)) \dots \oplus A(l(x-1), l(x)) \dots \oplus A(l(k-1), l(k))$  と,  $k$  個の部分パタ

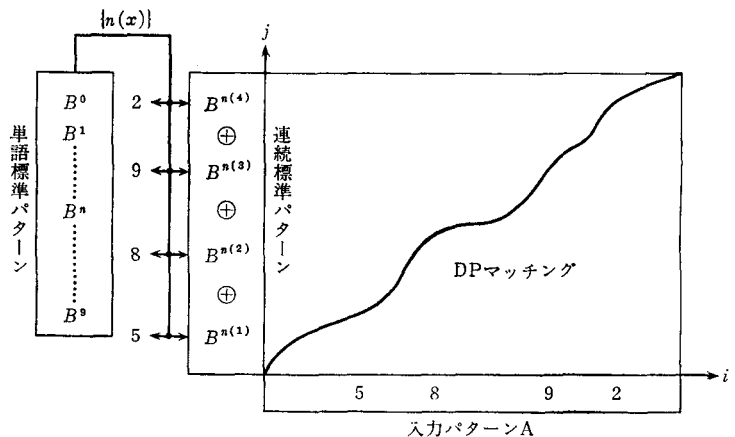


図 7 連続標準パターンの合成とDPマッチングによる連続単語認識の原理

ーの接続として表現される。ここに  $l(0)=0, l(k)=I$ 。これを(8)式に代入して整理すると

$$\min_{\{l(x)\}} \left[ \sum_{x=1}^k \min_{n(x)} [D(A(l(x-1)), l(x)), B^{n(x)})] \right] \quad (9)$$

と2段階の最小化に分解される。内側の最小化は  $n(x)$  に関する単純比較なので、一般的に次のように計算できる。

$$\hat{D}(l, m) = \min_n [D(A(l, m), B^n)] \quad (10)$$

この最小値を与える単語名

$$\hat{N}(l, m) = \operatorname{argmin}_n [D(A(l, m), B^n)] \quad (11)$$

は、部分パターン  $A(l(x-1), l(x))$  が1個の単語であると仮定して認識したときの結果である。

(10)式を用いると(9)式は

$$\min_{\{l(x)\}} \left[ \sum_{x=1}^k \hat{D}(l(x-1), l(x)) \right] \quad (12)$$

となるので、次のようなDPによって計算できる。

○初期条件

$$T(0, 0) = 0$$

○漸化式

$$T(x, m) = \min_{l < m} [\hat{D}(l, m) + T(x-1, l)] \quad (13)$$

この(13)式を  $x=5, m=I$  まで計算する。この間、(13)式の  $l$  の最適値を  $x, m$  に対応つけて記憶しておくバックトラックによって(12)式的最適解

$$\{\hat{l}(x) | x=0, 1, \dots, k\}$$

が定まる。これをもとに(11)式の  $\hat{N}$  のテーブルを参照することにより

$$\{\hat{N}(\hat{l}(x-1), \hat{l}(x))\}$$

として認識結果が得られる。

表2 2段DPマッチングによる連続数字認識実験結果

連続桁数		1	2	3	4	認識率
発 声 音	A	0	0	0	0	100%
	B	0	0	1	0	99.8%
	C	0	0	0	1	99.8%
	D	0	2	0	2	99.2%
	E	0	0	0	5	99.0%
平均認識率		100%	99.6%	99.9%	99.2%	99.6%

テストデータは1~4桁数字を100種ずつ発声したもの

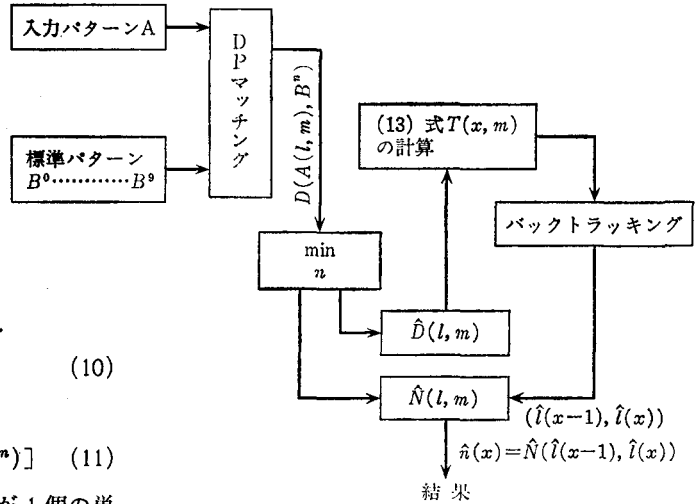


図8 2段DPマッチングのブロック図

以上の処理のブロック図を図8に示す。この方法では(10)式等の距離  $D(A(l, m), B^n)$  を計算するために前章で述べた単語単位でのDPが必要とされ、これに加えて(13)式のDPが計算される。これにちなんで2段DPマッチング法と呼んでいる。この方法では3.3節のDPの並列性が有効に利用できる。(10), (11)の計算のための距離  $D(A(l, m), B^n)$  は各  $(l, m)$  に対して独立に算出する必要はない。始端  $l$  を定めてDPマッチングを行なうと、整合窓内の  $m$  に対して並列的に  $D(A(l, m), B^n)$  が得られるからである。

2段DPマッチング法によって、連続単語認識がはじめて可能になった(1975)。表2に連続数字認識に適用した実験結果を示す。その後、種々の改良が行なわれたが、主として計算効率の向上を目的としたもので、連続標準パターンとのマッチングをDPによって処理するという原理は今も生きている。

## 5. ハードウェア, LSI

DPの処理を実行するための専用のハードウェアが作られ、LSI化まで実現しているのは音声認識の分野だけではなからうか? 実用音声認識装置となると、コンピュータ・シミュレーションと異なって、実時間動作が必須である。このため、

やや強引にハードウェア化が進められたとも言える。

1978年に発売され、世界で初のDP式音声認識装置となった日電のDP-100には、約400個のICを使ったDPプロセッサが使用された。その後、半導体技術が進歩し、現在では1個のLSIではほぼ同等の機能が実現されている。図9はその拡大写真である。マイクロプログラム制御を採用したので、DP特有の回路構成といった特徴は見えないが、マイクロインストラクション等にDP向けの工夫がなされている。

## 6. あとがき

以上、音声認識の分野におけるDPの応用の一端を紹介した。ここで述べたものは単語を単位としてパターンマッチングを行なう方式で技術的に完成度が高く、実用化されているものである。今後の方向として、大語彙化をめざして音素や音節を単位とした認識が必要と考えられ、研究が進められている。この中でもDPは基本的手法の1つとして位置づけられている。

音声認識の他に文字認識においてもDPが利用されている。これらを含めて、日経エレクトロニクス誌(1983年11月7日)に解説しているので興味をおもちの方は参照されたい。

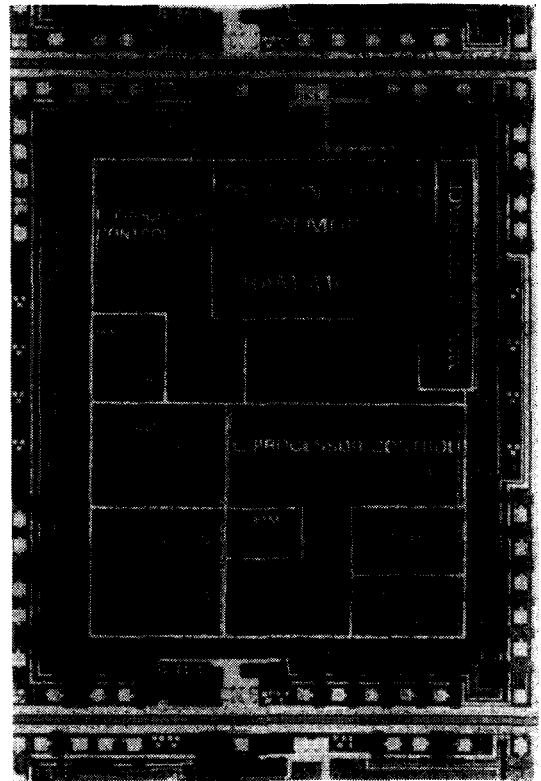


図9 DPマッチングLSI

1970年以来、音声認識の分野でDPが果たしてきた役割は非常に大きい。特に実用面での評価は高く、筆者らが何となく用いた“DPマッチング”という技術名が、多くの研究論文はもとより、各社の製品カタログでも技術的な裏づけを示すためのPR用語として、引用されている。筆者を含む多くの研究者が、DPによって大きな仕事をする事ができた。故ベルマン教授に誌上を借りて心からお礼を述べたい。

### 次号予告

#### 特集 事例研究——59年秋季研究発表会より

北洋漁業の地域経済におよぼす効果に関する研究

伊藤昭男・阿部秀明・佐藤博樹

バスパンチング発生要因に関する調査研究

定方希夫

ポンプ・ステーションの最適計画

石堂一成・南部和幸

配電系統の最適供給計画

青木兼一・一森哲男

一般ネットワークにおける複数施設の配置問題

川中子敬至・山城光雄・矢部 眞