

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6852478号
(P6852478)

(45) 発行日 令和3年3月31日 (2021.3.31)

(24) 登録日 令和3年3月15日 (2021.3.15)

(51) Int.Cl.	F I
G 1 O L 15/22 (2006.01)	G 1 O L 15/22 4 6 O Z
G 1 O L 15/30 (2013.01)	G 1 O L 15/30
G 1 O L 21/055 (2013.01)	G 1 O L 21/055
G 1 O L 15/00 (2013.01)	G 1 O L 15/00 2 O O A
H O 4 N 7/15 (2006.01)	H O 4 N 7/15 1 5 O
請求項の数 7 (全 17 頁) 最終頁に続く	

(21) 出願番号	特願2017-48205 (P2017-48205)	(73) 特許権者	000006747
(22) 出願日	平成29年3月14日 (2017.3.14)		株式会社リコー
(65) 公開番号	特開2018-151533 (P2018-151533A)		東京都大田区中馬込1丁目3番6号
(43) 公開日	平成30年9月27日 (2018.9.27)	(74) 代理人	100107766
審査請求日	令和2年1月16日 (2020.1.16)		弁理士 伊東 忠重
		(74) 代理人	100070150
			弁理士 伊東 忠彦
		(72) 発明者	中島 章敬
			東京都大田区中馬込1丁目3番6号 株式
			会社リコー内
		審査官	菊池 智紀
		最終頁に続く	

(54) 【発明の名称】 通信端末、通信プログラム及び通信方法

(57) 【特許請求の範囲】

【請求項 1】

集音装置により集音した音声データの送受信を行う通信端末であって、
 前記音声データを前記通信端末とネットワークを介して接続された音声認識装置へ送信し、前記音声認識装置から前記音声データの音声認識結果であるテキストデータを受信する通信部と、
 前記音声データを再生し、前記音声データが再生されている期間中に前記テキストデータを表示装置に表示させる出力部と、
 前記音声データに対する加工を行うか否かを判定する加工判定部と、
 前記音声データに対し、前記音声データの再生時間を延ばす加工を行う音声加工部と、
 を有し、
 前記加工判定部は、前記音声データの受信が開始されてから、前記テキストデータを受信するまでの期間が、所定の期間よりも長いとき、前記音声加工部による加工を行うものと判定する、通信端末。

【請求項 2】

前記通信部により、前記音声データの受信が開始されてから、前記テキストデータの受信が完了するまでの期間、受信した前記音声データを保持させるバッファ処理部を有し、
 前記出力部は、
 前記テキストデータの受信が完了した後に、前記音声データの再生と、前記テキストデータの表示と、を同時に開始する、請求項 1 記載の通信端末。

10

20

【請求項 3】

前記音声加工部は、

加工後の前記音声データの再生時間が、前記音声データの受信が開始されてから、前記テキストデータを受信するまでの期間よりも長くなるように、前記音声データの加工を行い、

前記出力部は、

加工後の前記音声データが再生されている期間中に前記テキストデータを表示装置に表示させる、請求項 1 又は 2 記載の通信端末。

【請求項 4】

前記所定の期間は、

前記音声データの受信が開始されてから、前記テキストデータを受信するまでの期間の平均である、請求項 2 又は 3 記載の通信端末。

【請求項 5】

前記音声データは、

前記表示装置に表示された画像の画像データ、前記表示装置に対して、入力された文字や画像を示すストローク情報の少なくとも何れかを含むコンテンツデータに含まれる、請求項 1 乃至 4 の何れか一項に記載の通信端末。

【請求項 6】

集音装置により集音した音声データの送受信を行う通信端末による通信方法であって、前記通信端末が、

前記音声データを前記通信端末とネットワークを介して接続された音声認識装置へ送信し、前記音声認識装置から前記音声データの音声認識結果であるテキストデータを受信する手順と、

前記音声データを再生し、前記音声データが再生されている期間中に前記テキストデータを表示装置に表示させる手順と、

前記音声データに対する加工を行うか否かを判定する手順と、

前記音声データに対し、前記音声データの再生時間を延ばす加工を行う手順と、を有し

、

前記判定する手順において、前記音声データの受信が開始されてから、前記テキストデータを受信するまでの期間が、所定の期間よりも長いとき、前記加工を行う手順による加工を行うものと判定する、通信方法。

【請求項 7】

集音装置により集音した音声データの送受信を行う通信端末において実行される通信プログラムであって、

前記音声データを前記通信端末とネットワークを介して接続された音声認識装置へ送信し、前記音声認識装置から前記音声データの音声認識結果であるテキストデータを受信する処理と、

前記音声データを再生し、前記音声データが再生されている期間中に前記テキストデータを表示装置に表示させる処理と、

前記音声データに対する加工を行うか否かを判定する処理と、

前記音声データに対し、前記音声データの再生時間を延ばす加工を行う処理と、通信端末に実行させ、

前記判定する処理において、前記音声データの受信が開始されてから、前記テキストデータを受信するまでの期間が、所定の期間よりも長いとき、前記加工を行う処理による加工を行うものと判定する、処理を前記通信端末に実行させる、通信プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、通信端末、通信プログラム及び通信方法に関する。

【背景技術】

【 0 0 0 2 】

従来から、音声データをテキストデータに変換する音声認識サービスを利用して、音声と対応したテキストデータを表示させる技術が普及している。

【 0 0 0 3 】

具体的には、例えば、会議システムにおいて、音声認識サービスにより発言者の音声データをテキストデータに変換し、聞き手が使用しているディスプレイに、発言者の音声データに対応したテキストデータを表示させる技術が知られている（特許文献１）。

【発明の概要】

【発明が解決しようとする課題】

【 0 0 0 4 】

従来の会議システムでは、音声データが入力された後に、この音声データが音声認識サービスに送信されてテキストデータに変換されるため、音声データが再生されるタイミングとテキストデータが表示されるタイミングにずれが生じる。

【 0 0 0 5 】

そのため、従来では、再生される音声データの内容と、表示されるテキストデータの内容とが対応せず、聞き手に違和感を与える、音声データの内容の理解を妨げる等の可能性があった。

【 0 0 0 6 】

開示の技術は、上記事情に鑑みてなされたものであり、再生中の音声データと、表示されるテキストデータとを対応させることを目的としている。

【課題を解決するための手段】

【 0 0 0 7 】

開示の技術は、集音装置により集音した音声データの送受信を行う通信端末であって、前記音声データを前記通信端末とネットワークを介して接続された音声認識装置へ送信し、前記音声認識装置から前記音声データの音声認識結果であるテキストデータを受信する通信部と、前記音声データを再生し、前記音声データが再生されている期間中に前記テキストデータを表示装置に表示させる出力部と、前記音声データに対する加工を行うか否かを判定する加工判定部と、前記音声データに対し、前記音声データの再生時間を延ばす加工を行う音声加工部と、を有し、前記加工判定部は、前記音声データの受信が開始されてから、前記テキストデータを受信するまでの期間が、所定の期間よりも長いとき、前記音声加工部による加工を行うものと判定する。

【発明の効果】

【 0 0 0 8 】

再生中の音声データと、表示されるテキストデータとを対応させることができる。

【図面の簡単な説明】

【 0 0 0 9 】

【図１】通信システムのシステム構成の一例を示す図である。

【図２】第一の実施形態の通信システムの動作の概略を説明するシーケンス図である。

【図３】比較例となる通信システムの動作の概略を説明するシーケンス図である。

【図４】第一の実施形態の通信端末のハードウェア構成の一例を示す図である。

【図５】第一の実施形態の通信端末の有する各装置の機能を説明する図である。

【図６】第一の実施形態の通信システムの動作を説明するシーケンス図である。

【図７】第一の実施形態の通信端末の動作を説明する図である。

【図８】第一の実施形態の通信端末の表示装置の表示例を示す図である。

【図９】第二の実施形態の通信端末の有する各装置の機能を説明する図である。

【図１０】第二の実施形態の通信システムの動作を説明するシーケンス図である。

【図１１】第二の実施形態の通信端末の動作を説明するフローチャートである。

【図１２】第二の実施形態の通信端末の動作を説明する図である。

【発明を実施するための形態】

【 0 0 1 0 】

(第一の実施形態)

以下に図面を参照して、第一の実施形態について説明する。図1は、通信システムのシステム構成の一例を示す図である。

【0011】

本実施形態の通信システム100は、通信端末200-1、200-2と、サーバ装置300と、を有する。通信システム100において、通信端末200-1、200-2、サーバ装置300のそれぞれは、ネットワークNを介して接続される。また、本実施形態の通信システム100は、ネットワークNを介して、音声データをテキストデータに変換する音声認識装置400と接続される。

【0012】

本実施形態の通信システム100において、通信端末200-1、200-2は、例えば電子黒板等であっても良く、サーバ装置300は、例えばテレビ会議を実現するためのテレビ会議用のサーバ装置等であっても良い。以下の説明において、通信端末200-1、200-2のそれぞれを区別しない場合には、通信端末200と呼ぶ。

【0013】

本実施形態の通信端末200は、マイク等の集音装置を有しており、集音装置によって集音された音声データを、サーバ装置300と、音声認識装置400とに送信する。また、本実施形態の通信端末200は、手書き入力された文字や画像等を示すストローク情報、画面をキャプチャした画像データ等を、サーバ装置300に送信する。さらに、本実施形態の通信端末200は、カメラ等の撮像装置を有しており、撮像装置によって撮像された画像データを、サーバ装置300に送信する。

【0014】

本実施形態の音声認識装置400は、例えば、人工知能により提供されるサービスである。音声認識装置400は、受信した音声データを音声認識機能によりテキストデータとし、サーバ装置300へ送信する。尚、本実施形態の音声データは、通信端末200の周辺で発話した人の声や、人の声以外の様々な音等、集音装置が集音した全ての音の音声データである。したがって、本実施形態では、通信端末200の周辺で発話した人の声を示す発話データは、音声データの一部である。

【0015】

本実施形態のサーバ装置300は、受信したストローク情報や画像データ、音声データ等を格納する。また、本実施形態のサーバ装置300は、音声認識装置400から送信されたテキストデータを、音声データと対応付けて格納する。以下の説明では、通信端末200からサーバ装置300に送信される各種のデータをコンテンツデータと呼ぶ。したがって、本実施形態のコンテンツデータは、音声データ、画像データ、ストローク情報等を含む。

【0016】

また、本実施形態のサーバ装置300は、例えば、ある会議において通信端末200が使用された場合、会議名と、会議中に取得したコンテンツデータと、音声データから変換されたテキストデータとが対応付けられて格納されても良い。言い換えれば、サーバ装置300では、通信端末200から取得したコンテンツデータが、会議毎に格納されても良い。

【0017】

本実施形態の通信システム100では、例えば、通信端末200-1の利用者と、通信端末200-2の利用者により、テレビ会議を行うことができる。この場合、サーバ装置300は、通信端末200-1、200-2のそれぞれから入力された情報を、通信端末200-1、200-2のそれぞれの画面に表示させ、情報を共有させる。

【0018】

具体的には、通信システム100の通信端末200は、一方の通信端末200において撮像された画像データと集音された音声データとを、サーバ装置300を介して、他方の通信端末200に送信する。

10

20

30

40

50

【 0 0 1 9 】

他方の通信端末 2 0 0 は、サーバ装置 3 0 0 から受信した画像データを表示装置に表示させ、音声データを再生する。また、他方の通信端末 2 0 0 は、受信した音声データを音声認識装置 4 0 0 へ送信してテキストデータとし、このテキストデータを表示装置に表示させる。

【 0 0 2 0 】

ここで、本実施形態の通信端末 2 0 0 では、サーバ装置 3 0 0 から受信した音声データが再生されている期間中に、この音声データの音声認識結果であるテキストデータを表示装置に表示させる。

【 0 0 2 1 】

本実施形態では、この処理により、通信端末 2 0 0 において、再生される音声データと、表示されるテキストデータとを対応させる。

【 0 0 2 2 】

尚、図 1 では、通信端末 2 0 0 の一例として、電子黒板としているが、これに限定されない。本実施形態の通信端末 2 0 0 は、集音装置と表示装置を有しており、外部の装置（サーバ装置 3 0 0、音声認識装置 4 0 0）と通信を行うことができる端末であれば良い。具体的には、本実施形態の通信端末 2 0 0 は、例えば、一般的なコンピュータ、タブレット型端末、スマートフォン等がある。また、その他にも、各種の電子機器に本実施形態を適用することができる。

【 0 0 2 3 】

以下に、図 2 及び図 3 を参照して、本実施形態の通信システム 1 0 0 の動作の概略について説明する。図 2 は、第一の実施形態の通信システムの動作の概略を説明するシーケンス図である。

【 0 0 2 4 】

本実施形態の通信システム 1 0 0 において、通信端末 2 0 0 - 1 は、撮像された画像データと集音された音声データとをサーバ装置 3 0 0 へ送信する（ステップ S 2 0 1）。サーバ装置 3 0 0 は、受信した画像データと音声データとを、通信端末 2 0 0 - 2 に送信する（ステップ S 2 0 2）。

【 0 0 2 5 】

通信端末 2 0 0 - 2 は、画像データと音声データとを受信すると、音声データを音声認識装置 4 0 0 へ送信する（ステップ S 2 0 3）。また、通信端末 2 0 0 - 2 は、受信した画像データと音声データとを、一時的に保持する（ステップ S 2 0 4）。

【 0 0 2 6 】

続いて、通信端末 2 0 0 - 2 は、音声認識装置 4 0 0 から、ステップ S 2 0 3 で送信した音声データの音声認識結果であるテキストデータを受信すると（ステップ S 2 0 5）、画像データと音声データを再生させ、受信したテキストデータを表示させる（ステップ S 2 0 6）。

【 0 0 2 7 】

このように、本実施形態の通信端末 2 0 0 では、音声データを受信した場合には、この音声データの音声認識結果のテキストデータを取得するまで、音声データの再生を行わずに待機する。

【 0 0 2 8 】

以下に、図 3 を参照して、本実施形態が適用されない通信システムの動作を説明する。図 3 は、比較例となる通信システムの動作の概略を説明するシーケンス図である。

【 0 0 2 9 】

図 3 のステップ S 3 0 1 からステップ S 3 0 3 までの処理は、図 2 のステップ S 2 0 1 からステップ S 2 0 3 までの処理と同様であるから、説明を省略する。

【 0 0 3 0 】

通信端末 2 - 2 は、ステップ S 3 0 3 において音声データを音声認識装置 4 0 0 に送信すると、サーバ装置 3 0 0 から受信した画像データと音声データとを再生する（ステップ

10

20

30

40

50

S 3 0 4)。続いて、通信端末 2 - 2 は、音声認識装置 4 0 0 からテキストデータを受信し（ステップ S 3 0 5）、受信したテキストデータを表示装置に表示させる（ステップ S 3 0 6）。

【 0 0 3 1 】

このように、図 3 の例では、通信端末 2 - 2 は、音声データを受信すると、テキストデータに変換される前に音声データの再生を開始する。したがって、図 3 の例では、テキストデータを受信して表示するまでの間に、音声データの再生が終了する可能性がある。この場合、通信端末 2 - 2 では、表示装置に表示されたテキストデータは、すでに再生が終了した音声データと対応するものとなる。よって、図 3 の例では、再生される音声データと表示されるテキストデータとが対応しない。

10

【 0 0 3 2 】

これに対し、図 2 に示す本実施形態の通信端末 2 0 0 - 2 は、音声データを受信した後に、この音声データの音声認識結果であるテキストデータを受信するまで、音声データの再生を保留する。そして、通信端末 2 0 0 - 2 は、テキストデータを受信した後に、音声データの再生とテキストデータの表示とを行う。したがって、本実施形態によれば、音声データが再生されている期間中に、この音声データの音声認識結果であるテキストデータを表示させることができ、音声データとテキストデータとを対応させることができる。

【 0 0 3 3 】

以下に、本実施形態の通信端末 2 0 0 について説明する。図 4 は、第一の実施形態の通信端末のハードウェア構成の一例を示す図である。

20

【 0 0 3 4 】

本実施形態の通信端末 2 0 0 は、入力装置 2 1 と、表示装置 2 2 と、外部 I / F 2 3 と、通信 I / F 2 4 と、R O M 2 5 (Read Only Memory) とを有する。また、本実施形態の通信端末 2 0 0 は、R A M (Random Access Memory) 2 6 と、C P U (Central Processing Unit) 2 7 と、H D D (Hard Disk Drive) 2 8 と、集音装置 2 9 と、撮像装置 3 0 と、を有する。これらの各ハードウェアは、それぞれがバス B 1 で接続されている。

【 0 0 3 5 】

入力装置 2 1 は、タッチパネル等であり、ユーザによる各種操作（例えば、音声テキスト変換（日本語）や音声テキスト変換（英語）等の機能の選択操作）を入力するのに用いられる。表示装置 2 2 は、ディスプレイ等であり、各種情報（例えば、音声テキスト変換（日本語）による変換結果を示すテキストや音声データ等）を表示する。尚、本実施形態では、タッチパネルが、入力装置と 2 1 と表示装置 2 2 の両方を兼ねていても良い。

30

【 0 0 3 6 】

外部 I / F 2 3 は、外部装置とのインターフェースである。外部装置には、記録媒体 2 3 a 等がある。これにより、通信端末 2 0 0 は、外部 I / F 2 3 を介して、記録媒体 2 3 a 等の読み取りや書き込みを行うことができる。なお、記録媒体 2 3 a には、例えば、U S B メモリや C D、D V D、S D メモリカード等がある。

【 0 0 3 7 】

通信 I / F 2 4 は、通信端末 2 0 0 をネットワーク N 1 等に接続するためのインターフェースである。これにより、通信端末 2 0 0 は、通信 I / F 2 4 を介して、他の装置（と通信を行うことができる。

40

【 0 0 3 8 】

H D D 2 8 は、プログラムやデータを格納している不揮発性の記憶装置である。H D D 2 8 に格納されるプログラムやデータには、通信端末 2 0 0 全体を制御する基本ソフトウェアである O S (Operating System)、O S 上において各種機能を提供するアプリケーションプログラム等がある。

【 0 0 3 9 】

また、H D D 2 8 は、格納しているプログラムやデータを所定のファイルシステム及び / 又は D B (データベース) により管理している。なお、通信端末 2 0 0 は、H D D 2 8 に代えて、記録媒体としてフラッシュメモリを用いるドライブ装置（例えばソリッドステ

50

ートドライブ：SSD)を有していても良い。

【0040】

ROM25は、電源を切ってもプログラムやデータを保持することができる不揮発性の半導体メモリである。ROM25には、通信端末200の起動時に実行されるBIOS(Basic Input/Output System)、OS設定、及びネットワーク設定等のプログラムやデータが格納されている。RAM26は、プログラムやデータを一時保持する揮発性の半導体メモリである。

【0041】

CPU27は、ROM25やHDD28等の記憶装置からプログラムやデータをRAM26上に読み出し、処理を実行することで、通信端末200全体の制御や機能を実現する演算装置である。

10

【0042】

集音装置29は、例えばマイクロフォン(マイク)等であり、通信端末200の周囲の音を集音する。

【0043】

撮像装置30は、例えばカメラ等であり、通信端末200の周辺の画像を撮像する。具体的には、例えば、撮像装置30は、通信端末200を用いて会議等を行っている様子等を撮像する。

【0044】

本実施形態の通信端末200は、図2に示すハードウェア構成を有することにより、後述するような各種処理を実現できる。

20

【0045】

次に、図5を参照して、本実施形態の通信端末200の機能について説明する。図5は、第一の実施形態の通信端末の有する各装置の機能を説明する図である。

【0046】

本実施形態の通信端末200の機能は、CPU27がRAM26等に格納されたプログラムを読み出して実行することで実現される。

【0047】

本実施形態の通信端末200は、集音部210、入力部220、出力部230、コンテンツ保持部240、バッファ処理部250、通信部260を有する。

30

【0048】

集音部210は、集音装置29に入力された音声を音声データとして取得する。入力部220は、通信端末200の有するタッチパネル(入力装置21、表示装置22)に対して手書き入力された文字や画像を示すストローク情報や、タッチパネルに表示された画像の画像データ等を取得する。尚、本実施形態のストローク情報とは、タッチパネルに対して手書き入力が行われた場合の、利用者による一画毎の軌跡を示す点群の座標報である。

【0049】

入力部220は、撮像装置30により撮影された画像データを取得する。尚、本実施形態の画像データは、動画データと静止画データの両方を含む。

【0050】

出力部230は、音声データや画像データを出力する。具体的には、出力部230は、例えば、表示装置22に対して画像データを表示させたり、音声データを再生させたりする。

40

【0051】

コンテンツ保持部240は、バッファ処理部250の指示により、音声データを一時的に保持する。本実施形態のコンテンツ保持部240は、例えば、通信部260によりサーバ装置300から受信したコンテンツデータを保持しても良いし、集音部210と入力部220により取得されたコンテンツデータを保持しても良い。また、コンテンツ保持部240は、コンテンツデータに含まれる音声データのみを保持しても良い。

【0052】

50

バッファ処理部 250 は、音声データが音声認識装置 400 に送信された場合、音声認識装置 400 から音声認識結果のテキストデータを受信するまで、コンテンツデータをコンテンツ保持部 240 に保持させる。

【0053】

通信部 260 は、サーバ装置 300 及び音声認識装置 400 との通信を行う。具体的には、通信部 260 は、集音部 210、入力部 220 により取得したコンテンツデータをサーバ装置 300 へ送信する。また、通信部 260 は、集音部 210 により取得した音声データを音声認識装置 400 へ送信し、音声認識結果のテキストデータを受信する。また、通信部 260 は、他の通信端末 200 から送信されたコンテンツデータを、サーバ装置 300 を介して受信する。

10

【0054】

次に、図 6 を参照して、本実施形態の通信システム 100 の動作について説明する。図 6 は、第一の実施形態の通信システムの動作を説明するシーケンス図である。

【0055】

図 6 では、通信端末 200 - 1 が取得したコンテンツデータを、サーバ装置 300 を介して通信端末 200 - 2 が受信する場合の動作を示している。

【0056】

本実施形態の通信システム 100 において、通信端末 200 - 1 は、集音部 210 - 1 により音声データを取得すると、通信部 260 - 1 へ渡す（ステップ S601）。また、通信端末 200 - 1 の入力部 220 - 2 により画像データを取得すると、通信部 260 - 1 へ渡す（ステップ S602）。通信端末 200 - 1 の通信部 260 - 1 は、音声データと画像データとを含むコンテンツデータをサーバ装置 300 へ送信する（ステップ S603）。

20

【0057】

サーバ装置 300 は、このコンテンツデータを、通信端末 200 - 2 へ送信する（ステップ S604）。

【0058】

通信端末 200 - 2 は、コンテンツデータを受信すると、通信部 260 - 2 により、コンテンツデータに含まれる音声データを音声認識装置 400 へ送信する（ステップ S605）。また、通信部 260 - 2 は、バッファ処理部 250 に、コンテンツデータを渡す（ステップ S606）。

30

【0059】

バッファ処理部 250 - 2 は、コンテンツデータを受けて、コンテンツデータをコンテンツ保持部 240 - 2 へ保持させる（ステップ S607）。

【0060】

続いて、通信端末 200 - 2 は、通信部 260 - 2 により、音声認識装置 400 から、ステップ S605 で送信した音声データの音声認識結果のテキストデータを受信する（ステップ S608）。

【0061】

続いて、通信部 260 - 2 は、バッファ処理部 250 - 2 に受信したテキストデータを渡す（ステップ S609）。バッファ処理部 250 - 2 は、テキストデータを受けると、コンテンツ保持部 240 - 2 からコンテンツデータを取得する（ステップ S610）。

40

【0062】

そして、バッファ処理部 250 - 2 は、出力部 230 - 2 に、コンテンツデータとテキストデータとを渡す（ステップ S611）。出力部 230 - 2 は、コンテンツデータとテキストデータとを同時に出力させる（ステップ S612）。

【0063】

以下に、図 7 を参照して、本実施形態の通信端末 200 の動作について、さらに具体的に説明する。図 7 は、第一の実施形態の通信端末の動作を説明する図である。

【0064】

50

図 7 の例では、本実施形態の通信端末 200 - 2 では、タイミング T 1 において、サーバ装置 300 からコンテンツデータの受信を開始し、タイミング T 2 において、コンテンツデータの受信を完了したものとする。

【0065】

このとき、通信端末 200 - 2 は、タイミング T 1 において、コンテンツデータの受信を開始した直後からコンテンツデータの再生を行わず、バッファ処理部 250 - 2 により、コンテンツ保持部 240 - 2 にコンテンツデータを保持させる。

【0066】

そして、通信端末 200 - 2 は、タイミング T 2 において、コンテンツデータの受信が完了し、タイミング T 3 において、テキストデータの受信が完了すると、コンテンツデータの再生と、テキストデータの表示を同時に開始する。したがって、本実施形態では、タイミング T 1 からタイミング T 3 までの期間 K 13 は、コンテンツ保持部 240 - 2 にコンテンツデータが保持されるコンテンツ保持期間となる。

【0067】

したがって、通信端末 200 - 2 では、タイミング T 1 からタイミング T 3 の間に、コンテンツデータに含まれる音声データを音声認識装置 400 へ送信し、音声認識結果のテキストデータを受信する。尚、本実施形態では、タイミング T 2 において、コンテンツデータの受信が完了した後に、音声認識装置 400 への音声データの送信を開始しても良いし、タイミング T 1 においてコンテンツデータの受信を開始したときから、音声認識装置 400 への音声データの送信を開始しても良い。

【0068】

図 7 の例では、コンテンツデータの再生時間は、タイミング T 3 からタイミング T 5 までの期間 K 35 であり、テキストデータの表示時間は、タイミング T 3 からタイミング T 4 までの期間 K 34 となる。

【0069】

したがって、本実施形態の通信端末 200 では、音声データを含むコンテンツデータが再生されている期間中に、この音声データと対応するテキストデータが表示装置 22 に表示させることができる。

【0070】

このように、本実施形態によれば、再生中の音声データと、表示されるテキストデータとを対応させることができるため、聞き手に対して違和感を与えることがない。また、本実施形態によれば、音声データの再生中にテキストデータを表示することによって、音声データの内容の理解を支援することができる。

【0071】

図 8 は、第一の実施形態の通信端末の表示装置の表示例を示す図である。図 8 に示す画面 81 は、例えば、拠点 A に設置された通信端末 200 - 1 と、拠点 B に設置された通信端末 200 - 2 とを用いてテレビ会議を行った場合の、通信端末 200 - 1、200 - 2 のそれぞれの表示装置 22 に表示される画面の例である。

【0072】

画面 81 は、通信端末 200 - 1 の入力部 220 - 1 により取得された画像データが表示される表示領域 82、通信端末 200 - 2 の入力部 220 - 2 により取得された画像データが表示される表示領域 83 を含む。また、画面 81 には、通信端末 200 - 1 の集音部 210 - 1 と、通信端末 200 - 2 の集音部 210 - 2 のそれぞれから取得された音声データから変換されたテキストデータが表示される表示領域 84 を含む。

【0073】

本実施形態では、例えば、表示領域 83 において表示された利用者の画像データと音声データが再生されている期間中に、表示領域 84 において、この音声データと対応するテキストデータが表示される。したがって、本実施形態によれば、音声データの内容を示すテキストデータが、音声データの再生中に表示されることになる。このため、本実施形態によれば、画面 81 を閲覧している利用者に対して、音声データの再生のタイミングと、

10

20

30

40

50

テキストデータの表示のタイミングのずれを感じさせることがなく、操作性を向上させることができる。

【 0 0 7 4 】

尚、本実施形態では、主に通信端末 2 0 0 - 2 の動作を説明したが、通信端末 2 0 0 - 1 と通信端末 2 0 0 - 2 は、同様の構成を有するものであり、通信端末 2 0 0 - 1 も、通信端末 2 0 0 - 2 と同様の動作を行うものである。

【 0 0 7 5 】

(第二の実施形態)

以下に図面を参照して、第二の実施形態について説明する。第二の実施形態では、コンテンツデータを受信してから音声データを再生するまでのコンテンツ保持期間に応じて、音声データに対する加工を行う点が、第一の実施形態と相違する。よって、以下の第二の実施形態の説明では、第一の実施形態との相違点についてのみ説明し、第一の実施形態と同様の機能構成を有するものには、第一の実施形態の説明で用いた符号を付与し、その説明を省略する。

【 0 0 7 6 】

図 9 は、第二の実施形態の通信端末の有する各装置の機能を説明する図である。本実施形態の通信端末 2 0 0 A は、第一の実施形態の通信端末 2 0 0 の有する各部に加え、保持期間取得部 2 7 0 、加工判定部 2 8 0 、音声加工部 2 9 0 を有する。

【 0 0 7 7 】

本実施形態の保持期間取得部 2 7 0 は、通信部 2 6 0 がコンテンツデータの受信を開始してから、コンテンツデータに含まれる音声データの音声認識結果のテキストデータを受信するまでのコンテンツ保持期間を算出して取得し、記憶する。

【 0 0 7 8 】

加工判定部 2 8 0 は、コンテンツ保持期間に基づき、音声データに対する加工を行うか否かを判定する。具体的には、加工判定部 2 8 0 は、コンテンツ保持期間が所定の期間より長い期間であるか否かを判定し、所定の期間より長い期間である場合には、音声データに対する加工を行うものと判定する。また、加工判定部 2 8 0 は、コンテンツ保持期間が所定の期間以内である場合には、音声データに対する加工はせずに、コンテンツ保持部 2 4 0 により、音声データを含むコンテンツデータを保持するものと判定する。

【 0 0 7 9 】

また、本実施形態の所定の期間とは、例えば、過去にコンテンツデータを受信した際のコンテンツ保持期間の平均である。本実施形態では、例えば、加工判定部 2 8 0 が、保持期間取得部 2 7 0 によりコンテンツ保持期間が記憶される度に、所定の期間となるコンテンツ保持期間の平均を算出し、保持していても良い。

【 0 0 8 0 】

音声加工部 2 9 0 は、前回のコンテンツデータの受信の際に、保持期間取得部 2 7 0 により記憶されたコンテンツ保持期間が、予め設定された所定の期間より長い期間であった場合に、コンテンツデータに含まれる音声データを引き伸ばす加工を行う。具体的には、音声加工部 2 9 0 は、音声データである波形を編集することで、音声データの再生時間を延ばすようにしても良い。

【 0 0 8 1 】

尚、本実施形態の音声加工部 2 9 0 は、コンテンツデータに含まれる画像データが静止画の画像データである場合には、音声データについてのみ、加工を行っても良い。また、本実施形態の音声加工部 2 9 0 は、コンテンツデータに含まれる画像データが動画データである場合には、動画データに対しても、音声データに対する加工と同様の加工を行う。

【 0 0 8 2 】

例えば、コンテンツデータが、動画データであり、フレームレートが 3 0 [f p s] である場合には、このフレームレートを 1 5 [f p s] とすれば、コンテンツデータの再生時間を 2 倍にすることができる。

【 0 0 8 3 】

以下に、図 1 0 を参照して、本実施形態の通信システム 1 0 0 の動作について説明する。図 1 0 は、第二の実施形態の通信システムの動作を説明するシーケンス図である。

【 0 0 8 4 】

図 1 0 では、通信端末 2 0 0 A - 1 からサーバ装置 3 0 0 に送信されたコンテンツデータを通信端末 2 0 0 A - 2 が受信した際に、コンテンツ保持期間が所定の期間より長い期間であった場合の動作を示している。

【 0 0 8 5 】

図 1 0 のステップ S 1 0 0 1 とステップ S 1 0 0 2 の処理は、図 2 のステップ S 2 0 1 とステップ S 2 0 2 の処理と同様であるから、説明を省略する。

10

【 0 0 8 6 】

通信端末 2 0 0 A - 2 は、通信部 2 6 0 - 2 により、コンテンツデータを受信すると、保持期間取得部 2 7 0 により、前回の保持されているコンテンツ保持期間を取得する（ステップ S 1 0 0 3 ）。

【 0 0 8 7 】

続いて、通信端末 2 0 0 A - 2 は、加工判定部 2 8 0 により、取得したコンテンツ保持期間が所定の期間より長い期間であると判定された場合に、音声加工部 2 9 0 により、コンテンツデータに含まれる音声データを加工する（ステップ S 1 0 0 4 ）。このとき、音声加工部 2 9 0 は、加工後の音声データを再生したときの再生時間が、コンテンツ保持期間よりも長くなるように、音声データを引き延ばす加工を行う。

20

【 0 0 8 8 】

続いて、通信端末 2 0 0 A - 2 は、音声データを音声認識装置 4 0 0 に送信する（ステップ S 1 0 0 5 ）。また、通信端末 2 0 0 A - 2 は、加工された音声データと画像データとの再生を開始する（ステップ S 1 0 0 6 ）。

【 0 0 8 9 】

続いて、通信端末 2 0 0 A - 2 は、音声認識装置 4 0 0 より、音声データの音声認識結果のテキストデータを受信し（ステップ S 1 0 0 7 ）、このテキストデータを表示装置 2 2 に表示させる（ステップ S 1 0 0 8 ）。

【 0 0 9 0 】

以上のように、本実施形態では、コンテンツ保持期間の長さに応じて、音声データの再生時間が、コンテンツ保持期間よりも長くなるように加工する。したがって、本実施形態によれば、加工後の音声データの再生中に、この音声データと対応するテキストデータを表示させることができる。

30

【 0 0 9 1 】

以下に、図 1 1 を参照して、本実施形態の通信端末 2 0 0 A の全体の動作を説明する。図 1 1 は、第二の実施形態の通信端末の動作を説明するフローチャートである。本実施形態の通信端末 2 0 0 A は、コンテンツデータを受信する度に、図 1 1 の処理を実行する。

【 0 0 9 2 】

本実施形態の通信端末 2 0 0 A は、通信部 2 6 0 により、コンテンツデータの受信を開始すると（ステップ S 1 1 0 1 ）、保持期間取得部 2 7 0 により、前回のコンテンツデータの受信において記憶されたコンテンツ保持期間を取得する（ステップ S 1 1 0 2 ）。

40

【 0 0 9 3 】

続いて、通信端末 2 0 0 A は、加工判定部 2 8 0 により、コンテンツ保持期間が所定の期間より長い期間であるか否かを判定する（ステップ S 1 1 0 3 ）。ステップ S 1 1 0 3 において、コンテンツ保持期間が所定の期間よりも長い場合、通信端末 2 0 0 A は、後述するステップ S 1 1 0 9 へ進む。

【 0 0 9 4 】

ステップ S 1 1 0 3 において、コンテンツ保持期間が所定の期間以内である場合、通信端末 2 0 0 A は、バッファ処理部 2 5 0 により、コンテンツデータをコンテンツ保持部 2 4 0 に保持させる（ステップ S 1 1 0 4 ）。続いて、通信端末 2 0 0 A は、通信部 2 6 0

50

により、音声データを音声認識装置 400 へ送信し（ステップ S 1105）、音声認識装置 400 からテキストデータを受信する（ステップ S 1106）。

【0095】

続いて通信端末 200A は、出力部 230 により、コンテンツデータを再生し、テキストデータを表示させる（ステップ S 1107）。続いて、通信端末 200A は、保持期間取得部 270 により、ステップ S 1101 でコンテンツデータの受信を開始してから、ステップ S 1106 でテキストデータを受信するまでのコンテンツ保持期間を取得して記憶し（ステップ S 1108）、処理を終了する。

【0096】

ステップ S 1103 において、コンテンツ保持期間が所定の期間より長い場合、通信端末 200A は、音声加工部 290 により、音声データの再生時間がコンテンツ保持期間よりも長くなるように、音声データを引き伸ばす加工を行う（ステップ S 1109）。

【0097】

続いて、通信端末 200A は、通信部 260 より、加工していない音声データを音声認識装置 400 へ送信する（ステップ S 1110）。また、通信端末 200A は、出力部 230 により、加工された音声データと、コンテンツデータに含まれる画像データと、の再生を開始する（ステップ S 1111）。

【0098】

続いて、通信端末 200A は、通信部 260 により、音声認識装置 400 からテキストデータを受信し（ステップ S 1112）、出力部 230 により、受信したテキストデータを表示装置 22 に表示させ（ステップ S 1113）、ステップ S 1108 へ進む。

【0099】

以下に、図 12 を参照して、実施形態の通信端末 200A の動作について、さらに説明する。図 12 は、第二の実施形態の通信端末の動作を説明する図である。

【0100】

図 12 の例では、本実施形態の通信端末 200A - 2 では、タイミング T 1 において、サーバ装置 300 からコンテンツデータの受信を開始し、タイミング T 2 において、コンテンツデータの受信を完了したものとする。また、図 12 では、通信端末 200A - 2 の保持期間取得部 270 により記憶されたコンテンツ保持期間が所定の期間より長いときの、通信端末 200A - 2 の動作を示している。

【0101】

図 12 の例では、タイミング T 1 において、コンテンツデータの受信を開始した直後からコンテンツデータの保持を行わず、コンテンツデータの再生を開始する。このとき、コンテンツデータに含まれる音声データは、コンテンツ保持期間よりも再生時間が長くなるように、加工されている。

【0102】

図 12 において、通信端末 200A - 2 は、タイミング T 2 において、コンテンツデータの受信が完了する。また、通信端末 200A - 2 は、タイミング T 3 において、テキストデータを受信すると、タイミング T 3 からタイミング T 4 までの期間 K 34 において、テキストデータを表示させる。

【0103】

このとき、コンテンツデータは、まだ再生中である。図 12 の例では、通信端末 200A - 2 は、タイミング T 5 において、加工後の音声データを含むコンテンツデータの再生が完了する。したがって、加工後の音声データを含むコンテンツデータの再生時間は、タイミング T 1 からタイミング T 5 までの期間 K 15 となる。

【0104】

したがって、本実施形態によれば、テキストデータが表示される期間 K 34 は、加工後の音声データを含むコンテンツデータの再生時間である期間 K 15 に含まれることになる。

【0105】

10

20

30

40

50

このように、本実施形態では、音声データ（コンテンツデータ）を引き延ばすことで、音声データの再生中に、この音声データの音声認識結果であるテキストデータが表示されるようにすることができる。

【0106】

以上、各実施形態に基づき本発明の説明を行ってきたが、上記実施形態に示した要件に本発明が限定されるものではない。これらの点に関しては、本発明の主旨をそこなわない範囲で変更することができ、その応用形態に応じて適切に定めることができる。

【符号の説明】

【0107】

- 100 通信システム
- 200、200A 通信端末
- 210 集音部
- 220 入力部
- 230 出力部
- 240 コンテンツ保持部
- 250 バッファ処理部
- 260 通信部
- 270 保持期間取得部
- 280 加工判定部
- 290 音声加工部
- 300 サーバ装置
- 400 音声認識装置

10

20

【先行技術文献】

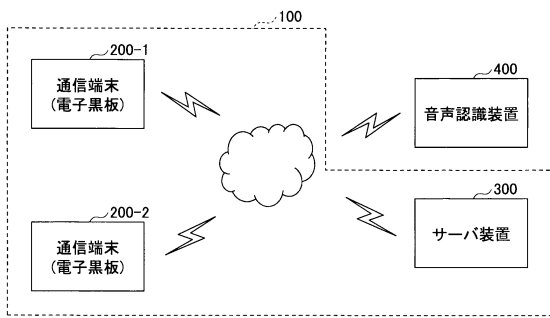
【特許文献】

【0108】

【特許文献1】特開2011-182125号公報

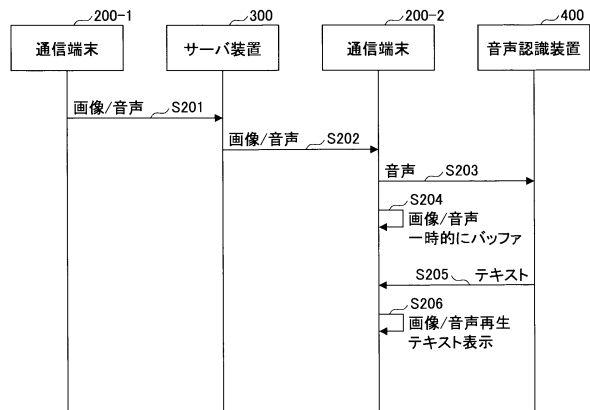
【図 1】

通信システムのシステム構成の一例を示す図



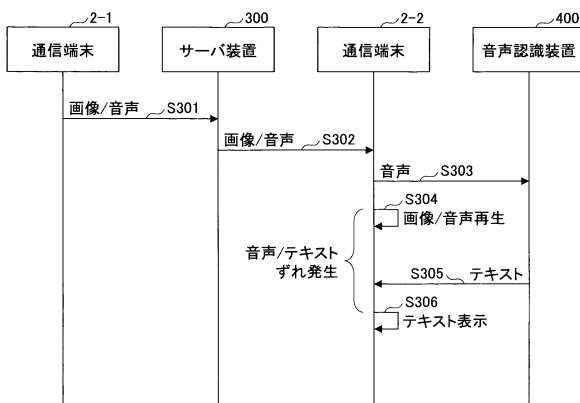
【図 2】

第一の実施形態の通信システムの動作を説明するシーケンス図



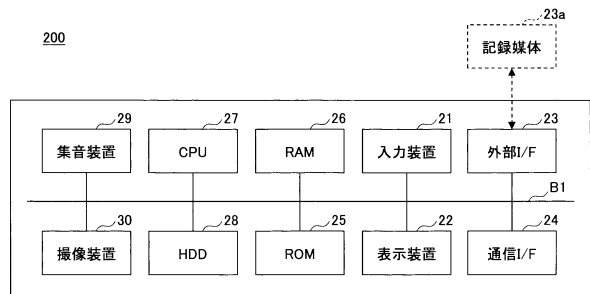
【図 3】

比較例となる通信システムの動作を説明するシーケンス図



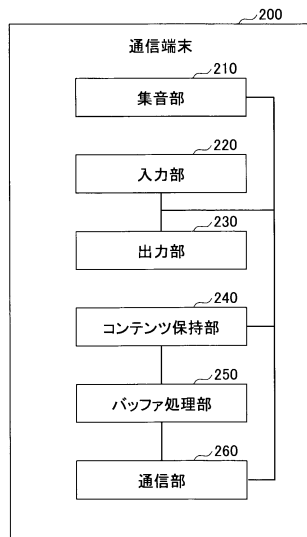
【図 4】

第一の実施形態の通信端末のハードウェア構成の一例を示す図



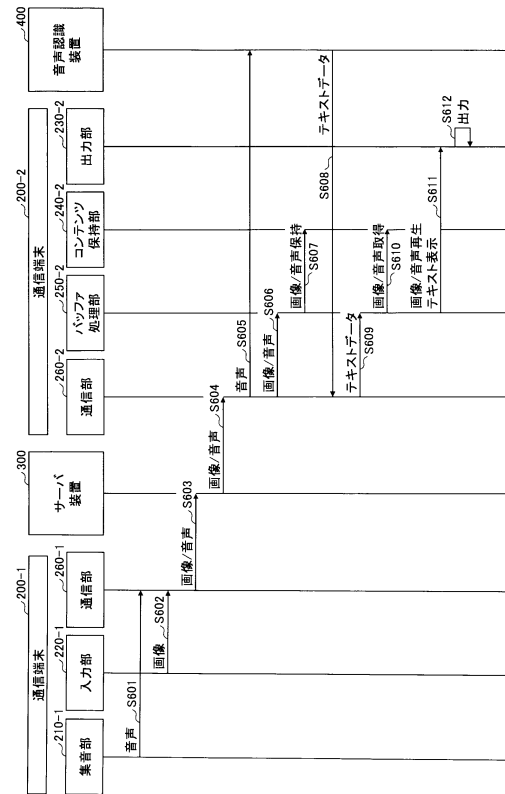
【 図 5 】

第一の実施形態の通信端末の有する各装置の機能を説明する図



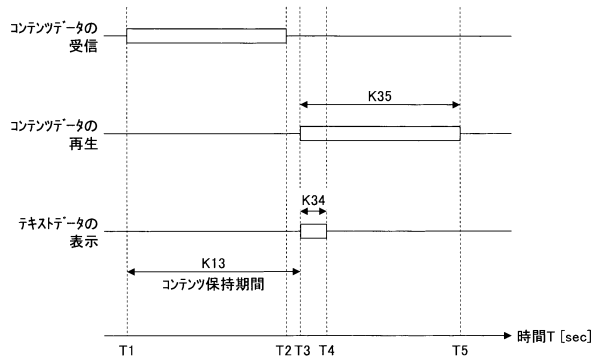
【 図 6 】

第一の実施形態の通信システムの動作を説明するシーケンス図



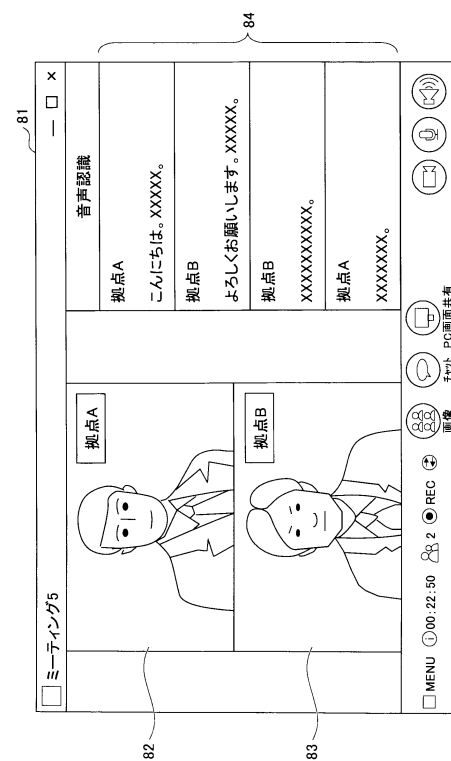
【 圖 7 】

第一の実施形態の通信端末の動作を説明する図



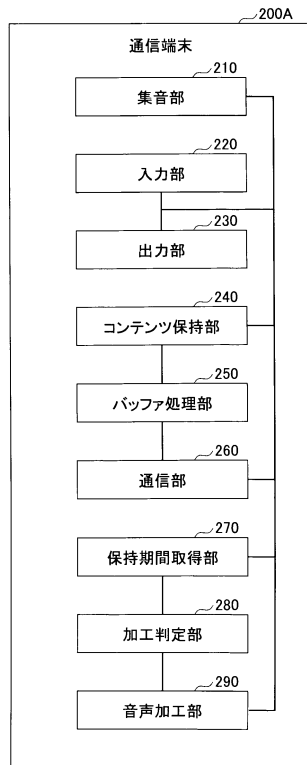
【 図 8 】

第一の実施形態の通信端末の表示装置の表示例を示す図



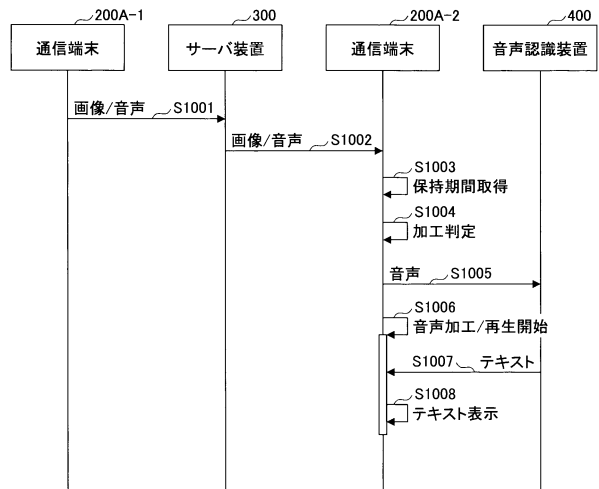
【図 9】

第二の実施形態の通信端末の有する各装置の機能を説明する図



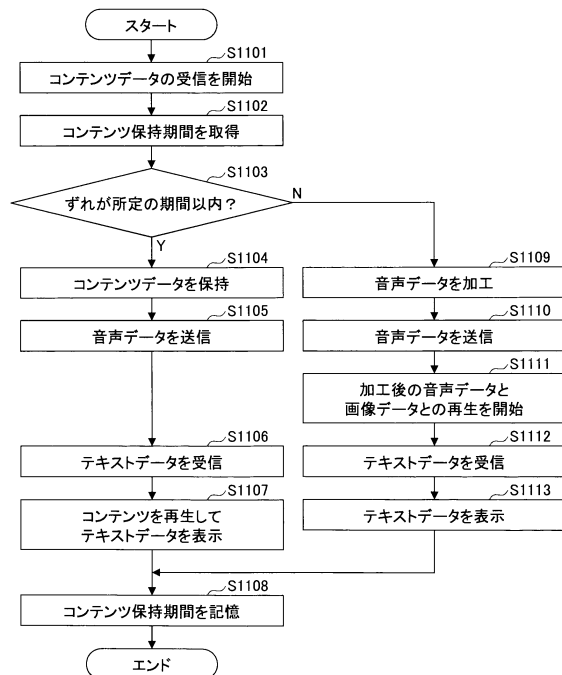
【図 10】

第二の実施形態の通信システムの動作を説明するシーケンス図



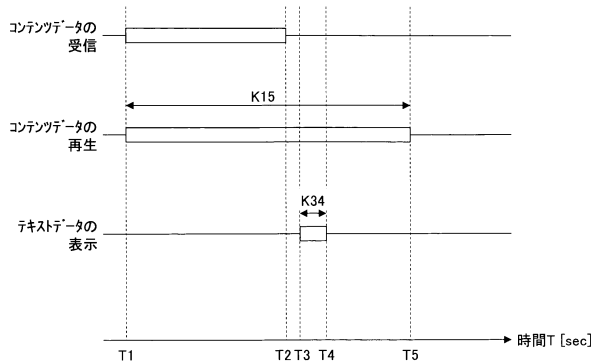
【図 11】

第二の実施形態の通信端末の動作を説明するフローチャート



【図 12】

第二の実施形態の通信端末の動作を説明する図



フロントページの続き

(51)Int.Cl. F I
H 0 4 M 3/42 (2006.01) H 0 4 M 3/42 R

(56)参考文献 特開 2 0 1 0 - 2 3 0 9 4 8 (J P , A)
特開 2 0 0 5 - 6 4 6 0 0 (J P , A)
特開 2 0 0 3 - 2 0 9 6 0 0 (J P , A)
特開 2 0 0 4 - 3 0 4 6 0 1 (J P , A)
特開 2 0 1 2 - 1 2 9 9 5 0 (J P , A)
特開 2 0 0 8 - 6 6 8 6 6 (J P , A)
米国特許出願公開第 2 0 1 4 / 0 1 9 2 1 3 8 (U S , A 1)

(58)調査した分野(Int.Cl. , D B 名)

G 1 0 L 1 5 / 0 0 - 1 7 / 2 6 ,
2 1 / 0 0 - 2 1 / 0 5 7
H 0 4 M 1 / 0 0 - 3 / 6 4
H 0 4 N 7 / 1 4 - 7 / 1 5