

WOZによるオンライン修正が可能な事例ベース音声対話システム

村尾 浩也^{†,††} 河口 信夫^{††} 松原 茂樹^{††} 山口由紀子^{†††} 武田 一哉^{††}
稲垣 康善^{††††}

† 三洋電機(株) デジタルシステム技術開発センター BU 〒 573-8534 大阪府枚方市走谷 1-18-13
 †† 名古屋大学統合音響情報研究拠点 〒 464-8603 名古屋市千種区不老町 1
 ††† 名古屋大学情報連携基盤センター 〒 464-8601 名古屋市千種区不老町 1
 †††† 愛知県立大学 〒 480-1198 愛知県愛知郡長久手町熊張字茨ヶ廻間 1522-3
 E-mail: †murao@hr.hm.rd.sanyo.co.jp

あらまし 筆者らはこれまで、柔軟な対話制御を行うことを目指し、対話事例を利用する音声対話システムにおいて、WOZ(Wizard Of Oz) システムと人間との対話のログ情報から事例データを生成する手法を提案してきた。今回我々は、事例ベース対話処理をコンテキスト依存発話を取り扱うことができるように拡張するとともに、対話システムとWOZシステムをオンライン接続し対話中にリアルタイムで事例を修正することにより、事例データを効率的に修正することができる新しい枠組み“GROW アーキテクチャ”を提案する。そしてさらに、その手法に基づいて動作する情報検索向け音声対話システムの概要と、その評価結果について述べる。

キーワード 音声対話, 音声認識, Wizard of OZ, 意図理解, 応答生成, 自動車, 対話コーパス, GROW アーキテクチャ

Example-based Spoken Dialogue System with Online Example Correction

Hiroya MURAO^{†,††}, Nobuo KAWAGUCHI^{††}, Shigeki MATSUBARA^{††},
Yukiko YAMAGUCHI^{†††}, Kazuya TAKEDA^{††}, and Yasuyoshi INAGAKI^{††††}

† Digital Systems Development Center BU, SANYO Electric Co., Ltd.
Hashiridani 1-18-13, Hirakata-shi, Osaka, 573-8534 Japan
 †† Center for Integrated Acoustic Information Research, Nagoya Univ.
Furo-cho, Chikusa-ku, Nagoya-shi, 464-8603 Japan
 ††† Information Technology Center, Nagoya Univ.
Furo-cho, Chikusa-ku, Nagoya-shi, 464-8601 Japan
 †††† Faculty of Information Science and Technology, Aichi Prefectural University,
Nagakute-cho, Aichi-gun, Aichi, 480-1198, Japan
 E-mail: †murao@hr.hm.rd.sanyo.co.jp

Abstract This paper proposes a new framework of spoken dialogue system, named “GROW architecture” that is based on the examples of dialogue between human and a Wizard-of-OZ (WOZ) system. Along with modeling of information retrieval dialogue, a system for shop information retrieval in a driving car has been designed. The system refers to the dialogue examples to find an example that is suitable to generate a query or a reply. Using the WOZ system for correcting dialogue examples via network, it becomes efficient to construct dialogue examples.

Key words spoken dialogue, speech recognition, Wizard of OZ, speech understanding, reply generation, car environment, dialogue corpus, GROW architecture

1. はじめに

近年、計算機能力の向上などを背景に、大量の音声やテキストのコーパスを利用した音声対話システム構築手法が盛んに研究されている ([1]~[3])。筆者らはこれまで、情報検索対話において柔軟な対話制御を行うことや、音声認識誤りに対してロバストな意図理解を行うことなどを旨とし、対話事例を利用する音声対話制御手法を提案してきた [4], [5]。この手法では発話データとその発話に伴う検索式などの付加情報の対を事例として用いることにより、システムの行動を決定する。このシステムを有効に動作させるためには、大量の事例データを収集することが必要であるが、これまでは、人間対人間の対話データを収集し、発話に対して人手で付加情報の付与を行っており、大きな労力を伴うという問題点があった。そこで我々は、WOZ (Wizard of OZ) 方式を利用した擬似音声対話システム (以下、WOZ システムと呼ぶ) と人間との間で行われた対話事例を事例データとして利用する新しい事例ベース対話システム構築手法を提案した [6]。この方法を用いると、WOZ システムが残したログ情報から発話データと検索式の対を容易に取り出すことができるため、事例データ収集の手間を大幅に軽減することができる。しかしながら、取り扱う対話事例が、対話コンテキストを扱わない発話、つまり対話の初期状態における対話のみであること、データを一度ログとして保存した後に事例データに加工する必要があり、集めた事例データをすぐに対話システムに反映できないこと、などの問題点があった。

本稿では、それらの問題点に鑑み、事例ベース対話処理をコンテキスト依存発話を取り扱うことができるように拡張するとともに、WOZ で収集したデータをリアルタイムで事例データとして利用することができる新しい枠組み“GROW アーキテクチャ”を提案する。そしてさらに、本手法に基づいて構築した情報検索向け音声対話システムの概要と、その評価結果について述べる。

2. 事例に基づく対話処理

まずはじめに、我々がこれまでに提案してきた事例に基づく対話処理の概要について述べる [4] [5] [6]。

2.1 情報検索対話のモデル化

人間のオペレータが情報データベースを検索し、ユーザに対して情報を提供する状況における、オペレータとユーザ間の対話は、図 1 のようにモデル化することができる。

(1) **要求** ユーザの要求発話を受けたオペレータは、ドメイン知識や現在の対話コンテキストを参照しながら検索式を生成する。検索式の生成は、一般にはコンピュータなどの検索用ツールを操作することで間接的に行われる。

(2) **検索** 生成された検索式により、検索が実行される。

(3) **検索結果** 検索結果が生成される。

(4) **応答** オペレータは、検索結果と対話コンテキストに基づいて応答を行う。

このように情報の流れを整理すると、次のように考えることができる。図 1 で、オペレータは対話の進行のために次の 2 つ

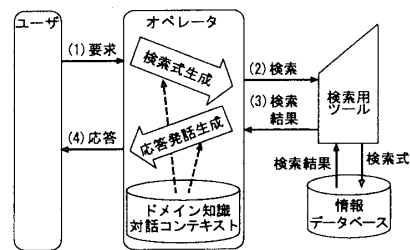


図 1 情報検索対話における対話情報の流れ

Fig. 1 Information flow of information retrieval dialogue

の判断を行っている。

判断 1 ユーザ発話と対話コンテキストに基づく検索式生成

判断 2 検索結果と対話コンテキストに基づく応答発話生成

熟練したオペレータは、ユーザ発話の表層情報だけでなく、ドメイン知識や対話コンテキスト、検索結果などの情報を使用して、その状況に最も適した検索式あるいは応答発話を生成している。つまりこの場合の対話処理とは、ユーザの発話やドメイン知識などの入力情報から検索式などの出力情報へのマッピング操作であるとみなすことができる。そこで我々は、このような対話進行のための「判断」を、熟練した人間のオペレータの行動事例を参照して行うことができるのではないかと考えた。つまり、大量の対話事例すなわち前述のようなマッピング情報を蓄積し、ある入力情報に対する出力情報を、類似した対話事例におけるマッピング情報を基に決定するのである。

2.2 事例に基づく検索式および応答発話の生成

2.2.1 事例データの構造

ユーザとオペレータが情報検索対話を行う際に行われる 2 つの「判断」は、以下のような入出力情報間のマッピングとしてあらわすことができる。

- 検索式生成のための「判断」における入出力情報
入力 ユーザの発話内容、対話コンテキスト
出力 検索式
- 応答発話生成のための「判断」における入出力情報
入力 ユーザの発話内容、対話コンテキスト、検索結果
出力 応答発話

従って事例データとしては、このような入出力情報を網羅するものを保持すればよい。具体的には、以下の情報を要素とする事例データを多数収集し、事例データベースを構築する。

- (1) ユーザの発話テキスト
- (2) 対話コンテキスト
- (3) 検索式
- (4) 検索結果
- (5) 応答発話テキスト

事例データベースを参照して上述の 2 つのマッピングを実行する過程は以下の通りである。

2.2.2 検索式生成の過程

入力発話およびその時点の対話コンテキストについて、事例データベースの中から最も類似した事例を抽出する。この最類似事例中の検索式を入力発話と対話コンテキストに合うように

	最低 価格	最高 価格	距離	ジャンル ジャンル	駐車場 駐車場	キーワード キーワード	検索 結果数	店名	ソート順 ORDER BY
発話前 スロット			<=5000	フード					距離
発話後 スロット					>=10			マクド ナルド	

図2 スロット情報の構成

Fig. 2 Configuration of dialogue slot

修正し、検索式を生成する。

2.2.3 応答発話生成の過程

入力発話とその時点の対話コンテキストおよび情報データベースから得られる検索結果について、事例データベースの中から最も類似した事例を抽出する。この最類似事例中の応答発話を現在の状況に合うように修正し、応答発話を生成する。

前稿 [6] では、対話コンテキストが不要な状況（対話の開始状態）、つまり事例データベースにコンテキスト情報を含まない条件での事例ベース対話処理について検討した。本稿ではそれをさらに拡張し、対話コンテキストを含む事例データベースを構築し、コンテキスト依存発話を取り扱えるようにした。

3. スロット情報を利用したコンテキスト依存対話処理

前述のように、対話コンテキスト依存の対話処理を行うためには、対話コンテキストを事例として保持することが必要である。対話コンテキストの情報を効率よく蓄積するために、我々は、スロット情報を利用することとした。

スロット情報とは、属性と値 (Attribute-Value) のペアの集合であり、対話システムでは、対話の進行で得られた情報を属性ごとに保存するために用いられる。スロットの値は対話中に獲得された、対話進行に必要な情報であり、その埋まり具合は対話の進行状況をあらわすため、対話コンテキストを反映した情報であると考えられる。我々は、以下の点について考慮し、事例データベース中でのスロットの保持方法を検討した。

- 検索式生成のための必要十分な情報を取り出せること
- ある発話が行われる時点の対話の状況を記述できること
- 発話前後の対話の状況の変化を記述できること

その結果、事例データベース中に持つスロット情報として、図2のような形式を採用した。1つの事例データ中には、スロットデータが計2つ保持される。各スロットデータは、「店名」「駐車場」などの属性と「マクドナルド」「>=10 (10台以上の意)」などの値が保持される。2つのスロットデータのうち、「発話前スロット」はユーザの発話が行われる時点の状況を、「発話後スロット」はユーザの発話が行われた後のスロットの変化値(差分)が保持される。

従って、事例データベースの要素は以下の5項目とした^(注1)

- (1) ユーザの発話テキスト
- (2) 発話前のスロット情報

(注1)：本来は、事例データとして検索式を保持する必要があるが、本システムではスロット情報から一意に検索式(検索エンジンとして今回新たにSQLを使用することとした)を生成できるようにしたため、事例データベースにスロット情報を保持することで検索式が表現できるようになっている。

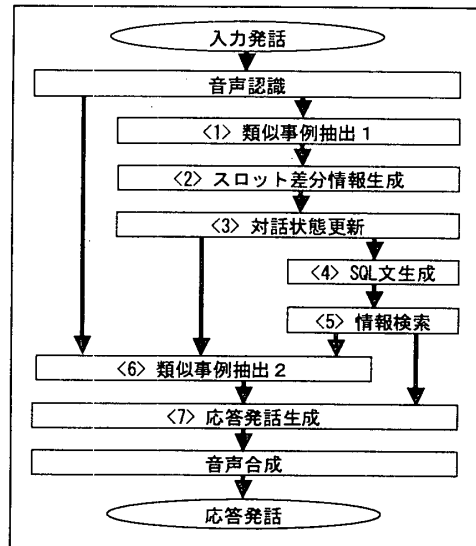


図3 対話処理の流れ

Fig. 3 Flow of dialogue process

- (3) 検索結果
- (4) 発話後のスロット情報
- (5) 応答発話テキスト

また、現在の対話の進行状況(対話状態)を保持するために対話システム内部にもスロット情報(内部スロット)を持つ。

スロット情報を利用した事例ベース対話処理における、処理の流れを図3に示す。

ある発話が入力されると、品詞情報を基にキーワードが抽出され、検索式生成用の類似事例が検索される(図3の<1>)。ここでは、以下のスコアの重み付き加算によって求められたスコア値を基に、もっとも良いスコアを示す事例が抽出される。

- 入力発話と事例中の発話との類似スコア
- 内部スロットと事例中の発話前スロットとの類似スコア

その後、抽出した事例中の発話後スロットを基に、入力発話の内容を参照しながら内部スロットの差分情報(スロット差分情報)が生成される(<2>)。スロット差分情報は、単語クラス情報(単語を意味的基準でクラス化したもの)を用いて、発話後スロット中の単語を入力発話に合うように入れ替えることによって生成される。

そして生成されたスロット差分情報を内部スロットに上書きして対話状態の更新が行われる(<3>)とともに、SQL文がスロット情報を基に生成され(<4>)、検索が行われる(<5>)。

検索結果が得られた後、再び類似事例の検索が行われる(<6>)。ここでは、<1>でのスコアに加え、以下の情報を重み付き加算したスコア値が用いられる。

- 検索結果数と事例中の検索結果数との類似スコア
- 更新後の内部スロットと事例中の発話後スロットとの類似スコア

抽出した事例中の応答発話を、検索結果の情報を基に単語を入れ替えることで修正し(<7>)、最終的な応答発話生成され、音声合成に渡される。

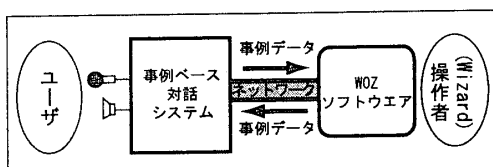


図4 GROWアーキテクチャの構成
Fig.4 Configuration of GROW architecture

4. GROWアーキテクチャ ～WOZによるオンライン事例修正～

これまでに我々は、事例として用いる付加情報付き対話データを効率的に収集するために、WOZ (Wizard of OZ) 方式を使用することを提案してきた [6]。WOZ 方式による対話とは、対話システムの構成要素の一部を人間が代行するが、使用者にはあたかも機械が全てを行っているかのように見せかけたものである。WOZ 方式によって人間との間で情報検索対話を行い、その際に生成された検索式や検索結果をログ情報として保存、利用することによって事例データベースを構築し、対話制御に利用する。

WOZ 方式を用いるメリットとして、

(1) WOZ システムの行動に埋め込まれた知識 (入力発話から検索式への変換, 入力発話から応答発話への変換など) を対話システムの知識として利用できる

(2) 実際の対話システム使用環境により近い、つまり人間対機械の対話に近い状況下の対話事例が収集できるなどが挙げられる。

前稿では、WOZ システムは単独で用いられ、すべての発話に対して WOZ システムが応答を行い、WOZ システムが残したログ情報を利用して事後的に事例データを作成していた。しかし実際には、WOZ による事例の生成は、対話システムが自動では正しい応答文を生成できない場合にのみ必要である。そこで我々は、事例ベース対話システムと WOZ システムをオンライン接続し、必要なとき、つまり対話システムが処理を誤るときのみ WOZ システムが介入する新しい枠組み **GROW** アーキテクチャを提案する。

4.1 システムの構成

音声対話による情報検索アプリケーションとして、自動車内の店情報検索をターゲットとし、事例ベース対話システムと WOZ システムとがオンラインで協調処理を行うシステムを構築した。今回構築したシステム全体の構成を図 4 に示す。事例ベース対話システムおよび WOZ ソフトウェア (操作者=ウィザード、が操作する PC 上のソフトウェア) は WindowsXP 上で動作する。使用プログラム言語は C++である。

GROW アーキテクチャでは、事例ベース対話システムと WOZ ソフトウェアはネットワーク接続されており、対話システムが自動生成したスロット情報や応答文などの事例データはネットワークを通して WOZ ソフトウェアに送られる。ウィザードは、送られてきたデータを見て修正の必要があれば画面を操作して修正を行い、修正後のデータを事例ベース対話シ

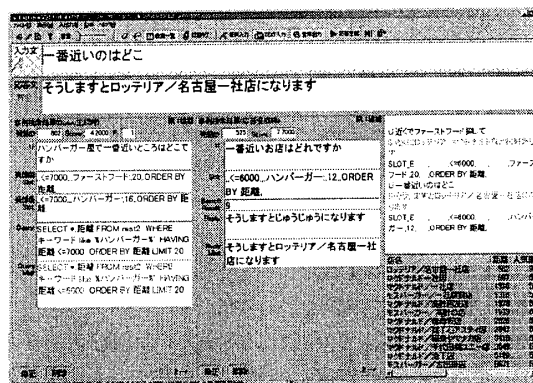


図5 事例ベース対話システムの画面例
Fig.5 An example of display of example-based dialogue system

テムに送り返す。修正の必要が無ければそのまま送り返す。対話システムは送られてきた事例データが修正されている場合にはそのデータを対話システムの事例データに追加する。

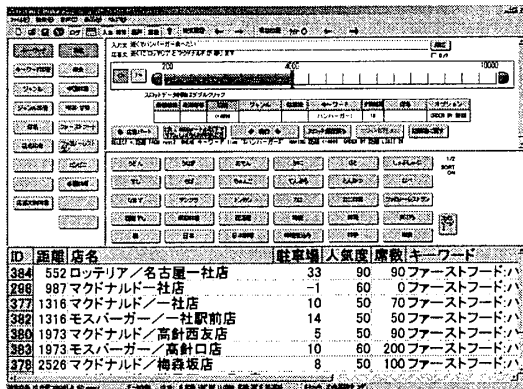
事例ベース対話システムは、3. 節で述べたアルゴリズムによってコンテキスト依存の対話処理を行うことができる。対話システムの動作状況をあらわす画面を図 5 に示す。

WOZ ソフトウェアは対話システムから受け取ったスロット情報を修正、再検索を行う画面 (入力パート) と、応答文を修正する画面 (応答パート) から成る。各パートの画面例を図 6, 7 に示す。操作者が簡単に操作を行うことができるようにタッチパネルを採用し、GUI 表示された画面中のボタンをタッチすることにより、スロットの充当、情報データベースの検索、応答発話の生成を行うことができる。

WOZ システムと人間の間で自然な対話を行うためには、WOZ ソフトウェアの高い操作性が求められる。我々は、事前に収集した人間対人間の対話コーパスの情報を用いて WOZ ソフトウェアの設計を行い、ユーザとの対話中に実時間で操作を行うことができるように以下のような種々の工夫を施している。

入力パート (図 6) では、内容別に木構造に整理されたキーワードの中から適切なキーワードをタッチパネル操作によって選択することでスロットを充当し、素早く店情報の検索が実行できるようにした。駐車場の台数、距離など、数値の範囲が充当されるスロットについては、スライドバー (画面上部) を利用した直感的な範囲選択を行うことができるようになっている。検索結果は画面下部にリスト形式で表示される。

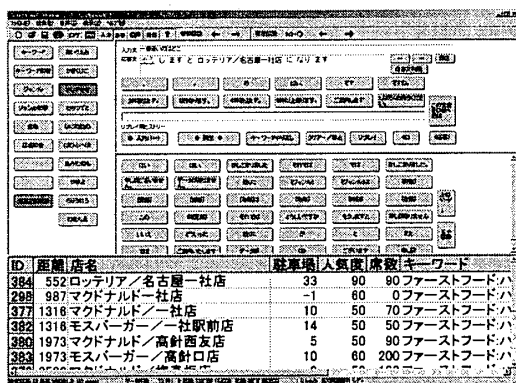
応答パート (図 7) には、応答発話生成用のテキスト入力ボタンと店情報の検索結果リストが表示される。テキスト入力ボタンは、応答発話生成に用いられる単語や文節、定型短文などに対応しており、そのボタンを順次押すことにより応答発話テキストが生成される。テキスト入力ボタンの配置は、人間対人間の対話コーパスから事前に求めた、各テキスト要素間の接続頻度 (応答文バイグラム) を基に決定した。すなわち、直前に入力されたテキストに後続する頻度の高いものから順にテキスト入力ボタンが表示され、応答発話テキストを素早く生成することができる。さらに、固有名詞である店情報などを素早



ID	距離	店名	駐車場	人気度	席数	キーワード
384	552	ロッテリア/名古屋一社店	33	90	90	ファーストフード/ハ
298	987	マクドナルド一社店	-1	60	0	ファーストフード/ハ
377	1316	マクドナルド一社店	10	50	70	ファーストフード/ハ
382	1316	モスバーガー/一社駅前店	14	50	50	ファーストフード/ハ
380	1973	マクドナルド/高針西友店	5	50	90	ファーストフード/ハ
383	1973	モスバーガー/高針口店	10	60	200	ファーストフード/ハ
378	2526	マクドナルド/梅森坂店	8	50	100	ファーストフード/ハ

図 6 WOZ システムの画面の例 (1) 検索式生成画面

Fig. 6 An example of display of Wizard of OZ System (1): Query generation part



ID	距離	店名	駐車場	人気度	席数	キーワード
384	552	ロッテリア/名古屋一社店	33	90	90	ファーストフード/ハ
298	987	マクドナルド一社店	-1	60	0	ファーストフード/ハ
377	1316	マクドナルド一社店	10	50	70	ファーストフード/ハ
382	1316	モスバーガー/一社駅前店	14	50	50	ファーストフード/ハ
380	1973	マクドナルド/高針西友店	5	50	90	ファーストフード/ハ
383	1973	モスバーガー/高針口店	10	60	200	ファーストフード/ハ

図 7 WOZ システムの画面の例 (2) 応答文生成画面

Fig. 7 An example of display of Wizard of OZ System (2): reply generation part

く入力できるようにするため、検索結果リストの店名の部分を直接タッチすることにより応答発話テキストに店名を含めることができるようにした。生成された応答発話テキストは、音声合成器 [8] で音声に変換され出力される。

2つのパート間は切り替えボタンにより自由に移動することができる。また、あいづちなど、被験者発話中に即座に提示すべき音声については入力パートにも発話生成ボタンを設け、自然な対話を行うことができるよう配慮している。

5. 評価

5.1 コンテキスト依存対話の評価手法

対話システムの評価はコンテキストに依存した発話を含めて評価対象とする必要がある。一般的に対話システムは被験者によるタスク達成率などの基準で評価されることが多いが、その場合客観的で再現性のある評価は難しい。発話データを用いた客観的な評価を行うためには対話の状況を評価時に再現できるようにすることが必要である。

我々は、スロット情報により対話コンテキストの情報を表現する提案システムの特徴を活かし、コンテキスト依存の対話処理に対する客観評価を行うことができると考えた。

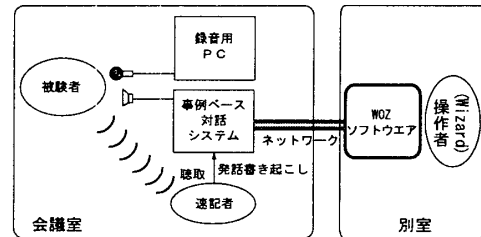


図 8 実験設備の配置図

Fig. 8 Arrangement of equipments for the experiment

提案システムでは、対話コンテキストの情報は全てスロット情報として表現されている。つまり、システムに対しスロット情報とその状況下での発話データを与えることにより、様々な対話の状態を再現し、その際のシステムの応答を客観的で再現性のある形で評価することが可能となる。本稿では以下のような手順でコンテキスト依存の対話処理の評価を行う。

(1) テストセットの収集 被験者と WOZ との間で対話を行い、発話データとそれに対応してウィザードが作成したスロット情報を記録する。これを整理し、発話情報と発話時のスロット情報との組を作成した。これをテストセットとして用いる。

(2) 対話処理の評価 まず、収集したテストセットのスロット情報をシステムに与え、対話の状況を再現する。その後、発話テキストを与えてその応答を評価する。本稿では、事例データベースが持つ事例の数を変化させ、生成された応答発話の妥当性で評価を行う。

以下、評価手順の詳細について述べる。

5.2 対話データの収集

対話データは、被験者を情報検索を行うユーザとし、名古屋大学内の会議室で収集した。WOZ ソフトウェアが動作する PC を別室に配置し、ウィザードは被験者から見えない位置で操作を行う。対話の内容は別の PC を用いて全て録音した。実験の状況を図 8 に示す。

データ収集は以下のような手順で、一人あたり約 15 分間行った。

(1) 被験者に対し、12 種の状況プレートから任意に 1 枚を選び提示する。(一例を図 9 に示す)

(2) 被験者はプレートの内容を基に、自由に店情報を問い合わせる対話を行う

(3) 被験者の発話内容を速記者がその場で書き起こし、対話システムへ入力する

(4) 対話システムで自動生成されたスロット情報と応答発話が WOZ に送られる

(5) ウィザードは必要に応じてスロット情報と応答発話の修正を行い対話システムに送り返す

(6) 被験者は応答発話 (音声合成音) を聴き対話を続ける
収集したデータについて表 1 に示す。データ収集は、テストセットを作成するためのセッションと、事例データを追加するセッションの 2 つを行った。総被験者数は 9 名で、そのうち 5 名は、上記の両方のセッションを行った。残りの 4 名については、事例データの追加セッションのみを行った。

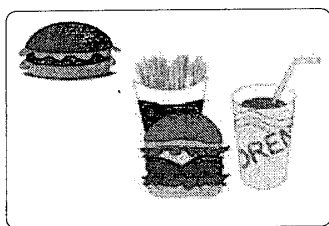


図9 状況プレートの例

Fig. 9 An example of prompting panel

表1 収集したWOZ対話データ

Table 1 Collected WOZ Data

収集したデータ	総被験者数	9
	テストセット用発話数	95
	事例データ追加用発話	356
実験に使用したデータ	テストセット発話数	66
	type1(コンテキスト非依存)	19
	type2(詳細情報問い合わせ)	36
	type3(店最終選択)	11
	追加用事例データ数(クローズ話者)	100
	追加用事例データ数(オープン話者)	120

5.3 対話処理の評価

テストセットとして集めたデータのうち、66発話を実験に使用した。テストセットはあらかじめ対話コンテキスト情報利用の有無や、発話の内容により type1 から type3 の3種類に分類し評価を行った。type1 は対話の初期状態に行われた発話データである。type2 は主に type1 の発話時に提示された店の情報をさらに詳細に問い合わせる発話である。type3 は最後に一つの店を決定し、そこへの案内を依頼するような発話である。type の分類は発話を見て主観的に決定した。

事例追加のセッションは、我々がこれまでに収集してきた車内音声対話データ [7] の書き起こしを基に作成した標準事例データベース(事例数 138) を使用して行った。そうして集めた事例データを、1) テストセット話者と同話者 5 名の事例(テストセットとは異なる状況プレートによる対話: 事例数 100) 2) テストセット話者と異なる話者 4 名の事例(事例数 120) に分け、それぞれ標準事例データベースに追加し評価に用いた。

評価結果を図 5.3 に示す。

評価は、生成された応答文の妥当性を主観的に判断し、正解率(%)で行った。正解の基準は、検索式、応答文の要素として必要十分な情報(キーワードなど)が含まれていることとした。また、必要な要素が一部不足しているが、余分な誤った情報は含まないものも正解として扱った。

評価結果を見ると、ほとんど全てのケースで、事例を追加することにより検索式生成性能、応答文生成性能が向上していることがわかる。特に話者クローズ条件での性能向上が大きい。しかし、応答生成率の値はまだ低く十分とはいえない。今回の実験で使用した事例数は高々250程度と小規模なものであり、数としては十分であるとはいえない。さらに事例数を増加し性能評価を行う必要がある。(以前の研究 [5] においては、事例数

表2 評価結果(応答文正解率)

Table 2 Evaluation result (Correct rate for reply generation)

		テストセットの分類			
		全データ (データ数:66)	type1 (19)	type2 (36)	type3 (11)
事例 データ ベース	標準(事例数:138)	27.3	42.1	11.1	54.5
	標準(138)+ 話者クローズ事例追加(100)	53.0	42.1	50.0	81.8
	標準(138)+ 話者オープン事例追加(120)	42.4	52.6	30.6	63.6
					単位(%)

500程度で検索式生成の評価を行っている。ここでは事例数を約390から約540まで増加させた場合でも性能の向上が確認されている。従って応答発話生成の評価においても少なくともこの程度の規模での評価が必要である。)しかし、今回の実験においても、特に type3 において80%以上の高い正解率が得られていることから、少なくとも本方式は、検索結果を一つに確定する発話、つまり「じゃあその店でお願いします」のような定型的な発話に対しては有効に働いているといえる。

6. おわりに

本稿では、事例ベース対話処理をコンテキスト依存発話を取り扱うことができるように拡張するとともに、WOZで収集したデータをリアルタイムで事例データとして利用する新しい枠組み **GROW** アーキテクチャを提案した。

評価実験を行った結果、小規模な事例データベースに対して提案手法により事例を追加することで正しい応答文を生成する割合が向上することがわかった。特に、事例データが評価話者と同じ話者のデータで追加された場合に性能向上が大きくなるという結果となった。今後は、1) さらに大規模な事例データベースを構築して評価を行う 2) 多数の被験者による評価を行う 3) 応答文生成に失敗する原因について詳しく調査し、応答文生成性能の向上を図る などを行う予定である。

文 献

- [1] E. Levin, R. Pieraccini, and W. Eckert. Using markov decision processes for learning dialogue strategies. Proc. ICASSP98 Vol. 1, pp. 201-204, 1998.
- [2] S. Young. Talking to machines (statistically speaking). Proc. of ICSLP-2002 pp. 9-16, 2002.
- [3] 松原茂樹, 河口信夫, 外山勝彦, 武田一哉. 音声対話コーパスの収集と利用. 人工知能学会誌, Vol. 17, No. 3, pp. 279-284, May 2002.
- [4] 村尾浩也, 河口信夫, 松原茂樹, 稲垣康善. 対話事例を利用した音声対話システム. 情報処理学会研究報告, 2000-SLP-34-34, 2000.
- [5] H. Murao, N. Kawaguchi, S. Matsubara, and Y. Inagaki. Example-based query generation for spontaneous speech. Proc. of 2001 IEEE Workshop on Automatic Speech Recognition and Understanding, 2001.
- [6] 村尾浩也, 河口信夫, 松原茂樹, 山口由紀子, 稲垣康善. WOZ システムのログ情報を利用した事例ベース音声対話システムの開発. 情報処理学会研究報告, 2002-SLP-44-23, 2002.
- [7] N. Kawaguchi, S. Matsubara, K. Takeda, and F. Itakura. Multi-dimensional data acquisition for integrated acoustic information research. Proc. of 3rd International Language Resources and Evaluation Conference (LREC-2002), pp. 2043-2046, 2002.
- [8] 余田直之, 平井裕之, 橋本誠, 大西宏樹. 自然音声合成ソフトウェア. Sanyo Technical Review, Vol. 33, No. 3, pp. 55-62, Dec. 2001.