



音声認識応用アプリ の現状と未来

2010年3月8日
NTTサイバースペース研究所
高橋 敏

- **長年の音声認識の研究開発により、技術レベルは着実に進歩し、応用アプリも広がってきた。しかし現状、「できるようになったこと」、「まだできないこと」の境界線は、一般に、正しく理解されているとは言い難い。**
- **本講演では、現状技術によって構築可能な音声認識応用アプリを紹介することを通して、音声認識技術の現状レベルを理解していただくことを目的とする。**
- **また、「高度な技術」と「使える機能」・「欲しい機能」とは必ずしもリンクしない。現状技術でもアイデア次第でAndroid端末向けのおもしろい音声応用アプリが考えられるだろう。発想のヒントを提供できれば幸いである**

- 約20年前の入社当時は、数百単語～数千単語の単語認識または定型文認識が研究テーマであった。
- 現在では、語彙サイズ10万語の連続音声認識が市販のPCソフトウェアで利用できるまでになった。

過去

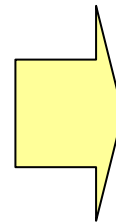
特定話者

事前登録あり

防音室録音

孤立単語発声

ハードウェア



現在

不特定話者

事前登録なし

日常環境録音

連続発声

ソフトウェア

NTT音声認識エンジン *VoiceRex*

VoiceRex 2000 *VoiceRex* 2003 *VoiceRex* 2006 *VoiceRex* 2008 次世代 *VoiceRex*



アナウンサー	一般ユーザ			
明瞭な発声	ゆっくり, 丁寧	発話速度が速い, ラフな発声		
文法に則した発話	語順の入れ替え, 不要語の挿入	オペレータの丁寧な話し言葉	一般ユーザのラフな話し言葉	
ニュース話題	タスク内話題 (レストラン検索)	コンタクトセンタでの話題	一般オフィス話題	広範囲な話題
雑音がないスタジオ	雑音が比較的少ない		雑音が多い	

コンピュータへの丁寧な発声

人と人とのコミュニケーションにおける自由な発声

- **音声インタフェース(音声コマンド)として使う**
 - スイッチの代替, もしくは, キー入力の代替
 - 人に優しいインタフェース(低リテラシー層向け)
 - ☆ 電話自動音声応答装置, 音声対話エージェントシステム
 - ☆ 音声カーナビ, (腕時計型PHS)

- **ディクテーション(口述筆記)として使う**
 - キーボードによる文書作成の代替
 - ☆ テレビ字幕作成, 議会録作成, 音声ワープロ

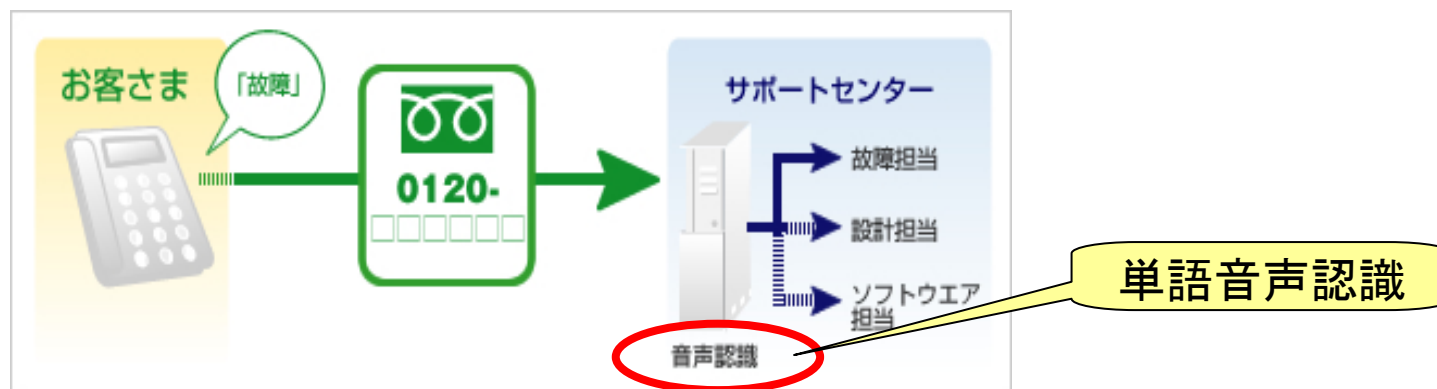
音声インタフェース (音声コマンド)

音声自動応答装置

お客様の音声を認識し、内容に合わせた着信先に自動的に振り分ける(商品名や問い合わせ内容など)



(<http://www.ntt.com/freedialin/index.html>)



得られるメリット

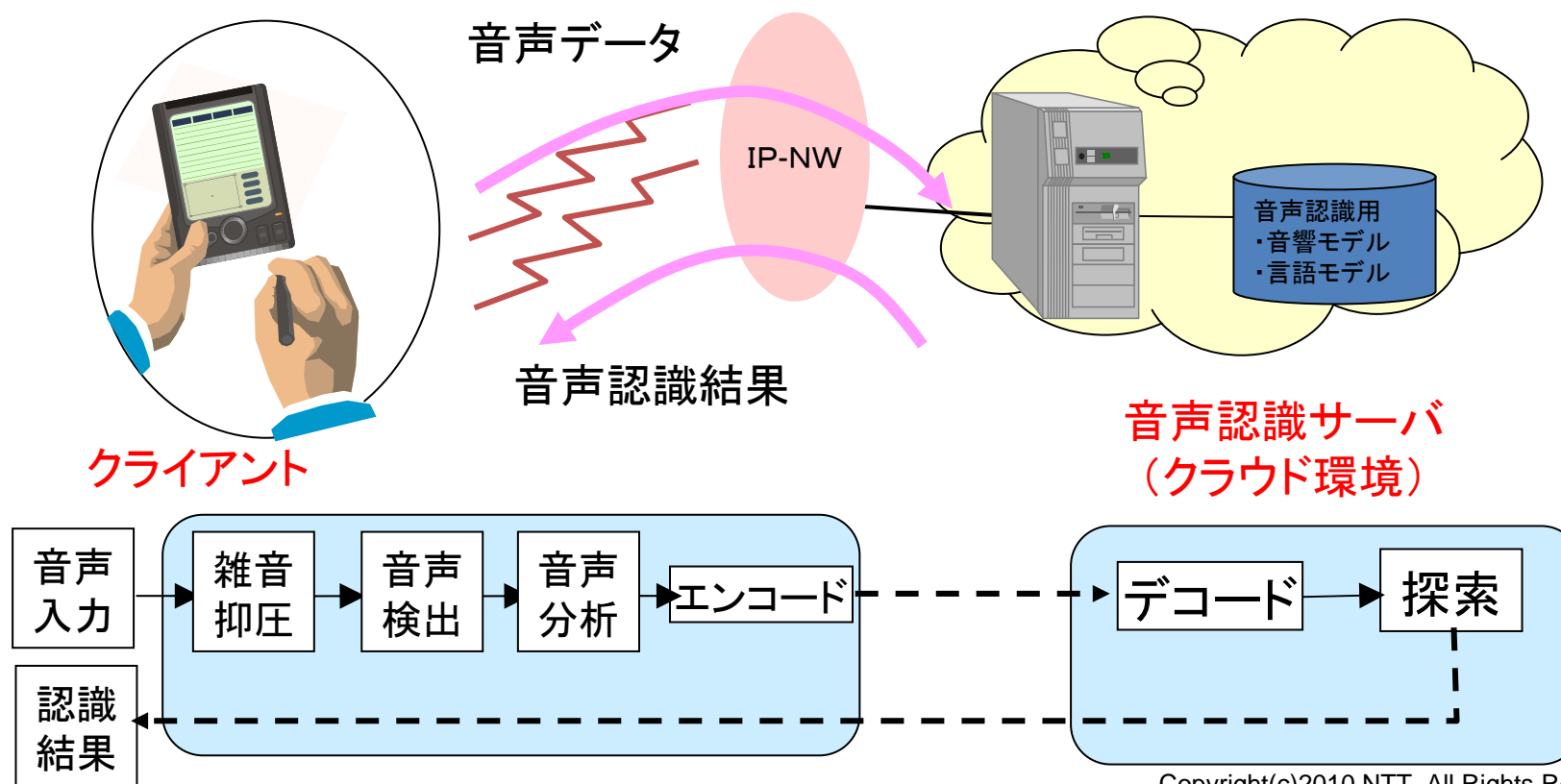
- 振り分け先を多数設定可能
- 低コストで案内
- 24時間対応

デメリットへの対応

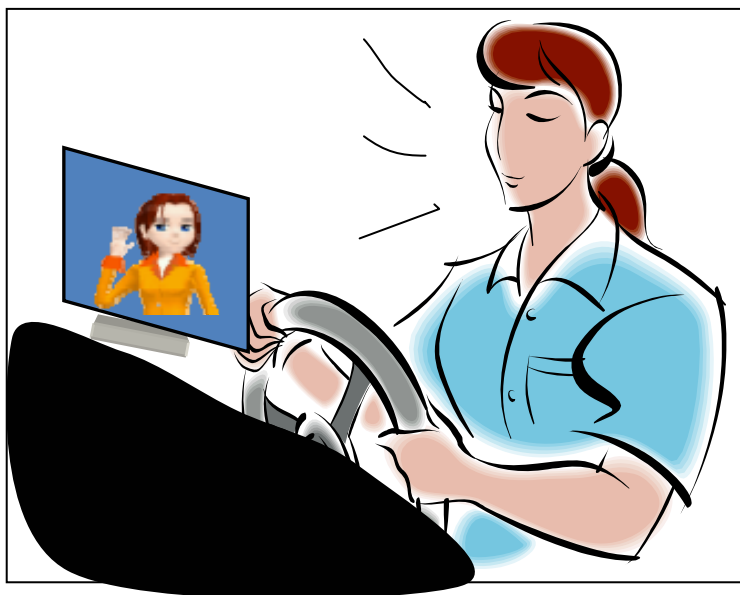
- 単語発声、PB入力も使える
- うまくいかない場合にオペレータ登場

小型端末における文字入力の煩わしさを音声入力により軽減. 通信を利用して入力音声をサーバに転送. バックエンドの豊富な計算機リソース, 知識リソースを活用して認識処理を実施

分散型音声認識 (Distributed Speech Recognition: DSR)



自動車運転時でも音声を使ってハンズフリー・アイズフリーでのカーナビ操作が可能



音声対話によって目的地をダイレクトに設定可能

- ・住所
- ・電話番号
- ・ランドマーク
- ・ジャンル

工場や倉庫での種分け作業，検査作業でもハンズフリー，アイズフリーが求められる

雑音下音声認識が課題

マルチモーダル音声対話エージェント



概要

日常会話に近い発声で、キャラクターエージェントと対話しながらタスクを実行する対話型サービスを提供します

特長

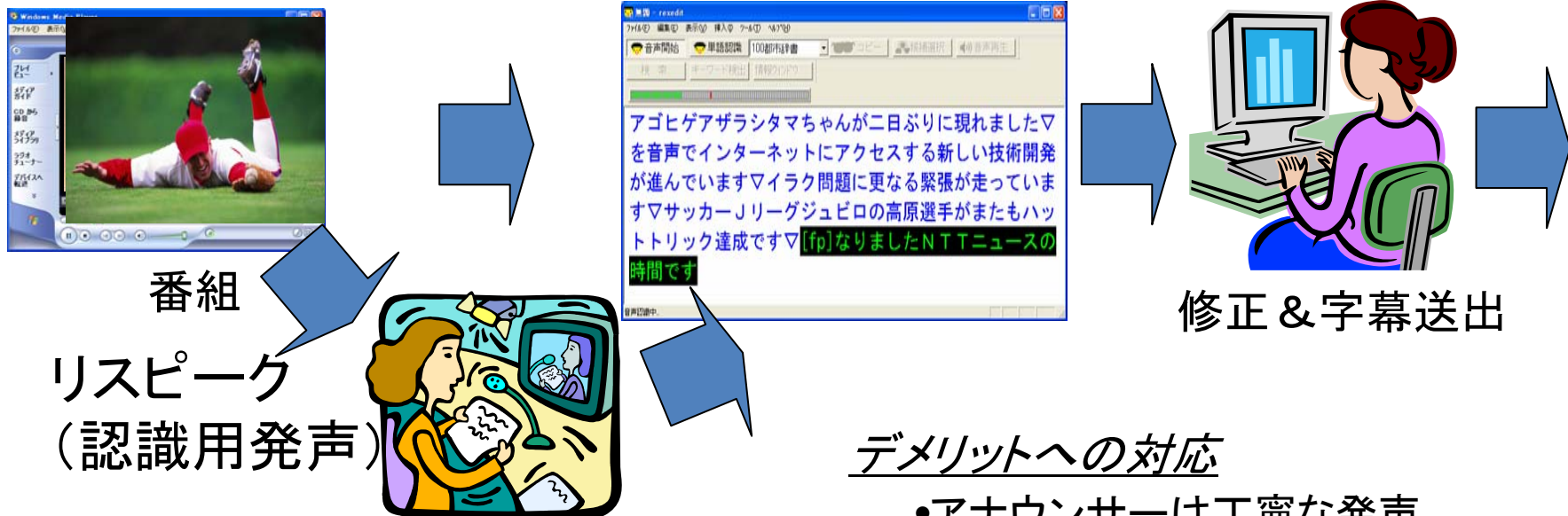
- 言い回し、語順を制限しない日常会話に近い発声を受け付け
- 自然な発話からキーワードを抽出して用件を理解
- 各種システムを容易に作成できるツール類を用意

The image shows two screenshots of a multi-modal voice dialogue agent interface. The top screenshot is for TV program reservation, displaying a grid of programs from NHK, Japan TV, and TV Asahi. A character agent asks, "So, starting at 8 PM, can I reserve the Friday drama?" A yellow box highlights "TV番組予約". The bottom screenshot is for a shop search, showing a search form with "Area (Station Name): Shinjuku" and "Genre: Italian". A character agent asks, "Um, Shinjuku, but there are no Italian shops, right?" A yellow box highlights "店舗検索". A map of the Kanto region shows station names and search paths.

ディクテーション (口述筆記)

テレビ字幕作成支援システム

アナウンサーのニュース音声を認識. または, 実況中継の生放送音声をリスピーカーが音声認識用に再発声してた音声を認識.
誤認識箇所は人手で修正して字幕送出



得られるメリット

- 字幕作成の負担軽減
- 特殊スキル者が実施するよりも低コスト
- 字幕付与義務化への準備

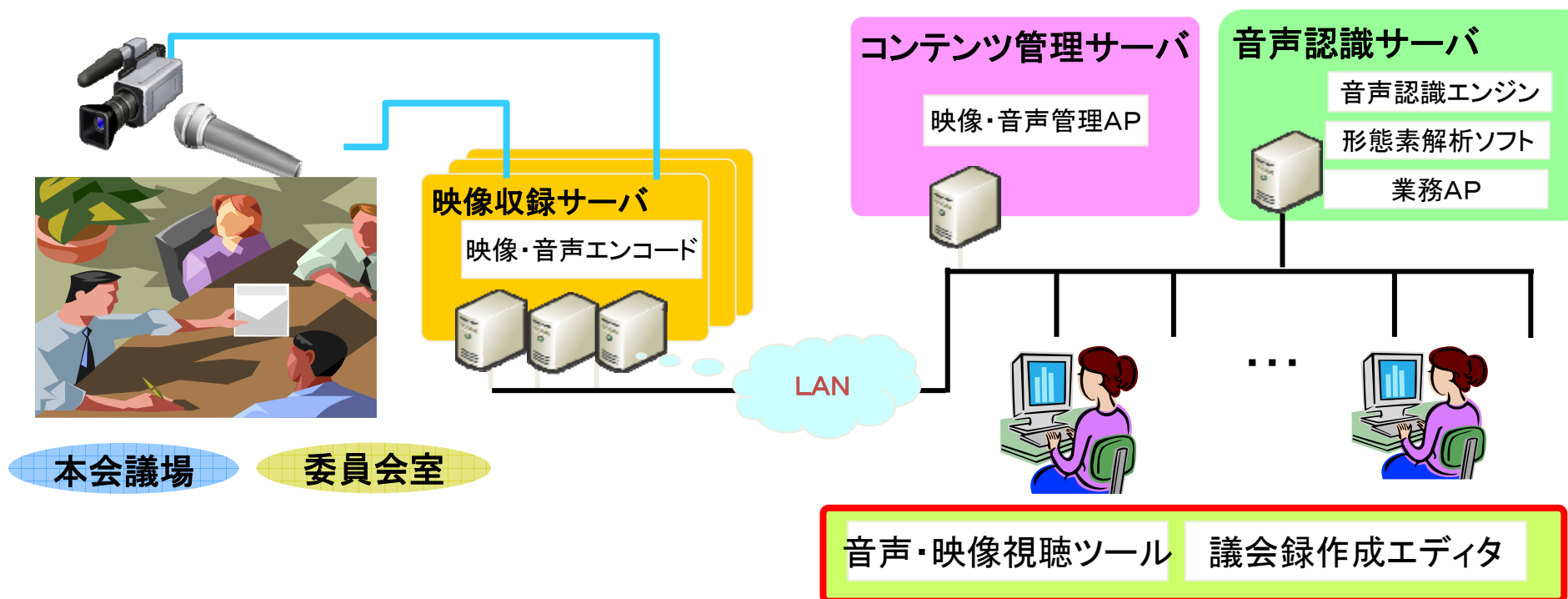
デメリットへの対応

- アナウンサーは丁寧な発声
- ニュースの話題は事前に把握可能
- アナウンサー(話者)は事前に把握可能
- 字幕向けに再発声
- 最後は人が修正

音声議会議録作成支援システム

議場で収録される音声を音声認識してテキスト化し、修正・編集することで、議会議録の作成を効率化するシステム

- ・音声認識サーバと連携した収録・管理・認識・編集の一連のシステム
- ・速記者など特殊技能者を必要としない



音声認識の利点と欠点

- **人に優しいインターフェースである**
 - 高齢者をはじめ，低リテラシー層でも分かりやすい入力手段である。
- **少ない労力で文字入力できる**
 - 音声入力に慣れれば，キーボードよりも高速に文字入力が可能である。
- **手がふさがっていても利用可能である**
 - 両手を使った作業中でも文字(コマンド)入力が可能である。
- **省スペース**
 - マイクとディスプレイ(またはスピーカ)さえあれば実現可能である。

- (人間がそうであるように,) 機械も聞き間違いをする
 - 音声認識率は100%ではなく, 誤認識が発生する.
 - 機械には知識がないので, とんでもない誤認識をする

- 入力音声に制約条件がある
 - 騒音下では認識率が低下する. 人の声の雑音に弱い.
 - 滑舌よく丁寧に発声される必要がある
(音素と音素の音声特徴が類似して混同してしまう)

- 音声認識辞書にない単語は認識できない
 - 発声する内容(話題)が, ある程度, 限定される必要がある
(単語の出現頻度や単語連鎖を学習するため)

■ 議会録作成

適用可能

- 原稿がある
- プレゼン形式(独話)
- 滑舌よく発声
- 静かな議場
- 専用マイクがある
- 発話に即した過去の議事録がある

- 発話内容をそのまま書き起こし残す

vs. (会社の) 会議録作成

適用困難

- 思いつくままにしゃべる
- 対話または多人数参加
- 発声が曖昧
- 雑音がある会議室
- 専用マイクなし(ICレコーダ)
- 議事録はない
(あっても要約した議事録)

- 要点をまとめて残す

- 認識性能が低い.
- なぜ誤認識したのかわからない.
- 同じ間違いを何度もする.
- 環境の変化に弱い(ロバストネスが低い)
- 何をしゃべったらよいかわからない, どこまで対応できるのかわからない.
- 新しい単語(最近登場した新語)が認識されない
- 発声しながら文章を考えるのは難しい(話し言葉と書き言葉は異なる)
- 誤認識箇所が発見しにくく, 修正がしにくい.

- **音声認識が適用しやすい領域(技術的観点から)**
 - 蓄積された大量テキストがある(新聞データ, 議会録, クエリーログ)
 - 何を発声すればよいかすぐわかる(住所, 電話番号)

- **誤認識への対応**
 - 誤認識が即アプリの致命的なエラーにならない
 - 修正する手段を備える
 - 徹底的に精度を高める(タスクチューニング)

- **認識精度向上のための対応**
 - 誤認識した要因をフィードバックする
 - 声の大きさ,
 - 発話のタイミング
 - 周囲雑音
 - 滑舌の悪さ
 - 語彙外発声

- **システムの透明性を高める**
 - 誤認識してもどういう状態にあるか, 次に何をすればよいかわかる

- **音声認識技術は発展途上. しかし着実に進展している.**
- **様々な応用アプリを構築し, 研究開発にフィードバックすることが重要.**
- **どのような領域でどう使われると真価を発揮できるか, 柔軟な発想が求められる.**
- **「端末に向かってしゃべる文化」をいかに醸成するか.**
- **あなたなら, ちょっと癖のある音声認識技術を何に応用しますか?**