

TR-I-0145

TDNNの構造の音韻認識率、  
シフトイバリアント性への影響

南 泰浩 † 沢井 秀文

Yasuhiro MINAMI † Hidefumi SAWAI

1990.2

概要

TDNNを連続音声認識に用いる場合、TDNNの連続音声に対する音韻認識率の向上が非常に重要となってきた。本報告ではTDNNの連続音声に対する音韻認識率をHMMと比較し、TDNNの連続音声に対する性能を調べた。この結果、TDNNはカテゴリ数の大きな連続音声の認識に対して、HMMに比べて弱いことが確認された。本報告ではさらに、TDNNの構造の変化による音韻認識率とシフトイバリアント性の比較を行い、TDNNの連続音声に対する音韻認識率の向上の可能性を示す。

本報告は、学外実習生南泰浩(慶応大学)が行った実習の報告書である。

ATR Interpreting Telephony Research Laboratories  
ATR 自動翻訳電話研究所

© ATR Interpreting Telephony Research Laboratories

© ATR 自動翻訳電話研究所

† Keio Univ.

† 慶応義塾大学

## 1. はじめに

A T R では従来より、T D N N ( T i m e - D e l a y N e u r a l N e t w o r k s ) を用いた音声認識を行ってきた<sup>(1)</sup>。T D N N は特定話者の / b d g / のタスクに対して 98.6% と高精度の音韻認識率を達成することが確認された<sup>(1)</sup>。しかし、T D N N を用いた連続音声認識システムを構築する過程で様々な問題が生じてきた。我々が構成したシステムは従来 A T R において構成した連続音声認識用システムの H M M ( H i d d e n M a r k o v M o d e l ) - L R の音声処理部 ( H M M 部 ) を T D N N の処理部に置き換えたものである<sup>(2)</sup>。このシステムを用いて連続音声認識を行った結果、H M M - L R に比べ文節認識率が低いことが確認された。この原因は T D N N の連続音声に対する音韻認識率の低下であると考えられる。

そこで本報告では 2 章で T D N N の文節発声データに対する音韻認識率を H M M と比較することによって、T D N N の連続音声に対する性能を調べる。そして、3 章、4 章で T D N N の構造による音韻認識や、ソフトインバリエント性を調べ、T D N N の改良の可能性を調べる。

## 2. T D N N と H M M の音韻認識率の比較

まず、T D N N が連続音声に対してどのような音韻認識率を示すかを H M M と比較する。以下に実験の条件を示す。3 章以降で述べる実験結果はすべてこの章で述べる実験条件と同じものを使用した。

以下に比較実験を行うタスク、認識に用いた音韻、学習に用いた音韻について述べる。

## 2. 1 比較対象タスク

比較に用いたタスクは以下の小カテゴリ（3カテゴリ）と中カテゴリ（6カテゴリ）と大カテゴリ（18カテゴリ）の3種類である。

(1) / b d g /

(2) / b d g m n N /

(3) 全子音（18子音）

18子音は / b / , / d / , / g / , / p / , / t / , / k / , / m / , / n / , / N / , / s / , / s h / , / h / , / z / , / c h / , / t s / , / r / , / w / , / y / であり音韻区間は以下の条件を満たすもののみを用いた。

(1) / c h / は / c h i / のみから切り出す。

(2) / s h / は / s h i / のみから切り出す。

(3) / z / は / z a / , / z e / , / z u / , / z o / のみから切り出す。

(4) / h / は / h a / , / h i / , / h e / , / h o / のみから切り出す。

(5) その他はラベルに従いその音韻の存在する区間を切り出す。

以上のラベルの表記はATRデータベースの表記法<sup>(3)</sup>に従い、音韻はこのラベルにより切り出したものである。

## 2. 2 認識対象音韻

認識データとしては発声様式の異なった以下の4種類のデータから切り出した音韻を用いる。認識に用いた音韻の数は全ての手法で同じ数とした。

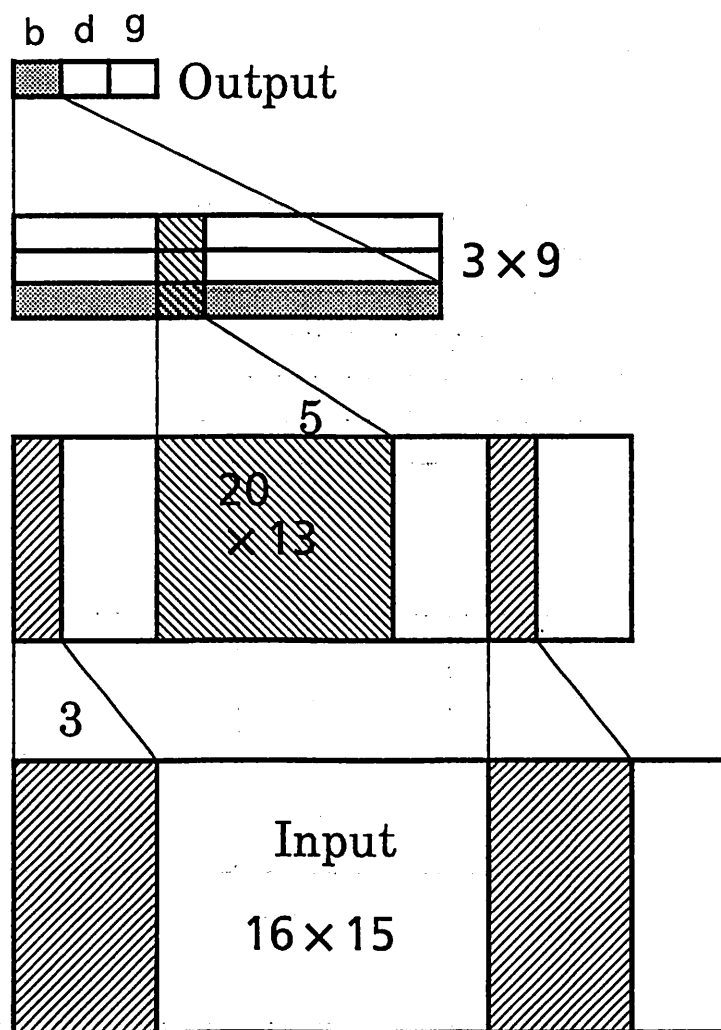


図1 /bdg/用TDNN

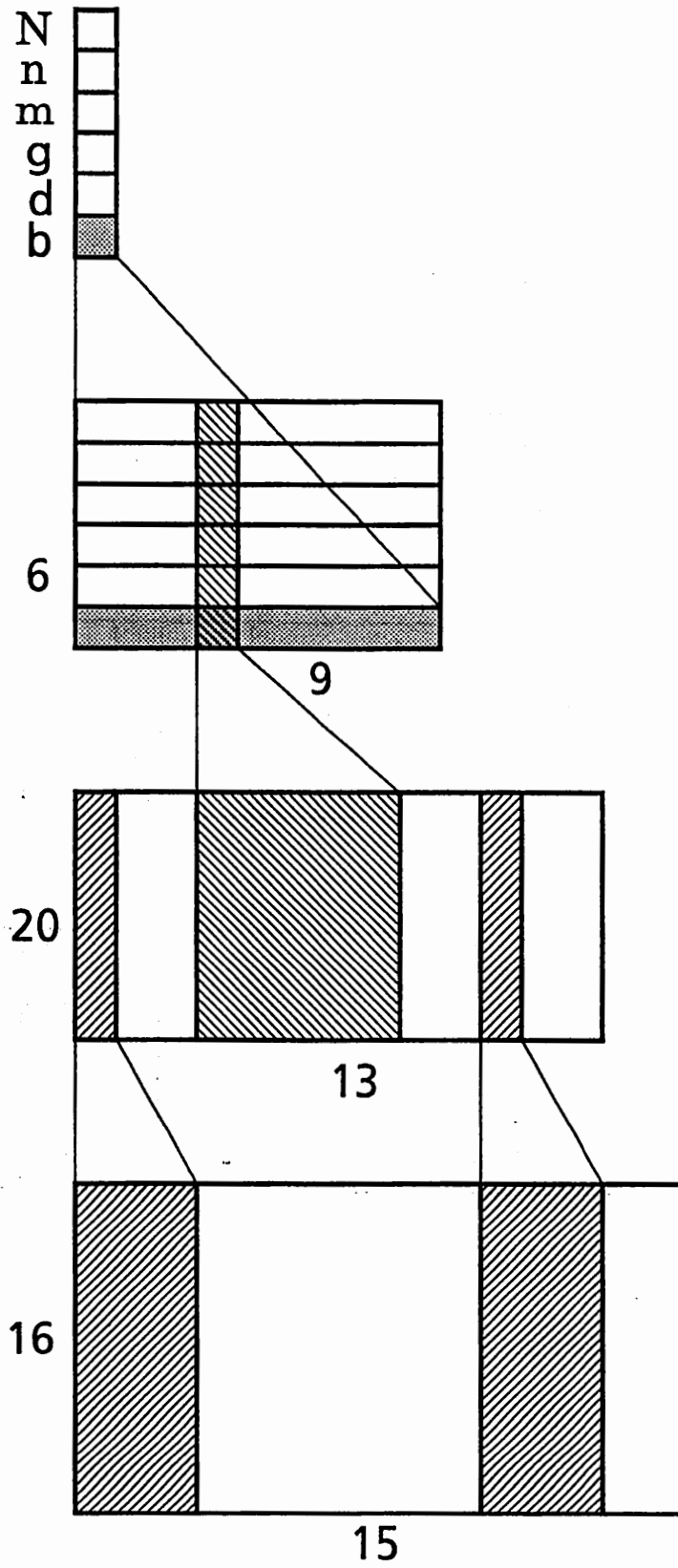


图2 /bdgmnN/用TDNN

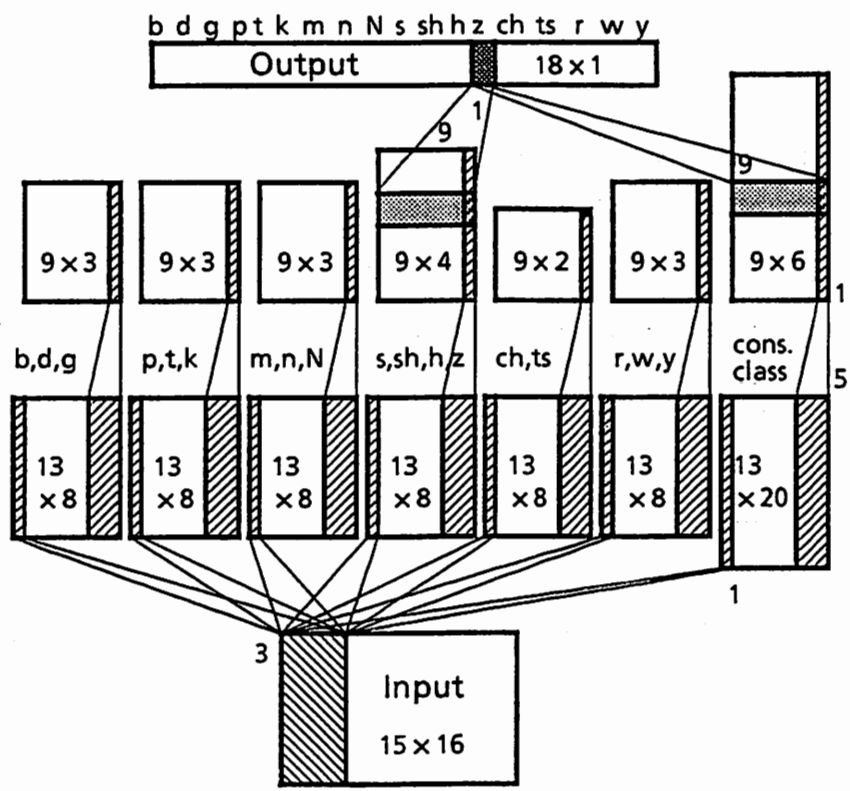


图3 18子音用TDNN

( 1 ) 5 2 4 0 単語の奇数番目のデータ ( 以下単語発声とする ) ( 5 . 6 8 モーラ / s e c )

( 2 ) 複合語を許さない文節発声 ( 以下短い文節発声とする ) ( 7 . 1 4 モーラ / s e c )

( 3 ) 文節発声 ( 7 . 7 2 モーラ / s e c )

( 4 ) 句切り指定を行わない発声 ( 以下自由発声とする ) ( 9 . 5 6 モーラ / s e c )

## 2 . 3 学習音韻

学習音韻は重要5240単語の偶数番目のデータから切り出した音韻を用いた。単語発声から切り出したデータだけで、どれだけ連続音声認識に対して有効であるかを調べるため学習データには文節発声や自由発声データを用いていない。

## 2 . 4 発声話者

発声はアナウンサー1名の発声した音声で特定話者の認識を行った。

## 2 . 5 認識率の計算

ここで求めた認識率は各音韻毎に認識率を求め、その値を各音韻の個数に対応する重みを掛けて総和したものを全体の認識率とした。

## 2 . 6 T D N N

T D N N を用いた / b d g / を認識するアーキテク

チャは図1のようになる<sup>(1)</sup>。また、 / b d g m n N / を識別するアーキテクチャを図2に示す(1)。

18子音に対しては図3に示すアーキテクチャを用いる<sup>(4)</sup>。これは図1に示したネットワークを多数カテゴリを識別するように構成したものである。このアーキテクチャは18音韻(18子音 / b / , / d / , / g / , / p / , / t / , / k / , / m / , / n / , / N / , / s / , / s h / , / h / , / z / , / c h / , / t s / , / r / , / w / , / y / )を識別する。このTDNNは6個の子音グループ内を識別するサブネットワーク" b d g " , " p t k " , " m n N " , " s s h h z " , " c h t s " , " r w y " と子音グループ間を識別するネットワーク" c - c l a s s " のモジュールで構成される。以上で述べた3つのアーキテクチャをBack-Propagation学習則<sup>(5)</sup>を用いて学習することにより、他カテゴリの出力値を考慮(側抑制)するようなネットワークを構成することができる。この際、同じ値を持つように設定された重みは、この設定を崩さないように学習される。18子音の学習は各モジュールに分割することなく、一括に学習した。

## 2. 7 TDNNにおける分析

TDNNで用いる1フレームの音声データは12KHzでサンプリングし、256ポイントのハミング窓で、FFTを計算した後、10ms毎に16次のメル尺度のフィルタバンクを通し、15フレーム内で平均0、±1に正規化したものである。

## 2. 8 複数コードブックを用いるHMM



本実験ではHMMとして複数のコードブックを利用するHMMを用いる。

ここで用いたパラメータは以下の3種類がある。

(1) スペクトル WLR 距離尺度による256個のベクトルからなるコードブック。

(2) パワー PWR のパワー項の距離尺度による64個のスカラからなるコードブック。

(3) スペクトルの動的特徴 LPC ケプストラム係数の1次から16次までの時間方向の線形回帰係数、のユークリッド距離による256個のベクトルからなるコードブック。

## 2. 9 HMM における分析

分析は音声データを12kHzでサンプリング後、窓長21.3msのハミング窓で、フレーム周期3msごとに切り出し、高域強調後12次のLPC分析を行った。

## 2. 10 比較結果

表1から表12に比較結果を示す。

## 2. 11 比較実験に対する検討

HMMとTDNNでは大カテゴリーの識別を行ったときに、音韻認識率の差が大きく現れる。この中で最も特徴的なのはTDNNが自由発声に発声様式が近づくにつれて認識率が極端に低下することである。

表1 単語発声データに対する累積音韻認識率(/bdg/)

順位	1	2	3
HMM	95.6%	99.4%	100%
TDNN	98.8%	99.9%	100%

表2 文節発声データに対する累積音韻認識率(/bdg/)

順位	1	2	3
HMM	86.8%	97.8%	100%
TDNN	90.8%	94.0%	100%

表3 短い文節発声データに対する累積音韻認識率(/bdg/)

順位	1	2	3
HMM	87.5%	96.8%	100%
TDNN	91.3%	95.0%	100%

表4 自由発声データに対する累積音韻認識率(/bdg/)

順位	1	2	3
HMM	80.8%	92.8%	100%
TDNN	85.3%	91.8%	100%

表5 単語発声データに対する累積音韻認識率(/bdgmnN/)

順位	1	2	3	4	5
HMM	89.5%	98.3%	99.5%	99.8%	99.9%
TDNN	95.1%	98.6%	99.7%	99.8%	100%

表6 文節発声データに対する累積音韻認識率(/bdgmnN/)

順位	1	2	3	4	5
HMM	73.7%	90.1%	97.0%	98.4%	99.6%
TDNN	72.5%	88.8%	94.6%	98.0%	99.2%

表7 短い文節発声データに対する累積音韻認識率(/bdg/mnN)

順位	1	2	3	4	5
HMM	75.4%	89.6%	95.5%	98.3%	99.3%
TDNN	75.7%	89.3%	95.2%	97.2%	99.1%

表8 自由発声データに対する累積音韻認識率(/bdgmnN/)

順位	1	2	3	4	5
HMM	64.6%	83.5%	91.8%	96.6%	99.4%
TDNN	64.7%	83.2%	90.4%	94.5%	97.9%

表9 単語発声データに対する累積音韻認識率(18子音)

順位	1	2	3	4	5
HMM	93.1%	98.4%	99.5%	99.7%	99.8%
TDNN	96.2%	98.9%	99.6%	99.8%	99.9%

表10 文節発声データに対する累積音韻認識率(18子音)

順位	1	2	3	4	5
HMM	79.8%	92.2%	96.6%	97.9%	99.0%
TDNN	72.4%	84.3%	90.5%	93.1%	94.9%

表11 短い文節発声データに対する累積音韻認識率(18子音)

順位	1	2	2	3	2
HMM	81.4%	91.7%	95.9%	98.0%	98.8%
TDNN	76.2%	86.3%	91.5%	94.7%	96.6%

表12 自由発声データに対する累積音韻認識率(18子音)

順位	1	2	2	3	2
HMM	71.6%	86.2%	92.4%	95.3%	97.2%
TDNN	56.6%	70.9%	78.5%	84.3%	87.8%

### 3. ネットワークの構造による

#### 認識率の違い（18子音用）

ネットワークの構造によってどのように認識率が変化するかを調べるために以下の基本的な2種類の実験を行った。ここで用いた音韻は2. で求めたものと同じものを用いた。

#### 3. 1 ネットワークのモジュール構成

T D N Nとしてどのような構成がよいのかを調べるために、以下の2つのT D N Nを用いて実験を行った。

##### (1) 分散型 T D N N

このネットワークを図4に示す。

##### (2) 集中型 T D N N

このネットワークを図5に示す。

表13に図3の従来型のT D N Nの各発声様式に対する累積音韻認識率を示す。以上の2つのT D N Nの各発声様式に対する累積認識率を表14と表15に示す。学習にはD Y N E T<sup>(6)</sup>を用いた。

#### ● 検討

分散型も集中型も図3のT D N Nに比べ認識率が低い。これは図3のクラス分けを行うネットワークが有効なのではないかと考えられる。

#### 3. 2 出力層と中間層の間の接続法

出力層と中間層の接続がどのようにシフトインバリアントに影響を及ぼすかを見極めるため出力層と中間

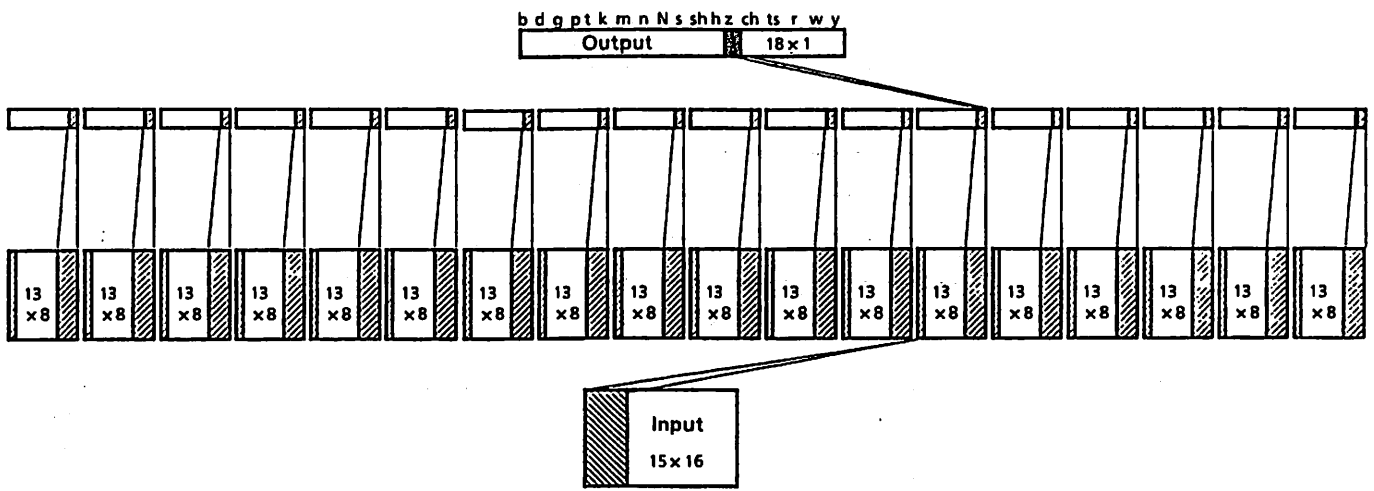


图 4 分散型18子音用TDNN

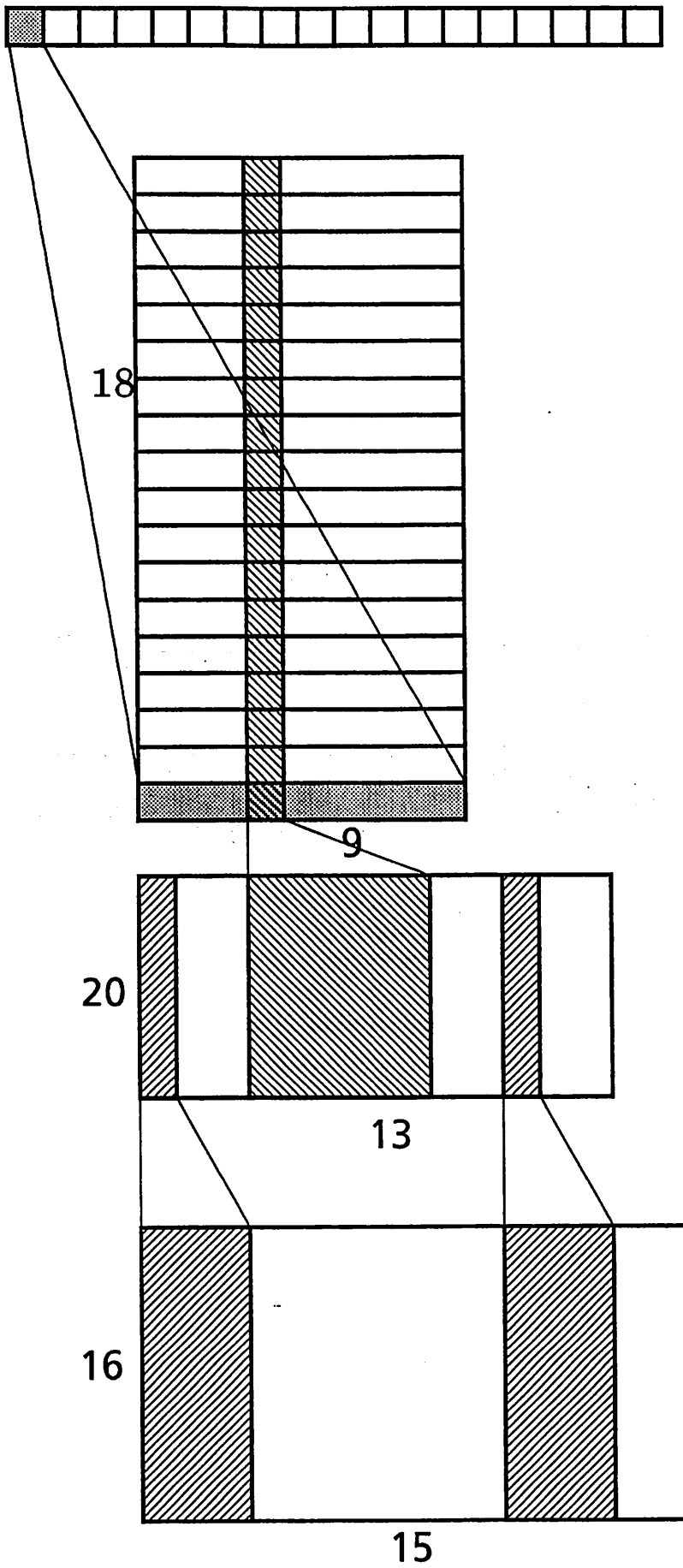


图5 集中型TD18子音用TDNN

表13 図3のTDNNの累積認識率

順位	1	2	3	4	5
単語発声	96.2%	98.9%	99.6%	99.8%	99.9%
文節発声	72.4%	84.3%	90.5%	93.1%	94.9%
短い文節発声	76.2%	86.3%	91.5%	94.7%	96.6%
自由発声	56.6%	70.9%	78.5%	84.3%	87.8%



表14 分散型TDNNの累積認識率

順位	1	2	3	4	5
単語発声	96.9%	99.2%	99.6%	99.8%	99.9%
文節発声	66.4%	82.3%	88.4%	91.5%	93.7%
短い文節発声	69.7%	83.8%	89.5%	92.8%	94.5%
自由発声	49.9%	67.9%	77.4%	83.0%	86.9%

表15 集約型TDNNの累積認識率

順位	1	2	3	4	5
単語発声	95.5%	98.7%	99.3%	99.6%	99.8%
文節発声					
短い文節発声	71.5%	84.5%	89.1%	92.3%	94.3%
自由発声	50.0%	65.9%	%	%	%

層との間の接続法を変化させたものを用いて、シフトインバリエントの効果調べた。ここでは 3. 2. 1, 3. 2. 2 に示す 2 種類の接続方法に対して評価を行った。この 2 つ接続法に対するシフトインバリエント性を図 3 の TDNN と比較するために、図 3 の TDNN のシフトインバリエント性について表 1 6 に示す。この表は TDNN の入力を 10 ms (1 フレーム) 毎にシフトし、結果を表にしたものである。認識に用いた音韻の発声様式は単語発声である。

### 3. 2. 1 自由な接続

図 3 に示す TDNN の中間層と出力層との接続はカテゴリ毎に全て同じ重みでつながっていた。ここではこの重みを全て自由にして、音韻の認識率とシフトインバリエント性を調べた。表 1 7, 1 8 に結果を示す。表 1 7 は各発声様式に対する音韻認識率を示し、表 1 8 は単語発声した音韻のシフトに対する認識率を示す。この学習には D Y N E T を用いた。

#### ● 検討

この表 1 7 と表 1 3 を比較するとかなりの音韻認識率の向上がみられる。しかし、表 1 8 と表 1 6 を比較するとシフトインバリエント性に対しては従来の TDNNの方がよいということが分かる。このことから従来の TDNN 出力層と中間層の間の結合法はシフトインバリエント性には有効であると考えられる。

### 3. 3. 2 分割接続

前節では全ての重みを自由にしていたが、シフトイ

表16 図のTDNNのシフトインバリエント性能(単語発声)

シフト幅(ms)	1	2	3	4	5
-30	61.2%	79.3%	88.0%	92.5%	95.1%
-20	82.8%	93.3%	96.7%	98.4%	99.2%
-10	93.2%	97.9%	99.2%	99.6%	99.8%
0	96.2%	98.9%	99.6%	99.8%	99.9%
10	93.0%	98.5%	99.4%	99.8%	99.8%
20	78.5%	92.0%	96.8%	98.7%	99.3%
30	60.0%	78.1%	87.2%	91.8%	94.3%

表17 出力層と中間層間の重みを自由にしたTDNNの累積認識率

順位	1	2	3	4	5
単語発声	95.7%	98.9%	99.6%	99.8%	99.8%
文節発声	78.7%	90.4%	94.1%	96.2%	97.2%
短い文節発声	82.2%	92.1%	94.9%	96.4%	97.7%
自由発声	65.2%	78.6%	84.3%	88.3%	90.9%

表18 出力層と中間層間の重みを自由にしたTDNNのシフトインバリエント性能(単語発声)

シフト幅(ms)	1	2	3	4	5
-30	33.5%	47.4%	56.4%	63.6%	69.5%
-20	52.9%	68.4%	76.4%	81.5%	85.5%
-10	84.7%	94.1%	96.7%	98.0%	98.8%
0	96.7%	99.1%	99.7%	99.8%	99.9%
10	78.6%	88.5%	92.8%	95.0%	96.6%
20	44.8%	57.3%	63.9%	68.8%	73.0%
30	28.9%	42.4%	48.5%	53.2%	58.2%

ンバリエーション性を高めるため、ここではこの重みを前部4個と後部の5個とに二分割、前部の重みは全て同じ値とし、後部の重みも全て同じ重みとした。表19, 20に結果を示す。表19は各発声様式に対する音韻認識率を示し、表20は単語発声の音韻の位置ずれ(シフト)に対する認識率を示す。この実験の学習にはDYNETを用いた。

#### ● 検討

表19は表17と表13からこの方式が従来型のTDNNと出力層と中間層の間の重みを自由にしたTDNNの中間程度の認識率を示すことが分かる。表20は表18と表16からシフトインバリエーション性についてもこの手法が2つのTDNNの間にあることが確認された。

#### 4. ネットワークの構造による

##### 認識率の違い ( / b d g m n N / )

18子音の学習は非常に多くの時間を必要とするのでここでは、6音韻 ( / b d g m n N / ) のTDNNを用いる。本章では、3. 章で述べた手法を考慮にいられて、様々な手法を提案し、その効果を調べる。ここで各手法に入る前に図2のTDNNの各発声様式に対する音韻認識率と、シフトインバリエーション性について表21と表22に示す。ここで用いた音韻は2. で求めたものと同じものを用いた。この実験の学習にはDYNETを用いた。

#### 4. 1 周波数方向の重みをぼかす手法

表19 出力層と中間層間の重みを中心に2つに分けたTDNNの累積認識率

順位	1	2	3	4	5
単語発声	96.7%	99.1%	99.7%	99.8%	99.9%
文節発声	77.9%	90.5%	94.3%	96.0%	96.9%
短い文節発声	80.8%	91.1%	95.0%	96.7%	98.2%
自由発声	62.0%	79.2%	85.2%	89.1%	92.0%

表20 出力層と中間層間の重みを中心に2つに分けたTDNNのシフトインバリエント性能(単語発声)

シフト幅(ms)	1	2	3	4	5
-30	46.2%	63.6%	73.6%	79.9%	84.2%
-20	75.7%	87.7%	92.4%	94.9%	96.4%
-10	93.0%	98.0%	99.1%	99.5%	99.7%
0	95.7%	98.9%	99.6%	99.8%	99.8%
10	92.0%	97.4%	98.7%	99.3%	99.7%
20	68.8%	81.5%	87.8%	91.9%	94.5%
30	38.5%	51.2%	60.4%	67.9%	73.7%

表21 図2のTDNNの累積認識率(/bdgmnN/)

順位	1	2	3
文節発声	72.5%	88.8%	94.6%
短い文節発声	75.7%	89.3%	95.2%
自由発声	64.7%	83.2%	90.4%

表22 図2のTDNNの  
シフトインバリエント性能(/bdgmnN/)(短い音節発声)

シフト幅(ms)	1	2	3
-30	44.9%	66.8%	86.7%
-20	60.3%	82.8%	92.6%
-10	72.5%	88.5%	94.7%
0	75.7%	89.3%	95.2%
10	74.1%	88.7%	94.4%
20	69.2%	86.1%	93.2%
30	56.6%	79.0%	89.7%



重みの周波数方向に対する感受性を軽減するために重みを学習させるときに周波数方向の隣の重みの変分量をこの重みに30%程度加えて学習させた。この効果によって周波数方向に対する位置ずれを吸収できると考えた。この結果を表23に示す。この実験の学習にはDCPを用いた。

● 検討

表23の結果と表21の結果を比較すると、この手法があまり有効でないことが分かる。原因としては、周波数領域がすでに16次元に情報が落とされてぼけているなどが挙げられる。

4. 2 出力層と中間層の間の重みを自由にする手法

3. 2で行ったように / b d g m n N / の識別に対しても出力層と中間層の接続がどのようにシフトインバリエントに影響を及ぼすかを見極めるため出力層と中間層との間の接続を自由にしたものを用いて音韻の認識と、シフトインバリエントの効果を調べた。この結果を表24、表25に示す。表24は各発声様式による、音韻認識率である。表25は各シフト(m s)に対する音韻認識率である。この実験の学習にはDYN E Tを用いた。

● 検討

表24、表25と表21と表22と比較すると次のことが分かる。18子音と同様この場合も音韻認識率はかなりよい値を示しているが、シフトインバリエント性はかなり悪化する。出力層と中間層の間の重みがシフトインバリエントに重要であることが確認された。

表23 周波数方向に重みをほかす  
TDNNの累積認識率(/bdgmnN/)

順位	1	2	3
文節発声	69.8%	87.4%	94.3%
短い文節発声	70.2%	85.7%	92.2%
自由発声	59.9%	82.0%	91.6%

表24 出力層と中間層間の重みを自由にした  
TDNNの累積認識率(/bdgmnN/)

順位	1	2	3
文節発声	85.1%	94.4%	96.9%
短い文節発声	83.1%	93.6%	97.0%
自由発声	80.8%	93.1%	97.5%

表25 出力層と中間層間の重みを自由にしたTDNNの  
シフトインバリエント性能(/bdgmnN/)(短い音節発声)

シフト幅(ms)	1	2	3
-30	19.7%	34.4%	43.7%
-20	33.7%	55.7%	67.0%
-10	67.6%	85.9%	92.1%
0	83.1%	93.6%	97.0%
10	71.1%	89.4%	94.5%
20	36.6%	60.7%	77.8%
30	22.4%	33.4%	51.2%

#### 4. 3 入力層と中間層の間の重みを自由にする手法

入力層と中間層間の重みを自由にして、入力層と中間層の間の接続がどのくらいシフトインバリエントに影響を及ぼすかを調べる。実験結果を表 2 6、表 2 7 に示す。表 2 6 は各発声様式による、音韻認識率である。表 2 7 は各シフト (ms) に対する音韻認識率である。この実験の学習には D Y N E T を用いた。

##### ● 検討

表 2 6、表 2 7 と表 2 1 と表 2 2 と比較すると次のことが分かる。音韻認識率はかなりよい値を示しているが、シフトインバリエント性はかなり悪化する。入力層と中間層の間の重みがシフトインバリエントに重要であることが確認された。表 2 7 を表 2 5 と比較すると、シフトインバリエントには中間層と出力層の接続の方が効果あると考えられる。この実験の学習には D Y N E T を用いた。

#### 4. 4 出力層と中間層の間の重みに窓を掛ける手法

出力層の直前の中間層は、各音韻毎 9 個のユニットを持っている。このユニットの周波数方向での一番外側の値は出力層に対してあまり有効に作用していない。そこで、出力の直前の中間層と出力層との間の重みに図 6 に示すような窓を掛けることを行った。この効果による各発声様式による認識率の比較を表 2 8 に示す。また、このときのシフトインバリエント性がどの様になるかを表 2 9 に示す。この実験の学習には D Y N E T を用いた。

##### ● 検討

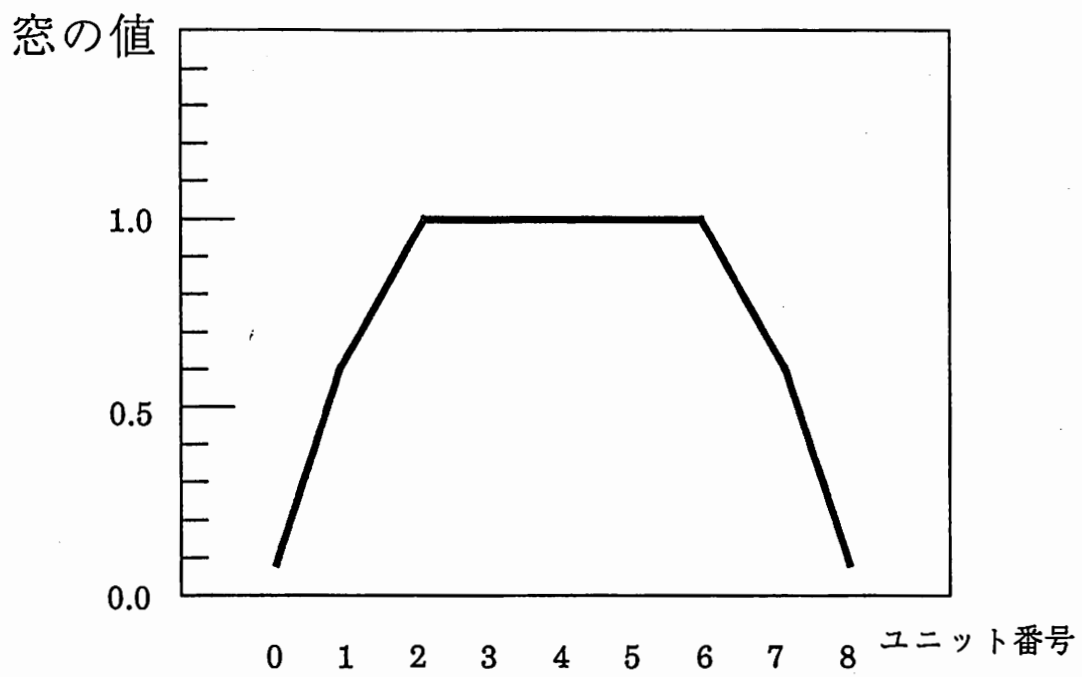


図6 出力層と中間層の間の結合に掛ける窓の値

表26 入力層と中間層の重みを自由にした  
TDNNの累積認識率(/bdgmnN/)

順位	1	2	3
文節発声	84.7%	94.9%	98.4%
短い文節発声	84.9%	94.6%	97.7%
自由発声	78.8%	92.0%	96.3%

表27 入力層と中間層の重みを自由にしたTDNNの  
シフトインバリエント性能(/bdgmnN/)(短い音節発声)

シフト幅(ms)	1	2	3
-30	30.0%	48.7%	64.3%
-20	44.0%	61.6%	74.9%
-10	71.7%	87.0%	94.0%
0	84.9%	94.6%	95.7%
10	79.1%	91.8%	96.8%
20	52.7%	75.4%	87.4%
30	28.3%	49.4%	68.3%

表 2 8、表 2 9 と表 2 1 と表 2 2 と比較すると次のことが分かる。音韻認識率はかなりよい値を示しており、シフトインバリエント性も比較的よい結果となっている。このことは出力層と中間層の間の重みと中間層の出力値をうまく設定すると認識率を向上させることが可能であるということを示している。言い換えると今の T D N N の出力層と中間層の間の接続には改良の余地があると考えられる。

#### 4. 5 出力層と中間層の間の重みを固定する手法

4. 4 で調べた結果から出力層の直前の中間層の出力はあまりよく学習されていないのではないかと考え出力層と最終の中間層の間の重みをすべて 0. 5 に固定し、出力層の閾値の値もすべて 0 に固定して学習を行った。このことにより出力層に接続する中間層の全てのユニットが 0. 5 または - 0. 5 になるように学習される。この効果による各発声様式による認識率の比較を表 3 0 に示す。また、このときのシフトインバリエント性がどの様になるかを表 3 1 に示す。

#### ● 検討

表 3 0、表 3 1 と表 2 1 と表 2 2 と比較すると次のことが分かる。音韻認識率は向上を示しており、シフトインバリエント性も向上している。

このネットワークの学習は出力層に接続する中間層に 0. 5 または - 0. 5 を直接教えることに相当する。この学習で著者が期待したのは次のことである。この中間層の 1 つのユニットの下には入力層の 7 フレームの窓内のユニットすべてがつながっている。もしこの窓内にユニットに対応する音韻の特徴が存在するので

表28 出力層と中間層の間の重みに窓をかけた  
TDNNの累積認識率(/bdgmnN/)

順位	1	2	3
短い文節発声	79.2%	92.5%	95.9%
自由発声	69.8%	86.4%	93.1%

表29 出力層と中間層間の重みに窓をかけたTDNNの  
シフトインバリエント性能(/bdgmnN/)(短い音節発声)

シフト幅(ms)	1	2	3
-30	37.8%	62.5%	82.2%
-20	57.2%	81.7%	92.2%
-10	72.6%	90.3%	95.0%
0	79.2%	92.5%	95.9%
10	80.8%	93.2%	96.3%
20	74.8%	89.5%	95.1%
30	59.1%	81.0%	90.8%



表30 出力層と中間層の間の重みを固定した  
TDNNの累積認識率(/bdgmnN/)

順位	1	2	3
文節発声	76.0%	92.4%	96.2%
短い文節発声	77.3%	90.7%	95.4%
自由発声	69.8%	88.4%	94.7%

表31 出力層と中間層の間の重みを固定したTDNNの  
シフトインバリエント性能(/bdgmnN/)(短い音節発声)

シフト幅(ms)	1	2	3
-30	59.9%	80.2%	90.0%
-20	68.1%	85.7%	92.9%
-10	73.3%	89.7%	89.7%
0	77.3%	90.7%	95.4%
10	76.8%	91.3%	95.5%
20	75.0%	89.2%	94.9%
30	71.8%	87.2%	93.8%

あれば、このユニットが0.5になり、そうでなければ-0.5になってほしいということである。つまり少しでも音韻を識別する特徴が入力層7フレームの窓内にあれば、それを見つけるようにネットワークを学習させるということになる。

#### 4. 6 クラス分けネットワークを付加する手法

3. の実験に於て、クラス分けのネットワークが効果があるのではないかと考えられる。そこで図に示すような / b d g / と、 / m n N / とを分けるクラス分けネットを用いて、各発声様式の違いによる音韻の認識率を表3.2に示す。この実験の学習にはDYNETを用いた。

##### ● 検討

3. の実験とこの実験からクラス分けのネットワークが有効であると考えられる。

#### 5. まとめ

本研究を通していろいろな実験を行った。この過程のなかで以下の4つのことがいえる。

- (1) クラス分けネットの有効である。
- (2) 出力層と中間層の接続の改良が必要である。
- (3) 従来のTDNNの入力層と中間層の間の接続がシフトインバリエントに有効である。
- (4) 従来のTDNNの出力層と中間層の間の接続がシフトインバリエントに有効である。

#### 謝辞

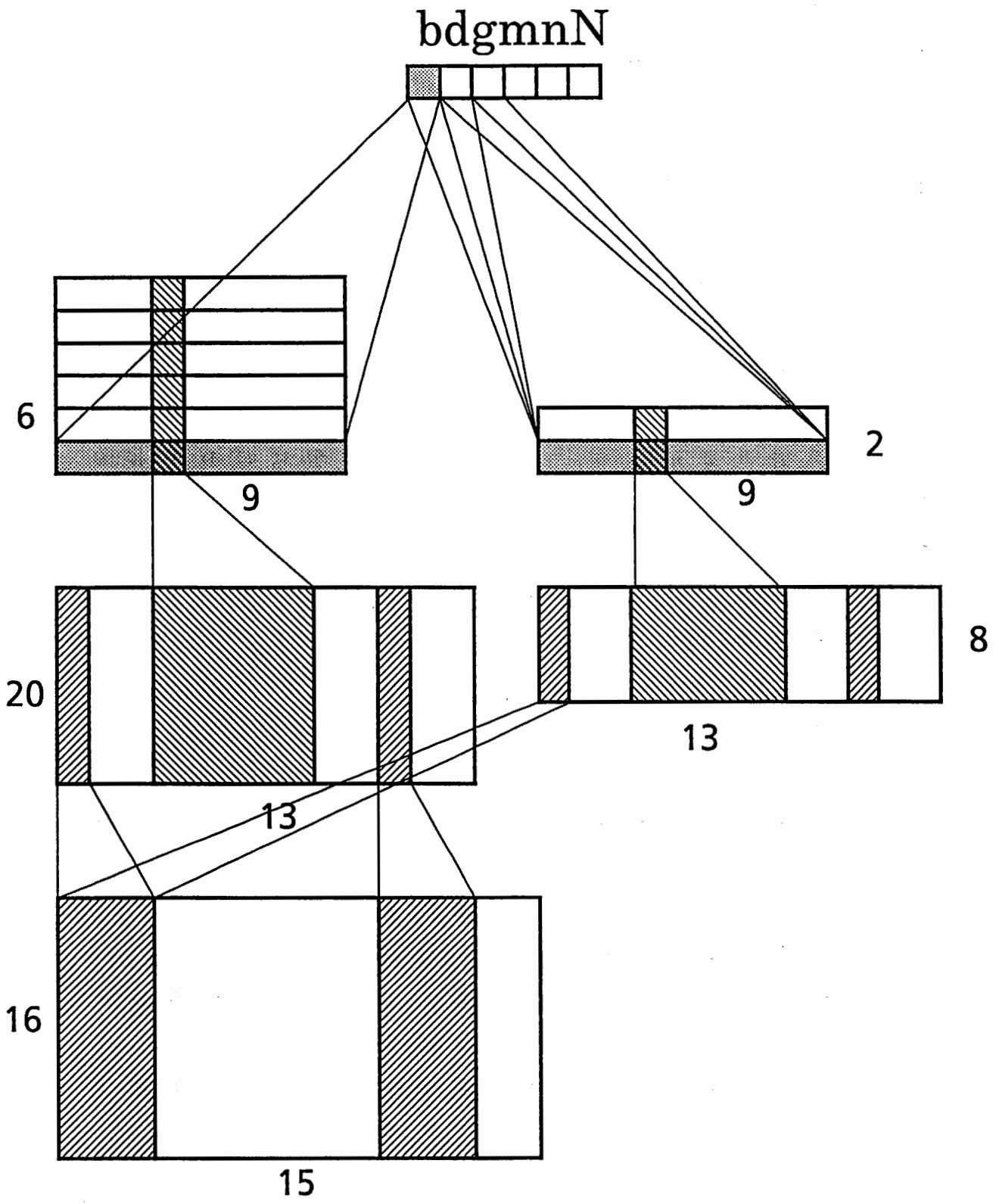


図7 クラス分けネットを持つ/bdgmnN/用TDNN

表32 クラス分けネットを付加した  
TDNNの累積認識率(/bdgmnN/)

順位	1	2	3
短い文節発声	80.7%	92.5%	96.6%
自由発声	67.6%	86.9%	94.3%

本研究の機会を与えて頂いたATR自動翻訳電話研究所榑松明社長に深謝致します。また、熱心に御討論頂いた音声情報処理研究室前室長の鹿野清宏博士、主任研究員の川端豪博士、田村震一氏をはじめとする音声情報処理研究室の皆様には感謝致します。

#### 参考文献

- (1) A. Waibel: "時間遅れ神経回路網(TDNN)による音韻認識"、信学会、信学技法SP87-100(1987.12)
- (2) 南泰浩、宮武正典、沢井秀文、鹿野清宏: "TDNN音韻スポッティングと拡張LRパーザを用いた文節音声認識"、音響講論集 3-1-11(1989.10).
- (3) 武田一哉、匂坂芳典、片桐滋、桑原尚夫: "研究用日本語音声データベースの構築"、音響学会誌, 44, 10, pp. 747-754(昭63.10).
- (4) A. Waibel, H. Sawai and K. Shikano, "Consonant Recognition by Modular Construction of Large Phonemic Time-Delay Neural Networks", ICASSP'89, pp. 112-115(1989).
- (5) D. E. Rumelhart and J. L. McClelland: "Parallel Distributed Processing; Explorations in the Microstructure of Cognition", Chap. 8, vol. II, MIT Press, Cambridge, MA(1986).

(6) P. Haffner, H. Sawai, A. Waibel and K. Shikano: "Fast Back-Propagation Learning Methods for Large Phonemic Neural Networks"、音響講論集 1-6-14 (1989.3) または P. Haffner, A. Waibel, H. Sawai and K. Shikano: "Fast Back-Propagation Learning Methods for Neural Networks in Speech" ATR Tech. Report TR-1-0058 (1988.11).