

ヒューマン・インタフェースにおける音声対話速度の適応化 (分担研究：相互作用と乳幼児の心理行動発達に 関する基礎的研究)

渡辺富夫*

要約 個人の発声速度に対応して機械がその個人に最適な音声出力速度で応答できれば、円滑な音声対話が図られると期待されよう。本報告では、個人の発声速度を時間率[平均ON区間/(平均ON区間+平均OFF区間)]に基づいて推定し、その発声速度に音声応答速度を適応化させる手法を提案している。

見出し語：ヒューマン・インタフェース，発声速度，適応化

1. 緒言

人間機械間の音声対話においては、人間は機械へ音声入力する話し手であると同時に機械からの音声応答の聞き手である。人間にとって理解しやすい音声応答速度は、個人の発声速度(無声区間をも含めたスピーチ区間でのモーラ/秒)¹⁾と密接な関係があると考えられる。したがって、個人の発声速度に対応して機械がその個人に最適な音声出力速度で応答できれば、個人と機械との円滑なコミュニケーションが図られ、人間に適合したヒューマン・インタフェースの実現に役立つと期待されよう。

著者は、前報¹⁾において音声分析システムを開発し、官能検査により、人間一般について話し手が理解しやすいと考える発声速度と聞き手として理解しやすい音声応答速度との相関関係を明らかにし、機械からの音声応答速度を人間の発声速度に適応化させることの有効性を確認した。²⁾

本報告では、留守番電話を想定した場合について、人間の発声速度を時間率[平均ON区間/

(平均ON区間+平均OFF区間)]に基づいて推定し、機械からの音声応答のOFF区間を伸縮してその音声入力された時間率に一致させることにより、音声応答速度を適応化させる手法を提案している。

2. スピーチのON区間とOFF区間の特性分析

無声子音の前のOFF区間(無声区間)は、発声器官のメカニズム上生じる区間で、実際にはON区間(有声区間)と考えられる。ここでは、図1に示すようにON区間にハングオーバー(ある時間だけON区間を伸ばすこと)とフィルイン(ある値より小さな継続時間のOFF区間のみをON区間に置換すること)を施して、上記OFF区間をON区間に置換し、呼気段落区分でのON区間とOFF区間(ポーズ区間)を計測する。しかしこの呼気段落区分に適したハングオーバーとフィルインの時間は、個人によって異なり、しかもポーズ区間が極端に短い場合には明確に決定できないので、任意のハングオーバーとフィルインを施した場合について、人間

* 山形大学工学部

(Faculty of Engineering, Yamagata Univ.)

注)

1 モーラは連続音声を仮名書きした場合の1文字に相当する。

一般に関する発声速度と時間率との相関関係を検討する。ここで時間率 (speech activity) は、代表的なスピーチの ON-OFF パラメータで、スピーチ区間で ON 区間の占める割合をさす。しかしスピーチ区間の短い対話においては、スピーチ区間で OFF 区間の回数が ON 区間の回数よりも 1 つ少ないことは時間率の算出において無視しえないので、本論文では時間率を平均 ON 区間 / (平均 ON 区間 + 平均 OFF 区間) と定義する。

2・1 実験 被験者は 19 才～23 才の男性 46 人である。被験者は留守番電話を想定したメッセージを、表 1 に示すキーワードを参考にして自分なりの言葉で伝達するという条件で発声した。被験者 46 人のモーラ数の平均値 ± 標準偏差は 128.9 ± 17.0 モーラであった。

2・2 測定 分析システム構成図を図 2 に示す。被験者の音声は、ビデオテープに PCM 録音した後、2.5 KHz のローパスフィルタを通

し、サンプリング周波数 5 KHz で A/D 変換され、本システムに収録される。その音声信号の ON 区間と OFF 区間の機械識別は、50 msec の平均雑音レベルに 12 dB 加えた値を臨界値として、ON-OFF の判定周期 10 msec (フレームと呼ぶ) についてこの値を越えた部分を ON、越えない部分を OFF とした。ハングオーバーとフィルインを施した場合の ON-OFF パターンの測定例を図 1 に示す。

2・3 スピーチの ON-OFF パラメータ

任意のフィルインとハングオーバーを施した場合の平均 OFF 区間、平均 ON 区間、時間率は OFF 区間の確率密度関数に基づいて以下のように計算される。³⁾

2・3・1 任意のフィルインを施した場合

フィルインを施さない、即ちフィルイン時間が 0 でのフレーム数 K における OFF 区間の確率密度関数を $f(k)$ とすると、フィルイン時間 m

表 1 キーワード表

“もしもし”	
“山形大学の_____です。”	
用件	講演の依頼
日時	11月3日(文化の日)
テーマ	「情報の価値について」
“今週末までに_____まで連絡して下さい。”	
“よろしくお願ひ致します。”	

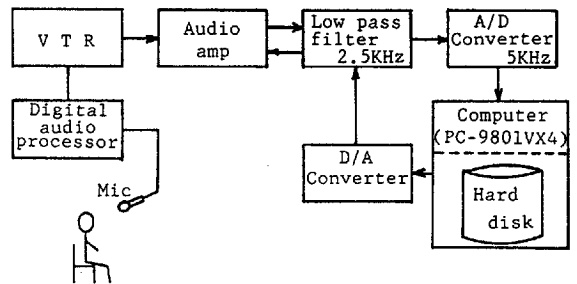


図 2 分析システム構成図

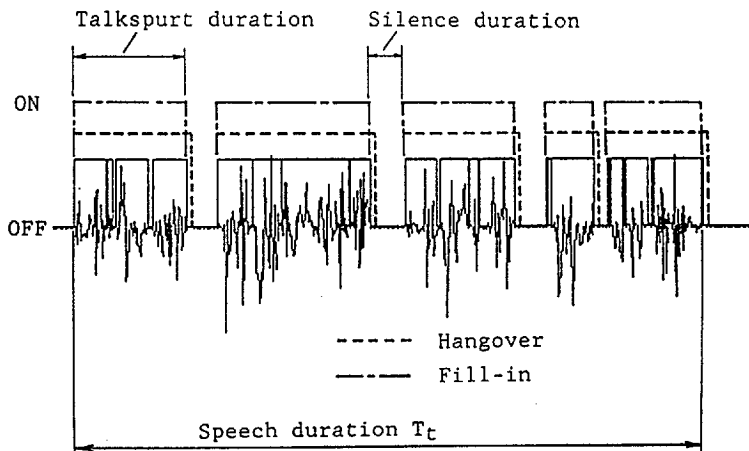


図 1 フィルインとハングオーバーを施した場合の ON-OFF パターンの測定例

の確率密度関数 $f_m(K)$ は、

$$f_m(K) = f(K) / \sum_j f(j) \quad (1)$$

$$j, K = m+1, m+2, m+3, \dots$$

$$m = 0, 1, 2, \dots$$

である。したがってフィルイン時間 m の平均 OFF 区間 $S_f(m)$ は、

$$S_f(m) = \sum_K (K \cdot f_m(K)) \quad (2)$$

$$m = 0, 1, 2, \dots$$

で算出される。一方、平均 ON 区間は以下のように求められる。 $m=0$ でスピーチ区間 T_t の ON 区間の回数を N_0 とすると、OFF 区間の回数が ON 区間より 1 つ少ないことを考慮して、フィルイン時間 m での ON 区間の回数 N_m は、

$$N_m = 1 + (N_0 - 1) \cdot \sum_K f(K) \quad (3)$$

$$K = m+1, m+2, \dots$$

である。

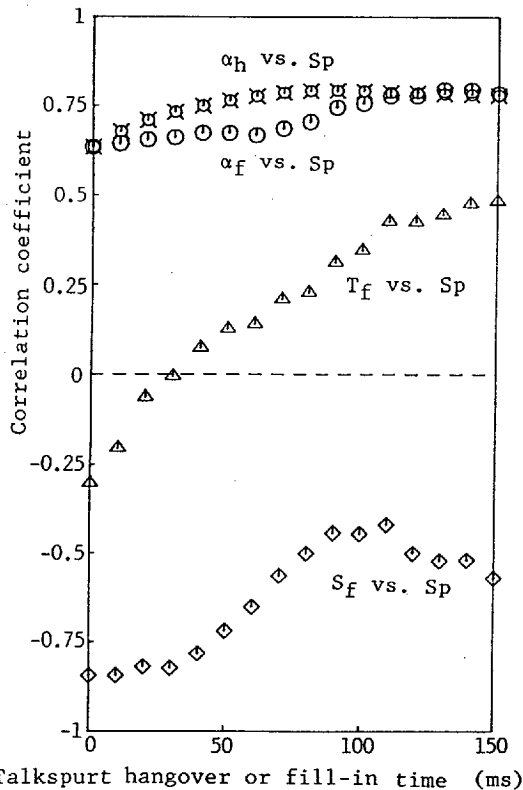


図3 フィルイン時間とハングオーバー時間の関数としてみた発声速度に対する時間率、平均 ON 区間及び平均 OFF 区間との相関係数

したがって平均 ON 区間 $T_f(m)$ は、

$$T_f(m) = T_t / N_m - (N_m - 1) \cdot S_f(m) / N_m \quad (4)$$

$$m = 0, 1, 2, \dots$$

で算出される。これは文献3)の近似計算式を修正したものである。時間率 $\alpha_f(m)$ は次式で定義される。

$$\alpha_f(m) = T_f(m) / (T_f(m) + S_f(m)) \quad (5)$$

発声速度 S_p は T_t とモーラ数 M_r より次式で算出される。

$$S_p = M_r / T_t \quad (6)$$

2.3.2 任意のハングオーバーを施した場合

ハングオーバーの性質上、フィルインと比べて平均 OFF 区間 $S_h(m)$ は m だけ短く、平均 ON 区間 $T_h(m)$ は m だけ長くなる。

$$S_h(m) = S_f(m) - m$$

$$T_h(m) = T_f(m) + m \quad (7)$$

時間率 $\alpha_h(m)$ は次式で与えられる。

$$\alpha_h(m) = T_h(m) / (T_h(m) + S_h(m))$$

$$= (T_f(m) + m) / (T_f(m) + S_f(m)) \quad (8)$$

2.4 分析 発声速度を、早口、普通、遅口の3種に変えた場合、ON 区間の伸縮に比べて OFF 区間の伸縮の比率が大きく、早口ほど時間率が高くなる。人間一般についても発声速度を推定するのに、時間率が有効な指標になると考えられる。

図3に被験者46人について、フィルイン時間とハングオーバー時間の関数としてみた発声速度に対する時間率、平均 ON 区間及び平均 OFF 区間との相関係数を示す。発声速度と時間率の相関が比較的高いのは、フィルインが90msec以上、ハングオーバーが80msec以上であった。最大の相関係数を示すのは、フィルイン時間が130msecで、ハングオーバー時間が90msecであり、その相関係数の値に差異はない。ハングオーバーを施した場合、その時間だけ呼気段落区分での OFF 区間が短くなるので、以後はフィルインについて検討する。またフィルイン時間0msecでの発声速度と平均 OFF 区間との相関係数は-0.84であり、フィルイン時間130msecでの発声速度と時間率との相関と同程度に高い。しかし、音声対話速度の適応化においては、無声子音の前の短い OFF 区間の伸縮は音声合成

上不適切であり、この呼気段落区分でのOFF区間が伸縮の対象になる。フィルインを施した場合の呼気段落区分での発声速度と平均OFF区間との相関は、時間率との相関よりも低く、したがって発声速度と最も相関が高いのは時間率である。

図4にフィルイン130msecでの発声速度と時間率との関係を示す。時間率 α より、発声速度 S_p は次式の回帰直線で推定される。

$$S_p = 10.5\alpha - 0.279 \quad (9)$$

フィルイン時間130msecでのモーラ当りのON区間とOFF区間との関係を図5に示す。両者の相関は低く、早口の人が必ずしもOFF区間が短いわけではない。したがって、最適な音声応答速度の制御は、厳密にはON区間と、OFF区間の両方について行なう必要がある。しかし、図5に示すようにモーラあたりのON区間との標準偏差/平均値の比は2:9で、ON区間の変動はOFF区間に比べて小さく、OFF区間のみの制御によってかなりの適応化が図られると考えられる。事実、フィルイン時間130msecでのモーラ当りのON区間と時間率との相関係数が-0.11とほとんど0に近いのに対し(図6)、モーラ当りのOFF区間と時間率との

相関係数は-0.95と極めて負の相関が強い(図7)。

註

平均ON区間及び平均OFF区間の発声速度に対する相関係数は、各々ハングオーバー時間とフィルイン時間の関数として等しい。

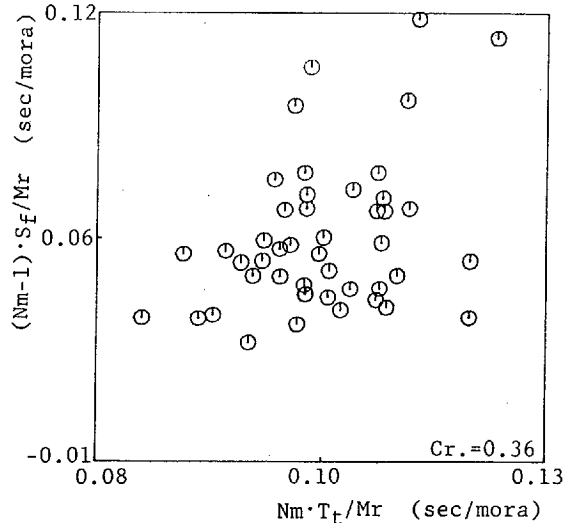


図5 モーラ当りのON区間とOFF区間との関係(フィルイン130msec)

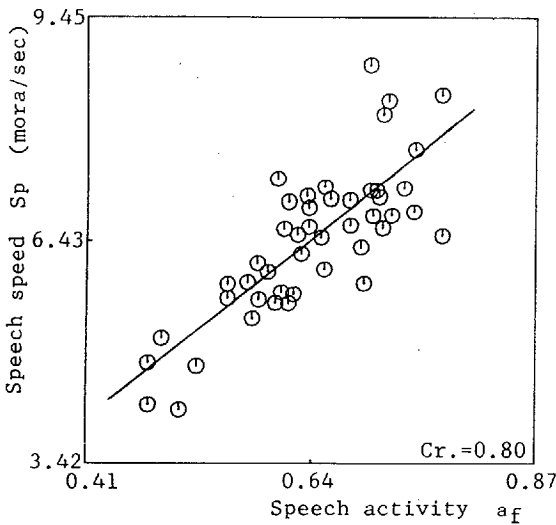


図4 発声速度と時間率との関係
(フィルイン130msec)

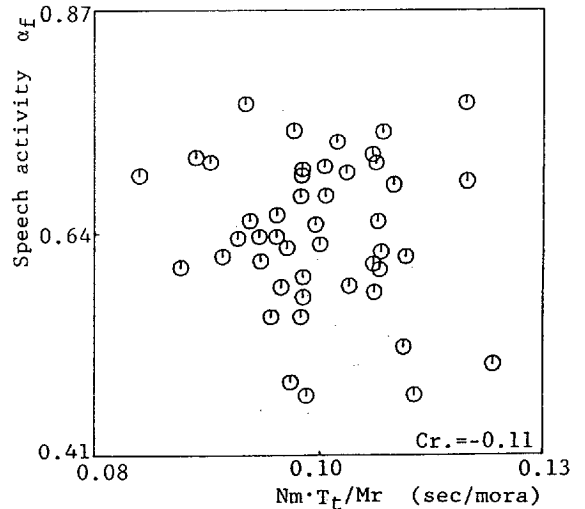


図6 時間率とモーラ当りのON区間との関係
(130msec)

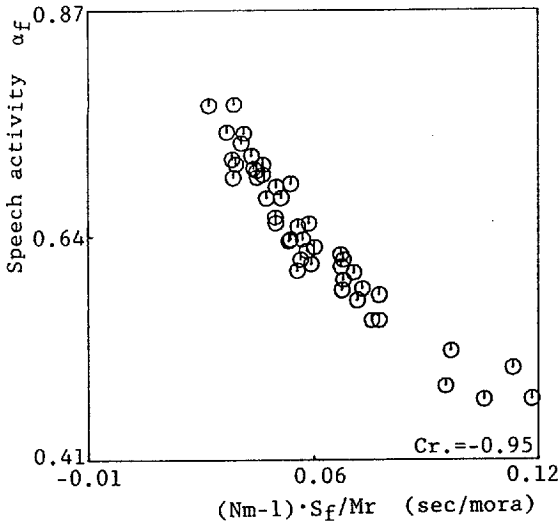


図7 時間率とモーラ当りのOFF区間との関係 (フィルイン130msec)

3. 音声対話速度の適応化

2章での分析結果より、音声応答速度の適応化は主として、OFF区間について行なえばよいことが判明した。また呼気段落区分での発声速度との相関は時間率が最も高く、音声応答速度を入力された発声速度に適応化させるには、音声応答のOFF区間を単にその入力音声の平均OFF区間に置換するのではなく、時間率に一致するように平均OFF区間を伸縮する必要がある。時間率を合わせることは、そのON-OFFパターンを合わせることであり、ON区間とOFF区間の繰り返しのリズムに適応化させることである。

いま、留守番電話のように予め録音されている音声を入音の発声速度に適応化させる場合を考える。機械応答音声と入力音声の時間率が各々 $\alpha_1 = T_1 / (T_1 + S_1)$ と $\alpha_2 = T_2 / (T_2 + S_2)$ のとき、両者の時間率が一致するように、

$$T_1 / (T_1 + S_1) = T_2 / (T_2 + S_2) \quad (10)$$

$$\therefore S' = T_1 \cdot S_2 / T_2$$

S_1 を $T_1 \cdot S_2 / T_2$ に伸縮することによって、即ち、フィルインを施した(呼気段落区分での)各OFF区間を $(T_1 \cdot S_2) / (T_2 \cdot S_1)$ の比率

で伸縮すれば、適応化が図られる。このとき、機械応答速度 $Sp_1 = Mr / (N_0 \cdot T_1 + (N_0 - 1) \cdot S_1)$ は、

$$Sp_1' = Mr / (N_0 \cdot T_1 + (N_0 - 1) \cdot T_1 \cdot S_2 / T_2) \quad (11)$$

Mr : スピーチ区間内のモーラ数

になる。

全被験者46人中、時間率が最大値、平均値、最小値の被験者a, b, c各々について、本手法を用いて本人以外の被験者45人の発声速度に適応化させた結果を図8に示す。適応化させた発声速度 Sp_1' と時間率 α_2 には、

$$Sp_1' = [Mr / (N_0 \cdot T_1)] \cdot [T_2 / (T_2 + (N_0 - 1) \cdot S_2 / N_0)] \\ = [Mr / (N_0 \cdot T_1)] \alpha_2 \quad (12)$$

の関係があり、図中の+印の傾きは $Mr / (N_0 \cdot T_1)$ 、即ちa, b, c各々のOFF区間を含まないON区間のみに基づく発声速度である。いま発声速度の適応化率 λ を次式で定義する。

$$\lambda = 1 - \sqrt{\sum \epsilon_0^2 / \sum \epsilon_n^2} \quad (13)$$

ϵ_0 : 適応化後の発声速度と目標の発声速度との差

ϵ_n : 適応化前の発声速度と目標の発声速度との差

λ は、a, b, cについて各々61.6%, 40.0%, 59.4%であった。bの λ が、a, cの λ よりも低

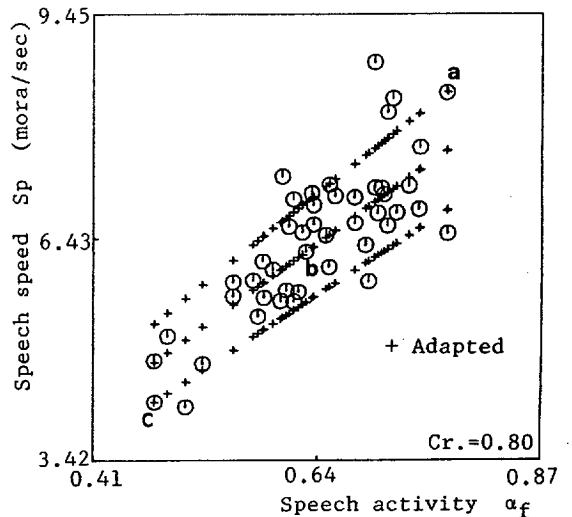


図8 全被験者46人中、時間率が最大値、平均値、最小値の被験者a, b, c各々について、本人以外の被験者45人の発声速度に適応化させた結果

いが、これは Σe_n^2 がもともと小さいためで、 $\sqrt{\Sigma e_n^2}/45$ は最も小さく 0.61 モーラ/秒であった。

4. 結言

本報告では、人間一般についての発声速度と時間率との関係を明らかにするとともに、個人の発声速度を時間率に基づいて推定し、その発声速度の時間率に一致するように機械からの音声応答のOFF区間を伸縮して音声応答速度を適応化させる手法を提案した。本手法で、音声応答の平均OFF区間を単に発声速度の平均OFF区間に一致させるのではなく、時間率に一致するように平均OFF区間を伸縮させることの根拠は、1)時間率とOFF区間とは強い相関関係があり、音声応答速度の適応化は主としてOFF区間について行なえばよいことが判明したこと、2)呼気段落区分での発声速度との相関は、時間率の方が平均OFF区間よりも高かったことによる。時間率を合わせることは、そ

のON-OFFパターンに、即ち話し手のリズムに適応化させることである。本手法を人間と機械との音声対話に適用することにより、個人に適合した円滑な対話が図られると考えられる。

文献

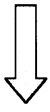
- 1) 渡辺富夫ら：乳幼児の泣き声収録分析システムの開発、厚生省心身障害研究「家庭保険と小児の成長・発達に関する総合的研究」昭和61年度研究報告書，51-53，1987.
- 2) 渡辺富夫：マン・マシン・インタフェースにおけるスピーチの理解しやすい速度に関する研究，日本機械学会論文集，53(496)，851-856，1987.
- 3) Gruber, J.G.: A Comparison of Measured and Calculated Speech Temporal Parameters Relevant to Speech Activity Detection, IEEE Trans. on Communications, 3(4), 728-738, 1982.

Abstract

Adaptation Method of Machine Conversational Speed to Speaker's Speed for Man-Machine Communication

Tomio WATANABE

This paper first discusses the relationship between the speech speed and the speech activity [mean talkspurt duration / (mean talkspurt duration + mean silence duration)] in 46 subjects for any value of "hangover" and "fill-in". Secondly, based on this relationship, an adaptation method of machine conversational speed to speaker's speed is proposed. The method controls the mean silence duration in the machine's speech activity in the same way as the speaker's speech activity. The simulated adaptation rate to speaker's speed in this method is approximately 60%. This method is effective in adaptation for the great differences between machine's speech activity and a speaker's speech activity, and is useful in realizing a man-machine interface with smooth information exchange.



検索用テキスト OCR(光学的文字認識)ソフト使用

論文の一部ですが、認識率の関係で誤字が含まれる場合があります



要約 個人の発声速度に対応して機械がその個人に最適な音声出力速度で応答できれば、円滑な音声対話が図られると期待されよう。本報告では、個人の発声速度を時間率[平均 ON 区間/(平均 ON 区間+平均 OFF 区間)]に基づいて推定し、その発声速度に音声応答速度を適応化させる手法を提案している。