

合成処理時間が設定可能な文音声合成ソフト

6B-3

原 義幸

新田 恒雄

小林 賢一郎†

(株)東芝 情報処理・機器技術研究所

†東芝AVE (株)

1. はじめに

近年、電子メール文の読み上げのように、漢字かな混じり文を音声に変換する「文音声合成 (Text-to-Speech; 以下TTSと略す)」技術利用の要求が高まってきている。このような背景のもと著者らは、先にケプストラム方式を用いた文音声合成ボードの試作について報告した⁽¹⁾⁽²⁾。

一方、現今のワークステーション (WS) は、処理能力が向上し (数十~数百MIPS)、同時に、オーディオデバイス (CODEC, スピーカ) を標準で搭載する機種が増えつつある。このようなWSを用いると、専用ハードウェアなしにソフトウェアのみでTTSを実行できる。しかし、サーバ/クライアント、あるいはマルチタスク処理環境のもとでは、TTSの実時間処理が困難となる場合を生ずる。こうした問題に対処するため、処理時間の設定が可能なTTSソフトをWS (AS4075) 上に構築したので、概要を述べる。

2. 処理の概要

本ソフトは図1に示すように大きく分けて言語処理部と音声合成部から成る。

2.1 言語処理部

入力された漢字かな混じり文は、最初、単語辞書を照合しながら総当たり法によって単語単位に分割され、読み・品詞・アクセント型などの情報に置き換えられる (形態素解析)。次に、読みの変形処理が行われ、続いてアクセント句結合規則により単語をアクセント句に結合してアクセント型を求めるとともに、ポーズ挿入規則に従ってポーズ記号が挿入される (音声記号列生成)。得られた音声記号列 (韻律情報と音韻系列) は、音声合成部に送られる。

2.2 音声合成部

ここでは、最初、音韻系列に対応する音声素片 (ケプストラムパラメータ) を取り出すとともに、音節時間配置パラメータをもとに音韻長を決定する。次に、補間結合規則に基づいて各素片間

を補間することにより、音韻パラメータ列を生成する。一方、韻律パラメータ生成では、アクセント型やポーズ長などの韻律情報と音韻長に基づきピッチパターンが与えられる。

合成器における音源生成は、有声としてインパルス系列を、無声としてM-系列を用いた。また、合成フィルタは、ケプストラムパラメータを直接その係数とするLMA (Log Magnitude Approximation) フィルタ⁽³⁾を用いている。

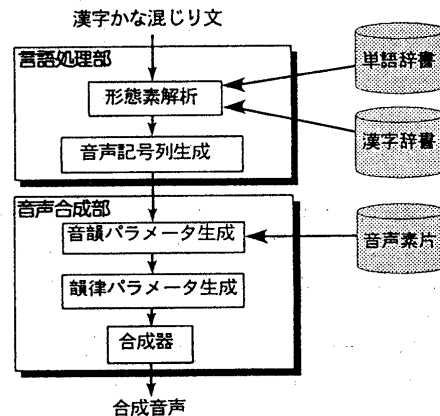


図1 TTS処理の流れ

3. 実時間処理のための手法

一般にWSはマルチタスク処理が可能なOS下で動作している。したがって、TTS以外のタスクが増えると、音声生成が実時間で行えなくなり、途中で途切れや異音を生じることになる。同様の問題は、TTSサーバとしての利用の際にも起こる。こうした問題に対処するため、音声生成に要する時間を可変にする方法と、実時間処理できない場合でも途切れによる聞きづらさを低減する方法の2つを採用した。

3.1 合成次数可変フィルタ

文を音声に変換する処理の中で、最も時間を要するのは合成フィルタ (LMAフィルタ) である。このフィルタは、サンプリング周期内に

TTS (Text-to-Speech Synthesis) Software with Controllable Processing Time.

Yoshiyuki Hara, Tsuneo Nitta, and Ken-ichiro Kobayashi †

Information Systems Engineering Lab., Toshiba Corp.

† Toshiba AVE Co., Ltd.

数十回の積和演算処理を行っている。したがって、この演算回数を可変にすることによって、音声生成に要する時間を増減することができる。これを実現するには、LMAフィルタにおける基礎フィルタの修正 $pad\ e$ 近似式の次数、つまりフィルタの内部構造を可変にする方法と、合成するときの音韻パラメータの次数を可変にする方法とがある⁽³⁾。今回は、簡単に実現できる後者の方法を採用した。

図2に、このフィルタのブロック図を示す。

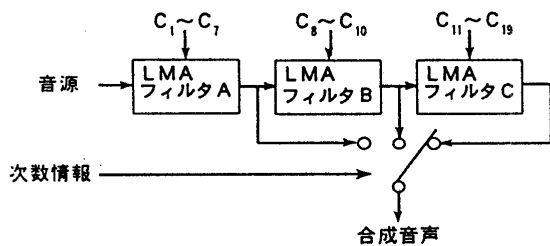


図2 合成次数可変フィルタ

図のように次数情報の入力により、合成次数を変えることが可能である。したがって、演算量を変えることができ、音声生成に要する時間を増減できる。

図3は3通りの合成次数Mで生成した音声のスペクトル包絡(「toshiba」)を示している。

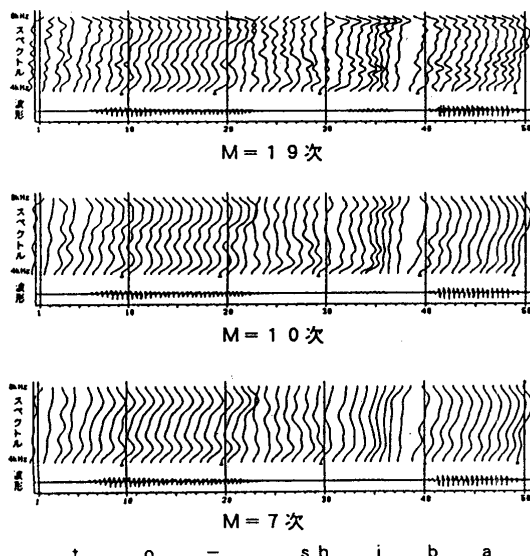


図3 「toshiba」のスペクトル包絡

19次、10次、7次と、スペクトルの山谷がなだらかになっているのが分かる。音質は、19次が歯切れの良いのに対して、次第になまった感じの音になるが、7次までは実用に耐えると考えている。

3.2 音声の途切れへの対応

通常、TTSを専用ハードで処理する際は1サンプルの音声データが生成される毎にDAコンバータに転送する方法が採られる。しかし、この方法では実時間処理が不可能な場合、音声の途切れも1サンプル毎になるため、非常に聞きづらくなる。そこで、無音区間を表す記号で挟まれたフレーズの音声データを生成した後、CODECに転送する方式を採用した。この方式では実時間に処理できない場合にも、無音の箇所に新たに無音が加わるだけのため、音声途中で途切れることによる聞きづらさが解消される。

4. 処理時間の例

0.5秒間の音声を生成するために必要な時間をAS4075上で測定したので、その結果を表1に示す。

表1 処理時間

合成次数	処理時間 [s]	
	A	B
19	0.58	0.7
10	0.39	0.48
7	0.31	0.38

A: 合成フィルタのみ
B: 全処理

この例では、次数が19次の場合、実時間で処理されていないが、本来はいるべき無音区間が少し長くなるだけのため、文理解に支障はない。

5. まとめ

合成処理時間が設定可能な文音声合成ソフトについて報告した。本ソフトは、電子メール、文章校正支援、ドキュメント読み上げ、音声による通知などに幅広く応用できると思われる。今後は、音質向上とともに応用ソフトの開発を行いたい。

【参考文献】

- (1) 原ほか：“文音声合成ボードの試作”
信学全大, A-231 (1992.3)
- (2) 小林ほか：“日本語解析処理を用いた音声読み上げ方式の開発”
情報全大, 7N-6 (1992.3)
- (3) 今井：“対数振幅近似(LMA)フィルタ”
信学論(A), J62-A, 12(1980.12)