

(11)特許出願公開番号

**特開2021-51172**

(P2021-51172A)

(43) 公開日 令和3年4月1日(2021.4.1)

(51) Int.Cl.

F I

テーマコード (参考)

**G 1 O L 13/00 (2006.01)**

G10L 13/00 100M

**G 1 O L 15/10 (2006.01)**

G10L 15/10 500Z

**G 1 O L 15/22 (2006.01)**

G 1 0 L 15/22 3 0 0 Z

**G 1 O L 13/02 (2013.01)**

G 1 0 L 13/02 1 3 0 Z

**G06F 3/16 (2006.01)**

G O 6 F      3/16      6 5 0

審査請求 未請求 請求項の数 14 O L (全 63 頁) 最終頁に続く

(21) 出願番号 特願2019-173551 (P2019-173551)

(22) 出願日 令和1年9月24日 (2019.9.24)

(71) 出願人 899000068

学校法人早稻田大学

東京都新宿区戸塚町1丁目104番地

(74) 代理人 100114638

弁理士 中野 寛也

(72) 発明者 小林 哲則

東京都新宿区戸塚町1丁目104番地 学  
校法人早稲田大学内

(72) 発明者 藤江 真也

東京都新宿区戸塚町1丁目104番地 学  
校法人早稲田大学内

[最終頁に続く](#)

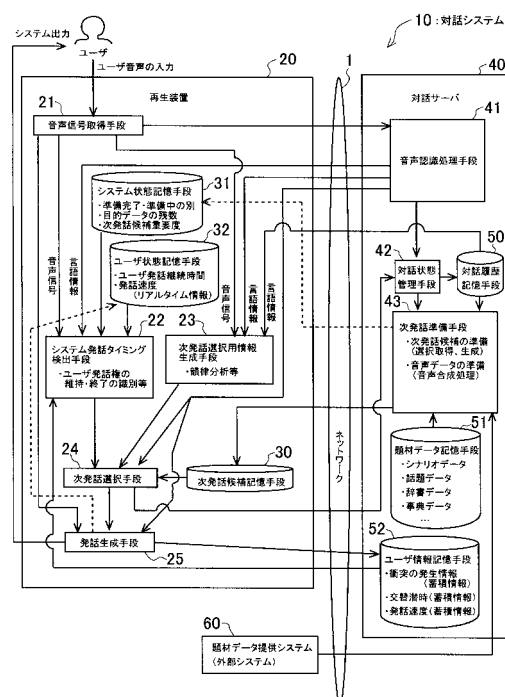
(54) 【発明の名称】 対話システムおよびプログラム

(57) 【要約】

【課題】システムの応答性を向上させ、衝突の発生を回避または抑制しつつ、不要に長いシステムの交替潜時の発生を回避または抑制することができる対話システムを提供する。

【解決手段】ユーザ発話の音声信号から抽出した音響特徴量を用いて、音声認識処理手段４１とは非同期で、ユーザ発話権の維持・終了を識別してシステム発話の開始タイミングを検出するシステム発話タイミング検出手段２２と、その検出前に、題材データ記憶手段５１等に記憶された題材データ、対話履歴情報や進行中のユーザ発話の途中までの音声認識処理の結果を用いて、システムの次発話を準備する次発話準備手段４３と、システム発話の開始タイミングの検出後に、次発話準備手段４３により準備された次発話を用いて、システム発話を再生する発話生成手段２５とを設け、対話システム１０を構成した。

【選択図】図1



**【特許請求の範囲】****【請求項 1】**

ユーザとの音声対話のための処理を実行するコンピュータにより構成された対話システムであって、

ユーザ発話の音声信号を取得する音声信号取得手段と、

この音声信号取得手段により取得したユーザ発話の音声信号についての音声認識処理を実行する音声認識処理手段と、

前記音声信号取得手段により取得したユーザ発話の音声信号から音響特徴量を抽出し、抽出した音響特徴量を用いるか、または、この音響特徴量に加え、前記音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報から抽出した言語特徴量を用いて、前記音声認識処理手段による音声認識処理の実行タイミングに依拠しない周期で、ユーザが発話する地位または立場を有していることを示すユーザ発話権の維持または終了を識別するパターン認識処理を繰り返し実行し、このパターン認識処理の結果を用いて、システム発話の開始タイミングを検出する処理を実行するシステム発話タイミング検出手段と、

このシステム発話タイミング検出手段による前記パターン認識処理の前記周期に依拠しないタイミングで、かつ、このシステム発話タイミング検出手段によりシステム発話の開始タイミングが検出される前に、題材データ記憶手段に記憶された題材データまたはネットワークを介して接続された外部システムに記憶された題材データを用いるとともに、ユーザとシステムとの間の対話履歴情報の少なくとも一部および/または前記音声認識処理手段による進行中のユーザ発話についての途中までの音声認識処理の結果を用いて、システムの次発話の内容データを取得または生成する準備処理を実行する次発話準備手段と、

前記システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、前記次発話準備手段による準備処理で得られた次発話の内容データを用いて、システム発話の音声信号の再生を含むシステム発話生成処理を実行する発話生成手段と

を備えたことを特徴とする対話システム。

**【請求項 2】**

前記次発話準備手段は、

次発話の候補となる複数の次発話候補の内容データを取得または生成する準備処理を実行する構成とされ、

前記システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、前記音声認識処理手段による音声認識処理の結果として得られた言語情報を用いて、前記次発話準備手段による準備処理で得られた複数の次発話候補の内容データの中から、前記発話生成手段で用いる次発話の内容データを選択する処理を実行する次発話選択手段を備えた

ことを特徴とする請求項 1 に記載の対話システム。

**【請求項 3】**

前記次発話準備手段は、

次発話の候補となる複数の次発話候補の内容データを取得または生成する準備処理を実行する構成とされ、

前記音声信号取得手段により取得したユーザ発話の音声信号から得られる韻律情報を用いるか、若しくは、この韻律情報に加えて、前記音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報を用いるか、またはこれらの韻律情報およびユーザ発話の言語情報に加えて、ユーザとシステムとの間の対話履歴情報のうちの直前のシステム発話の言語情報を用いて、質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図を識別するパターン認識処理を繰り返し実行する次発話選択用情報生成手段と、

前記システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、前記次発話選択用情報生成手段による処理で得られた前記ユーザ発話意図の識別結果を用いて、前記次発話準備手段による準備処理で得られた複数の次発話候補の内容デ

10

20

30

40

50

ータの中から、前記発話生成手段で用いる次発話の内容データを選択する処理を実行する次発話選択手段と

を備えたことを特徴とする請求項 1 に記載の対話システム。

【請求項 4】

前記次発話準備手段は、

次発話の候補となる複数の次発話候補の内容データを取得または生成する準備処理を実行する構成とされ、

前記音声信号取得手段により取得したユーザ発話の音声信号から得られる韻律情報を用いるか、若しくは、この韻律情報に加えて、前記音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報を用いるか、またはこれらの韻律情報およびユーザ発話の言語情報に加えて、ユーザとシステムとの間の対話履歴情報のうちの直前のシステム発話の言語情報を用いて、質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図を識別するパターン認識処理を繰り返し実行する次発話選択用情報生成手段と、

前記システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、前記次発話選択用情報生成手段による処理で得られた前記ユーザ発話意図の識別結果と、前記音声認識処理手段による音声認識処理の結果として得られた言語情報とを組み合わせ用いて、前記次発話準備手段による準備処理で得られた複数の次発話候補の内容データの中から、前記発話生成手段で用いる次発話の内容データを選択する処理を実行する次発話選択手段と

を備えたことを特徴とする請求項 1 に記載の対話システム。

【請求項 5】

前記次発話準備手段は、

次発話の候補となる複数の次発話候補の内容データを取得または生成する準備処理を実行する構成とされ、

前記システム発話タイミング検出手段は、

前記ユーザ発話権の維持または終了を識別するパターン認識処理を実行する際に、終了については、質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図のうちのいずれのユーザ発話意図で終了するのかを識別するパターン認識処理を実行する構成とされ、

前記システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、前記システム発話タイミング検出手段による処理で得られたユーザ発話意図の識別結果を用いて、前記次発話準備手段による準備処理で得られた複数の次発話候補の内容データの中から、前記発話生成手段で用いる次発話の内容データを選択する処理を実行する次発話選択手段を備えた

ことを特徴とする請求項 1 に記載の対話システム。

【請求項 6】

前記次発話準備手段による準備処理の状態を含むシステム状態を示す情報を記憶するシステム状態記憶手段を備え、

前記システム発話タイミング検出手段は、

前記ユーザ発話権の維持または終了を識別するパターン認識処理の結果および前記システム状態記憶手段に記憶されている前記システム状態を示す情報を用いて、システム発話の開始タイミングを検出する処理を実行する際に、

前記パターン認識処理の結果が前記ユーザ発話権の維持を示している場合には、システム発話の開始タイミングではないと判断し、

前記パターン認識処理の結果が前記ユーザ発話権の終了を示し、かつ、前記システム状態を示す情報が準備完了を示している場合には、システム発話の開始タイミングであると判断し、

前記パターン認識処理の結果が前記ユーザ発話権の終了を示し、かつ、前記システム状態を示す情報が準備中を示している場合には、前記次発話準備手段による準備中の処理内

10

20

30

40

50

容に応じ、直ぐに完了する処理内容として予め分類されている処理の準備中であるときには、準備完了になるまで待ってシステム発話の開始タイミングであると判断し、直ぐに完了しない処理内容として予め分類されている処理の準備中であるときには、システム発話の開始タイミングであると判断するとともに、フィラーの挿入タイミングである旨の情報を出力する処理を実行する構成とされている

ことを特徴とする請求項 1 ～ 5 のいずれかに記載の対話システム。

【請求項 7】

ユーザ発話継続時間を含むユーザ状態を示す情報を記憶するユーザ状態記憶手段を備え、

前記システム発話タイミング検出手段は、

10

前記ユーザ発話権の維持または終了を識別するパターン認識処理の結果および前記ユーザ状態記憶手段に記憶されている前記ユーザ状態を示す情報を用いて、システム発話の開始タイミングを検出する処理を実行し、この際の処理として、

( 1 ) 前記ユーザ状態記憶手段に記憶されている前記ユーザ発話継続時間が、予め定められた短時間判定用閾値以下または未満の場合には、前記パターン認識処理の結果として得られる尤度に対して設定されているユーザ発話権終了判定用閾値を標準値よりも高く設定し、予め定められた長時間判定用閾値以上または超過の場合には、前記ユーザ発話権終了判定用閾値を標準値よりも低く設定する処理と、

( 2 ) 前記ユーザ状態記憶手段に記憶されている前記ユーザ発話継続時間を用いて、前記パターン認識処理の結果として得られる尤度に対するユーザ発話権終了判定用閾値を、前記ユーザ発話継続時間が短いときには当該ユーザ発話権終了判定用閾値が高くなり、前記ユーザ発話継続時間が長いときには当該ユーザ発話権終了判定用閾値が低くなるように予め定められた関数により設定する処理と、

20

( 3 ) 前記ユーザ状態記憶手段に記憶されている前記ユーザ発話継続時間が、予め定められた短時間判定用閾値以下または未満の場合には、前記パターン認識処理の結果が前記ユーザ発話権の終了を示していても、システム発話の開始タイミングではないと判断し、予め定められた長時間判定用閾値以上または超過の場合には、前記パターン認識処理の結果が前記ユーザ発話権の維持を示していても、システム発話の開始タイミングであると判断する処理とのうちのいずれかの処理を実行する構成とされている

ことを特徴とする請求項 1 ～ 6 のいずれかに記載の対話システム。

30

【請求項 8】

システムによる発話開始に対する要求の強さの度合いを示すシステム発話意欲度の指標値として、対話目的を達成するためのシステムの最終の次発話候補の内容データとなり得る目的データの残数および / または前記次発話準備手段による準備処理で得られた次発話候補の内容データの重要度を含むシステム状態を示す情報を記憶するシステム状態記憶手段を備え、

前記システム発話タイミング検出手段は、

前記パターン認識処理の結果として得られる尤度に対するユーザ発話権終了判定用閾値を、前記システム状態記憶手段に記憶されている前記目的データの残数および / または前記重要度で定まる前記システム発話意欲度を用いて、前記システム発話意欲度が強いときには当該ユーザ発話権終了判定用閾値が低くなり、前記システム発話意欲度が弱いときには当該ユーザ発話権終了判定用閾値が高くなるように予め定められた関数により設定する処理を実行する構成とされている

40

ことを特徴とする請求項 2 ～ 5 のいずれかに記載の対話システム。

【請求項 9】

前記次発話準備手段は、

前記音声認識処理手段によるユーザ発話の音声認識処理の結果が新たに出力された場合には、新たに出力された当該音声認識処理の結果を用いて、次発話の候補となる複数の次発話候補の内容データの少なくとも一部を入れ替えるか否かを判定し、入れ替えると判定した場合には、次発話の候補となる別の複数の次発話候補の内容データを取得または生成

50

する準備処理を実行する構成とされている

ことを特徴とする請求項 2 ～ 5 のいずれかに記載の対話システム。

【請求項 10】

前記次発話準備手段は、

新たに出力された前記音声認識処理の結果を用いて、この結果に含まれる単語のうち予め定められた重要度の高い単語を用いて、ユーザの関心のある話題を決定し、前記題材データ記憶手段に記憶された題材データまたは前記外部システムに記憶された題材データの中から、決定した話題に関連付けられて記憶されている題材データを選択し、次発話の候補となる別の複数の次発話候補の内容データを取得または生成する準備処理を実行する構成とされている

10

ことを特徴とする請求項 9 に記載の対話システム。

【請求項 11】

前記発話生成手段は、

前記音声信号取得手段により取得したユーザ発話の音声信号と、再生中のシステム発話の音声信号との衝突の発生を検出し、検出した衝突の発生情報を、ユーザ識別情報と関連付けてユーザ情報記憶手段に記憶させるとともに、ユーザ発話の終了からシステム発話の開始までの交替潜時を計測し、計測した交替潜時を、ユーザ識別情報と関連付けて前記ユーザ情報記憶手段に記憶させる処理も実行する構成とされ、

前記システム発話タイミング検出手段は、

前記ユーザ情報記憶手段に記憶されている音声対話相手のユーザとの衝突の発生情報を取得して当該ユーザとの衝突の発生頻度または累積発生回数を算出し、算出した衝突の発生頻度または累積発生回数が上方調整用閾値以上または超過の場合には、前記ユーザ発話権の維持または終了を識別するパターン認識処理の結果として得られる尤度に対して設定されているユーザ発話権終了判定用閾値を標準値または前回調整値よりも高く設定し、

20

前記ユーザ情報記憶手段に記憶されている音声対話相手のユーザについてのユーザ発話の終了からシステム発話の開始までの複数の交替潜時を取得して当該ユーザについての交替潜時の長短の傾向を示す平均値若しくはその他の指標値を算出し、算出した交替潜時の指標値が下方調整用閾値以上または超過の場合には、前記ユーザ発話権終了判定用閾値を標準値または前回調整値よりも低く設定する処理も実行する構成とされている

ことを特徴とする請求項 1 ～ 10 のいずれかに記載の対話システム。

30

【請求項 12】

前記発話生成手段は、

前記音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報を用いて発話速度を算出し、算出した発話速度を、ユーザ識別情報と関連付けて前記ユーザ情報記憶手段に記憶させる処理も実行する構成とされ、

前記システム発話タイミング検出手段は、

前記ユーザ情報記憶手段に記憶されている音声対話相手のユーザについてのユーザ発話の終了からシステム発話の開始までの複数の交替潜時を取得して当該ユーザについての交替潜時の長短の傾向を示す平均値若しくはその他の指標値を算出し、算出した交替潜時の指標値が下方調整用閾値以上または超過の場合に、前記ユーザ発話権終了判定用閾値を標準値または前回調整値よりも低く設定する処理を実行する際に、

40

前記ユーザ情報記憶手段に記憶されている音声対話相手の複数の発話速度を取得して当該ユーザの発話速度の傾向を示す平均値若しくはその他の指標値を算出し、前記下方調整用閾値を、算出した前記発話速度の指標値を用いて、前記発話速度の指標値が大きいときには当該下方調整用閾値が小さくなり、前記発話速度の指標値が小さいときには当該下方調整用閾値が大きくなるように予め定められた関数により設定する処理を実行する構成とされている

ことを特徴とする請求項 11 に記載の対話システム。

【請求項 13】

ユーザのリアルタイムの発話速度を含むユーザ状態を示す情報を記憶するユーザ状態記

50

憶手段を備え、

前記発話生成手段は、

前記音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報を用いてリアルタイムの発話速度を算出し、算出したリアルタイムの発話速度を前記ユーザ状態記憶手段に記憶させる処理も実行する構成とされ、

前記システム発話タイミング検出手段は、

前記音声信号取得手段により取得したユーザ発話の音声信号から音響特徴量を抽出し、抽出した音響特徴量および前記ユーザ状態記憶手段に記憶されているリアルタイムの発話速度を用いるか、または、これらの音響特徴量およびリアルタイムの発話速度に加え、前記音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報から抽出した言語特徴量を用いて、前記音声認識処理手段による音声認識処理の実行タイミングに依拠しない周期で、前記ユーザ発話権の維持または終了を識別するパターン認識処理を繰り返し実行し、このパターン認識処理の結果を用いて、システム発話の開始タイミングを検出する処理を実行する構成とされている

ことを特徴とする請求項 1 ～ 12 のいずれかに記載の対話システム。

【請求項 14】

請求項 1 ～ 13 のいずれかに記載の対話システムとして、コンピュータを機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ユーザとの音声対話のための処理を実行するコンピュータにより構成された対話システムおよびプログラムに係り、例えば、ニュースやコラムや歴史等の各種の話題を記載した記事データから生成したシナリオデータを用いてユーザに対して記事の内容を伝達するニュース対話システム、ユーザに対して機器の使用方法的説明や施設の案内等を行うガイダンス対話システム、選挙情勢や消費者志向等の各種のユーザの動向調査を行うアンケート対話システム、ユーザが店舗・商品・旅行先・聞きたい曲等の情報検索を行うための情報検索対話システム、ユーザが家電機器や車等の各種の機器や装置等を操作するための操作対話システム、子供や学生や新入社員等であるユーザに対して教育を行うための教育対話システム、システムがユーザ属性等の情報を特定するための情報特定対話システム等に利用できる。

【背景技術】

【0002】

一般に、音声対話システムは、人であるユーザと、コンピュータシステムである自身との間で、互いに主に音声チャネルを通じた言語情報のやりとりを行うことにより、所望のタスクを実行し、その目的（例えば、ユーザへのニュース等の記事の内容の伝達、ユーザに対するガイダンス、ユーザへのアンケート、ユーザによる情報検索、ユーザによる機器等の操作、ユーザの教育、システムによる情報特定等）を達成するものである。

【0003】

より詳細には、従来の音声対話システムでは、まず、ユーザ発話の音声信号を取得し（音声信号取得）、連続的に得られる音声信号から、ユーザの発話が途切れたことを手がかりとして発話単位の音声信号を切り出す発話区間検出を行い（発話区間検出）、次に、得られた発話区間の音声信号を言語情報に変換する音声認識処理を行うことにより、検出したユーザ発話の意味を推定し（音声認識）、続いて、推定した意味に応じて次発話を決定し、すなわち得られたユーザの言語情報に適したシステム発話の内容を生成し（発話内容生成）、さらに、その発話内容を音声信号に変換する音声合成処理を行い（音声合成）、その後、システム発話の内容をユーザに伝達するため、生成したシステム発話の音声信号を再生する処理を行う（音声信号再生）。従来の音声対話システムは、これらの一連の処理を、原則的にはシーケンシャルに行うため、それぞれの処理における遅延が蓄積することで、ユーザが発話を完了してから、システムが応答するまでに長い遅延が生じることに

10

20

30

40

50

なる。

【0004】

音声対話における二者間の発話の間(ま)の長さを交替潜時と呼ぶが、人同士の円滑な対話における交替潜時は、平均的には0.6秒程度であり、長くとも1秒程度である。また、相手の発話が終了する前に、発話を開始することも多く、これを衝突と呼ぶ。一方、近年普及しているスマートスピーカ等の対話システムと人との対話においては、ユーザの発話終了からシステムの発話開始までの間(ま)(以下、特にユーザからシステムという方向性を持たせた交替潜時を指すときは、システムの交替潜時と呼ぶ。)が、1秒から数秒となることが多い。従来の研究によれば、一方の交替潜時が他方の交替潜時に影響を与えるとされているので、システムの交替潜時が不要に長くなると、これに影響されてユーザの間(ま)(システムの発話終了からユーザの応答開始までに要する時間)も長くなる。これにより、対話全体に要する時間が不要に長くなるため、タスク達成の効率や、ユーザ体験の観点から好ましくない。

10

【0005】

従って、システムの応答性を向上させることにより、上述した従来生じていたユーザ発話とシステム発話との間に生じる不要に長い無音の時間を短くするか、あるいは発生そのものを避けることが望ましく、それを実現するためには、システム発話の開始タイミングを適切に検出することが必要となる。なぜなら、システムの交替潜時を短くするためにシステム発話の開始タイミングを不当に早めるような方法で検出処理を行えば、衝突が発生する可能性が高くなるので、単純にシステム発話の開始タイミングが早まる方法を採用すればよいというものではないからである。

20

【0006】

より詳細には、従来の音声対話システムでは、ユーザ発話の終了時をシステム発話の開始タイミングとみなしていた。1対1の対話においては、これは極めて自然な考え方であるが、そもそもユーザ発話が終了する現象の定義が明確ではなかった。例えば、特定の長さ(例えば、100ミリ秒以上)のポーズで区切られた音声区間をInter-Pausal Unit(IPU)と呼び、音声分析や会話分析では音声区間の単位として広く用いられているが、100ミリ秒程度の無音区間は、1人の話者の発話区間内にも頻繁に生じるため、必ずしもその前後で話者交替が起こるわけではない。そのため、ユーザ発話の音声信号における短い無音区間をシステム発話の開始タイミングの検出に用いると、生成して再生を開始したシステム発話と、継続されたユーザ発話とがオーバーラップする衝突を起こし、対話を崩してしまう可能性がある。一方、より長い無音区間で区切ることで、オーバーラップ(衝突)を防ぐことはできるが、システム発話の開始タイミングは、無音区間の長さだけ遅れ、ユーザ発話とシステム発話との間の無音区間を短くすることができなくなる。

30

【0007】

また、従来の音声認識では、音声認識対象とする音声区間を決定するために音声区間検出(Voice Activity Detection; VAD)と呼ばれる処理を行う。音声信号の振幅やゼロ交差数を閾値処理する単純なものから、音声信号から得られる特徴量に基づき確率的に音声が含まれるか否かを決定するモデルなど、様々な手法が研究されてきた。しかし、システム発話の開始タイミングを早期に決定するということを意図した手法は提案されていなかった。

40

【0008】

さらに、システム発話の開始タイミングを決定するために、ユーザ発話の継続または終了、あるいはシステムが次にどのような行動をとるべきか(発話だけに限らず、相槌なども含む)を検出する技術も、本願発明者らにより研究されているが、ユーザ発話途中でのシステムの相槌・復唱の生成技術を除けば、これらは全て音声認識と同様にVADを前提としており、VAD処理による遅延の影響を排除することができない。

【0009】

これらの従来技術に対し、本願発明者らは、音声信号を逐次処理し、短い周期(例えば

50

、１０ミリ秒～１００ミリ秒）で音声信号から音響特徴量を抽出し、抽出した音響特徴量を用いて、システムが発話をすべきか否かの識別を行う技術、換言すれば、ユーザが発話する地位または立場を有していることを示すユーザ発話権の維持または終了（終了には、譲渡、放棄が含まれる。）を識別する技術を開発した（非特許文献１，２参照）。このようにすることで、音声区間検出処理（ＶＡＤ処理）による遅延なしにシステム発話の開始タイミングを決定することができる。

【００１０】

なお、本発明では、複数の次発話候補が準備された場合に、その中から次発話を選択する処理が行われるが、この選択処理を行うために必要となる情報を生成する技術としては、本願発明者らにより開発された、韻律分析によりユーザ発話意図を推定する技術が知られている（非特許文献３参照）。

10

【００１１】

また、本発明は、例えば、ニュース対話システム、ガイダンス対話システム、アンケート対話システム、情報検索対話システム、操作対話システム、教育対話システム等の各種の対話システムに適用することができるが、ユーザへの効率的な情報伝達を実現することができる対話システムとしては、本願発明者らにより開発された、主計画および副計画からなるシナリオデータを用いてユーザに対してニュース等の記事の内容を伝達するニュース対話システムが知られている（非特許文献４参照）。

【先行技術文献】

【非特許文献】

20

【００１２】

【非特許文献１】藤江真也、横山勝矢、小林哲則、“音声対話システムのためのユーザ発話終了タイミングの逐次予測”、日本音響学会講演論文集、２０１８

【非特許文献２】藤江真也、横山勝矢、小林哲則、“音声対話システムのためのユーザの発話権維持状態の逐次推定”、人工知能学会全国大会、２０１８、June 2018

【非特許文献３】高津弘明、横山勝矢、本田裕、藤江真也、小林哲則、“システム発話の文脈を考慮した発話意図理解”、言語処理学会 第２５回年次大会 発表論文集、pp. 320-323、2019

【非特許文献４】高津弘明、福岡維新、藤江真也、林良彦、小林哲則、“意図性の異なる多様な情報行動を可能とする音声対話システム”、人工知能学会論文誌、vol. 22、no. 1、p. DSH-C\_1-24、2018

30

【発明の概要】

【発明が解決しようとする課題】

【００１３】

従来の音声対話システムでは、前述したように、音声信号取得、発話区間検出、音声認識、発話内容生成、音声合成、音声信号再生という一連の処理を、シーケンシャルに行うため、それぞれの処理における遅延が蓄積するという問題があった。

【００１４】

また、前述した非特許文献１，２に記載された技術を用いれば、短い周期（例えば、１０ミリ秒～１００ミリ秒）で音声信号から抽出した音響特徴量を用いてユーザ発話権の維持または終了を識別するパターン認識処理を行うので、音声区間検出処理（ＶＡＤ処理）による遅延なしにシステム発話の開始タイミングを決定することができる。

40

【００１５】

しかし、非特許文献１，２に記載された技術を用いれば、システム発話の開始タイミングを、ＶＡＤ処理による遅延なしに早期に、かつ、衝突の発生を回避または抑制しながら適切に、決定することができるものの、その後の処理、すなわち、前述した一連の処理のうちの音声認識、発話内容生成、音声合成、音声信号再生の各処理を、従来通りにシーケンシャルに行うと、そこでの遅延が生じるという問題がある。

【００１６】

従って、非特許文献１，２に記載された技術を利用してシステム発話の開始タイミング

50



を早期かつ適切に決定しつつ、ユーザ発話とシステム発話との間に生じる不要に長い無音の時間を短くするか、あるいは発生を回避することができる技術の開発が望まれる。

【 0 0 1 7 】

本発明の目的は、システムの応答性を向上させることができ、衝突の発生を回避または抑制しつつ、不要に長いシステムの交替潜時の発生を回避または抑制することができる対話システムおよびプログラムを提供するところにある。

【 課題を解決するための手段 】

【 0 0 1 8 】

< 本発明の基本構成 >

【 0 0 1 9 】

本発明は、ユーザとの音声対話のための処理を実行するコンピュータにより構成された対話システムであって、

ユーザ発話の音声信号を取得する音声信号取得手段と、

この音声信号取得手段により取得したユーザ発話の音声信号についての音声認識処理を実行する音声認識処理手段と、

音声信号取得手段により取得したユーザ発話の音声信号から音響特徴量を抽出し、抽出した音響特徴量を用いるか、または、この音響特徴量に加え、音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報から抽出した言語特徴量を用いて、音声認識処理手段による音声認識処理の実行タイミングに依拠しない周期で、ユーザが発話する地位または立場を有していることを示すユーザ発話権の維持または終了を識別するパターン認識処理を繰り返し実行し、このパターン認識処理の結果を用いて、システム発話の開始タイミングを検出する処理を実行するシステム発話タイミング検出手段と、

このシステム発話タイミング検出手段によるパターン認識処理の周期に依拠しないタイミングで、かつ、このシステム発話タイミング検出手段によりシステム発話の開始タイミングが検出される前に、題材データ記憶手段に記憶された題材データまたはネットワークを介して接続された外部システムに記憶された題材データを用いるとともに、ユーザとシステムとの間の対話履歴情報の少なくとも一部および/または音声認識処理手段による進行中のユーザ発話についての途中までの音声認識処理の結果を用いて、システムの次発話の内容データを取得または生成する準備処理を実行する次発話準備手段と、

システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、次発話準備手段による準備処理で得られた次発話の内容データを用いて、システム発話の音声信号の再生を含むシステム発話生成処理を実行する発話生成手段と

を備えたことを特徴とするものである。

【 0 0 2 0 】

ここで、ユーザ発話権の「終了」には、放棄および譲渡の双方が含まれる。放棄は、自分の発話を終了させるだけの場合であり、譲渡は、相手への質問等のように、自分の発話を終了させるとともに、相手の発話開始を促す場合である。

【 0 0 2 1 】

また、「次発話準備手段」における「題材データ記憶手段に記憶された題材データまたはネットワークを介して接続された外部システムに記憶された題材データ」には、例えば、ニュース等の各種の話題をシナリオ化したシナリオデータ、シナリオ化されていない各種の話題データ、辞書データ、事典データ、機器の使用法や施設等のガイダンス用データ、アンケート調査用データ、機器や装置等の操作補助用データ、教育用データ等が含まれる。

【 0 0 2 2 】

さらに、「次発話準備手段」による「システムの次発話の内容データを取得または生成」の「取得」には、題材データ記憶手段やネットワークを介して接続された外部システムに記憶されている複数の題材データの中からの必要な題材データ（使用するか、または使用する可能性のある題材データ）の選択的な取得と、題材データ記憶手段や外部システムに記憶されている任意の1つの題材データの構成要素の中からの必要な構成要素（使用する

10

20

30

40

50

るか、または使用する可能性のある構成要素)の選択的な取得とが含まれる。

【0023】

また、上記の「次発話準備手段」における「生成」には、取得した言語情報(題材データまたはその構成要素であるテキストデータ)の加工(語尾等の部分的な変換調整、結合等)が含まれる。但し、題材データは、題材データ記憶手段や外部システムに記憶されている段階で、予め加工されていることが好ましい。そして、「生成」には、テキストデータから音声データへの変換(音声合成)も含まれる。なお、題材データ記憶手段や外部システムに、題材データまたはその構成要素として、音声データ(例えばwavファイル等)が既に用意されている場合には、「次発話準備手段」による音声合成処理は行わなくてもよい。

10

【0024】

さらに、「次発話準備手段」により準備される「システムの次発話の内容データ」は、テキストデータおよびこれに対応する音声データの場合と、テキストデータだけの場合とがある。但し、「発話生成手段」の処理負荷の軽減および遅延防止の観点からは、「次発話準備手段」により音声データも併せて準備することが好ましい。そして、対話中に、付帯的な情報として、映像(動画)や静止画を再生する場合には、「システムの次発話の内容データ」には、映像データや画像データが付随していてもよく、対話中に音楽を再生する場合には、「システムの次発話の内容データ」には、楽曲データが含まれていてもよい。

【0025】

また、「次発話準備手段」における「ユーザとシステムとの間の対話履歴情報の少なくとも一部および/または音声認識処理手段による進行中のユーザ発話についての途中までの音声認識処理の結果を用いて」の「対話履歴情報の少なくとも一部」を用いることには、対話履歴情報(システム発話、ユーザ発話)の全体を用いること、直前のシステム発話のみを用いること、直前のシステム発話を用いずにそれよりも前のシステム発話やユーザ発話を用いること(例えば、ユーザの「さっき言ったXXXのこと、もう少し詳しく聞きたいんだけど・・・」等の要求に応答する場合等)、直前のユーザ発話のみを用いること等が含まれる。そして、「音声認識処理手段による進行中のユーザ発話についての途中までの音声認識処理の結果」を用いることは、進行中のユーザ発話の部分的な内容(ユーザ発話の発話区間全体の内容ではなく、その途中までの部分的な内容)を用いることである。なお、部分的な音声認識処理の結果が得られた場合に、それがユーザ発話の発話区間の最後の部分であるか否かは、その時点では判らないことがあるが、結果的にそれがユーザ発話の発話区間の最後の部分であった場合には、直前のユーザ発話ということになり、「対話履歴情報の少なくとも一部」に該当する。

20

【0026】

さらに、「発話生成手段」における「システム発話の音声信号の再生を含むシステム発話生成処理」には、次発話のテキストデータについての音声合成が未だ済んでいない場合には、音声合成処理が含まれる。なお、前述した通り、「発話生成手段」の処理負荷の軽減および遅延防止の観点からは、音声データ(例えばwavファイル等)は、次発話準備手段による準備処理で用意することが好ましい。さらに、「次発話準備手段」により準備された「システムの次発話の内容データ」に映像データや画像データが付随している場合には、「発話生成手段」における「システム発話生成処理」には映像や画像の再生処理も含まれ、「システムの次発話の内容データ」に楽曲データが含まれている場合には、「システム発話生成処理」には音楽の再生処理も含まれる。

30

40

【0027】

<本発明の基本構成の作用・効果>

【0028】

このような本発明の対話システムにおいては、システム発話タイミング検出手段により、ユーザが自己の発話権を維持しているか、または、譲渡若しくは放棄により終了させたかをパターン認識処理により逐次推定するとともに、次発話準備手段により、システム発

50

話タイミング検出手段によるパターン認識処理とは非同期で、かつ、システム発話タイミング検出手段によりシステム発話の開始タイミングが検出される前に、ユーザ発話に対するシステムの次発話の内容データを準備する。すなわち、ユーザ発話に対するシステムの次発話の内容データを、当該ユーザ発話の進行中に、または、それよりも前の段階である当該ユーザ発話の開始前に準備しておく。

【0029】

このため、対話相手であるユーザが自己の発話権を譲渡若しくは放棄することによりユーザ発話権が終了し、システム発話タイミング検出手段により、このユーザ発話権の終了が捉えられ、システム発話の開始タイミングが検出された場合には、その検出直後に、発話生成手段により、タイミングよくシステム発話を開始させることが可能となるので、システムの応答性を向上させることが可能となる。

10

【0030】

また、システム発話タイミング検出手段は、音声認識処理手段による音声認識処理とは非同期で、ユーザ発話権の維持または終了を識別するパターン認識処理を繰り返し実行する構成とされているので、音声区間検出処理（VAD処理）を前提としない処理を実現することができるため、VAD処理による遅延なしに早期に、システム発話の開始タイミングを決定することができるとともに、ユーザ発話とシステム発話との衝突も回避または抑制することができる。

【0031】

以上より、本発明では、次発話準備手段により、システムが発話すべき内容を早期に確定したうえで、システム発話タイミング検出手段により、ユーザ発話権の終了が推定され、システム発話の開始タイミングが検出されるのを待って、発話生成手段により、システム応答を行うので、ユーザ発話の終了後、システム発話の開始までに、長い間（ま）が空くことを避けることができるうえ、両者の発話の衝突の発生も回避または抑制することができ、これらにより前記目的が達成される。

20

【0032】

また、本発明では、次発話準備手段により、次発話の候補となる複数の次発話候補の内容データを取得または生成する準備処理を実行し、この準備処理で得られた複数の次発話候補の内容データの中から、発話生成手段で用いる次発話の内容データを選択する処理を実行するようにしてもよい。これにより、様々な種別の対話に対応可能となる。具体的には、以下のような構成を採用することができる。

30

【0033】

< 音声認識処理の結果として得られた言語情報を用いて、複数の次発話候補の内容データの中から、次発話の内容データを選択する構成 >

【0034】

すなわち、前述した対話システムにおいて、

次発話準備手段は、

次発話の候補となる複数の次発話候補の内容データを取得または生成する準備処理を実行する構成とされ、

システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、音声認識処理手段による音声認識処理の結果として得られた言語情報を用いて、次発話準備手段による準備処理で得られた複数の次発話候補の内容データの中から、発話生成手段で用いる次発話の内容データを選択する処理を実行する次発話選択手段を備えた構成を採用することができる。

40

【0035】

このように音声認識処理の結果として得られた言語情報を用いて、複数の次発話候補の内容データの中から、次発話の内容データを選択する構成とした場合には、次発話選択手段により、ユーザ発話の内容に応じて、システムの次発話の内容データを選択することができる。

【0036】

50

このため、例えば、次発話準備手段により、直前のシステム発話の内容に基づき、またはそれまでの対話履歴（システム発話、ユーザ発話）の内容に基づき、システムの次発話での使用が想定される複数の次発話候補の内容データを準備しておき、ユーザ発話の内容に応じて、準備した複数の次発話候補の内容データの中から、次発話の内容データを選択することができる。

【0037】

また、例えば、次発話準備手段により、進行中のユーザ発話の途中までの部分的な内容（ユーザ発話の開始時点から途中の時点までの内容、あるいは、ユーザ発話の途中の時点から別の途中の時点までの内容）に応じ、システムの次発話での使用が想定される複数の次発話候補の内容データを準備しておき、その後の発話内容（途中の時点以降、あるいは別の途中の時点以降の発話内容）を含めたユーザ発話の発話区間全体の内容に応じて、準備した複数の次発話候補の内容データの中から、次発話の内容データを選択することもできる。

10

【0038】

< 韻律分析で推定したユーザ発話意図を用いて、複数の次発話候補の内容データの中から、次発話の内容データを選択する構成 >

【0039】

また、前述した対話システムにおいて、

次発話準備手段は、

次発話の候補となる複数の次発話候補の内容データを取得または生成する準備処理を実行する構成とされ、

20

音声信号取得手段により取得したユーザ発話の音声信号から得られる韻律情報を用いるか、若しくは、この韻律情報に加えて、音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報を用いるか、またはこれらの韻律情報およびユーザ発話の言語情報に加えて、ユーザとシステムとの間の対話履歴情報のうちの直前のシステム発話の言語情報を用いて、質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図を識別するパターン認識処理を繰り返し実行する次発話選択用情報生成手段と、

システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、次発話選択用情報生成手段による処理で得られたユーザ発話意図の識別結果を用いて、次発話準備手段による準備処理で得られた複数の次発話候補の内容データの中から、発話生成手段で用いる次発話の内容データを選択する処理を実行する次発話選択手段と

30

を備えた構成を採用することができる。

【0040】

ここで、「質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図」における質問、応答、相槌等は、ユーザ発話意図の例示列挙であり、ここに列挙されていない「その他」のユーザ発話意図を用意してもよい。また、質問、応答、相槌等は、例示列挙であるので、これらの各々は必須ではなく、別の定義のユーザ発話意図を用意してもよい。他の発明においても同様である。

【0041】

40

このように韻律分析で推定したユーザ発話意図を用いて、複数の次発話候補の内容データの中から、次発話の内容データを選択する構成とした場合には、前述した< 音声認識処理の結果として得られた言語情報を用いて、複数の次発話候補の内容データの中から、次発話の内容データを選択する構成 > の場合と同様な作用・効果が得られることに加え、ユーザ発話意図を用いるので、音声認識処理の結果を得ることなく、次発話の内容データを選択することが可能となるため、システムの応答性を向上させることが可能となる。

【0042】

なお、次発話選択用情報生成手段のユーザ発話意図の識別器と、システム発話タイミング検出手段のユーザ発話権の維持・終了の識別器とは、マルチタスクの識別器とすることにより一体化させてもよい。

50

## 【 0 0 4 3 】

また、次発話選択用情報生成手段で用いる韻律情報を得るための分析処理は、システム発話タイミング検出手段で用いる音響特徴量を抽出するための分析処理と共通の処理としてもよい。

## 【 0 0 4 4 】

＜韻律分析で推定したユーザ発話意図と、音声認識処理の結果として得られた言語情報とを組み合わせ用いて、複数の次発話候補の内容データの中から、次発話の内容データを選択する構成＞

## 【 0 0 4 5 】

さらに、前述した対話システムにおいて、

次発話準備手段は、

次発話の候補となる複数の次発話候補の内容データを取得または生成する準備処理を実行する構成とされ、

音声信号取得手段により取得したユーザ発話の音声信号から得られる韻律情報を用いるか、若しくは、この韻律情報に加えて、音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報を用いるか、またはこれらの韻律情報およびユーザ発話の言語情報に加えて、ユーザとシステムとの間の対話履歴情報のうちの直前のシステム発話の言語情報を用いて、質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図を識別するパターン認識処理を繰り返し実行する次発話選択用情報生成手段と、

システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、次発話選択用情報生成手段による処理で得られたユーザ発話意図の識別結果と、音声認識処理手段による音声認識処理の結果として得られた言語情報とを組み合わせ用いて、次発話準備手段による準備処理で得られた複数の次発話候補の内容データの中から、発話生成手段で用いる次発話の内容データを選択する処理を実行する次発話選択手段と

を備えた構成を採用することができる。

## 【 0 0 4 6 】

このように韻律分析で推定したユーザ発話意図と、音声認識処理の結果として得られた言語情報とを組み合わせ用いて、複数の次発話候補の内容データの中から、次発話の内容データを選擇する構成とした場合には、ユーザ発話意図を用いるだけでは、対応できないときでも、あるいは、音声認識処理の結果を用いるだけでは、対応できないときでも、次発話の内容データの選擇処理を行うことができるようになるので、あらゆるタイプの音声対話に対応可能となる。

## 【 0 0 4 7 】

＜システム発話タイミング検出手段によりユーザ発話意図の識別も行う構成＞

## 【 0 0 4 8 】

また、前述した対話システムにおいて、

次発話準備手段は、

次発話の候補となる複数の次発話候補の内容データを取得または生成する準備処理を実行する構成とされ、

システム発話タイミング検出手段は、

ユーザ発話権の維持または終了を識別するパターン認識処理を実行する際に、終了については、質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図のうちのいずれのユーザ発話意図で終了するのかを識別するパターン認識処理を実行する構成とされ、

システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された後に、システム発話タイミング検出手段による処理で得られたユーザ発話意図の識別結果を用いて、次発話準備手段による準備処理で得られた複数の次発話候補の内容データの中から、発話生成手段で用いる次発話の内容データを選択する処理を実行する次発話選択手段を備えた構成を採用することができる。

## 【 0 0 4 9 】

このようにシステム発話タイミング検出手段によりユーザ発話意図の識別も行う構成とした場合には、システム発話タイミング検出手段によりシステム発話の開始タイミングが検出された時点で、同時にユーザ発話意図の識別結果も得られているので、システムの応答性を向上させることが可能となる。

## 【 0 0 5 0 】

<システム発話タイミング検出手段により、システム状態を示す情報である準備完了・準備中の別を参照する構成>

## 【 0 0 5 1 】

また、以上に述べた対話システムにおいて、

次発話準備手段による準備処理の状態を含むシステム状態を示す情報を記憶するシステム状態記憶手段を備え、

システム発話タイミング検出手段は、

ユーザ発話権の維持または終了を識別するパターン認識処理の結果およびシステム状態記憶手段に記憶されているシステム状態を示す情報を用いて、システム発話の開始タイミングを検出する処理を実行する際に、

パターン認識処理の結果がユーザ発話権の維持を示している場合には、システム発話の開始タイミングではないと判断し、

パターン認識処理の結果がユーザ発話権の終了を示し、かつ、システム状態を示す情報が準備完了を示している場合には、システム発話の開始タイミングであると判断し、

パターン認識処理の結果がユーザ発話権の終了を示し、かつ、システム状態を示す情報が準備中を示している場合には、次発話準備手段による準備中の処理内容に応じ、直ぐに完了する処理内容として予め分類されている処理の準備中であるときには、準備完了になるまで待つてシステム発話の開始タイミングであると判断し、直ぐに完了しない処理内容として予め分類されている処理の準備中であるときには、システム発話の開始タイミングであると判断するとともに、フィラーの挿入タイミングである旨の情報を出力する処理を実行する構成としてもよい。

## 【 0 0 5 2 】

このようにシステム発話タイミング検出手段により、システム状態を示す情報である準備完了・準備中の別を参照する構成とした場合には、システム状態を考慮し、より適切なシステム発話の開始タイミングを検出することが可能となる。

## 【 0 0 5 3 】

<システム発話タイミング検出手段により、ユーザ状態を示す情報であるユーザ発話継続時間を用いて、ユーザ発話権終了判定用閾値の調整を行うか、またはシステム発話の開始タイミングであるか否かの判断を行う構成>

## 【 0 0 5 4 】

さらに、以上に述べた対話システムにおいて、

ユーザ発話継続時間を含むユーザ状態を示す情報を記憶するユーザ状態記憶手段を備え、

システム発話タイミング検出手段は、

ユーザ発話権の維持または終了を識別するパターン認識処理の結果およびユーザ状態記憶手段に記憶されているユーザ状態を示す情報を用いて、システム発話の開始タイミングを検出する処理を実行し、この際の処理として、

( 1 ) ユーザ状態記憶手段に記憶されているユーザ発話継続時間が、予め定められた短時間判定用閾値以下または未満の場合には、パターン認識処理の結果として得られる尤度に対して設定されているユーザ発話権終了判定用閾値を標準値よりも高く設定し、予め定められた長時間判定用閾値以上または超過の場合には、ユーザ発話権終了判定用閾値を標準値よりも低く設定する処理と、

( 2 ) ユーザ状態記憶手段に記憶されているユーザ発話継続時間を用いて、パターン認識処理の結果として得られる尤度に対するユーザ発話権終了判定用閾値を、ユーザ発話継

10

20

30

40

50

続時間が短いときには当該ユーザ発話権終了判定用閾値が高くなり、ユーザ発話継続時間が長いときには当該ユーザ発話権終了判定用閾値が低くなるように予め定められた関数により設定する処理と、

(3) ユーザ状態記憶手段に記憶されているユーザ発話継続時間が、予め定められた短時間判定用閾値以下または未満の場合には、パターン認識処理の結果がユーザ発話権の終了を示していても、システム発話の開始タイミングではないと判断し、予め定められた長時間判定用閾値以上または超過の場合には、パターン認識処理の結果がユーザ発話権の維持を示していても、システム発話の開始タイミングであると判断する処理とのうちのいずれかの処理を実行する構成としてもよい。

【0055】

ここで、「ユーザ発話権終了判定用閾値」の「標準値」は、別の情報に基づく別の趣旨での閾値調整が別途に行われている場合には、その別途の閾値調整後の値を指す。

【0056】

このようにシステム発話タイミング検出手段により、ユーザ状態を示す情報であるユーザ発話継続時間を用いて、ユーザ発話権終了判定用閾値の調整を行う(上記(1)、(2))か、またはシステム発話の開始タイミングであるか否かの判断を行う(上記(3))構成とした場合には、ユーザ発話継続時間の長短に応じ、システム発話の開始タイミングを調整することが可能となる。

【0057】

上記(1)、(2)では、ユーザ発話継続時間が短いときにはユーザ発話権終了判定用閾値が高くなり、ユーザ発話継続時間が長いときにはユーザ発話権終了判定用閾値が低くなるように設定することができるので、ユーザ発話の開始直後の時期には、ユーザ発話権が終了したという識別結果が出にくい設定状態とし、ユーザ発話の開始時点から比較的長時間が経過している時期には、ユーザ発話権が終了したという識別結果が出やすい設定状態とすることができる。

【0058】

上記(3)では、ユーザ発話継続時間を、ユーザ発話権終了判定用閾値に反映させるのではなく、ユーザ発話権終了判定用閾値を用いて維持・終了の識別結果を出した後におけるシステム発話の開始タイミングの判断処理に反映させることにより、上記(1)、(2)と同様な作用・効果を得る。

【0059】

<システム発話タイミング検出手段により、システム状態を示す情報であるシステム発話意欲度を用いてユーザ発話権終了判定用閾値を動的に調整する構成>

【0060】

また、前述した次発話準備手段により準備した複数の次発話候補の内容データの中から、次発話選択手段により次発話の内容データを選択する構成とした場合において、

システムによる発話開始に対する要求の強さの度合いを示すシステム発話意欲度の指標値として、対話目的を達成するためのシステムの最終の次発話候補の内容データとなり得る目的データの残数および/または次発話準備手段による準備処理で得られた次発話候補の内容データの重要度を含むシステム状態を示す情報を記憶するシステム状態記憶手段を備え、

システム発話タイミング検出手段は、

パターン認識処理の結果として得られる尤度に対するユーザ発話権終了判定用閾値を、システム状態記憶手段に記憶されている目的データの残数および/または重要度で定まるシステム発話意欲度を用いて、システム発話意欲度が強いときには当該ユーザ発話権終了判定用閾値が低くなり、システム発話意欲度が弱いときには当該ユーザ発話権終了判定用閾値が高くなるように予め定められた関数により設定する処理を実行する構成としてもよい。

【0061】

このようにシステム発話タイミング検出手段により、システム状態を示す情報であるシ

10

20

30

40

50

システム発話意欲度を用いてユーザ発話権終了判定用閾値を動的に調整する構成とした場合には、システム発話意欲度が強いときには、ユーザ発話権が終了したという識別結果が出やすくなる設定状態とし、システム発話意欲度が弱いときには、ユーザ発話権が終了したという識別結果が出にくい設定状態とすることが可能となる。

【0062】

< 音声認識処理の結果が新たに出力されたときに、その音声認識処理の結果を用いて、次発話候補の入替が可能な構成 >

【0063】

さらに、前述した次発話準備手段により準備した複数の次発話候補の内容データの中から、次発話選択手段により次発話の内容データを選択する構成とした場合において、

10

次発話準備手段は、

音声認識処理手段によるユーザ発話の音声認識処理の結果が新たに出力された場合には、新たに出力された当該音声認識処理の結果を用いて、次発話の候補となる複数の次発話候補の内容データの少なくとも一部を入れ替えるか否かを判定し、入れ替えると判定した場合には、次発話の候補となる別の複数の次発話候補の内容データを取得または生成する準備処理を実行する構成としてもよい。

【0064】

このように音声認識処理の結果が新たに出力されたときに、その音声認識処理の結果を用いて、次発話候補の入替が可能な構成とした場合には、進行中のユーザ発話の内容に応じて、既に準備されている複数の次発話候補の内容データの入替を行うことが可能となるので、ユーザ発話の内容に応じた適切な次発話候補の内容データを準備することが可能となる。

20

【0065】

< 音声認識処理の結果が新たに出力されたときに、この結果に含まれる重要度の高い単語を用いてユーザの関心のある話題を決定し、決定した話題に従って次発話候補の入替を行う構成 >

【0066】

また、上述した音声認識処理の結果が新たに出力されたときに、その音声認識処理の結果を用いて、次発話候補の入替が可能な構成とした場合において、

30

次発話準備手段は、

新たに出力された音声認識処理の結果を用いて、この結果に含まれる単語のうち予め定められた重要度の高い単語を用いて、ユーザの関心のある話題を決定し、題材データ記憶手段に記憶された題材データまたは外部システムに記憶された題材データの中から、決定した話題に関連付けられて記憶されている題材データを選択し、次発話の候補となる別の複数の次発話候補の内容データを取得または生成する準備処理を実行する構成としてもよい。

【0067】

このように音声認識処理の結果が新たに出力されたときに、この結果に含まれる重要度の高い単語を用いてユーザの関心のある話題を決定し、決定した話題に従って次発話候補の入替を行う構成とした場合には、進行中のユーザ発話の内容に応じて、既に準備されている複数の次発話候補の内容データの入替を行い、次発話により提示する話題を変更することが可能となる。

40

【0068】

< システム発話タイミング検出手段により、衝突の発生情報やシステムの交替潜時を用いて、ユーザ発話権終了判定用閾値を調整する構成 >

【0069】

また、以上に述べた対話システムにおいて、

発話生成手段は、

音声信号取得手段により取得したユーザ発話の音声信号と、再生中のシステム発話の音声信号との衝突の発生を検出し、検出した衝突の発生情報を、ユーザ識別情報と関連付け

50



てユーザ情報記憶手段に記憶させるとともに、ユーザ発話の終了からシステム発話の開始までの交替潜時を計測し、計測した交替潜時を、ユーザ識別情報と関連付けてユーザ情報記憶手段に記憶させる処理も実行する構成とされ、

システム発話タイミング検出手段は、

ユーザ情報記憶手段に記憶されている音声対話相手のユーザとの衝突の発生情報を取得して当該ユーザとの衝突の発生頻度または累積発生回数を算出し、算出した衝突の発生頻度または累積発生回数が上方調整用閾値以上または超過の場合には、ユーザ発話権の維持または終了を識別するパターン認識処理の結果として得られる尤度に対して設定されているユーザ発話権終了判定用閾値を標準値または前回調整値よりも高く設定し、

ユーザ情報記憶手段に記憶されている音声対話相手のユーザについてのユーザ発話の終了からシステム発話の開始までの複数の交替潜時を取得して当該ユーザについての交替潜時の長短の傾向を示す平均値若しくはその他の指標値を算出し、算出した交替潜時の指標値が下方調整用閾値以上または超過の場合には、ユーザ発話権終了判定用閾値を標準値または前回調整値よりも低く設定する処理も実行する構成としてもよい。

【0070】

ここで、「システム発話タイミング検出手段」における「標準値または前回調整値」は、別の情報に基づく別の趣旨での閾値調整が別途に行われている場合には、その別途の閾値調整後の値を指す。

【0071】

また、ここでの「衝突」は、ユーザ発話権が終了したという識別結果が出て、システム発話を開始したところ、実際にはユーザ発話権が維持されていて、両者の発話が重なった場合の衝突である。従って、ユーザ発話権が終了したものの、システム発話の開始が遅れたために、再び、ユーザ発話が開始されてしまい、ほぼ同時に両者の発話が開始されて重なった場合の衝突ではない。

【0072】

さらに、ここでの「交替潜時」は、ユーザ発話の終了からシステム発話の開始までの間（ま）であり、システムの交替潜時である。従って、「当該ユーザについての交替潜時」とされているが、これは、当該ユーザとの音声対話を行うときのシステムの交替潜時のことであり、システム発話の終了からユーザ発話の開始までの間（ま）のことではない。

【0073】

このようにシステム発話タイミング検出手段により、衝突の発生情報やシステムの交替潜時を用いて、ユーザ発話権終了判定用閾値を調整する構成とした場合には、各ユーザについて、衝突の発生が起きる傾向にあるときには、ユーザ発話権が終了したという識別結果が出にくい設定状態とし、システムの交替潜時が長い傾向にあるときには、ユーザ発話権が終了したという識別結果が出やすくなる設定状態とすることが可能となる。このため、ユーザ属性に応じたユーザ発話権終了判定用閾値の調整を実現することができる。

【0074】

<ユーザ発話権終了判定用閾値を下方調整することを決めるための下方調整用閾値を、ユーザの発話速度の関数とする構成>

【0075】

さらに、上述したシステム発話タイミング検出手段により、衝突の発生情報やシステムの交替潜時を用いて、ユーザ発話権終了判定用閾値を調整する構成とした場合において、発話生成手段は、

音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報を用いて発話速度を算出し、算出した発話速度を、ユーザ識別情報と関連付けてユーザ情報記憶手段に記憶させる処理も実行する構成とされ、

システム発話タイミング検出手段は、

ユーザ情報記憶手段に記憶されている音声対話相手のユーザについてのユーザ発話の終了からシステム発話の開始までの複数の交替潜時を取得して当該ユーザについての交替潜時の長短の傾向を示す平均値若しくはその他の指標値を算出し、算出した交替潜時の指標

10

20

30

40

50

値が下方調整用閾値以上または超過の場合に、ユーザ発話権終了判定用閾値を標準値または前回調整値よりも低く設定する処理を実行する際に、

ユーザ情報記憶手段に記憶されている音声対話相手の複数の発話速度を取得して当該ユーザの発話速度の傾向を示す平均値若しくはその他の指標値を算出し、下方調整用閾値を、算出した発話速度の指標値を用いて、発話速度の指標値が大きいときには当該下方調整用閾値が小さくなり、発話速度の指標値が小さいときには当該下方調整用閾値が大きくなるように予め定められた関数により設定する処理を実行する構成としてもよい。

【0076】

このようにユーザ発話権終了判定用閾値を下方調整することを決めるための下方調整用閾値を、ユーザの発話速度の関数とする構成とした場合には、各ユーザの発話速度の傾向に応じ、下方調整用閾値の設定を変更することが可能となる。このため、ユーザ属性に応じたユーザ発話権終了判定用閾値の下方調整を実現することができる。すなわち、システムの交替潜時が長い傾向にあるときには、ユーザ発話権終了判定用閾値を下方調整することにより、ユーザ発話権が終了したという識別結果が出やすくなる設定状態とし、システムの交替潜時が短くなるようにすることができるが、この際、システムの交替潜時が長い傾向にあるか否かは、ユーザ毎に異なり、各ユーザの発話速度の傾向と関係するので、下方調整用閾値をユーザの発話速度の関数とすることで、ユーザ属性に応じてユーザ発話権終了判定用閾値の下方調整を行うか否かを定めることができる。

【0077】

<ユーザ発話の音声信号から抽出した音響特徴量、およびリアルタイムのユーザの発話速度を用いて、ユーザ発話権の維持または終了を識別するパターン認識処理を行う構成>

【0078】

また、以上に述べた対話システムにおいて、

ユーザのリアルタイムの発話速度を含むユーザ状態を示す情報を記憶するユーザ状態記憶手段を備え、

発話生成手段は、

音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報を用いてリアルタイムの発話速度を算出し、算出したリアルタイムの発話速度をユーザ状態記憶手段に記憶させる処理も実行する構成とされ、

システム発話タイミング検出手段は、

音声信号取得手段により取得したユーザ発話の音声信号から音響特徴量を抽出し、抽出した音響特徴量およびユーザ状態記憶手段に記憶されているリアルタイムの発話速度を用いるか、または、これらの音響特徴量およびリアルタイムの発話速度に加え、音声認識処理手段による音声認識処理の結果として得られたユーザ発話の言語情報から抽出した言語特徴量を用いて、音声認識処理手段による音声認識処理の実行タイミングに依拠しない周期で、ユーザ発話権の維持または終了を識別するパターン認識処理を繰り返し実行し、このパターン認識処理の結果を用いて、システム発話の開始タイミングを検出する処理を実行する構成としてもよい。

【0079】

ここで、「リアルタイムの発話速度」における「リアルタイムの」という意味は、事後的に計算するのではなく、その場で計算されるという意味であり、「逐次得られる最新の」という意味である。発話速度の計算には、音声認識処理の結果が用いられるが、音声認識処理の結果は、若干の時間遅れで得られるため、厳密に言えば、ここでいう「リアルタイム」には、「略リアルタイム」が含まれる。

【0080】

このようにユーザ発話の音声信号から抽出した音響特徴量、およびリアルタイムのユーザの発話速度を用いて、ユーザ発話権の維持または終了を識別するパターン認識処理を行う構成とした場合には、システム発話タイミング検出手段の識別器は、ユーザ発話の各時点における発話速度（蓄積した発話速度から得られるユーザ属性としての発話速度の傾向ではなく、瞬間的な発話速度という意味）を用いた学習を行うことにより構築されるので

、その時々ユーザの発話速度を加味した識別結果を得ることが可能となる。このため、ユーザ毎に異なる発話速度の傾向（蓄積した発話速度から得られるユーザ属性）に応じてユーザ発話権終了判定用閾値を調整する必要がなくなる。なお、閾値調整と併用してもよく、その場合には、閾値調整量が少なくなる。

【0081】

＜プログラムの発明＞

【0082】

本発明のプログラムは、以上に述べた対話システムとして、コンピュータを機能させるためのものである。

【0083】

なお、上記のプログラムまたはその一部は、例えば、光磁気ディスク（MO）、コンパクトディスク（CD）、デジタル・バーサタイル・ディスク（DVD）、フレキシブルディスク（FD）、磁気テープ、読出し専用メモリ（ROM）、電氣的消去および書換可能な読出し専用メモリ（EEPROM）、フラッシュ・メモリ、ランダム・アクセス・メモリ（RAM）、ハードディスクドライブ（HDD）、ソリッドステートドライブ（SSD）、フラッシュディスク等の記録媒体に記録して保存や流通等させることが可能であるとともに、例えば、ローカル・エリア・ネットワーク（LAN）、メトロポリタン・エリア・ネットワーク（MAN）、ワイド・エリア・ネットワーク（WAN）、インターネット、イントラネット、エクストラネット等の有線ネットワーク、あるいは無線通信ネットワーク、さらにはこれらの組合せ等の伝送媒体を用いて伝送することが可能であり、また、搬送波に載せて搬送することも可能である。さらに、上記のプログラムは、他のプログラムの一部分であってもよく、あるいは別個のプログラムと共に記録媒体に記録されていてもよい。

【発明の効果】

【0084】

以上に述べたように本発明によれば、システム発話タイミング検出手段により、ユーザが自己の発話権を維持しているか、または、譲渡若しくは放棄により終了させたかをパターン認識処理により逐次推定するとともに、次発話準備手段により、システム発話タイミング検出手段によるパターン認識処理とは非同期で、かつ、システム発話タイミング検出手段によりシステム発話の開始タイミングが検出される前に、ユーザ発話に対するシステムの次発話の内容データを準備するので、システムの応答性を向上させることができ、衝突の発生を回避または抑制しつつ、不要に長いシステムの交替潜時の発生を回避または抑制することができるという効果がある。

【図面の簡単な説明】

【0085】

【図1】本発明の一実施形態の対話システムの全体構成図。

【図2】前記実施形態のシステム発話タイミング検出手段の詳細構成図。

【図3】前記実施形態の次発話選択用情報生成手段の詳細構成図。

【図4】前記実施形態の次発話準備手段の詳細構成図。

【図5】前記実施形態のユーザからシステムへの話者交替時の処理の流れを示すフローチャートの図。

【図6】前記実施形態のシステム発話、ユーザ発話、各処理の時間的な前後関係を示す説明図。

【図7】前記実施形態の各処理のタイミングとその結果との関係を示す説明図。

【図8】前記実施形態のシステム発話タイミング検出手段によるシステム発話の開始タイミングの判断処理のロジックを示すブロック図。

【図9】前記実施形態のシステム発話タイミング検出手段によるユーザ発話権終了判定用閾値のリアルタイム調整（その1）の説明図。

【図10】前記実施形態のシステム発話タイミング検出手段によるユーザ発話権終了判定用閾値のリアルタイム調整（その2）の説明図。

【図 1 1】前記実施形態のシステム発話タイミング検出手段によるユーザ発話権終了判定用閾値の事前調整（その 1）の説明図。

【図 1 2】前記実施形態のシステム発話タイミング検出手段によるユーザ発話権終了判定用閾値の事前調整（その 2）の説明図。

【図 1 3】前記実施形態のシナリオのデータ構成の具体例を示す図。

【図 1 4】前記実施形態のシナリオ再生（システム発話）とユーザの反応（ユーザ発話）との関係を示す説明図。

【図 1 5】前記実施形態の次発話候補の準備処理の具体例（1）を示す図。

【図 1 6】前記実施形態の次発話候補の準備処理の具体例（2）を示す図。

【図 1 7】前記実施形態の次発話候補の準備処理の具体例（3）を示す図。

10

【発明を実施するための形態】

【0086】

以下に本発明の一実施形態について図面を参照して説明する。図 1 には、本実施形態の対話システム 10 の全体構成が示されている。図 2 には、システム発話タイミング検出手段 22 の詳細構成が示され、図 3 には、次発話選択用情報生成手段 23 の詳細構成が示され、図 4 には、次発話準備手段 43 の詳細構成が示されている。また、図 5 には、ユーザからシステムへの話者交替時の処理の流れがフローチャートで示され、図 6 には、システム発話、ユーザ発話、各処理の時間的な前後関係が示され、図 7 には、各処理のタイミングとその結果との関係が示され、図 8 には、システム発話タイミング検出手段 22 によるシステム発話の開始タイミングの判断処理のロジックが示されている。さらに、図 9 ~ 図 12 は、システム発話タイミング検出手段 22 によるユーザ発話権終了判定用閾値の調整の説明図である。また、図 13 には、シナリオのデータ構成の具体例が示され、図 14 には、シナリオ再生（システム発話）とユーザの反応（ユーザ発話）との関係の具体例が示され、図 15 ~ 図 17 には、次発話候補の準備処理の具体例が示されている。

20

【0087】

< 対話システム 10 の全体構成 >

【0088】

図 1 において、対話システム 10 は、ユーザとの音声対話を行うシステムであり、1 台または複数台のコンピュータにより構成され、本実施形態では、一例として、再生装置 20 と、対話サーバ 40 とをネットワーク 1 で接続した構成とされている。また、ネットワーク 1 には、外部システムである題材データ提供システム 60 も接続されている。

30

【0089】

ここで、ネットワーク 1 は、主としてインターネットのような外部ネットワークであるが、これとイントラネットや LAN 等の内部ネットワークとの組合せ等でもよく、有線であるか、無線であるか、有線・無線の混在型であるかは問わない。また、ネットワーク 1 は、例えば、社内、工場内、事業所内、グループ企業内、学内、病院内、マンション内、建物内、公園・遊園地・動物園・博物館・美術館・博覧会場等の施設内、所定の地域内等に限定されたイントラネットや LAN 等の内部ネットワークであってもよい。

【0090】

再生装置 20 は、1 台または複数台のコンピュータにより構成され、音声信号取得手段 21 と、システム発話タイミング検出手段 22 と、次発話選択用情報生成手段 23 と、次発話選択手段 24 と、発話生成手段 25 と、次発話候補記憶手段 30 と、システム状態記憶手段 31 と、ユーザ状態記憶手段 32 とを備えている。この再生装置 20 は、例えば、スマートフォン、タブレット、モバイル PC（パーソナル・コンピュータ）等の携帯機器であってもよい。また、汎用機器ではなく、音声対話の専用機器としてもよい。

40

【0091】

このうち、音声信号取得手段 21（但し、マイクロフォンの部分を除く。）、システム発話タイミング検出手段 22（但し、ユーザ発話権終了判定モデル記憶手段 22E（図 2 参照）の部分を除く。）、次発話選択用情報生成手段 23（但し、第 1、第 2 の発話意図識別モデル記憶手段 23D、23G（図 3 参照）の部分を除く。）、次発話選択手段 24

50

、および発話生成手段 25（但し、スピーカやディスプレイの部分を除く。）は、再生装置 20 を構成するコンピュータ本体の内部に設けられた中央演算処理装置（CPU）、およびこの CPU の動作手順を規定する 1 つまたは複数のプログラムにより実現される。また、次発話候補記憶手段 30、システム状態記憶手段 31、ユーザ状態記憶手段 32、システム発話タイミング検出手段 22 を構成するユーザ発話権終了判定モデル記憶手段 22E（図 2 参照）、および次発話選択用情報生成手段 23 を構成する第 1、第 2 の発話意図識別モデル記憶手段 23D、23G（図 3 参照）としては、例えば、ハードディスクドライブ（HDD）、ソリッドステートドライブ（SSD）等を採用することができる。

#### 【0092】

対話サーバ 40 は、1 台または複数台のコンピュータにより構成され、音声認識処理手段 41 と、対話状態管理手段 42 と、次発話準備手段 43 と、対話履歴記憶手段 50 と、題材データ記憶手段 51 と、ユーザ情報記憶手段 52 とを備えている。

#### 【0093】

このうち、音声認識処理手段 41、対話状態管理手段 42、および次発話準備手段 43（但し、先行次発話候補情報記憶手段 43D（図 4 参照）の部分は除く。）は、対話サーバ 40 を構成するコンピュータ本体の内部に設けられた中央演算処理装置（CPU）、およびこの CPU の動作手順を規定する 1 つまたは複数のプログラムにより実現される。また、対話履歴記憶手段 50、題材データ記憶手段 51、ユーザ情報記憶手段 52、および次発話準備手段 43 を構成する先行次発話候補情報記憶手段 43D としては、例えば、ハードディスクドライブ（HDD）、ソリッドステートドライブ（SSD）等を採用することができる。なお、先行次発話候補情報記憶手段 43D（図 4 参照）は、主メモリ等の揮発性メモリとしてもよい。

#### 【0094】

題材データ提供システム 60 は、外部システムであり、1 台または複数台のコンピュータにより構成され、対話サーバ 40 を構成する題材データ記憶手段 51 に相当する外部題材データ記憶手段（不図示）を備えている。

#### 【0095】

なお、本実施形態では、図 1 に示すように、対話システム 10 は、再生装置 20 と、対話サーバ 40 とをネットワーク 1 で接続した構成とされているが、本発明の対話システムは、スタンドアローンのシステムとしてもよい。また、図 1 に示したネットワーク構成は、一例に過ぎないので、ネットワーク構成とする場合でも、各機能の分散形態として、図 1 の状態とは異なる様々な形態を採用することができる。

#### 【0096】

例えば、再生装置 20 は、音声対話相手であるユーザと音声によるやりとりを行うので、ユーザの近く（音声の届く範囲）に配置する必要があることから、これを本体と端末とに分割して無線または有線で通信を行うようにし、端末をユーザの近くに配置する一方、本体をユーザから比較的離れた位置（音声が届かない位置でもよい）に配置する構成とすることができる。この場合、例えば、端末は、再生装置 20 を構成する音声信号取得手段 21 またはその一部であるマイクロフォンの部分と、再生装置 20 を構成する発話生成手段 25 またはその一部であるスピーカの部分（映像や静止画像の再生を伴う場合には、ディスプレイの部分を含む。）とにより構成することができる。そして、例えば、本体を固定設置された機器とし、端末を移動機器とすること等ができるが、本体、端末のいずれについても、固定機器でも移動機器でもよい。また、本体と端末との個数の関係は、1 対 1 でも、1 対多でもよい。さらに、例えば、対話種別（ニュース対話、ガイダンス対話、アンケート対話、情報検索対話、操作対話、教育対話等の別）に対応させて異なる本体を設置する場合、新旧異なるタイプの本体を併用する場合、機能の異なる本体を使い分ける場合等には、本体と端末との個数の関係は、多対 1、多対多でもよく、この場合には、任意の 1 つの端末と、複数の本体から目的に応じて選択された 1 つの本体とが接続されることになる。また、ユーザとの関係では、1 つの端末は、同時使用でなければ、複数のユーザが交代して使用することができる。本体は、複数の端末と同時接続可能な構成とすれば、

10

20

30

40

50

複数のユーザの同時使用に対応可能な構成とすることができるが、複数のユーザの同時使用を許容しない構成としてもよい。

【0097】

また、再生装置20を構成するシステム発話タイミング検出手段22と、次発話選択用情報生成手段23と、次発話選択手段24と、次発話候補記憶手段30と、システム状態記憶手段31と、ユーザ状態記憶手段32とは、それぞれ別々のコンピュータに設けてもよく、適宜組み合わせで同じコンピュータに設けてもよい。

【0098】

さらに、対話サーバ40も同様であり、対話サーバ40を構成する各機能の部分は、それぞれ別々のコンピュータに設けてもよく、適宜組み合わせで同じコンピュータに設けてもよい。また、再生装置20を構成する1つまたは複数の機能の部分と、対話サーバ40を構成する1つまたは複数の機能の部分とを適宜組み合わせで同じコンピュータに設けてもよい。

【0099】

<再生装置20 / 音声信号取得手段21の構成>

【0100】

音声信号取得手段21は、ユーザ発話の音声信号を取得するものであり、音（ここでは、音声）をアナログの電気信号に変換する機器であるマイクロフォン、A/D変換手段、A/D変換で得られたデジタルの音声信号を保持する音声信号記憶手段、音声信号を各所に送信する送信手段等を含んで構成されている。

【0101】

<再生装置20 / システム発話タイミング検出手段22の構成>

【0102】

図2において、システム発話タイミング検出手段22は、音響特徴量抽出手段22Aと、言語特徴量抽出手段22Bと、ユーザ発話権終了判定用パターン認識器22Cと、システム発話開始タイミング判断手段22Fと、ユーザ発話権終了判定用閾値調整手段22Gとを含んで構成されている。このうち、音響特徴量抽出手段22Aと、言語特徴量抽出手段22Bと、ユーザ発話権終了判定用パターン認識器22Cについては、例えば、前述した非特許文献1, 2に記載された技術を採用することができる。

【0103】

このシステム発話タイミング検出手段22による処理は、音声認識処理手段41による音声認識処理の実行タイミングに依拠しない周期で、すなわち音声認識処理とは非同期で実行される。具体的には、例えば、10ms（ミリ秒）～100msという短い周期で実行される。図6の最下部に示した処理の周期Q（時間間隔）である。なお、音声区間（IPU）を形成する際のポーズは、通常は100ms以上であるから、周期Qは、そのIPU形成用の閾値よりも短いか、同等の周期ということになる。

【0104】

<再生装置20 / システム発話タイミング検出手段22 / 音響特徴量抽出手段22Aの構成>

【0105】

音響特徴量抽出手段22Aは、音声信号取得手段21により取得したユーザ発話の音声信号から音響特徴量を抽出する処理を実行するものである。前述した非特許文献1, 2に記載された技術を採用する場合には、狭帯域スペクトログラムを符号化、復号化する自己符号化器（オートエンコーダ）をニューラルネットワーク（CNN）により構築し、その中間層の出力を音響特徴量（ボトルネック特徴量）とする。具体的には、周波数分析により例えば10ms毎に得られる256点のパワースペクトルを10本分並べたものを入力とし、中間層の出力256次元を特徴量とする。すなわち、CNNオートエンコーダの入力は、例えば、フレームサイズ＝800サンプル（50ms）、フレームシフト＝160サンプル（10ms）で切り出した音声信号から生成したスペクトログラムを時系列に10本分（図6の下部に示したR本分）並べたものとし、そのサイズを10×256次元と

10

20

30

40

50

する。そして、この入力サイズを256次元に圧縮し、音響特徴量とする。

【0106】

このようにして音響特徴量を抽出する場合、256点のパワースペクトルを10本分並べた入力データを作成する際に、256点のパワースペクトルを1本ずつずらしていけば、図6の最下部に示した処理の周期Q（時間間隔）は、周波数分析のフレームシフト＝10msと同じになり、2本ずつずらしていけば、2倍の20msとなり、同様に10本ずつずらしていけば、10倍の100msとなる。従って、256点のパワースペクトルを用いる際に、ずらす本数を変えることにより、処理の周期Qを変更することができる。なお、ずらす本数を多くすることにより、使用する音声信号の区間に重なりがないようにしてもよいが、使用されない音声信号の区間が生じることは避ける必要がある。なお、図7

10

【0107】

また、音声信号からの音響特徴量の抽出処理は、上述した非特許文献1，2に記載された技術による抽出処理に限定されるものではなく、ユーザ発話権終了判定用パターン認識器22Cの入力に用いる音響特徴量は、ユーザ発話の音声信号から得られる音響特徴量であれば、いずれの特徴量でもよい。

20

【0108】

例えば、音響特徴量は、基本周波数（F0）や、メル周波数ケプストラム係数（MFCC）等でもよい。但し、特徴量を計算すること自体に遅延が生じないことが好ましい。なお、MFCC等の音響特徴量を用いると、処理遅れは無くなるが、韻律情報が失われるという欠点がある。

【0109】

<再生装置20 / システム発話タイミング検出手段22 / 言語特徴量抽出手段22Bの構成>

30

【0110】

言語特徴量抽出手段22Bは、音声認識処理手段41による音声認識処理の結果として得られたユーザ発話の言語情報から言語特徴量を抽出する処理を実行するものである。この言語特徴量抽出手段22Bの設置は省略してもよい。

【0111】

具体的には、例えば、LSTM言語モデルの中間出力（512次元）を言語特徴量とすることができる（非特許文献2参照）。なお、LSTM（Long short-term memory）は、リカレントニューラルネットワーク（RNN）の一種である。

【0112】

また、音声認識処理手段41による音声認識処理は、上述した音響特徴量抽出手段22Aの処理と非同期で実行されるため、音響特徴量抽出手段22Aにより音響特徴量が抽出されたときに、この言語特徴量抽出手段22Bによる言語特徴量の抽出が行われていない場合があるので、その場合には、言語特徴量はゼロベクトルとする。

40

【0113】

<再生装置20 / システム発話タイミング検出手段22 / ユーザ発話権終了判定用パターン認識器22Cの構成>

【0114】

ユーザ発話権終了判定用パターン認識器22Cは、音響特徴量抽出手段22Aにより抽出した音響特徴量を入力とするか（非特許文献1参照）、あるいは、この音響特徴量および言語特徴量抽出手段22Bにより抽出した言語特徴量を入力とし（非特許文献2参照）

50

、ユーザが発話する地位または立場を有していることを示すユーザ発話権の維持または終了（終了には、譲渡、放棄が含まれる。）を識別するパターン認識処理を繰り返し実行するものである。

【0115】

このユーザ発話権終了判定用パターン認識器22Cは、識別アルゴリズムによるパターン認識処理を実行するユーザ発話権終了判定用パターン認識処理手段22Dと、このパターン認識処理で用いるモデル（パラメータ）を記憶するユーザ発話権終了判定モデル記憶手段22Eとにより構成されている。

【0116】

具体的には、例えば、音響特徴量（256次元）および言語特徴量（512次元）を入力とし、ユーザ発話権の維持または終了を逐次推定するモデルをニューラルネットワークにより構築し、ユーザ発話権終了判定用パターン認識器22Cとすることができる（非特許文献2参照）。この際、ニューラルネットワークには、時系列情報を考慮するため、LSTM（RNNの一種）を用いることができる（非特許文献1参照）。

【0117】

このユーザ発話権終了判定用パターン認識器22Cは、ユーザ発話権が終了したことの確からしさを示す尤度を出力するので、その尤度が、予め定められたユーザ発話権終了判定用閾値（但し、この閾値は、ユーザ発話権終了判定用閾値調整手段22Gにより、事前に、またはリアルタイムで動的に調整されることがある。）以上であるか、またはこの閾値を超えているかを判定する閾値処理を行い（図7の最下部を参照）、ユーザ発話権終了判定用閾値以上または超過と判定した場合には、ユーザ発話権が終了したという識別結果を出力し、ユーザ発話権終了判定用閾値未満または以下と判定した場合には、ユーザ発話権が維持されているという識別結果を出力する。

【0118】

この際、尤度は、ユーザ発話権が終了したことの確からしさを示す尤度としているので、尤度の値が大きい程（1に近い程）、ユーザ発話権の終了の状態に近く、尤度の値が小さい程（0に近い程）、ユーザ発話権の維持の状態に近い（図7の最下部を参照）。従って、尤度がユーザ発話権終了判定用閾値以上になるか、超えれば、ユーザ発話権が終了したという識別結果が出力されることになるので、ユーザ発話権終了判定用閾値の上方調整というのは、ユーザ発話権が終了したという識別結果が出にくくなる方向への調整であり、下方調整というのは、ユーザ発話権が終了したという識別結果が出やすくなる方向への調整である。本願の請求項は、この場合の記載とされている。

【0119】

一方、尤度は、ユーザ発話権が維持されていることの確からしさを示す尤度としてもよく、この場合には、尤度の値が大きい程（1に近い程）、ユーザ発話権の維持の状態に近く、尤度の値が小さい程（0に近い程）、ユーザ発話権の終了の状態に近い。従って、尤度がユーザ発話権終了判定用閾値以下になるか、未満になれば、ユーザ発話権が終了したという識別結果が出力されることになるので、ユーザ発話権終了判定用閾値の上方調整というのは、ユーザ発話権が終了したという識別結果が出やすくなる方向への調整であり、下方調整というのは、ユーザ発話権が終了したという識別結果が出にくくなる方向への調整である。このため、本願の請求項は、この場合とは逆の記載とされているが（上方、下方が逆の表現となっているが）、両者は等価なことであり、また、1から尤度の値を減じれば、逆の意味の尤度になるので、本願の請求項は、いずれの場合も含むものである。

【0120】

また、閾値処理を行う際には、フィルタをかけた後の尤度（出力される直近の幾つかの尤度を用いて平準化した後の尤度）を用いてもよい。

【0121】

さらに、ユーザ発話権終了判定用パターン認識器22Cは、ユーザ発話権の維持または終了を識別するパターン認識処理を実行する際に、終了については、質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図のうちの



いずれのユーザ発話意図で終了するのかを識別するパターン認識処理を実行してもよい。この場合、識別器の学習段階で、終了については、質問、応答、相槌等のユーザ発話意図を含むラベル(タグ)を付すことになる。すなわち、「維持」、「質問で終了」、「応答で終了」、「相槌で終了」等のタグ付けを行った学習用データをそれぞれ多数用意して学習を行い、3クラス識別以上のユーザ発話権終了判定モデルを構築する。そして、運用段階では、ユーザ発話権が「維持」されていることの確からしさを示す尤度、「質問で終了」したことの確からしさを示す尤度、「応答で終了」したことの確からしさを示す尤度等のように、いずれのユーザ発話意図で終了したのかを示す情報が出力される。例えば、質問で終了の尤度 = 0.90、応答で終了の尤度 = 0.04、相槌で終了の尤度 = 0.03等のように出力される。従って、ユーザ発話権終了判定用閾値は、質問、応答、相槌等のユーザ発話意図毎に設定し、ユーザ発話意図毎に閾値処理を行う。但し、ユーザ発話意図毎に設定したユーザ発話権終了判定用閾値が、同じ値となってもよい。このようにした場合、次発話選択手段24は、ユーザ発話権終了判定用パターン認識器22Cによるユーザ発話意図の識別結果(いずれのユーザ発話意図で終了したのかという情報)を用いて、次発話準備手段43による準備処理で得られた複数の次発話候補の内容データの中から、発話生成手段25で用いる次発話の内容データを選択することができる。そして、以上のように質問、応答、相槌等のユーザ発話意図を含むタグ付けをした学習で得られたユーザ発話権終了判定用パターン認識器22Cは、次発話選択用情報生成手段23のユーザ発話意図の識別器(図3に示す第1、第2の発話意図識別器23B, 23E)と、いずれのユーザ発話意図で終了したのかというタグ付けをしないで単純に維持・終了を識別するための学習で得られた識別器とを、マルチタスクでまとめて一体化させたユーザ発話権終了判定用パターン認識器22Cとは、異なるものである。

10

20

#### 【0122】

なお、ユーザ発話の音声信号を逐次処理して短い周期で識別を行う技術として、スマートスピーカのウェイクワードのスポッティングが挙げられるが、ユーザ発話権終了判定用パターン認識器22Cは、特定の語のスポッティングではなく、ユーザ発話権の維持または終了を、その発話内容に依らずに検出することを目的とする点で異なる。

#### 【0123】

<再生装置20/システム発話タイミング検出手段22/システム発話開始タイミング判断手段22Fの構成>

30

#### 【0124】

システム発話開始タイミング判断手段22Fは、ユーザ発話権終了判定用パターン認識器22Cによるパターン認識処理の結果(維持または終了の識別結果)を用いるか、またはこのパターン認識処理の結果に加え、システム状態記憶手段31に記憶されているシステム状態を示す情報(準備完了・準備中の別)や、ユーザ状態記憶手段32に記憶されているユーザ状態を示す情報(ユーザ発話継続時間)を用いて、システム発話の開始タイミングを検出する処理を実行するものである。

#### 【0125】

具体的には、図8に示すように、システム発話開始タイミング判断手段22Fは、まず、ユーザ発話権終了判定用パターン認識器22Cによるパターン認識処理の結果が、ユーザ発話権の維持を示している場合(P1)と、終了を示している場合(P2)とに判断分岐する。

40

#### 【0126】

次に、維持を示している場合(P1)には、システム発話の開始タイミングではないと判断する(P7)。

#### 【0127】

但し、図8中および図2中の二点鎖線で示すように、ユーザ状態記憶手段32に記憶されているユーザ発話継続時間が、予め定められた長時間判定用閾値以上または超過の場合には、パターン認識処理の結果がユーザ発話権の維持を示していても(P1)、システム

50

発話の開始タイミングであると判断する処理（P 8）を行ってもよい。

【0128】

一方、終了を示している場合（P 2）において、システム状態記憶手段31に記憶されているシステム状態を示す情報が準備完了（ステータス＝「準備完了」）を示している場合（P 3）には、システム発話の開始タイミングであると判断する（P 8）。

【0129】

但し、図8中および図2中の二点鎖線で示すように、ユーザ状態記憶手段32に記憶されているユーザ発話継続時間が、予め定められた短時間判定用閾値以下または未満の場合には、パターン認識処理の結果がユーザ発話権の終了を示していても（P 2）、システム発話の開始タイミングではないと判断する処理（P 7）を行ってもよい。

10

【0130】

また、終了を示している場合（P 2）において、システム状態記憶手段31に記憶されているシステム状態を示す情報が準備中を示している場合（P 4）には、その準備中を示すステータスに応じ（次発話準備手段43による準備中の処理内容）に応じ、判断を分岐させる。

【0131】

そして、準備中を示している場合（P 4）において、その準備中の処理内容が、直ぐに完了する処理内容として予め分類されている処理の準備中である場合（P 5）には、準備完了になるまで待つてシステム発話の開始タイミングであると判断するために（但し、結果的に、直ぐに準備完了にならない場合もある。）、その時点では、システム発話の開始タイミングではないと判断する（P 9）。直ぐに完了する処理内容として予め分類されている処理の準備中とは、例えば、ステータス＝「自サーバ検索中」等である。

20

【0132】

一方、準備中を示している場合（P 4）において、その準備中の処理内容が、直ぐに完了しない処理内容として予め分類されている処理の準備中である場合（P 6）には、システム発話の開始タイミングであると判断するとともに、フィラーの挿入タイミングである旨の情報を出力する（P 10）。直ぐに完了しない処理内容として予め分類されている処理の準備中とは、例えば、ステータス＝「外部システムアクセス中」、「音声合成処理中」等である。フィラーの挿入タイミングである旨の情報には、どのような種別のフィラーを挿入するかの情報を含めてもよく、この場合、準備中のステータスの種別と、フィラーの種別との対応関係を、予め定めておけばよい。例えば、直ぐに完了しない処理内容にも、その程度があるので、かなり長時間の準備を要する場合には、「ちょっと待ってね、今、調べてるから。」、「少々お待ちください、処理中です。」等のフィラーを挿入することができ、そこまで長時間を要しない場合には、「えー。」、「あのね。」等のフィラーを挿入することができる。

30

【0133】

なお、準備中を示すステータスのうち、どのようなステータスが、直ぐに完了する処理内容なのか、直ぐに完了しない処理内容なのかは、システムの構築、運用、管理を行う者が適宜設計すればよく、対話の種別（ニュース対話、アンケート対話、情報検索対話、操作対話、教育対話等の別）に応じて定めてもよい。

40

【0134】

また、図8において、P 9の下流部分で点線で示されているように、P 9の判断を行って待った結果、システム状態が変化することもあるので、次回以降の判断時の状態に従って、判断分岐が行われることになる。

【0135】

すなわち、待った結果、直ぐに準備処理が完了した場合には、P 2 P 3 P 8という流れとなり、一方、直ぐに完了しない別の準備処理に移行した場合には、P 2 P 4 P 6 P 10という流れとなる。直ぐに完了しない別の準備処理に移行した場合とは、例えば、ステータス＝「自サーバ検索中」であったが、自サーバ内で目的の情報が得られなかったため、外部システムにアクセスし、ステータス＝「外部システムアクセス中」となっ

50

た場合等である。

【0136】

さらに、図8において、P10の判断に基づきフィラーの再生を開始した後、フィラーの再生を行っている間に準備が完了すれば、P2 P3 P8という流れとなり、フィラーの再生を中断するか、またはフィラーの再生終了後に、準備が完了した複数の次発話候補の中からの次発話の選択が行われ、選択された次発話の再生が行われることになる。一方、フィラーの挿入(P10)を行っても、未だ準備が続いていた場合には、直ぐに完了しない準備処理が続いていることになるので、P2 P4 P6 P10という流れとなり、再び、フィラーの挿入(P10)が行われる。なお、フィラーの再生中に、P10の判断が再びなされた場合には、再生中のフィラーを優先させて再生を続ける。新たなP10の判断を優先させると、例えば「ちょっと待っ」「ちょっと待っ」「ちょっと待っ」のような繰り返しをする不自然な発話になってしまうからである。

10

【0137】

<再生装置20/システム発話タイミング検出手段22/ユーザ発話権終了判定用閾値調整手段22Gの構成>

【0138】

ユーザ発話権終了判定用閾値調整手段22Gは、ユーザ発話権終了判定用パターン認識器22Cによる維持・終了の識別処理で用いるユーザ発話権終了判定用閾値の事前調整、ユーザ発話権終了判定用閾値の下方調整を行うことを決めるための下方調整用閾値の事前調整、およびユーザ発話権終了判定用閾値のリアルタイム調整の各処理を実行するものである。

20

【0139】

ここで、事前調整は、ユーザとの対話(その日またはその時における対話、あるいは、その週、その月、その季節、その年等の所定の区切りの期間における対話)を開始する前に行う調整であり、ユーザ情報記憶手段52に記憶されているユーザの属性情報(当該ユーザとの対話中における一時的な情報ではなく、当該ユーザとの複数回の対話を通じて得られた蓄積情報)を用いて行われる。一方、リアルタイム調整は、ユーザとの対話の開始後(特に、ユーザ発話の進行中)に行う調整であり、ユーザ状態記憶手段32に記憶されているユーザ状態を示す情報(対話中における一時的な情報)や、システム状態記憶手段31に記憶されているシステム状態を示す情報(対話中における一時的な情報)を用いて行われる。

30

【0140】

具体的には、図11に示すように、ユーザ発話権終了判定用閾値調整手段22Gは、対話相手のユーザについてのユーザ識別情報を用いてユーザ情報記憶手段52に記憶されている当該ユーザの衝突の発生情報(蓄積情報)を取得し、当該ユーザとの衝突の発生頻度または累積発生回数を算出する。この際、衝突の発生頻度は、例えば、1日、1週間、1か月等の所定の長さの期間における衝突の発生回数としてもよく、対話の総数に対する衝突の発生回数としてもよく、ユーザ発話からシステム発話への交替の総数に対する衝突の発生回数としてもよく、発生頻度の単位は、任意である。そして、算出した衝突の発生頻度または累積発生回数が、予め定められた上方調整用閾値以上または超過の場合には、ユーザ発話権終了判定用閾値を標準値または前回調整値よりも高く設定する上方調整を実行する。これにより、図11に示すように、システム発話の開始タイミングが遅れる方向、すなわち衝突回避方向に調整される。

40

【0141】

また、図12に示すように、ユーザ発話権終了判定用閾値調整手段22Gは、対話相手のユーザについてのユーザ識別情報を用いてユーザ情報記憶手段52に記憶されている当該ユーザについてのユーザ発話の終了からシステム発話の開始までの複数の交替潜時(システムの交替潜時の蓄積情報)を取得し、当該ユーザを対話相手とするときのシステムの交替潜時の長短の傾向を示す平均値若しくはその他の指標値を算出する。この際、交替潜時の長短の傾向を示す指標値は、複数の交替潜時をまとめた指標値であれば、いずれでも

50

よく、例えば、平均値、中央値、最頻値等とすることができ、中央値や最頻値とする場合は、交替潜時を幾つかに区分していずれかの区分に帰属させ、各区分の代表値の中のいずれかを中央値、最頻値とすること等ができる。そして、算出した交替潜時の指標値が、予め定められた下方調整用閾値以上または超過の場合には、ユーザ発話権終了判定用閾値を標準値または前回調整値よりも低く設定する下方調整を実行する。これにより、図 12 に示すように、システム発話の開始タイミングが早まる方向、すなわち交替潜時が短くなる方向に調整される。

#### 【0142】

さらに、図 12 に示すように、ユーザ発話権終了判定用閾値調整手段 22G は、対話相手のユーザについてのユーザ識別情報を用いてユーザ情報記憶手段 52 に記憶されている当該ユーザの複数の発話速度（蓄積情報）を取得し、当該ユーザの発話速度の傾向を示す平均値若しくはその他の指標値を算出する。なお、発話速度の単位は「モーラ/秒」等である。この際、ユーザの発話速度の傾向を示す指標値は、複数の発話速度をまとめた指標値であれば、いずれでもよく、例えば、平均値、中央値、最頻値等とすることができ、そして、下方調整用閾値を、算出した発話速度の指標値を用いて、発話速度の指標値が大きい（速い）ときには当該下方調整用閾値が小さくなり、発話速度の指標値が小さい（遅い）ときには当該下方調整用閾値が大きくなるように予め定められた関数により設定する。この関数は、上述した前提条件を満たす関数であれば、どのような関数でもよく、図 12 の例では、1 次関数とされているが、例えば、2 次以上の関数でもよく、1 段または多段のステップ関数等でもよい。これにより、早口のユーザについては、下方調整用閾値が小さくなり、比較的短い交替潜時でも、ユーザ発話権終了判定用閾値の下方調整を行うことができるようになり（下方調整の条件が成立するようになり）、交替潜時が短くなる方向へのシステム発話の開始タイミングの調整を行うことができるようになる。一方、ゆっくり発話するユーザについては、下方調整用閾値が大きくなり、比較的長い交替潜時でないと、ユーザ発話権終了判定用閾値の下方調整を行うことができないようになり（下方調整の条件が成立しなくなり）、交替潜時が短くなる方向へのシステム発話の開始タイミングの調整を行うことができないようになる。

#### 【0143】

また、図 9 中の実線で示すように、ユーザ発話権終了判定用閾値調整手段 22G は、ユーザ状態記憶手段 32 に記憶されている対話相手のユーザについてのユーザ発話継続時間（リアルタイム情報）を逐次取得し、取得したユーザ発話継続時間が、予め定められた短時間判定用閾値以下または未満の場合には、ユーザ発話権終了判定用閾値を標準値よりも高く設定し、予め定められた長時間判定用閾値以上または超過の場合には、ユーザ発話権終了判定用閾値を標準値よりも低く設定する処理を逐次実行する。これにより、ユーザ発話の開始直後には、ユーザ発話権が終了したという識別結果が出にくくなり、ユーザ発話の開始時からの経過時間が長くなると、ユーザ発話権が終了したという識別結果が出やすくなる。

#### 【0144】

また、図 9 中の二点鎖線で示すように、ユーザ発話権終了判定用閾値調整手段 22G は、ユーザ状態記憶手段 32 に記憶されている対話相手のユーザについてのユーザ発話継続時間（リアルタイム情報）を逐次取得し、ユーザ発話権終了判定用閾値を、取得したユーザ発話継続時間を用いて、ユーザ発話継続時間が短いときには当該ユーザ発話権終了判定用閾値が高くなり、ユーザ発話継続時間が長いときには当該ユーザ発話権終了判定用閾値が低くなるように予め定められた関数（図 9 中の実線で示された階段状の関数に限らず、それ以外の様々な関数）により逐次設定してもよい。この関数は、上述した前提条件を満たす関数であれば、どのような関数でもよく、例えば、1 次関数でもよく、2 次以上の関数でもよく、1 段のステップ関数や、図 9 中の実線で示された 2 段のステップ関数以外の多段（3 段以上）のステップ関数等でもよい。

#### 【0145】

さらに、図 10 に示すように、ユーザ発話権終了判定用閾値調整手段 22G は、システ

ム状態記憶手段31に記憶されている対話相手のユーザについての目的データの残数（対話目的を達成するためのシステムの最終の次発話候補の内容データとなり得る題材データである目的データの残数）および／または次発話候補の重要度（次発話準備手段43による準備処理で得られた複数の次発話候補の内容データの各々に付されている重要度）を取得する処理を逐次実行する。これらの目的データの残数および／または次発話候補の重要度は、システムによる発話開始に対する要求の強さの度合いを示すシステム発話意欲度の指標値である。

#### 【0146】

ここで、目的データの残数については、例えば、目的データの残数が1であれば、システム発話意欲度が強く、目的データの残数が2以上であれば、システム発話意欲度が弱い設定とすること等ができる。例えば、情報検索対話において、ユーザ発話の進行に伴ってユーザによる条件提示が進み、その条件提示の内容に応じて目的データの残数が1になった時点で、システム発話意欲度を強く設定することができる。具体例を挙げると、飲食店を検索するときに、ユーザが、システムの「食べる場所はどこ？」に対して「東京駅周辺のお店を探したい。」と答え、システムの「何が食べたいの？」に対して「中華料理が食べたい。」と答え、システムの「どんなお店がいいの？」に対して「おいしいと評判のお店がよくて、それと・・・」と答える等の条件提示を積み重ねていった結果、目的データ（情報提供する飲食店のデータ）が1つに絞り込まれる場合があり、この場合、ユーザは、それ以上、条件提示を行う必要はなく（つまり、「それと・・・」以降の条件提示を行う必要はなく）、1つに絞り込まれた目的データ（飲食店のデータ）を早く再生した方がよいという状況になるので、システム発話意欲度が強くなる。また、ユーザが「おいしいと評判のお店がよい。」と言った後に「待って、やっぱり評判はどうでもいいから、安いお店がいいな。」と訂正の発話を行い、それに基づき、再び、目的データ（飲食店のデータ）の残数が2以上になったときには、システム発話意欲度が弱くなる。システム発話意欲度の数値化方法は任意であり、例えば、1～10の10段階（段階数は任意）、0～1の連続値、0～100%の連続値等とすることができる。例えば、目的データの残数＝1の場合には、システム発話意欲度＝10段階のうちの10とし、目的データの残数＝2または3の場合には、システム発話意欲度＝10段階のうちの7とし、目的データの残数＝4以上の場合には、システム発話意欲度＝10段階のうちの2とすること等ができる。この対応関係は、予め定めておけばよい。なお、上記の例の対応関係では、10段階のうち使用されないシステム発話意欲度が存在するが、これは、下記の次発話候補の重要度により定まるシステム発話意欲度とのレベル合わせをしているからである。

#### 【0147】

また、次発話候補の重要度については、重要度が高ければ、システム発話意欲度が強くなり、重要度が低ければ、システム発話意欲度が弱くなる関係にある。この重要度は、記事データ（ニュースやコラムや歴史等の各種の話題を記載した記事の原文データ）を要約してシナリオデータを生成する際の元の記事データの各構成文の重要度と同じとしてもよいが、本実施形態では、それだけではなく、防災関連情報の緊急性や日常生活への影響の大きさ等を加味した重要度としている。例えば、ニュース対話において、重要度が、「XXXで大きな地震が発生しましたので、YYY沿岸地域の方は、すぐに高台に避難してください。」＝10、「XXX地方に大雨洪水警報が出ました。」＝8、「明日から消費税が10%となります。」＝6、「早稲田花子選手が女子100mの日本新記録を出しました。」＝4等のように、1～10の10段階の数値で設定されていれば、これらの数値をそのままシステム発話意欲度を示す数値とすること等ができる。また、重要度が8以上は、システム発話意欲度＝3段階のうちの3とし、重要度が7～5は、システム発話意欲度＝3段階のうちの2とし、重要度が4以下は、システム発話意欲度＝3段階のうちの1とすること等ができる。この対応関係は、予め定めておけばよい。そして、複数の次発話候補の内容データが次発話候補記憶手段30に記憶され、システム状態記憶手段31にそれらの複数の次発話候補の内容データの各々についての重要度が記憶されている場合には、複数の重要度の平均値、中央値、最頻値等を代表の重要度としてもよく、最も大きい重要

10

20

30

40

50

度や最も小さい重要度を代表の重要度としてもよい。

【0148】

なお、目的データの残数や、複数の次発話候補の内容データの各々の重要度をシステム状態記憶手段31に記憶させるのではなく、次発話準備手段43によりこれらを換算して求めたシステム発話意欲度を、システム状態記憶手段31に記憶させてもよい。また、複数の次発話候補の内容データの各々の重要度をシステム状態記憶手段31に記憶させるのではなく、次発話準備手段43により求めた代表の重要度を、システム状態記憶手段31に記憶させてもよい。さらに、目的データの残数により定まるシステム発話意欲度と、次発話候補の重要度により定まるシステム発話意欲度とを、対話の種別に応じて使い分けてもよいが、両者の平均値や加重平均値等を求めて統合して用いてもよい。

10

【0149】

そして、図10に示すように、ユーザ発話権終了判定用閾値調整手段22Gは、取得した目的データの残数および/または複数の次発話候補の内容データの各々の重要度からシステム発話意欲度を求め、ユーザ発話権終了判定用閾値を、求めたシステム発話意欲度を用いて、システム発話意欲度が強いときには当該ユーザ発話権終了判定用閾値が低くなり、システム発話意欲度が弱いときには当該ユーザ発話権終了判定用閾値が高くなるように予め定められた関数により設定する処理を逐次実行する。これにより、システム発話意欲度が強いときには、ユーザ発話権が終了したという識別結果が出やすくなり、システム発話意欲度が弱いときには、ユーザ発話権が終了したという識別結果が出にくくなる。

【0150】

20

<再生装置20/次発話選択用情報生成手段23の構成>

【0151】

図3において、次発話選択用情報生成手段23は、韻律特徴量抽出手段23Aと、第1の発話意図識別器23Bと、第2の発話意図識別器23Eとを含んで構成されている。この次発話選択用情報生成手段23には、例えば、前述した非特許文献3に記載された技術を採用することができる。

【0152】

韻律特徴量抽出手段23Aは、音声信号取得手段21により取得したユーザ発話の音声信号から韻律特徴量(韻律情報)を抽出する処理を実行するものである。この韻律特徴量抽出手段23Aは、システム発話タイミング検出手段22の音響特徴量抽出手段22Aと同様な構成を採用することができる。すなわち、音響特徴量抽出手段22Aで抽出された音響特徴量を、韻律特徴量(韻律情報)とすることができる。従って、この韻律特徴量抽出手段23Aと、システム発話タイミング検出手段22の音響特徴量抽出手段22Aとは、共通化することができる。従って、例えば、韻律特徴量抽出手段23Aで得られる韻律特徴量(韻律情報)は、CNNオートエンコーダの中間層から取り出した256次元のボトルネック特徴量とすることができる。

30

【0153】

第1の発話意図識別器23Bは、韻律特徴量抽出手段23Aで抽出した韻律特徴量(韻律情報)を用いて、質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図を識別する処理を実行する第1の発話意図識別処理手段23Cと、この識別処理で用いるモデル(パラメータ)を記憶する第1の発話意図識別モデル記憶手段23Dとにより構成されている。この第1の発話意図識別器23Bは、例えば、LSTM(RNNの一種)により構築することができる。

40

【0154】

具体的には、第1の発話意図識別器23Bは、例えば、CNNオートエンコーダの中間層から取り出した256次元の韻律特徴量(韻律情報)を逐次入力し、LSTMによるパターン認識処理を行って発話意図を識別し、その識別結果を出力する構成とすることができる(非特許文献3参照)。そして、この第1の発話意図識別器23Bから出力された発話意図を、次発話選択手段24に送ってもよい。

【0155】

50

第2の発話意図識別器23Eは、第1の発話意図識別器23Bで得られた韻律情報（例えば、LSTMの隠れ層の値）と、音声認識処理手段41による音声認識処理の結果として得られたユーザ発話の言語情報と、対話履歴記憶手段50に記憶されているユーザとシステムとの間の対話履歴情報のうちの直前のシステム発話の言語情報とを用いて、質問、応答、相槌、補足要求、反復要求、理解、不理解、無関心、若しくはその他のユーザ発話意図を識別する処理を実行する第2の発話意図識別処理手段23Fと、この識別処理で用いるモデル（パラメータ）を記憶する第2の発話意図識別モデル記憶手段23Gとにより構成されている。この第2の発話意図識別器23Eは、例えば、BERTにより構築することができる（非特許文献3参照）。BERTは、自然言語処理モデルであり、トランスフォーマのエンコーダ部分をユニットとする双方向トランスフォーマモデルである。この第2の発話意図識別器23Eから出力された発話意図は、次発話選択手段24に送られる。

10

#### 【0156】

また、次発話選択用情報生成手段23は、ユーザの顔画像やジェスチャー画像（身振り・手振りの画像）を取得し、顔の表情やジェスチャーの内容を解析し、その解析結果（表情の識別結果、身振り・手振りの意図の識別結果）を、次発話選択用情報として次発話選択手段24に送ってもよい。

#### 【0157】

<再生装置20/次発話選択手段24の構成>

#### 【0158】

次発話選択手段24は、システム発話タイミング検出手段22によりシステム発話の開始タイミングが検出された後に（システム発話タイミング検出手段22からシステム発話の開始タイミングであるという判断結果を受け取ったときに）、次発話選択用情報生成手段23による処理で得られたユーザ発話意図の識別結果と、音声認識処理手段41による音声認識処理の結果として得られた言語情報（文字列）とを組み合わせ用いて、次発話準備手段43による準備処理で得られて次発話候補記憶手段30に記憶されている複数（但し、1つの場合もある。）の次発話候補の内容データの中から、発話生成手段25で用いる次発話の内容データを選択し、選択した次発話の内容データを、発話生成手段25に送るとともに、選択した次発話の内容データまたはその識別情報（例えば、シナリオID、発話節ID等）を、ネットワーク1を介して対話状態管理手段42へ送信する処理を実行するものである。

20

30

#### 【0159】

なお、ユーザ発話意図の識別結果（例えば、質問、相槌等の別）と、音声認識処理の結果として得られた言語情報（文字列）とのうちのいずれか一方だけで、次発話の内容データを選択することができる場合には、これらを組み合わせ用いなくてもよい。また、次発話選択用情報生成手段23からではなく、システム発話タイミング検出手段22から、システム発話の開始タイミングであるという判断結果とともにユーザ発話意図の識別結果を受け取った場合（システム発話タイミング検出手段22において、どのようなユーザ発話意図で終了したのかを識別した場合）には、そのユーザ発話意図の識別結果を用いて、次発話の内容データを選択してもよい。さらに、次発話選択用情報生成手段23から、ユーザの顔の表情やジェスチャーの内容についての解析結果を受け取った場合には、それらの解析結果を用いて、またはそれらの解析結果と他の情報とを組み合わせ、次発話の内容データを選択してもよい。

40

#### 【0160】

また、次発話選択手段24は、システム発話タイミング検出手段22から、システム発話の開始タイミングであるという判断結果とともに、フィラーの挿入タイミングである旨の情報（どのような種別のフィラーを挿入するかの情報を含む）を受け取った場合には、指定された種別のフィラーの内容データ（音声データを含む）を、発話生成手段25に送る処理を実行する。この際、挿入するフィラーの内容データまたは当該フィラーの種別の識別情報を、ネットワーク1を介して対話状態管理手段42へ送信する処理を実行しても

50

よく、実行しなくてもよいが、実行した場合でも、対話状態管理手段 42 は、フィラーの挿入を、対話履歴上は、システム発話として取り扱うのではなく、システム発話の準備用繋ぎ発話として取り扱う。この点については、後述する対話状態管理手段 42 の説明で詳述するので、ここでは詳しい説明を省略する。なお、各種のフィラーの内容データ（音声データを含む）は、フィラーの種別の識別情報と対応付けて再生装置 20 に設けられたフィラー記憶手段（不図示）に記憶させておけばよいが、フィラーは常に次発話候補になり得ると考え、次発話候補記憶手段 30 に記憶させておいてもよい。後者とする場合、各種のフィラーの内容データ（音声データを含む）を、次発話準備手段 43 から毎回ネットワーク 1 を介して受信する必要はなく、次発話候補記憶手段 30 に固定的に準備されているデータとすればよい。なお、このようにフィラーを次発話候補であると考えて次発話候補記憶手段 30 に記憶させたとしても、上述したように、対話状態管理手段 42 は、フィラーの挿入が行われても、対話履歴上は、それをシステム発話として取り扱わないので、次発話候補記憶手段 30 に固定的に記憶しておくフィラーと、それ以外の複数の次発話候補の内容データ（頻繁に更新されるデータ）とは別のものであり、単に同じ次発話候補記憶手段 30 に記憶させるに過ぎない。

10

20

30

40

50

#### 【0161】

具体的には、次発話選択手段 24 は、音声認識処理手段 41 による音声認識処理の結果（文字列）を用いて次発話の選択を行う場合には、音声認識処理の結果に含まれる各単語と、次発話候補記憶手段 30 に記憶されている複数の次発話候補の内容データの各々に含まれる各単語とを用いて、キーワードマッチングを行い、マッチングした次発話候補の内容データを、次発話の内容データとして選択することができる。また、言語処理と機械学習とを合わせた複雑なマッチングを行ってもよい。さらに、音声認識処理の結果として得られた文字列と、複数の次発話候補の内容データ（文字列）の各々との類似度を、doc2vec 等により求め、類似度の高い次発話候補の内容データを、次発話の内容データとして選択してもよい。

#### 【0162】

また、電話の自動応答における音声対話等のように、システムがユーザに質問し、ユーザがそれに答えていく場合には、ユーザ発話の内容は、システムから与えられた選択肢等のように限られたものになるので、ユーザ発話の内容は予測することができる。この場合、次発話候補記憶手段 30 に記憶されている複数の次発話候補の内容データの各々に、対応するユーザ発話の予測データが付随していれば、その付随しているユーザ発話の予測データのうちのいずれが、音声認識処理の結果として得られた文字列と一致するのかを判断することにより、一致したユーザ発話の予測データに対応する次発話候補の内容データを、次発話の内容データとして選択することができる。

#### 【0163】

例えば、システム発話  $S(N)$  が「 $XXX$ 党と、 $YYY$ 党のどちらを支持しますか？」であり、システム発話  $S(N+1)$  の複数（2つ）の候補として、「 $XXX$ 党のどの政治家が総理大臣になると思いますか？」という内容データおよびそれに付随する「 $XXX$ 党」というユーザ発話  $U(N)$  の予測データと、「 $YYY$ 党のどの政治家が党首に相応しいですか？」という内容データおよびそれに付随する「 $YYY$ 党」というユーザ発話  $U(N)$  の予測データとを、次発話準備手段 43 により準備し、次発話候補記憶手段 30 に記憶させたとする。このとき、ユーザ発話  $U(N)$  の音声認識処理の結果が「 $XXX$ 党」であれば、「 $XXX$ 党」というユーザ発話  $U(N)$  の予測データと一致するので、それに対応する「 $XXX$ 党のどの政治家が総理大臣になると思いますか？」がシステム発話  $S(N+1)$  として選択され、発話生成手段 25 により再生される。

#### 【0164】

さらに、次発話選択手段 24 は、ユーザ発話意図の識別結果（質問、相槌等）を用いて次発話の選択を行う場合には、ユーザ発話意図に対応するシステム発話種別が、次発話候補記憶手段 30 に記憶されている複数の次発話候補の内容データの各々について定められているので、得られたユーザ発話意図の識別結果に対応するシステム発話種別である次発



話候補の内容データを、次発話の内容データとして選択することができる。

【0165】

例えば、ユーザ発話意図が「相槌」、「理解」等であれば、システム発話種別が主計画である次発話候補の内容データを選択し、ユーザ発話意図が「定義型質問」（用語の意義を問う質問）であれば、システム発話種別が副計画（定義）である次発話候補の内容データを選択し、ユーザ発話意図が「反復要求」、「不理解」であれば、システム発話種別が繰り返し用の主計画である次発話候補の内容データを選択し、ユーザ発話意図が「補足要求」等であれば、補足説明用の副計画（トリビア等）である次発話候補の内容データを選択する等のように、ユーザ発話意図と、システム発話種別との対応関係を予め定めておけばよい。この対応関係は、次発話選択手段24を構成するプログラム内に記述されていてもよく、再生装置20に設けられた発話意図・システム発話種別対応関係記憶手段（不図示）に記憶しておいてもよい。従って、シナリオデータの構成要素が、主計画要素であるか、副計画要素であるかも、システム発話種別に該当する。なお、主計画、副計画についての詳細は、図13および図14を用いて後述する。

10

【0166】

なお、ユーザ発話意図が「無関心」、「既知」であれば、長短2つ用意された主計画のうちの短い方の主計画である次発話候補の内容データを選択し、伝達情報量を減らすことができる。但し、これらのユーザ発話意図の場合は、次発話準備手段43の入替準備手段43Cにより、複数の次発話候補の内容データを入れ替える準備処理が進行しているか（ステータス＝準備中）、あるいは、既にその準備が完了し、次発話候補記憶手段30に、別の話題のシナリオデータ内の先頭の構成要素（主計画要素）が記憶されているか、同じシナリオデータ内の別の構成要素（主計画要素）が記憶されていることもある。その場合は、その主計画要素を選択すればよい。

20

【0167】

また、次発話選択手段24は、ユーザ発話意図の識別結果（質問、相槌等）と、音声認識処理の結果として得られた言語情報（文字列）とを組み合わせ用いて、次のように、次発話を選択することができる。

【0168】

例えば、システム発話S(N)が「早稲田太郎選手が4回転フリップを成功させたよ。」であり、システム発話S(N+1)の複数の候補として、「グランプリシリーズのカナダ大会で跳んだそうだ。」（主計画要素）と、「早稲田太郎選手は、...」という早稲田太郎の人物の説明データ（副計画要素の定義）と、「4回転フリップっていうのは、...」という4回転フリップの技の説明データ（副計画要素の定義）と、繰り返し用の「早稲田太郎選手が4回転フリップを成功させたよ。」（主計画要素）とが、次発話準備手段43により準備され、次発話候補記憶手段30に記憶されているとする。このとき、U(N)のユーザ発話意図が「相槌」、「理解」であったとすると、「グランプリシリーズのカナダ大会で跳んだそうだ。」（主計画要素）を次発話S(N+1)として選択すればよく、ユーザ発話意図が「反復要求」であったとすると、繰り返し用の「早稲田太郎選手が4回転フリップを成功させたよ。」（主計画要素）を選択すればよい。しかし、U(N)のユーザ発話意図が「質問」であったとすると、次発話候補記憶手段30に記憶されている複数の次発話候補の内容データの中には、定義型質問に対するシステム応答（副計画要素の定義）が2つ用意されているので、早稲田太郎選手について質問しているのか、4回転フリップについて質問しているのかが判明しないと、システム応答を行うことができないが、いずれの質問であるかは、ユーザ発話意図だけでは判断することができない。そこで、音声認識処理の結果として得られた言語情報（文字列）を用いて、どちらの質問であるかを判断し、どちらのシステム応答（副計画要素の定義）を選択するのかを判断する。一方、次発話候補記憶手段30に記憶されている複数の次発話候補の内容データの中に、定義型質問に対するシステム応答（副計画要素の定義）が1つしかない場合には、音声認識処理の結果を使用せずに（つまり、ユーザ発話意図だけで）、その1つのシステム応答（副計画要素の定義）を選択することができる。

30

40

50

## 【 0 1 6 9 】

また、逆に音声認識処理の結果だけでは、次発話を選択できない場合もある。例えば、「えっ？」というユーザ発話は、驚きなのか、質問なのか、聞き返し（反復要求）なのかは判断できないので、ユーザ発話意図を用いて、次発話を選択することができる。

## 【 0 1 7 0 】

なお、次発話候補の内容データが、次発話候補記憶手段 3 0 に 1 つも記憶されていない期間があるが、これは、次発話選択手段 2 4 による処理には影響しない。なぜなら、次発話候補の内容データが次発話候補記憶手段 3 0 に 1 つも記憶されていない期間は、次発話準備手段 4 3 による準備処理が完了していない期間（準備中の期間）であるが、図 8 の P 5 P 9 の流れの場合（直ぐに準備が完了する場合）には、システム発話の開始タイミングではないと判断されるので、次発話選択手段 2 4 による処理には進まず、一方、図 8 の P 6 P 1 0 の流れの場合（直ぐに準備が完了しない場合）には、次発話選択手段 2 4 による処理に進むものの、フィラーの挿入になるので、次発話候補の内容データは必要ないからである。

## 【 0 1 7 1 】

< 再生装置 2 0 / 発話生成手段 2 5 の構成 >

## 【 0 1 7 2 】

発話生成手段 2 5 は、システム発話タイミング検出手段 2 2 によりシステム発話の開始タイミングが検出された後に、次発話選択手段 2 4 で選択された次発話の内容データ（次発話準備手段 4 3 による準備処理で得られた複数の次発話候補の内容データの中から選択された次発話の内容データ）を用いて、システム発話の音声信号の再生を含むシステム発話生成処理を実行するものである。この発話生成手段 2 5 には、スピーカ、ディスプレイも含まれる。

## 【 0 1 7 3 】

この際、発話生成手段 2 5 は、次発話選択手段 2 4 から受け取った次発話の内容データに音声データ（例えば wav ファイル等）が含まれていない場合には、次発話選択手段 2 4 から受け取ったテキストデータから音声データを生成する音声合成処理も実行する。但し、音声合成処理は、システム応答の遅延防止の観点から、次発話準備手段 4 3 で実行するか、または題材データとして予め用意されていることが好ましい。

## 【 0 1 7 4 】

また、発話生成手段 2 5 は、システム発話の音声信号の再生処理を実行するとともに、次発話選択手段 2 4 から受け取った次発話の内容データに映像データや静止画データ、あるいは楽曲データが付随している場合には、ディスプレイでの動画や静止画の再生処理、あるいは音楽の再生処理も実行する。例えば、直前のシステム発話が「早稲田太郎選手が 4 回転フリップを跳びました。」であり、それに対するユーザの反応が「4 回転フリップってどんな技？」という質問だった場合に、4 回転フリップの技の説明用の映像を再生し、「早稲田太郎選手ってどんな選手なの？」という質問だった場合に、早稲田太郎選手の顔画像を再生すること等ができる。また、直前のシステム発話が「X X X ホールで第九が演奏されました。」であり、それに対するユーザの反応が「第九ってどんな曲なの？」という質問だった場合に、第九の楽曲データを再生すること等ができる。なお、システム発話中に、システム発話の音声信号の再生と同期または略同期させて、ディスプレイでシステム発話の内容を示すテキスト表示を行ってもよい。

## 【 0 1 7 5 】

さらに、発話生成手段 2 5 は、音声信号取得手段 2 1 により取得したユーザ発話の音声信号と、再生中のシステム発話の音声信号との衝突の発生を検出し、検出した衝突の発生情報を、ネットワーク 1 を介して対話サーバ 4 0 へ送信し、ユーザ識別情報と関連付けてユーザ情報記憶手段 5 2 に記憶させる処理を実行する。この際、衝突には 2 種類あり、ここで検出する衝突は、[ 1 ] ユーザ発話権が終了したという識別結果が出て、システム発話を開始したところ、実際にはユーザ発話権が維持されていて、両者の発話が重なった場合の衝突である。従って、[ 2 ] ユーザ発話権が終了したものの、システム発話の開始が

遅れたために、再び、ユーザ発話が開始されてしまい、ほぼ同時に両者の発話が開始されて重なった場合の衝突ではないので、この〔 2 〕の場合の衝突を排除する処理を実行する。例えば、衝突を起こしたときの直前の無音区間の長さ（衝突を起こしたユーザ発話の音声区間の開始時点とその直前のユーザ発話の音声区間の終了時点との間の時間間隔）が、予め定めた衝突種別判定用閾値以上または超過の場合に、〔 2 〕の場合の衝突であると判断し、排除すること等ができる。また、衝突の発生前後の関連するデータを全て保存しておき、事後的に〔 2 〕場合の衝突であるか否かを判断し、〔 2 〕の衝突を排除する処理を行ってもよい。関連するデータとは、例えば、衝突の直前のユーザ発話の音声区間の終了時刻、システム発話タイミング検出手段 2 2 によるシステム発話の開始タイミングの検出時刻、発話生成手段 2 5 によるシステム発話の音声信号の再生開始時刻、衝突を起こしたユーザ発話の音声区間の開始時刻、衝突の前後双方のユーザ発話の音声認識処理の結果としての言語情報、衝突を起こしたシステム発話の内容データ（テキストデータ）等である。さらに、これらの関連するデータを用いて学習を行い、〔 1 〕と〔 2 〕の衝突を識別する識別器を構築し、事後的に、またはリアルタイム若しくは略リアルタイムで、識別器による識別結果に従って、〔 2 〕の衝突を排除する処理を行ってもよい。なお、本発明は、予め次発話候補を準備する等、システム応答の遅延防止が図られているので、〔 2 〕の場合の衝突は、殆ど発生しないようになっている。

10

#### 【 0 1 7 6 】

また、発話生成手段 2 5 は、音声信号取得手段 2 1 により取得したユーザ発話の音声信号を用いてユーザ発話の終了時刻を検出するとともに、システム発話の音声信号の再生開始時刻を検出することにより、ユーザ発話の終了からシステム発話の開始までの交替潜時を計測し、計測したシステムの交替潜時を、ネットワーク 1 を介して対話サーバ 4 0 へ送信し、ユーザ識別情報と関連付けてユーザ情報記憶手段 5 2 に記憶させる処理を実行する。

20

#### 【 0 1 7 7 】

さらに、発話生成手段 2 5 は、音声信号取得手段 2 1 により取得したユーザ発話の音声信号を用いてユーザ発話の開始時刻を検出し、検出した開始時刻と現在時刻との差分によりシステム発話継続時間を逐次計測し、計測したシステム発話継続時間を、ユーザ状態記憶手段 3 2 に逐次記憶させる処理を実行する。

30

#### 【 0 1 7 8 】

そして、発話生成手段 2 5 は、音声認識処理手段 4 1 による音声認識処理の結果として得られたユーザ発話の言語情報をネットワーク 1 を介して逐次取得し、取得した言語情報およびその取得時刻（または、言語情報とともに取得した時刻情報若しくは時間情報）を用いて発話速度をリアルタイムで算出し、算出したリアルタイムの発話速度を、ユーザ状態記憶手段 3 2 に逐次記憶させる処理を実行する。また、発話生成手段 2 5 は、対話全体におけるユーザ発話についての発話速度を算出し、算出した発話速度を、ネットワーク 1 を介して対話サーバ 4 0 へ送信し、ユーザ識別情報と関連付けてユーザ情報記憶手段 5 2 に記憶させる処理を実行する。

#### 【 0 1 7 9 】

< 再生装置 2 0 / 次発話候補記憶手段 3 0 、システム状態記憶手段 3 1 、ユーザ状態記憶手段 3 2 の構成 >

40

#### 【 0 1 8 0 】

次発話候補記憶手段 3 0 は、次発話準備手段 4 3 からネットワーク 1 を介して送信されてきた複数の次発話候補の内容データを、それらのデータの識別情報（例えば、シナリオ ID、発話節 ID 等）と対応付けて記憶するものである。記憶する次発話候補の内容データには、テキストデータの他に、音声データ（例えば wav ファイル等）が含まれ、さらに映像データや静止画データ、あるいは楽曲データが付随している場合もある。なお、再生装置 2 0 にフィルタ記憶手段（不図示）を設けない場合には、各種のフィルタの内容データ（音声データを含む）を、既に述べたように固定的に準備されているデータとして、フィルタの種別の識別情報と対応付けて次発話候補記憶手段 3 0 に記憶しておいてもよい

50

。

#### 【0181】

また、複数の次発話候補の内容データの各々には、データの属性を示すシステム発話種別（例えば、シナリオデータにおける主計画・副計画の別等）が対応付けられて記憶されている。このシステム発話種別については、次発話選択手段24の説明で既に詳述しているので、ここでは詳しい説明を省略する。

#### 【0182】

システム状態記憶手段31は、システム状態を示す情報として、次発話準備手段43による準備処理の状態（準備完了・各種の準備中の別を示すステータス）、目的データの残数（対話目的を達成するためのシステムの最終の次発話候補の内容データとなり得る題材データである目的データの残数）、および次発話候補の重要度（次発話準備手段43による準備処理で得られた複数の次発話候補の内容データの各々についての重要度）を記憶するものである。このシステム状態は、現在進行しているユーザとの対話中に得られる一時的な情報（逐次更新されるリアルタイム情報）である。このうち、目的データの残数および次発話候補の重要度は、システム発話意欲度の指標値であるが、システム発話意欲度については、システム発話タイミング検出手段22のユーザ発話権終了判定用閾値調整手段22Gの説明で既に詳述しているので（図10参照）、ここでは詳しい説明を省略する。準備完了・各種の準備中の別を示すステータス、目的データの残数および次発話候補の重要度のいずれについても、次発話準備手段43からネットワーク1を介して送信されてきてシステム状態記憶手段31に記憶されて逐次更新される。

#### 【0183】

ユーザ状態記憶手段32は、ユーザ状態を示す情報として、進行中のユーザ発話についての発話開始からのユーザ発話継続時間、および進行中のユーザ発話についての発話速度を記憶するものである。このユーザ状態は、現在進行しているユーザとの対話中に得られる一時的な情報（逐次更新されるリアルタイム情報）であるため、ユーザ情報記憶手段52に記憶されているユーザの属性情報（複数回の対話を通じて得られた蓄積情報）とは異なる。本実施形態では、ユーザ発話継続時間および発話速度のいずれについても、発話生成手段25により計測され、ユーザ状態記憶手段32に記憶されて逐次更新される。

#### 【0184】

<対話サーバ40 / 音声認識処理手段41の構成>

#### 【0185】

音声認識処理手段41は、音声信号取得手段21により取得したユーザ発話の音声信号をネットワーク1を介して逐次取得し、取得したユーザ発話の音声信号についての音声認識処理を実行し、音声認識処理の結果として得られた言語情報を逐次出力し、出力した言語情報を、対話状態管理手段42を介して次発話準備手段43に逐次送るとともに、ネットワーク1を介してシステム発話タイミング検出手段22、次発話選択用情報生成手段23、次発話選択手段24、および発話生成手段25にも逐次送信する処理を実行するものである。

#### 【0186】

この音声認識処理手段41による音声認識処理は、システム発話タイミング検出手段22によるユーザ発話権の維持・終了を識別するパターン認識処理とは非同期で実行される。

#### 【0187】

具体的には、図6および図7に示すように、音声認識処理手段41は、ショートポーズセグメンテーションと呼ばれる機能により、音声信号取得手段21により取得したユーザ発話の音声信号の中に短時間の無音区間が現れるたびに音声信号を細かく区切り、音声認識対象とする区間を順次確定させていく。これにより、長時間の音声入力を自動的に区切りながら逐次的に音声認識処理を実行することができる。このショートポーズセグメンテーションでの音声認識対象とする音声信号の区間（対応する音声認識処理の時間長は、図6中および図7中の点線で示されている。）は、通常の音声区間検出（VAD）で決定さ

10

20

30

40

50

れる音声区間（対応する音声認識処理の時間長は、図6中および図7中の実線で示されている。）よりも短い。

【0188】

このショートポーズセグメンテーションの機能は、対話サーバ40の音声認識処理手段41の中に設けてもよく、あるいは、音声認識処理手段41により、図示されない外部サーバにアクセスしてストリーミング音声認識を行うようにしてもよい。後者の場合は、例えば、グーグル・クラウド・スピーチAPI（<https://cloud.google.com/speech/>）のストリーミング音声認識等を用いることができる（非特許文献2参照）。

【0189】

<対話サーバ40/対話状態管理手段42の構成>

【0190】

対話状態管理手段42は、進行中のユーザとシステムとの間の対話状態を管理する処理を実行するものである。ここで、システム発話とユーザ発話との時間的な前後関係の説明を容易にするため、図6の最上部に示すように、最初のシステム発話をS(1)、最初のユーザ発話をU(1)とし、以降、S(2)、U(2)、S(3)、U(3)、...、S(N-1)、U(N-1)と対話が進み、直前のシステム発話をS(N)、進行中のユーザ発話をU(N)とし、さらにU(N)を音声区間U(N,1)、U(N,2)、U(N,3)、...に分割し、U(N,K)まで進んでいるものとする。そして、未来の新たな音声区間をU(N,K+1)とする。なお、実際にはショートポーズセグメンテーションにより処理が進行するので、U(N,K+1)よりも細かい区間で、新たな出力が得られる。

【0191】

具体的には、図4に示すように、対話状態管理手段42は、次発話選択手段24からネットワーク1を介して送信されてくる選択結果を受信する処理を実行する。この選択結果は、次発話選択手段24により選択された次発話S(N+1)の内容データまたはその識別情報（例えば、シナリオID、発話節ID等）である。

【0192】

そして、対話状態管理手段42は、次発話選択手段24からの選択結果の受信により、システム発話の開始タイミングが検出されたこと、すなわちユーザ発話権が終了したことを把握することができるので、その時点まで対話状態管理手段42のメモリ（主メモリでよい）で保持していた言語情報、すなわちユーザ発話権が終了したユーザ発話U(N)の発話区間全体の内容データを、対話履歴記憶手段50に記憶させる処理を実行する。この時点の前までには、図6の最上部に示すように、対話履歴記憶手段50には直前のシステム発話S(N)までが保存されているので、これにユーザ発話U(N)が追加されることになる。

【0193】

また、対話状態管理手段42は、次発話選択手段24からの選択結果の受信により、システム発話の開始タイミングが検出されたこと、すなわち次発話S(N+1)の再生が開始されることを把握することができるので、選択結果として受信した次発話S(N+1)の内容データを、対話履歴記憶手段50に記憶させる処理を実行する。これにより、直前のシステム発話S(N)まで保存されていた対話履歴記憶手段50には、ユーザ発話U(N)およびシステム発話S(N+1)が追加されることになる。

【0194】

さらに、対話状態管理手段42は、次発話選択手段24からの選択結果の受信により、システム発話の開始タイミングが検出されたこと、すなわち次発話S(N+1)の再生が開始されることを把握することができるので、さらに次の次発話候補の準備処理を開始させるための準備開始指示情報を、次発話準備手段43に送る処理を実行する。これにより、次発話準備手段43によるシステム発話S(N+2)についての複数の候補の準備処理が開始されることになる。

【0195】

10

20

30

40

50

また、対話状態管理手段 4 2 は、ユーザ発話  $U(N)$  の進行中には、音声認識処理手段 4 1 から逐次出力される音声認識処理の結果としてのユーザ発話の言語情報を逐次受け取り、受け取った言語情報を、次発話準備手段 4 3 に入替要否判断のために逐次送るとともに、対話状態管理手段 4 2 のメモリ（主メモリでよい）に保持する。この際、音声認識処理手段 4 1 から逐次受け取る音声認識処理の結果は、ショートポーズセグメンテーションによる短い区間についての音声認識処理の分であるから、 $U(N, 1)$ 、 $U(N, 2)$ 、 $U(N, 3)$ 、... よりも細かい区間についての音声認識処理の結果である（図 6 参照）。なお、音声認識処理手段 4 1 から受け取った言語情報を、次発話準備手段 4 3 に逐次送る際には、受け取った言語情報だけを送ってもよく、対話状態管理手段 4 2 のメモリに保持しているその時点までのユーザ発話  $U(N)$  の全部の言語情報を送ってもよい。

10

【0196】

また、対話状態管理手段 4 2 は、上述したように、ユーザ発話  $U(N)$  の進行中には、音声認識処理手段 4 1 から逐次受け取った言語情報を、次発話準備手段 4 3 に入替要否判断のために逐次送るが、結果的に、それがユーザ発話  $U(N)$  の発話区間全体における最後の部分（または発話区間全体）であったとしても、次発話準備手段 4 3 に送る。そして、次発話準備手段 4 3 において、ユーザ発話  $U(N)$  の最後の部分を含めて入替要否判断を行い、システム発話  $S(N+1)$  の複数の候補を、新しい複数の候補に入れ替えると判定し、入替の準備を行った場合には、それらの新しい複数の次発話候補（システム発話  $S(N+1)$  の複数の候補）の内容データが、ネットワーク 1 を介して再生装置 20 へ送信され、複数の次発話候補として次発話候補記憶手段 30 に既に記憶されているシステム発話  $S(N+1)$  の複数の候補が、新しい複数の候補に更新される。その後、システム発話タイミング検出手段 2 2 によりシステム発話の開始タイミングが検出されたときには、次発話選択手段 2 4 によりシステム発話  $S(N+1)$  の新しい複数の候補のうちの 1 つが選択され、その選択結果が、ネットワーク 1 を介して対話状態管理手段 4 2 に送信されてくるので、対話状態管理手段 4 2 のメモリに保持されているユーザ発話  $U(N)$  の発話区間全体の内容データ、および選択結果として受信したシステム発話  $S(N+1)$  の内容データを、対話履歴記憶手段 50 に記憶させるとともに、システム発話  $S(N+2)$  の複数の候補を準備するための準備開始指示情報を、次発話準備手段 4 3 に送る。

20

【0197】

なお、上記において、仮に、システム発話  $S(N+1)$  の複数の候補を、新しい複数の候補に入れ替える準備処理に多少時間がかかった場合でも、システム状態記憶手段 31 に記憶されている準備状態を示すステータスが準備中になるので、システム発話の開始タイミングは検出されないことから（図 8 の P 9 参照）、次発話候補記憶手段 30 に記憶されているシステム発話  $S(N+1)$  の複数の候補は更新されないまま保たれ（但し、入替の準備を開始した時点でクリアしてもよい。）、新しい複数の候補への入替を待つことになる。一方、準備処理にかなりの時間がかかった場合には、フィラーが挿入されるが（図 8 の P 10 参照）、この場合も次発話候補記憶手段 30 に記憶されているシステム発話  $S(N+1)$  の複数の候補は更新されないまま保たれ（但し、入替の準備を開始した時点でクリアしてもよい。）、新しい複数の候補への入替を待つことになる。この場合、フィラーの挿入情報は、次発話選択手段 2 4 から対話状態管理手段 4 2 へ送信してもよく、送信しなくてもよいが、送信した場合でも、フィラーの挿入情報を受信した対話状態管理手段 4 2 は、フィラーの挿入を、選択されたシステム発話  $S(N+1)$  として取り扱うわけではないので、システム発話  $S(N+2)$  の準備のための準備開始指示情報を次発話準備手段 4 3 に送る処理は行わない。新しいシステム発話  $S(N+1)$  の準備処理に時間がかかっているため、フィラーを挿入したのに、そのフィラーの挿入をもって、さらに次のシステム発話  $S(N+2)$  の準備を開始するための処理を行うのは不合理だからである。但し、挿入したフィラーの情報を、システム発話  $S(N+1)$  として取り扱うのではなく、システム発話  $S(N+1)$  の準備用繋ぎ発話  $S(N+1: 準備)$  として対話履歴記憶手段 50 に記憶させてもよい。

30

40

【0198】

50

また、対話状態管理手段 4 2 は、対話の開始時には、最初のシステム発話 S ( 1 ) についての準備開始指示情報を次発話準備手段 4 3 に送るが、S ( 1 ) については、複数の候補を準備する必要はないので、次発話準備手段 4 3 は、S ( 1 ) の準備開始指示情報を受け取った場合には、1 つの次発話 ( システム発話 S ( 1 ) ) の内容データを、ネットワーク 1 を介して再生装置 2 0 へ送信し、次発話候補記憶手段 3 0 に S ( 1 ) の内容データを記憶させればよい。S ( 1 ) の内容データを記憶させた時点で、ユーザは発話していないので、すぐにシステム発話タイミング検出手段 2 2 によりシステム発話の開始タイミングが検出され、次発話選択手段 2 4 により S ( 1 ) が選択され、発話生成手段 2 5 により ( 1 ) の再生が開始されることになる。また、次発話選択手段 2 4 により S ( 1 ) が選択されると、その選択結果がネットワーク 1 を介して対話状態管理手段 4 2 へ送信されるので、対話状態管理手段 4 2 は、選択結果として受信したシステム発話 S ( 1 ) を対話履歴記憶手段 5 0 に記憶させるとともに、システム発話 S ( 2 ) の複数の候補の準備のための準備開始指示情報を次発話準備手段 4 3 へ送る。なお、この時点で、ユーザは未だ発話していないので、対話履歴記憶手段 5 0 へのユーザ発話の保存はない。

10

【 0 1 9 9 】

< 対話サーバ 4 0 / 次発話準備手段 4 3 の構成 >

【 0 2 0 0 】

次発話準備手段 4 3 は、システム発話タイミング検出手段 2 2 によるパターン認識処理の周期に依拠しないタイミングで、かつ、システム発話タイミング検出手段 2 2 によりシステム発話の開始タイミングが検出される前に、題材データ記憶手段 5 1 に記憶された題材データまたはネットワーク 1 を介して接続された外部システムである題材データ提供システム 6 0 に記憶された題材データを用いるとともに、ユーザとシステムとの間の対話履歴情報の少なくとも一部および / または音声認識処理手段 4 1 による進行中のユーザ発話についての途中までの音声認識処理の結果を用いて、システムの次発話の内容データ ( 本実施形態では、複数の次発話候補の内容データ ) を取得または生成する準備処理を実行するものである。

20

【 0 2 0 1 】

より詳細には、図 4 に示すように、次発話準備手段 4 3 は、次発話候補初期準備手段 4 3 A と、入替要否判断手段 4 3 B と、入替準備手段 4 3 C と、先行次発話候補情報記憶手段 4 3 D とを含んで構成されている。

30

【 0 2 0 2 】

次発話候補初期準備手段 4 3 A は、対話状態管理手段 4 2 からの準備開始指示情報を受け取ったときに、システムの複数の次発話候補の内容データを取得または生成する準備処理を実行するものである。

【 0 2 0 3 】

入替要否判断手段 4 3 B は、次発話候補初期準備手段 4 3 A により準備した複数の次発話候補の内容データ、または入替要否判断手段 4 3 B 自身により前回準備した複数の次発話候補の内容データを、別の複数の次発話候補の内容データに入れ替えるか否かを判断する処理を実行するものである。

【 0 2 0 4 】

入替準備手段 4 3 C は、入替要否判断手段 4 3 B により入替が必要であると判断した場合に、現在、次発話候補記憶手段 3 0 に記憶されている最新の複数の次発話候補の内容データとは別の複数の次発話候補の内容データを取得または生成する準備処理を実行するものである。

40

【 0 2 0 5 】

先行次発話候補情報記憶手段 4 3 D は、次発話候補初期準備手段 4 3 A や入替準備手段 4 3 C による準備処理を行って得られた複数の次発話候補の内容データ、すなわちネットワーク 1 を介して再生装置 2 0 へ送信し、現在、次発話候補記憶手段 3 0 に記憶されている最新の複数の次発話候補の内容データについての情報 ( 先行情報 ) を記憶するものである。この先行情報は、入替要否判断手段 4 3 B による判断処理を行う際に、先行する複数

50

の次発話候補の内容を把握するために利用される。

【0206】

具体的には、次発話候補初期準備手段43Aは、システム発話S(N+1)の準備開始指示情報を受け取ったときに、対話履歴記憶手段50に記憶されているそれまでの対話履歴情報(システム発話S(N)までの対話履歴情報)の少なくとも一部、すなわちS(1)、U(1)、S(2)、U(2)、...、S(N)(図6の最上部を参照)の少なくとも一部を用いて、題材データ記憶手段51または題材データ提供システム60に記憶されている題材データやその構成要素の中から、システムの複数の次発話候補S(N+1)の内容データを選択取得する処理を実行する。但し、最初のシステム発話S(1)の準備開始指示情報を受け取ったときには、選択取得するS(1)は、1つだけでよい。また、常に複数の次発話候補を選択取得しなければならないわけではなく、選択取得した次発話候補が、結果的に1つになる場合があってもよい。

10

【0207】

一方、入替準備手段43Cは、次発話候補記憶手段30に複数の次発話候補S(N+1)の内容データが既に記憶されている状態において、入替要否判断手段43Bにより入替が必要であると判断した場合に、対話状態管理手段42から逐次送られてくる音声認識処理手段41による現在進行中のユーザ発話U(N)の音声認識処理の結果である言語情報を用いて、題材データ記憶手段51または題材データ提供システム60に記憶されている題材データやその構成要素の中から、システムの別の複数の次発話候補S(N+1)の内容データを選択取得する処理を実行する。但し、常に複数の次発話候補を選択取得しなければならないわけではなく、選択取得した次発話候補が、結果的に1つになる場合があってもよいのは、上述した次発話候補初期準備手段43Aの場合と同様である。また、上記の説明では、現在進行中のユーザ発話U(N)の音声認識処理の結果を用いるとしているが、準備処理を行う時点では、対話状態管理手段42から逐次送られてくる音声認識処理の結果がユーザ発話U(N)の発話区間全体における最後の部分(または発話区間全体)であるか否かは判らない場合があるので(あるいは、最後の部分であるか否かの区別をする必要はないので)、対話状態管理手段42の説明で既に詳述したように、対話状態管理手段42から受け取った音声認識処理の結果が、結果的に、ユーザ発話U(N)の発話区間全体における最後の部分(または発話区間全体)であった場合でも、この入替準備手段43Cによる準備処理は実行される。この場合、結果的には、進行中のユーザ発話U(N)ではなく、対話履歴情報として対話履歴記憶手段50に記憶されることになる発話終了後のユーザ発話U(N)の情報を用いていることになる。なお、入替準備手段43Cは、ユーザ発話U(N)だけではなく、次発話候補初期準備手段43Aの場合と同様に、対話履歴情報S(1)、U(1)、...、S(N)を用いてもよい。

20

30

【0208】

また、用意されている題材データには、様々な状態のものがあ、例えば、テキストデータだけの場合、テキストデータおよびそれに対応する音声データがある場合、それらのテキストデータや音声データに、映像データや静止画データ、あるいは楽曲データが付随している場合、付随させる映像データや静止画データ、あるいは楽曲データだけの場合等があり、更にはテキストデータにも様々な語調のものがある。このため、次発話候補初期準備手段43Aおよび入替準備手段43Cは、必要な場合には、テキストデータの加工調整(例えば、語尾の調整、一部削除、結合・分割・組替・その他の編集等)、テキストデータから音声データ(例えばwavファイル等)を生成する音声合成、動画や静止画の画質・サイズ調整といった各種の生成処理も行う。但し、システムの応答性を向上させる観点からは、次発話候補の準備処理自体に時間がかかることを避ける必要がある。準備処理は、ユーザ発話中に行い、原則として、ユーザ発話が終了する前に準備が完了していることが好ましいからである。従って、テキストデータの語調等の加工調整、音声合成処理、動画や静止画の画質・サイズ調整等は、予め実行しておき、それらの処理を実行済の題材データを、題材データ記憶手段51または題材データ提供システム60に用意しておくことが好ましい。

40

50



## 【0209】

そして、次発話候補初期準備手段43Aによるシステム発話S(N+1)の複数の候補の選択取得では、通常は、直前のシステム発話S(N)の内容が最も重要な選択用判断材料となるが、それよりも前のS(1)、U(1)、...、S(N-1)、U(N-1)も使用されることがある。例えば、シナリオデータ内における各構成要素(本実施形態では、主計画要素、副計画要素がある。)を予め定めた順序で再生していくときに、ユーザ発話の内容に応じて各構成要素の再生順序を変更する場合がある。この場合、例えば、1回再生した構成要素については、2度目の再生は行わないというルールがあれば、それまでにいずれの構成要素が再生されたのかを把握する必要があるので、それまでの対話履歴情報の全部を使用する必要がある。

10

## 【0210】

また、入替準備手段43Cによるシステム発話S(N+1)の複数の候補の選択取得では、ユーザ発話U(N)が、例えば「さっき言っていたXXX選手について、別の情報が知りたいな。」「その話は知っているから、さっきのYYY事件の話を詳しく聞きたいな。」等であれば、選択用判断材料として、そのユーザ発話U(N)の情報を使用することは勿論であるが、直前のシステム発話S(N)を使用せず、それよりも前の情報S(1)、U(1)、...、S(N-1)、U(N-1)を使用する場合もある。つまり、少し前(例えば数分前等)の対話履歴情報に基づき、XXX選手やYYY事件について、どこまで話していたのかを把握し、それとは別の情報を、題材データ記憶手段51または題材データ提供システム60から選択取得する場合等がある。

20

## 【0211】

また、題材データ記憶手段51または題材データ提供システム60に記憶されている題材データには、様々な種類のデータがあり、情報量の多少も異なっている。例えば、題材データがシナリオデータであれば、複数の構成要素により構成され、一方、題材データの中には、シナリオデータ内の1つの構成要素に相当するような比較的短い題材データも存在する。従って、次発話候補初期準備手段43Aや入替準備手段43Cにより複数の次発話候補を「選択」することには、複数(多数)のシナリオデータの中から1つのシナリオデータを選択し、かつ、選択した1つのシナリオデータの中から1つ(S(1)の場合)または複数の構成要素を選択すること、既に選択されている1つのシナリオデータの中から複数の構成要素を選択すること、複数(多数)の比較的短い題材データの中から1つ(S(1)の場合)または複数の題材データを選択すること等が含まれる。

30

## 【0212】

また、次発話候補初期準備手段43Aおよび入替準備手段43Cは、準備処理で得られた複数(結果的に1つの場合もあり、また、S(1)の場合は1つである。)の次発話候補の内容データまたはそれらに加えてそれらの識別情報(例えば、シナリオID、発話節ID等)を、ネットワーク1を介して再生装置20へ送信し、次発話候補記憶手段30に記憶させる処理も実行する。2回目以降は、更新処理である。この更新により、次発話候補記憶手段30に記憶されている次発話候補の内容データの数は、ユーザ発話の進行に伴って、例えば、図7中の中央部に示すように、N1、N2、N3のように変化する。また、入替要否判断手段43Bにより入替が必要であるという判断結果が出た場合に、入替準備手段43Cによる準備処理が開始されるが、この準備期間中は、図7中の中央部に示すように、次発話候補記憶手段30に記憶されている複数の次発話候補の内容データを削除し、次発話候補の内容データの数をゼロにクリアしてもよく、あるいは、削除せずに維持し、ゼロにクリアしない処理を行ってもよい。

40

## 【0213】

さらに、次発話候補初期準備手段43Aおよび入替準備手段43Cは、準備状態を示すステータス、目的データの残数および次発話候補の重要度(準備処理で得られた複数の次発話候補の内容データの各々の重要度)を、ネットワーク1を介して再生装置20へ送信し、システム情報記憶手段31に記憶させる処理も実行する。2回目以降は、更新処理である。

50

## 【 0 2 1 4 】

なお、目的データの残数は、対話目的を達成するためのシステムの最終の次発話候補の内容データとなり得る目的データの残数であるが、次発話候補記憶手段 30 に記憶させる次発話候補の内容データの数は異なる。例えば、情報検索対話で、ユーザが自分の利用する飲食店を探すときには、飲食店のデータが目的データとなる。しかし、条件提示による絞り込みが進んでいない段階では、目的データ（例えば、飲食店のデータ等）は多数存在し、それらの目的データの全部を、次発話候補の内容データとして次発話候補記憶手段 30 に記憶させるわけではなく、情報検索対話の初期の段階や途中の段階では、次発話候補記憶手段 30 には、「何を食べたいですか?」、「費用はどれくらいですか?」等が記憶されるだけである。そして、絞り込みが進んだ段階や絞り込みが完了した最終段階で、目的データ（例えば、飲食店のデータ等）は、次発話候補の内容データとして次発話候補記憶手段 30 に記憶されることになる。従って、目的データの残数は、潜在的な次発話候補の内容データの数である。

10

## 【 0 2 1 5 】

先行次発話候補情報記憶手段 43D には、例えば、次発話候補初期準備手段 43A および入替準備手段 43C による準備処理で得られた複数の次発話候補の内容データ（現在、次発話候補記憶手段 30 に記憶されている最新の複数の次発話候補の内容データ）、それらの内容データについての各分野（例えば、IT・科学、テニス、野球等）、分野以外の属性（例えば、男性向け、10代～30代向け等）、それらの内容データに含まれる1つまたは複数の重要度の高い単語等が記憶されている。

20

## 【 0 2 1 6 】

入替要否判断手段 43B は、対話状態管理手段 42 から逐次送られてくる音声認識処理手段 41 によるユーザ発話 U(N) の音声認識処理の結果である言語情報を受け取り、受け取った言語情報と、先行次発話候補情報記憶手段 43D に記憶されている再生装置 20 へ送信済の複数の次発話候補の内容データ（現在、次発話候補記憶手段 30 に記憶されている最新の複数の次発話候補の内容データ）についての情報（先行情報）とを用いて、次発話の候補となる複数の次発話候補の内容データの少なくとも一部を入れ替えるか否かを逐次判定し、入れ替えると判定した場合には、その結果を入替準備手段 43C に送る処理を実行する。

## 【 0 2 1 7 】

具体的には、入替要否判断手段 43B は、現在までに、図 6 に示す U(N, K) までの音声認識処理の結果を用いた入替要否判断処理およびそれに伴う入替準備処理が行われていたとすると、例えば、新たに出力された U(N, K + 1)（但し、ショートセグメンテーションであるから、正確には、その一部）の音声認識処理の結果である言語情報の中に、重要度の高い単語が含まれているか否かを判断する。ここで、単語の重要度としては、例えば、TF (Term Frequency: 文書における単語の出現頻度) および IDF (Inverse Document Frequency: 逆文書頻度) による TF - IDF、Okapi - BM25 等を採用することができ、予め算出して単語重要度記憶手段（不図示）に記憶しておけばよい。

30

## 【 0 2 1 8 】

そして、例えば、U(N, K + 1)（正確には、その一部）の中に、単語  $w_i$  が含まれていたとすると、これらの単語  $w_i$  の全ての重要度が、予め定めた重要度判定用閾値以下または未満であった場合（単語  $w_i$  がいずれも重要度の高い単語ではなかった場合）には、入替は不要であると判断すること等ができる。一方、単語  $w_i$  の中に重要度の高い単語が含まれていた場合には、その重要度の高い単語（例えば単語  $w_j$ ）が、先行次発話候補情報記憶手段 43D に記憶されている重要度の高い単語の中に含まれているか否かを判断し、含まれていれば、入替は不要であると判断すること等ができる。あるいは、単語  $w_i$  の中に重要度の高い単語が含まれていた場合には、その重要度の高い単語（例えば単語  $w_j$ ）と、先行次発話候補情報記憶手段 43D に記憶されている重要度の高い単語の各々の類似度を、例えば word2vec や GloVe 等により求め、求めた各類似度のうちのいずれかが類似度判定用閾値以上または超過であった場合

40

50

(類似する重度語の高い単語があった場合)には、入替は不要であると判断すること等ができる。また、上記のように単語 , , の中の重要度の高い単語(例えば単語 )と、先行次発話候補情報記憶手段43Dに記憶されている1つまたは複数の重要度の高い単語とを用いて判断するのではなく、単語 , , の中の重要度の高い単語(例えば単語 )と、先行次発話候補情報記憶手段43Dに記憶されている複数の次発話候補の内容データの全体(それらに含まれる全ての単語)とを用いて判断してもよい。

【0219】

また、上記の例の単語 , , の中の重要度の高い単語(例えば単語 )が、いずれの分野に属する単語であるかを判断し、判断した分野が(複数の分野でもよく、その場合には、いずれかの分野が)、先行次発話候補情報記憶手段43Dに記憶されている1つまたは複数の分野(通常は1つの分野であることが多い。)の中に含まれていれば、入替は不要であると判断すること等ができる。なお、各単語(重要度の高い単語)と各分野(例えば、IT・科学、テニス、ゴルフ、エンタメ、政治経済、国際等)との対応関係は、予め定めて単語帰属分野記憶手段(不図示)に記憶しておけばよく、1つの単語が複数の分野に帰属していてもよい。この対応関係は、例えば、各分野の文書における各単語の出現頻度や、累積出現回数等により定めることができる。

【0220】

なお、題材データまたはその構成要素には、分野の識別情報が関連付けられている。分野の粒度は、システム設計者が適宜定めればよく、例えば、テニス、ゴルフ、野球等を別々の分野とするか、スポーツで1つの分野にまとめるか、あるいは、政治、経済を別々の分野とするか、1つにまとめるか等は任意である。1つの題材データまたはその構成要素は、複数の分野に帰属していてもよい。また、題材データまたはその構成要素が、女子プロゴルフの話題のみである場合に、例えば、女子プロゴルフ<ゴルフ<スポーツのように、包含関係にある分野の識別情報を全て関連付けるようにしてもよい。

【0221】

また、音声対話には、各種の目的の対話(例えば、ニュース対話、アンケート対話、ガイド対話、情報検索対話、操作対話、教育対話、情報特定対話等)があり、対話の進行も各種のタイプのものがある。対話の進行のタイプとの関係では、次のようになる。

【0222】

次発話候補初期準備手段43Aは、シナリオデータ(主計画および副計画を有する複雑な分岐を行うシナリオに限らず、より単純なシナリオも含む。)があり、シナリオとして予め定められた順序に従って対話を進めていく場合には、そのシナリオの順序に従って、複数の次発話候補を選択していく。この場合、入替準備手段43Cにより、予め定められた順序が変更された場合には、その変更を反映させ、例えば、1回再生したシナリオ構成要素については、2度目の再生は行わないというルールがあれば、そのルールに従いつつ、当初の順序をなるべく維持した順序で複数の次発話候補を選択していく。

【0223】

また、シナリオデータがなく、対話の進行や分岐のパターンが予め定まっているわけではないが、システム発話の内容については、予定外の情報を外部システムから取得しなければならない場合を除き、予め用意されていて、毎回のユーザ発話の内容に従って、その都度、次のシステム発話の内容を定める場合がある。このような場合には、直前のシステム発話S(N)で、次のシステム発話S(N+1)の複数の候補が定まることは少ない。なぜなら、S(N)でS(N+1)の候補が定まるということは、結局、広い意味で、または部分的にシナリオが形成されていると考えることができるので、シナリオがない場合に該当しないからである。

【0224】

例えば、自動車のカーナビ操作のための操作対話において、システム発話S(N) = 「住所で目的地を設定しますか?」であった場合、次発話候補初期準備手段43Aにより、システム発話S(N+1) = 「最初に都道府県を教えてください。」、「市町村を教えてください。」、「何丁目何番地ですか?」等を準備して次発話候補記憶手段30に記憶さ

10

20

30

40

50

せておく。そして、ユーザ発話  $U(N)$  が「はい。」であれば、次発話選択手段 24 により「最初に都道府県を教えてください。」を選択して発話生成手段 25 によりそれを再生し、「はい、東京都です。」であれば、「市町村を教えてください。」を選択して再生し、「はい、東京都新宿区です。」であれば、「何丁目何番地ですか？」を選択して再生する。この際、ユーザ発話  $U(N)$  が「はい、東京都新宿区です。」の途中の「はい、東京都...」まで進行した段階で、入替準備手段 43C により、システム発話  $S(N+1) =$  「市町村を教えてください。」、「何丁目何番地ですか？」等への入替が行われる場合（「最初に都道府県を教えてください。」が次発話候補から除かれている場合）もある。このような場合、部分的にシナリオが形成されていると考えることができ、そのシナリオに従って次発話候補初期準備手段 43A による準備処理が行われているが、入替準備手段 43C による役割も大きい。

10

#### 【0225】

一方、シナリオデータがない場合は、次発話候補初期準備手段 43A による準備処理よりも、入替準備手段 43C による準備処理が中心となる。最初はシナリオがあり、その後、フリートークに近い状態になる場合の後半の処理も同様である。そして、シナリオデータがない場合、次発話候補初期準備手段 43A は、対話状態管理手段 42 からの  $S(N+1)$  の準備開始指示情報を受け取ったときに、システム発話  $S(N+1)$  の複数の候補を定めることができなければ、準備中のステータス（ステータス＝次発話候補検討中）を、ネットワーク 1 を介して再生装置 20 へ送信してシステム状態記憶手段 31 に記憶させ、入替準備手段 43C に準備処理を任せることができる。このようにした場合は、入替準備手段 43C による準備処理は、入替の準備というより初期データの準備となるので、入替準備手段 43C は、入替要否判断手段 43B からの入替が必要であるという判断結果を受け取って準備処理を開始するのではなく、次発話候補初期準備手段 43A から転送されてくる準備開始指示情報を受け取って準備処理を開始することになる。この場合、入替要否判断手段 43B による判断処理は行われないので、入替準備手段 43C は、対話状態管理手段 42 から逐次送られてくる音声認識処理の結果を受け取り、入替要否判断手段 43B による重要度の高い単語の抽出処理に相当する処理を実行するが、この際の重要度判定用閾値は低く設定してもよい。なお、重要度判定用閾値を低く設定しても、重要度の高い単語が抽出されない場合には、進行中のユーザ発話  $U(N)$  の中に、未だシステム発話  $S(N+1)$  の複数の候補の決定をするのに十分な情報（単語）が現れていないことになるので、ユーザ発話  $U(N)$  の進行を待つことになる。入替準備手段 43C は、以上のような次発話候補初期準備手段 43A から転送されてくる準備開始指示情報を受け取った場合の初期データの準備処理を行い、複数の次発話候補の内容データを次発話候補記憶手段 30 に記憶させた後には、通常通りの入替の準備処理（重要度判定用閾値も通常の設定とする。）を実行する。

20

30

#### 【0226】

また、シナリオデータがない場合は、フリートークの状態に近いと考え、次発話候補初期準備手段 43A は、とりあえず、題材データ記憶手段 51 や題材データ提供システム 60 に記憶されている題材データの中からランダムに選択取得した複数の次発話候補の内容データを、ネットワーク 1 を介して再生装置 20 へ送信して次発話候補記憶手段 30 に記憶させてもよい。ランダムな選択取得を行っても、その後、ユーザ発話  $U(N)$  が進行すると、入替準備手段 43C による準備処理が行われ、複数の次発話候補の内容データが適切なものに入れ替えられる。仮に、ユーザ発話  $U(N)$  が進行しても、ランダムに選択取得した複数の次発話候補の内容データがそのまま維持されていたとすると、そのランダムな選択取得が適切であったということになる。なお、ランダムに選択取得する題材データが存在する（選択取得する範囲が定まっている）ということは、完全なフリートークではなく、システム発話の内容は、想定範囲外の情報を外部システムから取得しなければならない場合を除き、予め用意されていることになる。

40

#### 【0227】

前述したように、入替要否判断手段 43B は、対話状態管理手段 42 から逐次送られて

50

くる新たな音声認識処理の結果を受け取り、この結果に含まれる単語のうち予め定められた重要度の高い単語を抽出する処理を実行するので、入替準備手段43Cは、入替要否判断手段43Bから、抽出された重要度の高い単語を受け取る。そして、入替準備手段43Cは、この重要度の高い単語を用いて、予め定められた各単語と各分野との対応関係（単語帰属分野記憶手段（不図示）に記憶されている情報）から、ユーザの関心のある話題（分野）を決定し、題材データ記憶手段51または題材データ提供システム60に記憶されている題材データの中から、決定した話題（分野）に関連付けられて記憶されている題材データを選択し、次発話の候補となる別の複数の次発話候補の内容データを取得または生成する準備処理を実行することができる。

#### 【0228】

例えば、システム発話S(N)が「早稲田太郎選手が4回転フリップを成功させたよ。」であり、システム発話S(N+1)の複数の候補として、「グランプリシリーズのカナダ大会で跳んだそうだ。」（主計画要素）、「早稲田太郎選手は、...」という早稲田太郎の人物の説明データ（副計画要素）、「4回転フリップっていうのは、...」という4回転フリップの技の説明データ（副計画要素）が、次発話候補記憶手段30に記憶されているとする。このとき、ユーザ発話U(N)が「フィギュアスケートは興味がないので、野球の話が聞きたいんだけど...」、「つまらない、野球の方がおもしろいから...」であった場合には、入替準備手段43Cは、再生中のシナリオデータ（分野=アイススケート、または、分野=スポーツ、アイススケート）の中に野球の話は全くないので、シナリオデータ自体を別の分野（この例では、野球の分野）に入れ替え、その入替後のシナリオデータ内の先頭の構成要素を、次発話候補とすることができる。なお、この場合は、S(N+1)の候補ではあるが、シナリオデータの先頭からの再生となるので、S(1)と同等であるから、次発話候補は1つでよい。また、この場合、「興味がない」、「つまらない」、「もう飽きた」、「くだらない」、「話題を変えてほしい」、「その話はもういい」、「その話はやめて」、「ところで」、「話は変わるけど」、「そういえば」等の話題転換要求を伴っているため、シナリオデータ自体を入れ替える話題転換処理を行っている。また、明確な話題転換要求が無くても、例えば、ユーザ発話U(N)が「来週、日米野球があるけど、高校野球の...」のように、再生中のシナリオデータ内にない単語が繰り返される場合にシナリオデータ自体を入れ替える話題転換処理を行ってもよい。

#### 【0229】

一方、上記の例において、ユーザ発話U(N)が「早稲田次郎選手の方が好きなんですけど、早稲田次郎選手の成績は...」であった場合には、入替準備手段43Cは、再生中のシナリオデータ内に早稲田次郎選手についての構成要素（主計画要素）も含まれているので、シナリオデータの入替は行わずに、同じシナリオデータ内での再生順序の変更を行う。例えば、「早稲田次郎」によるキーワードマッチングで「早稲田次郎選手は、4回転アクセルに挑戦したけど失敗したんだ。」（主計画要素）を選択するとともに、「早稲田次郎選手は、...」という早稲田次郎の人物の説明データ（副計画要素）、「4回転アクセルっていうのは、...」という4回転アクセルの技の説明データ（副計画要素）を選択し、次発話候補とする。

#### 【0230】

また、次発話候補初期準備手段43Aや入替準備手段43Cによる準備処理において、位置データや時刻データを反映させてもよい。例えば、博物館や遺跡等の案内を行うガイド対話では、ユーザの位置データ（例えば、再生装置20に設置されたGPS受信機や、再生装置20が本体と端末とに分割されている場合の端末に設置されたGPS受信機で得られる位置データ等）を用いて、複数の次発話候補の内容データが定まるようにしてもよい。例えば、博物館のガイド対話において、予め登録されて対話サーバ40のメモリに記憶されている展示物Xの位置データと、再生装置20からネットワーク1を介して送信されてきたユーザの位置データとを用いて、ユーザが展示物Xのそばに近づいたことを検出し、さらに時刻データを用いて12時近くであることを検出した時点で、「そろそろ展示物Xが見えてきます。」というシステム発話S(N)を行い、次発話候補S(N

10

20

30

40

50

+ 1)として「展示物 X は、・・・」、「食堂をご案内しましょうか。」等を用意して次発話候補記憶手段 30 に記憶させておく。そして、ユーザ発話 U ( N ) が「展示物 X の説明を聞きたいな。」であった場合には、「展示物 X は、・・・」を選択して再生し、「お腹すいた。」であった場合には、「食堂をご案内しましょうか。」を選択して再生する。また、時刻データを用いて、「外が暗くなってきたから、そろそろ帰り支度を始めましょう。」、「閉館時間が迫っているから、素早く回ろうね。」等を準備して次発話候補記憶手段 30 に記憶させておくこともできる。

#### 【 0 2 3 1 】

さらに、次発話候補初期準備手段 4 3 A や入替準備手段 4 3 C による準備処理において、位置データや時刻データ以外の状態データ（変化の速度の大小の相違はあるが、原則として、時々刻々と変化するデータ）、例えば、温度データ、湿度データ、天候データ、高度データ等を用いて次発話候補の内容データを準備してもよい。例えば、「今日は暑いね。」、「今日は蒸すね。」、「今日は天気がいいね。」、「空気が薄くなってきたけど、大丈夫？」等の次発話候補を準備し、次発話候補記憶手段 30 に記憶させておくことができる。なお、上記の例の天候データは、選択用判断材料としての天候データ（ユーザが操作する再生装置 20 の所在地における晴・雨・曇り等のデータ）であるから、題材データとして用意されている「台風 28 号が沖縄地方に接近しています。」、「XX 地方に大雨洪水警報が出ていますので、YY 川の氾濫に注意してください。」等の警報データとは異なる。つまり、例えば、雨という天候データに基づき、「今日は雨だけど、東京ドームは屋根があるから、野球の観戦はできるよ。」等の題材データが選択取得され、次発話候補記憶手段 30 に記憶されることになる。

#### 【 0 2 3 2 】

また、入替準備手段 4 3 C による準備処理は、必ずしも複数の次発話候補の全部を入れ替える必要はなく、少なくとも一部の入替が行われればよい。例えば、最初の複数の次発話候補の内容データ（初期データ）または前回の入替後のデータが、次発話候補 A , B , C であったとすると、入替後の次発話候補は、D , E , F のように全部が入れ替わっていてもよく、A , B , D のように一部が入れ替わっていてもよい。また、入替後の次発話候補は、A , B , C , D , E のように候補が追加されて増えた状態となってもよく、A , B のように一部削除された状態となってもよい。

#### 【 0 2 3 3 】

さらに、対話状態管理手段 4 2 は、ユーザ発話 U ( N ) の進行中において、各時点において、それまでの U ( N ) の全部を保持しているので、対話状態管理手段 4 2 から入替要否判断手段 4 3 B に送る入替要否判断に用いるための音声認識処理の結果の長さは、自在に調整することができる。従って、ショートポーズセグメンテーションの単位の最新の音声認識結果だけとしてもよく、最新の音声認識結果を含めたある程度の時間長の音声認識結果としてもよく、対話状態管理手段 4 2 に保持されている U ( N ) の音声認識結果の全部としてもよい。

#### 【 0 2 3 4 】

< 対話サーバ 40 / 対話履歴記憶手段 50 の構成 >

#### 【 0 2 3 5 】

対話履歴記憶手段 50 は、システムとユーザとの間の対話履歴情報を記憶するものである。具体的には、図 6 の最上部に示すように、システム発話 S ( 1 ) の内容データ（テキストデータ）、ユーザ発話 U ( 1 ) の内容データ（テキストデータ）、同様に、S ( 2 )、U ( 2 )、S ( 3 )、U ( 3 )、... の各内容データ（テキストデータ）を、対話の順番に記憶する。ユーザ発話から始まってもよい。なお、進行中のユーザ発話 U ( N ) は、本実施形態では、対話状態管理手段 4 2 のメモリ（主メモリでよい）に記憶され、発話の終了後に、発話区間全体が対話履歴記憶手段 50 に記憶される。

#### 【 0 2 3 6 】

< 対話サーバ 40 / 題材データ記憶手段 51 の構成 >

#### 【 0 2 3 7 】

題材データ記憶手段 5 1 は、題材データを記憶するものである。題材データは、例えば、シナリオデータ（主計画および副計画を有する複雑な分岐を行うシナリオに限らず、より単純なシナリオも含む。）、シナリオが形成されていない最近のトピックを集めた各種の話題データの集合（但し、話題データの 1 つ 1 つが、独立した題材データであり、それぞれ比較的短いデータである。）、辞書データ、事典データ、機器の使用方法や施設等のガイダンス用データ、アンケート調査用データ、機器や装置等の操作補助用データ、教育用データ等である。これらの題材データまたはその構成要素には、分野（例えば、IT・科学、政治・経済、国際、エンタメ、相撲、ゴルフ等）の識別情報が関連付けられている。なお、分野が定められていない題材データまたはその構成要素が混在していてもよいが、その場合は、キーワードマッチングにより、必要な情報を選択取得する。題材データ提供システム 6 0 も同様である。

10

【0238】

&lt;対話サーバ 4 0 / ユーザ情報記憶手段 5 2 の構成&gt;

【0239】

ユーザ情報記憶手段 5 2 は、ユーザ発話とシステム発話との衝突の発生情報、システムの交替潜時、およびユーザの発話速度を、ユーザ識別情報と関連付けて記憶するものである。このユーザ情報記憶手段 5 2 に記憶される情報は、各ユーザとの複数回の対話を通じて得られたユーザ毎の蓄積情報であるから、ユーザの属性情報である。従って、ユーザとの対話中における一時的な情報ではないので、ユーザ状態記憶手段 3 2 に記憶される情報とは異なる。これらの衝突の発生情報、システムの交替潜時、およびユーザの発話速度は、いずれも発話生成手段 2 5 により得られて記録されたものである。

20

【0240】

&lt;ユーザからシステムへの話者交替時の処理の流れ：図 5 &gt;

【0241】

このような本実施形態においては、以下のようにしてユーザからシステムへの話者交替が行われる。

【0242】

図 5 において、先ず、対話開始前に、システム発話タイミング検出手段 2 2 のユーザ発話権終了判定用閾値調整手段 2 2 G（図 2 参照）により、ユーザ情報記憶手段 5 2 から、対話相手のユーザについてのユーザ識別情報（ユーザ ID）を用いて、衝突の発生情報（蓄積情報）、交替潜時（蓄積情報）、および発話速度（蓄積情報）を取得し、衝突の発生情報（蓄積情報）によるユーザ発話権終了判定用閾値の事前調整（図 1 1 参照）、発話速度（蓄積情報）による下方調整用閾値の事前調整およびその下方調整用閾値を用いた交替潜時（蓄積情報）によるユーザ発話権終了判定用閾値の事前調整（図 1 2 参照）を行う（ステップ S 1）。

30

【0243】

次に、対話開始後においては、音声信号取得手段 2 1 により取得したユーザの音声信号を用いて、システム発話タイミング検出手段 2 2 の音響特徴量抽出手段 2 2 A（図 2 参照）により、周波数分析等を行って音響特徴量を抽出する（ステップ S 2）。また、必要な場合には、音声認識処理手段 4 1 により得られた音声認識処理の結果である言語情報を用いて、言語特徴量抽出手段 2 2 B（図 2 参照）により言語特徴量を抽出する。

40

【0244】

続いて、システム発話タイミング検出手段 2 2 のユーザ発話権終了判定用閾値調整手段 2 2 G（図 2 参照）により、ユーザ状態記憶手段 3 2 からユーザ発話継続時間を取得し、ユーザ発話権終了判定用閾値のリアルタイム調整（図 9 参照）を行うとともに、システム状態記憶手段 3 1 からシステム発話意欲度の指標値（目的データの残数および / または次発話候補の重要度）を取得し、ユーザ発話権終了判定用閾値のリアルタイム調整（図 1 0 参照）を行う（ステップ S 3）。なお、これらの 2 種類のリアルタイム調整は、いずれか一方の調整を行ってもよく、双方の調整を行ってもよく、双方の調整を行う場合は、どちらの調整を先に行ってもよい。

50

## 【0245】

それから、システム発話タイミング検出手段22のユーザ発話権終了判定用パターン認識器22C(図2参照)により、ステップS2で抽出した音響特徴量、またはこれに加えて言語特徴量を用いて、ユーザ発話権の維持または終了を識別するパターン認識処理を実行し、このパターン認識処理で得られる尤度を用いたユーザ発話権終了判定用閾値による閾値判定を行い、維持または終了の識別結果を出力する(ステップS4)。

## 【0246】

また、これと並行して、次発話選択用情報生成手段23により、ステップS2で得られた音響特徴量を用いた韻律分析を行い、ユーザ発話意図の識別処理を行う(ステップS5)。なお、図5中の2点鎖線で示すように、音声信号取得手段21により取得したユーザの音声信号を用いて、ステップS2とは別途に韻律特徴量を抽出し、その韻律特徴量を用いた韻律分析を行い、ユーザ発話意図の識別処理を行ってもよい。

## 【0247】

そして、前述したステップS4の識別結果が維持または終了のいずれであるかを判断し(ステップS6)、維持であった場合には、ステップS2の処理に戻る。一方、識別結果が終了であった場合には、システム発話タイミング検出手段22のシステム発話開始タイミング判断手段22F(図2参照)により、システム状態記憶手段31から、準備の状態を示すステータス(準備完了・各種の準備中の別)を取得し、図8に示す流れで、システム発話の開始タイミングであるか否かの判断結果を出力する(図5のステップS7)。

## 【0248】

ここで、出力された判断結果が、システム発話の開始タイミングであるか否かに応じ(ステップS8)、システム発話の開始タイミングではなかった場合には、ステップS2の処理に戻る。一方、システム発話の開始タイミングであった場合(当該タイミングが検出された場合)には、次発話選択手段24により、前述したステップS5で次発話選択用情報生成手段23により得られたユーザ発話意図の識別結果(質問、相槌等の別)および/または音声認識処理手段41による音声認識処理の結果である言語情報(文字列)を用いて、次発話候補記憶手段30に記憶されている複数(但し、1つである場合もある)の次発話候補の内容データの中から、次発話の内容データを選択するとともに、選択した次発話の内容データまたはその識別情報(シナリオID、発話節ID等)を、ネットワーク1を介して対話サーバ40の対話状態管理手段42へ送信する(ステップS9)。

## 【0249】

このステップS9の処理を行う際、次発話候補記憶手段30には、図5の流れとは非同期で行われる次発話準備手段43の準備処理により、次発話選択手段24によるステップS9の処理に先んじて用意された複数(但し、1つである場合もある)の次発話候補の内容データが既に記憶されている状態である。なお、準備中の場合には、フィラーの挿入が行われる(図8のP10参照)。

## 【0250】

また、次発話選択手段24による選択結果が、ネットワーク1を介して対話状態管理手段42へ送信されると、対話履歴記憶手段50に記憶された対話履歴情報が更新されるとともに、次発話準備手段43に対してさらにその次の次発話候補の準備開始指示情報が出され、対話履歴記憶手段50の情報をういたさらにその次の次発話候補の準備処理が開始されるので、この意味では、図5の流れと、次発話準備手段43の処理とは一連の流れのように見える。しかし、図4に示すように、次発話準備手段43は、図5の中心であるシステム発話タイミング検出手段22の処理とは非同期で行われる音声認識処理手段41による音声認識処理の結果(言語情報)を用いた準備処理を実行するので、結局、次発話準備手段43の処理は、図5の一連の処理の流れの中に記載することはできない。

## 【0251】

その後、発話生成手段25により、次発話選択手段24により選択された次発話の内容データを用いて、システム発話の音声信号の再生処理が行われる(図5のステップS10)。また、付随する映像データ、静止画データ、楽曲データがあれば、それらの再生処理



も行われる。なお、必要な場合には、ここでの音声合成処理も行われるが、次発話候補記憶手段30に記憶された次発話候補の内容データには、音声合成処理で得られた音声データも含まれていることが好ましい。

【0252】

そして、対話終了であるか否かを判断し（ステップS11）、対話終了でない場合には、ステップS2の処理に戻る。

【0253】

<シナリオのデータ構成、シナリオの再生、および次発話候補の準備処理の流れ：図13～図17>

【0254】

<シナリオのデータ構成：図13>

【0255】

図13には、主計画および副計画からなるシナリオのデータ構成の具体例が示されている。このようなシナリオデータは、非特許文献4に記載されたシナリオデータと同様のものであり、対話システム10で利用することができる題材データの一種として題材データ記憶手段51に記憶されている。

【0256】

より詳細には、このシナリオデータは、ニュースやコラムや歴史等の各種の話題を記載した記事データ（文書データ）から生成したものであり、元の文書データを構成する要素を、元の文書データの内容の要約となる主計画要素と、この主計画要素を補完する副計画要素と、これら以外の要素（省略要素）とに分割し、これらの3種類の要素のうちの主計画要素および副計画要素、並びに、発話計画情報（主計画要素の再生順序および副計画要素への分岐を定めた情報）を含むように構成したものである。なお、元の文書データからのシナリオ生成時に、結果的に省略要素が発生しなくてもよい。つまり、主計画要素を除いた残り全ての要素が、副計画要素に割り当てられてもよい。

【0257】

主計画要素は、元の文書データを要約し、口語化することにより生成される。文書の要約は、重要文抽出、整列、文圧縮の処理を経て行われる。まず、重要文抽出で、文書の要点となる情報を文単位で大まかに抽出する。次に、整列を行い、抽出した重要文の提示順序を決定する。そして、文圧縮を行い、文自体を短く縮約する。最後に、口語化処理を行い、書き言葉から会話表現への書き換えを行う。なお、このシナリオ生成時における重要文抽出で考慮される文の重要度は、前述したシステム発話タイミング検出手段22のユーザ発話権終了判定用閾値調整手段22Gの説明で既に詳述した通り、システム状態記憶手段31に記憶される次発話候補の内容データの重み度とは異なるものであり、防災関連情報の緊急性や日常生活への影響の大きさ等を加味した重要度ではない。

【0258】

副計画要素は、主計画要素の情報を補うためのシステム発話の計画要素である。この副計画要素には、主計画要素で省かれた内容に基づく補足説明データ、予想される質問に対する回答データが含まれる。ユーザ発話の内容に応じて、副計画要素が再生されることになる。副計画要素についても、文圧縮と口語化の処理を行う。

【0259】

図13において、シナリオを構成するデータ（カラム）には、元の記事（文書）についての文書ID、元の文書を構成する段落についての段落ID、元の文書を構成する文の重要度（シナリオ生成時に考慮した文の重要度に、防災関連情報の緊急性や日常生活への影響の大きさ等を加味した重要度であり、システム状態記憶手段31に記憶される対象となる重要度である。）、元の文書を構成する文の内容を伝達したか否かの情報（未伝達・伝達済の別）、元の文書を構成する段落内の文についての文ID、元の文書を構成する段落内の文を構成する文節についての文節ID、シナリオの構成要素として選択されたか否かを示す情報（選択文節）、元の文書を構成する文内での文節提示順序、シナリオ再生を行うための発話節ID、リンクする発話節の合成音声ファイル（wavファイル等）の名称

10

20

30

40

50

、口語表現、文節間の間（ま）、元の文節の内容、ユーザの定義型質問に対する応答用の定義の文字情報、リンクする定義の合成音声ファイル（wavファイル等）の名称、トリビアの文字情報等が含まれる。また、図13での図示は省略されているが、リンクするトリビアの合成音声ファイル（wavファイル等）の名称も含まれ、さらに、口語表現は、複数段階の表現（例えば、伝聞口調・断定口調を使い分ける「標準」、伝聞口調だけの「伝聞」、断定口調だけの「断定」、ですます調だけの「敬体」等の口調の別を含む）が用意されている。

#### 【0260】

なお、合成音声ファイルの名称は、上記のように、シナリオを構成するデータであるが、その他に、合成音声ファイル自体（自体とは、ファイルの名称を示す文字情報ではなく、音声データを記録しているファイルという意味）を、予め生成してシナリオデータに含めてもよく、そうすることにより、次発話準備手段43による準備時や、発話生成手段25による再生時に音声合成処理を行う必要がなくなるので、システムの応答性を向上させることができる。

#### 【0261】

また、上記の例では、重要度（システム状態記憶手段31に記憶される対象となる重要度）は、元の文書を構成する文を単位とする重要度とされているが、元の文書を構成する文の単位ではなく、より細かく発話節毎に設定してもよい。

#### 【0262】

さらに、元の文書を構成する文の内容を伝達したか否かの情報（未伝達・伝達済の別）は、対話履歴記憶手段50に記憶されている対話履歴情報（図6の最上部のS（1）、S（2）、S（3）、...）と同期して更新されるが、この情報も、発話節毎に持たせてもよい。この未伝達・伝達済の別は、対話の進行に伴って逐次更新されるので、題材データ記憶手段51に記憶されているシナリオデータを直接に書き換えるわけではなく、対話状態管理手段42のメモリ（主メモリでよいが、不揮発性メモリでもよい。）にコピーされて保持されているシナリオデータを書き換える。題材データ記憶手段51に記憶されている当該シナリオデータは、同時期に他のユーザとの対話で使用されることもあるからである。主計画および副計画を備えていない他のタイプのシナリオデータの場合も同様である。

#### 【0263】

<シナリオデータを用いた音声対話の進行の概要：図14>

#### 【0264】

図14には、図13のシナリオデータを用いてシステムとユーザとの間で行われる音声対話の進行の概要が示されている。但し、次発話準備手段43による次発話候補の準備処理等の詳細は、図15および図16を用いて後述するので、ここでは表面的に表れる発話だけで対話の進行を説明する。

#### 【0265】

まず、1番目の主計画要素である「 社が3DS向けにSuicaとかと連携するゲームソフトを開発してるらしいよ」という発話節（文書ID=1、段落ID=1、文ID=1における発話節ID=1；合成音声ファイル=1-1-1-1.wav）が再生される。この発話節の途中で（例えば「 社が」の再生直後に）、ユーザから「 社って、どんな会社なの？」という定義型質問（割込み）があった場合には、「 社は」という元の文節に対して予め用意されている定義型質問応答の副計画要素「 社は、・・・」が再生される。また、発話節の途中における別の位置で（例えば「ゲームソフトを」の再生直後に）、「 社って、どんな会社なの？」という割込みがあった場合でも、「 社は」という元の文節に対して予め用意されている定義型質問応答の副計画要素「 社は、・・・」が再生される。ユーザの割込みを受けた後のシステム発話の戻りの再生開始位置は、図14中の実線で示すように、割込みを受けた位置でもよく、図14中の点線で示すように、幾つか前の文節からの再開でもよく、発話節の先頭からの再開でもよい。なお、発話節の再生の終了直後に、ユーザから「 社って、どんな会社なの？」という定義型質問（割込み）があった場合も、同様に定義型質問応答の副計画要素「 社は、・・・」が再生される

10

20

30

40

50

ので、途中であっても、終了直後であっても、同じシステム応答となる。

【0266】

次に、ユーザ発話が「楽しみだね。」であったとすると、2番目の主計画要素である「ICカードから読み取った乗車履歴を基に、」という発話節（文書ID = 1、段落ID = 1、文ID = 2における発話節ID = 1：合成音声ファイル = 1 - 1 - 2 - 1 . wav）が再生され、その再生の終了直後に、ユーザから「ICカードって、何？」という定義型質問があった場合には、「ICカードから」という元の文節に対して予め用意されている定義型質問応答の副計画要素「ICカードっていうのは、・・・」が再生される。

【0267】

続いて、3番目の主計画要素である「ゲーム内で使えるポイントが手に入るんだって」という発話節（文書ID = 1、段落ID = 1、文ID = 2における発話節ID = 2：合成音声ファイル = 1 - 1 - 2 - 2 . wav）が再生され、その再生の終了直後に、ユーザから「へー、そうなんだ。」という反応があった場合には、4番目の主計画要素（図13での図示は省略）が再生される。

【0268】

さらに、副計画要素の再生中に、ユーザからの割込みがあれば、別の副計画要素が再生されるので、副計画要素の再生は、ユーザの反応次第で階層的になることがある。また、ユーザ発話の内容次第では、副計画要素としてシナリオ内に用意していない計画外の回答を行うこともある。

【0269】

以上が対話の進行の概要であるが、以上のような対話を実現するために、対話システム10は、具体的には、例えば、図15～図17に示すような各処理を実行する。但し、図17は、図13のシナリオデータではなく、主計画および副計画からなる同型のシナリオデータを用いている。

【0270】

<次発話候補の準備処理の具体例（1）：図15>

【0271】

図15において、次発話準備手段43は、対話状態管理手段42からのシステム発話S（1）の準備開始指示情報を受け取り、次発話候補（但し、ここでは最初の発話）の準備処理を行う。対話開始時であるから、複数の次発話候補を選択取得するのではなく、シナリオデータ内から1番目の主計画要素を選択取得する。従って、次発話準備手段43は、図13のシナリオデータから、S（1）＝「社が3DS向けにSuicaとかと連携するゲームソフトを開発してるらしいよ。」を選択取得し、これを次発話候補記憶手段30に記憶させる。

【0272】

続いて、ユーザ発話は未だ無い状態なので、システム発話タイミング検出手段22により、直ぐにシステム発話の開始タイミングが検出され、次発話選択手段24により、次発話候補記憶手段30に記憶されているS（1）＝「社が3DS向けにSuicaとかと連携するゲームソフトを開発してるらしいよ。」が選択されるとともに、その選択結果がネットワーク1を介して対話状態管理手段42に送信される。対話状態管理手段42は、受け取ったS（1）を対話履歴記憶手段50に保存するとともに、S（2）の準備開始指示情報を次発話準備手段43に送る。

【0273】

それから、発話生成手段25により、次発話選択手段24により選択されたS（1）＝「社が3DS向けにSuicaとかと連携するゲームソフトを開発してるらしいよ。」の再生が開始される。また、これと並行して、次発話準備手段43により、S（2）の準備処理が進む。ここで準備されるS（2）の次発話候補は、2番目の主計画要素である「ICカードから読み取った乗車履歴を基に、」と、現在再生中の1番目の主計画要素に対するユーザの定義型質問に答えるために用意する定義型質問応答の副計画要素である「社は、...」、「社は、...」、「社3DSっていうのは、...」、「交通系ICカードっ

10

20

30

40

50

ていうのは、...」、「S u i c a っていうのは、...」、「連携っていうのは、...」、「ゲームソフトっていうのは、...」、「開発っていうのは、...」、「発表っていうのは、...」と、現在再生中の1番目の主計画要素に対するユーザの補足要求に応えるために用意する補足説明用の副計画要素(トリビア)である「S u i c a の名称は...」、「開発は、もともと...」等であり、これらが次発話候補記憶手段30に記憶される。なお、「社は、...」という定義型質問応答の副計画要素は、1番目の主計画要素の発話節として選択されなかった元の文節「社の」について用意された情報であるが、ユーザの連想が及ぶ範囲であるため、ここでは次発話候補としている。

#### 【0274】

その後、ユーザ発話U(1) = 「楽しみだね。」であったとすると、この場合は、次発話選択用情報生成手段23により、ユーザ発話意図として、例えば「理解」等の識別結果が得られるので、次発話選択手段24により選択される次発話は、S(2) = 2番目の主計画要素である「ICカードから読み取った乗車履歴を基に、」となり、この選択結果がネットワーク1を介して対話状態管理手段42に送信される。対話状態管理手段42は、選択結果としてS(2)を受け取ると、自身のメモリに保持しているU(1)(音声認識処理手段41による音声認識処理の結果として受け取り、保持している文字列)と、受け取ったS(2)とを対話履歴記憶手段50に保存するとともに、S(3)の準備開始指示情報を次発話準備手段43に送る。

#### 【0275】

それから、発話生成手段25により、次発話選択手段24により選択されたS(2) = 2番目の主計画要素である「ICカードから読み取った乗車履歴を基に、」の再生が開始される。また、これと並行して、次発話準備手段43により、S(3)の準備処理が進む。ここで準備されるS(3)の次発話候補は、3番目の主計画要素である「ゲーム内で使えるポイントが手に入るんだって。」と、現在再生中の2番目の主計画要素に対するユーザの定義型質問に答えるために用意する定義型質問応答の副計画要素である「ICカードっていうのは、...」、「基にっていうのは、...」と、現在再生中の2番目の主計画要素に対するユーザの補足要求に応えるために用意する補足説明用の副計画要素(トリビア)である「ICカードは、国際...」等であり、これらが次発話候補記憶手段30に記憶される。

#### 【0276】

< 次発話候補の準備処理の具体例(2) : 図16 >

#### 【0277】

図16には、ユーザ発話U(1)が定義型質問となり、システム発話S(2)として定義型質問応答の副計画要素が再生される場合の具体例が示されている。システム発話S(1)の再生、S(2)の複数の候補の準備までは、前述した図15の場合と同様である。

#### 【0278】

図16において、ユーザ発話U(1) = 「社って、どんな会社なの?」(割込でもよい)であったとすると、この場合は、次発話選択用情報生成手段23により、ユーザ発話意図として、例えば「質問」等の識別結果が得られるが、このユーザ発話意図だけでは、いずれの質問なのか判らないので、次発話選択手段24は、音声認識処理手段41による音声認識処理の結果(言語情報)を用いたキーワードマッチング等により、いずれの質問なのか判別し、社についての質問であることを把握する。従って、次発話選択手段24により選択される次発話は、S(2) = 定義型質問応答の副計画要素である「社は、...」となり、この選択結果がネットワーク1を介して対話状態管理手段42に送信される。対話状態管理手段42は、選択結果としてS(2) = 「社は、...」を受け取ると、自身のメモリに保持しているU(1) = 「社って、どんな会社なの?」(音声認識処理手段41による音声認識処理の結果として受け取り、保持している文字列)と、受け取ったS(2) = 「社は、...」とを対話履歴記憶手段50に保存するとともに、S(3)の準備開始指示情報を次発話準備手段43に送る。

#### 【0279】

それから、発話生成手段 25 により、次発話選択手段 24 により選択された  $S(2)$  = 定義型質問応答の副計画要素である「社は、...」の再生が開始される。また、これと並行して、次発話準備手段 43 により、 $S(3)$  の準備処理が進む。この準備処理では、次発話準備手段 43 は、対話履歴記憶手段 50 を参照し、未だ 2 番目の主計画要素である「ICカードから読み取った乗車履歴を基に、」が再生されていないことを確認することができる。従って、ここで準備される  $S(3)$  の次発話候補は、2 番目の主計画要素である「ICカードから読み取った乗車履歴を基に、」と、再生を終えている 1 番目の主計画要素に対するユーザの定義型質問に答えるために用意する定義型質問応答の副計画要素である「社は、...」、「社は、...」、「社 3 D S っていうのは、...」、「交通系 IC カードっていうのは、...」、「S u i c a っていうのは、...」、「連携っていうのは、...」、「ゲームソフトっていうのは、...」、「開発っていうのは、...」、「発表っていうのは、...」と、再生を終えている 1 番目の主計画要素に対するユーザの補足要求に応えるために用意する補足説明用の副計画要素（トリビア）である「S u i c a の名称は...」、「開発は、もともと...」等であり、これらが次発話候補記憶手段 30 に記憶される。従って、結果的に、副計画を再生した場合は、次発話候補を維持することになる。

10

20

30

40

50

#### 【0280】

この際、再生を終えている 1 番目の主計画要素に対するユーザの質問を想定した準備を行うのは、 $U(1)$  = 「社って、どんな会社なの？」というユーザの定義型質問に対し、システムが、 $S(2)$  = 定義型質問応答の副計画要素である「社は、...」を再生した後に、さらにユーザが、「S u i c a って、何？」という定義型質問をする場合があるからである。また、上述したように、結果的に次発話候補を維持するだけでもよいが、 $S(2)$  = 定義型質問応答の副計画要素である「社は、...」を再生すると、その後、ユーザから「社は、...」の中の用語について、さらに定義型質問が行われる場合があるので、シナリオデータ内に、「社は、...」という定義型質問応答の副計画要素の中の用語について、更なる定義型質問応答の副計画要素が用意されていれば、それを  $S(3)$  の次発話候補に含めて準備してもよい。

#### 【0281】

続いて、ユーザ発話  $U(2)$  = 「なるほど。」であったとすると、この場合は、次発話選択用情報生成手段 23 により、ユーザ発話意図として、例えば「理解」、「相槌」等の識別結果が得られるので、次発話選択手段 24 により選択される次発話は、 $S(3) = 2$  番目の主計画要素である「ICカードから読み取った乗車履歴を基に、」となり、この選択結果がネットワーク 1 を介して対話状態管理手段 42 に送信される。対話状態管理手段 42 は、選択結果として  $S(3)$  を受け取ると、自身のメモリに保持している  $U(2)$  = 「なるほど。」（音声認識処理手段 41 による音声認識処理の結果として受け取り、保持している文字列）と、受け取った  $S(3)$  とを対話履歴記憶手段 50 に保存するとともに、 $S(4)$  の準備開始指示情報を次発話準備手段 43 に送る。

#### 【0282】

それから、発話生成手段 25 により、次発話選択手段 24 により選択された  $S(3)$  = 2 番目の主計画要素である「ICカードから読み取った乗車履歴を基に、」の再生が開始される。また、これと並行して、次発話準備手段 43 により、 $S(4)$  の準備処理が進む。ここで次発話準備手段 43 により準備される  $S(4)$  の次発話候補は、3 番目の主計画要素である「ゲーム内で使えるポイントが手に入るんだって。」と、現在再生中の 2 番目の主計画要素に対するユーザの定義型質問に答えるために用意する定義型質問応答の副計画要素である「ICカードっていうのは、...」、「基にっていうのは、...」と、現在再生中の 2 番目の主計画要素に対するユーザの補足要求に応えるために用意する補足説明用の副計画要素（トリビア）である「ICカードは、国際...」等であり、これらが次発話候補記憶手段 30 に記憶される。

#### 【0283】

< 次発話候補の準備処理の具体例 (3) : 図 17 >

#### 【0284】

図 17 には、次発話候補の入替が行われる具体例が示されている。但し、図 17 は、図 13 のシナリオデータではなく、同型の別のシナリオデータ（不図示）を用いている。

【0285】

図 17 において、次発話準備手段 43 は、対話状態管理手段 42 からのシステム発話 S (1) の準備開始指示情報を受け取り、次発話候補（但し、ここでは最初の発話）の準備処理を行う。対話開始時であるから、複数の次発話候補を選択取得するのではなく、シナリオデータ内から 1 番目の主計画要素を選択取得する。ここでは、次発話準備手段 43 は、シナリオデータから、S (1) = 「早稲田太郎選手が 100 m 平泳ぎで優勝したんだ。」を選択取得し、これを次発話候補記憶手段 30 に記憶させる。

【0286】

続いて、ユーザ発話は未だ無い状態なので、システム発話タイミング検出手段 22 により、直ぐにシステム発話の開始タイミングが検出され、次発話選択手段 24 により、次発話候補記憶手段 30 に記憶されている S (1) = 「早稲田太郎選手が 100 m 平泳ぎで優勝したんだ。」が選択されるとともに、その選択結果がネットワーク 1 を介して対話状態管理手段 42 に送信される。対話状態管理手段 42 は、受け取った S (1) を対話履歴記憶手段 50 に保存するとともに、S (2) の準備開始指示情報を次発話準備手段 43 に送る。

【0287】

それから、発話生成手段 25 により、次発話選択手段 24 により選択された S (1) = 「早稲田太郎選手が 100 m 平泳ぎで優勝したんだ。」の再生が開始される。また、これと並行して、次発話準備手段 43 により、S (2) の準備処理が進む。ここで準備される S (2) の次発話候補は、2 番目の主計画要素である「オーストラリアで開催された国際大会での快挙なんだ。」と、現在再生中の 1 番目の主計画要素に対するユーザの定義型質問に答えるために用意する定義型質問応答の副計画要素である「早稲田太郎選手は、...」と、現在再生中の 1 番目の主計画要素に対するユーザの補足要求に応じるために用意する補足説明用の副計画要素（トリビア）である「去年の優勝者は、...」等であり、これらが次発話候補記憶手段 30 に記憶される。

【0288】

その後、ユーザ発話 U (1) = 「僕は、平泳ぎよりも、バタフライの方が興味あるんだよ。バタフライの・・・」であったとすると、次発話準備手段 43 は、例えば、U (1) = 「僕は、平泳ぎよりも、バタフライの方が興味あるんだよ。バタフライの」という途中までの情報（音声認識処理手段 41 による音声認識処理の結果である文字列）を、対話状態管理手段 42 から得ることになる。従って、次発話準備手段 43 は、このような U (1) の途中までの音声認識処理の結果に基づき、平泳ぎからバタフライへの話題の転換要求（2 回のバタフライという単語の出現、あるいは、「... よりもバタフライの方が興味ある」）を捉え、次発話候補の入替が必要であると判断し、次発話候補の入替のための準備処理を実行する。そして、次発話準備手段 43 は、使用中のシナリオデータ内にバタフライのデータ（構成要素であり、主計画要素でも副計画要素でもよい）が存在する場合には、それを選択取得する。また、使用中のシナリオデータ内にバタフライのデータが存在しない場合には、題材データ記憶手段 51 に記憶されている別のシナリオデータや、シナリオになっていない別の題材データの中からバタフライのデータを探し、それでも該当データが見つからない場合には、ネットワーク 1 を介して外部システムである題材データ提供システム 60 にアクセスし、該当データを探す。その間は、ステータス = 準備中となる。この際、次発話準備手段 43 は、題材データ記憶手段 51 および題材データ提供システム 60 のいずれを検索する場合でも、分野が関係付けられている題材データについては、先ず分野（例えば、スポーツおよび / または水泳）を用いた絞り込み検索を行うことができ、分野が関係付けられていない題材データについては、バタフライの語によるキーワード検索を行う。

【0289】

そして、バタフライの結果が見つかった場合には、次発話準備手段 43 により準備され

10

20

30

40

50

るS(2)の次発話候補は、S(2) = 「100mバタフライでは、早稲田次郎選手が6位入賞だったんだ。」、「200mバタフライでは、早稲田三郎選手が残念ながら予選落ちしたんだ。」等となり、これらが次発話候補記憶手段30に記憶される。一方、見つからなかった場合には、例えば、S(2) = 「ごめんね。バタフライの結果は知らないんだ。」、「バタフライじゃなくて、背泳ぎの結果なら知ってるよ。」等を準備し、次発話候補記憶手段30に記憶させる。

#### 【0290】

最終的に、ユーザ発話U(1)が終了し、U(1) = 「僕は、平泳ぎよりも、バタフライの方が興味あるんだよ。バタフライの結果を教えてくださいかな。早稲田次郎選手はどうだったの？」であった場合には、システム発話タイミング検出手段22によりシステム発話の開始タイミングが検出された後、次発話選択手段24は、次発話候補記憶手段30に記憶されているS(2)の複数の候補の中からの選択処理を行う。例えば、U(1)に含まれる単語である「早稲田次郎」によるキーワードマッチングを行う。従って、ここで次発話選択手段24により選択される次発話は、S(2) = 「100mバタフライでは、早稲田次郎選手が6位入賞だったんだ。」となる。そして、次発話選択手段24により、その選択結果がネットワーク1を介して対話状態管理手段42に送信される。対話状態管理手段42は、選択結果としてS(2)を受け取ると、自身のメモリに保持しているU(1) = 「僕は、平泳ぎよりも、バタフライの方が興味あるんだよ。バタフライの結果を教えてくださいかな。早稲田次郎選手はどうだったの？」(音声認識処理手段41による音声認識処理の結果として受け取り、保持している文字列)と、受け取ったS(2)とを対話履歴記憶手段50に保存するとともに、S(3)の準備開始指示情報を次発話準備手段43に送る。

10

20

#### 【0291】

<本実施形態の効果>

#### 【0292】

このような本実施形態によれば、次のような効果がある。すなわち、対話システム10は、システム発話タイミング検出手段22を備えているので、ユーザが自己の発話権を維持しているか、または、譲渡若しくは放棄により終了させたかをパターン認識処理により逐次推定することができる。また、次発話準備手段43を備えているので、システム発話タイミング検出手段22によるパターン認識処理とは非同期で、かつ、システム発話タイミング検出手段22によりシステム発話の開始タイミングが検出される前に(すなわち、ユーザ発話の進行中に、または、それよりも前の段階であるユーザ発話の開始前に)、ユーザ発話に対するシステムの次発話の内容データを準備することができる。

30

#### 【0293】

このため、対話相手であるユーザが自己の発話権を譲渡若しくは放棄することによりユーザ発話権が終了し、システム発話タイミング検出手段22により、このユーザ発話権の終了が捉えられ、システム発話の開始タイミングが検出された場合には、その検出直後に、発話生成手段25により、タイミングよくシステム発話を開始させることができるので、システムの応答性を向上させることができる。

40

#### 【0294】

また、システム発話タイミング検出手段22は、音声認識処理手段41による音声認識処理とは非同期で、ユーザ発話権の維持または終了を識別するパターン認識処理を繰り返し実行する構成とされているので、音声区間検出処理(VAD処理)を前提としない処理を実現することができるため、VAD処理による遅延なしに早期に、システム発話の開始タイミングを決定できるとともに、ユーザ発話とシステム発話との衝突も回避または抑制することができる。

#### 【0295】

以上より、対話システム10では、次発話準備手段43により、システムが発話すべき内容(本実施形態では、複数の次発話候補の内容)を早期に確定したうえで、システム発話タイミング検出手段22により、ユーザ発話権の終了が推定され、システム発話の開始

50

タイミングが検出されるのを待って、発話生成手段 25 により、システム応答を行うことができる。このため、ユーザ発話の終了後、システム発話の開始までに、長い間（ま）が空くことを避けることができるうえ、両者の発話の衝突の発生も回避または抑制することができる。

#### 【0296】

また、対話システム 10 は、次発話準備手段 43 による準備処理で取得または生成した複数の次発話候補の内容データ中から、次発話選択手段 24 が、発話生成手段 25 で用いる次発話の内容データを選択する構成とされているので、様々な種別の対話に柔軟に対応することができる。すなわち、各種の対話の中には、ユーザ発話の内容が確定する前に、そのユーザ発話に対するシステムの次発話の内容が 1 つに定まらない種別の対話も多いが、そのような場合でも、システムの応答性の向上を図ることができる。

10

#### 【0297】

そして、次発話選択手段 24 は、異なる処理で得られた複数の種類の情報を用いて、次発話の選択処理を行うことができるので、この点でも、様々な種別の対話に柔軟に対応することができる。

#### 【0298】

具体的には、次発話選択手段 24 は、音声認識処理手段 41 による音声認識処理の結果として得られた言語情報（文字列）と、次発話選択用情報生成手段 23 により得られた次発話選択用情報（主としてユーザ発話意図の識別結果であるが、ユーザの顔画像から得られた表情の識別結果や、ユーザのジェスチャー画像から得られた身振り・手振りの意図の識別情報を加えてもよい。）とのうちのいずれか一方の情報を用いて、次発話の選択処理を行うことができ、また、これらの情報を組み合わせて用いて、次発話の選択処理を行うこともできる。さらに、システム発話タイミング検出手段 22 で得られたユーザ発話意図の識別結果を用いることができる場合もある。従って、様々な種別の対話において、ユーザ発話の内容（必ずしも言語情報に限らず、ユーザ発話意図等も含めた内容）に応じて、システムの次発話の内容データを選択することができる。

20

#### 【0299】

また、上記において、次発話選択手段 24 が、韻律分析で推定したユーザ発話意図だけを用いて、複数の次発話候補の内容データの中から、次発話の内容データを選択することができる場合は、音声認識処理の結果を得る必要はないので、システムの応答性を、より一層向上させることができる。

30

#### 【0300】

また、システム発話タイミング検出手段 22 は、システム状態記憶手段 31 に記憶されている準備完了・準備中の別を示すステータスを参照する構成とされているので（図 8、図 2 参照）、システム状態を考慮し、より適切なシステム発話の開始タイミングを検出することができる。

#### 【0301】

さらに、システム発話タイミング検出手段 22 は、ユーザ状態記憶手段 32 に記憶されているユーザ発話継続時間を用いて、ユーザ発話権終了判定用閾値の調整を行うことができるので（図 9 参照）、ユーザ発話継続時間の長短に応じ、システム発話の開始タイミングを調整することができる。

40

#### 【0302】

また、システム発話タイミング検出手段 22 は、システム状態記憶手段 31 に記憶されているシステム発話意欲度の指標値（目的データの残数、次発話候補の重要度）を用いてユーザ発話権終了判定用閾値を動的に調整することができるので（図 10 参照）、システム発話意欲度が強いときには、ユーザ発話権が終了したという識別結果が出やすくなる設定状態とし、システム発話意欲度が弱いときには、ユーザ発話権が終了したという識別結果が出にくい設定状態とすることができる。

#### 【0303】

さらに、システム発話タイミング検出手段 22 は、ユーザ情報記憶手段 52 に記憶され

50



ている衝突の発生情報（蓄積情報）やシステムの交替潜時（蓄積情報）を用いて、ユーザ発話権終了判定用閾値を事前調整することができるので（図 1 1、図 1 2 参照）、各ユーザについて、衝突の発生が起きる傾向にあるときには、ユーザ発話権が終了したという識別結果が出にくい設定状態とし、システムの交替潜時が長い傾向にあるときには、ユーザ発話権が終了したという識別結果が出やすくなる設定状態とすることができる。このため、ユーザ属性に応じたユーザ発話権終了判定用閾値の調整を実現することができる。

【0304】

この際、システム発話タイミング検出手段 2 2 は、ユーザ情報記憶手段 5 2 に記憶されているユーザの発話速度（蓄積情報）を用いて、ユーザ発話権終了判定用閾値を下方調整することを決めるための下方調整用閾値を、ユーザの発話速度の関数として設定することができるので（図 1 2 参照）、各ユーザの発話速度の傾向に応じ、下方調整用閾値の設定を変更することができる。このため、ユーザ属性に応じたユーザ発話権終了判定用閾値の下方調整を実現することができる。すなわち、システムの交替潜時が長い傾向にあるときには、ユーザ発話権終了判定用閾値を下方調整することにより、ユーザ発話権が終了したという識別結果が出やすくなる設定状態とし、システムの交替潜時が短くなるようにすることができるが、この際、システムの交替潜時が長い傾向にあるか否かは、ユーザ毎に異なり、各ユーザの発話速度の傾向と関係するので、下方調整用閾値をユーザの発話速度の関数とすることで、ユーザ属性に応じてユーザ発話権終了判定用閾値の下方調整を行うか否かを決めることができる。

【0305】

また、次発話準備手段 4 3 は、入替要否判断手段 4 3 B および入替準備手段 4 3 C を備えているので（図 4 参照）、進行中のユーザ発話の内容を逐次反映させ、既に準備されている複数の次発話候補の内容データの入替を行うことができる。このため、ユーザ発話の内容に応じた適切な次発話候補の内容データを準備することができる。

【0306】

例えば、次発話準備手段 4 3 は、逐次得られる音声認識処理の結果に含まれる重要度の高い単語を用いて、ユーザの関心のある話題を決定し、題材データ記憶手段 5 1 または題材データ提供システム 6 0 に記憶された題材データの中から、決定した話題に関連付けられて記憶されている題材データを選択し、次発話の候補となる別の複数の次発話候補の内容データを取得または生成する準備処理を実行することができる。従って、次発話により提示する話題を変更することができる。

【0307】

< 変形の形態 >

【0308】

なお、本発明は前記実施形態に限定されるものではなく、本発明の目的を達成できる範囲内での変形等は本発明に含まれるものである。

【0309】

例えば、前記実施形態の対話システム 1 0 は、次発話準備手段 4 3 により、複数の次発話候補の内容データを準備して次発話候補記憶手段 3 0 に記憶させ、次発話選択手段 2 4 により、複数の次発話候補の内容データの中から、次発話の内容データを選択する構成とされていたが、本発明の対話システムは、このような構成に限定されるものではなく、次発話準備手段 4 3 により、次発話の内容データを 1 つだけ準備し、次発話選択手段 2 4 を設置しない構成としてもよい。但し、様々な種別の対話に対応できるようにするという観点で、前記実施形態のように、次発話選択手段 2 4 を設け、次発話準備手段 4 3 により複数の次発話候補の内容データを準備する構成とすることが好ましい。

【0310】

また、前記実施形態では、システム発話タイミング検出手段 2 2 は、ユーザ情報記憶手段 5 2 に記憶されているユーザの発話速度（蓄積情報）を用いて、ユーザ発話権終了判定用閾値を下方調整することを決めるための下方調整用閾値を、ユーザの発話速度の関数として設定する構成とれていたが（図 1 2 参照）、蓄積されたユーザの発話速度から得られ

るユーザ属性（発話速度の傾向）を用いるのではなく、ユーザ状態記憶手段 3 2 に記憶されているユーザのリアルタイムの発話速度（その時々発話速度）を用いて、ユーザ発話権の維持または終了を識別するパターン認識処理を行う構成としてもよい。

#### 【0311】

このような構成とする場合、ユーザ発話権の維持または終了を識別する識別器を構築する際には、ユーザ発話の音声信号から抽出した音響特徴量と、音声認識処理手段 4 1 による音声認識処理の結果として得られた言語情報から抽出した言語特徴量（但し、省略してもよい）と、ユーザ発話における対応する各時点でのユーザの発話速度とを入力して識別器の学習を行う。そして、運用時には、音響特徴量と、言語特徴量（但し、省略してもよい）と、逐次得られるリアルタイムの発話速度とを、識別器に入力する。これにより、ユーザのリアルタイムの発話速度を加味した識別結果を得ることができる。このため、ユーザ毎に異なる発話速度の傾向（蓄積された発話速度から得られるユーザ属性）に応じてユーザ発話権終了判定用閾値を調整する必要がなくなる。なお、閾値調整と併用してもよく、その場合には、閾値調整量が少なくなる。

#### 【産業上の利用可能性】

#### 【0312】

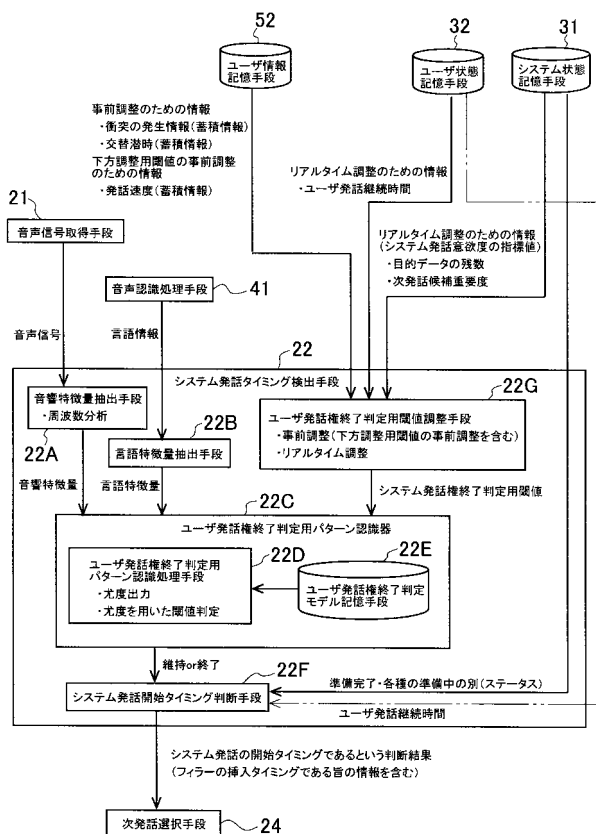
以上のように、本発明の対話システムおよびプログラムは、例えば、ニュースやコラムや歴史等の各種の話題を記載した記事データから生成したシナリオデータを用いてユーザに対して記事の内容を伝達するニュース対話システム、ユーザに対して機器の使用方法的説明や施設の案内等を行うガイダンス対話システム、選挙情勢や消費者志向等の各種のユーザの動向調査を行うアンケート対話システム、ユーザが店舗・商品・旅行先・聞きたい曲等の情報検索を行うための情報検索対話システム、ユーザが家電機器や車等の各種の機器や装置等を操作するための操作対話システム、子供や学生や新入社員等であるユーザに対して教育を行うための教育対話システム、システムがユーザ属性等の情報を特定するための情報特定対話システム等に用いるのに適している。

#### 【符号の説明】

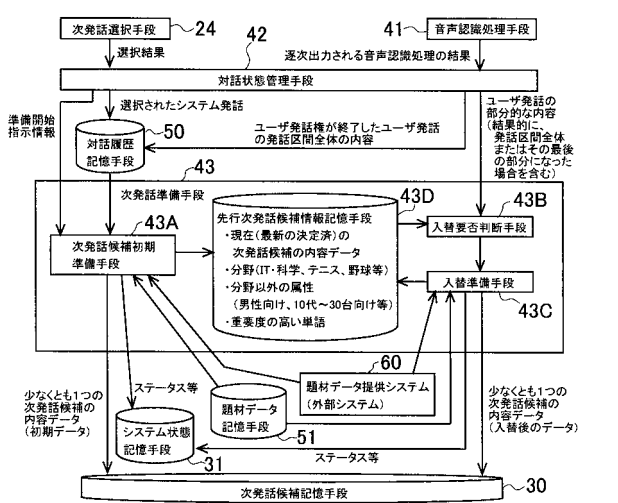
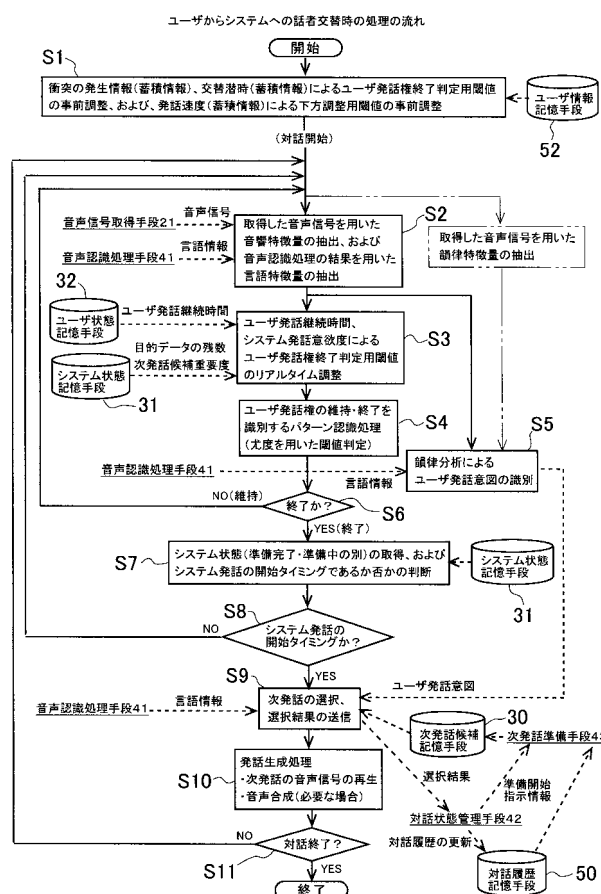
#### 【0313】

- 1 ネットワーク
- 10 対話システム
- 21 音声信号取得手段
- 22 システム発話タイミング検出手段
- 23 次発話選択用情報生成手段
- 24 次発話選択手段
- 25 発話生成手段
- 30 次発話候補記憶手段
- 31 システム状態記憶手段
- 32 ユーザ状態記憶手段
- 41 音声認識処理手段
- 42 対話状態管理手段
- 43 次発話準備手段
- 50 対話履歴記憶手段
- 51 題材データ記憶手段
- 52 ユーザ情報記憶手段
- 60 外部システムである題材データ提供システム

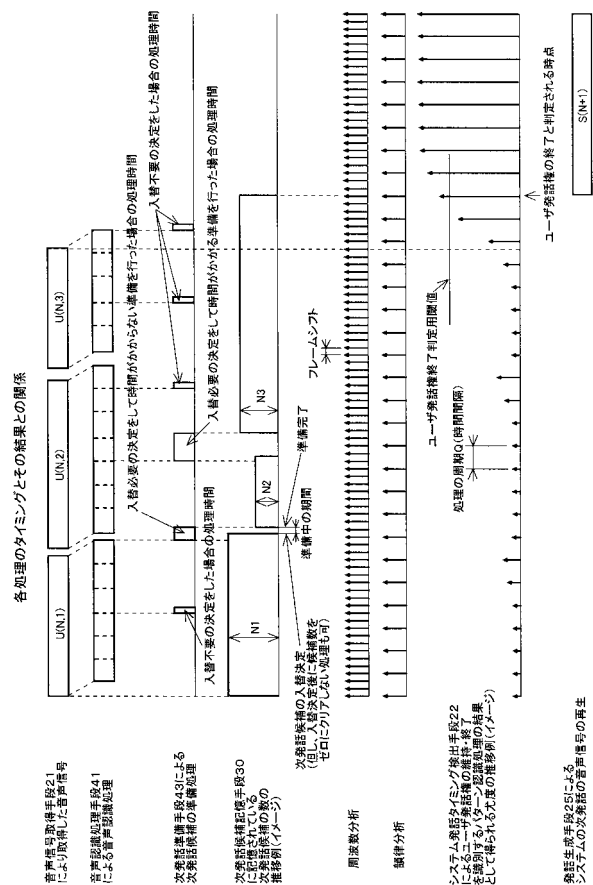
【 図 2 】



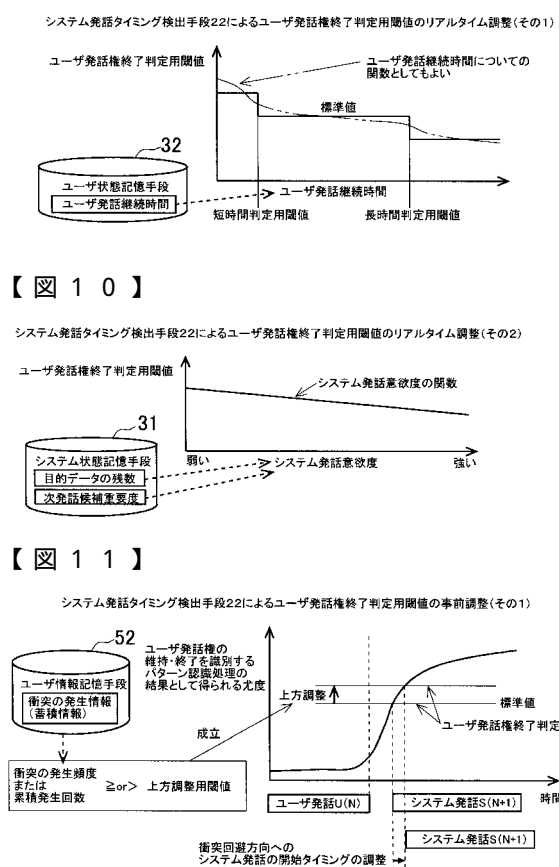
【 図 5 】



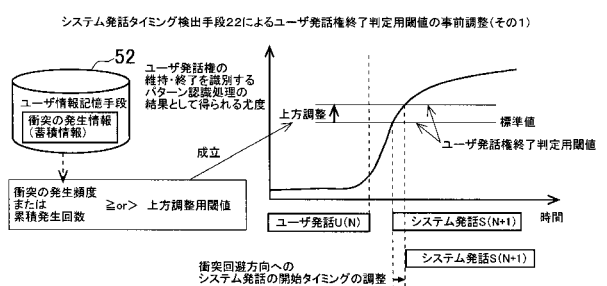
【 圖 7 】



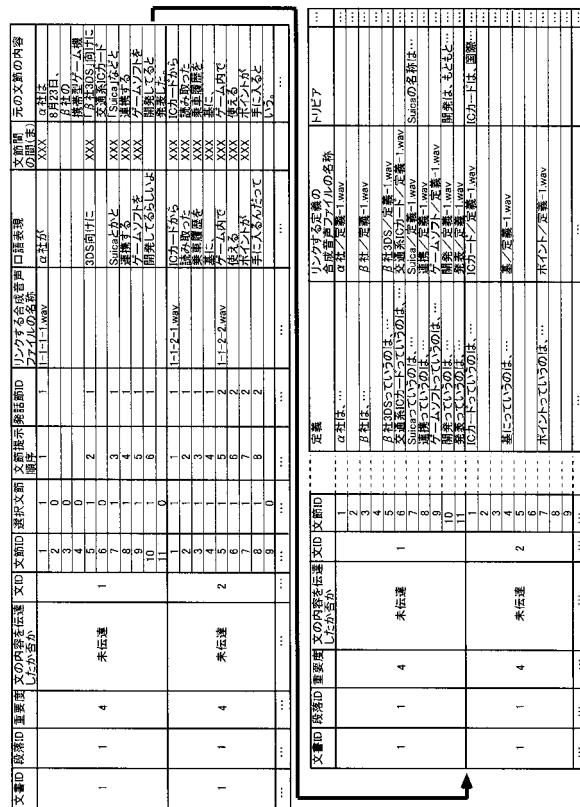
【圖 9】



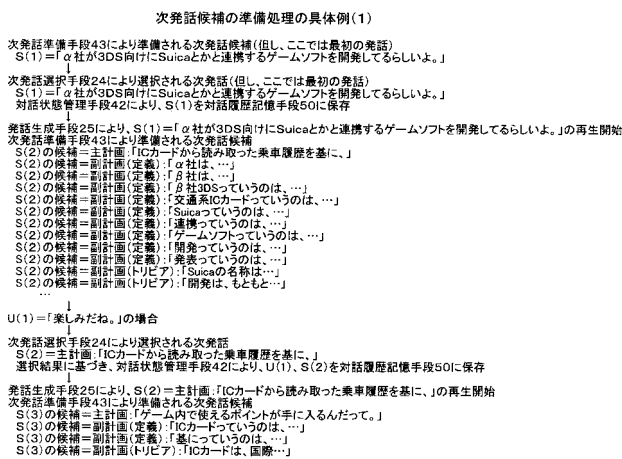
【 図 1 1 】



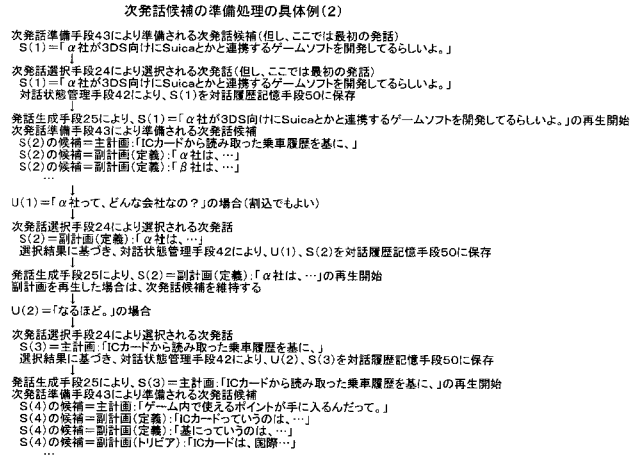
【 図 1 3 】



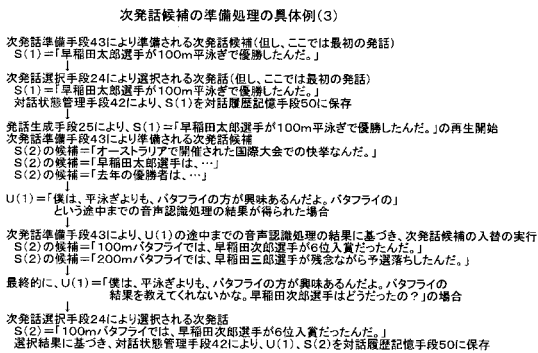
【 図 1 5 】



## 【図 16】



## 【図 17】



## フロントページの続き

|             |                          |            |
|-------------|--------------------------|------------|
| (51)Int.Cl. | F I                      | テーマコード(参考) |
|             | G 0 6 F    3/16    6 9 0 |            |

特許法第30条第2項適用申請有り (1) 令和1年7月12日掲載 (目次) [https://ipsj.ixsq.nii.ac.jp/ej/index.php?action=pages\\_view\\_main&active\\_action=repository\\_view\\_main\\_\\_item\\_snippet&index\\_\\_id=9859&pn=1&count=20&order=7&lang=japanese&page\\_\\_id=13&block\\_\\_id=8](https://ipsj.ixsq.nii.ac.jp/ej/index.php?action=pages_view_main&active_action=repository_view_main__item_snippet&index__id=9859&pn=1&count=20&order=7&lang=japanese&page__id=13&block__id=8) (論文へのアクセス用記事) <http://id.nii.ac.jp/1001/00197954/> を通じて論文抄録及び論文を発表 (2) 令和1年7月19日発表 一般社団法人情報処理学会の第128回音声言語情報処理研究会(SIG-SLP)の研究発表会、新潟県月岡温泉、風鈴屋にて、スライドを用いて発表するとともに、発表当日に来場者にダウンロード形式で論文を電子的に配布