

TR-A-0104

0009

固定Q型蝸牛フィルタを用いた単語認識
—耐雑音性、耐残響性、話者変動の評価—

小原 和昭 平原 達也

1991. 3.19

ATR 視聴覚機構研究所

〒619-02 京都府相楽郡精華町乾谷 ☎07749-5-1411

ATR Auditory and Visual Perception Research Laboratories

Inuidani, Sanpeidani, Seika-cho, Soraku-gun, Kyoto 619-02 Japan

Telephone: +81-7749-5-1411

Facsimile: +81-7749-5-1408

Telex: 5452-516 ATR J

目次

- 第1章 はじめに
- 第2章 特徴抽出のための信号処理
- 第3章 実験と実験結果
- 第4章 考察
- 第5章 まとめと今後の課題

第1章 はじめに

人間の聴覚系は音声認識装置として非常に優れた性能を持っている。例えば代表的な音声認識装置はSN比が25 [dB]以下のノイズ環境下ではその認識性能は大きく劣化してしまう。一方我々はSNが5~10 [dB]程度の悪い環境下であっても話者の発声内容を容易に認識できる[1]。この様な聴覚系での信号処理機構を反映した特徴抽出を音声認識フロントエンドとして用い音声認識の性能を改善しようとする試みがなされてきている。

本報告は聴覚固定Q型蝸牛フィルタを音声認識のフロントエンドとして用いた場合の雑音、残響、話者変動に対する認識耐性をDTWによる単語認識を通してDFTフロントエンドと比較したものである。第2章では聴覚モデルを音声認識のフロントエンドとして用いた従来の認識実験結果をまとめるとともに、本実験で用いた実験条件、音声データの作成方法について述べる。第3章では実験内容とその結果について述べる。第4章では実験に用いたフロントエンドについて考察する。第5章では本研究のまとめと今後の課題について整理する。

第2章 特徴抽出のための信号処理

2-1 聴覚末梢系のモデルをフロントエンドとした音声認識

聴覚末梢系のモデルをフロントエンドとして音声認識を行なう試みがいくつかなされている。ここではそれらの試みを概観し、その結果をまとめた。聴覚モデルをフロントエンドとして音声認識を行った従来の実験について第2-1図にまとめた [2] [3] [4] [5] [6] [7] [8] [9]。これらの結果を整理すると聴覚末梢系のモデルをフロントエンドとして用いて特徴抽出を行える効果については、

- 1) 耐雑音性の向上 [3] [4] [7]。
- 2) 認識率の改善 [5] [6] [8]

等が報告されている。

雑音下での聴覚モデルをフロントエンドとした音声認識の検討はいくつかなされているが、残響環境下での認識性能や話者依存性について、このようなフロントエンドを用いた検討は見あたらない。本研究は様々な雑音、残響環境下および話者変動下において、固定Q型蝸牛フィルタをベースにしたモデルによる特徴抽出と従来の代表的なDF Tメルスケールによる特徴抽出とを、DTWによる単語認識を通して比較検討したものである。

2-2 DFTによる特徴抽出

DFTによる特徴抽出のブロック図を第2-2図に示した。メルスケールDFTとバークスケールDFTによる特徴抽出をおこなった。メルスケールDFTでは20 [kHz] サンプリングされた音声を20 [ms] のハミング窓を10 [ms] 毎シフトし512点FFTし、10 [kHz] の帯域をメル周波数軸上で64 [ch] に分割した各帯域に含まれるFFT離散スペクトルパワーを加算し、各帯域に含まれるFFT離散スペクトルパワー数で平均して対数をとることにより64 [ch] の特徴ベクトルを求めた。Hz周波数からメル周波数への変換は次式を用いた。

$$\text{MEL} = \frac{1000}{\log_{10}(2.0)} * \log_{10}\left(1.0 + \frac{\text{Hz}}{1000.0}\right)$$

バークスケールDFTでは20 [kHz] サンプリングされた音声を20 [ms] のハミング窓を10 [ms] 毎シフトし512点FFTし、19.5 [Bark] 帯域をバーク周波数軸上で64 [ch] に分割した各帯域に含まれるFFT離散スペクトルパワーを加算し、各帯域に含まれるFFT離散スペクトルパワー数で平均して対数をとることにより64 [ch] の特徴ベクトルを求めた。Hz周波数からバーク周波数への変換は以下に示したZwickerの変換式[13]を用いた。

$$\text{Bark} = 13.0 * \text{atan}(0.76 * \text{kHz}) + 3.5 * \text{atan}\left(\frac{\text{kHz}}{7.5}\right)^2$$

2-3 固定Q型蝸牛フィルタ

本実験で用いた固定Q型蝸牛フィルタの特性の一例を第2-3図に示した。低域の減衰特性が高域のそれよりもゆるやかなことが固定Q型蝸牛フィルタの特徴である。1.5 [Bark] から19.5 [Bark] までを64分割して求めたBark周波数B [1] ~ B [64] をもとに64 [ch] の固定Q型蝸牛フィルタを求めた。特性はノッチフィルタとバンドパスフィルタの従続接続によって生成している [10]。

ノッチフィルタの伝達関数N [s] は

$$N[s] = \frac{\omega_p}{\omega_z} * \frac{s^2 + (\omega_z/Q_z)*s + \omega_z^2}{s^2 + (\omega_p/Q_p)*s + \omega_p^2} \quad (2-1)$$

で与えた。ここで ω_p, ω_z はそれぞれ極と零点の角周波数でB [1] ~ B [64] をHz周波数に変換した後、pre-warpingした周波数により設定した。

すなわち

$$\omega_p = \frac{\tan(F_d[n]*T_s*\pi)}{(T_s*\pi)} * \sqrt{R} \quad (n=1 \text{ to } 64)$$
$$\omega_z = \frac{\tan(F_d[n]*T_s*\pi)}{(T_s*\pi)} / \sqrt{R} \quad (n=1 \text{ to } 64)$$

ここで

Fd[n]は [Bark] 周波数B [n] に対応する [Hz] 周波数、
Tsはサンプリング周期 (1/20 [kHz])、
Rは極と零点の角周波数の比 ω_p/ω_z で0.975と設定した。

またQz, Qpはそれぞれ極と零点のQの値で今回の設計ではQz=7、Qp=5と固定した。(2-1)式を双一次変換して2次のデジタルノッチフィルタDN [z] を求めた。すなわち

$$DN[z] = \frac{B_0*z^2 + B_1*z + B_2}{z^2 + A_1*z + A_2}$$

但し

$$B_0 = \left(\frac{\omega_p}{\omega_z} \right)^2 * \left\{ \left(\frac{2}{T_s} \right)^2 + \left(\frac{\omega_z}{Q_z} \right)^2 * \left(\frac{2}{T_s} \right) + \left(\frac{2}{\omega_z} \right) \right\} / A$$

$$B_1 = \left(\frac{\omega_p}{\omega_z} \right)^2 * 2 * \left\{ \omega_z^2 - \left(\frac{2}{T_s} \right)^2 \right\} / A$$

$$B_2 = \left(\frac{\omega_p}{\omega_z} \right)^2 * \left\{ \left(\frac{2}{T_s} \right)^2 - \left(\frac{\omega_z}{Q_z} \right)^2 * \left(\frac{2}{T_s} \right) + \left(\frac{2}{\omega_z} \right) \right\} / A$$

$$A_1 = 2 * \left\{ \omega_p^2 - \left(\frac{2}{T_s} \right)^2 \right\} / A$$

$$A_2 = \left\{ \left(\frac{2}{T_s} \right)^2 + \left(\frac{\omega_p}{Q_p} \right)^2 * \left(\frac{2}{T_s} \right) + \left(\frac{2}{\omega_p} \right) \right\} / A$$

$$A = \left\{ \left(\frac{2}{T_s} \right)^2 + \left(\frac{\omega_p}{Q_p} \right)^2 * \left(\frac{2}{T_s} \right) + \omega_p^2 \right\}$$

である。

また バンドパスフィルタの伝達関数H [s] は

$$H[s] = \frac{\omega_b}{Q_b} * \frac{s^2 + (\omega_b / Q_z) * s + \omega_z^2}{s^2 + (\omega_p / Q_p) * s + \omega_p^2} \quad (2-2)$$

で与えた。ここで ω_b はバンドパスフィルタの共振角周波数でB [n] (nはチャンネル番号でn=1~64)をHz周波数に変換した周波数Fd [n]をprewarpingした周波数により設定した。

すなわち

$$\omega_b = \frac{\tan(F_d[n] * T_s * \pi)}{(T_s * \pi)} * \sqrt{R} \quad (n=1 \text{ to } 64)$$

またQbはバンドパスフィルタのQの値で今回の設計では各BPFの中心周波数B [n]から低域カットオフ0.5 [Bark]、高域カットオフ0.5 [Bark]それぞれ離れた角周波数BL [n]とBH [n]とをHz周波数に変換した後prewarpingした周波数により設定した。すなわち

$$Q_b = \frac{F_c[n]}{F_H[n] - H_L[n]}$$

$F_c[n]$ は $B[n]$ に対応する [Hz] 周波数を prewarping した周波数、 $F_L[n]$ は $(B[n] - 0.5)$ に対応する [Hz] 周波数を prewarping した周波数、 $F_H[n]$ は $(B[n] + 0.5)$ に対応する [Hz] 周波数を prewarping した周波数である。

(2-2) 式を双一次変換して2次のデジタルバンドパスフィルタ $DH[z]$ を求めた。すなわち

$$DH[z] = \frac{B_0 * z^2 + B_1 * z + B_2}{z^2 + A_1 * z + A_2}$$

ここで

$$B_0 = \left(\frac{\omega_p}{Q_b} \right) * \left(\frac{2}{T_s} \right) / A$$

$$B_1 = 0.0$$

$$B_2 = \left(\frac{\omega_p}{Q_b} \right)^2 * \left(\frac{2}{T_s} \right) / A$$

$$A_1 = 2 * \left\{ \omega_b^2 - \left(\frac{2}{T_s} \right)^2 \right\} / A$$

$$A_2 = \left\{ \left(\frac{2}{T_s} \right)^2 - \left(\frac{\omega_b}{Q_b} \right) * \left(\frac{2}{T_s} \right) + \left(\frac{2}{\omega_p} \right) \right\} / A$$

但し

$$A = \left\{ \left(\frac{2}{T_s} \right)^2 + \left(\frac{\omega_b}{Q_b} \right) * \left(\frac{2}{T_s} \right) + \omega_p^2 \right\}$$

である。

以上の設計でBark周波数からHz周波数への変換には以下の関係式を用いた。

```
if (Bark <= 5) then
    Hz = 100 * Bark
else if (Bark > 10.04) then
    Hz = 1000. * exp { (Bark - 8.85) / 6. }
else
    Hz = 1000. * (Bark - 1.5) / 7.
```

2-4 モデル1

モデル1は固定Q型蝸牛フィルタの特徴抽出とDFTによる特徴抽出で認識率の差異がどの程度生じるかを検討するモデルである。すなわちモデル1は64 [ch] 固定Q型蝸牛フィルタでフィルタリングしてその出力をウインド長20 [ms] シフト10 [ms] のハミングウインドウで切出し、フレーム内に含まれるデータの平均パワーの対数を取り特徴ベクトルとするモデルである。

2-5 モデル2

モデル2はモデル1の時間軸方向の変動を滑らかにする平均化処理に加えて周波数方向の変動を滑らかにする処理を行なったモデルである。周波数軸方向の処理は64 [ch] 出力に対しLPF処理を行った(LPリフタリング)。このLPF処理はスペクトログラム上の周期の速い変動を減衰させることを目的に行った。LPリフタリングはBark周波数方向にLPF処理を行うことで行った(周波数方向のデータを76 [KHz] サンプリングとしてカットオフ周波数を18 [KHz] に設定した)。

モデル1、2での信号処理を第2-4図にまとめた。

2-6 音声データ

音声の仕様を第2-5図にまとめた。サンプリングは20 [kHz] で音声の開始点から終了点までをラベルデータに基づき切り出しを行った。単語セットとしてはATR音声データベースの音韻バランス単語216語を使用した。話者は男性2名 (MST, MYS) 女性2名 (FST, FNY) を用いた。

2-7 雑音の重畳

雑音による単語認識の性能劣化を調べるために音声データに種々のSNの雑音を重畳し雑音音声データを作成した。雑音源として20 [kHz] でサンプリングされた5秒のピンク雑音を用いた。初めに音声のトータルパワーをもとめる。その音声と同じ時間長のノイズを開始点がランダムになるよう雑音源から切り出し、設定されたSNからノイズの振幅を設定し、音声に重畳した。ここでのSNは音声のトータルパワーと雑音のトータルパワーの比で定義した (グローバルSN)。SNが ∞ (ピュア), 40, 20, 10, 5, 0 [dB] の雑音音声データを作成した。すなわち音声データを $S[n]$, ランダムに切り出したノイズデータを $N[n]$, 雑音重畳された音声データを $S_n[n]$ とすると以下の式によって雑音重畳された音声データを作成した。

$$S_n[n] = S[n] + N[n] \frac{\sqrt{\sum_n (S[n]*S[n])}}{\sqrt{\sum_n (N[n]*N[n])}} * 10^{-\frac{SN}{20}}$$

2-8 残響音声の生成

残響音声は原音声とATR可変残響室で収録したインパルスレスポンスの畳み込みを行うことで生成した。インパルスレスポンスとしては残響時間の異なる5種類を50 [ms] の時間長で切り出して用いた。各々のインパルスレスポンスを第2-6図に示した。各々のインパルスレスポンスの500 [Hz] における残響時間はそれぞれ、約200、210、460、620、1100 [ms] であった。残響下での音声データ作成のブロック図を第2-7図に示した。音声と残響のインパルスレスポンスを2048サンプルのFFTとoverlap-add法によって残響音声を生じた [11]。畳み込みによって生成された音声は原音声と同じ時間長になるように音声の終端部を切り捨てた。

第3章 認識実験

第2章で述べたモデル1、モデル2、DFTにより音声の特徴抽出を行いDTWによる単語認識を行った。DTWとしてはスタガードアレーDP [12]を用いた。整合窓 r はテンプレートのフレーム長を L_t 、認識単語のフレーム長を L_w とすると

$$r = \min(L_t, L_w) / 4 + 3$$

とした。格子点の間引き間隔を3として、文献[12]のDP3-1のアルゴリズムによって認識を行った。認識実験の構成図を第3-1図に示す。認識実験結果を以下に示す。

3-1 認識実験とその結果

実験1) ウィンドウの影響について (1話者=MST)

特徴抽出に与えるウィンドウの影響を調べるために方形、ハミングの2種類の切り出しウィンドウを用いて、モデル2でのノイズによる認識率変動を調べた。ハミングウィンドウ長は20 [ms]で10 [ms]シフトで、また方形ウィンドウ長は10 [ms]で10 [ms]シフトで特徴抽出を行った。結果を第3-2図に示した。ノイズによる認識率の変動はウィンドウによらず同様な傾向を示した。また認識率には大きな差はみられなかった。以下の実験にはDFTとの比較のためにハミングウィンドウを用いた。

実験2) モデル1、2とDFTの耐雑音性の比較 (4話者)

各話者の一回目の発声をテンプレートとし二回目の発声に雑音を重畳して、その単語の認識実験を行った。4話者の単語認識でのDFT、モデル1、2の認識結果を第3-3図~第3-5図に示した。また4話者認識結果の平均を第3-6図に示した。SN40 [dB]以上ではDFTの性能がモデルの性能と同等か5 [%]程度上回ったが、SN20 [dB]以下ではモデル1、2の性能がDFTによる認識率を5~25 [%]上回った。SNの劣化に対してモデル1、2の認識率の劣化の方がDFTの認識率の劣化よりも少なかった。またバークスケールDFTのほうがメルスケールDFTよりも低SNでの劣化が少なかった。

実験3) モデル1、2とDFTの耐残響性の比較 (1話者 MST)

モデル1、2とDFTによる耐残響性能を調べるために、種々の残響時間を持つ音声データを作成して認識実験を行った。テンプレートは残響付加もノイズ付加もないクリーンな環境での発声を用いた。その結果を第3-7図に示した。残響環境下では残響時間によらずDFTによる特徴抽出の性能がモデルの性能を上回った。また残響時間の変化に対しDFTの認識率の変動はモデルの変動よりも少なく、実験に用いた条件内ではクリーンな環境から約5%の認識率の劣化があった。一方モデル1、2の認識率の劣化はクリーンな環境から約16%であった。今回用いた最も残響の長い環境ではDFTの方がモデル1、2に比べ12[%]程度認識率が高かった。またバークスケールDFTのほうがメルスケールDFTよりも劣化が大きかった。

実験4) モデル1、2とDFTの耐雑音性、耐残響性の比較 (1話者 MST)

雑音、残響下でのDFTとモデル1、2の認識率をクリーンなテンプレートを用いて比較した。残響インパルスR1 (第2-6図 R1) を用いて残響下での音声を生成して、それに雑音を加えSNが ∞ (ピュア), 40, 20, 10, 5, 0 [dB] の残響付加雑音音声データを作成した。認識結果を第3-8図に示した。残響下かつ雑音下での認識ではモデル1、2のほうがDFTよりも20 [dB] 以下の低SN下では認識性能が良いが、絶対的な認識性能はSNが20 [dB] で80 [%] 程度で、どちらのモデルでも大きく劣化した。

実験5) 正規化特徴量による雑音下認識 (1話者 MST)

各話者の一回目の発声をテンプレートとし二回目の発声に雑音を重畳して、その単語の認識実験を行った。各単語の特徴ベクトルの最大値で特徴ベクトルの正規化を行い認識を行った。すなわち各単語の特徴ベクトルを $S(m,n)$ とすると正規化特徴ベクトル $S'(m,n)$ を次式で求めた。

$$S'(m,n) = \frac{S(m,n)}{\max\{S(m,n)\}}$$

$$\left[\begin{array}{l} \max\{S(m,n)\}: \text{maximum of } \{S(m,n)\} \\ \text{for all } m \text{ and } n \end{array} \right]$$

メルDFT、モデル1、2による認識結果を第3-9図に示した。特徴量の正規化をおこなうことによりモデル1、2の認識率はSNによらずDFTよりも高くなった。また低SNでの認識率の大きな改善がみられた。DFTは特徴量の正規化を行うことで低SNでの認識率が低下したがモデル2では全てのSN

で認識率が向上した。

実験6) 正規化特徴量による残響下認識 (1話者 MST)

各話者の一回目の発声をテンプレートとし二回目の発声に残響付加を行った後、各単語の特徴ベクトルの最大値で特徴ベクトルの正規化を行い認識を行った。メルDF T、モデル1,2による認識結果を第3-10図に示した。特徴量の正規化をおこなうことによりモデル1,2の耐残響性が向上し認識率はDF Tと同程度になった。DF Tでも特徴量の正規化を行うことで認識率が数%向上した。

実験7) テンプレートと異なる話者による単語認識率

男性話者MSTの一回目の発声をテンプレートとし、他の話者の発声した単語の認識実験を行った (テンプレート MST、話者MYS、FST、FNY)。各単語の特徴ベクトルの最大値で特徴ベクトルを正規化した場合としない場合とについて、DF T、モデル1、2による認識結果を第3-11、12図に示した。異なる話者のテンプレートを用いた場合、DF Tの特徴量による劣化10 [%] に対しモデル特徴量による劣化は30 [%] 程度でDF T特徴量による認識の方がモデル1、2の特徴量による認識率よりも性能がよかった。特徴量の正規化をおこなうことによりモデル1、2の認識率は向上したがDF Tによる認識率の向上は少なかった。正規化を行ってもDF T特徴量による認識率の方がモデルによる認識率よりもよかった。モデル1、2の特徴量による正規化なしの場合の特徴量による認識ではとくに女性話者に対する認識率の劣化が男性話者の10 [%] に対し30 [%] 程度と大きかった。特徴量を正規化することにより話者間の変動が押えられた。

第4章 考察

この章ではDFTによる特徴抽出とモデル1、2による特徴抽出の差異についてまとめた。

4-1 モデルとDFTの特徴抽出の差

本実験で用いた4つの特徴抽出モデル（DFTメル、DFTバーク、モデル1、モデル2）の1000 [Hz]でのスペクトログラムを第4-1図に示した。これらを比べると、モデル2はDFTに比べかなり滑らかな特徴抽出になっておりまた応答チャンネルがDFTより広くなっている。またモデル1はモデル2との差は大きくはないがモデル2はモデル1よりも周波数方向の変動が少ない特徴抽出となっている。

モデルとDFTの特徴抽出の差を以下に示した。

(A) フィルタ特性：

モデルの周波数分析は基底膜のフィルタリング特性を反映して低域と高域の減衰特性が非対象であるが、DFTは通過帯域内での特性は一様でフラットである。また固定Q型蝸牛フィルタの帯域とDFTの帯域を比べるとDFTのほうが狭帯域である。固定Q型蝸牛フィルタの周波数軸はbarkスケールであり分析帯域は約5.7 [kHz]でありメルDFTの周波数軸はメルスケールで分析帯域が10 [kHz]となっている。そのため最大応答を示すチャンネルが異なっているがモデル1、2の固定Q型蝸牛フィルタのほうが広帯域になっているのが分かる。またバークスケールDFTでは分析帯域はモデル1、2とほぼ等しいが各チャンネルの周波数分解能はバークスケールDFTのほうがよい。

(B) フィルタ帯域ののオーバーラップ：

モデル1、2の固定Q型蝸牛フィルタは各チャンネルが約0.28 [bark]間隔に設定してあり隣接チャンネル間で通過帯域が重なっている。一方DFTでは隣接チャンネルの重なりはない。200 [Hz] + 500 [Hz] + 1000 [Hz]の入力に対するDFT、モデル1、2の出力を第4-2図に示した。DFTに比べ固定Q型蝸牛フィルタの周波数分析はブロードになっている。

(C) 時間積分による通過帯域：

モデル1、モデル2では20 [ms]のハミングウィンドウで音声を取り出す。

しその区間の400サンプルの時間軸の加算平均を10 [ms] 毎に行った。そのため各フィルタ出力のカットオフ周波数が50 [Hz] 程度となっている。またDFTでもフレーム周期を10 [ms] としているため各チャンネル出力のカットオフ周波数は50 [Hz] となっているので時間分解能は等しくなっている。

4-2 ノイズによるテンプレートと認識単語との距離の変動

ノイズによる認識性能の劣化の原因を調べるために、ノイズ重畳によるテンプレートと認識単語の距離の変動を調べた。クリーンな環境での1回目の発声と2回目の発声の距離を基準1として、テンプレートとノイズ付加による認識単語の距離の変化を測定した。その結果を第4-3図に示した。距離の変動はメルDFTによる特徴抽出の方がSNが20 [dB] 程度まではモデル1、2での距離変動よりも少なかった。これはDFTに比べモデル1、2の各チャンネルの周波数分析が広帯域であるためSNがDFTよりも悪いためであると考えられる。SNが比較的良いときにはこのSNの変化に対する距離の変化がSNの劣化に伴う認識率の変化に現れていると考えられる。すなわちSNが40 [dB] 程度の時にはDFTでの距離の変動対し、モデル1、2での距離の変化は大きく、認識率の劣化はモデル1、2の方が大きくなっている。

一方SNが20 [dB] 以下になるとモデル2での距離の変動はDFTの距離変動よりも少なくなっている。このためにSNが20 [dB] 以下ではモデル2の認識率がメルDFTよりも良くなっていると考えられる。またモデル1では距離の変動がメルDFTよりも大きくなっているが認識率はDFTより良くなっている。この点を調べるためにSNの変化に対するDFTとモデル1、2のスペクトログラムの変化を第4-4、5、6図に示した。これらの図からモデル1、2のスペクトログラムではSNが20 [dB] 程度になっても大まかなスペクトログラム上の特徴は保存されているが、DFTではスペクトログラム上の山谷の特徴がつぶれてきているのがわかる。モデル1、2がノイズにあまり影響されない比較的滑らかな特徴量を取り出しているの対しDFTによる特徴量はノイズによって大きく変化してしまっていると考えられる。すなわちDFTの特徴抽出は細かな特徴をよく取り出している反面ノイズによってその特徴が失われやすい。一方モデル1、2の特徴抽出は比較的大まかな特徴を取り出しているためノイズによってもその特徴が失われ難くなっていると考えられる。そのためにSNが20 [dB] 以下ではモデル1、2による認識率がDFTによるそれよりも良くなっていると考えられる。

すなわちSNが悪くなったときモデル1,2の特徴ベクトルのクラス間距離がDFTよりも大きくなっていると考えられる。また特徴量の正規化によってテンプレートと認識単語間の距離がどう変化するかを第4-7図に示した。正規化によってモデル1、2の距離変動がDFTよりも押えられているのがわかる。このことによりモデル1、2の認識率が正規化により向上したと考えられる。一方DFT特徴量を正規化すると20 [dB]以下の低SNではモデル1、2よりも距離変動が大きくなっているのがわかる。

4-3 残響付加によるスペクトログラムの変化

第4-8, 4-9, 4-10図に残響付加によるスペクトログラムの変動とその差分の一例を示した。モデル1、2に比べメルDFTは残響付加によりスペクトログラムの変動が少なくなっている。この点からも残響環境下でのDFTの認識率の劣化がモデル1、2よりも少ないことが推察できる。

4-4 特徴抽出にかかる時間の比較

モデル1、モデル2での処理時間はDFT処理時間の1.7倍であった。

4-5 誤りパターン

第4-11図にSN40 [dB]における単語認識の誤りの一例を示した。この図は例えばDFTによる特徴抽出において75番目の単語（きゅーげき 付録A参照）が第1候補から第5候補までそれぞれ201番目（とーひょー）、120番目（によーぼー）、18番目（ひょーじゅん）、53番目（いらしゃる）、23番目（わがまま）の単語と誤認識されたことを示している。モデル1、2、DFTとも同じ単語について似たような認識誤りをおかしているのがわかる。

第5章 まとめと今後の課題

本研究では以下の点が明らかになった。

1. モデル1、2とDF Tの特徴抽出を単語認識を通じて比較検討した結果、認識すべき単語のSNが40 [dB] 以上では両者の認識率の差はほとんどなかった。しかしノイズや残響がある環境での認識率はSNが20 [dB] 程度以下では固定Q型蝸牛フィルタを用いたモデル1、2のほうがDF Tよりも性能が良かった。
2. モデル1、2のSNによる認識率の劣化は特徴量のパワー正規化によって大きく改善され、SN10 [dB] においてDF Tの認識率16 [%] に対して73 [%] の認識率を得た。この認識率の改善は正規化によりテンプレートと認識単語間の距離変動を押えたことによる改善である。
3. 残響環境下の認識率はモデル1、2いずれよりもDF Tの方が良好であった。また残響環境の変化に対してもDF Tの認識率の変動がモデルの変動よりも少なかった。これはモデルによる特徴抽出よりもDF Tによる特徴抽出が残響環境での特徴量の変動が少ないためである。
4. モデル1,2の特徴量を正規化することにより対残響性は向上したがDF Tの性能は越えなかった。
5. テンプレートと話者が異なる場合の認識率の変動はDF Tの性能がモデル1、2の性能よりも良かった。この場合にもモデル1、2の特徴量を正規化することにより認識率は改善されたがDF Tの性能をこえることはなかった。
6. 総じてDF Tによる特徴抽出は今回の考慮した範囲での残響環境の変動に対し良好であるが低SNでの認識率の劣化が著しい。またモデル1、2による特徴抽出は低SNで効果が認められたが、その効果は低SNでの認識時に限られ、残響環境下での認識率の劣化はDF Tよりも大きかった。
7. モデル特徴量の正規化により耐雑音性,耐残響性が向上し,話者変動による認識率の変動も押えられることを確認した。

DFTやモデル1、2による特徴量が環境によって大きく変動してしまうため環境変動によって大きな認識率の劣化が生じる。今回のモデルで考慮した末梢系での信号処理のみでは人間の認識能力に匹敵する特徴抽出を行なうことは難しい。今後の課題としては、音声の認識においてどのような特徴量を人が利用して様々な環境の変化に対応しているのかを検討し、このような環境変動の影響を受けずらい音声特徴量の抽出方法とその特徴量を用いた認識の検討が必要である。

参考文献

- [1] Pollack,Pickett:"Masking of speech by noise at high sound levels", J.Acoust.Soc.Am.30, pp127-130(1958)
- [2] Blomberg,M.,et.al(1984):"Auditory Models and isolated word recognition," Q Prog.Stat. Rep.,SpeechTransmiss. Lab.(Royal Institute of Technology,Stockhom),pp.1-15
- [3] Ghiza,O.(1988):"Temporal Non-place information in the auditory-nerve firing patterns as a front-end for speech recognition in a noisy environment,"J.of Phonetics Vol.16,No.1 pp.109-123
- [4] 佐久嶋、磯、吉田、渡辺："聴覚モデル分析を用いた音声認識"、音学講論,1-3-9, 1988. 10
- [5] Hunt,M.J,C.Lefebvre(1988):"Speaker Dependent and Independent Speech Recognition Experiments with an Auditory Model," Proc. IEEE Int. Conf. Acoustics,Speech&Signal Processing,ICASSP-88,pp215-218(1988)
- [6] Cohen,J.R.(1989):"Application of an Auditory model to speech recognition,"J.Acoust.Soc.Am.85(6),pp.2623-2629
- [7] Paterson,R. Hirahara,T.(1989):"HMM Speech Recognition using DFT and Auditory Spectrogram,"ATR Tech.Report TR-A-0063
- [8] H.Hamada,T.Hirahara,A.Imamura,T.Matsuoka,R.Nakatu(1989):"Auditory-Based Filter-Bank Analysis as a front-end Processor for Speech Recognition," Proc. Eurospeech 89, Vol.2 pp.396-399
- [9] T.Hirahara,H.Iwamida(1990):"Auditory Spectrograms in HMM Phoneme Recognition,"1990 Int. Conf.on Spoken Language Processing,ICSLP-90, pp.381-384
- [10] 駒木根,平原："蝸牛の周波数分解機能を模擬するフィルターバンクの一構成法", 信学技報 SP87-45 (昭61)
- [11] A.V.Oppenheim & R.W.Schafer: "デジタル信号処理 (下)", (伊達訳) コロナ社 (昭53)
- [12] 鹿野、相川："Staggered Array DP マッチング"、信学論誌,61-D,9, pp.657-544 (昭57)
- [13] Zwicker, et al.(1980):"Analytical expression for critical-band rate and critical band width as a function of frequency", J.Acoust.Soc.Am.74(3),pp1523-1525

Blomberg et al.(1984)[2]

Model: 5 different spectral scaling
(FFT, Bark, Phon, sone, domin)
Task: 9 Swedish vowels and 18 Swedish consonants
Recognition: linear normalization of time
Results: FFT gives the best score

Ghiza(1988) [3]

Model: 85 ch log scale spacing cochlea filter+EIH
Task: 39-word alpha-degits(2male2female)
Recognition: DTW
Results: Same recognition score comparing DFT front end ,Better noise immunity

Sakushima et al.(1988)[4]

Model: Seneff model
Task: speaker dependent 250 Words
Recognition: DTW
Result: Better noise immunity

Hunt et al.(1988)[5]

Model: CBF+ Onset detector ,Periodicity detector
+Linear Discriminant analysis
Recognition: DTW
Task: talker independent and dependent, Number
Result: In all cases ,Better score comparing with
Mel Cepstrum

Cohen(1989)[6]

Model: 20 ch CBF+loudness scaling+ adaptation
Recognition: HMM (Training 100 Words, Test 50 Words)
Task: Words
Results: Better score comparing to Filter Bank Only
(Error 6.3-[%]>3.9[%])

Patterson et al.(1989)[7]

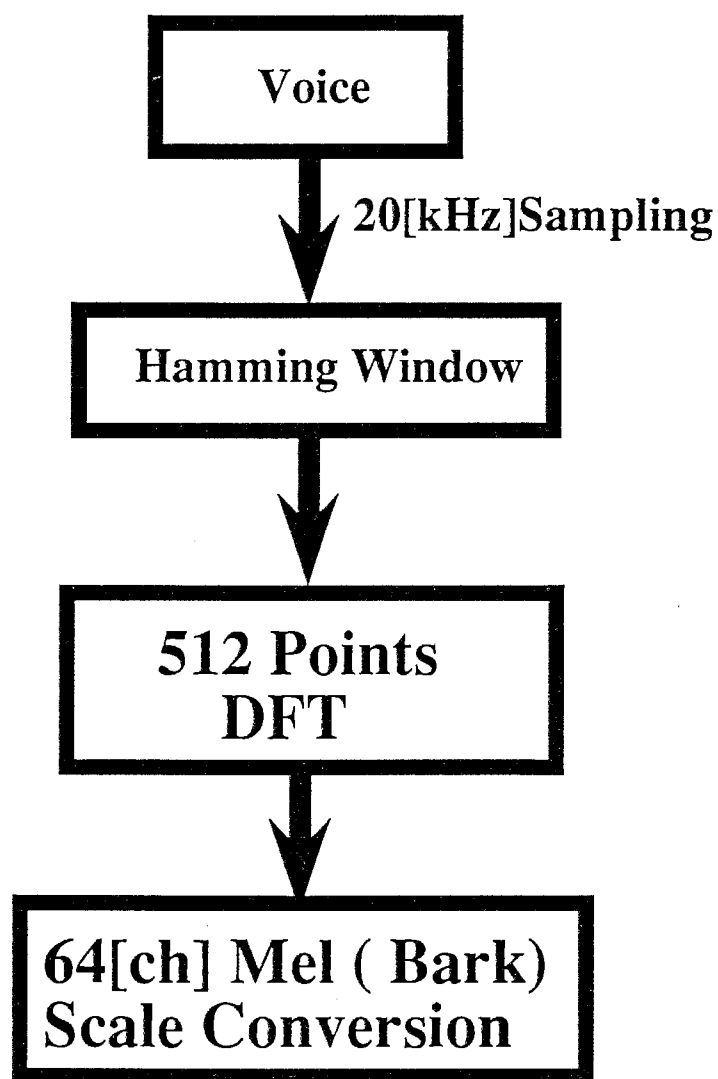
Model: S AS
Recognition: HMM
Task: 14 consonants
Results: When code book size is small,
better recognition score comparing DFT front
end

Hamada et al.(1989)[8]

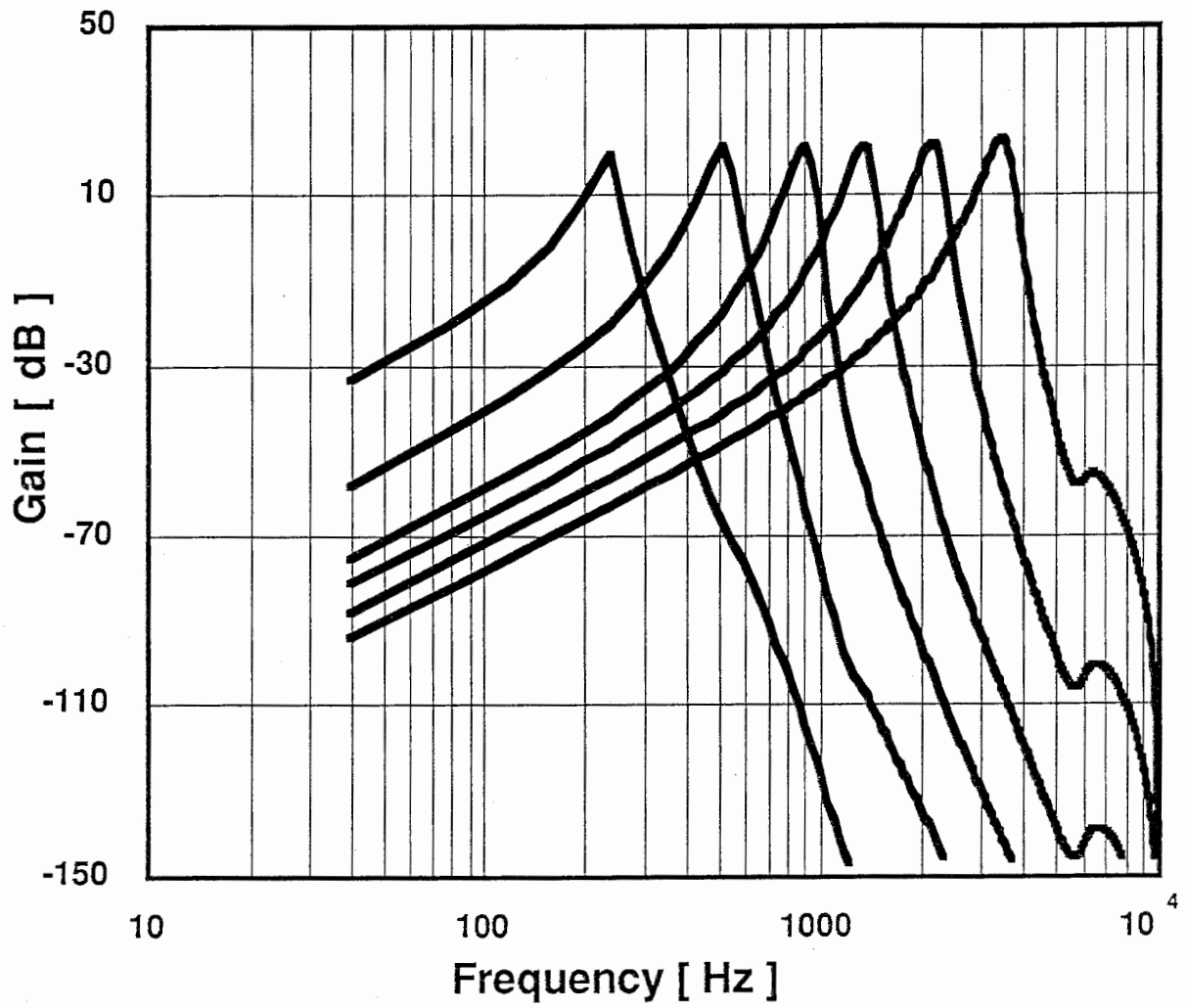
Model: 35 ch CBF+HWR+Integration+Log+LIN
Recognition: DTW, HMM
Task: 14 consonants
Results: Same or better recognition score comparing
LPC analysis

Hirahara et al.(1990)[9]

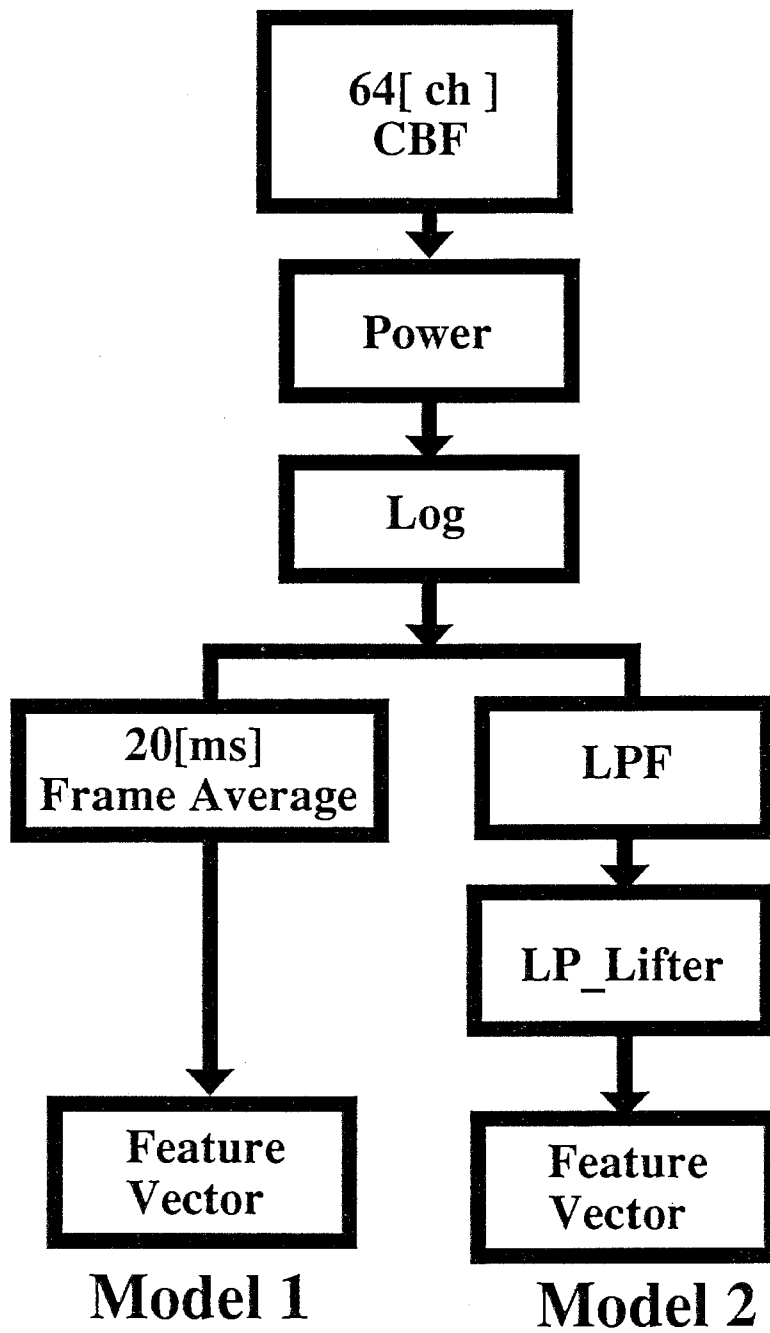
Model: Adaptive Q Chchlea Filter+Adaptation Model+LIN
Recognition: HMM
Task: 18 phoneme
Results: Without LINH,DFT gives the best score
With LINH ,,Adaptive Q Model gives the best score



第2-2図DFTによる特徴抽出



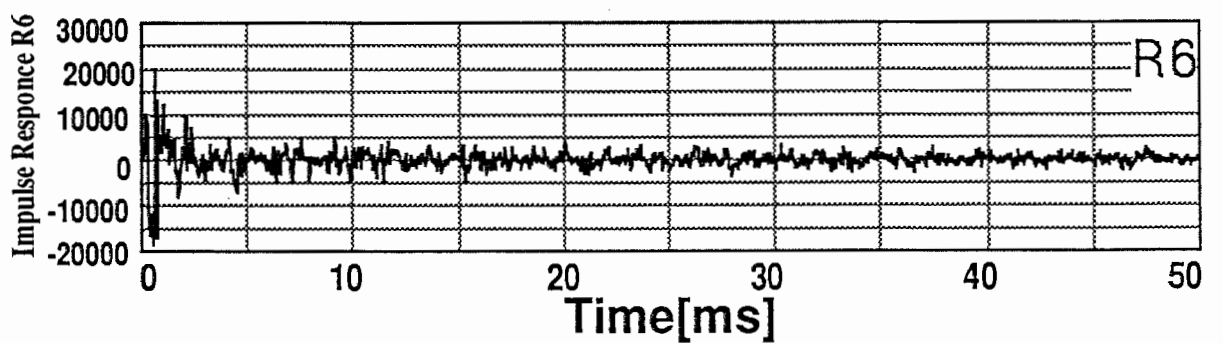
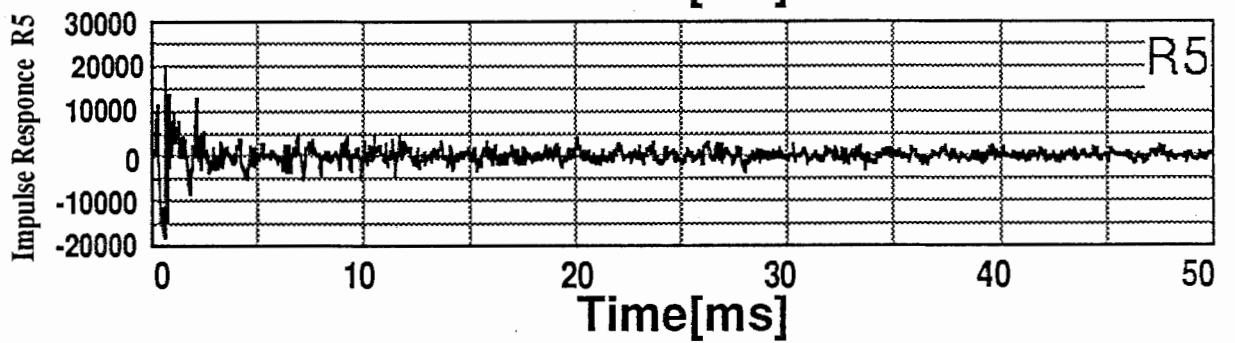
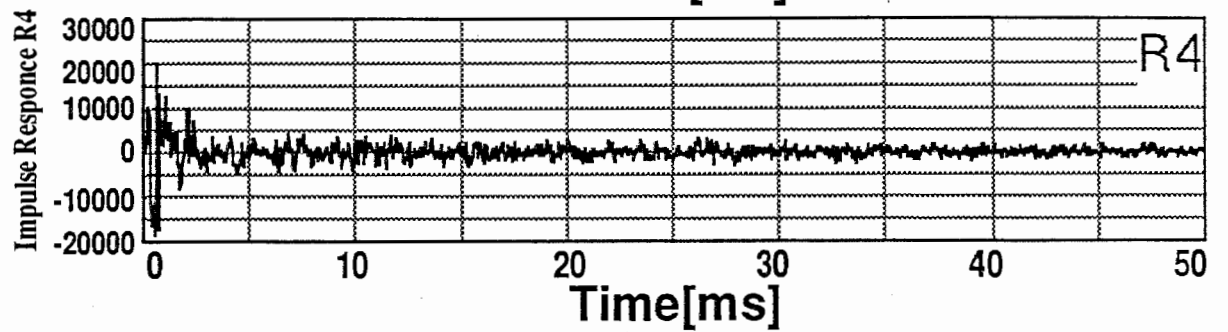
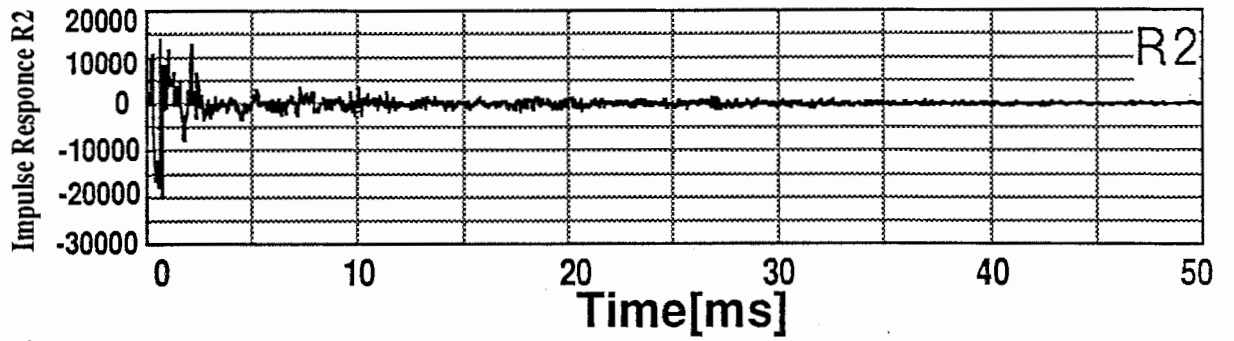
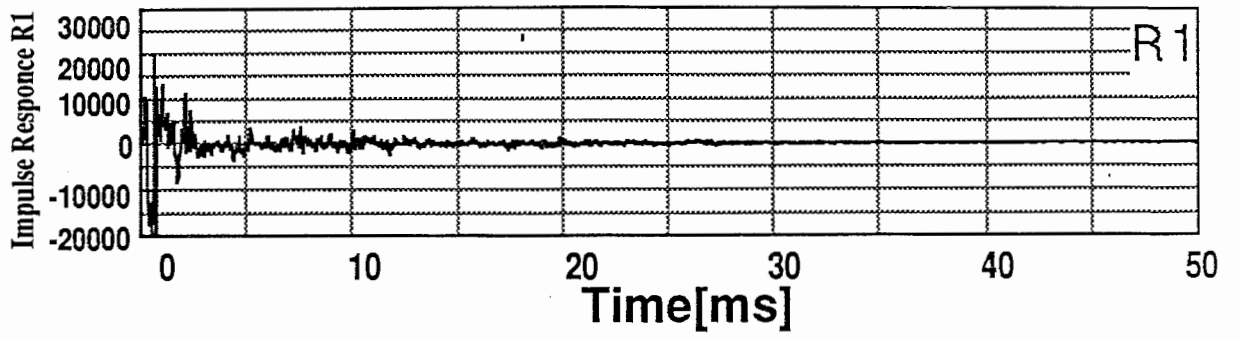
第2-3図 基底膜フィルタの特性の一例



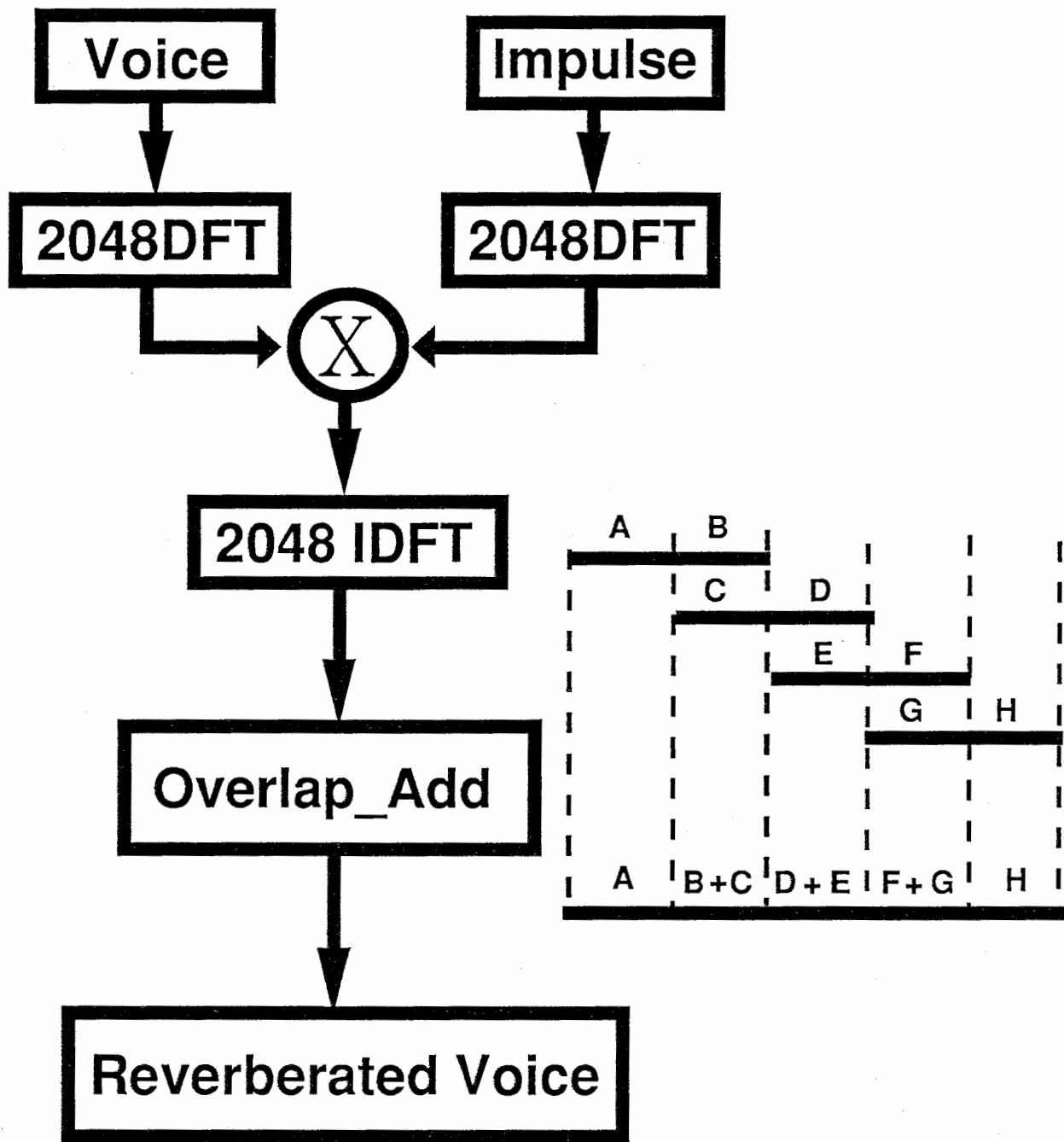
第2-4図 モデル1とモデル2による特徴抽出

サンプリング	20 [kHz]
話者	男性：MST, MYS 女性FST, FNY
単語セット	ATR 音声データベース 音韻バランス単語216語
リエンファシス	なし
雑音音声生成	雑音の音声への加算
残響音声生成	残響インパルスと音声の畳み込み
特徴抽出	モデル1 モデル2 メルスケールDFT
距離尺度	ユークリッド距離
認識部	スタグガードアレイDTW

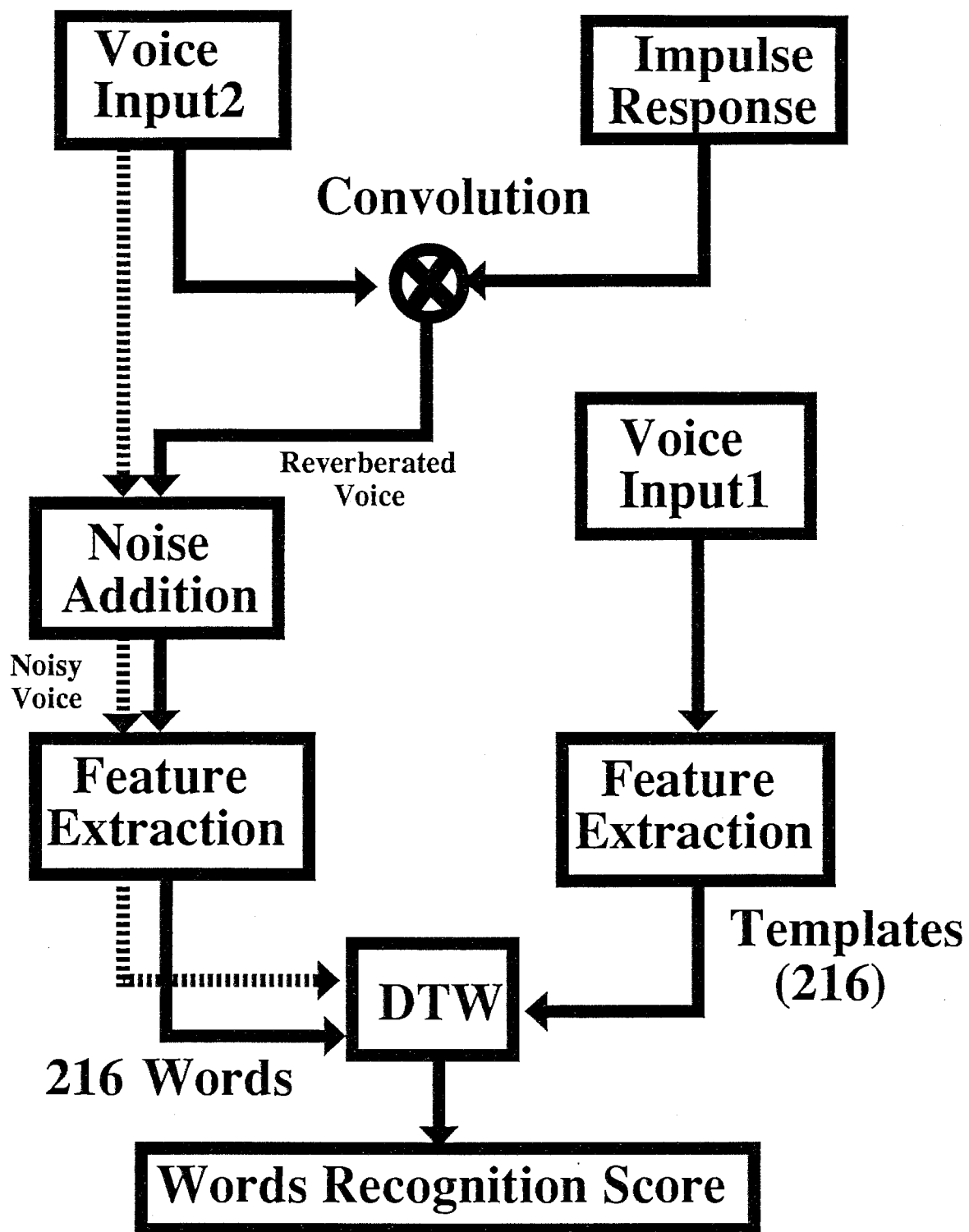
第2-5図 音声データの仕様



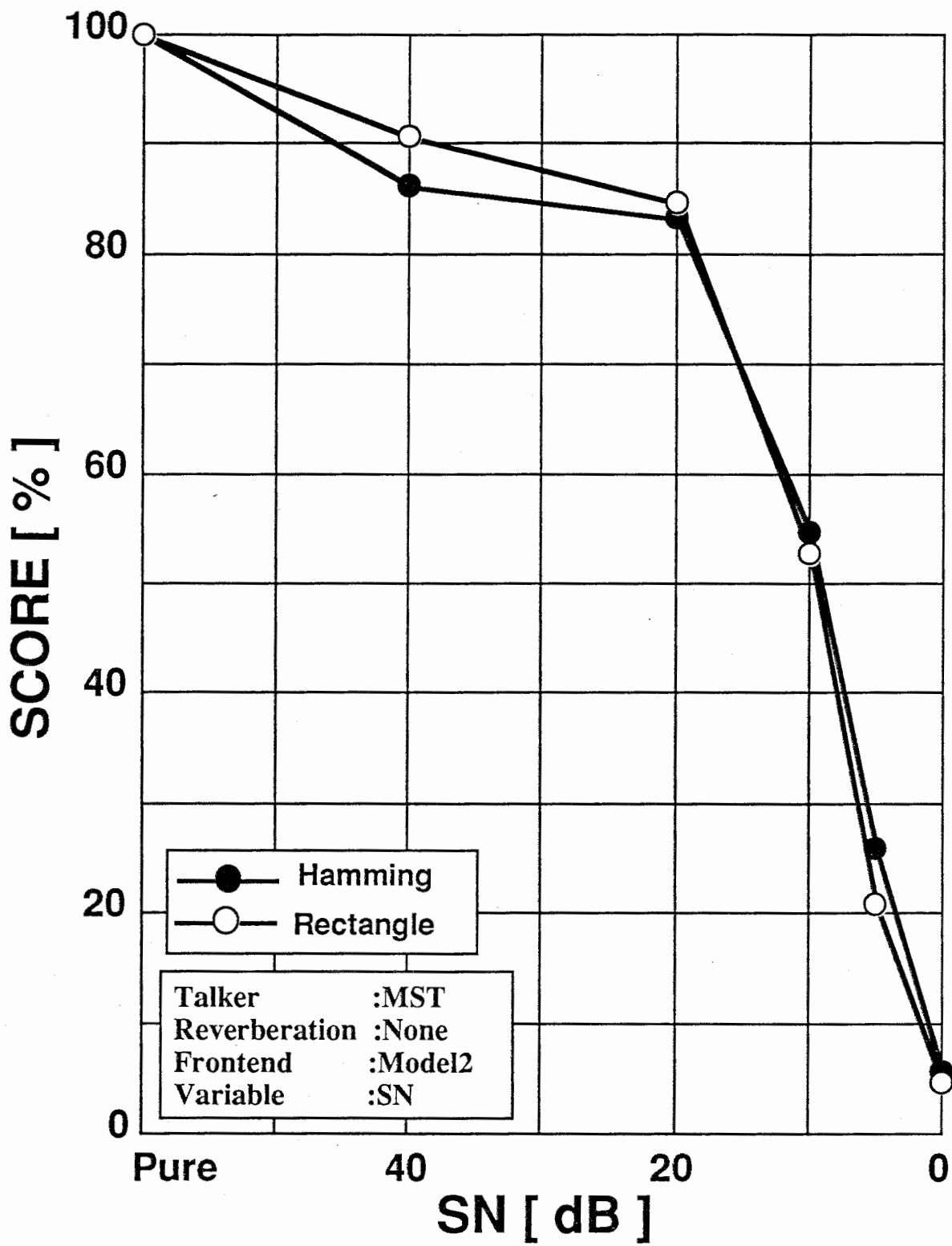
第2-6図 残響インパルスデータ



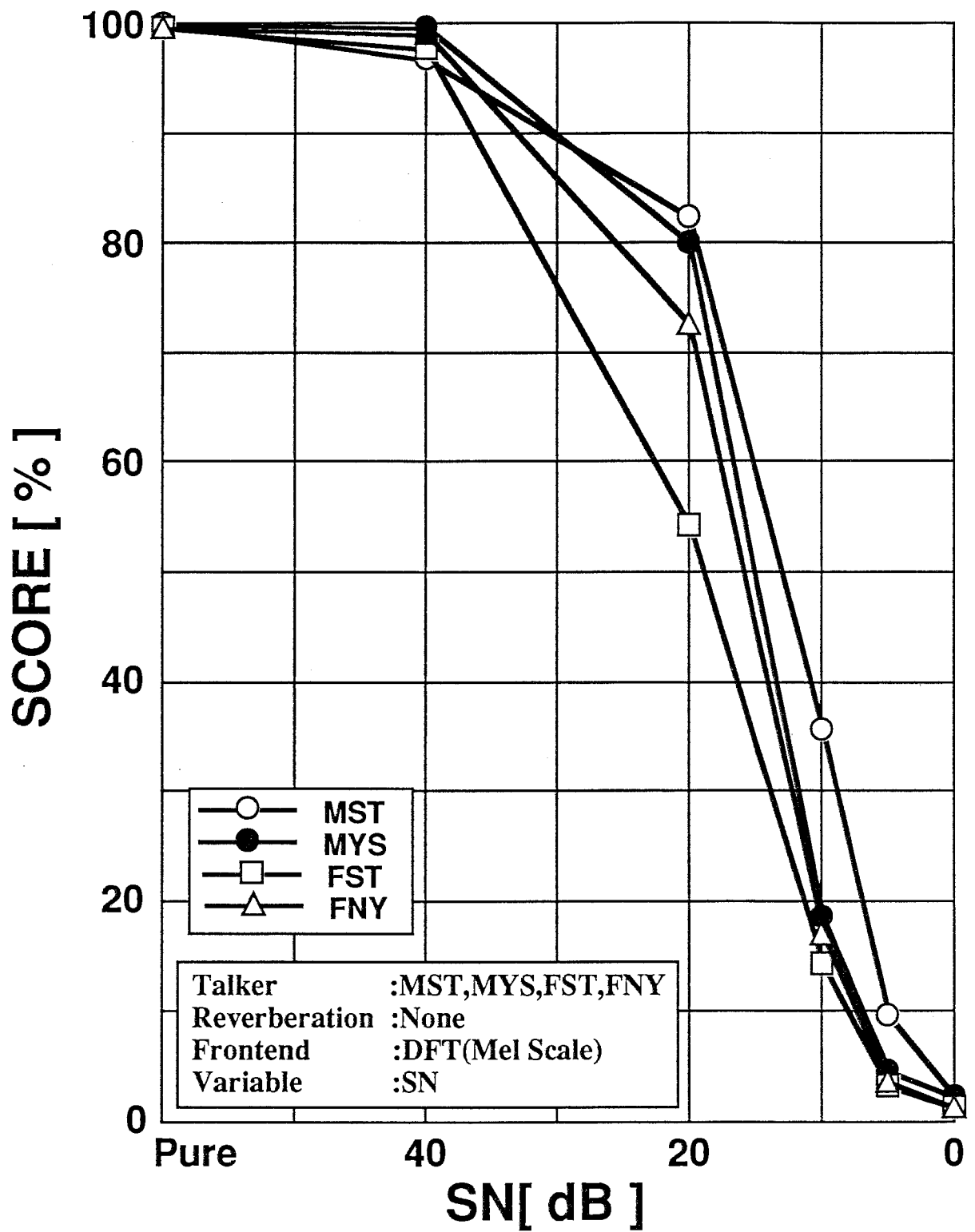
第2-7図 残響特性の畳み込み方法



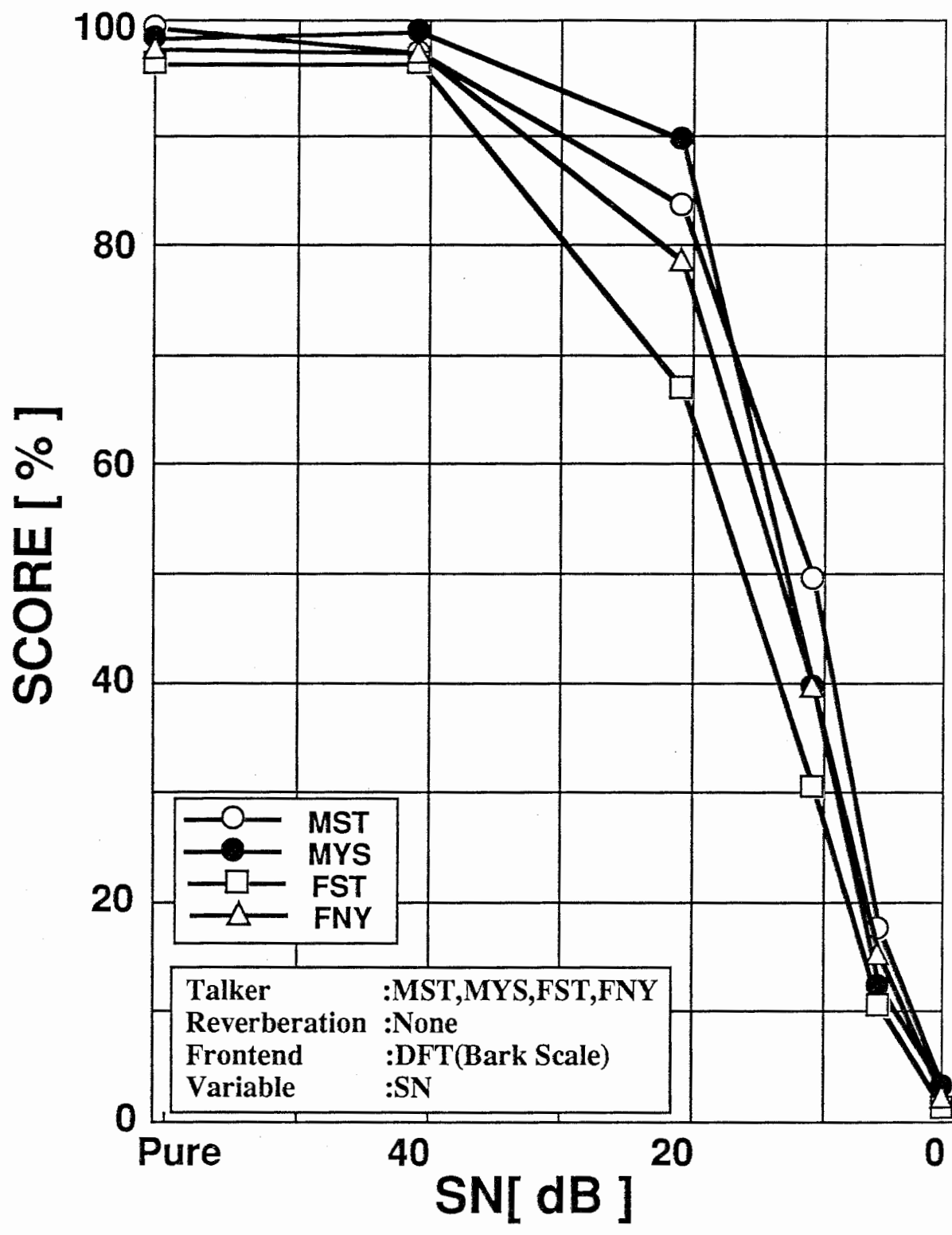
第3-1図 実験構成図



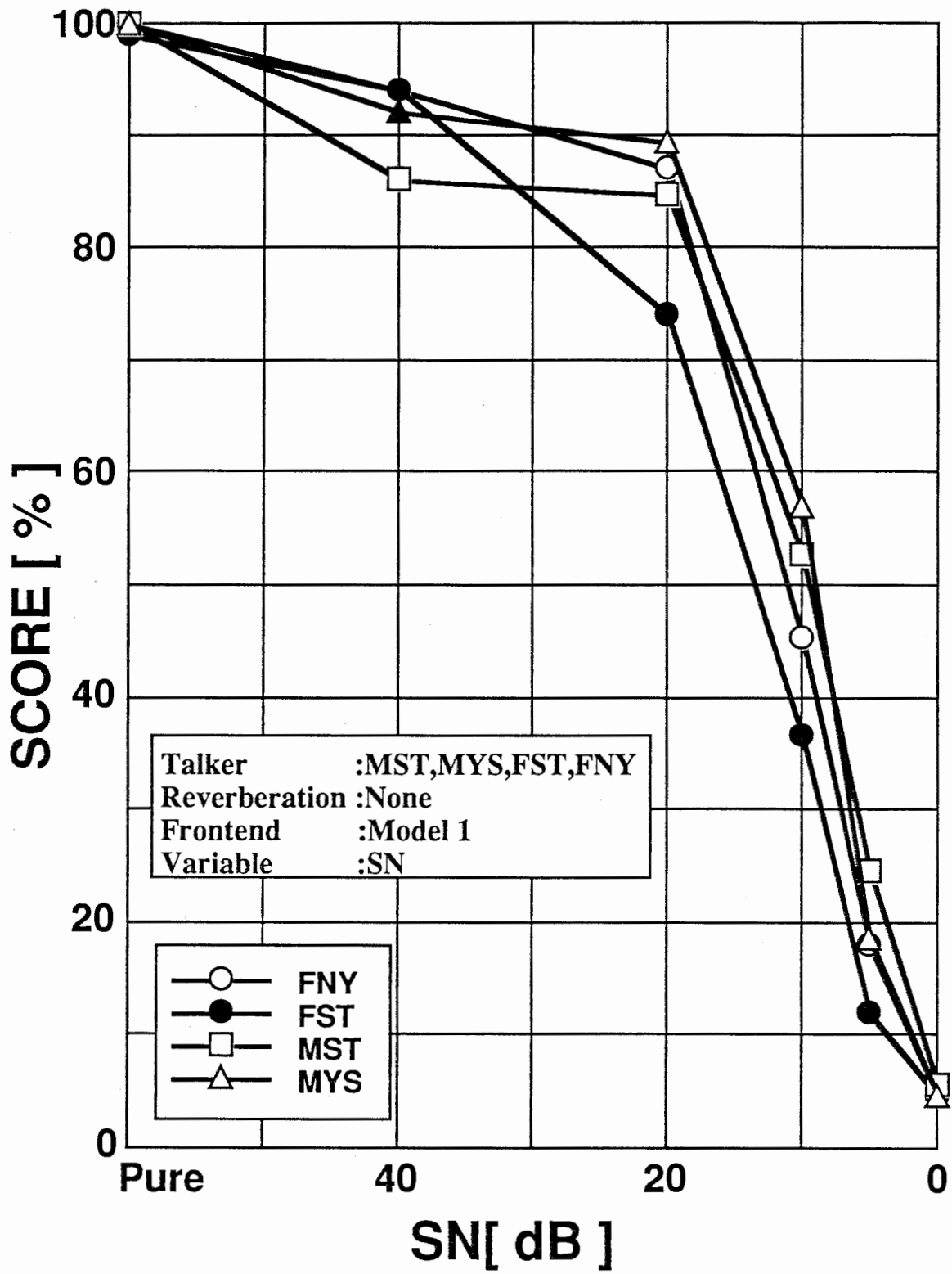
第3-2図 モデル2でのウィンドウの認識率への影響



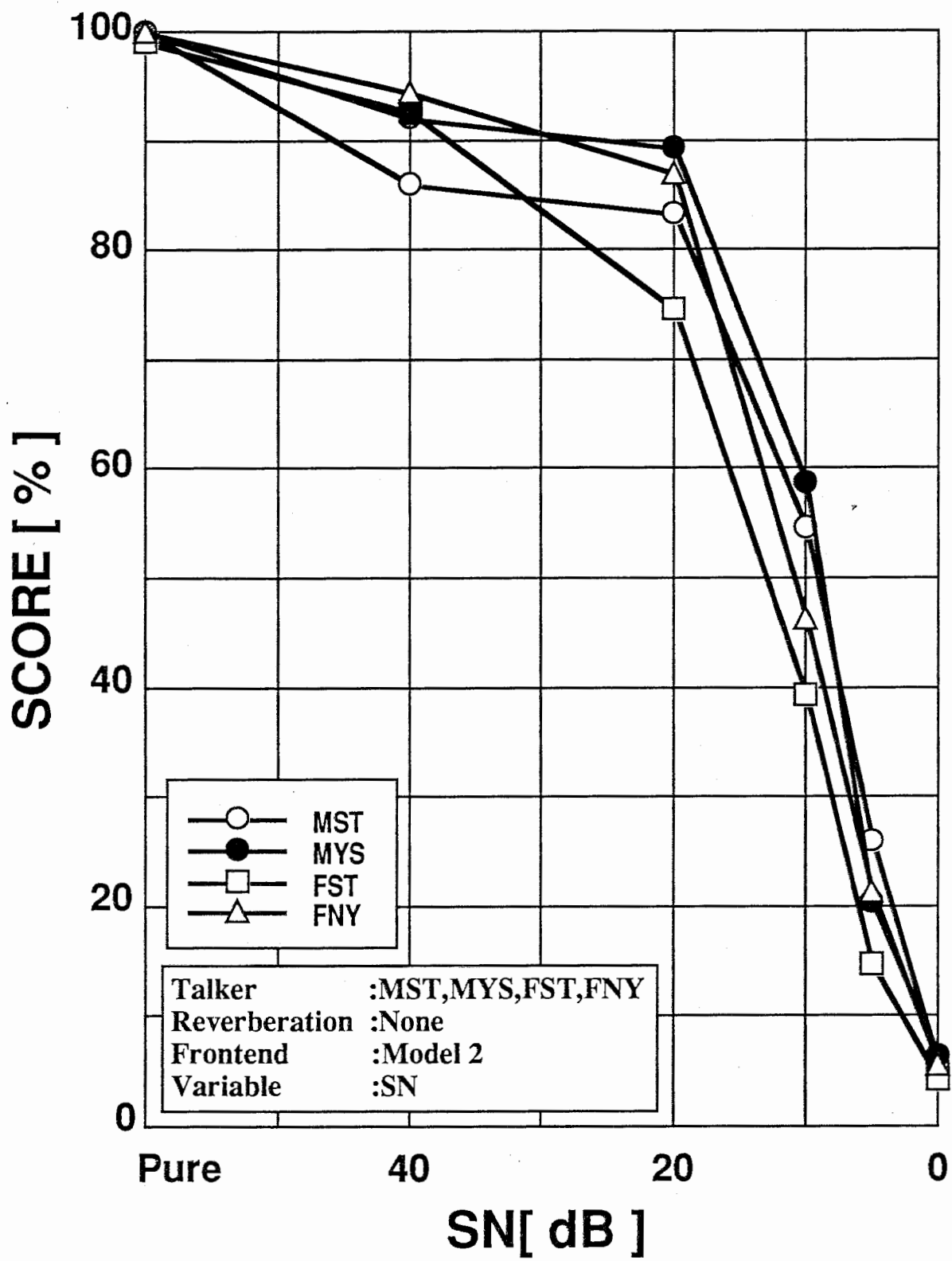
第3-3-1図 DFT(Mel) 特徴量による単語認識



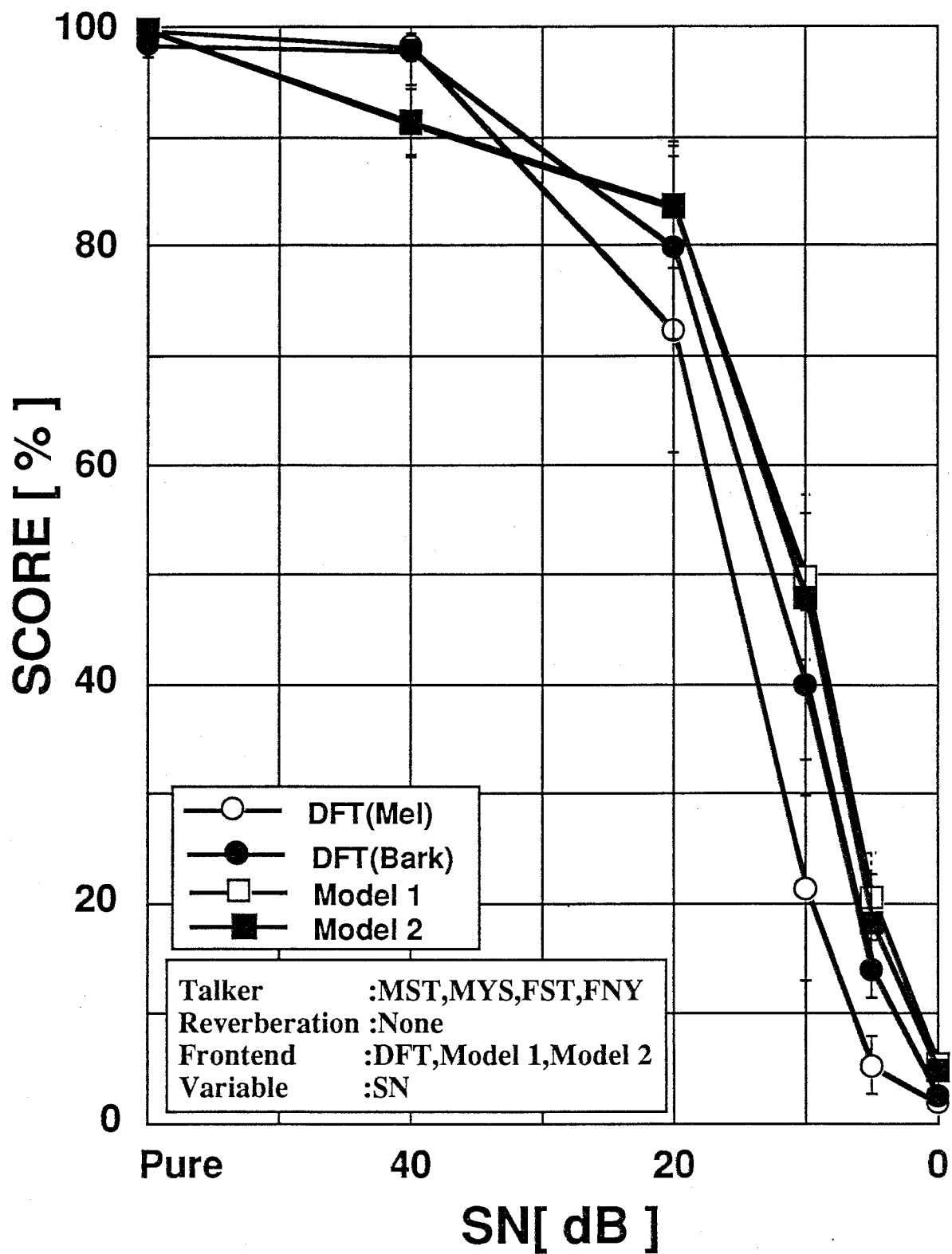
第3-3-2図 DFT(Bark) 特徴量による単語認識



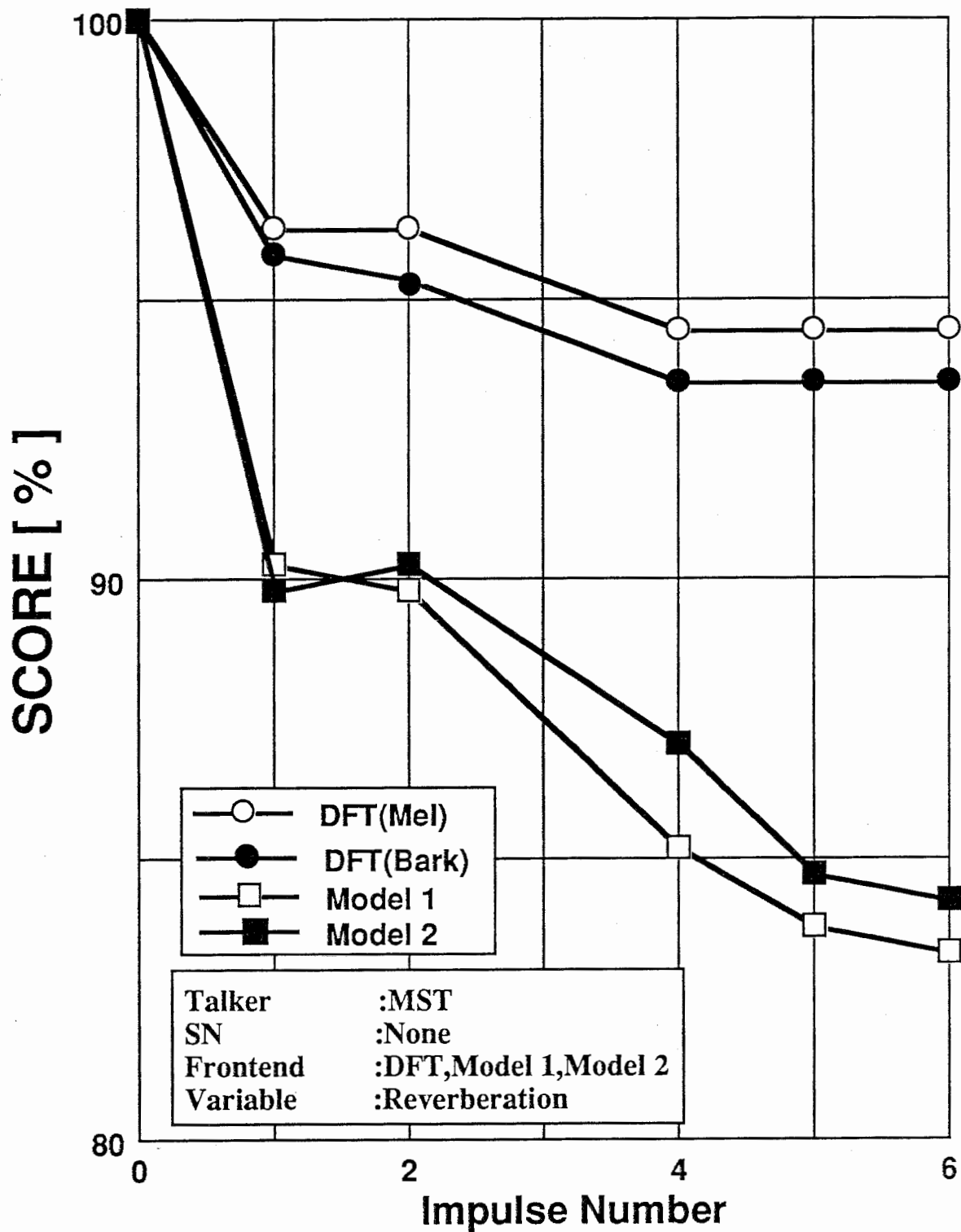
第3-4図 モデル1特徴量による単語認識



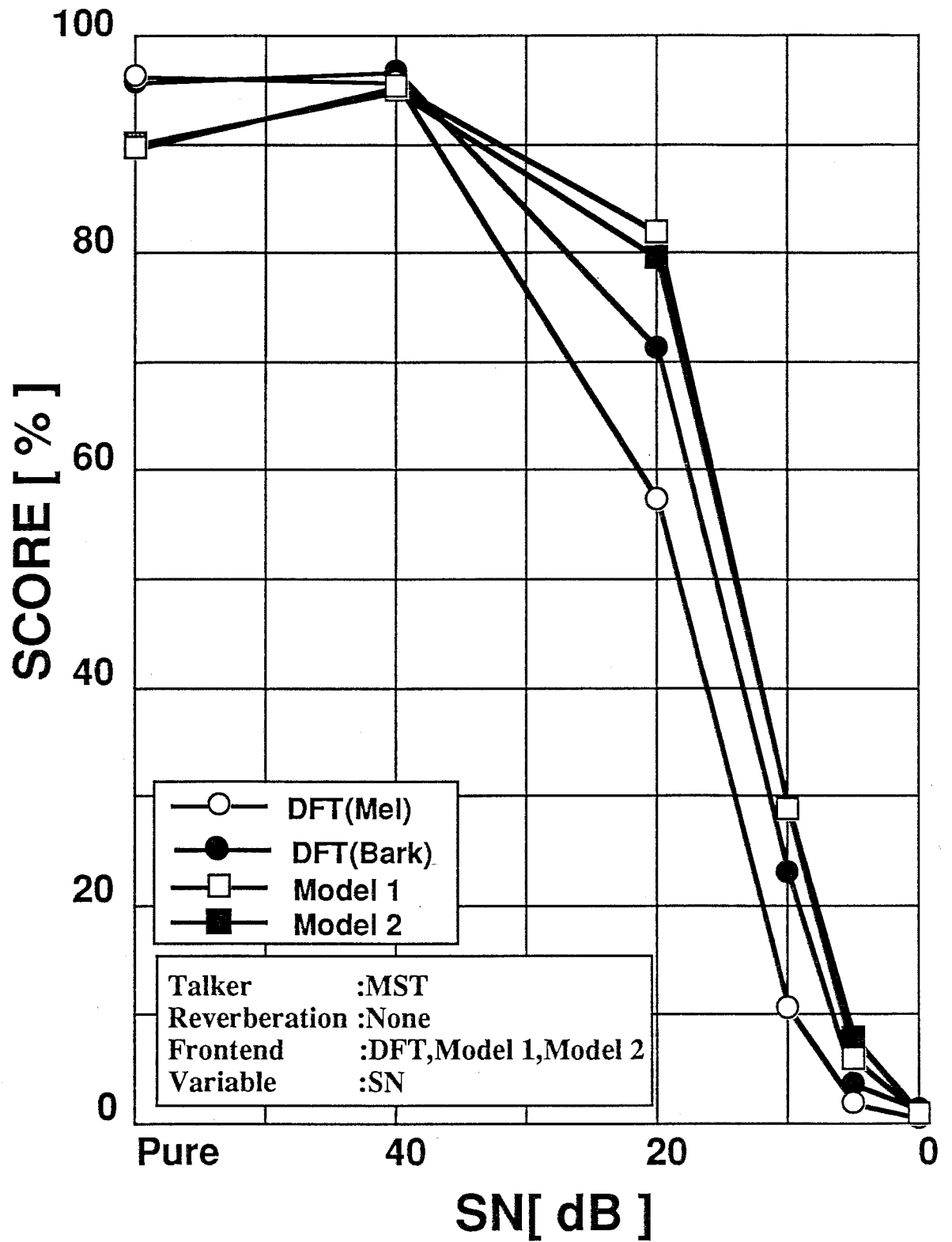
第3-5図 モデル2特徴量による単語認識



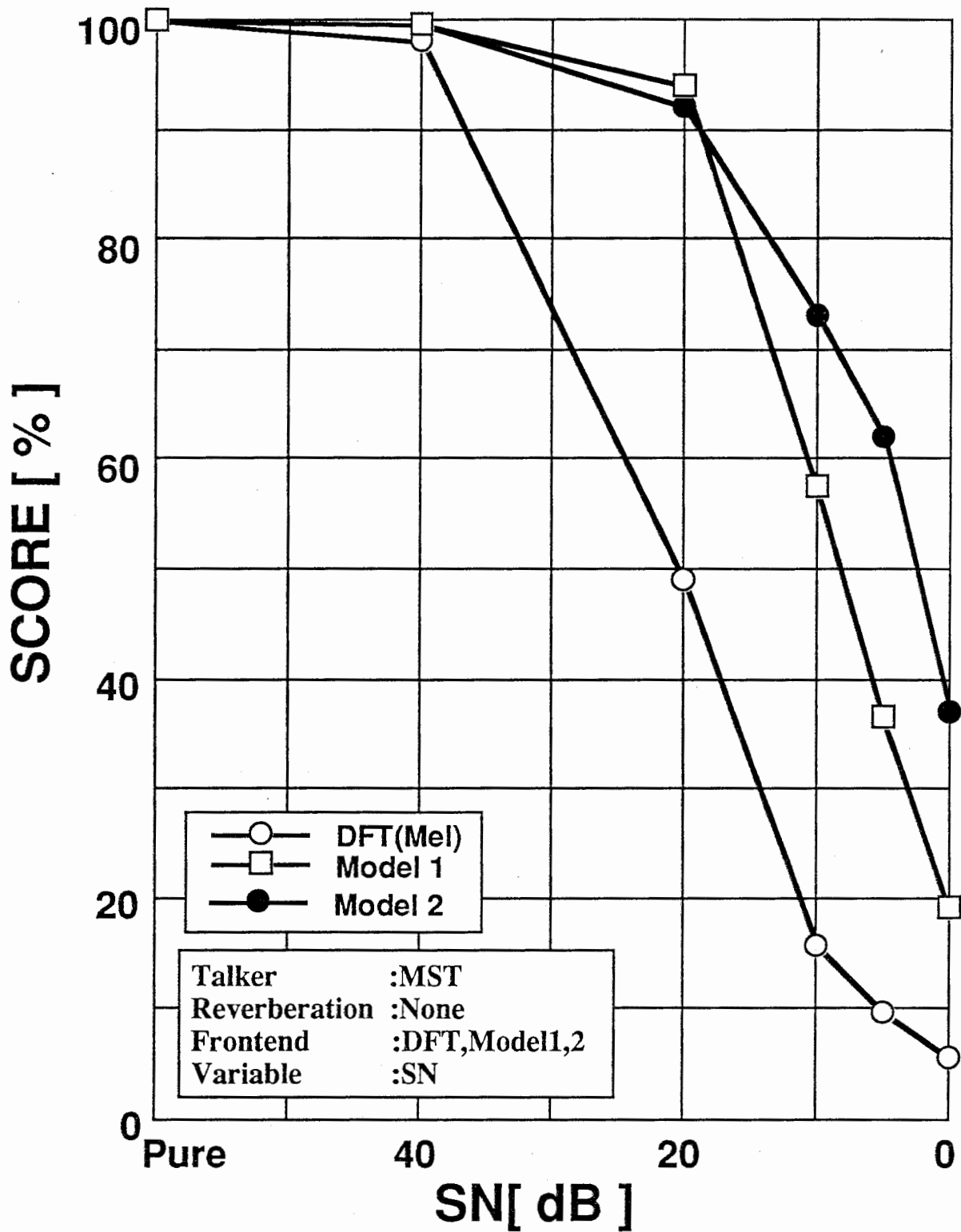
第3-6図 雑音下での異なるフロントエンド特徴量による単語認識(4話者平均)



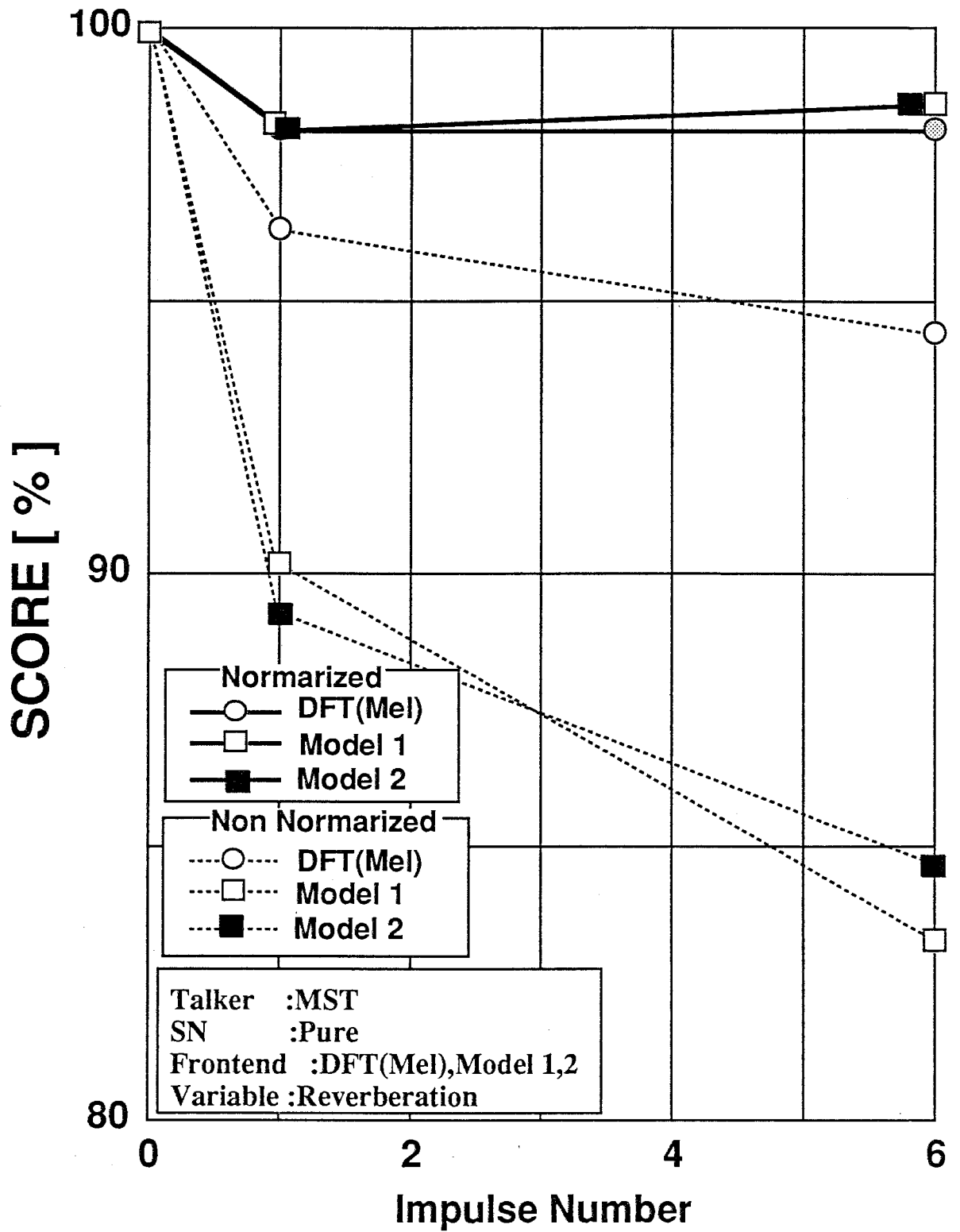
第3-7図 残響下での異なるフロントエンド特徴量による単語認識(4話者平均)



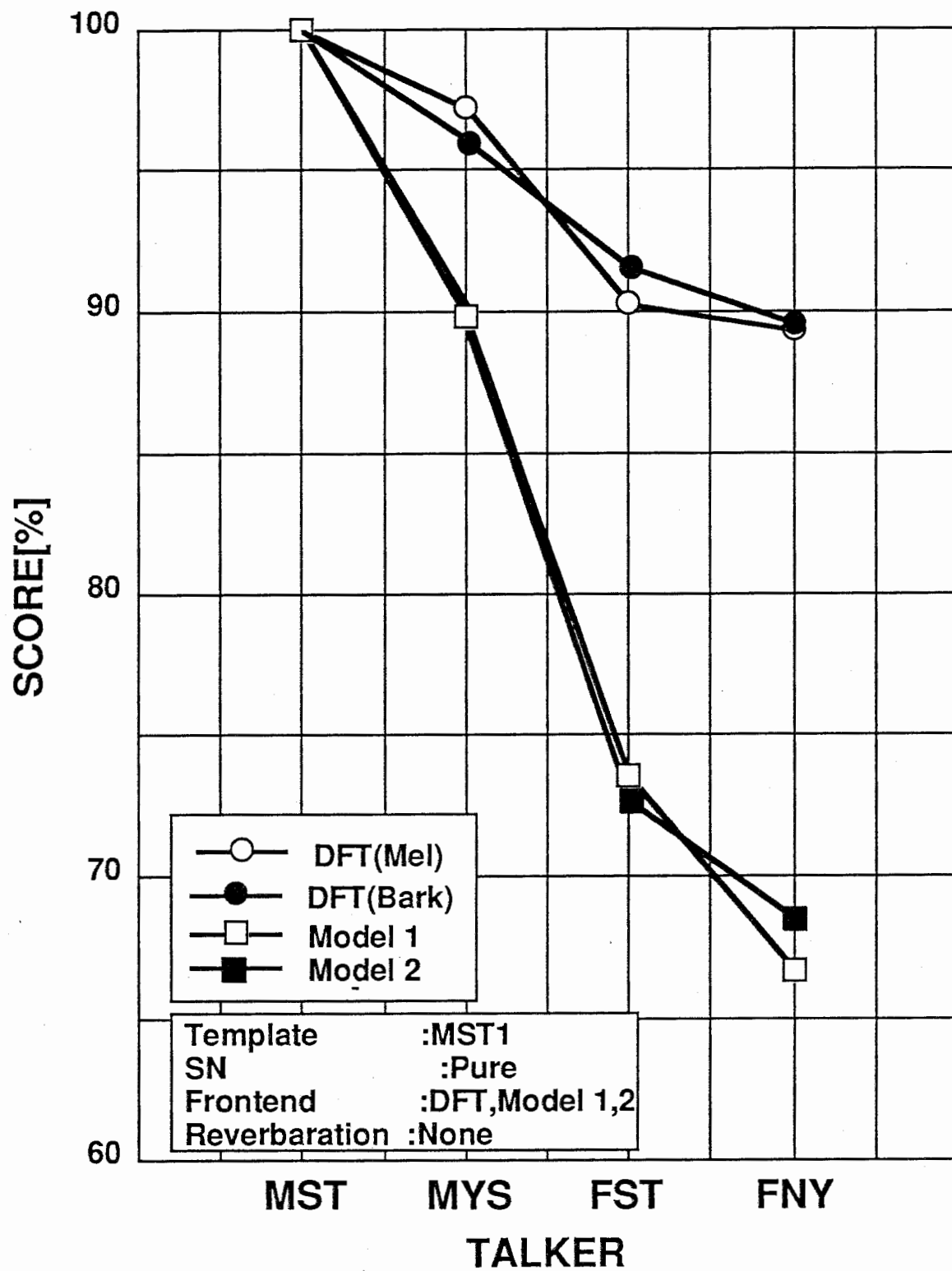
第3-8図 残響下,雑音下での異なるフロントエンド特徴量による単語認識(4話者平均)



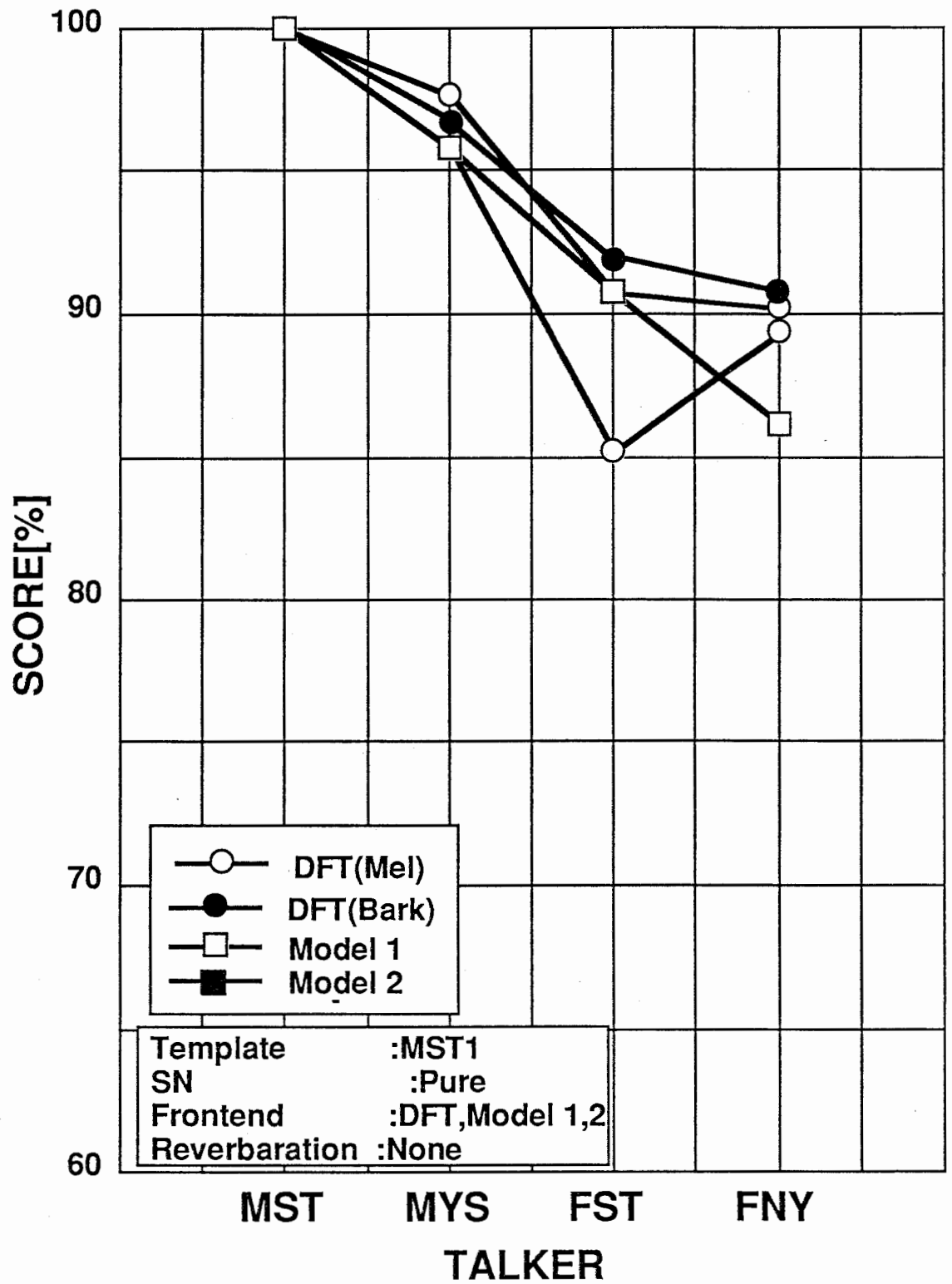
第3-9図 パワー正規化特徴量による単語認識



第 3 - 1 0 図 残響下での単語認識(正規化)

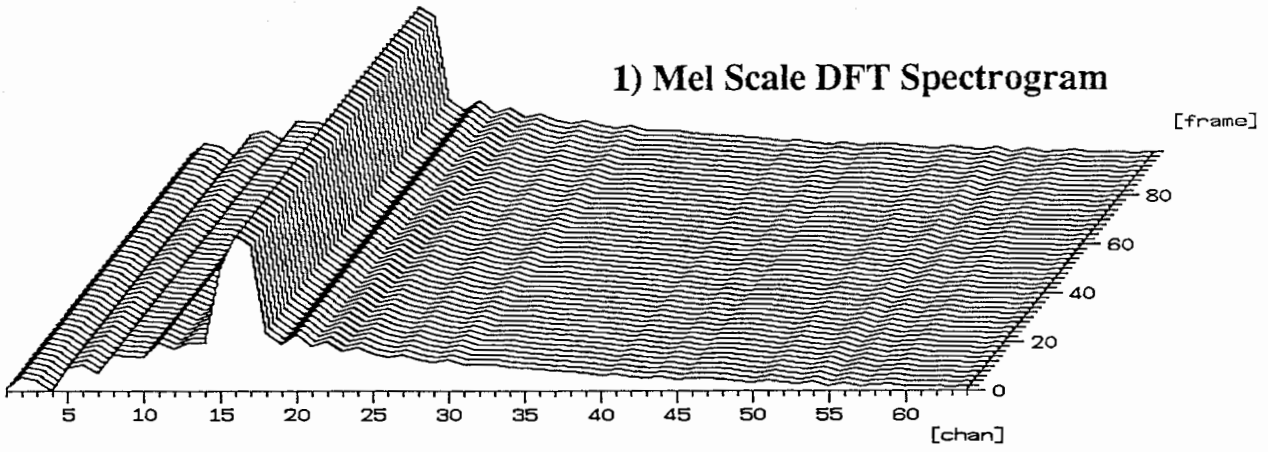


第3-11図 異なる話者間の単語認識

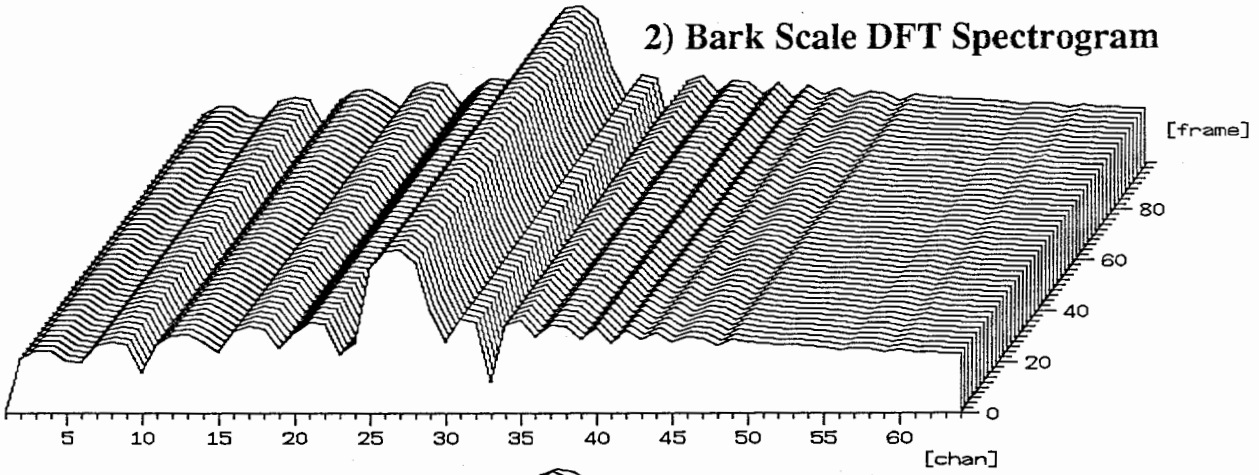


第3-12図 異なる話者間の単語認識 (特徴量正規化)

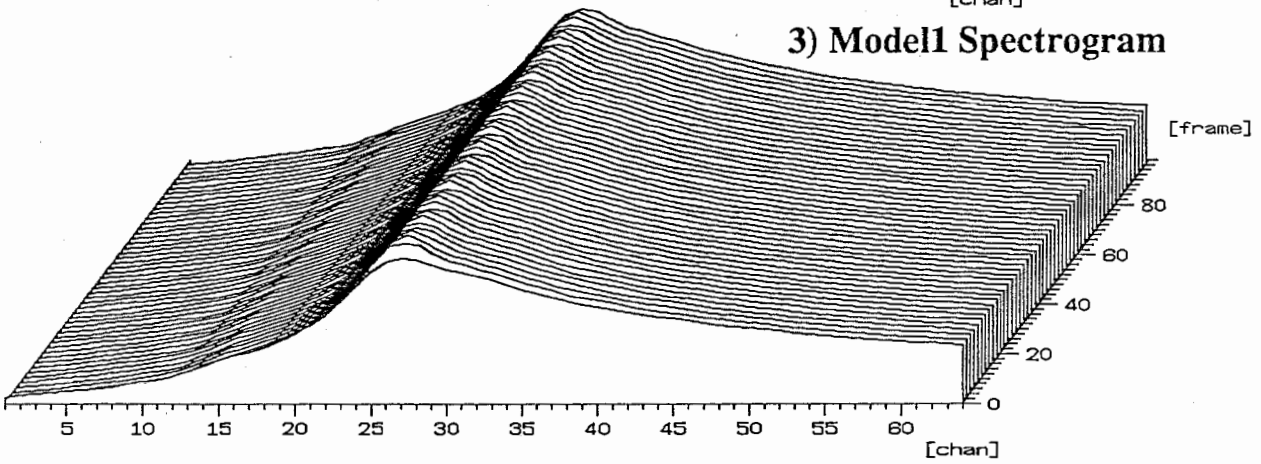
1) Mel Scale DFT Spectrogram



2) Bark Scale DFT Spectrogram

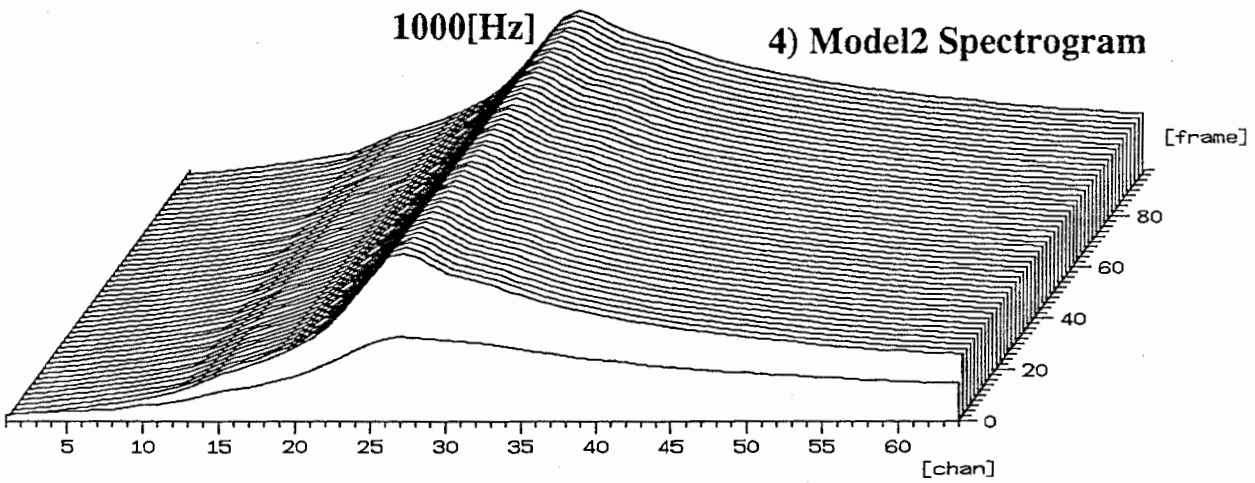


3) Model1 Spectrogram

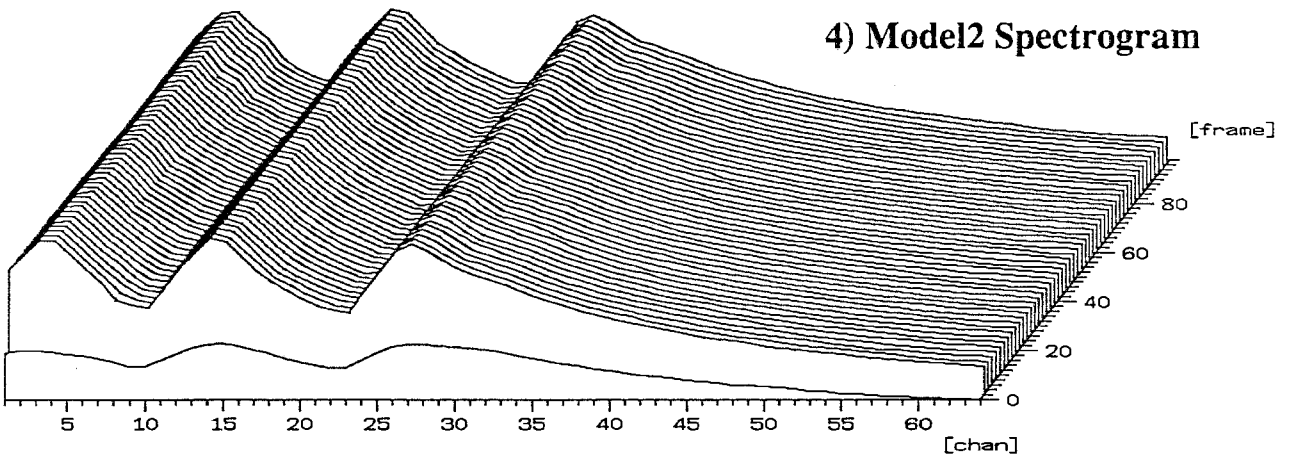
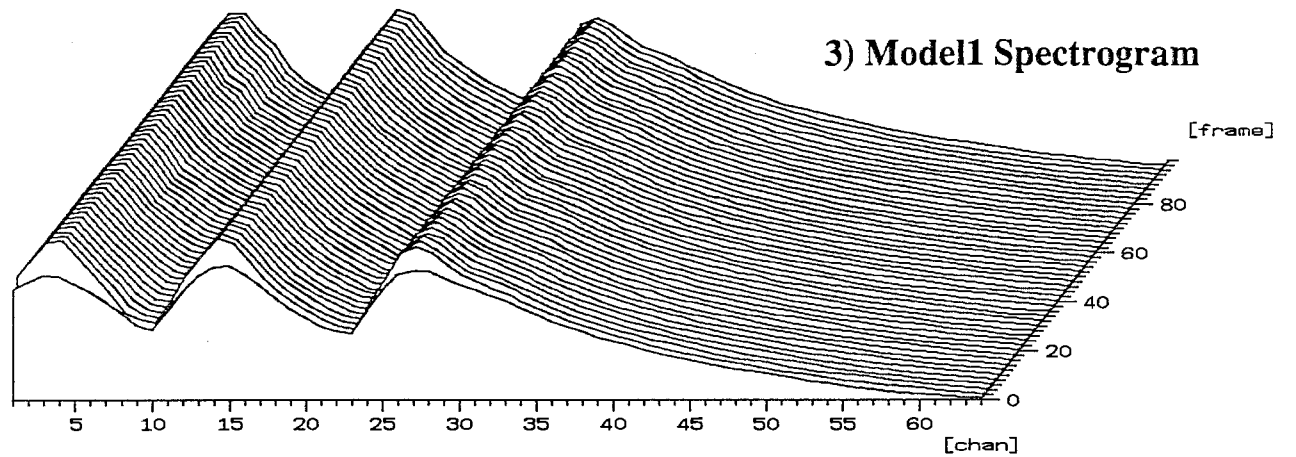
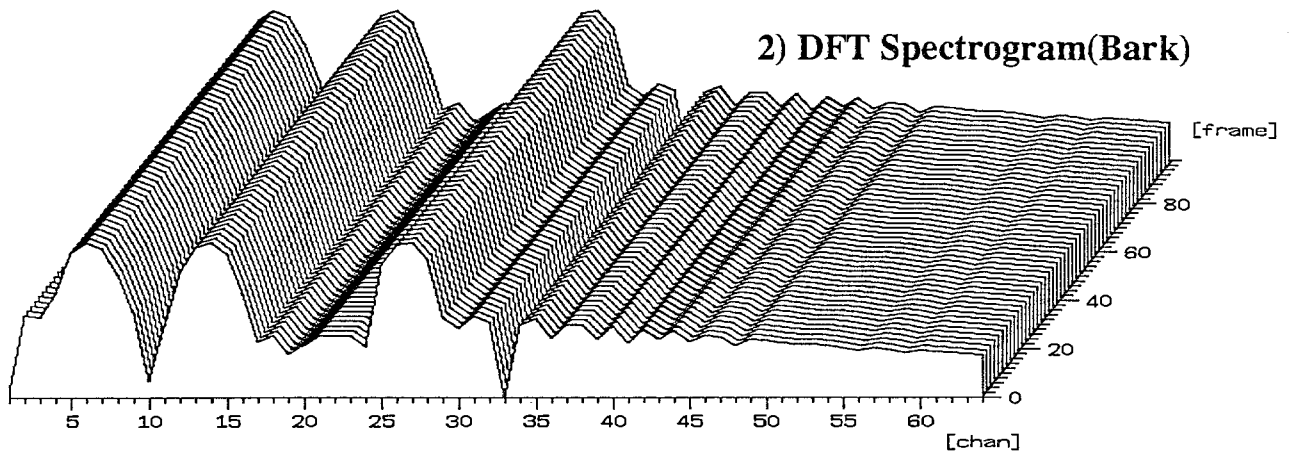
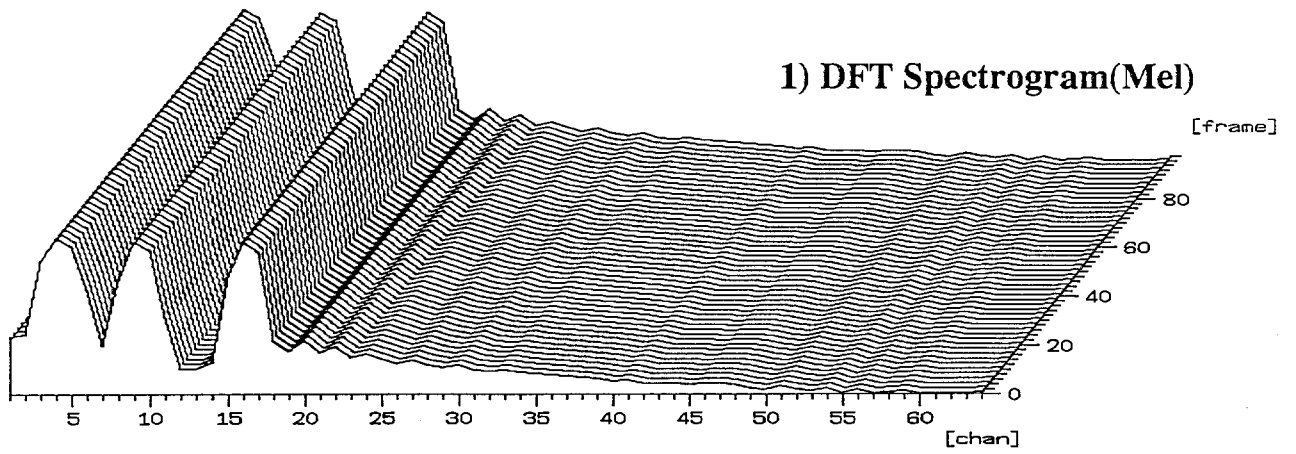


1000[Hz]

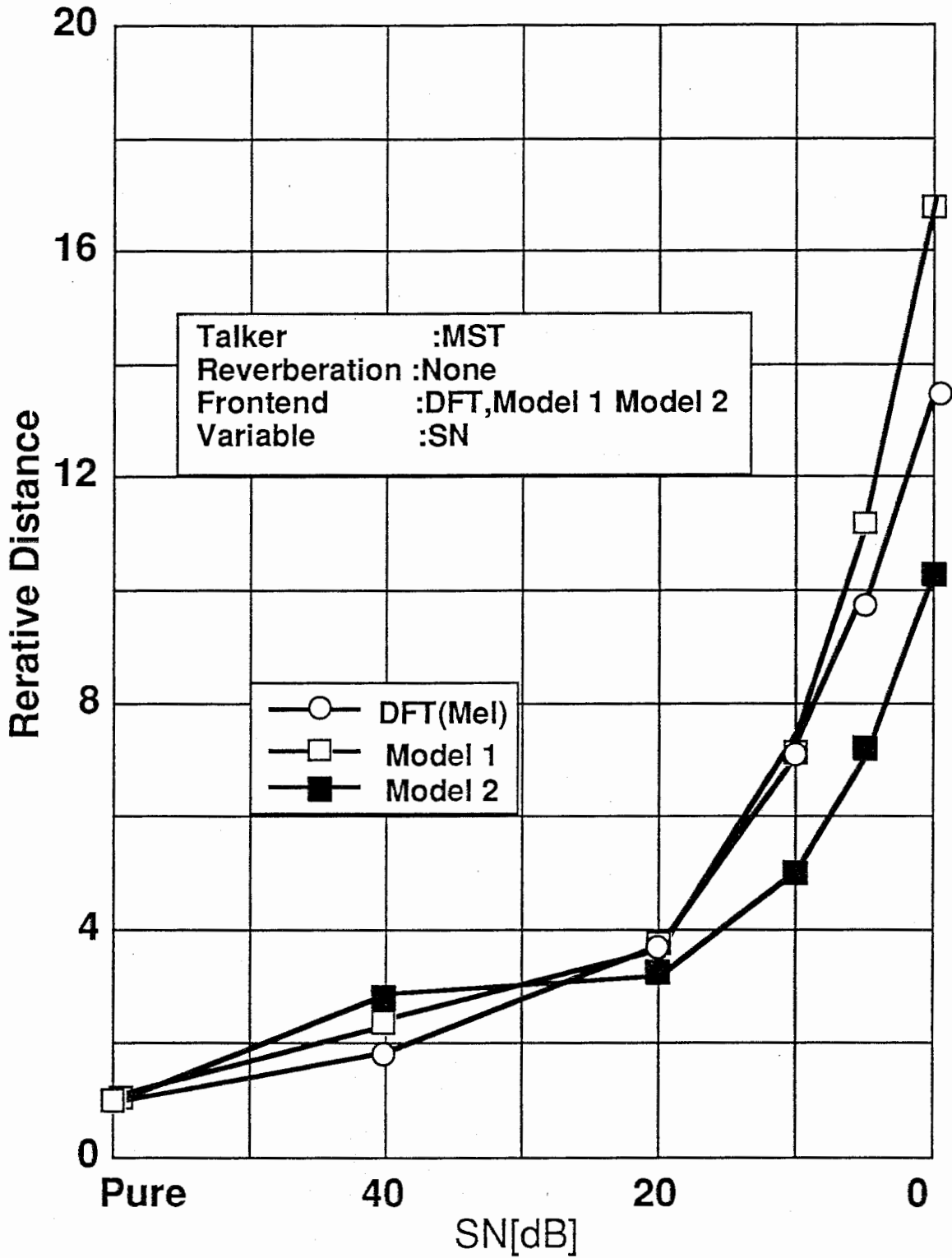
4) Model2 Spectrogram



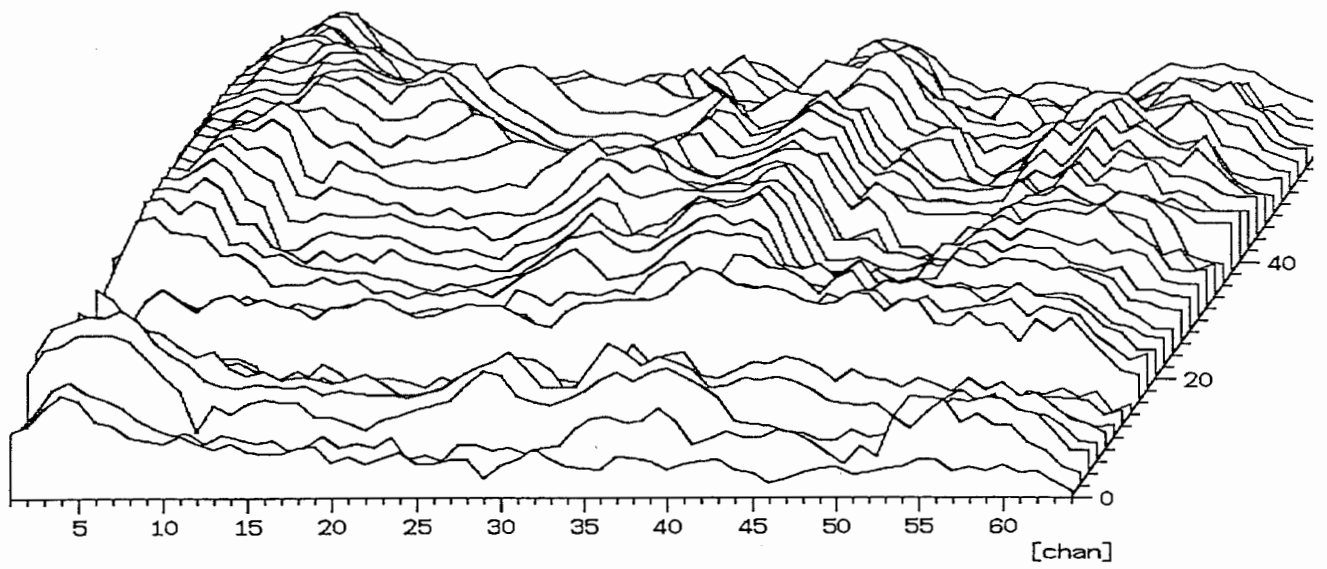
第4-1図 DFT、モデル1、モデル2スペクトログラム



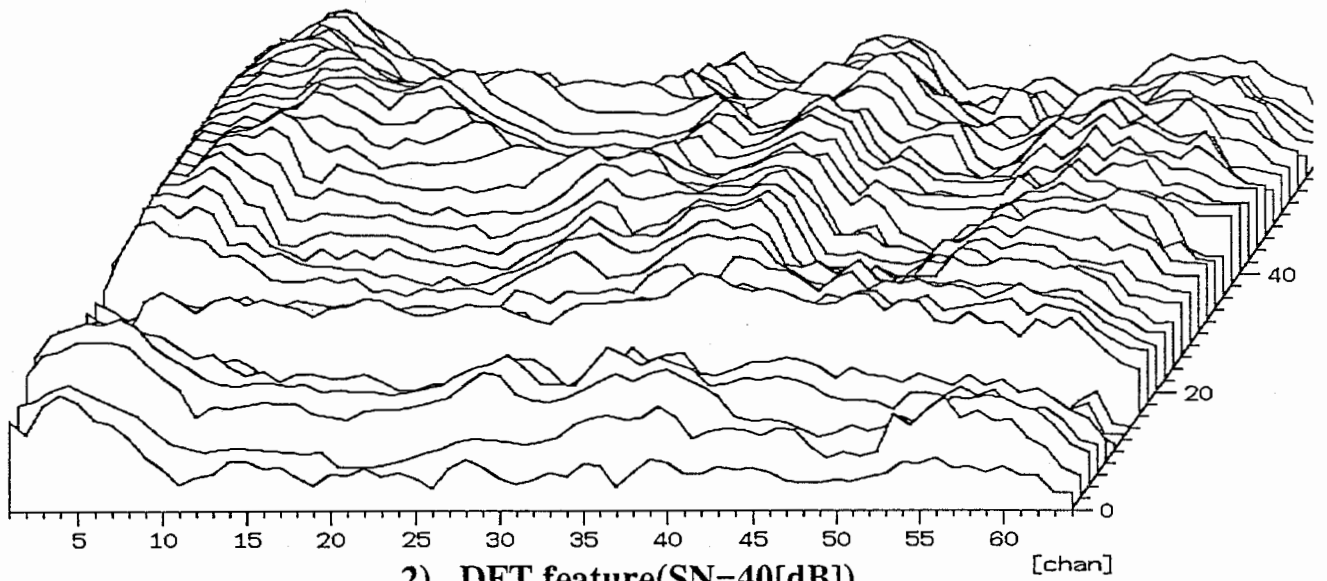
第4-2図 200+500+1000 [Hz] のスペクトログラム



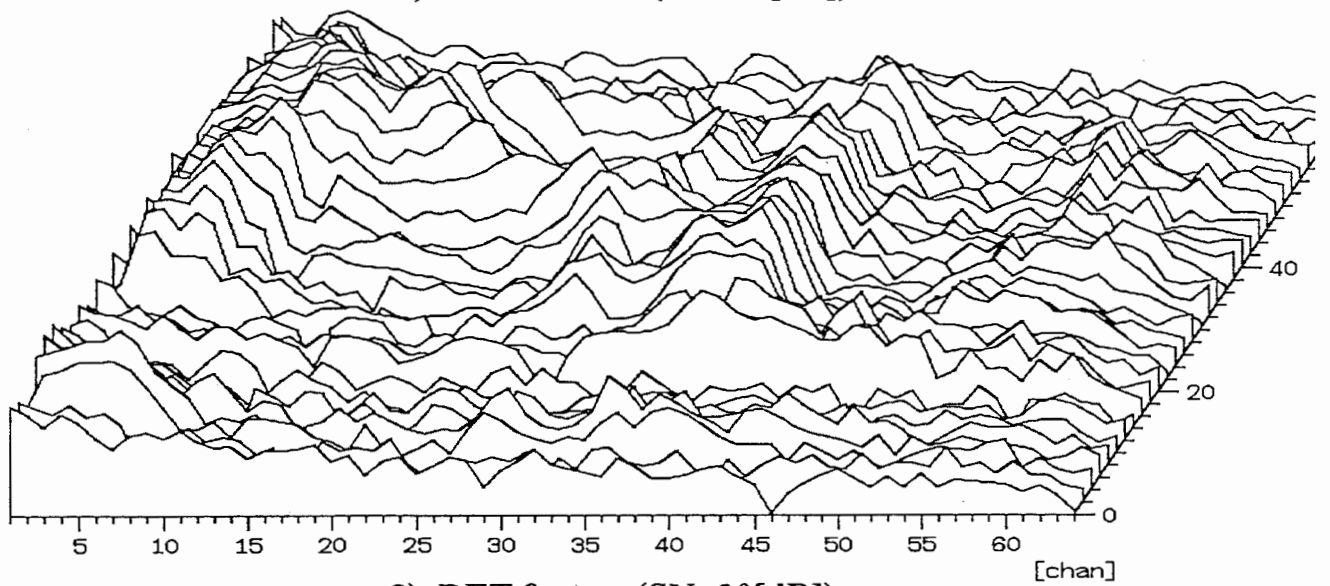
第4-3図 ノイズによるテンプレート単語間距離変動



1) DFT feature(Pure)

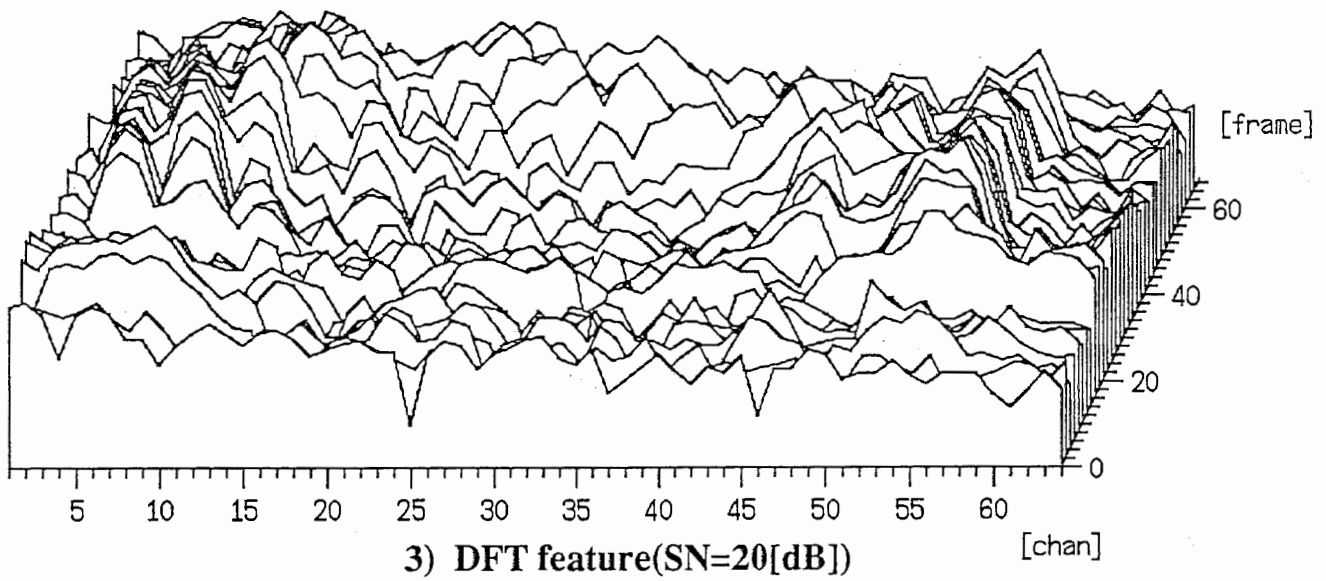
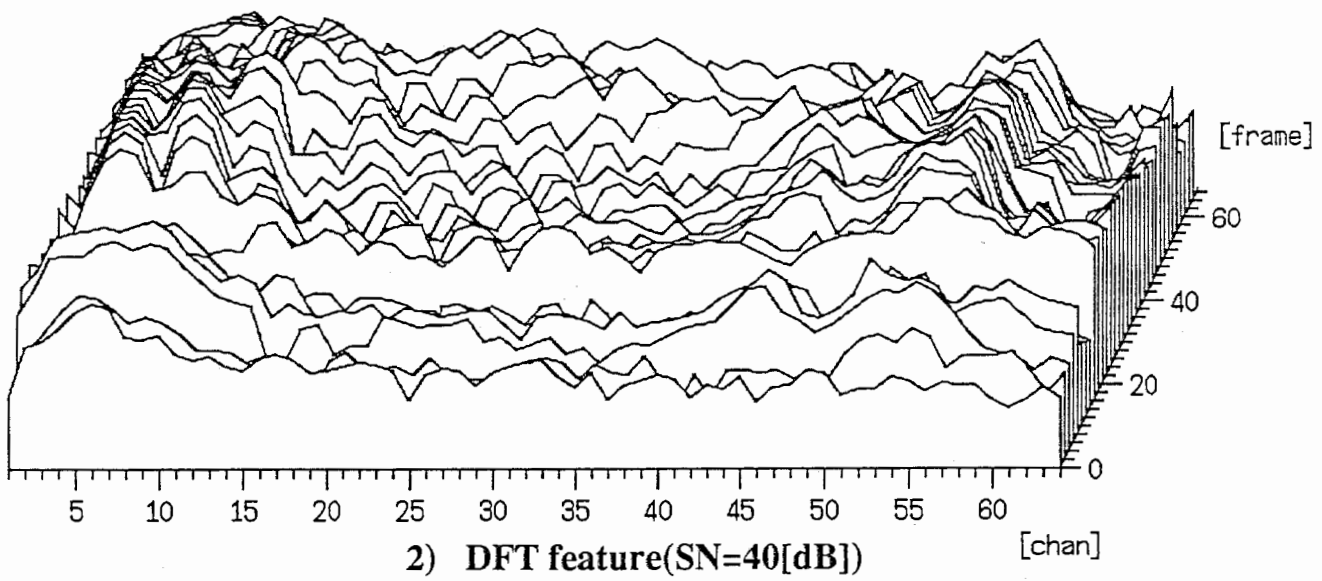
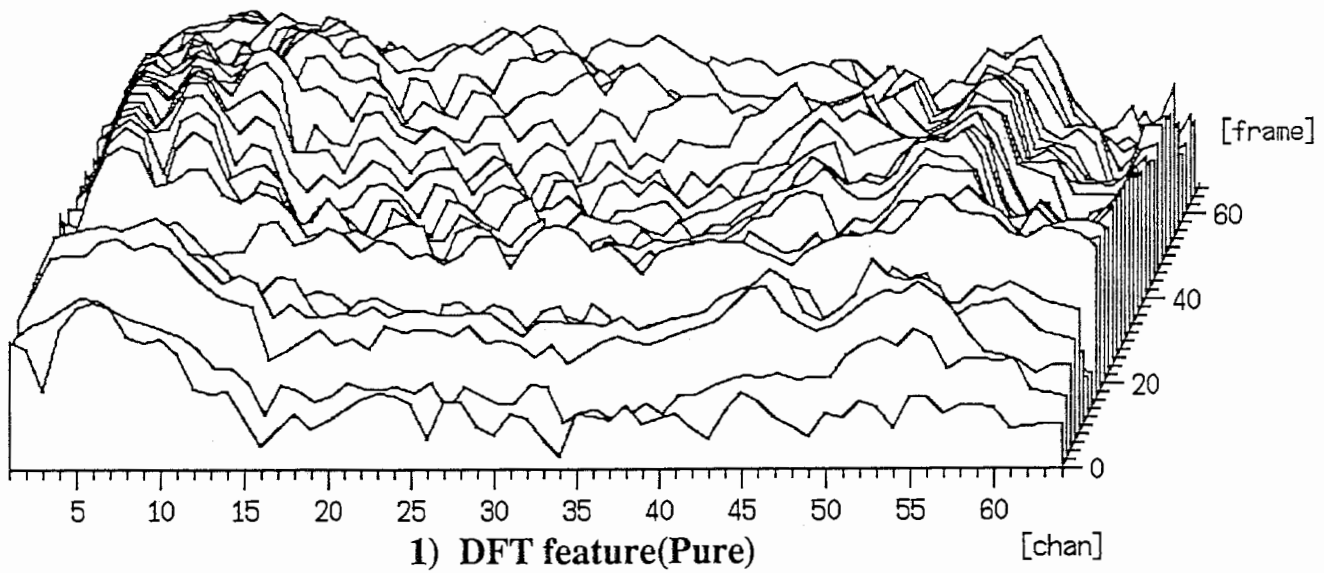


2) DFT feature(SN=40[dB])

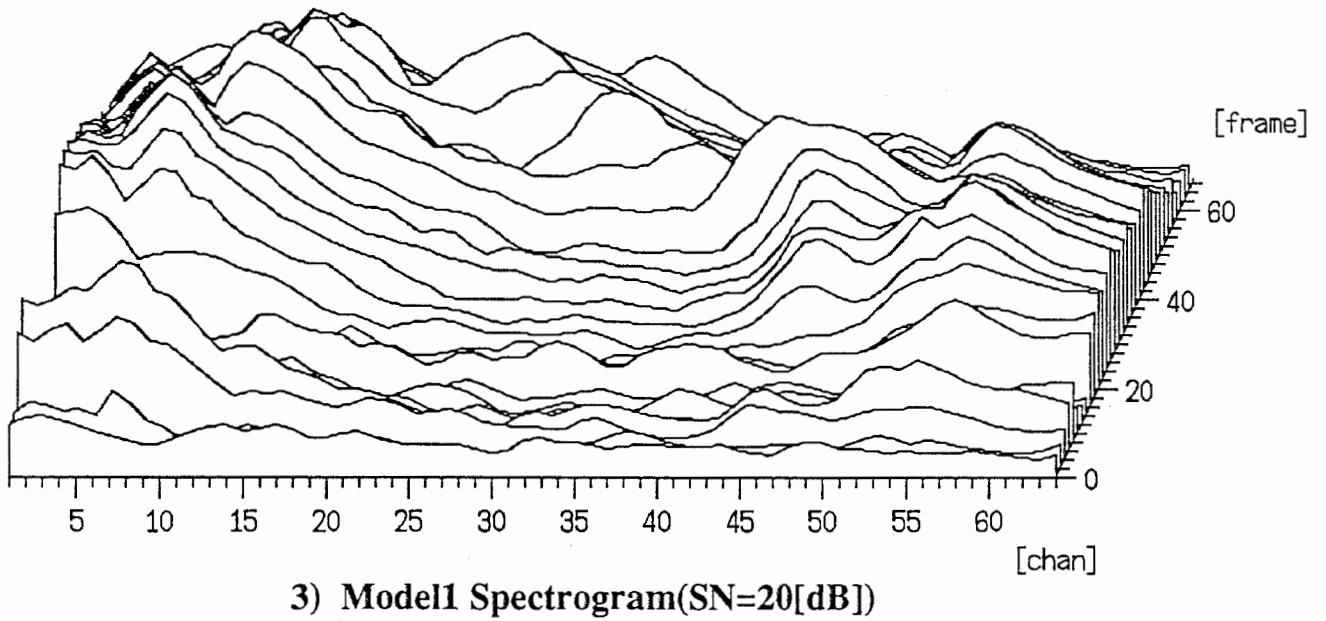
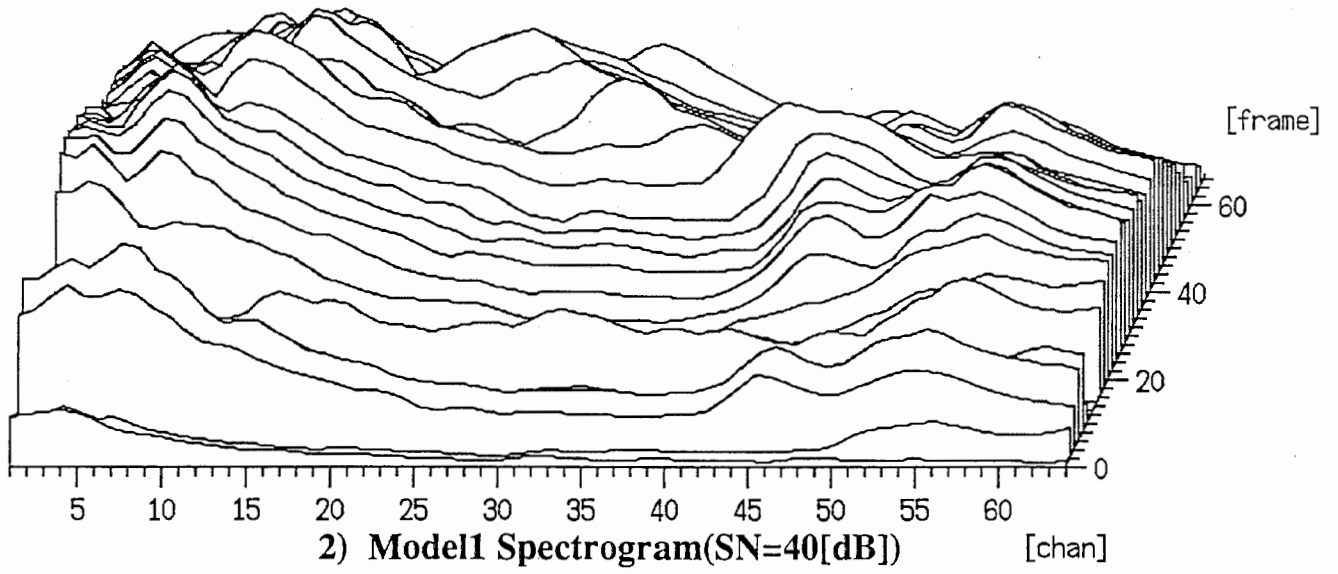
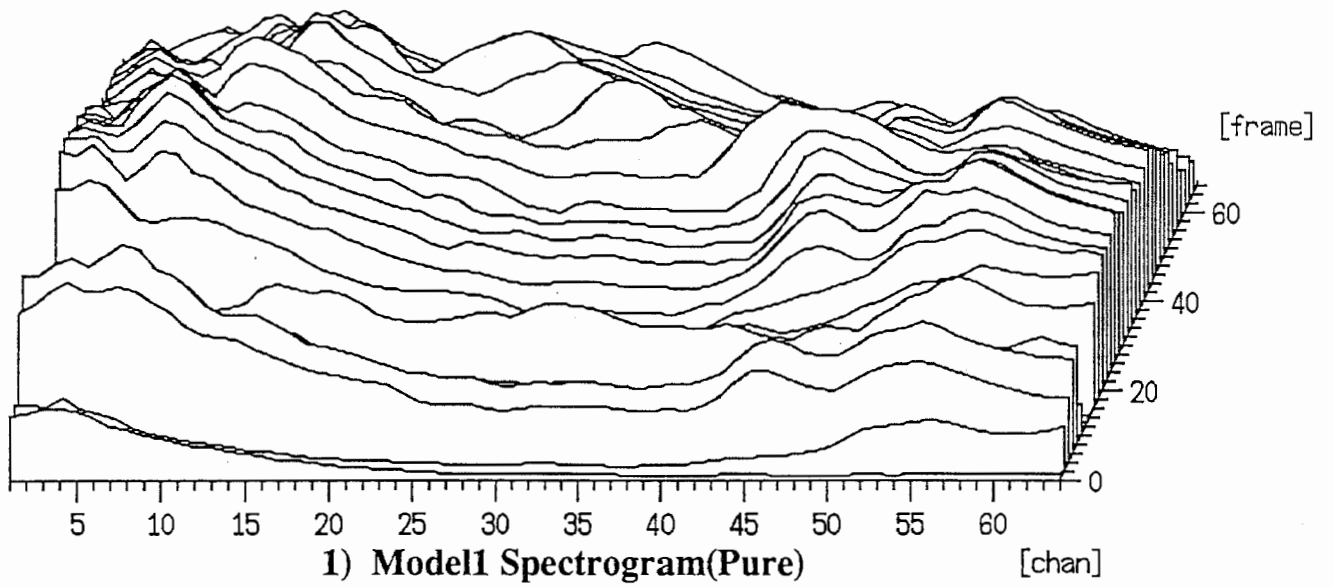


3) DFT feature(SN=20[dB])

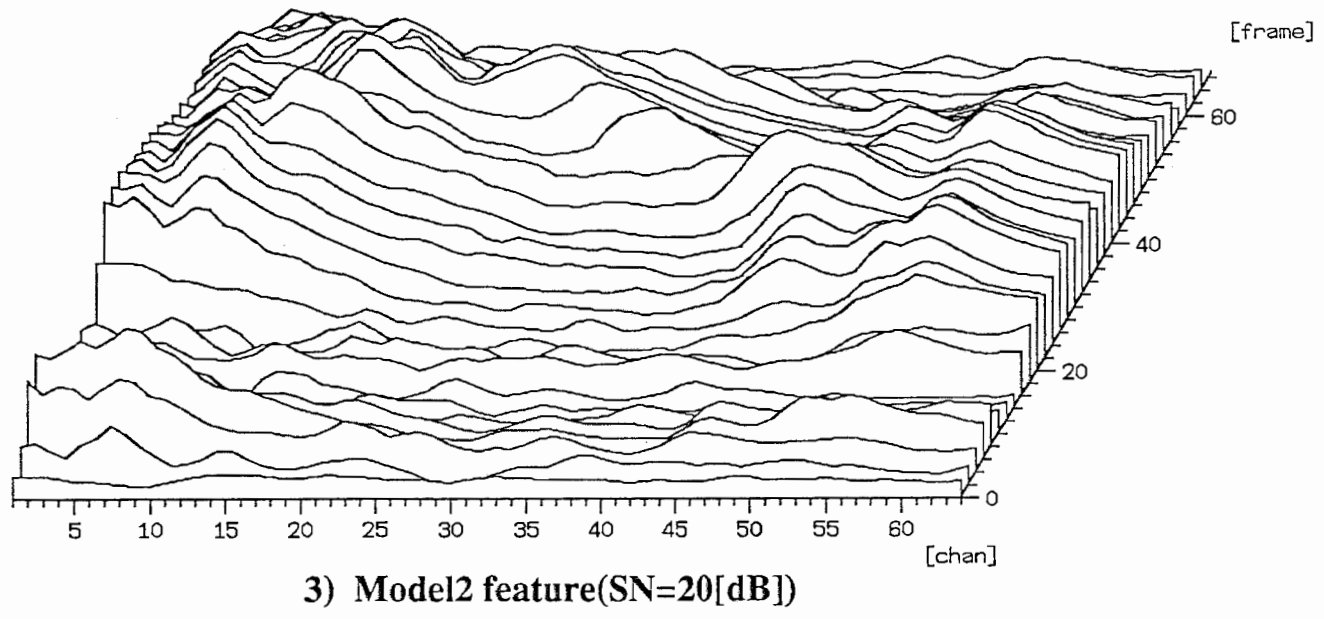
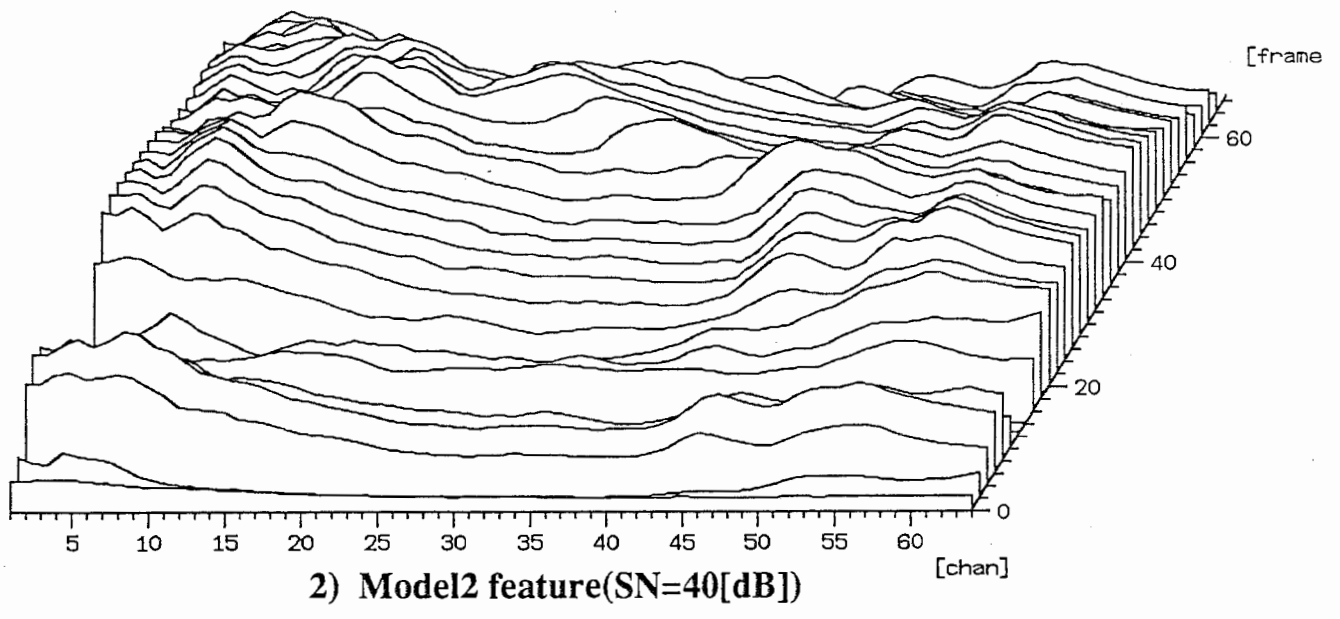
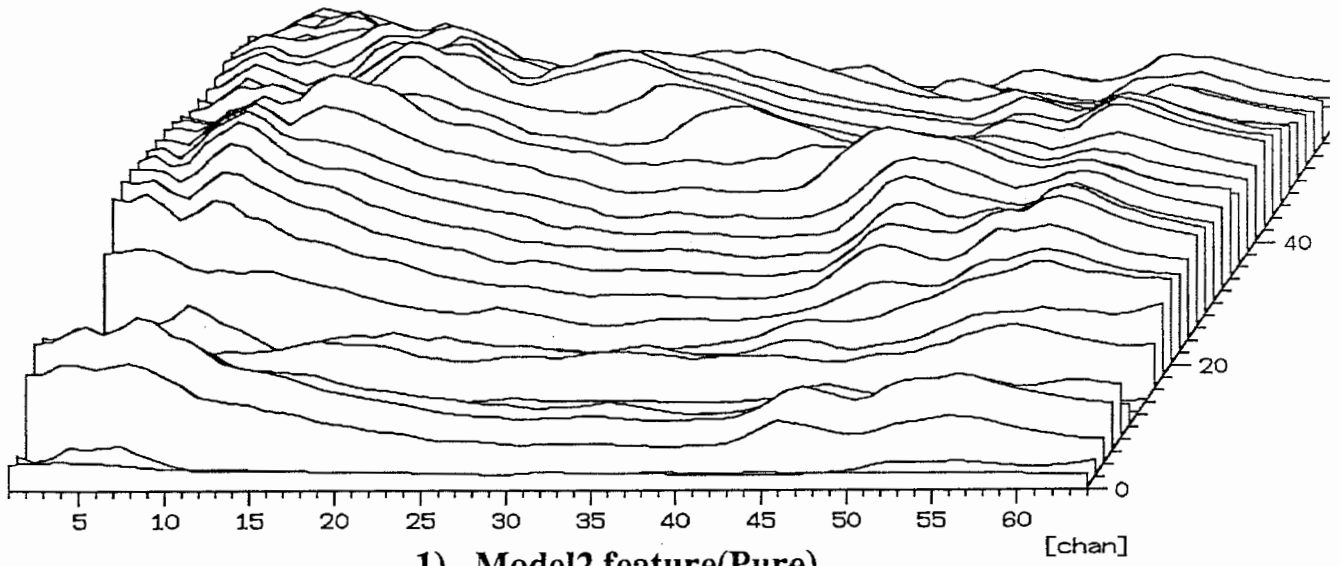
第4-4-1図 SN変動によるスペクトルの変化 (Mel DFT)



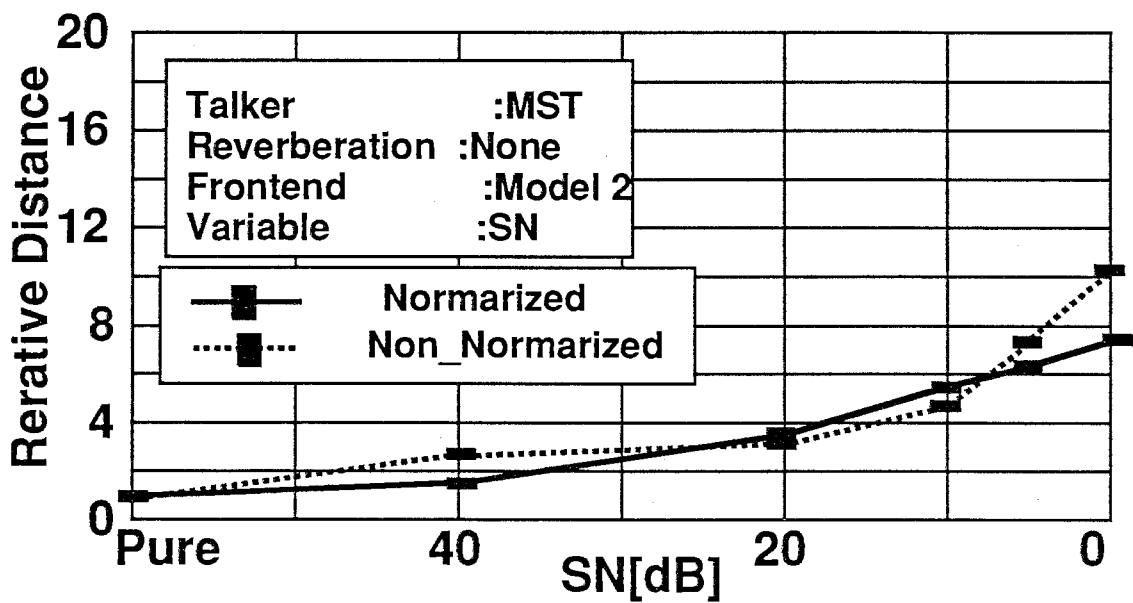
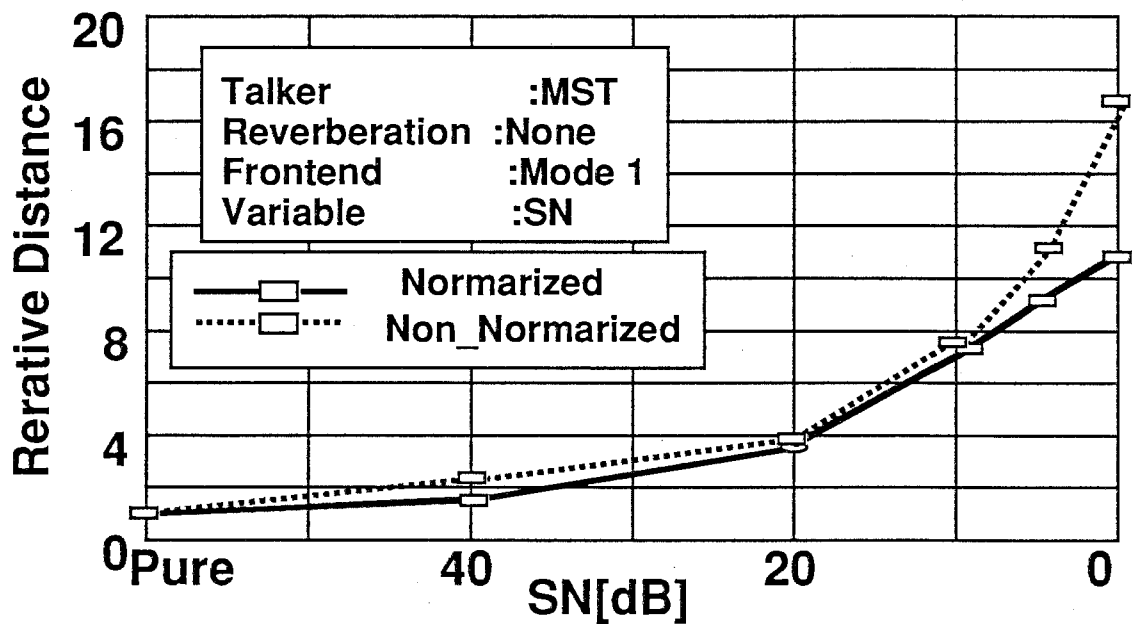
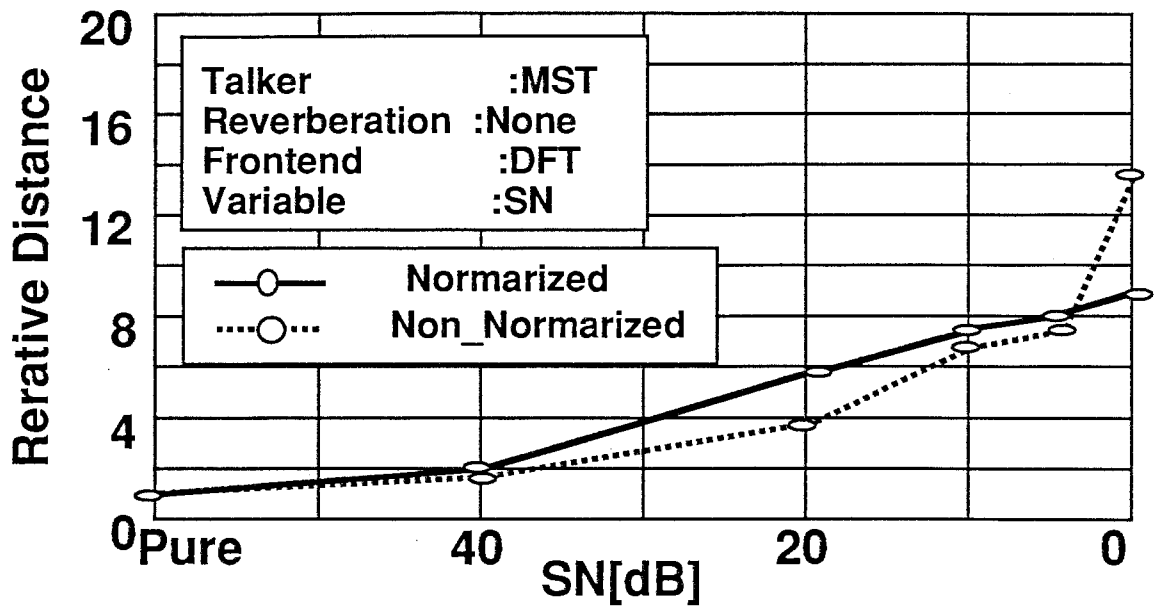
第4-4-2図 SN変動によるスペクトルの変化 (Bark DFT)



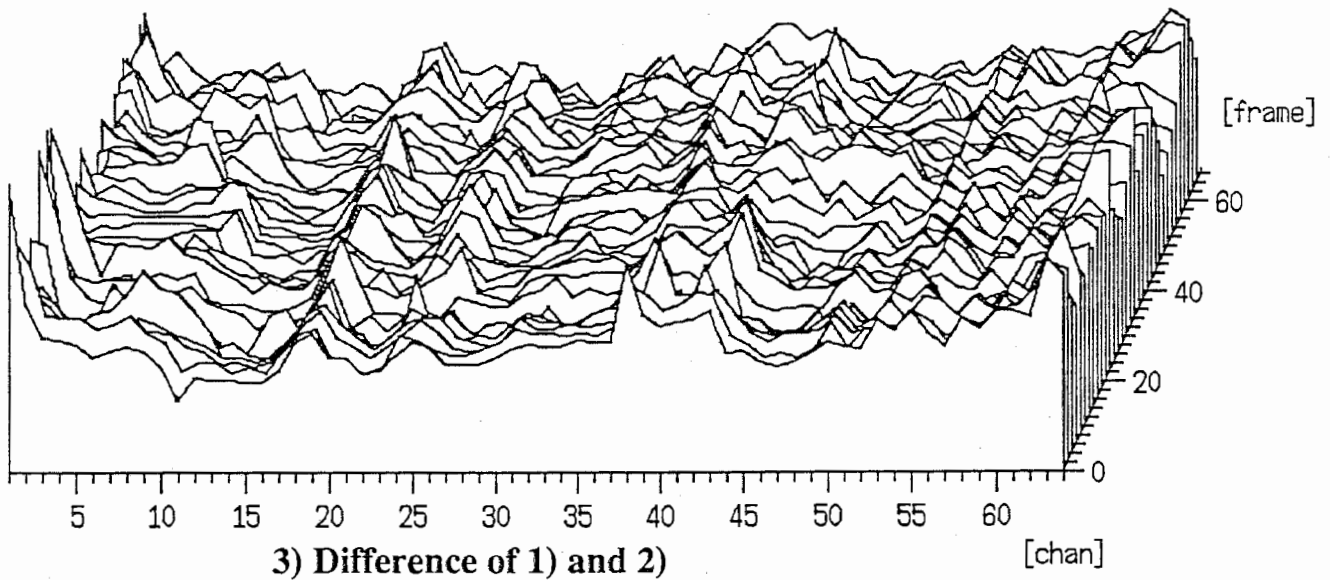
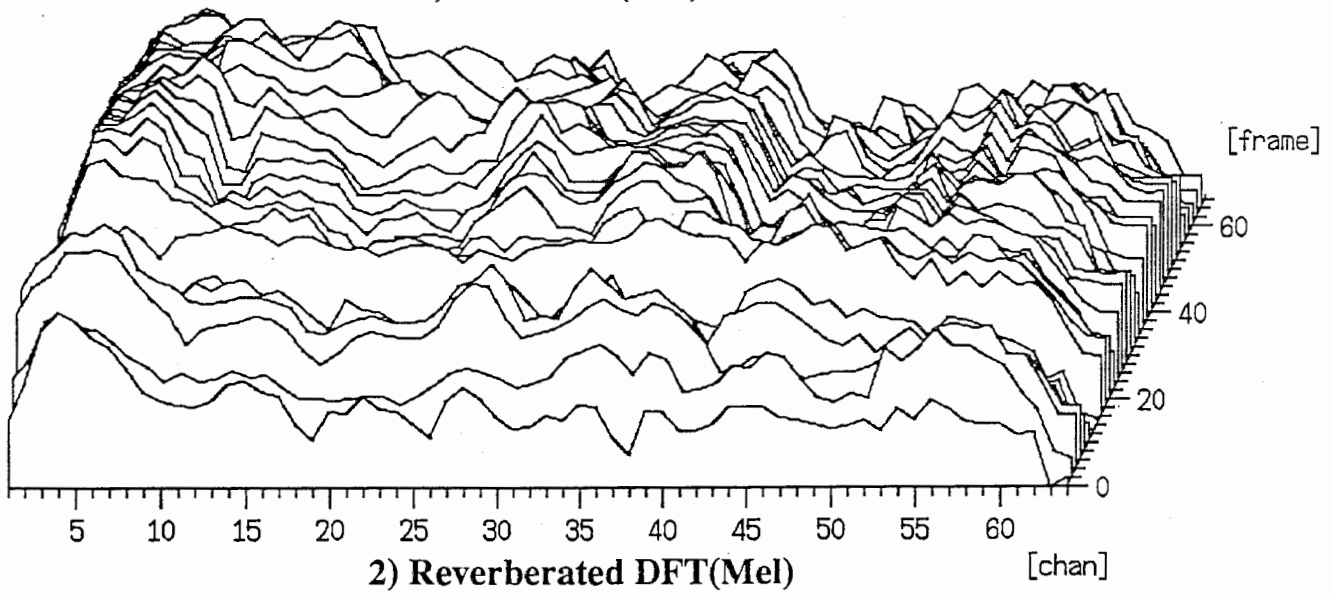
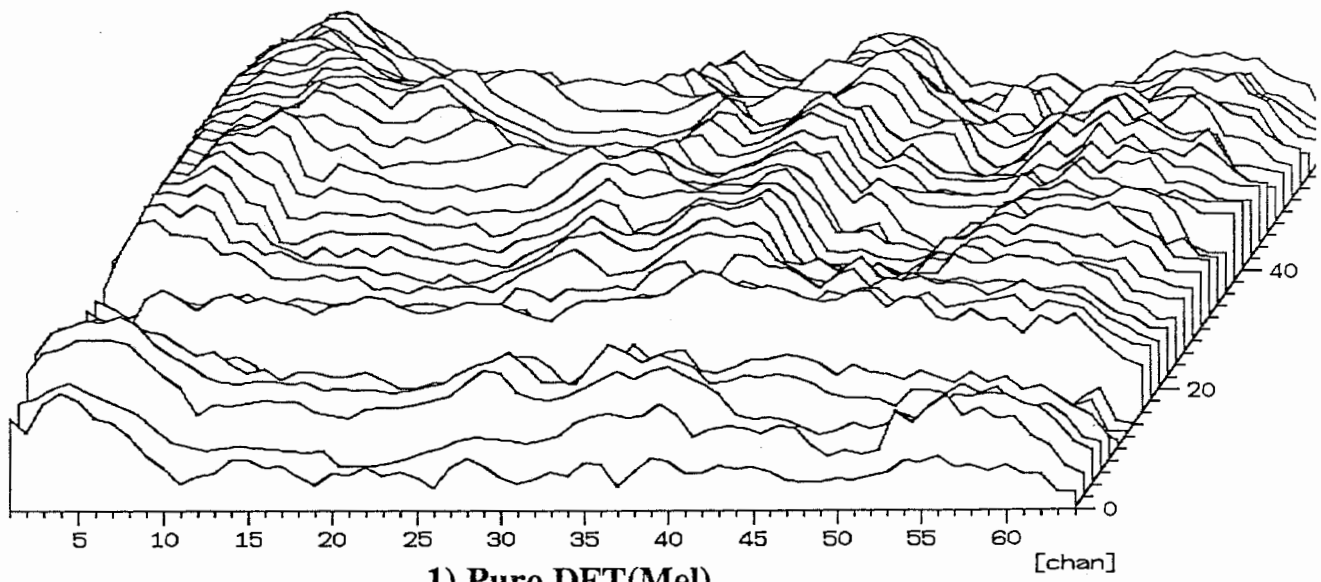
第4-5図 SN変動によるスペクトルの変化 (モデル1)



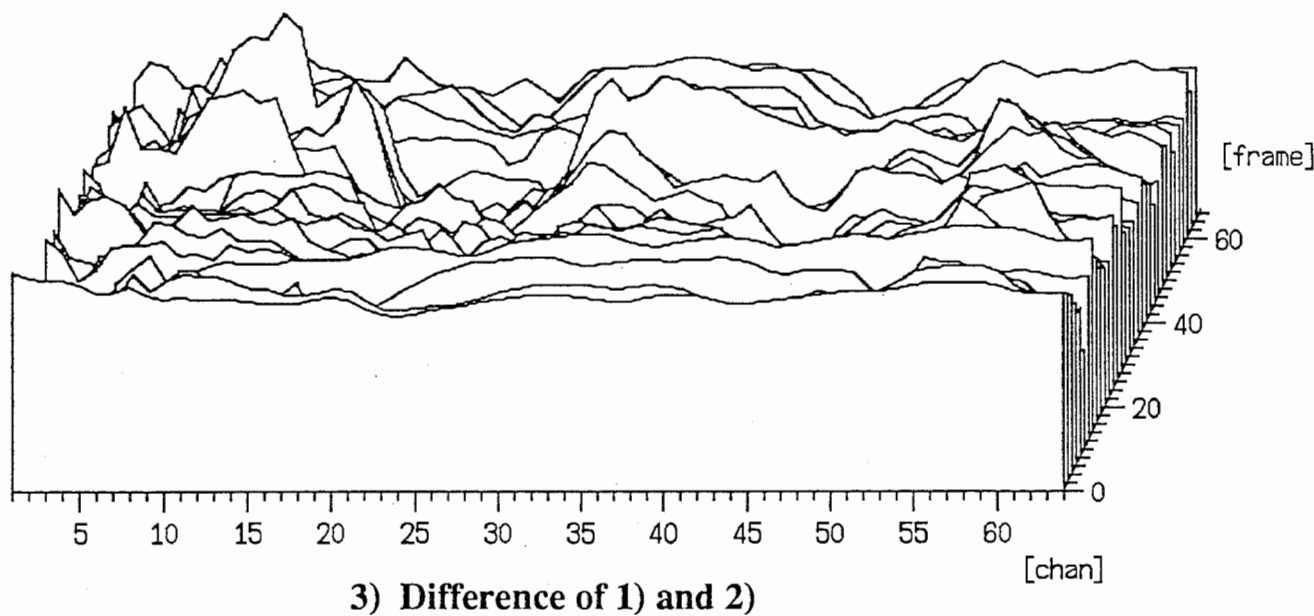
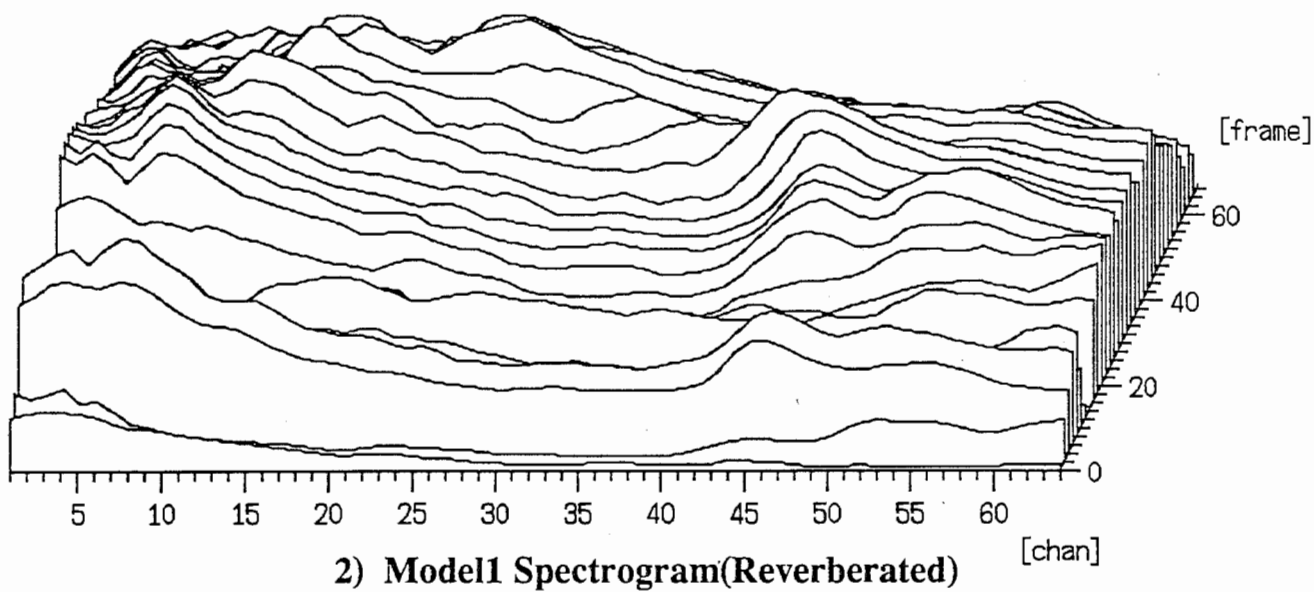
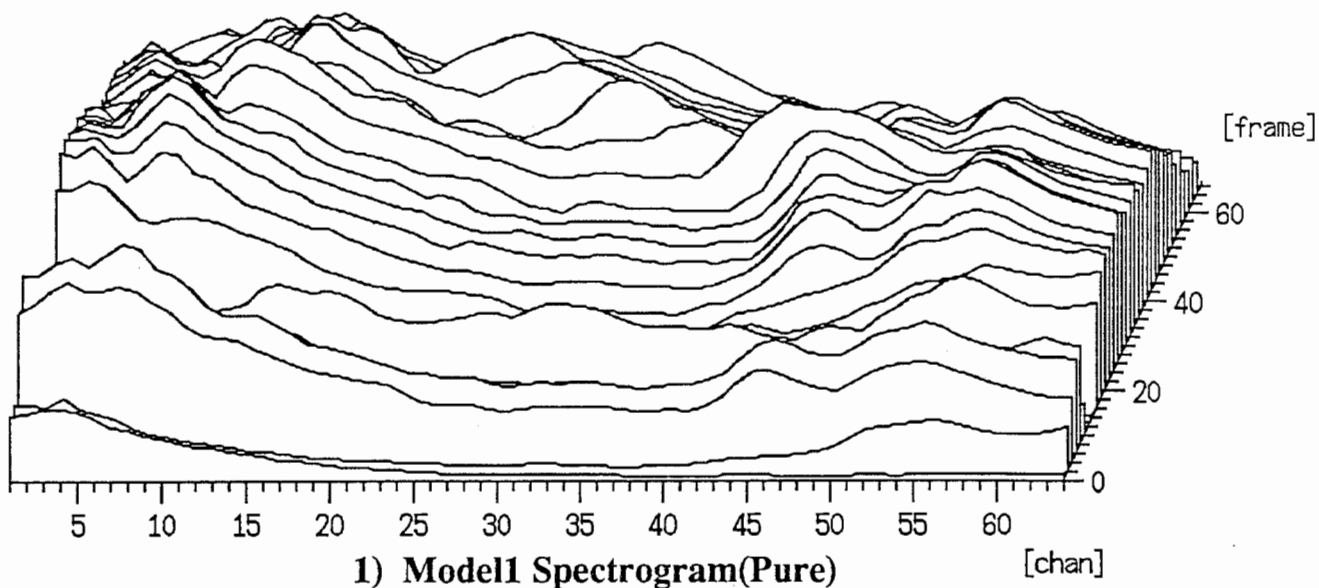
第4-6図 SN変動によるスペクトルの変化 (モデル2)



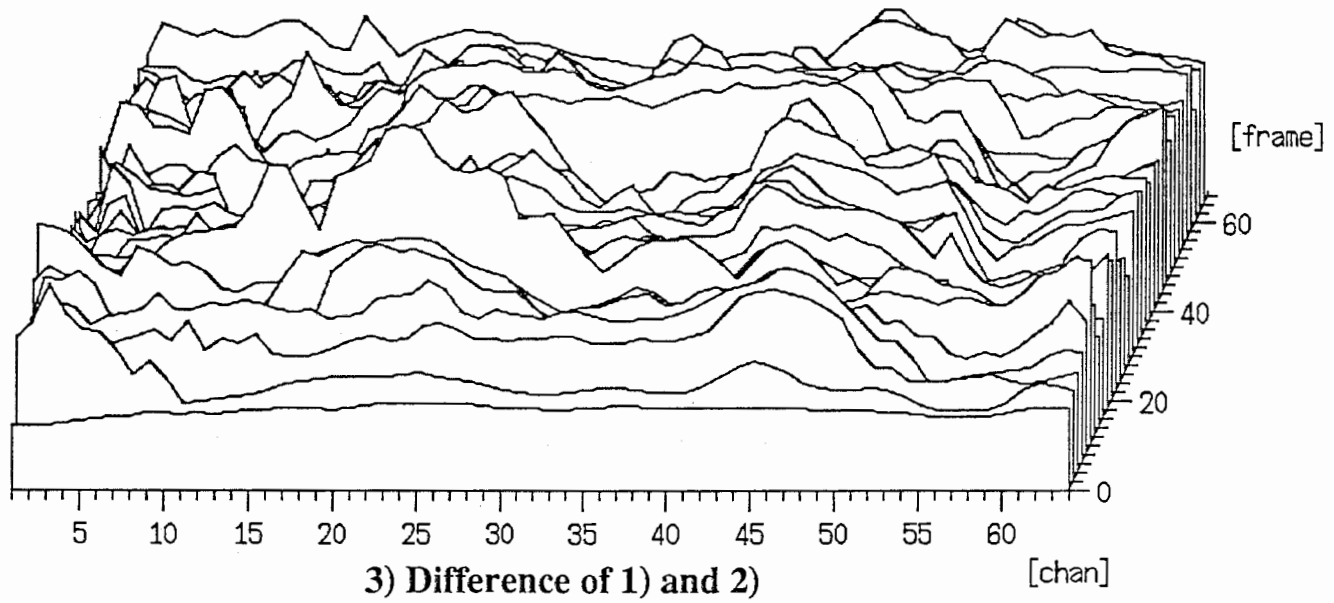
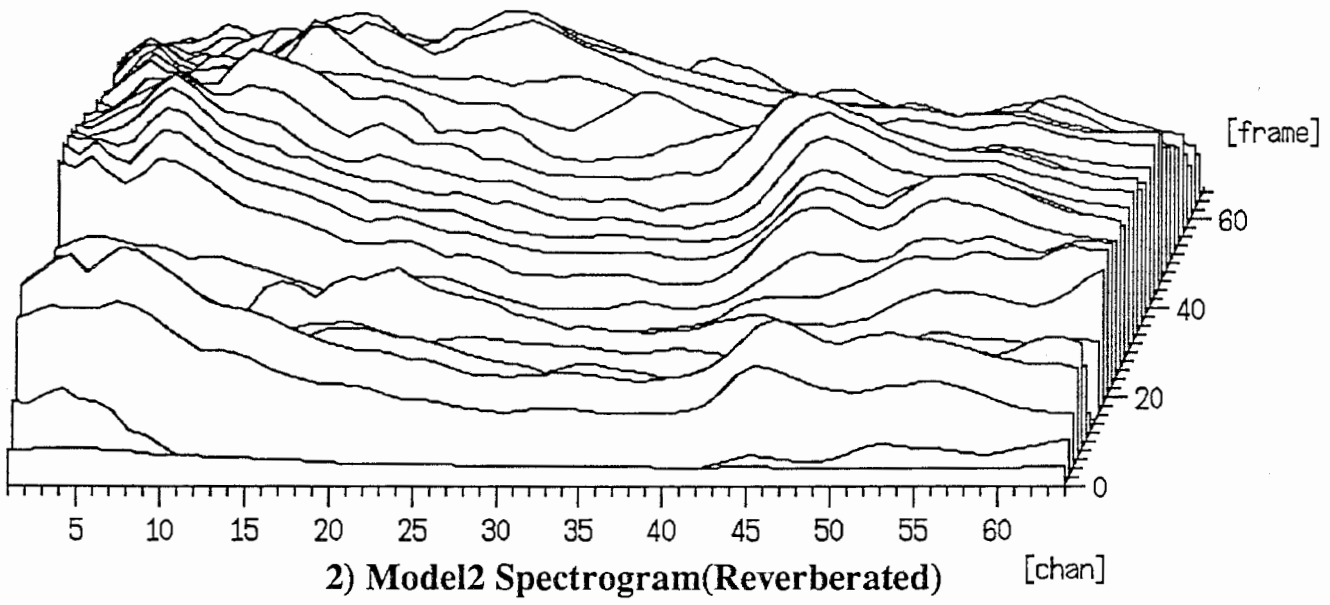
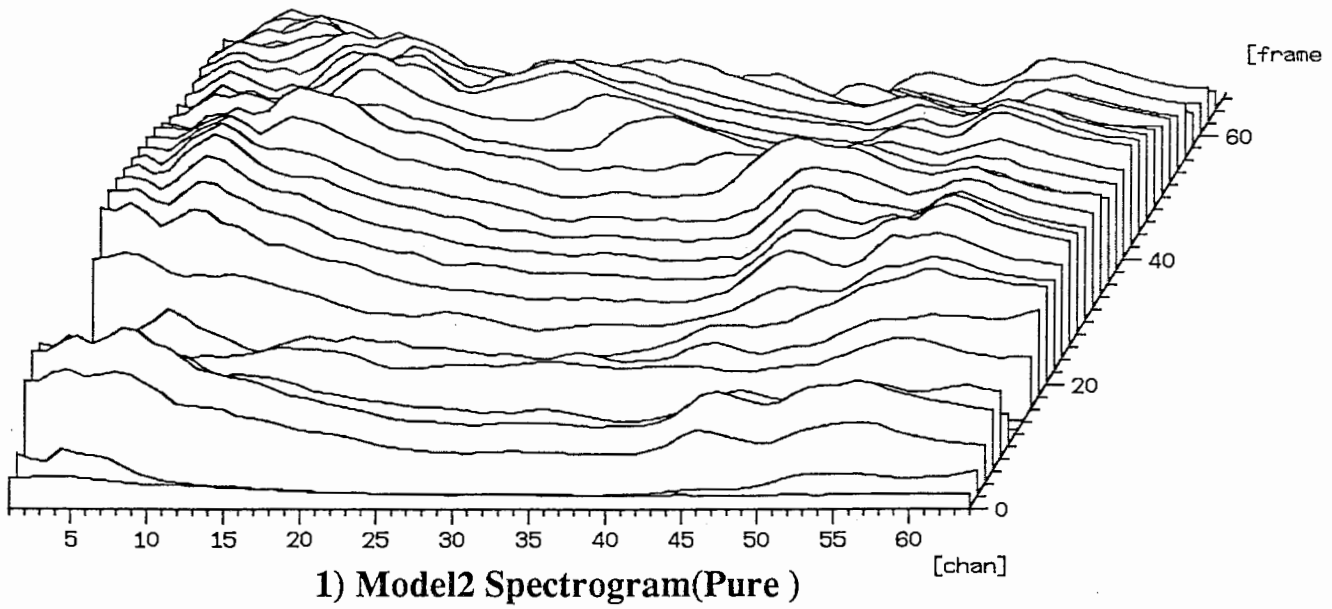
第4-7図 ノイズによるテンプレート単語間距離変動



第4-8図 残響付加によるスペクトル変化 (Mel DFT)



第4-9図 残響付加によるスペクトルの変動 (モデル1)



第4-10図 残響付加によるスペクトル変化 (モデル2)

Feature Extraction:DFT

Score:96.296295 Error 8

Word 75 recognized as 201 120 18 53 23
Word 76 recognized as 73 201 13 204 76
Word 134 recognized as 110 159 134 18 104
Word 154 recognized as 82 135 1 154 86
Word 155 recognized as 117 94 191 199 39
Word 195 recognized as 178 118 209 201 181
Word 204 recognized as 78 201 204 202 203
Word 212 recognized as 208 212 57 192 213

Feature Extraction:Model 1

Score 89.814812 Error 22

Word 7 recognized as 130 94 7 26 103
Word 36 recognized as 16 36 187 78 56
Word 50 recognized as 33 50 166 197 8
Word 66 recognized as 38 66 173 174 210
Word 73 recognized as 203 13 73 184 201
Word 75 recognized as 18 20 201 36 184
Word 76 recognized as 73 130 201 78 204
Word 90 recognized as 8 78 90 56 93
Word 91 recognized as 30 91 132 13 27
Word 134 recognized as 116 126 36 110 20
Word 147 recognized as 157 207 201 184 56
Word 154 recognized as 82 55 16 135 21
Word 155 recognized as 117 191 200 39 94
Word 186 recognized as 13 186 203 27 204
Word 187 recognized as 184 187 36 38 120
Word 191 recognized as 201 184 191 117 51
Word 195 recognized as 118 20 82 181 51
Word 204 recognized as 78 204 203 202 21
Word 205 recognized as 117 202 203 205 8
Word 208 recognized as 57 17 101 208 117
Word 211 recognized as 17 211 91 27 130
Word 212 recognized as 57 212 208 193 192

Feature Extraction:Model 2

Score: 90.27777 Error 21

Word 7 recognized as 130 12 94 205 106
Word 50 recognized as 33 50 98 197 75
Word 75 recognized as 201 16 120 53 18
Word 76 recognized as 201 73 184 13 78
Word 90 recognized as 21 90 8 104 53
Word 98 recognized as 16 98 36 116 58
Word 116 recognized as 16 116 36 58 187
Word 134 recognized as 116 16 110 36 38
Word 147 recognized as 207 201 120 157 56
Word 154 recognized as 135 82 1 55 21
Word 155 recognized as 117 94 39 200 191
Word 186 recognized as 13 202 27 30 120
Word 187 recognized as 184 187 38 201 120
Word 191 recognized as 201 191 5 53 184
Word 195 recognized as 118 20 21 201 178
Word 196 recognized as 173 184 182 38 187
Word 197 recognized as 210 197 200 180 121
Word 203 recognized as 202 203 73 201 13
Word 204 recognized as 202 201 120 21 78
Word 211 recognized as 17 58 100 91 93
Word 212 recognized as 57 208 17 212 135

付録A ATR音声データベース 音韻バランス単語一覧

単語番号	単語	単語
1	いきおい	勢い
2	いよいよ	いよいよ
3	うらやましい	羨ましい
4	おもしろい	面白い
5	ぐあい	具合
6	ざいりょー	材料
7	じゅーいちがつ	十一月
8	しゅーきょー	宗教
9	じゅんばん	順番
10	すいちよく	垂直
11	だいじょーぶ	大丈夫
12	だいどころ	台所
13	ちゃんと	ちゃんと
14	ちゅーおー	中央
15	とりあつかう	取り扱う
16	ばしよ	場所
17	びょーいん	病院
18	ひょーじゅん	標準
19	ぶんめー	文明
20	ぽけっと	ポケット
21	ぼんやり	ぼんやり
22	めがね	眼鏡
23	わがまま	我儘
24	うけあう	請け合う
25	うちあわせ	打ち合わせ
26	かれんだー	カレンダー
27	しゅーてん	終点
28	せーい	誠意
29	そびえる	そびえる
30	ちちおや	父親
31	ちゅーねん	中年
32	ておち	手落ち
33	でしゃばる	でしゃばる
34	にせもの	贋物
35	ははおや	母親
36	ひやくしよー	百姓
37	びょーにん	病人
38	ふらふら	ふらふら
39	ほほえむ	微笑む
40	れじゃー	レジャー
41	わすれもの	忘れ物
42	あかちゃん	赤ちゃん
43	あけがた	明け方
44	あじわう	味わう
45	あたりまえ	当たり前
46	あばーと	アパート
47	あるみにゆうむ	アルミニウム
48	あんけーと	アンケート
49	いえで	家出
50	いちじるしい	著しい

単語番号	単語	単語
51	いちぶぶん	一部分
52	いまごろ	今頃
53	いらっしゃる	いらっしゃる
54	うでまえ	腕前
55	うめぼれる	自惚れる
56	うめあわせる	埋め合わせる
57	うやまう	敬う
58	えーきゅー	永久
59	えーゆー	英雄
60	えすかれーたー	エスカレーター
61	えねるぎー	エネルギー
62	えぶろん	エプロン
63	えれべーたー	エレベーター
64	おいはらう	追い払う
65	おごそか	厳か
66	おしゃべり	お喋り
67	おぼえ	覚え
68	おもちゃ	おもちゃ
69	おもわず	思わず
70	がいしゅつ	外出
71	ぎこちない	ぎこちない
72	きむずかしい	気難しい
73	ぎやくたい	虐待
74	きゅーぎょー	休業
75	きゅーげき	急激
76	ぎゅーにゅー	牛乳
77	きゅーりょー	給料
78	ぎょーせー	行政
79	ぎよぎょー	漁業
80	きよくたん	極端
81	きよぜつ	拒絶
82	げひん	下品
83	げんえき	現役
84	こぎって	小切手
85	こころよい	快い
86	ことづて	言伝
87	このあいだ	この間
88	こびー	コピー
89	ごぶさた	御無沙汰
90	さきほど	先程
91	じぐざぐ	ジグザグ
92	しゅーへん	周辺
93	しゅっせ	出世
94	じゅんじゅんに	順々に
95	じゅんちょー	順調
96	じょーきやく	乗客
97	しよーふだ	正札
98	しよーりやく	省略
99	じんいん	人員
100	すびーど	スピード

単語番号	単語	単語
101	ぜんしゅー	全集
102	ぜんてー	前提
103	それぞれ	それぞれ
104	それでは	それでは
105	ぞんざい	ぞんざい
106	だいぶぶん	大部分
107	ちゅーしゃ	駐車
108	ちよーえつ	超越
109	つけくわえる	付け加える
110	できごと	出来事
111	てぬぐい	手拭
112	てのひら	手のひら
113	でばーと	デパート
114	とりあえず	取り敢えず
115	なおさら	なおさら
116	なかなかおり	仲直り
117	にゅーいん	入院
118	にゅーじょー	入場
119	にゅーよく	入浴
120	によーぼー	女房
121	ねあげ	値上げ
122	ねさげ	値下げ
123	のみこむ	飲み込む
124	ぱいぷ	パイプ
125	はなはだ	はなはだ
126	はなばなしい	華々しい
127	ばんぐみ	番組
128	ひきうける	引き受ける
129	ひっくりかえす	ひっくりかえす
130	びょーしゃ	描写
131	ひょーじょー	表情
132	びょーどー	平等
133	ひょーほん	標本
134	ぶろぐらむ	プログラム
135	ぶんるい	分類
136	ページ	ページ
137	べっど	ベッド
138	ぼんぷ	ポンプ
139	みうしなう	見失う
140	みすぼらしい	みすぼらしい
141	みせびらかす	見せびらかす
142	みやく	脈
143	みよーじ	名字
144	みよーにち	明日
145	めうえ	目上
146	もちぬし	持ち主
147	ものすごい	物凄い
148	ゆーもあ	ユーモア
149	ゆきずまる	行き詰まる
150	ゆびさす	指差す

単語番号	単語	単語
151	よこづな	横綱
152	よつかど	四つ角
153	よっぱらい	酔っ払い
154	りやくする	略する
155	りゅーいき	流域
156	りゅーちよー	流暢
157	りゅっくさっく	リュックサック
158	りよーがえ	両替
159	れくりえーしょん	レクリエーション
160	れんあい	恋愛
161	ろくおん	録音
162	ろっかー	ロッカー
163	ろんそー	論争
164	わざわざ	わざわざ
165	わりあてる	割り当てる
166	わりびき	割り引き

人為的な追加分

単語番号	単語	単語
167	せきらんうん	積乱雲
168	たんおんひょーじ	単音表示
169	かんふる	カンフル
170	まえひょーばん	前評判
171	おんなみよーり	女冥利
172	ふせーみやく	不整脈
173	かりゅー	下流
174	がびよー	画鋏
175	けびよー	仮病
176	ほびゅらー	ポピュラー
177	とつきよ	特許
178	そっちよく	率直
179	はんぎやく	反逆
180	まじよ	魔女
181	ざひょー	座標
182	こんにやく	こんにやく
183	かにゅー	加入
184	めにゅー	メニュー
185	こーみよー	巧妙
186	みゅーじっく	ミュージック
187	きにゅー	記入
188	こーにゅー	購入
189	さんみやく	山脈
190	はっぴやく	八百
191	けわしい	険しい
192	めーりよー	明瞭
193	かんびよー	干瓢
194	ついぎゅー	追及
195	おう	遑う
196	いんりよく	引力
197	ごびゅー	誤謬
198	てちよー	手帳
199	さぎよー	作業
200	だいひょー	代表

単語番号	単語	単語
201	とーひよー	投票
202	えーぎよー	営業
203	びみよー	微妙
204	じゅみよー	寿命
205	じぎよー	事業
206	じゅーびよー	重病
207	かんびよー	看病
208	でんわ	電話
209	けんきゅー	研究
210	けいびやく	啓白
211	だきよー	妥協
212	ひゅーひゅー	ひゅーひゅー
213	びゅーびゅー	びゅーびゅー
214	おやびゅーま	親ビューマ
215	たいびゅーた	タイビュータ
216	かくれびゅーりたん	隠れビューリタン