



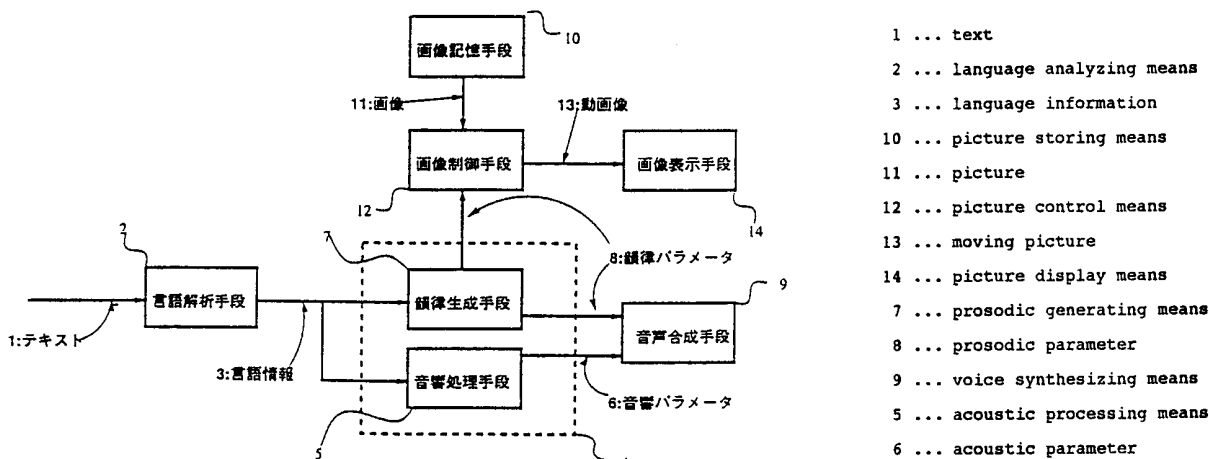
PCT

特許協力条約に基づいて公開された国際出願

<p>(51) 国際特許分類6 G06T 13/00, G10L 5/02</p>	<p>A1</p>	<p>(11) 国際公開番号 WO99/46732</p> <p>(43) 国際公開日 1999年9月16日(16.09.99)</p>
<p>(21) 国際出願番号 PCT/JP98/01025</p> <p>(22) 国際出願日 1998年3月11日(11.03.98)</p> <p>(71) 出願人 (米国を除くすべての指定国について) 三菱電機株式会社 (MITSUBISHI DENKI KABUSHIKI KAISHA)[JP/J] 〒100-8310 東京都千代田区丸の内二丁目2番3号 Tokyo, (JP)</p> <p>(72) 発明者; および</p> <p>(75) 発明者/出願人 (米国についてのみ) 海老原充(EBIHARA, Takashi)[JP/J] 石川 泰(ISHIKAWA, Yasushi)[JP/J] 〒100-8310 東京都千代田区丸の内二丁目2番3号 三菱電機株式会社内 Tokyo, (JP)</p> <p>(74) 代理人 弁理士 曾我道照, 外(SOGA, Michiteru et al.) 〒100-0005 東京都千代田区丸の内三丁目1番1号 国際ビルディング8階 曾我特許事務所 Tokyo, (JP)</p>	<p>(81) 指定国 JP, US, 欧州特許 (AT, BE, CH, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE)</p> <p>添付公開書類 国際調査報告書</p>	

(54) Title: MOVING PICTURE GENERATING DEVICE AND IMAGE CONTROL NETWORK LEARNING DEVICE

(54) 発明の名称 動画画像生成装置および画像制御ネットワーク学習装置



(57) Abstract

A moving picture generating device is provided with a language analyzing means (2) which produces language information (3) on the pronunciation of words and "bunsetsu", parts of speech, accent types, etc., by analyzing an inputted text (1), an acoustic processing means (5) which determines acoustic parameters (6) based on the language information (3), a prosody generating means (7) for determining prosodic parameters (8) such as the pitch, duration, pose, etc., based on the information (3), a voice synthesizing means (9) which generates synthesized voices from the acoustic parameters (6) and prosodic parameters (8), a picture storing means (10) which stores the picture of the face and body of a person corresponding to the prosodic parameters, a picture control means (12) which determines a moving picture (13) to be outputted referring to the picture storing means (10) based on the prosodic parameters (8), and a picture display means (14) which displays the moving picture.

(57)要約

入力されたテキスト1を解析して単語や文節の読み、品詞、アクセント型などの言語情報3を得る言語解析手段2、言語情報3を基に音響パラメータ6を決定する音響処理手段5、言語情報3を基にピッチや継続時間長やポーズ等の韻律パラメータ8を決定する韻律生成手段7、音響パラメータ6と韻律パラメータ8により合成音を生成する音声合成手段9、韻律パラメータに対応した顔や人体の画像を記憶する画像記憶手段10、韻律パラメータ8に基づき画像記憶手段10を参照して出力する動画像13を決める画像制御手段12、動画像を表示する画像表示手段14を備える。

PCTに基づいて公開される国際出願のパンフレット第一頁に掲載されたPCT加盟国を同定するために使用されるコード(参考情報)

AE アラブ首長国連邦	DM ドミニカ	KZ カザフスタン	SD スーダン
AL アルバニア	EE エストニア	LC セントルシア	SE スウェーデン
AM アルメニア	ES スペイン	LI リヒテンシュタイン	SG シンガポール
AT オーストリア	FI フィンランド	LK スリ・ランカ	SI スロヴェニア
AU オーストラリア	FR フランス	LR リベリア	SK スロヴァキア
AZ アゼルバイジャン	GA ガボン	LS レソト	SL シェラ・レオネ
BA ボスニア・ヘルツェゴビナ	GB 英国	LT リトアニア	SN セネガル
BB バルバドス	GD グレナダ	LU ルクセンブルグ	SZ スワジランド
BE ベルギー	GE グルジア	LV ラトヴィア	TD チャード
BF ブルキナ・ファソ	GH ガーナ	MC モナコ	TG トーゴ
BG ブルガリア	GM ガンビア	MD モルドヴァ	TJ タジキスタン
BJ ベナン	GN ギニア	MG マダガスカル	TZ タンザニア
BR ブラジル	GW ギニア・ビサウ	MK マケドニア旧ユーゴスラヴィア共和国	TM トルクメニスタン
BY ベラルーシ	GR ギリシャ	ML マリ	TR トルコ
CA カナダ	HR クロアチア	MN モンゴル	TT トリニダード・トバゴ
CF 中央アフリカ	HU ハンガリー	MR モーリタニア	UA ウクライナ
CG コンゴ	ID インドネシア	MW マラウイ	UG ウガンダ
CH スイス	IE アイルランド	MX メキシコ	US 米国
CI コートジボアール	IL イスラエル	NE ニジェール	UZ ウズベキスタン
CM カメルーン	IN インド	NL オランダ	VN ヴィエトナム
CN 中国	IS アイスランド	NO ノールウェー	YU ユーゴスラビア
CR コスタ・リカ	IT イタリア	NZ ニュー・ジーランド	ZA 南アフリカ共和国
CU キューバ	JP 日本	PL ポーランド	ZW ジンバブエ
CY キプロス	KE ケニア	PT ポルトガル	
CZ チェッコ	KG キルギスタン	RO ルーマニア	
DE ドイツ	KP 北朝鮮	RU ロシア	
DK デンマーク	KR 韓国		

明 細 書

動画像生成装置および画像制御ネットワーク学習装置

技術分野

この発明は、規則合成や音声対話システムなどにおいて、画像と音声を同時に出力提示する動画像生成装置および画像制御ネットワーク学習装置に関するものである。

背景技術

規則合成技術は、任意の文字列からなるテキストを音声に変換して提示する技術であり、音声による情報サービス、身障者向け朗読器、新聞校閲などの用途に用いられている。

また、音声によるマンマシンインタフェースを実現する技術においても、規則合成技術が適用されており、近年、音声の合成と同様に任意の文章情報からそれを喋ったときの口唇の動きやそれを含む人体の動画像を生成する技術の開発も行われている。これを音声合成技術と併用することで、より自然なマンマシンインタフェースが実現可能となる。

このように文章情報に基づいて音声と同時に画像を合成し出力する動画像生成装置の従来技術には、文字列に対応した口形状変化を有する顔動画像を生成する画像合成装置として、例えば特開平2-234285号公報の「画像合成方法及びその装置」（以下、文献1とする）に開示されている技術がある。

図11は文献1で示される手法に基づく従来の動画像生成装置の構成図である。

図11において、1はテキスト（文字列）、41はテキスト1を音素列42に分割して出力する音声合成部、43は音声合成部41からの各音素42毎の音韻特徴を口形特徴44に変換する変換部、45は種々の口形特徴と具体的な口形状

パラメータとの対応付けが記述されている変換テーブル46を記憶している変換テーブル記憶部である。

また、47は変換テーブル46を参照して口形特徴44に対応する口形状パラメータ48を取得する口形状パラメータ取得部、49は音素列42の各音素に対応する音素継続時間長50を出力する時間調整部、51は音素継続時間長50の期間口形状パラメータ48を出力するゲート、52はゲート51を介して入力される口形状パラメータに基づいて画像を生成する画像生成部である。

次に、上述した図1に示す従来例の動画像生成装置について説明する。

テキスト1を音声合成部41で音素列42に分割し、変換部43で各音素毎の音韻特徴から口形特徴44へ変換し、種々の口形特徴と具体的な口形状パラメータとの対応付けが記述されている変換テーブル46を変換テーブル記憶部45から参照して、口形状パラメータ取得部47で口形状パラメータ48を取り出し、一定時間間隔の画像系列として与えられる動画像を生成するために、時間調整部49で口形状パラメータの出力を各音素毎の継続時間50に従って制御して、画像生成部51で画像を生成する。

また、上述した動画像生成装置のほかに、音声と同期とした動画像生成を実現する手法として、例えば特開平4-155381号公報に記載の「音声による映像制御システム」（以下、文献2とする）がある。

文献2では、音または声を入力する音声入力部は、音または声の強弱及び高低を検出するフィルタを備え、音または声の有無、強弱、あるいは高低に対応して映像処理部の映像信号生成処理を制御し、映像処理部ではその制御に従って予め記憶した画像データを必要に応じて合成または加工して上記映像信号を生成する方法を提案している。

例えば物音や話し声、音楽等に合わせて映像を制御し、変化させることができ、TVモニタに映し出したキャラクタが聴覚を持つかのように振舞わせたりすることが可能となり、電話機を介して通話中の話し声に合わせた映像制御の実現を目標とする。

しかしながら、上述した文献1に示された従来技術では、音声特徴の種類によって画像の制御をすることで、音声と時間的に同期した動画像生成ができるものの、口唇画像のみの制御であり、韻律の違いや言語的文脈によって異なるような多様な画像を提供したり、口唇以外の表情や身振りなどの画像生成を実現することができない。

また、文献2に示される従来技術では、音声のピッチやパワーなどの韻律情報を用いて動画像を制御する手法であるが、対象が通話音声等の実音声に限定されている。すなわち、文の朗読や音声対話処理などの規則合成技術に基づく音声合成に対する考察がなされていない。

この発明の目的は、これらの課題を克服し、入力文から得られた音響的情報のみならず、韻律的情報や、入力文における言語的文脈の情報を積極的に利用し、人間が実際に音声を話しているようなより自然な動画像の提供することができる動画像生成装置を得ることを目的とする。

また、文の言語情報または韻律情報とそれを発生する人間の表情や動作との相関をモデル化し、文の文脈的特徴を有効に利用することにより自然な情報提供を行う動画像生成装置の実現を可能とするネットワークの学習を、簡便な学習方式によって行うことができる動画像制御ネットワーク学習装置を得ることを目的とするものである。

発明の開示

この発明に係る動画像生成装置は、入力されたテキストを解析して、単語や文節の読み、品詞、アクセント型などの言語情報を得る言語解析手段と、上記言語情報を基に音響パラメータを決定する音響処理手段と、上記言語情報を基にピッチや継続時間長やポーズ等の韻律パラメータを決定する韻律生成手段と、上記音響パラメータと上記韻律パラメータにより合成音を生成する音声合成手段と、韻律パラメータに対応した顔や人体の画像を記憶する画像記憶手段と、上記韻律生成手段から出力される韻律パラメータを用いて上記画像記憶手段を参照して出力

画像を決定する画像制御手段と、上記出力画像を表示する画像表示手段とを備えたものである。

また、上記画像記憶手段は、音韻と韻律パラメータに対応した顔や人体の画像を記憶すると共に、上記画像制御手段は、上記韻律生成手段から出力される韻律パラメータと上記言語解析手段から出力される言語情報の音韻とを用いて上記画像記憶手段を参照して出力画像を決定することを特徴とするものである。

また、上記言語解析手段から出力される上記言語情報を基に文中の言語情報のコンテキストを記憶するコンテキスト記憶手段をさらに備え、上記画像記憶手段は、コンテキストに対応した顔や人体の画像を記憶すると共に、上記画像制御手段は、上記コンテキスト記憶手段から出力されるコンテキストを用いて上記画像記憶手段を参照して出力画像を決定することを特徴とするものである。

また、上記言語解析手段から出力される上記言語情報を基に文中の言語情報のコンテキストを記憶するコンテキスト記憶手段をさらに備え、上記韻律生成手段は、当該コンテキスト記憶手段からの言語情報のコンテキストを基にピッチや継続時間長やポーズ等の韻律パラメータを決定することを特徴とするものである。

また、人間の表情や動作に関する状態と、制御カテゴリーとしての韻律パラメータと、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワークとして記憶するネットワーク記憶手段と、上記韻律生成手段からの韻律パラメータを入力して上記ネットワーク記憶手段を参照して上記韻律パラメータに対するネットワークを読み出して、ネットワークが該韻律パラメータを制御カテゴリーとして受理した際にネットワークの状態を遷移し、遷移時に画像制御パラメータを出力するネットワーク参照手段とをさらに備え、上記画像記憶手段は、画像制御パラメータに対応した顔や人体の画像を記憶すると共に、上記画像制御手段は、上記ネットワーク参照手段からの画像制御パラメータを用いて上記画像記憶手段を参照して出力画像を決定することを特徴とするものである。

また、他の発明に係る動画像生成装置は、入力されたテキストを解析して、単

語や文節の読み、品詞、アクセント型などの言語情報を得る言語解析手段と、上記言語情報を基に音響パラメータを決定する音響処理手段と、上記言語情報を基にピッチや継続時間長やポーズ等の韻律パラメータを決定する韻律生成手段と、上記音響パラメータと上記韻律パラメータにより合成音を生成する音声合成手段と、上記言語情報と、人間の顔や人体の画像の対応関係を予め学習して記憶し、上記言語解析手段から得た言語情報によりこの対応関係を参照して画像を出力する学習画像生成手段と、出力画像を表示する画像表示手段とを備えたものである。

また、上記学習画像生成手段は、言語情報と人間の顔や人体の画像との対応関係をテーブルとして記憶する画像テーブル記憶手段と、上記画像テーブル記憶手段を参照して上記言語情報に対応する画像を出力する画像テーブル参照手段とを備えたことを特徴とするものである。

また、上記学習画像生成手段は、言語情報と人間の表情や動作を示す画像の様態を示す画像制御パラメータとの対応関係をテーブル化して記憶するパラメータテーブル記憶手段と、上記パラメータテーブル記憶手段を参照して上記言語情報に対応する画像制御パラメータを出力するパラメータテーブル参照手段と、画像制御パラメータに対応した顔や人体の画像を記憶する画像記憶手段と、上記パラメータテーブル参照手段から出力される画像制御パラメータを用いて上記画像記憶手段を参照して出力画像を決定する画像制御手段とを備えたことを特徴とするものである。

また、上記一つの言語情報に対応して複数の動作が同時に起きる動作同期確率を記述したテーブルを記憶する動作同期確率記憶手段をさらに備え、上記パラメータテーブル参照手段は、上記言語解析手段から出力される言語情報に基づき上記パラメータテーブル記憶手段および上記動作同期確率記憶手段を参照して複数の画像制御パラメータを出力することを特徴とするものである。

また、上記学習画像生成手段は、人間の表情や動作に関する状態と、制御カテゴリーとしての言語情報と、制御カテゴリーが受理され状態遷移する際の遷移先

状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワークとして記憶するネットワーク記憶手段と、上記ネットワーク記憶手段を参照して上記言語情報に対するネットワークを読み出して、ネットワークが該言語情報を制御カテゴリーとして受理した際にネットワークの状態を遷移し、遷移時に画像制御パラメータを出力するネットワーク参照手段と、画像制御パラメータに対応した顔や人体の画像を記憶する画像記憶手段と、上記ネットワーク参照手段からの画像制御パラメータを用いて上記画像記憶手段を参照して出力画像を決定する画像制御手段とを備えたことを特徴とするものである。

さらに、この発明に係る画像制御ネットワーク学習装置は、文の言語情報または韻律パラメータを制御カテゴリーとして、その制御パラメータと文を発声している人間の顔や人体の画像パラメータとを対応づけて学習画像データとして記憶する学習画像データ記憶手段と、人間の表情や動作に関する状態と、制御カテゴリーとしての言語情報または韻律パラメータと、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態の画像パラメータとの組をネットワークとして記憶するネットワーク記憶手段と、上記学習画像データ記憶手段から読み出した学習画像データを入力とし、上記ネットワーク記憶手段から学習画像データの制御カテゴリーによって現在保持しているネットワークの状態からの状態遷移を読み出し、状態を遷移先状態に更新して保持し、上記状態遷移に含まれる画像パラメータと学習画像データの画像パラメータとの間の誤差を計算して出力する誤差計算手段と、状態遷移毎の上記誤差を基に学習対象の状態遷移を一つ決定し、学習対象アークとする学習対象アーク決定手段と、上記学習対象アークが受理可能とする制御カテゴリーの集合を分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する遷移分割手段と、上記学習対象アークの遷移元状態または遷移先状態を分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する状態分割手段と、上記遷

移分割手段から出力される誤差と上記状態分割手段から出力される誤差とを判定し、いずれかの分割手段により更新された状態遷移を選択して、上記ネットワーク記憶手段に記憶されているネットワーク内の状態遷移を書き換える最適分割決定手段とを備えたものである。

図面の簡単な説明

- 図1は、この発明の実施の形態1に係る動画像生成装置を示す構成図、
図2は、この発明の実施の形態2に係る動画像生成装置を示す構成図、
図3は、この発明の実施の形態3に係る動画像生成装置を示す構成図、
図4は、この発明の実施の形態4に係る動画像生成装置を示す構成図、
図5は、この発明の実施の形態5に係る動画像生成装置を示す構成図、
図6は、この発明の実施の形態6に係る動画像生成装置を示す構成図、
図7は、この発明の実施の形態7に係る動画像生成装置を示す構成図、
図8は、この発明の実施の形態8に係る動画像生成装置を示す構成図、
図9は、この発明の実施の形態9に係る動画像生成装置を示す構成図、
図10は、この発明の実施の形態10に係る画像制御ネットワーク学習装置を示す構成図、
図11は、従来例に係る動画像生成装置を示す構成図、
図12は、この発明の実施の形態8及び9に係る動画像生成装置におけるネットワークの一例を示す説明図、
図13は、この発明の実施の形態10に係る動画像生成装置における遷移分割手段から出力される遷移分割事後学習対象アークを含むネットワークの一例を示す説明図、
図14は、この発明の実施の形態10に係る動画像生成装置における状態分割手段の動作の一例を示す説明図、
図15は、この発明の実施の形態10に係る動画像生成装置における状態分割手段から出力される状態分割事後学習対象アークを含むネットワークの一例を示す説明図、

す説明図である。

発明を実施するための最良の形態

実施の形態 1.

図 1 は実施の形態 1 に係る動画像生成装置を示す構成図である。

図 1 において、1 はテキスト、2 は入力されたテキスト 1 を解析して、単語や文節の読み、品詞、アクセント型などの言語情報 3 を得る言語解析手段、4 は合成パラメータ生成手段、5 は言語情報 3 を基に音響辞書を参照して音響パラメータ 6 を決定する音響処理手段、7 は言語情報 3 を基に各種の規則を適用してピッチや継続時間長やポーズ等の韻律パラメータ 8 を決定する韻律生成手段である。

また、9 は音響パラメータ 6 と韻律パラメータ 8 により合成音を生成する音声合成手段、10 は韻律パラメータ 8 に対応する顔や人体の画像を記憶する画像記憶手段、12 は画像記憶手段 10 を参照して韻律パラメータ 8 により画像 11 を検索して出力する動画像 13 を決める画像制御手段、14 は出力画像を表示する画像表示手段である。

次に動作について説明する。

テキスト 1 は、日本語漢字仮名混じり文からなり、言語解析手段 2 へ入力される。このテキスト 1 の入力、キーボードからのキー入力や、音声認識装置が出力する認識結果などどのような入力形態でも構わない。

言語解析手段 2 は、単語の読み、品詞、アクセント型などを記憶する辞書を参照して、テキストに含まれる文節を決定し、文節毎の読み、品詞、アクセント型などを抽出する形態素解析処理を行い言語情報 3 を出力する。

言語解析手段 2 から出力される言語情報 3 は、規則合成を行う為に必要な単語や文節の読み、品詞、アクセント型などの情報から成る。

合成パラメータ生成手段 4 は、音響処理手段 5 と韻律生成手段 7 から成る。

ここで、音響処理手段 5 は、上記言語情報 3 を基に音響辞書を参照して、単語や文節の読みに対応した音響パラメータ 6 を抽出する。音響パラメータ 6 は、音

声波形の形態あるいは音声からあらかじめ抽出したLPC係数、メルケプストラム、LSP係数などの特徴量の形態で、CV（子音・母音の組み）、VCV（母音・子音・母音の組み）などの単位からなる音声素片の形で用意され、音響辞書に記憶されている。

また、韻律生成手段7は、上記言語情報3を基にピッチパターン生成規則、継続時間長生成規則、ポーズ挿入規則などの韻律制御規則を適用して、ピッチパターン、継続時間長、ポーズ挿入位置などの韻律パラメータ8を決定し、出力する。

音声合成手段9は、音響パラメータ6と韻律パラメータ8により合成音声を生成し、スピーカ等により出力する。

画像記憶手段10は、人間の表情や動作を示す画像11が韻律パラメータ8に対応して記憶されている。例えば高いピッチに対しては口の開きが大きく、目を見開いている画像、低ピッチに対しては口や目が閉じ気味の画像、ポーズに対しては口を閉じている画像などが記憶されている。

画像制御手段12は、上記韻律生成手段7から出力された韻律パラメータ8に基づき、上記画像記憶手段10を参照して韻律パラメータ8に合致した画像11を読み出す。そして、韻律パラメータ8の入力時間に同期させて画像11を提示することで、動画像13を出力する。

画像表示手段14は、上記動画像13をモニター等に表示する。

以上のように、この動画像生成方式は、テキストから得られた韻律パラメータにより画像制御を行い、合成音声のピッチの高低等に対応する動画像を生成するものである。

このような構成をとることにより、テキスト1から得られた韻律パラメータ8を利用して動画像生成を行うことができ、合成音声の韻律的特徴と同期した動画像を生成することで、より自然な情報提示を実現できる。

実施の形態2.

次に、図2は実施の形態2に係る動画像生成装置を示す構成図である。

図2において、図1に示す実施の形態1と同一部分は同一符号を付して示し、その説明は省略する。新たな符号として、15は言語解析手段2から出力される言語情報3の単語や文節の読み（音韻）を示す。

また、この実施の形態2において、画像記憶手段10は、音韻と韻律パラメータに対応した顔や人体の画像を記憶しており、画像制御手段12は、韻律生成手段7から出力される韻律パラメータ8と言語解析手段2から出力される言語情報3の読み15とを用いて画像記憶手段10を参照して出力画像を決定するようになされている。

次に動作について説明する。

日本語漢字仮名混じり文からなる入力されたテキスト1は言語解析手段2へ入力される。テキストは実施の形態1と同じものとする。言語解析手段2は、辞書を参照して、テキストに含まれる文節を決定し、文節毎の読み、品詞、アクセント型などを抽出する形態素解析処理を行う。

言語情報3は言語解析手段2から出力され、規則合成を行う為に必要な単語や文節の読み、品詞、アクセント型などの情報から成る。合成パラメータ生成手段4は音響処理手段5と韻律生成手段7から成る。

音響処理手段5は上記言語情報3を基に音響辞書を参照して、単語や文節の読みに対応した音響パラメータ6を抽出する。音響パラメータ6は、実施の形態1と同じものとする。韻律生成手段7は上記言語情報3を基に韻律制御規則を適用して、韻律パラメータ8を決定し、出力する。

音声合成手段9では音響パラメータ6と韻律パラメータ8により合成音声を生成し、スピーカ等により出力する。

画像記憶手段10では人間の表情や動作を示す画像11が、全ての音韻と韻律パラメータ8に対応して記憶されている。例えば一つの音韻“a”について、高いピッチに対しては口の開きが大きく、目を見開いている画像、低ピッチに対しては口や目が閉じ気味の画像、ポーズに対しては口を閉じている画像などが記憶

されている。

画像制御手段12では、上記韻律生成手段7から出力された韻律パラメータ8と上記言語解析手段2から得られたテキストの読み15に基づき、上記画像記憶手段10を参照して韻律パラメータ8に合致した画像11を読み出す。そして、韻律パラメータ8の入力時間に同期させて画像11を提示することで、動画像13を出力する。

画像表示手段14は上記動画像13をモニター等に表示する。

以上のように、この動画像生成装置は、テキストの読み15とテキストから得られた韻律パラメータ8により画像制御を行い、合成音声のピッチの高低等に対応する動画像を生成するものである。

すなわち、この装置において、画像記憶手段10は、音韻の種類と韻律パラメータに対応する顔や人体の画像を記憶する。画像制御手段12は、画像記憶手段10を参照して韻律パラメータ8とテキストの読み15により画像を検索し、出力する動画像を決める。

このような構成をとることにより、テキストの読み15とテキストから得られた韻律パラメータ8を利用して動画像生成を行うことができ、合成音声の音韻性及び韻律的特徴と同期した動画像を生成することで、より自然な情報提示を実現できる。

実施の形態3.

次に、図3は実施の形態3に係る動画像生成装置を示す構成図である。

図3において、図1に示す実施の形態1と同一部分は同一符号を付して示し、その説明は省略する。新たな符号として、16は言語解析手段2から出力される言語情報3を基に文中の言語情報のコンテキストを記憶するコンテキスト記憶手段であり、17は言語的コンテキストである。

また、この実施の形態3において、画像記憶手段10は人間の表情や動作を示す画像11が言語的コンテキスト17に対応して記憶されており、画像制御手段

12は、画像記憶手段10を参照してコンテキスト記憶手段16から出力される言語的コンテキスト17を用いて画像を検索し、出力する動画像13を決めるようになされている。

次に動作について説明する。

日本語漢字仮名混じり文からなる入力されたテキスト1は言語解析手段2へ入力される。テキストは実施の形態1と同じものとする。言語解析手段2は、辞書を参照して、テキストに含まれる文節を決定し、文節毎の読み、品詞、アクセント型などを抽出する形態素解析処理を行う。言語情報3は言語解析手段2から出力され、規則合成を行う為に必要な単語や文節の読み、品詞、アクセント型などの情報から成る。

合成パラメータ生成手段4は音響処理手段5と韻律生成手段7から成る。音響処理手段5は上記言語情報3を基に音響辞書を参照して、単語や文節の読みに対応した音響パラメータ6を抽出する。音響パラメータ6は、実施の形態1と同じものとする。韻律生成手段7は上記言語情報3を基に韻律制御規則を適用して、韻律パラメータ8を決定し、出力する。音声合成手段9では音響パラメータ6と韻律パラメータ8により合成音声を生成し、スピーカ等により出力する。

コンテキスト記憶手段16は、一文単位のテキスト1から言語解析手段2において得られた言語情報3の系列を言語的コンテキスト17として記憶する。言語的コンテキスト17は、一文の先頭から文末までの文節毎の読み、品詞、アクセント型の組を時系列順に記憶したものとする。

画像記憶手段10では人間の表情や動作を示す画像11が言語的コンテキスト17に対応して記憶されている。例えば一つの文節における先行文節および当該文節の言語情報に対してそれに対応する画像が一つ与えられている。具体的には、品詞が名詞から名詞へ変化している場合にはやや上向きの画像、名詞から動詞へ変化する場合はややうつむき加減の画像、文末の単語では下を向いている画像などである。

画像制御手段12では、上記コンテキスト記憶手段16から出力された言語コ

ンテキスト17に基づき上記画像記憶手段10を参照して言語的情報の変遷韻律に沿った画像11を読み出す。そして、韻律パラメータ8の入力時間に同期させて画像11を提示することで、動画像13を出力する。画像表示手段14は上記動画像13をモニター等に表示する。

以上のように、この動画像生成装置は、テキストから得られた言語情報3のコンテキストにより画像制御を行い、文節間の言語的特徴の文脈から動画像を生成するものである。

すなわち、この装置において、コンテキスト記憶手段16は、一文のテキストの言語情報を記憶する。画像記憶手段10は、言語的コンテキスト17に対応する顔や人体の画像を記憶する。また、画像制御手段12は、画像記憶手段10を参照して言語的コンテキスト17により画像を検索し、出力する動画像13を決める。

このような構成をとることにより、テキストの文節間の言語的特徴の文脈から動画像を生成することで、文の文脈的特徴を反映した、より自然な情報提示を実現できる。

実施の形態4.

次に、図4は実施の形態4に係る動画像生成装置を示す構成図である。

図4において、図3に示す実施の形態3と同一部分は同一符号を付して示し、その説明は省略する。この実施の形態4においては、コンテキスト記憶手段16は、言語解析手段2と韻律生成手段7との間に設けられていて、一文単位のテキスト1から言語解析手段2において得られた言語情報3の系列を言語的コンテキスト17として記憶する。

また、韻律生成手段7は、言語的コンテキスト17を基に韻律制御規則を適用してピッチや継続時間やポーズ等の韻律パラメータ8を決定して出力し、画像記憶手段10は、言語的コンテキスト17に対応する顔や人体の画像を記憶し、画像制御手段12は、画像記憶手段10を参照して韻律パラメータ8により画像を

検索して出力する動画像13を決めるようになされている。

次に動作について説明する。

日本語漢字仮名混じり文からなる入力されたテキスト1は言語解析手段2へ入力される。テキストは実施の形態1と同じものとする。言語解析手段2は、辞書を参照して、テキストに含まれる文節を決定し、文節毎の読み、品詞、アクセント型などを抽出する形態素解析処理を行う。言語情報3は言語解析手段2から出力され、規則合成を行う為に必要な単語や文節の読み、品詞、アクセント型などの情報から成る。

コンテキスト記憶手段16は、一文単位のテキスト1から言語解析手段2において得られた言語情報3の系列を言語的コンテキスト17として記憶する。言語コンテキスト17は実施の形態3と同じものとする。

合成パラメータ生成手段4は音響処理手段5と韻律生成手段7から成る。音響処理手段5は上記言語情報3を基に音響辞書を参照して単語や文節の読みに対応した音響パラメータ6を抽出する。音響パラメータ6は、実施の形態1と同じものとする。

韻律生成手段7は上記言語的コンテキスト17を基に韻律制御規則を適用して、韻律パラメータ8を決定し、出力する。例えば、数量化I類などの統計的学習手法により言語コンテキスト17と点ピッチモデルにおけるアクセント成分の強さとの相関を予め求めておき、入力された言語コンテキスト17に対するアクセント成分の強さを決定する。

音声合成手段9では音響パラメータ6と韻律パラメータ8により合成音声を生成し、スピーカ等により出力する。

画像記憶手段10では人間の表情や動作を示す画像11が韻律パラメータ8に対応して記憶されている。韻律パラメータ8が記憶される形式は、実施の形態1と同じとする。

画像制御手段12では、上記韻律生成手段7から出力された韻律パラメータ8に基づき上記画像記憶手段10を参照して韻律パラメータ8に合致した画像11

を読み出す。そして、韻律パラメータ 8 の入力時間に同期させて画像 1 1 を提示することで、動画像 1 3 を出力する。

画像表示手段 1 4 は上記動画像 1 3 をモニター等に表示する。

以上のように、この動画像生成装置は、テキスト 1 から得られた言語情報のコンテキストにより画像制御を行い、文節間の言語的特徴の文脈から韻律と同時に動画像を生成するものである。

すなわち、この装置において、コンテキスト記憶手段 1 6 は、一文のテキストの言語情報を記憶する。画像記憶手段 1 0 は、言語的コンテキスト 1 7 に対応する顔や人体の画像を記憶する。画像制御手段 1 2 は、画像記憶手段 1 0 を参照して韻律パラメータ 8 により画像を検索し、出力する動画像 1 3 を決める。

このような構成をとることにより、テキストの文節間の言語的特徴の文脈により韻律を制御し、その韻律により動画像を生成することで、文の文脈的特徴を反映した韻律と動画像の生成を同時に行うことができ、より自然な情報提示を実現できる。

実施の形態 5.

次に、図 5 は実施の形態 5 に係る動画像生成装置を示す構成図である。

図 5 において、図 1 に示す実施の形態 1 と同一部分は同一符号を付して示し、その説明は省略する。新たな符号として、1 8 は統計的手法によりテキストの言語情報と人間の顔や人体の画像との対応関係を画像テーブル 1 9 として記憶する画像テーブル記憶手段、2 0 は画像テーブル記憶手段 1 8 を参照して言語情報に対応する画像テーブル 1 9 を検索して出力する動画像 1 3 を決める画像テーブル参照手段である。

次に動作について説明する。

日本語漢字仮名混じり文からなる入力されたテキスト 1 は言語解析手段 2 へ入力される。テキストは実施の形態 1 と同じものとする。言語解析手段 2 は、辞書を参照して、テキストに含まれる文節を決定し、文節毎の読み、品詞、アクセン

ト型などを抽出する形態素解析処理を行う。言語情報3は言語解析手段2から出力され、規則合成を行う為に必要な単語や文節の読み、品詞、アクセント型などの情報から成る。

合成パラメータ生成手段4は音響処理手段5と韻律生成手段7から成る。音響処理手段5は上記言語情報3を基に音響辞書を参照して、単語や文節の読みに対応した音響パラメータ6を抽出する。音響パラメータ6は、実施の形態1と同じものとする。韻律生成手段7は上記言語情報3を基に韻律制御規則を適用して、韻律パラメータ8を決定し、出力する。

音声合成手段9では音響パラメータ6と韻律パラメータ8により合成音声を生成し、スピーカ等により出力する。

画像テーブル記憶手段18は、言語情報3と、人間の表情や動作を示す画像11の対応関係をテーブル化した画像テーブル19を記憶している。画像テーブル19の形式としては、例えば一つの文節における言語情報3に対応する画像が一つ与えられている。画像テーブル19は言語情報3と人間の表情や動作との関連を統計的に学習して得るものとする。学習方法としては、会話をしている一人の人間の画像を文節毎に蓄積し、各画像を文節の言語情報と共に記憶しておく。そして、一つの言語情報に対応した画像集合から一つの画像を選択し、言語情報との対応関係をテーブルの形式で保存する。選択画像は画像集合を平均化したものを用いる。

画像テーブル参照手段20は、言語解析手段2において得られた言語情報3により、画像テーブル記憶手段18に記憶されている画像テーブル19を参照して画像11を提示する。そして、画像11を言語情報の提示順に生成することで、動画像13を出力する。

画像表示手段14は上記動画像13をモニター等に表示する。

以上のように、この動画像生成装置は、テキストから得られた言語情報により統計的手法によって作成された画像テーブルを参照して画像制御を行い、統計的見地を反映した動画像を提供可能とするものである。

すなわち、この装置において、画像テーブル記憶手段18は、統計的手法によりテキストの言語情報と人間の顔や人体の画像との対応関係をテーブルとして記憶する。画像テーブル参照手段20は、画像テーブル記憶手段18を参照して言語情報により画像テーブル19を検索し、出力する画像を決める。

このような構成をとることにより、統計的手法によってテキストの文節の言語的特徴と対応付けられた画像を提示することが可能となり、統計的見地を反映した動画像を提供することにより、より自然な情報提示を実現できる。

実施の形態6.

次に、図6は実施の形態6に係る動画像生成装置を示す構成図である。

図6において、図1に示す実施の形態1と同一部分は同一符号を付して示し、その説明は省略する。新たな符号として、21は統計的手法によりテキストの言語情報3と人間の顔や人体の様態を示す画像制御パラメータ24との対応関係をパラメータテーブルとして記憶するパラメータテーブル記憶手段、23はパラメータテーブル記憶手段21を参照して言語情報によりパラメータテーブルを検索して画像制御パラメータ24を出力するパラメータテーブル参照手段である。

また、この実施の形態6において、画像記憶手段10は、表情や動作の様態の画像制御パラメータ対応した人間の顔や人体の画像を記憶し、画像制御手段12は、画像記憶手段10を参照して画像制御パラメータに対応した画像を検索して出力するようになっている。

次に動作について説明する。

日本語漢字仮名混じり文からなる入力されたテキスト1は言語解析手段2へ入力される。テキストは実施の形態1と同じものとする。言語解析手段2は、辞書を参照して、テキストに含まれる文節を決定し、文節毎の読み、品詞、アクセント型などを抽出する形態素解析処理を行う。言語情報3は言語解析手段2から出力され、規則合成を行う為に必要な単語や文節の読み、品詞、アクセント型などの情報から成る。

合成パラメータ生成手段4は音響処理手段5と韻律生成手段7から成る。音響処理手段5は上記言語情報3を基に音響辞書を参照して、単語や文節の読みに対応した音響パラメータ6を抽出する。音響パラメータ6は、実施の形態1と同じものとする。韻律生成手段7は上記言語情報3を基に韻律制御規則を適用して、韻律パラメータ8を決定し、出力する。

音声合成手段9では音響パラメータ6と韻律パラメータ8により合成音声を生成し、スピーカ等により出力する。

パラメータテーブル記憶手段21は、言語情報3と、人間の表情や動作を示す画像の様態を示すパラメータとの対応関係をテーブル化したパラメータテーブル22を記憶している。形式としては、例えば一つの文節における言語情報に対応する人間の表情や動作の様態のパラメータが与えられている。このパラメータテーブル22は、言語情報3と人間の表情や動作の様態との関連を統計的に学習して得るものとする。人間の表情や動作の様態は、例えば人間の頭が前に傾く角度、口唇の幅、上下のまぶたの間隔などを数量化したものとする。学習方法としては、会話をしている一人の人間の表情や動作の様態のパラメータを文節毎に蓄積し、各パラメータを文節の言語情報と共に記憶しておく。そして、一つの言語情報に対応したパラメータ集合から一つを選択し、言語情報との対応関係をテーブルの形式で保存する。パラメータの選択基準としては、パラメータ集合の平均を用いる。

パラメータテーブル参照手段23は言語解析手段2において得られた言語情報3により、パラメータテーブル記憶手段21に記憶されているパラメータテーブル22を参照して画像制御パラメータ24を出力する。

画像記憶手段10は人間の表情や動作の画像11を、動作や表情の様態のパラメータと共に記憶している。画像制御手段12では、パラメータテーブル参照手段23から出力された画像制御パラメータ24によって上記画像記憶手段10に記憶されている画像11を読み出す。そして、画像11を言語情報の提示順に生成することで、動画像13を出力する。

画像表示手段14は上記動画像13をモニター等に表示する。

以上のように、この動画像生成装置は、テキストから得られた言語情報により統計的手法によって作成されたパラメータテーブルを参照して画像制御を行い、統計的見地を反映しかつ顔や人体の部位などの詳細な動きを制御する動画像を提供可能とするものである。

すなわち、この装置において、パラメータテーブル記憶手段21は、統計的手法によりテキストの言語情報と人間の顔や人体の様態を示すパラメータとの対応関係をテーブルとして記憶する。パラメータテーブル参照手段23は、パラメータテーブル記憶手段21を参照して言語情報によりパラメータテーブル22を検索し、画像制御パラメータ24を出力する。画像記憶手段10は、人間の顔や人体の画像を、表情や動作の様態のパラメータと共に記憶する。画像制御手段12は、画像記憶手段10を参照して画像制御パラメータ24により画像を検索して出力する。

このような構成をとることにより、統計的手法によってテキストの文節の言語的特徴と人間の動作や表情の様態とを対応付けて、より詳細な画像を提示することが可能となり、統計的見地を反映した動画像を提供することにより、より自然な情報提示を実現できる。

実施の形態7.

次に、図7は実施の形態7に係る動画像生成装置を示す構成図である。

図7において、図6に示す実施の形態6と同一部分は同一符号を付してその説明は省略する。新たな符号として、25は1つの言語情報3に対応して複数の動作が同時に起きる動作同時確立26を記述したテーブルを記憶する動作同時確立記憶手段である。そして、この実施の形態7において、パラメータテーブル参照手段23は、動作同期確率記憶手段25およびパラメータテーブル記憶手段21を参照して、言語情報3により複数の人体のパーツにおける動作同期確率26を求め、パラメータテーブル22を検索して上記動作同期確率の閾値が一定値を越

すパーツについて画像制御パラメータ 2 4 を出力するようになされている。

次に動作について説明する。

日本語漢字仮名混じり文からなる入力されたテキスト 1 は言語解析手段 2 へ入力される。テキストは実施の形態 1 と同じものとする。言語解析手段 2 は、辞書を参照して、テキストに含まれる文節を決定し、文節毎の読み、品詞、アクセント型などを抽出する形態素解析処理を行う。言語情報 3 は言語解析手段 2 から出力され、規則合成を行う為に必要な単語や文節の読み、品詞、アクセント型などの情報から成る。

合成パラメータ生成手段 4 は音響処理手段 5 と韻律生成手段 7 から成る。音響処理手段 5 は上記言語情報 3 を基に音響辞書を参照して、単語や文節の読みに対応した音響パラメータ 6 を抽出する。音響パラメータ 6 は、実施の形態 1 と同じものとする。韻律生成手段 7 は上記言語情報 3 を基に韻律制御規則を適用して、韻律パラメータ 8 を決定し、出力する。

音声合成手段 9 では音響パラメータ 6 と韻律パラメータ 8 により合成音声を生成し、スピーカ等により出力する。

動作同期確率記憶手段 2 5 は、一つの言語情報が与えられた際に、複数の動作が同時に起きる動作同時確率 2 6 を記述したテーブルを記憶する。例えばテーブルには頭の動き、目の動き、口の動きなどの動作が詳細にカテゴライズされ、言語情報に対しての各カテゴリーの動作が行なわれる確率が記述されている。このテーブルの学習は会話をしている人間の画像を収集し、言語情報毎に人間の動作や表情に関わる顔や人体の部位のデータを蓄積し、動作や表情の平均発生確率を求める手法による。

パラメータテーブル記憶手段 2 1 は、言語情報 3 と、人間の表情や動作を示す画像の様態を示すパラメータとの対応関係をテーブル化したパラメータテーブル 2 2 を記憶している。パラメータテーブルの形式は実施の形態 6 と同じとする。

パラメータテーブル参照手段 2 3 は、言語解析手段 2 において得られた言語情報 3 により、まず、動作同期確率記憶手段 2 5 を参照し、言語情報 3 に対して人

体の種々のパーツが動作を起こす確率、すなわち動作同期確率 26 を読み出す。次に、パラメータテーブル記憶手段 21 に記憶されているパラメータテーブル 22 を参照し、動作同期確率 26 がある閾値を越すパーツについての画像制御パラメータ 24 を引出し、出力する。

画像記憶手段 10 は人間の表情や動作の画像 11 を、動作や表情の様態のパラメータと共に記憶している。

画像制御手段 12 ではパラメータテーブル参照手段 23 から出力された画像制御パラメータ 24 によって上記画像記憶手段 10 に記憶されている画像 11 を読み出す。そして、画像 11 を言語情報の提示順に生成することで、動画像 13 を出力する。

画像表示手段 14 は上記動画像 13 をモニター等に表示する。

以上のように、この動画像生成方式は、テキストから得られた言語情報により統計的手法によって作成されたパラメータテーブルを参照して画像制御を行ない、統計的見地を反映しかつ顔や人体の部位などの詳細な複数の動きを同時に制御する動画像を提供可能とするものである。

すなわち、この装置において、動作同期確率記憶手段 25 は、テキストの言語情報に対して人間の複数の動作や表情が起きる動作同時確率 26 を統計的に計算し、記憶する。パラメータテーブル記憶手段 21 は、統計的手法によりテキストの言語情報と人間の顔や人体の様態を示すパラメータとの対応関係をテーブルとして記憶する。パラメータテーブル参照手段 23 は、動作同期確率記憶手段 25 およびパラメータテーブル記憶手段 21 を参照して、言語情報 3 により複数の人体のパーツにおける動作同期確率 26 を求め、パラメータテーブル 22 を検索して上記動作同期確率 26 の閾値が一定値を越すパーツについて画像制御パラメータ 24 を出力する。

画像記憶手段 10 は、人間の顔や人体の画像を、表情や動作の様態のパラメータと共に記憶する。

画像制御手段 12 は、画像記憶手段 10 を参照して画像制御パラメータ 24 に

より画像を検索して出力する。

このような構成をとることにより、統計的手法によってテキストの文節の言語的特徴と人間の動作や表情の様態とを対応付けて、より詳細な画像を提示することが可能となり、さらに複数の動作や表情の画像の制御を同時に行うことができ、統計的見地を反映した動画像を提供することにより、より自然な情報提示を実現できる。

実施の形態 8.

次に、図 8 は実施の形態 8 に係る動画像生成装置を示す構成図である。

図 8 において、図 1 に示す実施の形態 1 と同一部分は同一符号を付してその説明は省略する。新たな符号として、27 は、人間の表情や動作に関する状態と、制御カテゴリーとしての言語情報 3 と、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワーク 28 として記憶するネットワーク記憶手段、29 は、ネットワーク記憶手段 27 を参照してテキストの言語情報 3 に対するネットワーク 28 を読み出して、ネットワーク 28 が該言語情報 3 を制御カテゴリーとして受理した際にネットワークの状態を遷移し、遷移時に画像制御パラメータ 24 を出力するネットワーク参照手段である。

そして、この実施の形態 8 において、画像記憶手段 10 は、人間の顔や人体の画像を、表情や動作の様態の画像制御パラメータ 24 に対応して記憶し、画像制御手段 12 は、画像記憶手段 10 を参照して画像制御パラメータ 24 により画像を検索して出力するようになされている。

次に動作について説明する。

日本語漢字仮名混じり文からなる入力されたテキスト 1 は言語解析手段 2 へ入力される。テキストは実施の形態 1 と同じものとする。言語解析手段 2 は、辞書を参照して、テキストに含まれる文節を決定し、文節毎の読み、品詞、アクセント型などを抽出する形態素解析処理を行う。言語情報 3 は言語解析手段 2 から出

力され、規則合成を行う為に必要な単語や文節の読み、品詞、アクセント型などの情報から成る。

合成パラメータ生成手段4は、音響処理手段5と韻律生成手段7から成る。音響処理手段5は上記言語情報3を基に音響辞書を参照して、単語や文節の読みに対応した音響パラメータ6を抽出する。音響パラメータ6は、実施の形態1と同じものとする。韻律生成手段7は上記言語情報3を基に韻律制御規則を適用して、韻律パラメータ8を決定し、出力する。

音声合成手段9では音響パラメータ6と韻律パラメータ8により合成音声を生成し、スピーカ等により出力する。

ネットワーク記憶手段27は、人間の表情や動作に関する状態と、制御カテゴリーとしての言語情報3と、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワーク28として記憶する。

ネットワークの構成例を図12に示す。

図12において、状態番号p、受理される制御カテゴリーとしての言語情報C、遷移先状態番号q、遷移時に出力される出力パラメータ α とする。同図(a)はネットワーク内部を示す表であり、同図(b)は状態間の関係を示す模擬図である。図12(b)中の小円が状態に、状態間の弧は状態遷移に相当する。ここでは、状態1は文頭、状態6は文末を示す。

ネットワーク参照手段29では、ネットワーク記憶手段27を参照して、入力された言語情報3と現在保持されている状態番号に対応するネットワーク28を検索する。そして、ネットワーク28に含まれる遷移先状態番号と出力パラメータを読み出し、上記遷移先状態番号を新たな状態番号へと更新して保持する。そして、上記出力パラメータを画像制御パラメータ24として出力する。例として、図12のネットワークに{c3、c2、c5}の順で言語情報系列が入力された場合は、2、6、9のネットワークが順番に参照され、出力される画像制御パラメータの系列は{ α 2、 α 6、 α 9}の順で出力される。状態番号が文末状態

の6となった場合は状態を文頭に戻し、状態番号は1となる。

画像記憶手段10は、人間の表情や動作の画像11を、動作や表情の様態のパラメータと共に記憶している。

画像制御手段12ではパラメータテーブル参照手段23から出力された画像制御パラメータ24によって上記画像記憶手段10に記憶されている画像11を読み出す。そして、画像11を言語情報の提示順に生成することで、動画像13を出力する。

画像表示手段14は、上記動画像13をモニター等に表示する。

以上のように、この動画像生成装置は、テキストから得られた言語情報によりネットワークを参照して画像制御を行い、簡便なモデルにより言語情報の文脈を反映した動画像を提供可能とするものである。

すなわち、この装置において、ネットワーク記憶手段27は、人間の表情や動作に関する状態と、制御カテゴリーとしての言語情報3と、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワーク28として記憶する。

ネットワーク参照手段29は、上記ネットワーク記憶手段27を参照してテキスト1の言語情報3に対するネットワーク28を読み出して、ネットワーク28が該言語情報3を制御カテゴリーとして受理した際にネットワークの状態を遷移し、遷移時に画像制御パラメータを出力する。

画像記憶手段10は、人間の顔や人体の画像を、表情や動作の様態のパラメータと共に記憶する。

画像制御手段12は、画像記憶手段10を参照して画像制御パラメータ24により画像を検索して出力する。

このような構成をとることにより、テキストの文節の言語的文脈の影響を反映した動画像の生成を行うことができ、簡便なモデルにより文の文脈的特徴を反映した、より自然な情報提示を実現できる。

実施の形態 9.

次に、図 9 は実施の形態 9 に係る動画像生成装置を示す構成図である。

図 9 において、図 8 に示す実施の形態 8 と同一部分は同一符号を付してその説明は省略する。

この実施の形態 9 においては、ネットワーク記憶手段 27 は、人間の表情や動作に関する状態と、制御カテゴリーとしての韻律パラメータ 8 と、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワーク 28 として記憶する。

また、ネットワーク参照手段 29 は、上記ネットワーク記憶手段 27 を参照してテキストの韻律パラメータ 8 に対するネットワーク 28 を読み出して、ネットワーク 28 が該韻律パラメータ 8 を制御カテゴリーとして受理した際にネットワークの状態を遷移し、遷移時に画像制御パラメータ 24 を出力するようになっている。

次に動作について説明する。

日本語漢字仮名混じり文からなる入力されたテキスト 1 は言語解析手段 2 へ入力される。テキストは実施の形態 1 と同じものとする。言語解析手段 2 は、辞書を参照して、テキストに含まれる文節を決定し、文節毎の読み、品詞、アクセント型などを抽出する形態素解析処理を行う。言語情報 3 は言語解析手段 2 から出力され、規則合成を行う為に必要な単語や文節の読み、品詞、アクセント型などの情報から成る。

合成パラメータ生成手段 4 は、音響処理手段 5 と韻律生成手段 7 から成る。音響処理手段 5 は上記言語情報 3 を基に音響辞書を参照して、単語や文節の読みに対応した音響パラメータ 6 を抽出する。音響パラメータ 6 は、実施の形態 1 と同じものとする。韻律生成手段 7 は上記言語情報 3 を基に韻律制御規則を適用して、韻律パラメータ 8 を決定し、出力する。

音声合成手段9では、音響パラメータ6と韻律パラメータ8により合成音声を生成し、スピーカ等により出力する。

ネットワーク記憶手段27は、人間の表情や動作に関する状態と、制御カテゴリーとしての韻律パラメータ8と、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワーク28として記憶する。

ネットワークの構成例を図12に示す。

ネットワークの構成は、受理される制御カテゴリーCが韻律パラメータ8であることの他は実施の形態8と同じである。制御カテゴリーとして韻律パラメータ8を用いるには、アクセントの大小により異なるインデックスを付けたり、ポーズ挿入を明示するなどの方法がある。

ネットワーク参照手段29では、ネットワーク記憶手段27を参照して、入力された韻律パラメータ8と現在保持されている状態番号に対応するネットワーク28を検索する。そして、ネットワーク28に含まれる遷移先状態番号と出力パラメータを読み出し、上記遷移先状態番号を新たな状態番号へと更新して保持する。そして、上記出力パラメータを画像制御パラメータ24として出力する。

画像記憶手段10は、人間の表情や動作の画像11を、動作や表情の様態のパラメータと共に記憶している。

画像制御手段12では、パラメータテーブル参照手段23から出力された画像制御パラメータ24によって上記画像記憶手段10に記憶されている画像11を読み出す。そして、画像11を言語情報の提示順に生成することで、動画像13を出力する。

画像表示手段14は、上記動画像13をモニター等に表示する。

以上のように、この動画像生成装置は、テキストから得られた韻律情報によりネットワークを参照して画像制御を行い、簡便なモデルにより言語情報の文脈を反映しかつ韻律情報と同期した動画像を提供可能とするものである。

すなわち、この装置において、ネットワーク記憶手段27は、人間の表情や動

作に関する状態と、制御カテゴリーとしての韻律パラメータ8と、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワーク28として記憶する。

ネットワーク参照手段29は、上記ネットワーク記憶手段27を参照してテキストの韻律パラメータ8に対するネットワーク28を読み出して、ネットワーク28が該韻律パラメータ8を制御カテゴリーとして受理した際にネットワークの状態を遷移し、遷移時に画像制御パラメータを出力する。

画像記憶手段10は、人間の顔や人体の画像を、表情や動作の様態のパラメータと共に記憶する。

画像制御手段12は、画像記憶手段10を参照して画像制御パラメータ24により画像を検索して出力する。

このような構成をとることにより、韻律情報の文脈と同期した動画像の生成を行なうことができ、簡便なモデルにより文の文脈的特徴を反映した、より自然な情報提示を実現できる。

実施の形態10.

次に、図10は実施の形態10に係る動画像生成装置における画像制御ネットワーク学習装置を示す構成図である。

図10において、27は実施の形態8及び9と同様なネットワーク記憶手段であり、人間の表情や動作に関する状態と、制御カテゴリーとしての言語情報3と、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワーク28として記憶している。

また、新たな符号として、30は文を発声している人間の顔や人体の画像を、発声している文の言語情報あるいは韻律情報と対応づけて学習画像データ31として記憶する学習画像データ記憶手段である。

32は上記学習画像データ記憶手段30から読み出した学習画像データ31を入力とし、上記ネットワーク記憶手段27から学習画像データの制御カテゴリから現在保持しているネットワーク28の状態からの状態遷移を読み出し、状態を遷移先状態に更新して保持し、上記状態遷移に含まれる画像パラメータと学習画像データの画像パラメータとの間の誤差データ33を計算して出力する誤差計算手段である。

34は全ネットワーク中の上記誤差データが最小となる状態遷移を一つ選択し、学習対象アーク35として決定する学習対象アーク決定手段である。

36は上記学習対象アーク35が受理可能とする制御カテゴリの集合を分割事後の誤差の和が最小となる組合せで分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する遷移分割手段である。

38は上記学習対象アーク35の分割事後の誤差の和が最小となるように、学習対象アーク35の遷移元状態へ遷移する状態遷移のカテゴリを二分することによって遷移元状態を分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する状態分割手段である。

40は上記遷移分割手段36から出力される誤差と上記状態分割手段38から出力される誤差とを判定し、いずれかの修正状態遷移を選択して、上記ネットワーク記憶手段27に記憶されている状態遷移を書き換える最適分割決定手段である。

なお、37は遷移分割後学習対象アーク、39は状態分割後学習対象アーク、41は書き換えネットワークである。

次に動作について説明する。

学習画像データ記憶手段30は、ネットワークの学習に用いる学習画像データ31を記憶する。学習画像データ31は、文を発声している人間の顔や人体の画像を、発声している文の言語情報あるいは韻律情報と対応づけてられている。

ネットワーク記憶手段27は、人間の表情や動作に関する状態と、制御カテゴ

リーとしての言語情報3と、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワーク28として記憶する。

誤差計算手段32は、上記学習画像データ記憶手段27から読み出した学習画像データ31を入力とし、上記ネットワーク記憶手段27から学習画像データの制御カテゴリーから現在保持しているネットワーク28の状態からの状態遷移を読み出し、状態を遷移先状態に更新して保持する。そして、上記状態遷移に含まれる画像パラメータと学習画像データの画像パラメータとの間の平均誤差を計算して、ネットワークと誤差の組からなる誤差データ33を出力する。

学習対象アーク決定手段34は、全ネットワーク中の誤差データ33を検索して、上記誤差が最大となる状態遷移を判定し、学習対象アーク35として出力する。

遷移分割手段36は、上記学習対象アークが受理可能とする制御カテゴリーの集合を分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する。

例えば、図12において、4番の状態遷移が学習対象アークとして判定された場合、まず、カテゴリ集合Cを可能な全ての分割法で2つのカテゴリ集合に分離する。4番の状態遷移では{c1}と{c2、c3}、{c2}と{c1、c3}、{c3}と{c1、c2}の3通りの分割となる。

次に、これら全ての分割法について、学習対象アークの誤差を分離された2つのカテゴリ集合毎に再計算する。誤差の再計算は、2つのカテゴリ集合内の画像パラメータの平均値を求めてそれぞれ $\alpha 4'$ 、 $\alpha 4''$ とし、カテゴリ集合毎に平均値と学習画像データの画像パラメータ値の差を集計し、その平均値を誤差として保持する。そして、誤差の和が最小の分割法を選択し、学習対象アークのカテゴリ集合を分割後の一方のカテゴリ集合に修正し、画像制御パラメータも分割後のカテゴリ集合のものに更新する。

また、もう一方の分離後のカテゴリ集合を受理する状態遷移を新たに作成し、

画像制御パラメータももう一方の分割後のカテゴリ集合のものにする。

例えば、カテゴリ集合を {c 1、c 3} と {c 2} に分離したときの誤差の和が最小であれば、図 1 3 のようにネットワークが更新される。

この様に更新された学習対象アークを遷移分割後学習対象アーク 3 7 として出力する。

状態分割手段 3 8 は、上記学習対象アークの遷移元状態を分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する。

例えば、図 1 2 において、4 番の状態遷移が学習対象アークとして判定された場合、まず、図 1 4 のように学習対象アークの状態番号に相当する状態、すなわち遷移元状態である状態 2 を分割対象とし、新たな状態 2' を用意する。ここで、状態 2 に遷移するすべての状態遷移を、カテゴリ集合の組合せにより二分し、状態 2 へ遷移する状態遷移と状態 2' へ遷移する状態遷移に分割する。

この状態の分割により、学習対象アークは、状態 2 から状態 4 への遷移と、状態 2' から状態 4 への遷移へと分離される。この時に、全ての状態の分割法について、学習対象アークの誤差を遷移元状態により分離された 2 つのグループ毎に再計算する。誤差の再計算は、2 つのグループ内の画像パラメータの平均値を求めてそれぞれ $\alpha 4'$ 、 $\alpha 4''$ とし、グループ毎に平均値と学習画像データの画像パラメータ値の差を集計し、その平均値を誤差として保持する。

そして、誤差の和が最小の状態分割法を選択し、分離された学習対象アークの画像制御パラメータをそれぞれ分割後の一方のグループのものに更新する。

例えば、ここでは先行カテゴリ集合は {c 1} と {c 2} に唯一分離したときの誤差の和が最小であれば、図 1 5 のようにネットワークが更新される。

この様に更新された学習対象アークを状態分割後学習対象アーク 3 9 として出力する。

最適分割決定手段 4 0 は、上記遷移分割手段 3 6 から出力される遷移分割後学習対象アーク 3 7 における誤差と、上記状態分割手段 3 8 から出力される状態分

割後学習対象アーク 38 における誤差を判定し、誤差の和が小さい方の分割後の学習対象アークを選択し、選択された分割手法により更新された状態遷移の集合を書き換えネットワーク 41 とし、上記ネットワーク記憶手段 27 に記憶されているネットワーク 28 を書き換える。

以上の動作を一定回数繰り返すことによって、画像制御ネットワークの学習を行う。

以上のように、この画像制御ネットワーク学習装置は、実施の形態 9 または実施の形態 10 の動画像生成装置で用いるネットワークの学習を行うものである。

すなわち、この装置において、学習画像データ記憶手段 30 は、文を発声している人間の顔や人体の画像を、発声している文の言語情報あるいは韻律情報と対応づけて共に記憶する。

ネットワーク記憶手段 27 は、人間の表情や動作に関する状態と、制御カテゴリーとしての言語情報または韻律パラメータと、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワークとして記憶する。

誤差計算手段 32 は、上記学習画像データ記憶手段 30 から読み出した学習画像データを入力とし、上記ネットワーク記憶手段 27 から学習画像データの制御カテゴリーから現在保持しているネットワークの状態からの状態遷移を読み出し、状態を遷移先状態に更新して保持し、上記状態遷移に含まれる画像パラメータと学習画像データの画像パラメータとの間の誤差を計算して出力する。

学習対象アーク決定手段 34 は、全ネットワーク中の上記誤差が最小となる状態遷移を一つ選択し、学習対象アークとして決定する。

遷移分割手段 36 は、上記学習対象アークが受理可能とする制御カテゴリーの集合を分割事後の誤差の和が最小となる組合せで分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する。

状態分割手段 38 は、上記学習対象アークの分割事後の誤差の和が最小となるように、学習対象アークの遷移元状態へ遷移する状態遷移のカテゴリーを二分する

ことによって遷移元状態を分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する。

最適分割決定手段40は、上記遷移分割手段36から出力される誤差と上記状態分割手段38から出力される誤差とを判定し、いずれかの修正状態遷移を選択して、上記ネットワーク記憶手段27に記憶されている状態遷移を書き換える。

このような構成をとることにより、文の言語情報または韻律情報とそれを発声する人間の表情や動作との相関をモデル化し、文の文脈的特徴を有効に利用したより自然な情報提示を行なう動画像生成装置の実現を可能とするネットワークの学習を、簡便な学習装置によって行うことができる。

その他の実施の形態A.

実施の形態1～9の入力であるテキスト1は、日本語の漢字仮名混じり文に限らず、仮名からなる文、あるいは英語文などの外国語文とすることもできる。

その他の実施の形態B.

実施の形態1～9における言語解析手段2は、形態素解析処理のみならず、文法を適用して文の統語的構造を抽出する統語解析処理や、意味的構造を抽出する意味解析処理を含めることができる。

その他の実施の形態C.

実施の形態2における画像記憶手段10は、全ての音韻種類に対応した画像を用意するのではなく、音韻種類を少数の音韻グループに分類し、その音韻グループおよび韻律パラメータに対応した画像を用意することもできる。

その他の実施の形態D.

実施の形態3および4におけるコンテキスト記憶手段16は、文頭を含む呼気段落の先頭の文節から文末を含む呼気段落の終端の文節までの文節毎の読み、品

詞、アクセント型の組を時系列順に記憶することができる。

その他の実施の形態E.

実施の形態3における画像記憶手段10は、一つの文節における当該文節および後続文節の言語情報に対して、それに対応する画像が一つ与えることができる。

その他の実施の形態F.

実施の形態3における画像記憶手段10は、一つの文節における先行文節、当該文節および後続文節の言語情報に対して、それに対応する画像が一つ与えることができる。

その他の実施の形態G.

実施の形態5における画像テーブル記憶手段18としては、画像テーブルの学習時において、一つの言語情報に対応した画像集合から一つの画像を選択する手法として、画像集合の中で比較的差分の小さい画像同士を小集合として収集し、その中で要素最大となる小集合から代表画像を一つ選択することもできる。

その他の実施の形態H.

実施の形態5において、画像テーブル記憶手段18における画像テーブルは画像と実施の形態3における言語的コンテキスト17との対応関係を示すものとし、画像テーブル参照手段20は、上記言語的コンテキスト17を入力として画像テーブルを検索することもできる。

その他の実施の形態I.

実施の形態6におけるパラメータテーブル記憶手段21としては、パラメータテーブルの学習時において、一つの言語情報に対応したパラメータ集合から代表

値を決定する手法として、言語情報を制御要因として数量化I類に代表される多変量解析法により得られた出力係数を用いることもできる。

その他の実施の形態J.

実施の形態6および7において、パラメータテーブル記憶手段21におけるパラメータテーブルは、画像制御パラメータ24と実施の形態3における言語的コンテキスト17との対応関係を示すものとし、パラメータテーブル参照手段23は、上記言語的コンテキスト17を入力としてパラメータテーブルを検索することもできる。

その他の実施の形態K.

実施の形態10におけるネットワーク記憶手段27は、実施の形態9で述べている韻律パラメータを制御カテゴリとして用いることができる。

その他の実施の形態L.

実施の形態10における状態分割手段38は、学習対象アークの遷移元状態を分割対象とするのみならず、学習対象アークの遷移先状態に後続する状態遷移の制御カテゴリーを二分することによる遷移先状態の分割も行い、最適分割決定手段40は、遷移分割手段38による学習対象アークの分割と、状態分割手段36による遷移元状態の分割と遷移先状態の分割の3つ分割手法の中から、最適な分割手法を選択することもできる。

その他の実施の形態M.

実施の形態10の画像制御ネットワーク学習装置において、学習終了の時点として一定回数の学習の終了時と設定する代わりに、学習対象アークの平均誤差が別途設定したしきい値を下回った時点を用いることができる。

その他の実施の形態N.

実施の形態10の画像制御ネットワーク学習装置において、学習終了の時点として一定回数の学習の終了時と設定する代わりに、全ネットワークの平均誤差が別途設定したしきい値を下回った時点を用いることができる。

その他の実施の形態O.

実施の形態10の誤差計算手段32において、平均誤差が最大の状態遷移を学習対象アークとする代わりに、状態遷移毎の学習画像データとの累積誤差が最大となる状態遷移を学習対象アークとすることもできる。

その他の実施の形態P.

実施の形態10の遷移分割手段36および状態分割手段38において、分割後の学習対象アークの誤差の総和が最小となるように分割する代わりに、分割後の最大誤差が最小となるように分割することもできる。

産業上の利用の可能性

以上のように、この発明に係る動画像生成装置は、入力文から得られた音響的情報のみならず、韻律的情報や、入力文における言語的文脈の情報を積極的に利用し、人間が実際に音声を話しているようなより自然な動画像を生成して、より自然な情報提示を実現できる。また、この発明に係る画像制御ネットワーク学習装置は、文の言語情報または韻律情報とそれを発声する人間の表情や動作との相関をモデル化し、文の文脈的特徴を有効に利用したより自然な情報提示を行う動画像生成装置の実現を可能とするネットワークの学習を、簡便な学習装置によって実現できる。

請 求 の 範 囲

1. 入力されたテキストを解析して、単語や文節の読み、品詞、アクセント型などの言語情報を得る言語解析手段と、

上記言語情報を基に音響パラメータを決定する音響処理手段と、

上記言語情報を基にピッチや継続時間長やポーズ等の韻律パラメータを決定する韻律生成手段と、

上記音響パラメータと上記韻律パラメータにより合成音を生成する音声合成手段と、

韻律パラメータに対応した顔や人体の画像を記憶する画像記憶手段と、

上記韻律生成手段から出力される韻律パラメータを用いて上記画像記憶手段を参照して出力画像を決定する画像制御手段と、

上記出力画像を表示する画像表示手段と

を備えた動画像生成装置。

2. 請求項1記載の動画像生成装置において、上記画像記憶手段は、音韻と韻律パラメータに対応した顔や人体の画像を記憶すると共に、上記画像制御手段は、上記韻律生成手段から出力される韻律パラメータと上記言語解析手段から出力される言語情報の音韻とを用いて上記画像記憶手段を参照して出力画像を決定することを特徴とする動画像生成装置。

3. 請求項1記載の動画像生成装置において、上記言語解析手段から出力される上記言語情報を基に文中の言語情報のコンテキストを記憶するコンテキスト記憶手段をさらに備え、上記画像記憶手段は、コンテキストに対応した顔や人体の画像を記憶すると共に、上記画像制御手段は、上記コンテキスト記憶手段から出力されるコンテキストを用いて上記画像記憶手段を参照して出力画像を決定することを特徴とする動画像生成装置。

4. 請求項1記載の動画像生成装置において、上記言語解析手段から出力される上記言語情報を基に文中の言語情報のコンテキストを記憶するコンテキスト記

憶手段をさらに備え、上記韻律生成手段は、当該コンテキスト記憶手段からの言語情報のコンテキストを基にピッチや継続時間長やポーズ等の韻律パラメータを決定することを特徴とする動画像生成装置。

5. 請求項1記載の動画像生成装置において、人間の表情や動作に関する状態と、制御カテゴリーとしての韻律パラメータと、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワークとして記憶するネットワーク記憶手段と、上記韻律生成手段からの韻律パラメータを入力して上記ネットワーク記憶手段を参照して上記韻律パラメータに対するネットワークを読み出して、ネットワークが該韻律パラメータを制御カテゴリーとして受理した際にネットワークの状態を遷移し、遷移時に画像制御パラメータを出力するネットワーク参照手段とをさらに備え、上記画像記憶手段は、画像制御パラメータに対応した顔や人体の画像を記憶すると共に、上記画像制御手段は、上記ネットワーク参照手段からの画像制御パラメータを用いて上記画像記憶手段を参照して出力画像を決定することを特徴とする動画像生成装置。

6. 入力されたテキストを解析して、単語や文節の読み、品詞、アクセント型などの言語情報を得る言語解析手段と、

上記言語情報を基に音響パラメータを決定する音響処理手段と、

上記言語情報を基にピッチや継続時間長やポーズ等の韻律パラメータを決定する韻律生成手段と、

上記音響パラメータと上記韻律パラメータにより合成音を生成する音声合成手段と、

上記言語情報と、人間の顔や人体の画像の対応関係を予め学習して記憶し、上記言語解析手段から得た言語情報によりこの対応関係を参照して画像を出力する学習画像生成手段と、

出力画像を表示する画像表示手段と

を備えた動画像生成装置。

7. 請求項6記載の動画像生成装置において、上記学習画像生成手段は、言語情報と人間の顔や人体の画像との対応関係をテーブルとして記憶する画像テーブル記憶手段と、上記画像テーブル記憶手段を参照して上記言語情報に対応する画像を出力する画像テーブル参照手段とを備えたことを特徴とする動画像生成装置。

。

8. 請求項6記載の動画像生成装置において、上記学習画像生成手段は、言語情報と人間の表情や動作を示す画像の様態を示す画像制御パラメータとの対応関係をテーブル化して記憶するパラメータテーブル記憶手段と、上記パラメータテーブル記憶手段を参照して上記言語情報に対応する画像制御パラメータを出力するパラメータテーブル参照手段と、画像制御パラメータに対応した顔や人体の画像を記憶する画像記憶手段と、上記パラメータテーブル参照手段から出力される画像制御パラメータを用いて上記画像記憶手段を参照して出力画像を決定する画像制御手段とを備えたことを特徴とする動画像生成装置。

9. 請求項8記載の動画像生成装置において、上記一つの言語情報に対応して複数の動作が同時に起きる動作同期確率を記述したテーブルを記憶する動作同期確率記憶手段をさらに備え、上記パラメータテーブル参照手段は、上記言語解析手段から出力される言語情報に基づき上記パラメータテーブル記憶手段および上記動作同期確率記憶手段を参照して複数の画像制御パラメータを出力することを特徴とする動画像生成装置。

10. 請求項6記載の動画像生成装置において、上記学習画像生成手段は、人間の表情や動作に関する状態と、制御カテゴリーとしての言語情報と、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態のパラメータとの組をネットワークとして記憶するネットワーク記憶手段と、上記ネットワーク記憶手段を参照して上記言語情報に対するネットワークを読み出して、ネットワークが該言語情報を制御カテゴリーとして受理した際にネットワークの状態を遷移し、遷移時に画像制御パラメータを出力するネットワーク参照手段と、画像制御パラメータに対応した顔や

人体の画像を記憶する画像記憶手段と、上記ネットワーク参照手段からの画像制御パラメータを用いて上記画像記憶手段を参照して出力画像を決定する画像制御手段とを備えたことを特徴とする動画像生成装置。

11. 文の言語情報または韻律パラメータを制御カテゴリーとして、その制御パラメータと文を発声している人間の顔や人体の画像パラメータとを対応づけて学習画像データとして記憶する学習画像データ記憶手段と、

人間の表情や動作に関する状態と、制御カテゴリーとしての言語情報または韻律パラメータと、制御カテゴリーが受理され状態遷移する際の遷移先状態と、状態遷移毎に与えられている値としての人間の表情や動作の様態の画像パラメータとの組をネットワークとして記憶するネットワーク記憶手段と、

上記学習画像データ記憶手段から読み出した学習画像データを入力とし、上記ネットワーク記憶手段から学習画像データの制御カテゴリーによって現在保持しているネットワークの状態からの状態遷移を読み出し、状態を遷移先状態に更新して保持し、上記状態遷移に含まれる画像パラメータと学習画像データの画像パラメータとの間の誤差を計算して出力する誤差計算手段と、

状態遷移毎の上記誤差を基に学習対象の状態遷移を一つ決定し、学習対象アークとする学習対象アーク決定手段と、

上記学習対象アークが受理可能とする制御カテゴリーの集合を分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する遷移分割手段と、

上記学習対象アークの遷移元状態または遷移先状態を分割し、それに応じて状態遷移を更新し、分割後の学習対象アークの間の画像パラメータの誤差を計算する状態分割手段と、

上記遷移分割手段から出力される誤差と上記状態分割手段から出力される誤差とを判定し、いずれかの分割手段により更新された状態遷移を選択して、上記ネットワーク記憶手段に記憶されているネットワーク内の状態遷移を書き換える最適分割決定手段と

を備えた画像制御ネットワーク学習装置。

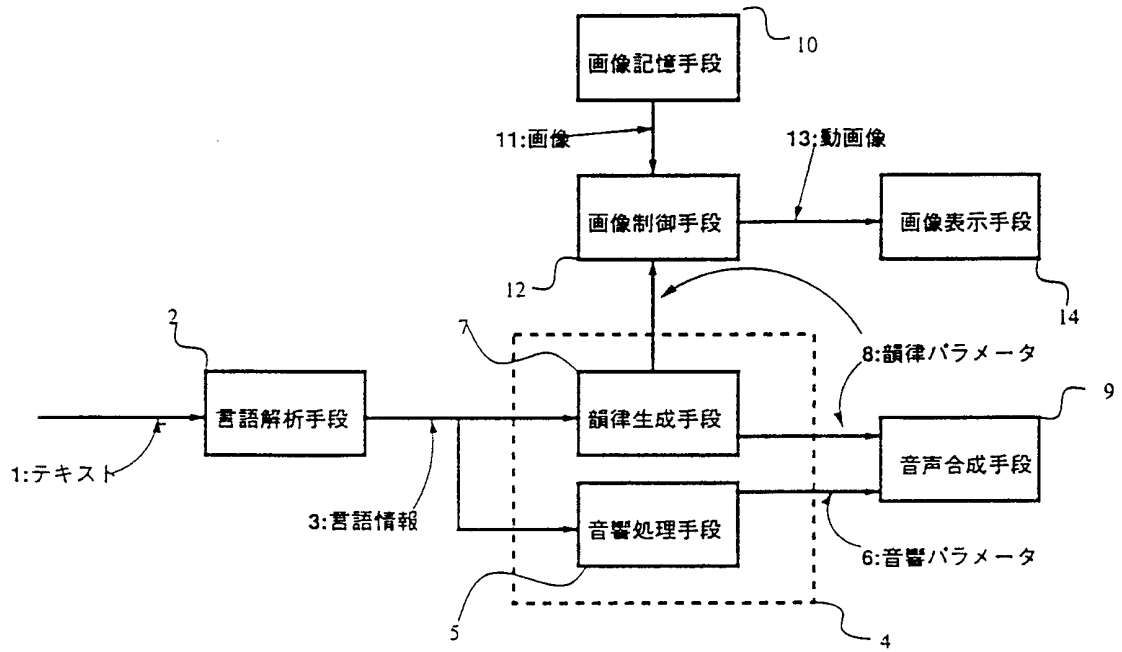


図 1

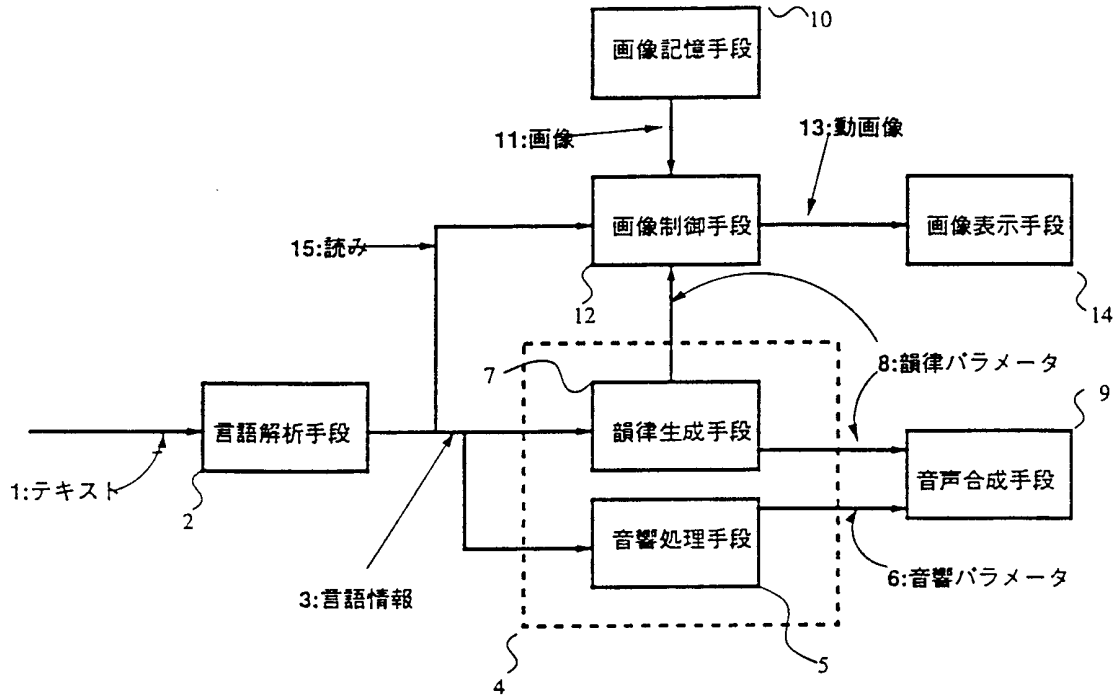


図 2

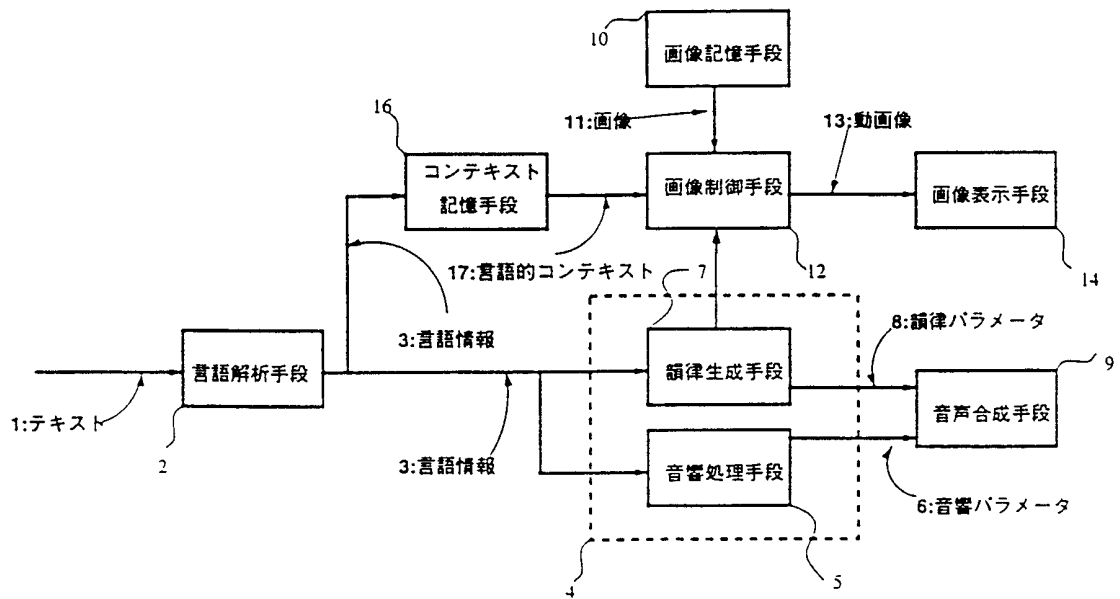


図 3

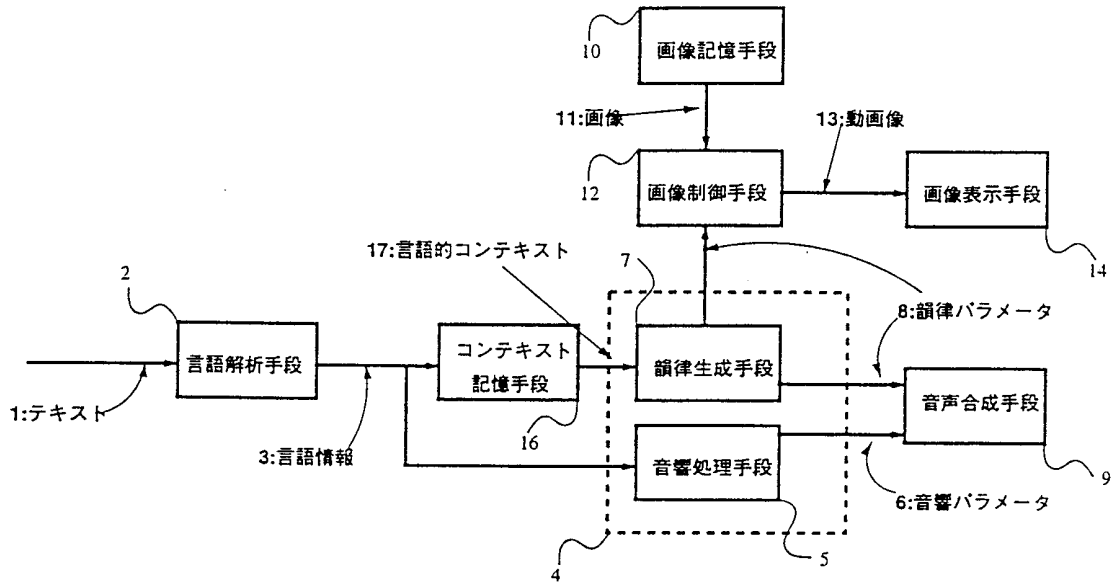


図 4

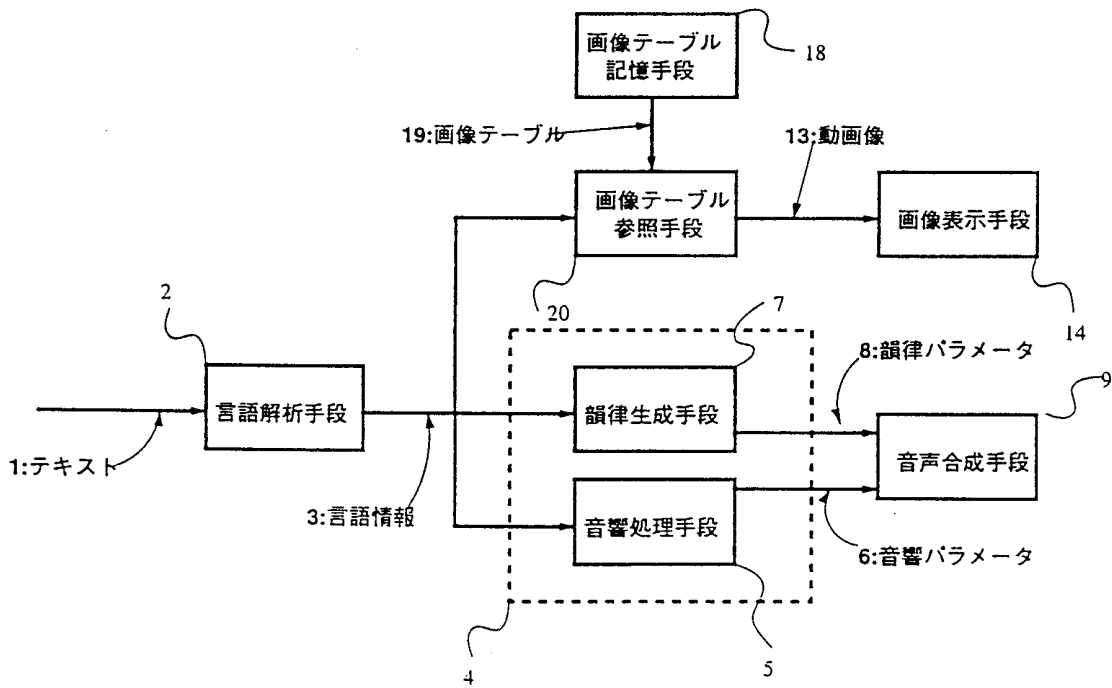


図 5

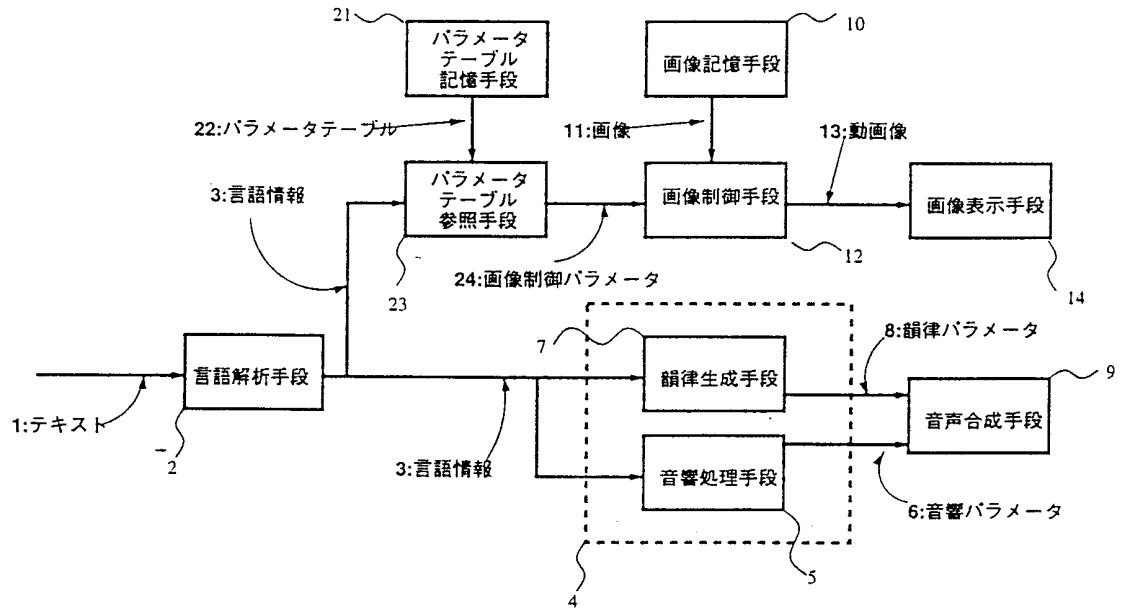


図 6

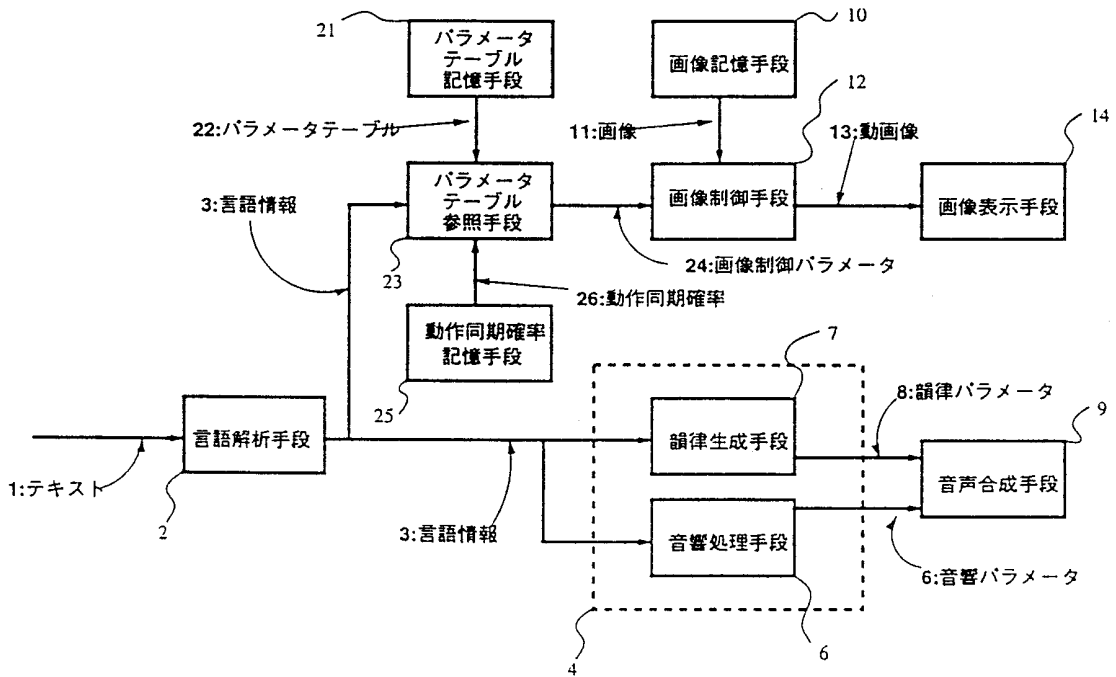


図 7

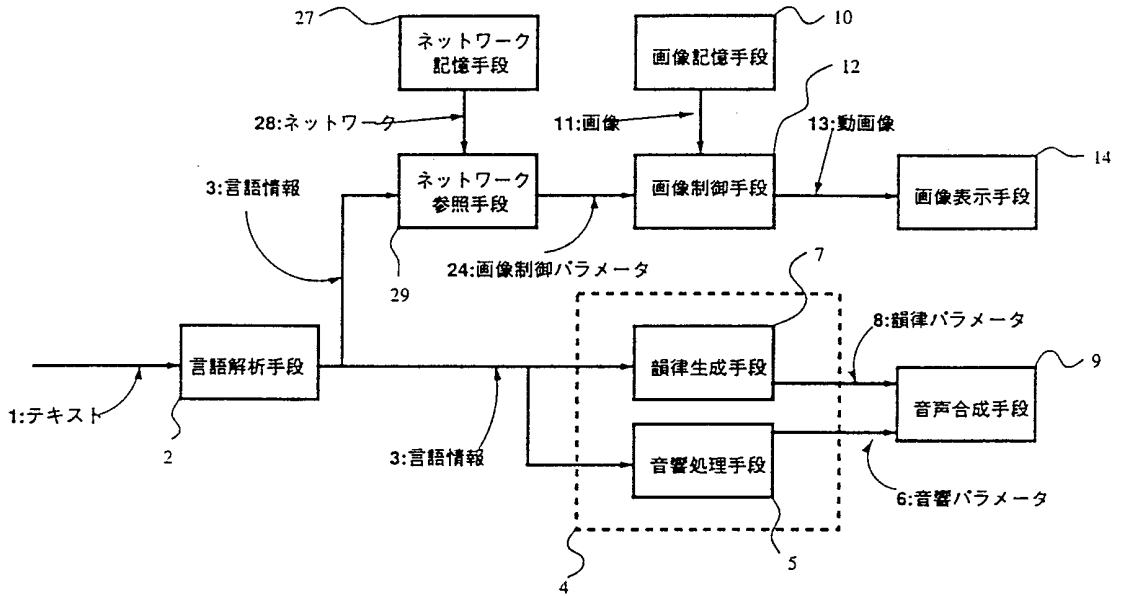


図 8

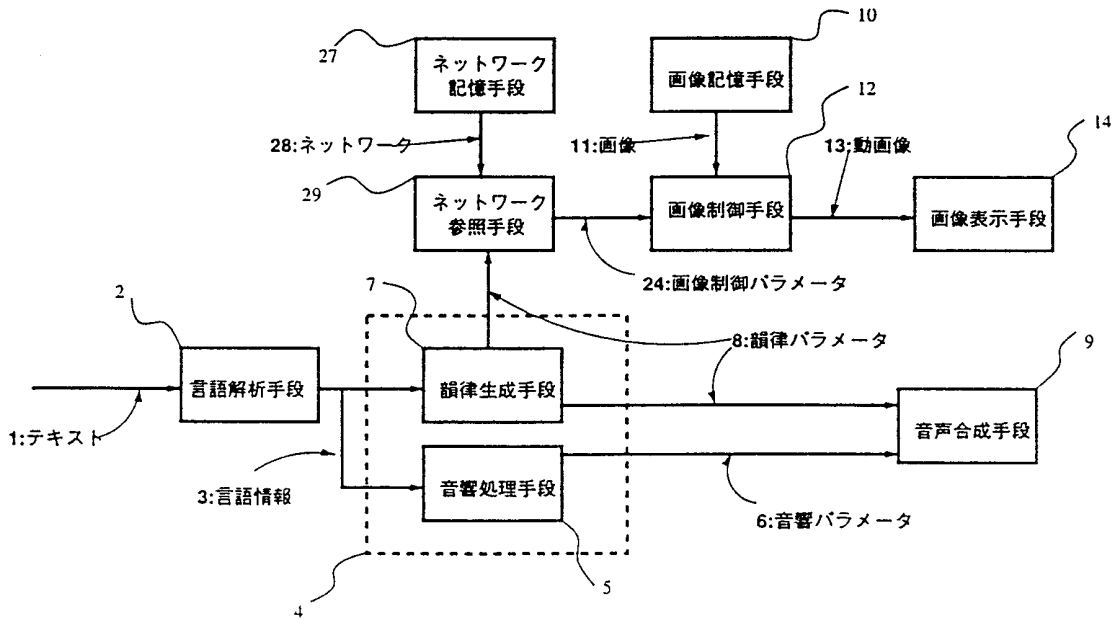


図 9

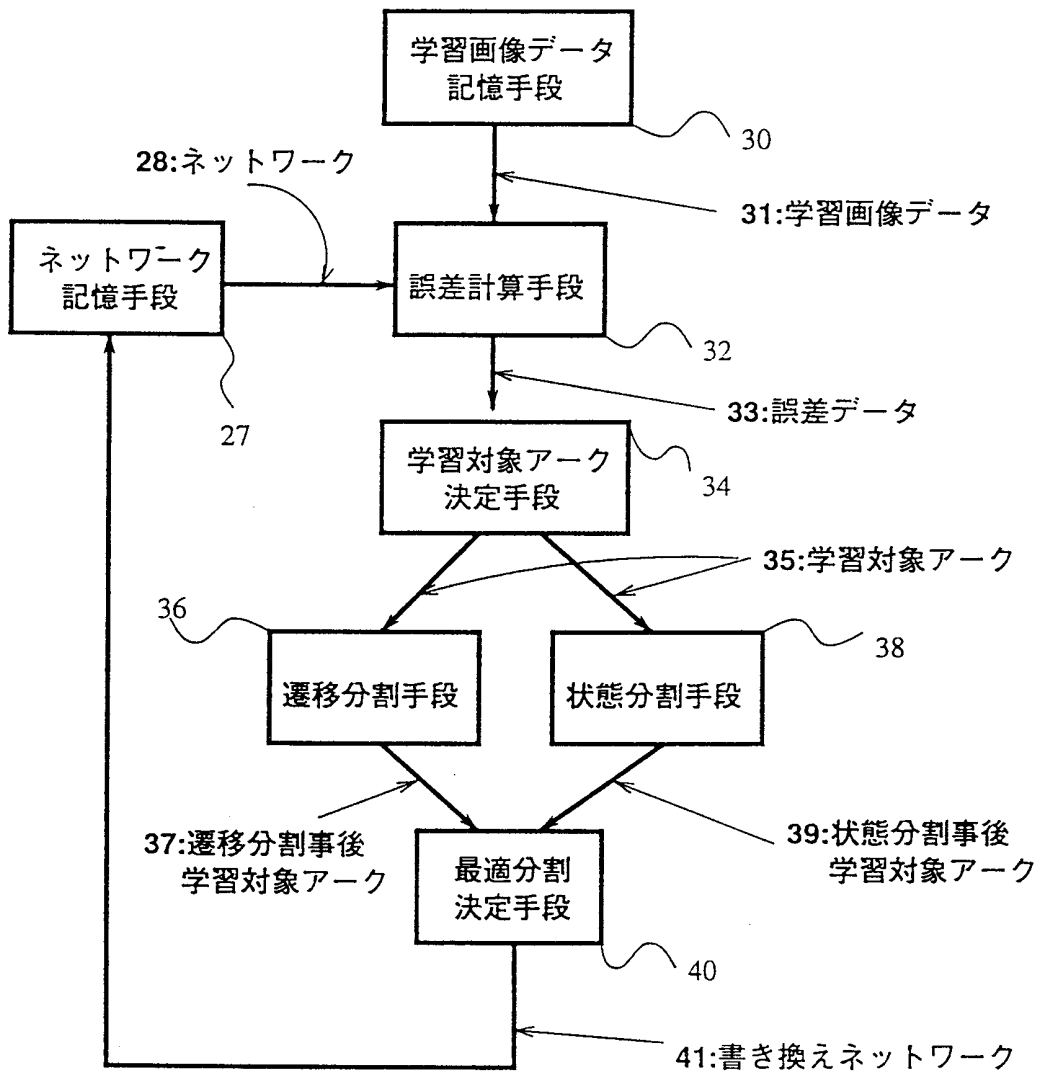


図 10

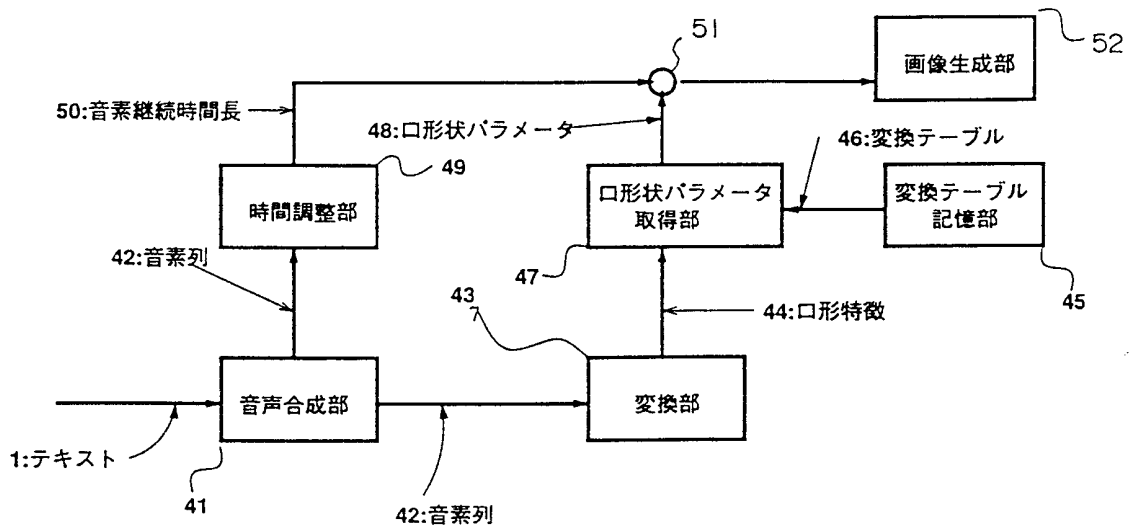
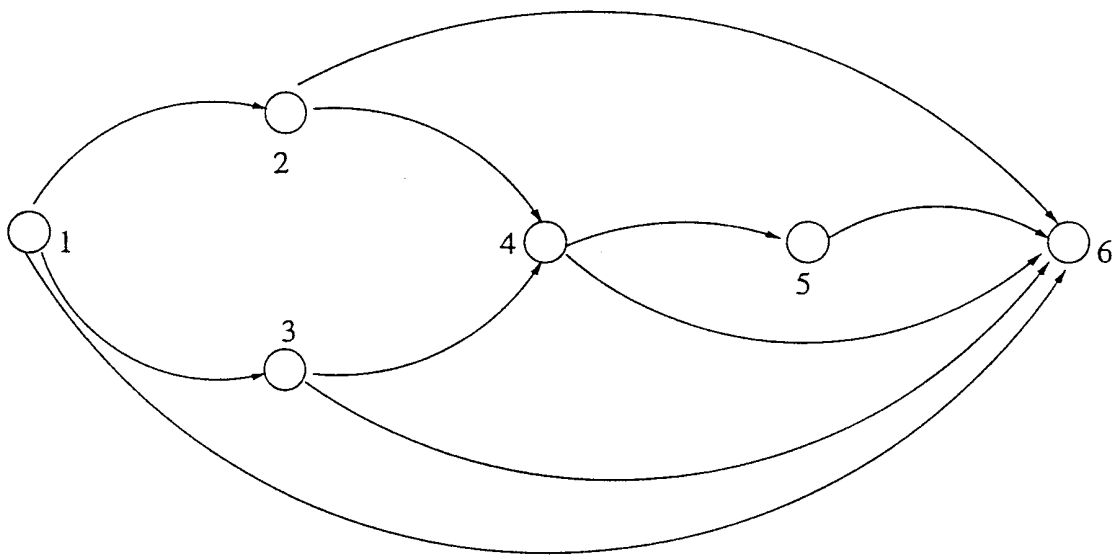


図 11

状態遷移番号	状態番号 p	遷移先状態番号 q	カテゴリ集合 C	出力パラメータ α
1	1	2	c1,c2	α 1
2	1	3	c3	α 2
3	1	6	c4,c5	α 3
4	2	4	c1,c2,c3	α 4
5	2	6	c4,c5	α 5
6	3	4	c1,c2	α 6
7	3	6	c5	α 7
8	4	5	c1,c3	α 8
9	4	6	c4,c5	α 9
10	5	6	c5	α 10

(a)

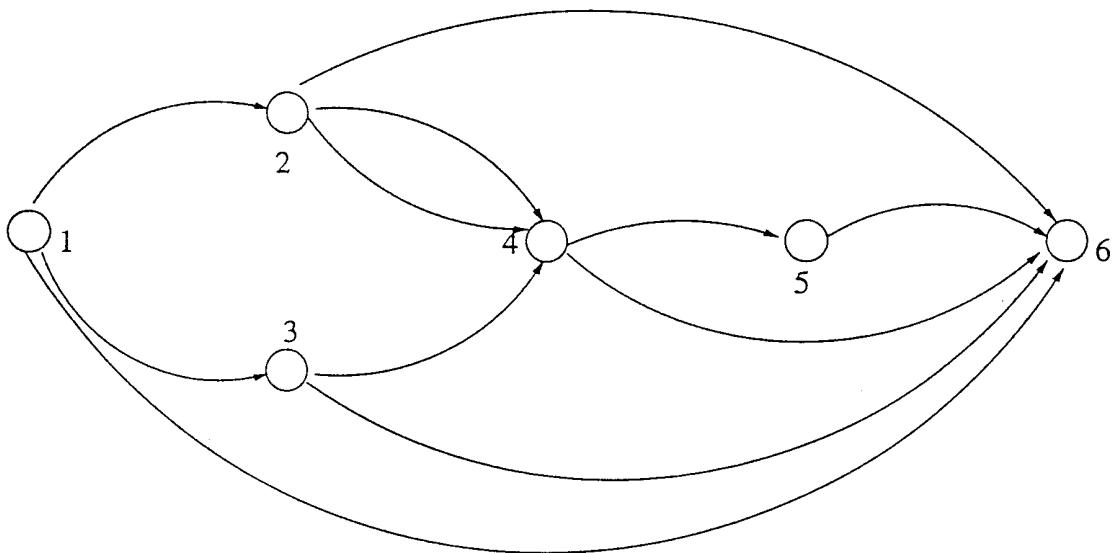


(b)

図 12

状態遷移番号	状態番号 p	遷移先状態番号 q	カテゴリ集合 C	出力パラメータ α
1	1	2	c1,c2	$\alpha 1$
2	1	3	c3	$\alpha 2$
3	1	6	c4,c5	$\alpha 3$
4	2	4	c1,c3	$\alpha 4'$
5	2	6	c4,c5	$\alpha 5$
6	3	4	c1,c2	$\alpha 6$
7	3	6	c5	$\alpha 7$
8	4	5	c1,c3	$\alpha 8$
9	4	6	c4,c5	$\alpha 9$
10	5	6	c5	$\alpha 10$
11	2	4	c2	$\alpha 4''$

(a)



(b)

図 13

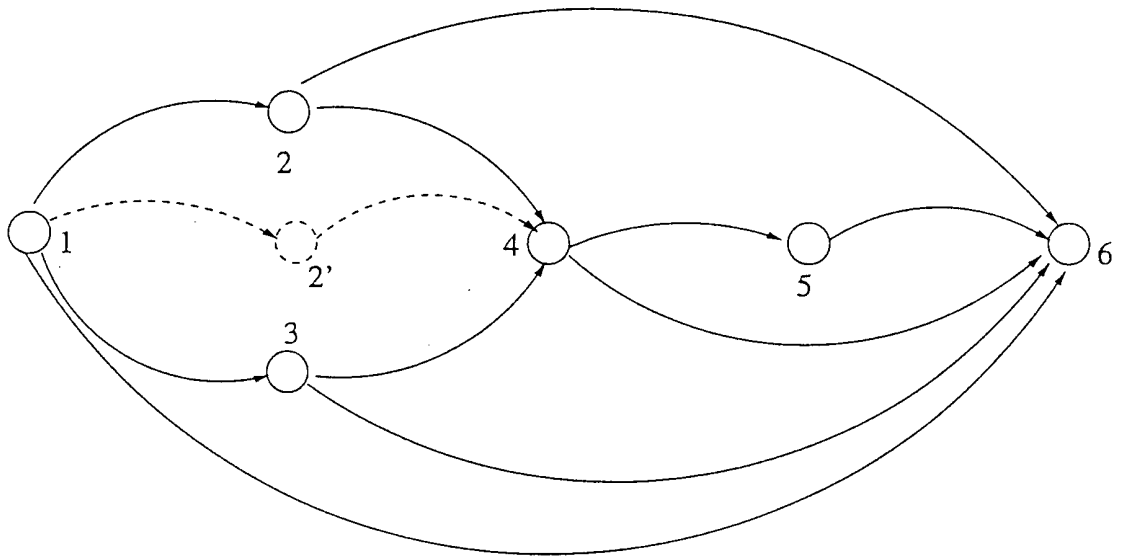
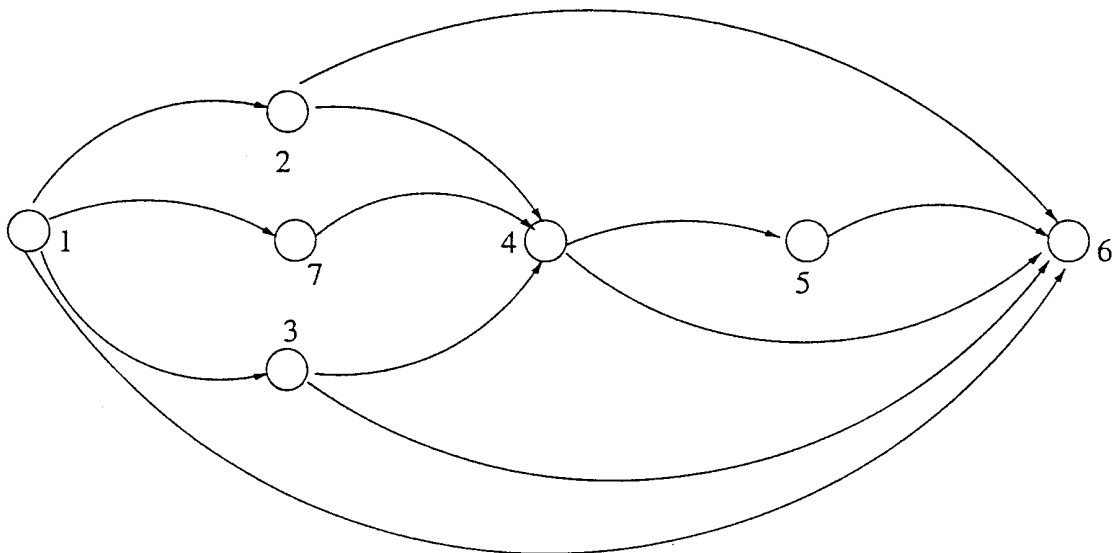


图 14

状態遷移番号	状態番号 p	遷移先状態番号 q	カテゴリ集合 C	出力パラメータ α
1	1	2	c1	$\alpha 1'$
2	1	3	c3	$\alpha 2$
3	1	6	c4,c5	$\alpha 3$
4	2	4	c1,c2,c3	$\alpha 4'$
5	2	6	c4,c5	$\alpha 5$
6	3	4	c1,c2	$\alpha 6$
7	3	6	c5	$\alpha 7$
8	4	5	c1,c3	$\alpha 8$
9	4	6	c4,c5	$\alpha 9$
10	5	6	c5	$\alpha 10$
11	1	7	c2	$\alpha 1''$
12	7	4	c1,c2,c3	$\alpha 4''$

(a)



(b)

図 15

INTERNATIONAL SEARCH REPORT

International application No.
PCT/JP98/01025

A. CLASSIFICATION OF SUBJECT MATTER
Int.Cl⁶ G06T13/00, G10L5/02

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
Int.Cl⁶ G06T13/00, G10L5/02

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched
Jitsuyo Shinan Koho 1926-1998
Kokai Jitsuyo Shinan Koho 1971-1998

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)
JICST File (JOIS)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	JP, 5-153581, A (Seiko Epson Corp.), June 18, 1993 (18. 06. 93)	1-11
Y	JP, 4-359299, A (Sony Corp.), December 11, 1992 (11. 12. 92)	1-11
Y	Noriharu Sakamoto and five others "Multi-Model Experiment System based on Interactive Database (in Japanese)", Gazo-Rabo (Image Lab), 1997, Vol. 8, No. 10, pages 20 to 23	1-11

Further documents are listed in the continuation of Box C. See patent family annex.

<p>* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier document but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed</p>	<p>"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family</p>
--	---

Date of the actual completion of the international search April 28, 1998 (28. 04. 98)	Date of mailing of the international search report May 19, 1998 (19. 05. 98)
--	---

Name and mailing address of the ISA/ Japanese Patent Office	Authorized officer
Facsimile No.	Telephone No.

A. 発明の属する分野の分類 (国際特許分類 (IPC))

Int. Cl^o G06T13/00
Int. Cl^o G10L5/02

B. 調査を行った分野

調査を行った最小限資料 (国際特許分類 (IPC))

Int. Cl^o G06T13/00
Int. Cl^o G10L5/02

最小限資料以外の資料で調査を行った分野に含まれるもの

日本国実用新案公報1926-1998年
日本国公開実用新案公報1971-1998年

国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)

JICSTファイル (JOIS)

C. 関連すると認められる文献

引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
Y	JP, 5-153581, A (セイコーエプソン株式会社) 18. 6月. 1993 (18. 06. 93)	1-11
Y	JP, 4-359299, A (ソニー株式会社) 11. 12月. 1992 (11. 12. 92)	1-11
Y	坂本憲治, 外5名, "対話データベースに基づくマルチモーダル実験システム", 画像ラボ, 1997, Vol. 8, No. 10, 第20~23頁	1-11

C欄の続きにも文献が列挙されている。

パテントファミリーに関する別紙を参照。

* 引用文献のカテゴリー

「A」特に関連のある文献ではなく、一般的技術水準を示すもの
「E」先行文献ではあるが、国際出願日以後に公表されたもの
「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)
「O」口頭による開示、使用、展示等に言及する文献
「P」国際出願日前で、かつ優先権の主張の基礎となる出願

の日の後に公表された文献

「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの
「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの
「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの
「&」同一パテントファミリー文献

国際調査を完了した日

28. 04. 98

国際調査報告の発送日

1905.98

国際調査機関の名称及びあて先

日本国特許庁 (ISA/J P)
郵便番号100-8915
東京都千代田区霞が関三丁目4番3号

特許庁審査官 (権限のある職員)

岩間直純



5H

9287

電話番号 03-3581-1101 内線 3533