

東芝における最近の音声合成・認識の応用

松浦博 正井康之 田中信一 原義幸 桃崎浩平 新田恒雄 *太田治徳 **小林賢一郎
(株)東芝 マルチメディア技術研究所 *青梅工場 **東芝AVE(株)

1. 東芝音声システム⁽¹⁾

東芝音声システムは日本語 Windows^(注1) 95で動作し、TTSを用いたユーザインタフェースを提供するソフトウェア環境である。本システムの主な特長を以下にあげる。

- ・直交化残差方式とケプストラム合成方式によって音質を向上させている。
- ・約12万語の単語辞書、2千語登録可能なユーザ辞書により正確なアクセントを付ける。
- ・知識処理を用いた解析によって、「質の良い宝石を質に入れる」の「質(しつ)」と「質(しち)」のような同字異音語を読み分ける。
- ・豊富なアプリケーションを備えている。

本システムは図1に示すように主にユーザインタフェース部とTTSエンジンからなる。

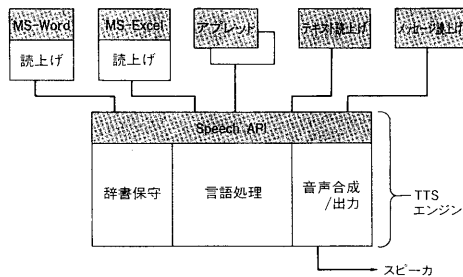


図1 東芝音声システムの構成

ユーザインタフェース部はユーザがTTSを利用するためのツール群である。現状ではTTSを標準にサポートしている市販のアプリケーションが少ないため、独自に作成した。TTSエンジンはTTS処理を行うモジュールであり、それと共に上位ソフトウェアに対して、TTS機能を使用するためのAPI(Application Programming Interface)を提供する。本APIはMicrosoft^(注2)社が提案するSpeech APIのサブセットになっている。

1.1 TTSエンジン

TTSエンジンは、実際のTTS処理を行う文音声合成エンジン部と、APIを提供するSpeech API部の2層からなる。文音声合成エンジン部は、言語処理、音声合成/出力および辞書保守の各処理を実行するモジュールで、32ビットDLL(Dynamic Link Library)になっている。

Speech API部は、文音声合成エンジン部の機能をユーザに提供するためのインタフェースモジュールである。ユーザはこのAPIを介してだけ文音声合成エンジンの機能を使用することができる。テキストあるいは音声記号列から音声への変換、音声出力の停止/再開、音声属性の各種設定、ユーザ辞書登録などの機能を提供する。APIはOLE2(Object Linking and Embedding)に基づいて定義されており、API部の上層に実装されたラッパーによって、TTSエンジンをオブジェクトとして扱うことができる。

1.2 ユーザインタフェース部

TTSを利用するための支援ツール群として以下のものを提供する。

(1) テキスト読上げ

アプリケーションプログラムであり、Speech APIを使用して、テキストファイルの読上げ、テキストファイルの作成、クリップボードの読上げ、OLEサーバ機能の動作を行う。OLEサーバ機能とは読上げ文書の埋込みおよびリンクのサポートである。さらに、読み上げ音声の種類の設定、読み上げモードの各種設定、ユーザ辞書登録/削除/更新/検索を行う。

(2) アプレット

入力したキーおよび計算結果を読み上げる音声電卓アプレット、時刻を読み上げる音声時

刻アプレット、キーボードから入力されたコードを読み上げるキーボード読上げアプレットがある。

(3) MS-Excel 読上げ

Microsoft Excel のデータを読み上げるためのマクロ機能を提供する。

(4) MS-Word 読上げ

Microsoft Wordの文書を読み上げるためのマクロ機能を提供する。

(5) メッセージ読上げ

システムの状態を常時監視して、システムからのエラーメッセージを読み上げたり、システムの開始/終了時などに特定のメッセージを読み上げるプログラムである。

2. 音声認識を利用した地図検索システム

音声認識を利用するメリットとして、多項目からの検索が迅速かつ容易に行えることがあげられる。我々は警察や消防の通信指令室などで使用される、目的地の周辺地図を音声認識を用いて迅速に検索する地図検索システムを開発した。今回、作成したシステムには地名(字名まで)3万4230件、目標物名7万4937件、公衆電話ボックス名2347件の計11万1514件が登録されている。このように大規模なシステムを構築できるようになった理由の一つは、単語の読みを入力するだけで、任意の単語を認識できるようになったためである。

2.1 本システムの操作例

地名「福島市森合字森下」を検索する場合を例にとって、本システムの操作方法を説明する。まず「フクシマ」を発声すると、「森合」などの福島市の大字名が表示される。次に「モリアイ」と発声すると、「森下」などの森合に属する小字名が表示される。さらに、「モリシタ」と発声すると、目的地が特定され相当する地図が表示される。

目標物「東芝福島工場」を検索する場合について説明する。まず市町村を確定し、

「カイシャ」などの分類名を発声し決定する。次に目標物名の「トウシバフクシマコウジョウ」を発声する。会社名の場合、3千単語を越えるような大語彙となるが、本システムではこれを一括で認識する。認識単語数を大きく増加させることによって、地図検索など多くの項目を入力する必要があるシステムへの音声認識の応用が可能になった。さらに、「市町村-大字-(小字)」、「市町村-分類-目標物」のように認識対象を階層化することによって、総検索単語数を十万件以上にすることが可能となった。ここで、括弧で囲われた項目は存在しない場合があることを示す。また、公衆電話ボックスは分類の一つとして扱っている。

2.2 システム構成

本システムは図2に示すように音声認識装置とワークステーションから構成されている。

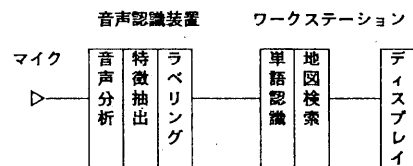


図2 音声認識地図検索システムの構成

音声認識装置は音声分析、特徴抽出、ラベリングを行う。一方、ワークステーションは単語単位の認識を行い5位までの認識結果を求めディスプレイに表示させると共に、目的地が特定された場合、1位の結果に対する地図を表示する地図検索ソフトを備えている。ここで、2位から5位の認識結果をマウスの操作によって指定すれば、対応した目的地の地図が改めて表示される。

文献(1) 太田治徳、高橋勉、原義幸：「パソコンにおける文音声合成を利用したヒューマンインタフェース」東芝レビュー, pp14-17, Vol. 51, No. 1(1996)

(注1) Windows はMicrosoft 社の商標。

(注2) Microsoft はMicrosoft 社の商標。