

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7621607号
(P7621607)

(45)発行日 令和7年1月27日(2025.1.27)

(24)登録日 令和7年1月17日(2025.1.17)

(51)国際特許分類

F I

H 0 4 N 21/435(2011.01)

H 0 4 N 21/435

H 0 4 N 21/854(2011.01)

H 0 4 N 21/854

請求項の数 17 (全35頁)

(21)出願番号	特願2023-547248(P2023-547248)	(73)特許権者	000002185
(86)(22)出願日	令和4年2月8日(2022.2.8)		ソニーグループ株式会社
(65)公表番号	特表2024-505988(P2024-505988		東京都港区港南1丁目7番1号
	A)	(74)代理人	100092093
(43)公表日	令和6年2月8日(2024.2.8)		弁理士 辻居 幸一
(86)国際出願番号	PCT/US2022/070572	(74)代理人	100109070
(87)国際公開番号	WO2022/170368		弁理士 須田 洋之
(87)国際公開日	令和4年8月11日(2022.8.11)	(74)代理人	100067013
審査請求日	令和5年8月3日(2023.8.3)		弁理士 大塚 文昭
(31)優先権主張番号	17/170,695	(74)代理人	
(32)優先日	令和3年2月8日(2021.2.8)		上杉 浩
(33)優先権主張国・地域又は機関	米国(US)	(74)代理人	100120525
			弁理士 近藤 直樹
		(72)発明者	キャンデロア ブラント
			アメリカ合衆国 カリフォルニア州 9 2
			最終頁に続く

(54)【発明の名称】 シーン説明の再生制御

(57)【特許請求の範囲】

【請求項1】

メディアレンダリング装置であって、
一連の撮影シーンと、前記一連の撮影シーンにおける撮影シーンを説明するテキストベースのビデオ説明情報及びタイミング情報を含むテキスト情報とを含むメディアコンテンツを取得し、
前記撮影シーンの前記テキスト情報から、前記ビデオ説明情報を再生するためのタイミング情報を抽出し、
前記撮影シーンの連続する音声部分間の自然な休止に対応する第1の時間間隔を前記タイミング情報から抽出し、
前記一連の撮影シーンにおける前記撮影シーンのオーディオ部分を再生するための時間間隔をそれぞれが示す、前記撮影シーンの一連の第2の時間間隔を決定し、
前記撮影シーンの前記ビデオ説明情報のオーディオ表現の第3の時間間隔を決定し、
前記決定された第3の時間間隔と前記第1の時間間隔の比率に基づいて乗数を決定し、
前記乗数と前記ビデオ説明情報の実際の再生速度とに基づいて、前記ビデオ説明情報のオーディオ表現を再生する速度を決定し、
前記決定された速度に基づいて、前記ビデオ説明情報の再生を、前記撮影シーンの前記抽出されたタイミング情報によって示される第1の時間間隔においてテキスト表現又はテキスト及びオーディオ表現で制御する、
ように構成された回路を備える、

ことを特徴とするメディアレンダリング装置。

【請求項 2】

前記回路は、

それぞれが前記一連の第 2 の時間間隔とは異なる、前記撮影シーンの一連の第 4 の時間間隔を決定し、

前記一連の第 4 の時間間隔から、時間間隔閾値よりも高い前記第 1 の時間間隔を選択する、

ようにさらに構成される、請求項 1 に記載のメディアレンダリング装置。

【請求項 3】

前記決定される速度は、前記オーディオ表現の実際の再生速度よりも低い、
請求項 1 に記載のメディアレンダリング装置。

10

【請求項 4】

前記決定される速度は、前記オーディオ表現の実際の再生速度よりも高い、
請求項 1 に記載のメディアレンダリング装置。

【請求項 5】

前記回路は、前記メディアレンダリング装置に関連する規定の速度設定に基づいて、前記ビデオ説明情報の前記オーディオ表現を再生する前記速度を決定するようにさらに構成され、

前記規定の速度設定は、前記ビデオ説明情報の前記オーディオ表現の最大再生速度を示す、

20

請求項 1 に記載のメディアレンダリング装置。

【請求項 6】

前記回路は、前記テキスト情報と共に速度情報を受け取り、前記決定された速度及び前記規定の速度設定に基づいて、前記撮影シーンの画像部分又はオーディオ部分の一方の再生を制御するようにさらに構成される、
請求項 5 に記載のメディアレンダリング装置。

【請求項 7】

前記回路は、

前記メディアコンテンツがレンダリングされている対象であるユーザのプロファイル情報を示す第 2 のユーザ入力を受け取り、

30

前記受け取った第 2 のユーザ入力に基づいて、前記ビデオ説明情報の前記オーディオ表現を再生するための速度設定を決定する、
ようにさらに構成される、請求項 5 に記載のメディアレンダリング装置。

【請求項 8】

前記回路は、

前記一連の撮影シーンのうちの 1 つの撮影シーンの説明に対応する第 1 のユーザ入力を受け取り、

前記受け取った第 1 のユーザ入力を、前記一連の撮影シーンの各々に関連する前記ビデオ説明情報内で検索し、

前記検索に基づいて、前記メディアコンテンツを再生するための再生タイミング情報を決定し、

40

前記決定された再生タイミング情報に基づいて前記メディアコンテンツの前記再生を制御する、

ようにさらに構成される、請求項 1 に記載のメディアレンダリング装置。

【請求項 9】

前記第 1 の時間間隔は、前記撮影シーンの第 1 のせりふと第 2 のせりふとの間である、
請求項 1 に記載のメディアレンダリング装置。

【請求項 10】

前記第 1 のせりふは、前記撮影シーンの第 1 のショットの最後の単語であり、前記第 2 のせりふは、前記撮影シーンの第 2 のショットの最初の単語であり、

50

前記第 1 のショット及び前記第 2 のショットは、前記撮影シーンの連続するショットである、
請求項 9 に記載のメディアレンダリング装置。

【請求項 1 1】

訓練済み機械学習 (ML) モデルを記憶するように構成されたメモリをさらに備え、前記テキスト情報は速度情報をさらに含み、前記回路は、

前記撮影シーンの少なくとも 1 つの特性の分析に基づいて前記撮影シーンのコンテキスト情報を決定し、

前記撮影シーンの前記決定されたコンテキスト情報に対する前記訓練済み ML モデルの適用に基づいて、前記ビデオ説明情報のオーディオ表現を再生するためのオーディオ特性を決定し、

10

前記速度情報及び前記決定されたオーディオ特性に基づいて、前記撮影シーンの前記抽出されたタイミング情報によって示される前記第 1 の時間間隔において前記ビデオ説明情報の前記オーディオ表現を再生するように制御する、
ようにさらに構成される、請求項 1 に記載のメディアレンダリング装置。

【請求項 1 2】

前記撮影シーンを説明する前記ビデオ説明情報は、前記撮影シーン内に存在する生物オブジェクト又は無生物オブジェクトに関する認知情報を含み、

前記回路は、前記撮影シーンの前記ビデオ説明情報に含まれる前記認知情報の再生を制御するようにさらに構成される、

20

請求項 1 に記載のメディアレンダリング装置。

【請求項 1 3】

前記メディアレンダリング装置は、前記ビデオ説明情報のテキスト表現を再生するように構成されたディスプレイ装置をさらに備える、

請求項 1 に記載のメディアレンダリング装置。

【請求項 1 4】

前記メディアコンテンツは、前記一連の撮影シーンの各々のオーディオ部分を表すクローズドキャプション情報をさらに含み、

前記一連の撮影シーンの各々を説明する前記ビデオ説明情報は、前記クローズドキャプション情報と共に前記メディアコンテンツ内に符号化される、

30

請求項 1 に記載のメディアレンダリング装置。

【請求項 1 5】

前記回路は、前記メディアレンダリング装置に関連するオーディオレンダリング装置を、前記ビデオ説明情報のオーディオ表現及び前記撮影シーンのオーディオ部分を再生するように制御するようにさらに構成される、

請求項 1 に記載のメディアレンダリング装置。

【請求項 1 6】

メディアレンダリング装置において、

一連の撮影シーンと、前記一連の撮影シーンにおける撮影シーンを説明するテキストベースのビデオ説明情報及びタイミング情報を含むテキスト情報とを含むメディアコンテンツを取得することと、

40

前記撮影シーンの前記テキスト情報から、前記ビデオ説明情報を再生するためのタイミング情報を抽出することと、

前記撮影シーンの連続する音声部分間の自然な休止に対応する第 1 の時間間隔を前記タイミング情報から抽出することと、

前記一連の撮影シーンにおける前記撮影シーンのオーディオ部分を再生するための時間間隔をそれぞれが示す、前記撮影シーンの一連の第 2 の時間間隔を決定することと、

前記撮影シーンの前記ビデオ説明情報のオーディオ表現の第 3 の時間間隔を決定することと、

前記決定された第 3 の時間間隔と前記第 1 の時間間隔の比率に基づいて乗数を決定する

50

ことと、

前記乗数と前記ビデオ説明情報の実際の再生速度とに基づいて、前記ビデオ説明情報のオーディオ表現を再生する速度を決定することと、

前記決定された速度に基づいて、前記ビデオ説明情報の再生を、前記撮影シーンの前記抽出されたタイミング情報によって示される第1の時間間隔においてテキスト表現又はテキスト及びオーディオ表現で制御することと、
を含むことを特徴とする方法。

【請求項17】

コンピュータ実行可能命令を記憶した非一時的コンピュータ可読媒体であって、前記コンピュータ実行可能命令は、メディアレンダリング装置によって実行された時に、前記メディアレンダリング装置に、

一連の撮影シーンと、前記一連の撮影シーンにおける撮影シーンを説明するテキストベースのビデオ説明情報及びタイミング情報を含むテキスト情報とを含むメディアコンテンツを取得することと、

前記撮影シーンの前記テキスト情報から、前記ビデオ説明情報を再生するためのタイミング情報を抽出することと、

前記撮影シーンの連続する音声部分間の自然な休止に対応する第1の時間間隔を前記タイミング情報から抽出することと、

前記一連の撮影シーンにおける前記撮影シーンのオーディオ部分を再生するための時間間隔をそれぞれが示す、前記撮影シーンの一連の第2の時間間隔を決定することと、

前記撮影シーンの前記ビデオ説明情報のオーディオ表現の第3の時間間隔を決定することと、

前記決定された第3の時間間隔と前記第1の時間間隔の比率に基づいて乗数を決定することと、

前記乗数と前記ビデオ説明情報の実際の再生速度とに基づいて、前記ビデオ説明情報のオーディオ表現を再生する速度を決定することと、

前記決定された速度に基づいて、前記ビデオ説明情報の再生を、前記撮影シーンの前記抽出されたタイミング情報によって示される第1の時間間隔においてテキスト表現又はテキスト及びオーディオ表現で制御することと、
を含む動作を実行させる、ことを特徴とする非一時的コンピュータ可読媒体。

【発明の詳細な説明】

【技術分野】

【0001】

〔関連出願との相互参照／引用による組み入れ〕

なし

【0002】

本開示の様々な実施形態は、メディア再生制御に関する。具体的には、本開示の様々な実施形態は、メディアレンダリング装置及びシーン説明の再生制御方法に関する。

【背景技術】

【0003】

近年のメディアコンテンツ再生分野の進歩は、メディアコンテンツの様々な部分を制御するための様々な技術の発展をもたらした。通常、メディアコンテンツ（例えば、映画）は、1又は2以上の視聴者のためにメディアレンダリング装置上で同時にレンダリングできるビデオトラック及び対応するオーディオトラックなどの異なる部分を含むことができる。いくつかの状況では、視覚障害者又は認知障害者などの視聴者が、映画を理解する上でメディアコンテンツのシーンにおける要素、文脈、筋書き又は感情を視覚化できないという問題に直面することがある。メディアコンテンツの中には、視覚障害又は認知障害視聴者のメディアコンテンツ体験をさらに強化するために、ビデオトラック及びオーディオトラックと共に、メディアコンテンツ内にビデオ説明付きオーディオを代替オーディオトラックとして含むことができるものもある。いくつかのシナリオでは、ビデオ説明がオー

10

20

30

40

50

ディオベースであってビデオの説明に使用され、従って「ビデオ説明 (video description)」と呼ばれる。しかしながら、米国では、連邦通信委員会 (FCC) が、2020年11月30日に発行された21世紀における通信とビデオアクセシビリティに関する2010年法、FCC 20-155 (2020年) によってこの用語を「オーディオ説明 (audio description)」に変更した。本文書では、古い用語である「ビデオ説明」を引き続き使用する。このナレーション付きの説明は、視覚障害者又は認知障害者などの視聴者にとってのメディアコンテンツの利用しやすさを強化するものである。これらのビデオ説明は、事前録画されたメディアコンテンツのオーディオトラック (例えば、せりふ) 間の自然な途切れに挿入される。自然な途切れにおけるビデオ説明の挿入に関するいくつかのシナリオでは、対応する自然な途切れの期間内にビデオ説明が収まるように、対応するビデオ説明の1又は2以上の関連部分を削除し、又は自然な途切れの期間を増加させる編集が行われる。このようなシナリオでは、ビデオ説明の関連部分の削除又はメディアコンテンツのオーディオトラックの期間の増大が望ましくない場合もあり、視聴者のコンテンツ体験が不快で低品質なものになってしまう恐れがある。さらに、メディアコンテンツのせりふの自然な途切れにビデオ説明が挿入されるので、認知障害者は、自然な老化過程の一部としてよく理解することができず、従ってビデオ説明を理解できないことが多い。従って、視聴者 (例えば、視覚障害者又は認知障害者) のメディアコンテンツ体験を改善するようにビデオ説明を効果的に制御できる強化された装置が必要とされている。

10

【発明の概要】

20

【発明が解決しようとする課題】

【0004】

当業者には、説明したシステムと、本出願の残り部分において図面を参照しながら示す本開示のいくつかの態様とを比較することにより、従来の慣習的な手法のさらなる限界及び不利点が明らかになるであろう。

【課題を解決するための手段】

【0005】

実質的に少なくとも1つの図に関連して図示及び/又は説明し、特許請求の範囲にさらに完全に示すような、シーン説明の再生制御のためのメディアレンダリング装置及び方法を提供する。

30

【0006】

全体を通じて同じ要素を同じ参照符号によって示す添付図面を参照しながら本開示の以下の詳細な説明を検討することにより、本開示のこれらの及びその他の特徴及び利点を理解することができる。

【図面の簡単な説明】

【0007】

【図1】本開示の実施形態による、シーン説明の再生制御のための例示的なネットワーク環境を示すブロック図である。

【図2】本開示の実施形態による、シーン説明の再生制御のための例示的なメディアレンダリング装置を示すブロック図である。

40

【図3A】本開示の実施形態による、シーン説明の再生制御のための例示的なシナリオを図3Bと合わせて示す図である。

【図3B】本開示の実施形態による、シーン説明の再生制御のための例示的なシナリオを図3Aと合わせて示す図である。

【図4】本開示の実施形態による、シーン説明の再生制御のための別の例示的なシナリオを示す図である。

【図5】本開示の実施形態による、シーン説明の再生制御のための例示的な動作を示す第1のフローチャートである。

【図6】本開示の実施形態による、シーン説明の再生制御のための例示的な動作を示す第2のフローチャートである。

50

【発明を実施するための形態】

【0008】

開示する（視覚障害者又は認知障害者ユーザなどの）視聴者のメディアコンテンツ体験を強化するシーン説明の再生の動的制御のためのメディアレンダリング装置及び方法では、後述する実装を見出すことができる。本開示の例示的な態様は、一連の撮影シーンを含むことができるメディアコンテンツ（例えば、映画）を検索するように構成できるメディアレンダリング装置（例えば、テレビ）を提供する。メディアコンテンツは、ビデオ説明情報（例えば、一連の撮影シーンにおける撮影シーンを説明できるビデオ、筋書き又はシーン説明）と、ビデオ説明情報を再生するためのタイミング情報とを含むことができるテキスト情報をさらに含むことができる。タイミング情報は、ビデオ説明情報のテキスト表現又はオーディオ表現、或いはこれらの組み合わせを収めることができる空白又は途切れ（すなわち、メディアコンテンツのオーディオ部分の空白）に関する情報を含むことができる。メディアレンダリング装置は、ビデオ説明情報を再生するために撮影シーンのテキスト情報からタイミング情報を抽出することができる。メディアレンダリング装置は、ビデオ説明情報の再生を、第1の時間間隔（すなわち、撮影シーンの抽出されたタイミング情報によって示される第1の時間間隔）においてオーディオ表現、テキスト表現、又はテキスト表現及びオーディオ表現で制御するように構成することができる。

10

【0009】

別の実施形態では、テキスト情報が、ビデオ説明情報を再生するための速度情報をさらに含むことができる。速度情報は、タイミング情報に対応するビデオ説明情報のオーディオ表現を再生するための再生速度に関する情報を含むことができる。メディアレンダリング装置は、ビデオ説明情報のオーディオ表現を再生するために撮影シーンのテキスト情報から速度情報を抽出することができる。メディアレンダリング装置は、抽出された速度情報に基づいて、第1の時間間隔（すなわち、抽出された撮影シーンのタイミング情報によって示される第1の時間間隔）においてビデオ説明情報のオーディオ表現の再生を制御するように構成することができる。

20

【0010】

別の実施形態では、メディアレンダリング装置を、一連の撮影シーンのみを含むことができるメディアコンテンツと、一連の撮影シーンにおける撮影シーンを説明することはできるがタイミング情報及び速度情報を含まないビデオ説明情報とを検索するように構成することができる。メディアレンダリング装置は、撮影シーンにおけるオーディオ部分（例えば、せりふ）を再生するための時間間隔をそれぞれが示すことができる、撮影シーンの一連の第2の時間間隔を決定するように構成することができる。メディアレンダリング装置は、撮影シーンのビデオ説明情報（すなわち、シーン説明）のオーディオ表現をレンダリングする期間に対応できる第3の時間間隔を決定するようにさらに構成することができる。メディアレンダリング装置は、ビデオ説明情報のオーディオ表現を含めるために、一連の第2の時間間隔の合間の第1の時間間隔（すなわち、空白又は途切れ）を決定し、決定された一連の第2の時間間隔及び決定された第3の時間間隔に基づいて、含められるビデオ説明情報のオーディオ表現の再生速度をさらに制御するように構成することができる。

30

【0011】

別の実施形態では、メディアレンダリング装置が、メディアレンダリング装置102に関連するディスプレイ装置上にビデオ説明情報を（例えば、テキストフォーマット又は表現で）直接レンダリングすることができる。ビデオ説明情報のテキストは、ディスプレイ装置上にレンダリングできる検索されたメディアコンテンツ上に、又は検索されたメディアコンテンツの外部にオーバーレイ表示することができる。別の実施形態では、ビデオ説明情報のテキストを、任意にクローズドキャプション情報（すなわち、メディアコンテンツのオーディオ部分又はせりふに関連するクローズドキャプション）と共に表示することができる。このことは、ビデオ説明情報が長く、ユーザがディスプレイ装置上にレンダリングされたビデオ説明情報を読むためにさらなる時間を必要とする場合に、メディアレンダリング装置のユーザがメディアコンテンツを手動で制御（一時停止及び再生）すること

40

50

ができるため有利である。

【 0 0 1 2 】

メディアコンテンツの途切れ / 空白にシーン説明を含めるためにビデオ / シーン説明の関連部分を削除し、又は途切れ / 空白の長さの期間を増加させることがある従来の解決策とは対照的に、開示するメディアレンダリング装置は、メディアコンテンツと共に検索できる、又は撮影シーンのせりふ間の検出された空白 / 途切れの期間に基づいて動的に決定できる速度に基づいて、ビデオ説明情報（すなわち、シーン又はビデオ説明）のオーディオ表現を再生することができる。メディアレンダリング装置は、撮影シーンにおいて識別される（単複の）自然な途切れ又は空白の期間と、メディアレンダリング装置に関連する規定の速度設定とに基づいて、ビデオ説明情報のオーディオ表現の再生速度を増加 / 減少させることができる。従って、シーン / ビデオ説明の全体的な再生品質が損なわれず、これによって視聴者（視覚障害者又は認知障害者）のコンテンツ体験をリアルタイムベースでさらに強化することができる。

10

【 0 0 1 3 】

さらに、ビデオ又はシーン説明をオーディオ形態で受け取ることができる従来の解決策と比べて、開示するメディアレンダリング装置は、シーン説明をテキストフォーマットで受け取り、又はテキストフォーマットでメディアコンテンツに含め、さらにシーン説明のテキスト情報をオーディオ表現に変換するように構成することができる。なお、任意に、ビデオ説明情報は、ディスプレイ装置上に直接レンダリングされるようにテキストとして保持する（すなわち、メディアコンテンツ上に、又はメディアコンテンツへの影響が大きい場合にはメディアコンテンツの外部にオーバーレイ表示する）こともできる。従って、シーン説明をテキストフォーマットで含め又は伝えることで、2つの装置間でビデオ説明をオーディオフォーマットで送信するのに必要とされる適切な帯域幅を節約することができる。従って、開示するメディアレンダリング装置は、ビデオ説明情報を含むオーディオトラックとビデオ説明情報を含まないオーディオトラックとを基本的に重複させる従来の解決策と比べて帯域幅を効率的に利用することができる。また、テキスト版のビデオ説明は、映画又はTV番組内の特定のシーンを検索するための単語検索を可能にすることもできる。

20

【 0 0 1 4 】

図1は、本開示の実施形態による、シーン説明の再生制御のための例示的なネットワーク環境を示すブロック図である。図1にはネットワーク環境100を示す。ネットワーク環境100は、メディアレンダリング装置102、ディスプレイ装置104、オーディオレンダリング装置106、サーバ108、及び通信ネットワーク110を含むことができる。メディアレンダリング装置102は、通信ネットワーク110を介してディスプレイ装置104、オーディオレンダリング装置106、サーバ108に通信可能に結合することができる。メディアレンダリング装置102は、アンテナに接続された時にメディアコンテンツ112を受信できるように無線地上波チューナ（図示せず）と共に構成することができる。図1では、メディアレンダリング装置102及びディスプレイ装置104を2つの独立した装置として示しているが、いくつかの実施形態では、本開示の範囲から逸脱することなく、ディスプレイ装置104の機能全体をメディアレンダリング装置102に含めることもできる。

30

40

【 0 0 1 5 】

さらに、図1では、オーディオレンダリング装置106をメディアレンダリング装置102及び / 又はディスプレイ装置104から分離して示しているが、本開示はこのように限定されるものではない。いくつかの実施形態では、本開示の範囲から逸脱することなく、オーディオレンダリング装置106をメディアレンダリング装置102及び / 又はディスプレイ装置104に統合することもできる。図1には、一連の撮影シーン114、オーディオ部分116及びテキスト情報118を含むことができるメディアコンテンツ112をさらに示す。図1に示すように、一連の撮影シーン114は、第1の撮影シーン114A、第2の撮影シーン114B、及び第Nの撮影シーン114Nを含むことができる。一

50

連の撮影シーン 114 の各々は、対応する撮影シーンを形成するように構築できる複数のショットを含むことができる。テキスト情報 118 は、ビデオ説明情報 118 A 及びタイミング情報 118 B を含むこともできる。いくつかの実施形態では、テキスト情報が速度情報 118 C を含むこともできる。ビデオ説明情報 118 A は、一連の撮影シーン 114 における少なくとも 1 つの撮影シーン（例えば、第 1 の撮影シーン 114 A）に関する説明を含むことができる。いくつかの実施形態では、複数のショットを含む一連の撮影シーン 114 の各々が、オーディオ部分 116 に関連する 1 又は 2 以上の画像フレーム又は部分をさらに含むことができる。さらに、メディアレンダリング装置 102 に関連することができるユーザ 120 も示す。例えば、ユーザ 120 は、メディアコンテンツ 112 の視聴者であることができ、視覚障害又は認知障害視聴者とすることができる。

10

【0016】

メディアレンダリング装置 102 は、（サーバ 108 などの）リモートソース又はメディアレンダリング装置 102 のメモリ（すなわち、図 2 のメモリ 204）からメディアコンテンツ 112 を検索するように構成できる好適なロジック、回路、インターフェイス及び/又はコードを含むことができる。いくつかの実施形態では、地上波チューナを利用して無線でメディアコンテンツ 112 を検索することができる。いくつかのシナリオでは、高度テレビシステム委員会（ATSC）又は ATSC 3.0 標準を使用して、メディアコンテンツ 112 をデジタルで受信することができる。

【0017】

メディアコンテンツ 112 は、ビデオ説明情報 118 A、タイミング情報 118 B 及び速度情報 118 C を含むことができるテキスト情報 118 を含むことができる。ビデオ説明情報 118 A は、一連の撮影シーン 114 のうちの（第 1 の撮影シーン 114 A などの）撮影シーンを説明することができる。メディアレンダリング装置 102 は、第 1 の撮影シーン 114 A のテキスト情報 118 からタイミング情報 118 B を抽出するように構成することができる。タイミング情報 118 B は、メディアレンダリング装置 102 がビデオ説明情報を再生するために使用することができる。いくつかの実施形態では、メディアレンダリング装置が、ビデオ説明情報を再生するためにタイミング情報 118 B と共に速度情報 118 C を使用することもできる。メディアレンダリング装置 102 は、ビデオ説明情報の再生を、撮影シーンの抽出されたタイミング情報によって示される第 1 の時間間隔においてテキスト表現、オーディオ表現、又はテキスト表現及びオーディオ表現の両方で制御するようにさらに構成することができる。他のいくつかの実施形態では、メディアレンダリング装置 102 を、抽出された速度情報 118 C に基づいて、撮影シーンの抽出されたタイミング情報によって示される第 1 の時間間隔においてビデオ説明情報のオーディオ表現の再生を制御するようにさらに構成することができる。

20

30

【0018】

別の実施形態では、メディアレンダリング装置 102 が、（サーバ 108 などの）リモートソース又はメディアレンダリング装置 102 のメモリ（すなわち、図 2 のメモリ 204）からメディアコンテンツ 112 を検索することができる。メディアコンテンツは、一連の撮影シーン 114 の（第 1 の撮影シーン 114 A などの）撮影シーンを説明できるビデオ説明情報 118 A を含むことができる。メディアレンダリング装置 102 は、第 1 の撮影シーン 114 A のオーディオ部分 116 を再生するための一連の第 2 の時間間隔を決定し、ビデオ説明情報 118 A のオーディオ表現を再生するための第 3 の時間間隔を決定するように構成することができる。メディアレンダリング装置 102 は、決定された一連の第 2 の時間間隔及び第 3 の時間間隔に基づいて、ビデオ説明情報 118 A のオーディオ表現の再生速度を決定するようにさらに構成することができる。メディアレンダリング装置 102 の例としては、以下に限定するわけではないが、デジタルメディアプレーヤ（DMP）、スマートテレビメディアプレーヤ、オーバーザトップ（OTT）プレーヤ、デジタルメディアストリーマ、メディアエクステンダ/レギュレータ、デジタルメディアハブ、メディアコンテンツコントローラ、テレビ、コンピュータワークステーション、メインフレームコンピュータ、ハンドヘルドコンピュータ、携帯電話機、スマートフォン、セル

40

50

ラー電話機、スマート家電、携帯情報端末（PDA）、スマートスピーカ、スマートメガネ、サウンドシステム、ヘッドマウント装置（HMD）、ヘッドセット、スマートヘッドホン、及び／又はオーディオ・ビデオレンダリング能力を有するその他のコンピュータ装置を挙げることができる。

【0019】

ディスプレイ装置104は、検索されたメディアコンテンツ112内に存在する一連の撮影シーン114を表示するように構成できる好適なロジック、回路及びインターフェイスを含むことができる。ディスプレイ装置104は、ビデオ説明情報118Aをテキストフォーマットで表示するようにさらに構成することができる。ディスプレイ装置104は、ユーザがディスプレイ装置104を介してユーザ入力を提供することを可能にするタッチ画面とすることができる。タッチ画面は、抵抗膜式タッチ画面、静電容量式タッチ画面、又は感熱式タッチ画面のうちの少なくとも1つとすることができる。ディスプレイ装置104は、以下に限定するわけではないが、液晶ディスプレイ（LCD）ディスプレイ、発光ダイオード（LED）ディスプレイ、プラズマディスプレイ、又は有機LED（OLED）ディスプレイ技術のうちの少なくとも1つ、或いはその他のディスプレイ装置などの複数の既知の技術を通じて実現することができる。ある実施形態によれば、ディスプレイ装置104は、ヘッドマウント装置（HMD）のディスプレイ画面、スマートメガネ装置、シースルーディスプレイ、投影型ディスプレイ、エレクトロクロミックディスプレイ、又は透明ディスプレイを意味することができる。

【0020】

オーディオレンダリング装置106は、ビデオ説明情報118A（すなわち、シーン又はビデオ説明）のオーディオ表現を再生又はプレイバックするように構成できる好適なロジック、回路及びインターフェイスを含むことができる。オーディオレンダリング装置106は、第1の撮影シーン114A又は一連の撮影シーン114のオーディオ部分116（例えば、せりふ）を再生するようにさらに構成することができる。オーディオレンダリング装置106の例としては、以下に限定するわけではないが、ラウドスピーカ、壁埋め込み型／天井取り付け型スピーカ、サウンドバー、ウーファ又はサブウーファ、サウンドカード、ヘッドフォン、ヘッドセット、ワイヤレススピーカ、及び／又はオーディオ再生能力を有するその他のコンピュータ装置を挙げることができる。

【0021】

サーバ108は、メディアコンテンツ112を記憶するように構成できる好適なロジック、回路、インターフェイス及びコードを含むことができる。サーバ108は、メディアレンダリング装置102から、サーバ108に記憶されているメディアコンテンツ112を検索するための要求を受け取ることができる。いくつかの実施形態では、サーバ108を、ビデオ説明情報118A（すなわち、シーン説明）のオーディオ表現の最大再生速度を示すことができる規定の速度設定を記憶するように構成することができる。サーバ108は、ウェブアプリケーション、クラウドアプリケーション、HTTPリクエスト、リボジトリ操作及びファイル転送などを通じて動作を実行できるクラウドサーバとして実装することができる。サーバ108の他の例としては、以下に限定するわけではないが、データベースサーバ、ファイルサーバ、ウェブサーバ、メディアサーバ、アプリケーションサーバ、メインフレームサーバ、クラウドサーバ、又はその他のタイプのサーバを挙げることができる。1又は2以上の実施形態では、サーバ108を、当業者に周知の複数の技術を使用することによって複数の分散型クラウドベースリソースとして実装することができる。当業者であれば、本開示の範囲は、サーバ108及びメディアレンダリング装置102を独立エンティティとして実装することに限定されるものではないと理解するであろう。いくつかの実施形態では、本開示の範囲から逸脱することなく、サーバ108の機能を全体的に又は少なくとも部分的にメディアレンダリング装置102に組み込むこともできる。

【0022】

通信ネットワーク110は、メディアレンダリング装置102、ディスプレイ装置10

10

20

30

40

50

4、オーディオレンダリング装置106及びサーバ108が互いに通信できるようにする通信媒体を含むことができる。通信ネットワーク110は、有線通信ネットワーク又は無線通信ネットワークとすることができる。通信ネットワーク110の例としては、以下に限定するわけではないが、インターネット、クラウドネットワーク、ワイヤレスフィデリティ(Wi-Fi)ネットワーク、パーソナルエリアネットワーク(PAN)、ローカルエリアネットワーク(LAN)、又はメトロポリタンエリアネットワーク(MAN)を挙げることができる。ネットワーク環境100内の様々な装置は、様々な有線及び無線通信プロトコルに従って通信ネットワーク110に接続するように構成することができる。このような有線及び無線通信プロトコルの例としては、以下に限定するわけではないが、伝送制御プロトコル及びインターネットプロトコル(TCP/IP)、ユーザデータグラム
10
プロトコル(UDP)、ハイパーテキスト転送プロトコル(HTTTP)、ファイル転送プロトコル(FTP)、ZigBee、EDGE、IEEE802.11、ライトフィデリティ(Li-Fi)、802.16、IEEE802.11s、IEEE802.11g、マルチホップ通信、無線アクセスポイント(AP)、装置間通信、セルラー通信プロトコル、Bluetooth(BT)通信プロトコルを挙げることができる。

【0023】

動作時には、開示するメディアレンダリング装置102が、ユーザ120からメディアコンテンツ112を再生するための要求を受け取ることができる。メディアコンテンツ112の例としては、以下に限定するわけではないが、ビデオクリップ、映画、広告、オーディオ-ビデオコンテンツ、ゲームコンテンツ、又はスライドショークリップを挙げること
20
ができる。メディアレンダリング装置102は、この要求に基づいて、(サーバ108などの)リモートソース又はメディアレンダリング装置102の(図2のメモリ204などの)メモリからメディアコンテンツ112を検索することができる。メディアコンテンツ112は、一連の撮影シーン114、オーディオ部分116、及びテキスト情報118を含むことができる。テキスト情報118は、一連の撮影シーン114のうちの撮影シーン(例えば、第1の撮影シーン114A)を説明することができる、テキストフォーマットであることができるビデオ説明情報118Aを含むことができる。いくつかの実施形態では、ビデオ説明情報118Aが、メディアコンテンツ112内に存在する一連の撮影シーン114の各々を説明することができる。ある実施形態では、メディアレンダリング装置102を、第1の撮影シーン114Aの(例えば、テキストフォーマットの)ビデオ説
30
明情報118Aをビデオ説明情報118Aのオーディオ表現に変換するようにさらに構成することができる。テキスト情報118は、タイミング情報118Bを含むこともできる。タイミング情報118Bは、ビデオ説明情報118Aのオーディオ表現を収めて再生できる第1の時間間隔を示すことができる。別の実施形態では、テキスト情報118が速度情報118Cをさらに含むことができる。速度情報118Cは、タイミング情報118Bによって示される(第1の時間間隔などの)特定の時間間隔中にビデオ説明情報118Aのオーディオ表現を再生する再生速度を示すことができる。メディアレンダリング装置102は、第1の撮影シーン114Aのテキスト情報118からタイミング情報118Bを抽出するようにさらに構成することができる。メディアレンダリング装置102は、ビデオ説明情報118Aの再生を、一連の撮影シーン114の第1の撮影シーン114Aの抽出されたタイミング情報118Bによって示される第1の時間間隔においてテキスト表現、オーディオ表現、又はテキスト表現及びオーディオ表現で制御するようにさらに構成
40
することができる。他のいくつかの実施形態では、メディアレンダリング装置102を、速度情報118Cを抽出するようにさらに構成することができる。このような事例では、メディアレンダリング装置102を、抽出された速度情報118Cに基づいて、一連の撮影シーン114のうちの第1の撮影シーン114Aの抽出されたタイミング情報118Bによって示される第1の時間間隔においてビデオ説明情報118Aのオーディオ表現の再生を制御するようにさらに構成することができる。

【0024】

別の実施形態では、メディアレンダリング装置102が、一連の撮影シーン114にお

10

20

30

40

50

ける第1の撮影シーン114Aのオーディオ部分116（すなわち、せりふ）を再生するための時間間隔をそれぞれが示すことができる、第1の撮影シーン114Aの一連の第2の時間間隔を決定することができる。メディアレンダリング装置102は、第1の撮影シーン114Aのビデオ説明情報118Aのオーディオ表現を再生するために必要な第3の時間間隔を決定するようにさらに構成することができる。第3の時間間隔は、第1の撮影シーン114Aのビデオ説明情報118Aのオーディオ表現を再生するのにかかる時間又はそのために必要な期間に対応することができる。一連の第2の時間間隔及び第3の時間間隔の詳細については、例えば図4で説明する。

【0025】

メディアレンダリング装置102は、ビデオ説明情報118Aのオーディオ表現を再生する速度を決定するようにさらに構成することができる。決定される速度は、例えば第1の撮影シーン114Aの再生中にユーザ120のためにビデオ説明情報118Aのオーディオ表現を再生できる速度とすることができる。ビデオ説明情報118Aのオーディオ表現の再生速度は、決定された一連の第2の時間間隔及び決定された第3の時間間隔に基づいて決定することができる。いくつかの実施形態では、決定される速度が、ビデオ説明情報118Aのオーディオ表現の実際の再生速度よりも低いことができる。他のいくつかの実施形態では、決定される速度が、ビデオ説明情報118Aのオーディオ表現の実際の再生速度よりも高いことができる。決定された一連の第2の時間間隔及び決定された第3の時間間隔に基づくビデオ説明情報118Aのオーディオ表現の再生速度の決定の詳細については、例えば図4で説明する。

【0026】

メディアレンダリング装置102は、決定された速度に基づいて、ビデオ説明情報118Aのオーディオ表現の再生を制御するようにさらに構成することができる。ビデオ説明情報118Aのオーディオ表現は、第1の時間間隔（例えば、第1の撮影シーン114Aのせりふ間の空白）において再生することができる。第1の時間間隔は、一連の第2の時間間隔とは異なることができる。いくつかの実施形態では、第1の時間間隔を、第1の撮影シーン114Aの第1のせりふと第2のせりふとの間の間隔（すなわち、空白）とすることができる。第1のせりふは、第1の撮影シーン114Aのあるショット（例えば、第1のショット）の最後の単語に対応することができ、第2のせりふは、第1の撮影シーン114Aの次のショット（例えば、第2のショット）の最初の単語に対応することができる。第1のショット及び第2のショットは、第1の撮影シーン114Aの連続するショットとすることができる。別の実施形態では、第1の時間間隔を、第1の撮影シーン114Aの開始と第1の撮影シーン114Aの第1のせりふとの間の間隔（すなわち、空白）とすることができる。ある実施形態では、第1の時間間隔（すなわち、空白）が第3の時間間隔よりも短い場合、メディアレンダリング装置102が、ビデオ説明情報118Aのオーディオ表現の再生速度をビデオ説明情報118Aのオーディオ表現の実際の再生速度よりも高くなるように決定することができる。別の実施形態では、第1の時間間隔（すなわち、空白）が第3の時間間隔よりも長い場合、メディアレンダリング装置102が、ビデオ説明情報118Aのオーディオ表現の再生速度をオーディオ表現の実際の再生速度よりも低くなるように決定することができる。ビデオ説明情報118Aのオーディオ表現（すなわち、第1の撮影シーン114Aの説明）の再生速度を増減することで、（説明内容の短縮のような）著しい修正を伴わずにシーン説明を再生することができ、視覚障害又は認知障害視聴者などのユーザ120のためにシーン/ビデオ説明の品質をさらに維持することができる。

【0027】

図1には、本開示の範囲から逸脱することなく修正、追加又は省略を行うことができる。例えば、ネットワーク環境100は、本開示において図示し説明する要素よりも多くの又は少ない要素を含むことができる。例えば、いくつかの実施形態では、ネットワーク環境100が、メディアレンダリング装置102を含んでディスプレイ装置104を含まないこともできる。また、いくつかの実施形態では、本開示の範囲から逸脱することなく、

10

20

30

40

50

各オーディオレンダリング装置 106 の機能をディスプレイ装置 104 に組み込むこともできる。

【0028】

図 2 は、本開示の実施形態による、シーン説明の再生制御のための例示的なメディアレンダリング装置を示すブロック図である。図 2 の説明は、図 1 の要素に関連して行う。図 2 には、メディアレンダリング装置 102 のブロック図 200 を示す。メディアレンダリング装置 102 は、シーン又はビデオ説明の再生を制御する動作を実行できる回路 202 を含むことができる。メディアレンダリング装置 102 は、メモリ 204、入力/出力 (I/O) 装置 206、テキスト - スピーチ変換器 208、ネットワークインターフェイス 210、ディスプレイ装置 104、及びオーディオレンダリング装置 106 をさらに含むことができる。メモリ 204 は、機械学習 (ML) モデル 212 を含むことができる。回路 202 は、メモリ 204、I/O 装置 206、テキスト - スピーチ変換器 208、ネットワークインターフェイス 210、ディスプレイ装置 104 及びオーディオレンダリング装置 106 に通信可能に結合することができる。

【0029】

回路 202 は、メディアレンダリング装置 102 によって実行される異なる動作に関連するプログラム命令を実行するように構成できる好適なロジック、回路及びインターフェイスを含むことができる。例えば、これらの動作の一部は、メディアコンテンツ 112 の検索、タイミング情報 118B 及び/又は速度情報 118C の抽出、及び抽出されたタイミング情報 118B 及び速度情報 118C に基づく第 1 の時間間隔におけるビデオ説明情報 118A のテキスト表現又はオーディオ表現又はテキスト表現及びオーディオ表現の再生を含むことができる。回路 202 は、独立したプロセッサとして実装できる 1 又は 2 以上の特殊処理ユニットを含むことができる。ある実施形態では、1 又は 2 以上の特殊処理ユニットを、1 又は 2 以上の特殊処理ユニットの機能をまとめて実行する統合プロセッサ又はプロセッサ群として実装することができる。回路 202 は、当業で周知の複数のプロセッサ技術に基づいて実装することができる。回路 202 の実装例は、X86 ベースのプロセッサ、グラフィックプロセッシングユニット (GPU)、縮小命令セットコンピューティング (RISC) プロセッサ、特定用途向け集積回路 (ASIC) プロセッサ、複合命令セットコンピューティング (CISC) プロセッサ、マイクロコントローラ、中央処理装置 (CPU)、及び/又はその他の制御回路とすることができる。

【0030】

メモリ 204 は、回路 202 によって実行される命令を記憶するように構成できる好適なロジック、回路、インターフェイス及び/又はコードを含むことができる。メモリ 204 は、メディアコンテンツ 112、テキスト情報 118、及びビデオ説明情報 118A のオーディオ表現の最大再生速度を示すことができる規定の速度設定を記憶するように構成することができる。メモリ 204 は、(第 1 の撮影シーン 114A などの) 撮影シーンのコンテキスト情報に基づいてオーディオ特性を決定するように構成できる訓練済み機械学習 (ML) モデル 212 を記憶するようにさらに構成することができる。ML モデル 212 の詳細な機能については、例えば図 4 で説明する。メモリ 204 は、ユーザのプロファイル情報を記憶するようにさらに構成することができる。メモリ 204 の実装例としては、以下に限定するわけではないが、ランダムアクセスメモリ (RAM)、リードオンリメモリ (ROM)、電氣的に消去可能なプログラマブルリードオンリメモリ (EEPROM)、ハードディスクドライブ (HDD)、固体ドライブ (SSD)、CPU キャッシュ、及び/又はセキュアデジタル (SD) カードなどを挙げることができる。

【0031】

I/O 装置 206 は、入力を受け取り、受け取った入力に基づいて出力を提供するように構成できる好適なロジック、回路及びインターフェイスを含むことができる。I/O 装置 206 は、撮影シーンの説明に対応する第 1 のユーザ入力を受け取るように構成することができる。I/O 装置は、ディスプレイ装置 104 及びオーディオレンダリング装置 106 を介してメディアコンテンツ 112 の再生を制御するようにさらに構成することがで

きる。I/O装置は、回路202と通信するように構成できる様々な入力及び出力装置を含むことができる。I/O装置206の例としては、以下に限定するわけではないが、ディスプレイ装置104、オーディオレンダリング装置106、タッチ画面、キーボード、マウス、ジョイスティック、及びマイクを挙げることができる。

【0032】

テキスト-スピーチ変換器208は、一連の撮影シーン114のうちの少なくとも第1の撮影シーン114Aを説明するビデオ説明情報118Aをオーディオレンダリング装置106による再生に適合できるオーディオフォーマットに変換するように構成できる好適なロジック、回路、インターフェイス及び/又はコードを含むことができる。本明細書では、変換されたオーディオをビデオ説明情報118Aのオーディオ表現と呼ぶことができ、オーディオレンダリング装置106上でレンダリングすることができる。テキスト-スピーチ変換器208は、当業で周知の数多くのプロセッサ技術に基づいて実装することができる。プロセッサ技術の例としては、以下に限定するわけではないが、中央処理装置(CPU)、x86ベースのプロセッサ、縮小命令セットコンピューティング(RISC)プロセッサ、特定用途向け集積回路(ASIC)プロセッサ、複合命令セットコンピューティング(CISC)プロセッサ、及びその他のプロセッサを挙げることができる。

【0033】

ネットワークインターフェイス210は、通信ネットワーク110を介して回路202とサーバ108との間の通信を容易にするように構成できる好適なロジック、回路及びインターフェイスを含むことができる。ネットワークインターフェイス210は、メディアレンダリング装置102と通信ネットワーク110との有線又は無線通信をサポートする様々な既知の技術を使用して実装することができる。ネットワークインターフェイス210は、以下に限定するわけではないが、アンテナ、無線周波数(RF)トランシーバ、1又は2以上の増幅器、チューナ、1又は2以上の発振器、デジタルシグナルプロセッサ、コーダ-デコーダ(CODEC)チップセット、加入者IDモジュール(SIM)カード、又はローカルバッファ回路を含むことができる。ネットワークインターフェイス210は、インターネット、イントラネットなどのネットワーク、又はセルラー電話ネットワーク、無線ローカルエリアネットワーク(LAN)及びメトロポリタンエリアネットワーク(MAN)などの無線ネットワークと無線通信を介して通信するように構成することができる。無線通信は、グローバルシステムフォーモバイルコミュニケーションズ(GSM)、拡張データGSM環境(EDGE)、広帯域符号分割多重アクセス(W-CDMA)、ロングタームエボリューション(LTE)、符号分割多重アクセス(CDMA)、時分割多重アクセス(TDMA)、Bluetooth、(IEEE802.11a、IEEE802.11b、IEEE802.11g又はIEEE802.11nなどの)ワイヤレスフィデリティ(WiFi)、ボイスオーバーインターネットプロトコル(VoIP)、ライトフィデリティ(Li-Fi)、ワールドワイド・インターオペラビリティ・フォー・マイクロウェーブ・アクセス(Wi-MAX)、電子メール用プロトコル、インスタントメッセージ、及びショートメッセージサービス(SMS)などの複数の通信標準、プロトコル及び技術のうちの1つ又は2つ以上を使用するように構成することができる。

【0034】

図3Aは、本開示の実施形態による、シーン説明の再生制御のための例示的なシナリオを示す図である。図3Aの説明は、図1及び図2の要素に関連して行う。図3Aには例示的なシナリオ300を示す。例示的なシナリオ300には、メディアレンダリング装置302(すなわち、メディアレンダリング装置102と同様のもの)を示す。図3Aには、メディアレンダリング装置302に関連するディスプレイ装置104及びオーディオレンダリング装置106をさらに示す。メディアレンダリング装置302は、ディスプレイ装置104及びオーディオレンダリング装置106を、メディアコンテンツをレンダリングするように制御することができる。メディアコンテンツの例としては、以下に限定するわけではないが、ビデオクリップ、映画、オーディオビデオコンテンツ、ゲームコンテンツ、広告、又はスライドショークリップを挙げることができる。メディアコンテンツは、(

図 3 A に示すような) ディスプレイ装置 1 0 4 上に表示された撮影シーン 3 0 4 を含むことができる(図 1 の一連の撮影シーン 1 1 4 などの)一連の撮影シーンを含むことができる。

【 0 0 3 5 】

なお、図 3 A に示す撮影シーン 3 0 4 は車のシーンの一例として提示するものにすぎない。本開示は、(以下に限定するわけではないが、アクションシーン、ドラマシーン、ロマンチックなシーン、感情的なシーン、ダンスシーン、音楽シーン、ホラーシーン、又はラブシーンなどの)他のタイプの撮影シーンにも適用可能である。他のタイプの撮影シーンの説明については、簡潔にするために本開示からは省略する。さらに、図 3 A に示すディスプレイ装置 1 0 4 はテレビの一例として提示するものにすぎない。本開示は、例えば図 1 で説明したような他のタイプのディスプレイ装置にも適用可能である。このような他のタイプのディスプレイ装置の説明については、簡潔にするために本開示からは省略する。さらに、図 3 A に示すオーディオレンダリング装置 1 0 6 はスピーカの一例として提示するものにすぎない。本開示は、例えば図 1 で説明したような他のタイプのオーディオレンダリング装置にも適用可能である。このような他のタイプのオーディオレンダリング装置の説明については、簡潔にするために本開示からは省略している。

【 0 0 3 6 】

ある実施形態では、メディアコンテンツの一連の撮影シーンの各々が、オーディオ部分、画像部分、及びテキスト情報 3 0 6 を含むことができる。オーディオ部分は、オーディオレンダリング装置 1 0 6 を介してレンダリングできるオーディオフォーマットでの、メディアコンテンツの一連の撮影シーンの各々の 1 又は 2 以上のせりふを含むことができる。各シーンの画像部分は、ディスプレイ装置 1 0 4 上にレンダリングできる 1 又は 2 以上の画像フレームを含むことができる。

【 0 0 3 7 】

テキスト情報 3 0 6 は、ビデオ説明情報 3 0 6 A、タイミング情報、及び/又は速度情報をさらに含むことができる。ビデオ説明情報 3 0 6 A は、一連の撮影シーンのうちの(撮影シーン 3 0 4 又は第 1 の撮影シーン 1 1 4 A などの)撮影シーンを説明することができる。いくつかの実施形態では、ビデオ説明情報 3 0 6 A が、一連の撮影シーンの各撮影シーンを説明することができる。ビデオ説明情報 3 0 6 A 又は撮影シーンの説明の例としては、以下に限定するわけではないが、撮影シーン内の 1 又は 2 以上の装飾品の説明、シーン内の照明条件の説明、撮影シーン内の場所の説明、撮影シーン内のカメラモーションの説明、撮影シーン内の背景情報の説明、撮影シーン内の環境条件の説明、撮影シーン内のショット推移の説明、撮影シーンに含まれるテキストの説明、撮影シーンに描かれるキャラクタの説明、撮影シーンに描かれるキャラクタの態度/感情の説明、撮影シーンに描かれるキャラクタ間の空間的関係の説明、撮影シーンに描かれるキャラクタの身体的属性の説明、撮影シーンに描かれるキャラクタの肉体的表現の説明、撮影シーンに描かれるキャラクタの表情の説明、撮影シーンに描かれるキャラクタの動きの説明、撮影シーンに描かれるキャラクタの職業又は役割の説明、撮影シーンに描かれるキャラクタの服装の説明などを挙げる事ができる。

【 0 0 3 8 】

ある実施形態によれば、回路 2 0 2 は、メディアレンダリング装置 3 0 2 のメモリ 2 0 4 から、一連の撮影シーン及びテキスト情報 3 0 6 を含むメディアコンテンツを検索するように構成することができる。いくつかの実施形態では、回路 2 0 2 を、メディアコンテンツを求める要求をサーバ 1 0 8 に送信するように構成することができる。送信された要求に基づいて、サーバ 1 0 8 から(一連の撮影シーン及びテキスト情報 3 0 6 を含むことができる)メディアコンテンツを受け取ることができる。テキスト情報 3 0 6 内に存在するビデオ説明情報 3 0 6 A は(オーディオフォーマットではなく)テキストフォーマットで受け取ることができ、これによりサーバ 1 0 8 とメディアレンダリング装置 3 0 2 との間におけるテキスト情報 3 0 6 の通信中の帯域幅をさらに節約することができる。テキス

10

20

30

40

50

トフォーマットでのテキスト情報 306 は、メモリ 204 又はサーバ 108 の記憶スペースをさらに節約することができる。ある実施形態では、メディアコンテンツの再生前に、メディアコンテンツから一連の撮影シーン及びテキスト情報 306 の各々を検索することができる。他のいくつかの実施形態では、回路 202 が、特定の撮影シーン（例えば、現在の撮影シーン）の再生時に、さらなる処理（例えば、次の撮影シーンのビデオ説明情報 306 A をオーディオ再生するためのタイミング情報及び速度情報の抽出又は速度の計算）のために次の撮影シーンのテキスト情報 306 を検索することができる。

【0039】

ある実施形態では、撮影シーン 304 が、第 1 のショット及び第 2 のショットなどの複数のショットを含むことができる。撮影シーン 304 は、複数の生物オブジェクト (animated objects) 及び無生物オブジェクト (in-animated objects) を含むことができる。例えば、図 3 A では、生物オブジェクトが、以下に限定するわけではないが、（例えば、「George」という名前の）第 1 の人物 308 及び（例えば、「Jack」という名前の）第 2 の人物 310 を含むことができる。図 3 A に示す無生物オブジェクトは、限定するわけではないが車 312 を含むことができる。図 3 A にはタイムライン 314 をさらに示す。タイムライン 314 は、撮影シーン 304 を再生できる（例えば、秒、分又は時間単位の）時間間隔を示すことができる。タイムライン 314 に示される合計時間は、撮影シーン 304 に関連する全ての画像フレーム及び/又はオーディオフレームをレンダリングするための再生時間とすることができる。

【0040】

図 3 A に示すように、撮影シーン 304 は、撮影シーン 304 のオーディオ部分 316 として第 1 のオーディオ部分 316 A 及び第 2 のオーディオ部分 316 B を含むことができる。第 1 のオーディオ部分 316 A 及び第 2 のオーディオ部分 316 B の各々は、撮影シーン 304 に取り込まれた第 1 の人物 308 及び/又は第 2 の人物 310 に対応する（図 3 A に示す「George: おい Jack、どこに向かっているんだ? (Hey Jack! Where are you heading)」及び「Jack: 仕事だよ (I am going to work)」などの）セリふを含むことができる。一例として、図 3 A に示すように、撮影シーン 304 に関連するタイムライン 314 には、時点 $t_0 \sim t_1$ に、第 1 の途切れ 318 A に対応できる自然な途切れが存在することができる。時点 $t_1 \sim t_2$ には、第 1 の人物 308 が、セリふ又は第 1 のオーディオ部分 316 A として「おい Jack、どこに向かっているんだ?」と発話することができる。さらに、時点 $t_2 \sim t_3$ には、第 2 の途切れ 318 B に対応できる別の自然な途切れが存在することができる。時点 $t_3 \sim t_4$ には、第 2 の人物 310 が、セリふ又は第 2 のオーディオ部分 316 B として、例えば「仕事だよ」というセリふで第 1 の人物 308 に返答することができる。時点 $t_4 \sim t_5$ には、第 3 の途切れ 318 C に対応できる別の自然な途切れが存在することができる。

【0041】

一例として、ビデオ説明情報 306 A、及びタイミング情報を含むことができる検索されたテキスト情報 306 を以下の表 1 に示す。

S. No	ビデオ説明	時間間隔
1.	Jack が車を運転中、George が Jack を見ている	$t_2 \sim t_3$
2.	ビデオ説明 1	$t_A \sim t_B$
3.	ビデオ説明 2	$t_C \sim t_D$

表 1：テキスト情報

【0042】

なお、テキスト情報 306 内の行数は一例として提示するものにすぎない。テキスト情報 306 は、撮影シーン 304 に含まれるビデオ説明の数に基づいてこれよりも多くの又は少ない数の行を含むことができる。

【 0 0 4 3 】

回路 2 0 2 は、撮影シーン 3 0 4 のテキスト情報 3 0 6 からタイミング情報を抽出するようにさらに構成することができる。タイミング情報は、テキスト情報 3 0 6 のビデオ説明情報 3 0 6 A を再生するために抽出することができる。タイミング情報は、ビデオ説明情報 3 0 6 A のテキスト表現又はオーディオ表現、或いはテキスト表現及びオーディオ表現の両方を再生のために収めることができる、タイムライン 3 1 4 内の第 1 の時間間隔（例えば、第 2 の途切れ 3 1 8 B としての時間間隔 $t_2 \sim t_3$ ）を示すことができる。

【 0 0 4 4 】

別の実施形態では、回路 2 0 2 を、撮影シーン 3 0 4 のテキスト情報 3 0 6 から速度情報を抽出するようにさらに構成することができる。タイミング情報と同様に、速度情報も、テキスト情報 3 0 6 のビデオ説明情報 3 0 6 A を再生するために抽出することができる。速度情報は、タイミング情報 1 1 8 B によって示される第 1 の時間間隔（すなわち、第 2 の途切れ 3 1 8 B）中にビデオ説明情報 3 0 6 A のオーディオ表現を再生する再生速度を示すことができる。一例として、ビデオ説明情報 3 0 6 A、タイミング情報及び速度情報を含むことができる検索されたテキスト情報 3 0 6 を以下の表 2 に示す。

S. No	ビデオ説明	時間間隔	再生速度
1.	J a c k が車を運転中、G e o r g e が J a c k を見ている	$t_2 \sim t_3$	1.6倍
2.	ビデオ説明 1	$t_A \sim t_B$	0.5倍
3.	ビデオ説明 2	$t_C \sim t_D$	2.0倍

表 2：テキスト情報

【 0 0 4 5 】

なお、テキスト情報 3 0 6 内の行数は一例として提示するものにすぎない。テキスト情報 3 0 6 は、撮影シーン 3 0 4 に含まれるビデオ説明の数に基づいてこれよりも多くの又は少ない数の行を含むことができる。

【 0 0 4 6 】

テキスト表現の場合には、回路 2 0 2 を、撮影シーン 3 0 4 の抽出されたタイミング情報によって示される第 1 の時間間隔（すなわち、第 2 の途切れ 3 1 8 B）において（テキスト情報 3 0 6 内に存在する）ビデオ説明情報 3 0 6 A をディスプレイ装置 1 0 4 上にレンダリングするように構成することができる。ビデオ説明情報 3 0 6 A のテキスト再生に関する詳細については、例えば図 3 B で説明する。

【 0 0 4 7 】

オーディオ表現の場合には、回路 2 0 2 を、撮影シーン 3 0 4 のテキスト情報 3 0 6 内に存在する検索されたビデオ説明情報 3 0 6 A をビデオ説明情報 3 0 6 A のオーディオ表現に変換するようにテキスト - スピーチ変換器 2 0 8 を制御するようさらに構成することができる。回路 2 0 2 は、撮影シーン 3 0 4 の抽出されたタイミング情報によって示される第 1 の時間間隔（すなわち、第 2 の途切れ 3 1 8 B）においてビデオ説明情報 3 0 6 A のオーディオ表現の再生を制御することができる。ビデオ説明情報 3 0 6 A のオーディオ表現の再生は、抽出された速度情報に基づくことができる。

【 0 0 4 8 】

テキスト表現及びオーディオ表現の両方の場合には、オーディオレンダリング装置 1 0 6 を介してビデオ説明情報 3 0 6 A のオーディオ表現をレンダリングできる第 1 の時間間隔（すなわち、 $t_2 \sim t_3$ ）中に、ビデオ説明情報 3 0 6 A をディスプレイ装置 1 0 4 上に（例えば、テキストフォーマットで）レンダリングすることができる。表 2 によれば、回路 2 0 2 は、第 1 の時間間隔（すなわち、 $t_2 \sim t_3$ ）中に、ビデオ説明情報 3 0 6 A のオーディオ表現（例えば、「J a c k が車を運転中、G e o r g e が J a c k を見ている（

George is looking at Jack while Jack is driving the car)」)の再生を、撮影シーン304のビデオ説明情報306Aのオーディオ表現の実際の再生速度の1.6倍の速度で制御することができる。実際の再生速度は、メディアコンテンツのオーディオをレンダリングできるレート又は速度(すなわち、1倍速)に対応することができる。実際の再生速度は、撮影シーン404の取り込み時にオーディオ部分116が録音されたレート又は速度とすることができる。ビデオ説明情報306Aのオーディオ表現を再生するための第1の時間間隔(すなわち、図3Aに示す $t_2 \sim t_3$)は、テキスト情報306に含まれるタイミング情報によって示すことができ、ビデオ説明情報306Aのオーディオ表現を再生できる速度(すなわち、1.6倍)は、テキスト情報306に含まれる速度情報によって示すことができる。

10

【0049】

限定ではなく一例として、表1によれば、回路202は、時間間隔 $t_A \sim t_B$ 中に、ビデオ説明情報306A(「ビデオ説明1」)のテキスト表現、又はテキスト表現及びオーディオ表現の両方の再生を制御することができる。限定ではなく別の例として、表1によれば、回路202は、時間間隔 $t_A \sim t_B$ 中に、ビデオ説明情報306A(「ビデオ説明1」)のオーディオ表現の再生を、撮影シーン304のビデオ説明情報306Aのオーディオ表現の実際の再生速度の0.5倍の速度で制御することができる。従って、開示するメディアレンダリング装置302は、ディスプレイ装置104及びオーディオレンダリング装置106を介して再生できるメディアコンテンツのテキスト情報306に(例えば、テキスト形態で)含まれるタイミング情報及び/又は速度情報に基づいて、ビデオ説明情報306Aの(テキスト表現、オーディオ表現、又はテキスト表現及びオーディオ表現の両方での)再生のタイミング及び/又は速度を制御することを可能にすることができる。

20

【0050】

ある実施形態では、回路202を、メディアコンテンツのレンダリング前又はその最中にユーザ112に対してディスプレイ装置104上に一連の選択肢を表示するように構成することができる。一連の選択肢のうちの第1の選択肢は、ビデオ説明情報のオーディオ表現(すなわち、ビデオ説明情報をオーディオフォーマットでレンダリングすること)の選択に対応することができる。一連の選択肢のうちの第2の選択肢は、ビデオ説明情報のテキスト表現(すなわち、ビデオ説明情報をテキストフォーマットでレンダリングすること)の選択に対応することができる。同様に、一連の選択肢のうちの第3の選択肢は、ビデオ説明情報のオーディオ表現及びテキスト表現の選択(すなわち、ビデオ説明情報をオーディオ表現及びテキスト表現の両方で同時にレンダリングすること)に対応することができる。いくつかの実施形態では、回路202が、ユーザ120のユーザプロファイルからビデオ説明情報の再生のためのユーザ選好を決定することができる。回路202は、このユーザ選好に基づいてビデオ説明情報の再生(テキストフォーマット、オーディオフォーマット、又はこれらの両方)をさらに制御することができる。

30

【0051】

図3Bには、ディスプレイ装置104及びオーディオレンダリング装置106をさらに含むことができるメディアレンダリング装置302を示す。メディアレンダリング装置302は、ディスプレイ装置104及びオーディオレンダリング装置106を、メディアコンテンツをレンダリングするように制御することができる。メディアコンテンツは、(図3Aに示すような)ディスプレイ装置104上に表示された撮影シーン304を含むことができる(図1の一連の撮影シーン114などの)一連の撮影シーンを含むことができる。

40

【0052】

ある実施形態では、メディアコンテンツの一連の撮影シーンの各々が、オーディオ部分、画像部分、テキスト情報306、及びクローズドキャプション情報320を含むことができる。オーディオ部分は、オーディオレンダリング装置106を介してレンダリングできるオーディオフォーマットでの、メディアコンテンツの一連の撮影シーンの各々の1又は2以上のせりふを含むことができる。各シーンの画像部分は、ディスプレイ装置104上にレンダリングできる1又は2以上の画像フレームを含むことができる。クローズドキ

50

ャプション情報 320 は、撮影シーン 304 の再生中に（図 3 B に示すような）ディスプレイ装置 104 上にレンダリングできるテキストフォーマットでの、撮影シーン 304 のオーディオ部分 116 を表すことができる。クローズドキャプション情報 320 は、撮影シーン 304 のオーディオ部分の転写とみなすことができる。いくつかの実施形態では、ビデオ説明情報 306 A（すなわち、シーン説明）をクローズドキャプション情報 320 と共にメディアコンテンツ内に符号化することができる。

【0053】

ある実施形態では、撮影シーン 304 が、第 1 のショット及び第 2 のショットなどの複数のショットを含むことができる。撮影シーン 304 は、複数の生物オブジェクト及び無生物オブジェクトを含むことができる。例えば、図 3 B では、生物オブジェクトが、以下

10

【0054】

ある実施形態では、図 3 B に示すように、第 1 の時間間隔（すなわち、図 3 A に示す $t_2 \sim t_3$ ）中に、「Jack が車を運転中、George が Jack を見ている」というビデオ説明情報 306 A をディスプレイ装置 104 上にテキストフォーマットでレンダリングすることができる。別の実施形態では、表 2 に従って、回路 202 が、第 1 の時間間隔（ $t_2 \sim t_3$ ）中に、「Jack が車を運転中、George が Jack を見ている」というビデオ説明情報 306 A のオーディオ表現の再生を、撮影シーン 304 のビデオ説明情報 306 A のオーディオ表現の実際の再生速度（すなわち、1.0 倍速）の 1.6 倍の速度で制御することができる。ある実施形態では、図 3 B に示すように、第 1 の時間間隔（すなわち、図 3 A に示す $t_2 \sim t_3$ ）中に、ビデオ説明情報 306 A をクローズドキャプション情報 320 の表示と共にディスプレイ装置 104 上にテキストフォーマットでレンダリングしながら、オーディオレンダリング装置 106 を介してビデオ説明情報 306 A のオーディオ表現をレンダリングすることもできる。図 3 B に示すように、ビデオ説明情報 306 A 及びクローズドキャプション情報 320 は、撮影シーン 304 の表示時に（画像フレームなどの）画像部分にオーバーレイ表示できるテキストフォーマットでディスプレイ装置 104 上にレンダリングすることができる。いくつかの実施形態では、ビデオ説明情報 306 A のオーディオ表現を再生する代わりに、第 1 の時間間隔（ $t_2 \sim t_3$ ）中にビデオ説明情報 306 A 及びクローズドキャプション情報 320 を同時にディスプレイ装置 104 上にレンダリングすることができる。

20

30

【0055】

図 4 は、本開示の実施形態による、シーン説明の再生制御のための別の例示的なシナリオを示す図である。図 4 の説明は、図 1、図 2、図 3 A 及び図 3 B の要素に関連して行う。図 4 には例示的なシナリオ 400 を示す。例示的なシナリオ 400 には、メディアレンダリング装置 402（すなわち、メディアレンダリング装置 102 と同様のもの）を示す。図 4 には、メディアレンダリング装置 402 に関連するディスプレイ装置 104 及びオーディオレンダリング装置 106 をさらに示す。メディアレンダリング装置 402 は、ディスプレイ装置 104 及びオーディオレンダリング装置 106 を、メディアコンテンツをレンダリングように制御することができる。メディアコンテンツは、ディスプレイ装置 104 上に表示された撮影シーン 404 を含むことができる（図 1 の一連の撮影シーン 114 などの）一連の撮影シーンを含むことができる。

40

【0056】

ある実施形態では、メディアコンテンツの一連の撮影シーンの各々が、オーディオ部分、画像部分、及びビデオ説明情報 406 を含むことができる。いくつかの実施形態では、一連の撮影シーンの各々が、（例えば、図 3 A で説明したようなビデオ説明情報 406 を含むことができるテキスト情報 306 などの）テキスト情報を含むことができる。オーディオ部分は、オーディオレンダリング装置 106 を介してレンダリングできるオーディオフォーマットでの、メディアコンテンツの一連の撮影シーンの各々の 1 又は 2 以上のせり

50

ふを含むことができる。各シーンの画像部分は、ディスプレイ装置 104 上にレンダリングできる 1 又は 2 以上の画像フレームを含むことができる。ビデオ説明情報 406 A は、一連の撮影シーンのうちの（撮影シーン 404 又は第 1 の撮影シーン 114 A などの）撮影シーンを説明することができ、撮影シーン 404 はディスプレイ装置 104 上に表示することができる。いくつかの実施形態では、ビデオ説明情報 406 が、一連の撮影シーンの各撮影シーンを説明することができる。

【0057】

ある実施形態によれば、回路 202 は、メディアレンダリング装置 402 のメモリ 204 から（一連の撮影シーン及びビデオ説明情報 406 を含むことができる）メディアコンテンツを検索するように構成することができる。いくつかの実施形態では、回路 202 を、メディアコンテンツを求める要求をサーバ 108 に送信するように構成することができる。送信された要求に基づいて、サーバ 108 から（一連の撮影シーン及びビデオ説明情報 406 を含むことができる）メディアコンテンツを受け取ることができる。ビデオ説明情報 406 A は（オーディオフォーマットではなく）テキストフォーマットで受け取ることができ、これによりサーバ 108 とメディアレンダリング装置 402 との間におけるビデオ説明情報 406 の通信中の帯域幅をさらに節約することができる。テキストフォーマットでのビデオ説明情報 406 は、メモリ 204 又はサーバ 108 の記憶スペースをさらに節約することができる。ある実施形態では、メディアコンテンツの再生前に、メディアコンテンツから一連の撮影シーン及びビデオ説明情報 406 の各々を検索することができる。他のいくつかの実施形態では、回路 202 が、特定の撮影シーン（例えば、現在の撮影シーン）の再生時に、さらなる処理（例えば、次の撮影シーンのビデオ説明情報 406 をオーディオ再生するための速度の計算）のために次の撮影シーンのビデオ説明情報 406 を検索することができる。

【0058】

回路 202 は、撮影シーン 404 の検索されたビデオ説明情報 406 をビデオ説明情報 406 のオーディオ表現に変換するようにテキスト - スピーチ変換器 208 を制御するようにさらに構成することができる。いくつかの実施形態では、撮影シーン 404 に関する情報がクロズドキャプション情報も含む。例えば図 3 B で説明したように、クロズドキャプション情報は、撮影シーン 304 の表示時に（画像フレームなどの）画像部分にオーバーレイ表示できる、テキストフォーマットでの撮影シーン 404 のオーディオ部分 116 を表すことができる。いくつかの実施形態では、ビデオ説明情報 406（すなわち、シーン説明）をクロズドキャプション情報と共にメディアコンテンツ内に符号化することができる。

【0059】

ある実施形態では、撮影シーン 404 が、第 1 のショット及び第 2 のショットなどの複数のショットを含むことができる。撮影シーン 404 は、複数の生物オブジェクト及び無生物オブジェクトを含むことができる。例えば、図 4 では、生物オブジェクトが、以下に限定するわけではないが、（例えば、「George」という名前の）第 1 の人物 408 及び（例えば、「Jack」という名前の）第 2 の人物 410 を含むことができる。図 4 に示す無生物オブジェクトは、限定するわけではないが車 312 を含むことができる。図 4 にはタイムライン 414 をさらに示す。タイムライン 414 は、撮影シーン 404 を再生できる（例えば、秒、分又は時間単位の）時間間隔を示すことができる。タイムライン 414 に示される合計時間は、撮影シーン 404 に関連する全ての画像フレーム及び／又はオーディオフレームをレンダリングするための再生時間とすることができる。タイムライン 414 は、撮影シーン 404 における第 1 の人物 408 と第 2 の人物 410 との間の会話中に発せられるせりふに対応できる一連の第 2 の時間間隔 416 を含むことができる。

【0060】

図 4 に関しては、メディアコンテンツ又は（図 1 に示すテキスト情報 118 などの）テキスト情報がタイミング情報及び速度情報（すなわち、例えば図 3 A で説明したものを）を含んでいないと仮定することができる。従って、開示するメディアレンダリング装置 10

10

20

30

40

50

2 は、ビデオ説明情報 4 0 6 のオーディオ表現を再生するための速度及び第 1 の時間間隔を決定することができる。ある実施形態によれば、回路 2 0 2 は、(第 1 のオーディオ部分 4 1 6 A 及び第 2 のオーディオ部分 4 1 6 B などの) オーディオ部分 1 1 6 を含むことができる撮影シーン 4 0 4 の一連の第 2 の時間間隔 4 1 6 を決定するようにさらに構成することができる。一連の第 2 の時間間隔 4 1 6 の各々は、一連の撮影シーンにおける撮影シーン 4 0 4 のオーディオ部分 1 1 6 を再生するための時間間隔を示すことができる。例えば、図 4 に示すように、撮影シーン 4 0 4 は、撮影シーン 4 0 4 のオーディオ部分 1 1 6 として第 1 のオーディオ部分 4 1 6 A 及び第 2 のオーディオ部分 4 1 6 B を含むことができる。第 1 のオーディオ部分 4 1 6 A 及び第 2 のオーディオ部分 4 1 6 B の各々は、撮影シーン 4 0 4 に取り込まれた第 1 の人物 4 0 8 及び / 又は第 2 の人物 4 1 0 に対応する (図 4 に示す「George: おい Jack、どこに向かってるんだ?」及び「Jack: 仕事だよ」などの) セリふを含むことができる。回路 2 0 2 は、撮影シーン 4 0 4 に含まれる各オーディオフレームのオーディオ分析に基づいて、撮影シーン 4 0 4 における一連の第 2 の時間間隔 4 1 6 を決定するように構成することができる。オーディオ分析では、回路 2 0 2 が、各オーディオフレーム内のオーディオ音量又はピッチをオーディオ閾値 (dB 単位) と比較して、撮影シーン 4 0 4 に関連するセリふ又は音楽を含むことができる一連の第 2 の時間間隔 4 1 6 を決定することができる。

【 0 0 6 1 】

ある実施形態では、回路 2 0 2 を、撮影シーン 4 0 4 のビデオ説明情報 4 0 6 のオーディオ表現の第 3 の時間間隔 4 1 8 (すなわち、図 4 に示すような「t₀₀」~「t₀₁」の時間間隔) を決定するようにさらに構成することができる。第 3 の時間間隔 4 1 8 は、ビデオ説明情報 4 0 6 のオーディオ表現をその実際の再生速度でプレイバック又は再生するために必要な期間 (例えば、数秒単位) に対応することができる。この時間間隔は、ユーザ 1 2 0 がビデオ説明情報 4 0 6 を表示する選択肢を選択した場合に (図 3 B に示すような) ディスプレイ装置 1 0 4 上にビデオ説明情報 4 0 6 のテキストフォーマットを表示できる期間であることもできる。実際の再生速度は、メディアコンテンツのオーディオをレンダリングできるレート又は速度 (すなわち、1 倍速) に対応することができる。実際の再生速度は、撮影シーン 4 0 4 の取り込み時にオーディオ部分 1 1 6 が録音されたレート又は速度とすることができる。ある実施形態では、第 3 の時間間隔 4 1 8 が、ビデオ説明情報 4 0 6 のサイズに基づくことができる。例えば、撮影シーン 4 0 4 を説明するためにより多くの数の単語がビデオ説明情報 4 0 6 に含まれている場合には、ビデオ説明情報 4 0 6 のオーディオ表現を実際の再生速度で再生するための第 3 の時間間隔 4 1 8 の期間も長くなることができる。

【 0 0 6 2 】

ある実施形態によれば、回路 2 0 2 は、撮影シーン 4 0 4 の一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C を決定するようにさらに構成することができる。一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C の各々は、一連の第 2 の時間間隔 4 1 6 とは異なることができ、撮影シーン 4 0 4 のタイムライン 4 1 4 内の自然な途切れ (又は空白) に対応できる全ての間隔を含むことができる。図 4 に示すように、一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C は、第 1 の途切れ 4 2 0 A、第 2 の途切れ 4 2 0 B、及び第 3 の途切れ 4 2 0 C を含むことができる。回路 2 0 2 は、撮影シーン 4 0 4 に含まれる各オーディオフレームのオーディオ分析に基づいて、撮影シーン 4 0 4 内の自然な途切れ又は空白 (すなわち、一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C に対応する途切れ又は空白) を決定するように構成することができる。オーディオ分析では、回路 2 0 2 が、各オーディオフレーム内のオーディオ音量又はピッチをオーディオ閾値 (dB 単位) と比較することができる。オーディオフレーム内のオーディオ音量又はピッチが (例えば、dB 単位の) オーディオ閾値よりも小さい場合には、対応するオーディオフレームを撮影シーン 4 0 4 内の自然な途切れ又は空白として決定することができる。回路 2 0 2 は、撮影シーン 4 0 4 に含まれる第 1 の途切れ 4 2 0 A、第 2 の途切れ 4 2 0 B 又は第 3 の途切れ 4 2 0 C などの決定された途切れ又は空白を再生するための一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C 又は期間を決定するように

さらに構成することができる。

【 0 0 6 3 】

一例として、図 4 に示すように、撮影シーン 4 0 4 に関連するタイムライン 4 1 4 には、時点 $t_0 \sim t_1$ に、第 1 の途切れ 4 2 0 A に対応できる自然な途切れが存在することができる。時点 $t_1 \sim t_2$ には、第 1 の人物 4 0 8 が、せりふ又はオーディオ部分 1 1 6 として「おい」ack、どこに向かっているんだ？」と発話することができる。さらに、時点 $t_2 \sim t_3$ には、第 2 の途切れ 4 2 0 B に対応できる別の自然な途切れが存在することができる。時点 $t_3 \sim t_4$ には、第 2 の人物 4 1 0 が、例えば「仕事だよ」というせりふで第 1 の人物 4 0 8 に返答することができる。時点 $t_4 \sim t_5$ には、第 3 の途切れ 4 2 0 C に対応できる別の自然な途切れが存在することができる。従って、図 4 に示すように、一連の第 2 の時間間隔 4 1 6 は、時点 t_1 から t_2 に及ぶことができる第 1 のオーディオ部分 4 1 6 A、及び時点 t_3 から t_4 に及ぶことができる第 2 のオーディオ部分 4 1 6 B を含むことができる。一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C は、時点 t_0 から t_1 に及ぶことができる第 1 の途切れ 4 2 0 A、時点 t_2 から t_3 に及ぶことができる第 2 の途切れ 4 2 0 B、及び時点 t_4 から t_5 に及ぶことができる第 3 の途切れ 4 2 0 C を含むことができる。

10

【 0 0 6 4 】

回路 2 0 2 は、撮影シーン 4 0 4 の一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C から第 1 の時間間隔 4 2 2 を選択するようにさらに構成することができる。第 1 の時間間隔 4 2 2 は、時間間隔閾値の期間よりも長い期間を有することができる時間間隔であることができ、ビデオ説明情報 4 0 6（すなわち、シーン説明）のオーディオ再生のための潜在的空白とみなすことができる。時間間隔閾値は、第 1 の人物 4 0 8 又は第 2 の人物 4 1 0 が特定のせりふを発話している間に発生し得る短い途切れ又は空白をフィルタ除去するために利用される（例えば、ミリ秒又は数秒単位の）所定の時間値とすることができる。例えば、時間間隔閾値は、第 1 の人物 4 0 8 又は第 2 の人物 4 1 0 が複数のせりふ間に息を吸う / 吐くために要する時間を示すことができる。

20

【 0 0 6 5 】

ある実施形態では、回路 2 0 2 が、一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C の各々と時間間隔閾値との比較に基づいて第 1 の時間間隔 4 2 2 を選択することができる。時間間隔閾値は、ビデオ説明情報 4 0 6 のオーディオ再生が不可能と考えられる間隔の値に対応することができる。換言すれば、時間間隔閾値は、それ未満ではビデオ説明情報 4 0 6 のオーディオ再生がメディアコンテンツのレンダリング対象であるユーザ 1 2 0 に対して十分に詳細なシーン説明を提供できないと考えられるタイミング値に対応することができる。

30

【 0 0 6 6 】

例えば、第 1 の途切れ 4 2 0 A の期間が 0 . 7 5 秒であり、第 2 の途切れ 4 2 0 B の期間が 1 秒であり、第 3 の途切れ 4 2 0 C の期間が 0 . 5 秒であり、時間間隔閾値が 1 秒である場合、回路 2 0 2 は、一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C 内の各途切れの期間と時間間隔閾値とを比較し、時間間隔閾値以上の期間を有する第 2 の途切れ 4 2 0 B を第 1 の時間間隔 4 2 2 として選択することができる。いくつかの実施形態では、期間が長くなるとビデオ説明情報 4 0 6（すなわち、シーン説明）の再生速度が実際の再生速度と同じになり、従ってビデオ説明情報 4 0 6 のオーディオ再生の品質を維持することができるので、回路 2 0 2 は、（第 1 の途切れ 4 2 0 A、第 2 の途切れ 4 2 0 B、又は第 3 の途切れ 4 2 0 C のうちの）最も長い期間を有する途切れを第 1 の時間間隔 4 2 2 として選択することができる。

40

【 0 0 6 7 】

ある実施形態によれば、回路 2 0 2 は、ビデオ説明情報 4 0 6 のオーディオ表現を再生する再生速度を決定するようにさらに構成することができる。再生速度は、ビデオ説明情報 4 0 6 のオーディオ表現の再生速度に対応することができる。いくつかの実施形態では、回路 2 0 2 が乗算係数（multiplication factor）を計算し、計算された乗算係数及びビデオ説明情報 4 0 6 のオーディオ表現の実際の再生速度に基づいて再生速度を決定することができる。乗算係数は、決定された第 3 の時間間隔 4 1 8 及び

50

選択された第 1 の時間間隔 4 2 2 に基づいて計算することができる。

【 0 0 6 8 】

ある例では、撮影シーン 4 0 4 内の第 1 の途切れ 4 2 0 A (時点 $t_0 \sim t_1$) の期間が 2 秒であり、第 2 の途切れ 4 2 0 B (時点 $t_2 \sim t_3$) の期間が 3 秒であり、第 3 の途切れ 4 2 0 C (時点 $t_4 \sim t_5$) の期間が 2 秒である。第 3 の時間間隔 4 1 8 の期間が 5 秒である場合、このような期間は、一連の第 4 の時間間隔 4 2 0 A ~ 4 2 0 C (すなわち、第 1 の途切れ 4 2 0 A、第 2 の途切れ 4 2 0 B、及び第 3 の途切れ 4 2 0 C) の各々又は選択された第 1 の時間間隔 4 2 2 に対応する時間間隔中にビデオ説明情報 4 0 6 を実際の再生速度で聞き取れるように再生するには不十分と考えられる。回路 2 0 2 は、以下の方程式 (1) を使用して乗算係数を決定するように構成することができる。

10

$$\text{乗算係数} = \frac{\text{第 3 の時間間隔}}{\text{第 1 の時間間隔}}$$

(1)

【 0 0 6 9 】

回路 2 0 2 は、計算された乗算係数及び実際の再生速度に基づいて、以下の方程式 (2) を使用することによって、ビデオ説明情報 4 0 6 のオーディオ表現を再生する再生速度を決定するようにさらに構成することができる。

$$\text{再生速度} = \text{乗算係数} \times \text{実際の再生速度} (2)$$

20

【 0 0 7 0 】

上述した例を参照すると、回路 2 0 2 は、方程式 (1) を使用することにより、乗算係数を 1 . 6 6 (すなわち、5 秒である第 3 の時間間隔 4 1 8 と、3 秒である第 2 の途切れ 4 2 0 B として選択された第 1 の時間間隔 4 2 2 との比率) であると決定するように構成することができる。乗算係数が 1 . 0 よりも大きい (すなわち、第 3 の時間間隔 4 1 8 が第 1 の時間間隔 4 2 2 よりも大きい) 場合、回路 2 0 2 は、ビデオ説明情報 4 0 6 のオーディオ表現の実際の再生速度を乗算係数によって増加させるように構成することができる。例えば、乗算係数が 1 . 6 6 である場合、回路 2 0 2 は、撮影シーン 4 0 4 のビデオ説明情報 4 0 6 のオーディオ表現の実際の再生速度の 1 . 6 6 倍を再生速度として決定することができる。その他の事例では、乗算係数が 1 . 0 未満である場合 (すなわち、第 3 の時間間隔 4 1 8 が第 1 の時間間隔 4 2 2 よりも小さい場合)、回路 2 0 2 は、ビデオ説明情報 4 0 6 のオーディオ表現の実際の再生速度を乗算係数によって減少させるように構成することができる。例えば、乗算係数が 0 . 8 である場合、回路 2 0 2 は、撮影シーン 4 0 4 のビデオ説明情報 4 0 6 のオーディオ表現の実際の再生速度の 0 . 8 倍を再生速度として決定することができる。いくつかの実施形態では、乗算係数が 1 . 0 未満である場合、回路 2 0 2 は実際の再生速度を変更せず、ビデオ説明情報 4 0 6 のオーディオ表現の再生速度は実際の再生速度と同じままであることができる (例えば、乗算係数が 0 . 9 5 である場合には実質的に 1 . 0 に近いと考えることができる)。他のいくつかの実施形態では、乗算係数が 1 . 0 に等しい場合 (すなわち、第 3 の時間間隔 4 1 8 が第 1 の時間間隔 4 2 2 に等しい場合)、回路 2 0 2 は、ビデオ説明情報 4 0 6 のオーディオ表現の実際の再生速度を再生速度として決定するように構成することができる。

30

40

【 0 0 7 1 】

回路 2 0 2 は、決定された再生速度に基づいて、ビデオ説明情報 4 0 6 のオーディオ表現の再生を第 1 の時間間隔 4 2 2 において制御するようにさらに構成することができる。第 1 の時間間隔 4 2 2 (すなわち、途切れのうちの 1 つ) は、一連の第 2 の時間間隔 4 1 6 (すなわち、撮影シーン 4 0 4 のオーディオ部分を含む第 2 の時間間隔) とは異なることができる。いくつかの実施形態では、第 1 の時間間隔 4 2 2 を、撮影シーン 4 0 4 の第 1 のせりふ (例えば、第 1 のオーディオ部分 4 1 6 A) と第 2 のせりふ (例えば、第 2 のオーディオ部分 4 1 6 B) との間とすることができる。例えば、図 4 に示すように、ビデオ説明情報 4 0 6 のオーディオ表現 (すなわち、「 J a c k が車を運転中、 G e o r g e

50

が「Jackを見ている」というシーン説明)は、第1の時間間隔422において、決定された再生速度で(例えば、第3の時間間隔418が5秒であり、第1の時間間隔422が3秒である場合には1.66倍で)再生することができる。従って、回路202は、ビデオ説明情報406の一部(例えば、特定の文字、テキスト又は単語)を短縮又は削除することなく、オーディオセリふの空白(すなわち、第1の時間間隔422)間のビデオ説明情報406(すなわち、シーン説明)のオーディオ再生速度を増加させることができる。この速度の増加により、ビデオ説明情報406の第3の時間間隔418よりも短い期間である第1の時間間隔422内にビデオ説明情報406のオーディオ表現を効果的に組み込み又は収めることができる。従って、たとえ決定された空白(すなわち、撮影シーン404内の特定の空白の第1の時間間隔422)が第3の時間間隔418(すなわち、シーン/ビデオ説明を聞き取れるように再生するのに必要な時間)より短い場合でも、ビデオ説明情報406の再生品質が維持される。

10

【0072】

いくつかの実施形態では、第1のセリふを撮影シーン404の第1のショットの最後の単語とすることができ、第2のセリふを撮影シーン404の第2のショットの最初の単語とすることができ、第1のショット及び第2のショットは、撮影シーン404の連続するショットとすることができ、他のいくつかの実施形態では、第1の時間間隔422を、撮影シーン404の開始と撮影シーン404の(第1のオーディオ部分416Aなどの)第1のセリふとの間とすることができ、このような場合、第1の時間間隔422は、図4に示すような第1の途切れ420Aに対応することができる。

20

【0073】

なお、図4に示す撮影シーン404、及び複数の生物オブジェクト又は無生物オブジェクトは、一例として提示するものにすぎない。本開示は、他のタイプの撮影シーン(例えば、以下に限定するわけではないが、アクションシーン、恋愛シーン、ドラマシーン、ダンスシーン又は音楽シーン)及び複数の生物オブジェクト又は無生物オブジェクトにも適用可能である。他のタイプの撮影シーン404及び複数の生物オブジェクト又は無生物オブジェクト、或いはこれらの例の説明については、簡潔にするために本開示からは省略する。

【0074】

ある実施形態では、回路202を、一連の撮影シーン114の各々について、対応する撮影シーンのオーディオ部分を再生するための時間間隔をそれぞれが示すことができる一連の第2の時間間隔を決定するように構成することができる。回路202は、一連の撮影シーン114のうちの対応する撮影シーンのビデオ説明情報のオーディオ表現の第3の時間間隔を決定するようにさらに構成することができる。撮影シーン404に関して上述したように、回路202は、各シーンの決定された一連の第2の時間間隔及び決定された第3の時間間隔に基づいて、ビデオ説明情報406のオーディオ表現を再生する速度を決定するようにさらに構成することができる。回路202は、決定された速度に基づいて、一連の撮影シーン114の各撮影シーンのビデオ説明情報のオーディオ表現の再生を第1の時間間隔(すなわち、一連の第2の時間間隔とは異なる時間間隔)において制御するようにさらに構成することができる。従って、開示するメディアレンダリング装置402は、対応する撮影シーン又は以前の撮影シーン(すなわち、対応する撮影シーンの直前のシーン)の再生中に、メディアコンテンツ内の各撮影シーンを処理し、対応する撮影シーンの第1の時間間隔422を選択し、撮影シーンに関連するビデオ説明情報406の再生速度を決定することができる。さらに、メディアレンダリング装置402は、一連の撮影シーン114内の各撮影シーンの決定された再生速度に基づいて、対応するビデオ説明情報のオーディオ表現(すなわち、シーン説明)の再生を動的に制御することができる。従って、開示するメディアレンダリング装置402は、例えば視覚障害者又は認知障害者などのユーザ120のコンテンツ体験を強化することができる。

30

40

【0075】

ある実施形態では、回路202を、ユーザ120からI/O装置206を介して第1の

50

ユーザ入力を受け取るようにさらに構成することができる。第1のユーザ入力はテキストフォーマットであることができ、ビデオ説明情報406、又は一連の撮影シーン114のうちの1つの撮影シーンのシーン説明に対応することができる。回路202は、メディアコンテンツの再生中又はメディアコンテンツの再生開始前に第1のユーザ入力を受け取ることができる。第1のユーザ入力は、一連の撮影シーン114のうちの撮影シーンのうちの1つの撮影シーンのビデオ説明情報406に含めることができるテキスト単語又は表現とすることができる。例えば、図4に示すビデオ説明情報406は、「Jackが車を運転中、GeorgeがJackを見ている」であることができる。受け取られた第1のユーザ入力は、ビデオ説明情報406の一部であることができる単語又は表現（例えば、「GeorgeがJackを見ている」）を含むことができる。

10

【0076】

回路202は、一連の撮影シーン114の各々に関連する記憶されたビデオ説明情報406内で、受け取られた第1のユーザ入力を検索するようにさらに構成することができる。いくつかの実施形態では、第1のユーザ入力を受け取られたテキスト説明が、一連の撮影シーン114のうちの1つの撮影シーンのビデオ説明情報406と全く同じものであることができる。他の実施形態では、第1のユーザ入力が、ビデオ説明情報406の一部であることができる。回路202は、検索に基づいて、メディアコンテンツを再生するための再生タイミング情報を決定するようにさらに構成することができる。回路202は、検索に基づいて再生タイミング情報を決定するために、撮影シーン（例えば、撮影シーン404）、及び第1のユーザ入力を含む対応するビデオ説明情報406を決定することができる。このような場合、再生タイミング情報は、決定された撮影シーンの再生タイミングであることができる。他のいくつかの実施形態では、第1のユーザ入力を受け取られたテキスト説明が、一連の撮影シーン114の各々に関連するビデオ説明情報406と全く同じではないことがある。このようなシナリオでは、回路202を、第1のユーザ入力において受け取られたテキスト説明と、一連の撮影シーン114の各々に関連するビデオ説明情報406との間の類似性スコアを決定するように構成することができる。類似性スコアは、テキスト説明と対応する撮影シーンのビデオ説明情報406の部分との一致に基づいて決定することができる。いくつかの実施形態では、類似性スコアを、メディアコンテンツの一連の撮影シーン114の各々に関連する人気度スコアに基づいて計算することができる。回路202は、サーバ108から各撮影シーンの人気スコアを検索することができる。ある実施形態では、サーバ108又はメモリ204から検索されたメディアコンテンツに各撮影シーン的人气スコアを含めることができる。例えば、第1のユーザ入力（すなわち、説明）が「GeorgeがJackを見ている」という単語であり、この単語が、撮影シーン404を含む複数の撮影シーンに関連するビデオ説明情報406内に存在し得るものとする。このような場合、回路202は、複数の撮影シーンの中の各撮影シーン的人气スコアを抽出し、どのシーンが人々の間で人気が高く、ユーザ120が人気の高い撮影シーンのビデオ説明情報406の説明を検索したいと思っている確率が高い（例えば、撮影シーン404）のはどのシーンであるかを識別することができる。回路202は、受け取られた説明（又は第1のユーザ入力）の類似度スコアが高い識別された撮影シーン（例えば、撮影シーン404）の再生タイミング情報を決定するようにさらに構成することができる。回路202は、決定された再生タイミング情報（ t_0 ）に基づいて、識別された撮影シーンからのメディアコンテンツの再生を制御するようにさらに構成することができる。従って、開示するメディアレンダリング装置402は、ユーザ120がメディアコンテンツの一連の撮影シーン114の各々の記憶されたビデオ説明情報406（すなわち、シーン説明）内の単語又はテキストを検索し、従って検索に基づいて識別できる識別された撮影シーンの再生タイミングを制御（すなわち、早送り又は巻き戻し）することを可能にすることができる。従って、メディアレンダリング装置402は、メディアコンテンツ内の1又は2以上のシーンに対応する説明をユーザ120が検索できるようにする検索エンジン機能を提供することができる。

20

30

40

【0077】

50

ある実施形態では、メディアレンダリング装置 402 を、一定期間（例えば、最後の 1 日又は 1 週間）内に第 1 のユーザ入力で受け取られた以前の検索説明に基づいて、新たなメディアコンテンツの個人化された推奨を提供するようにさらに構成することができる。一例として、ユーザ 120 が特定の期間内に「アクション」という単語を検索した頻度が高い場合、回路 202 は、「アクション」ジャンルに関連し得る他の又は新たなメディアコンテンツの推奨を提供することができる。従って、開示するメディアレンダリング装置 402 は、ユーザ 120 が頻繁に検索していると考えられるシーン又はビデオ説明に関連するメディアコンテンツを推奨することができる。

【0078】

ある実施形態では、メディアレンダリング装置 402 を、メディアレンダリング装置 402 に関連する第 1 の規定の速度設定をメモリ 204 に記憶するように構成することができる。第 1 の規定の速度設定は、ビデオ説明情報 406（すなわち、シーン説明）のオーディオ表現の最大再生速度を示すことができる。第 1 の規定の速度設定によって示される最大速度は、メディアコンテンツのレンダリング対象であるユーザ 120 がビデオ説明情報 406 のオーディオ表現を正しく理解できる速度とすることができる。例えば、最大速度は、実際の再生速度の 2 倍とすることができる。いくつかの実施形態では、第 1 の規定の速度設定が、再生速度を決定できる基になる乗算係数の最大値（例えば、2.0）を示すことができる。

【0079】

別の実施形態では、メディアレンダリング装置 402 を、メディアレンダリング装置 402 に関連する第 2 の規定の速度設定をメモリ 204 に記憶するように構成することができる。第 2 の規定の速度設定は、ビデオ説明情報 406（すなわち、シーン説明）のオーディオ表現の最小再生速度を示すことができる。第 2 の規定の速度設定によって示される最小速度は、メディアコンテンツのレンダリング対象であるユーザ 120 がビデオ説明情報 406 のオーディオ表現を正しく理解できる速度とすることができる。例えば、最小速度は、実際の再生速度の 0.5 倍とすることができる。いくつかの実施形態では、第 2 の規定の速度設定が、再生速度を決定できる基になる乗算係数の最小値（例えば、0.5）を示すことができる。

【0080】

ある実施形態によれば、回路 202 を、ビデオ説明情報 406 のオーディオ表現の決定された再生速度、及び第 1 / 第 2 の規定の速度設定に基づいて、撮影シーン 404 の画像部分又はオーディオ部分の一方の再生を制御するようにさらに構成することができる。撮影シーン 404 の画像部分又はオーディオ部分（すなわち、せりふ）の一方の再生制御は、自然な途切れ（すなわち、第 1 の時間間隔 422）が、決定された再生速度及び第 1 又は第 2 の規定の速度設定に基づいてビデオ説明情報 406 のオーディオ表現を収めることができるほど十分に長い場合の、撮影シーンの画像部分及び / 又はオーディオ部分のレンダリングの時間遅延又は一時停止に対応することができる。

【0081】

一例として、第 1 の規定の速度設定（すなわち、最大速度）がビデオ説明情報 406 のオーディオ表現の実際の再生速度の 2 倍であり、第 3 の時間間隔 418 が 7 秒であり、第 1 の時間間隔 422 の期間が 3 秒である場合、方程式（1）によれば、決定された再生速度は 2.33 倍となる。決定された再生速度が最大速度（すなわち、2 倍）よりも高いので、回路 202 は、一連の第 4 の時間間隔 420A ~ 420C から選択された第 1 の時間間隔 422 を廃棄することができる。このような場合、回路 202 は、ビデオ説明情報 406 のオーディオ表現をレンダリングするために撮影シーン 404 の画像部分又はオーディオ部分（すなわち、図 4 の第 2 のオーディオ部分 416B などのせりふ）を一時停止することができる。別の事例では、回路 202 が、レンダリングされているメディアコンテンツの品質を維持するために、ビデオ説明情報 406 を（2 倍のような）最大速度で聞き取れるようにレンダリングし、撮影シーンの画像部分又はオーディオ部分を（第 3 の時間間隔 418 が 7 秒であり、第 1 の時間間隔 422 の期間が 3 秒である場合の残りの 1 秒な

10

20

30

40

50

どの) 残りの時間にわたって一時停止することができる。

【0082】

ある実施形態では、回路202を、I/O装置206を介してユーザ120から第2のユーザ入力を受け取るようにさらに構成することができる。第2のユーザ入力は、メディアコンテンツをレンダリングできる対象であるユーザ120のプロファイル情報を示すことができる。プロファイル情報は、ビデオ説明情報406を聞き取れるようにレンダリングするためのユーザ120の過去の速度選好を含むことができる。いくつかの実施形態では、プロファイル情報が、ユーザ120に関連し得る一意の識別番号(例えば、以下に限定するわけではないが、社会保障番号(SSN)、電話番号、又は保険証券番号)を示すことができる。回路202は、受け取られた一意の識別番号に基づいて、サーバ108又はメモリ204からユーザ120の年齢を検索するようにさらに構成することができる。いくつかの実施形態では、回路202を、ユーザ120に関連する一意の識別番号に基づいてユーザ120の健康状態を決定するようにさらに構成することができる。健康状態は、ビデオ説明情報406のオーディオ表現又は撮影シーンのオーディオ部分(すなわち、せりふ)を特定の再生速度で理解するためのユーザ120の聞き取り能力の欠如を示すことができる。回路202は、受け取られた第2のユーザ入力に基づいて、ビデオ説明情報406のオーディオ表現を再生する再生速度を決定するようにさらに構成することができる。

10

【0083】

一例として、ユーザ120の年齢が65歳(すなわち、老齢)として決定された場合、回路202は、ビデオ説明情報406のオーディオ表現の実際の再生速度の1.5倍を再生速度として決定することができる。いくつかの実施形態では、回路202が、決定された年齢に基づいて(例えば、1.5倍を最大速度とする)第1の速度設定を定めることができる。別の例として、ユーザ120の健康状態によってユーザ120が過去の所定の期間内(例えば、過去6ヶ月以内)に耳の手術を受けたことが示される場合、回路202は、ビデオ説明情報406のオーディオ表現の実際の再生速度の1.2倍を第1の速度設定として定め、又は再生速度として決定することができる。従って、開示するメディアレンダリング装置402は、ユーザ120の(年齢又は健康状態などの)プロファイル情報に基づいて、視覚障害又は聴覚障害問題の一方又は両方を有する可能性がある異なるユーザにとってオーディオシーン説明の再生品質が維持されるように、シーン/ビデオ説明を再生するための再生速度又は速度設定(例えば、最大又は最小)を制御することができる。

20

30

【0084】

ある実施形態では、メディアレンダリング装置402のメモリ204に(図2に示す)訓練済み機械学習(ML)モデル212を記憶することができる。訓練済みMLモデル212は、撮影シーン404のコンテキスト情報(すなわち、コンテキストを示す情報)に基づいて、ビデオ説明情報406のオーディオ表現を再生するためのオーディオ特性を決定又は出力することができる。コンテキスト情報は、訓練済み機械学習(ML)モデル212への入力であることができる。機械学習(ML)モデル212は、入力(すなわち、コンテキスト情報)と出力(すなわち、オーディオ特性)との間の関係を識別するように訓練することができる。MLモデル212は、例えば重みの数、コスト関数、入力サイズ及び層の数などのハイパーパラメータによって定めることができる。MLモデル212のハイパーパラメータは、MLモデル212のコスト関数の大域的最小点に近づくように調整することができ、重みもそのように更新することができる。MLモデル212は、MLモデル212の訓練データセット内の特徴に基づく数エポックの訓練後に一連の入力(すなわち、コンテキスト情報)に対して予測結果(例えば、オーディオ特性)を出力するように訓練することができる。

40

【0085】

MLモデル212は、例えばソフトウェアプログラム、ソフトウェアプログラムのコード、ライブラリ、アプリケーション、スクリプト、或いは回路202などの処理装置によって実行されるその他のロジック又は命令などの電子データを含むことができる。MLモ

50

デル 2 1 2 は、メディアレンダリング装置 4 0 2 などのコンピュータ装置がコンテキスト情報に基づいてオーディオ特性を決定するための 1 又は 2 以上の動作を実行することを可能にするように構成されたコード及びルーチンを含むことができる。これに加えて又はこれに代えて、ML モデル 2 1 2 は、プロセッサ、（例えば、1 又は 2 以上の動作を実行し又は実行を制御する）マイクロプロセッサ、フィールドプログラマブルゲートアレイ（FPGA）、又は特定用途向け集積回路（ASIC）を含むハードウェアを使用して実装することもできる。或いは、いくつかの実施形態では、ハードウェアとソフトウェアとの組み合わせを使用して ML モデル 2 1 2 を実装することもできる。

【0086】

ある実施形態によれば、回路 2 0 2 は、撮影シーン 4 0 4 のコンテキスト情報を決定するように構成することができる。コンテキスト情報の例としては、以下に限定するわけではないが、アクション、格闘、冒険、アニメーション、コメディ、ダンス、ミュージカル、犯罪、叙事詩、エロティカ、ファンタジー、ホラー、ミステリー、哲学、政治、宗教、ロマンス、SF、スリラー、都市、戦争、伝記、又は悲劇を挙げることができる。コンテキスト情報は、撮影シーン 4 0 4 の少なくとも 1 つの視覚的特性の分析に基づいて決定することができる。撮影シーン 4 0 4 の視覚的特性としては、以下に限定するわけではないが、少なくとも 1 つのフレーム内で認識される物体（例えば、図 4 の車 4 1 2 ）、少なくとも 1 つのフレーム内で認識される（図 4 の第 1 の人物 4 0 8 又は第 2 の人物 4 1 0 などの）人物、少なくとも 1 つのフレーム内の少なくとも 1 つのオブジェクトの（幸福状態、悲しみ状態、怒り状態、混乱状態、ストレス状態、又は興奮状態などの）感情状態、少なくとも 1 つのフレームの背景情報、少なくとも 1 つのフレーム内の周囲照明条件、少なくとも 1 つのフレーム内の動き情報（すなわち、静止又は移動）、少なくとも 1 つのフレーム内の少なくとも 1 つのオブジェクトに関連する（ダンスジェスチャ又はアクションジェスチャなどの）ジェスチャ、又は少なくとも 1 つのフレームに関連するジャンル情報を挙げることができる。いくつかの実施形態では、回路 2 0 2 を、（撮影シーン 4 0 4 などの）撮影シーンの視覚的特徴及びコンテキスト情報を決定するために、当業で周知の様々な画像処理法、シーンマイニング法、又はシーン理解法を実装するように構成することができる。

【0087】

回路 2 0 2 は、撮影シーン 4 0 4 の決定されたコンテキストに対する訓練済み ML モデル 2 1 2 の適用に基づいて、ビデオ説明情報 4 0 6 のオーディオ表現を再生するためのオーディオ特性を決定するようにさらに構成することができる。オーディオ特性は、以下に限定するわけではないが、ラウドネスパラメータ、ピッチパラメータ、トーンパラメータ、発話速度パラメータ、声質パラメータ、音声学的パラメータ、イントネーションパラメータ、倍音の強度、音声変調パラメータ、発音パラメータ、韻律パラメータ、音色パラメータ、或いは 1 又は 2 以上の音響心理的パラメータを含むことができる。オーディオ特性は、撮影シーン 4 0 4 の決定されたコンテキスト情報に対する訓練済み ML モデル 2 1 2 の適用に基づいて決定することができる。

【0088】

回路 2 0 2 は、決定された速度及び決定されたオーディオ特性に基づいて、ビデオ説明情報 4 0 6 のオーディオ表現の再生を第 1 の時間間隔 4 2 2 において制御するようにさらに構成することができる。一例として、撮影シーン 4 0 4 のコンテキスト情報が格闘シーンとして決定された場合、回路 2 0 2 は、メディアコンテンツ及びビデオ説明情報 4 0 6（すなわち、シーン説明）をレンダリングできる対象であるユーザ 1 2 0 にリアルなユーザ体験を提供するために、ビデオ説明情報 4 0 6 の（音量などの）ラウドネスパラメータ、及び倍音パラメータ（すなわち、オーディオ特性）の強度を高めるように構成することができる。このような場合、回路 2 0 2 は、コンテキスト情報が格闘シーンとして決定されたことに基づいて、決定されたコンテキスト情報に対する訓練済み ML モデル 2 1 2 の適用に基づいて（音量などの）ラウドネスパラメータをオーディオ特性として決定することができる。

10

20

30

40

50

【 0 0 8 9 】

別の実施形態では、回路 2 0 2 を、ビデオ説明情報 4 0 6 のオーディオ表現、並びに撮影シーン 4 0 4 又は一連の撮影シーンの各撮影シーンの（第 1 のオーディオ部分 4 1 6 A 及び第 2 のオーディオ部分 4 1 6 B などの）オーディオ部分を聞き取れるように再生するようにオーディオレンダリング装置 1 0 6 を制御するようさらに構成することができる。オーディオレンダリング装置 1 0 6 は、（図 2 に示すような）メディアレンダリング装置 4 0 2 に関連することができ、又はメディアレンダリング装置 4 0 2 内に統合することができる。

【 0 0 9 0 】

ある実施形態では、撮影シーン 4 0 4 のビデオ説明情報 4 0 6 が、撮影シーン 4 0 4 内に存在する生物オブジェクト及び／又は無生物オブジェクトに関する認知情報を含むことができる。生物オブジェクトは（人間、動物又は鳥などの）生物を含むことができる。無生物オブジェクトは無生物を含むことができる。オブジェクト（生物又は無生物）に関する認知情報は、撮影シーン 4 0 4 のコンテキストに関連することも又はしないこともあるオブジェクトの徹底的な詳細を提供することができる。認知情報は、撮影シーン 4 0 4 内に存在するオブジェクトに関する一般的知識又は情報をユーザ 1 2 0 に提供することができる。いくつかの実施形態では、認知情報が、撮影シーン内に存在するオブジェクトに関連する画像又はアイコンに対応することができ、或いはオブジェクトに関連するオーディオトーンに対応することができる。ある実施形態では、回路 2 0 2 を、ディスプレイ装置 1 0 4 又はオーディオレンダリング装置 1 0 6 のいずれかによる認知情報の再生を制御するようさらに構成することができる。

【 0 0 9 1 】

図 5 は、本開示の実施形態による、シーン説明の再生制御のための例示的な動作を示す第 1 のフローチャートである。図 5 の説明は、図 1、図 2、図 3 A、図 3 B 及び図 4 の要素に関連して行う。図 5 にはフローチャート 5 0 0 を示す。5 0 2 ~ 5 0 8 の動作は、例えばメディアレンダリング装置 1 0 2 又は回路 2 0 2 などのいずれかのコンピュータ装置上で実施することができる。動作は 5 0 2 から開始して 5 0 4 に進むことができる。

【 0 0 9 2 】

5 0 4 において、メディアコンテンツを検索することができる。メディアコンテンツは、一連の撮影シーン 1 1 4 及びテキスト情報 1 1 8 を含むことができる。テキスト情報 1 1 8 は、ビデオ説明情報 1 1 8 A 及びタイミング情報 1 1 8 B を含むことができる。ビデオ説明情報 1 1 8 A は、一連の撮影シーン 1 1 4 の撮影シーンを説明することができる。1 又は 2 以上の実施形態では、回路 2 0 2 を、一連の撮影シーン 1 1 4 及びテキスト情報 1 1 8 を含むことができるメディアコンテンツ 1 1 2 を検索するように構成することができる。テキスト情報 1 1 8 は、ビデオ説明情報 1 1 8 A 及びタイミング情報 1 1 8 B をさらに含むことができる。ビデオ説明情報 1 1 8 A は、一連の撮影シーン 1 1 4 の撮影シーンを説明することができる。

【 0 0 9 3 】

5 0 6 において、撮影シーンのテキスト情報 1 1 8 から、ビデオ説明情報 1 1 8 A を再生するためのタイミング情報 1 1 8 B を抽出することができる。1 又は 2 以上の実施形態では、回路 2 0 2 を、撮影シーンのテキスト情報 1 1 8 からタイミング情報 1 1 8 B（すなわち、ビデオ説明情報 1 1 8 A を再生するためのタイミング情報）を抽出するように構成することができる。

【 0 0 9 4 】

5 0 8 において、ビデオ説明情報 1 1 8 A（テキスト表現、オーディオ表現、又はテキスト表現及びオーディオ表現の両方）の再生を制御することができる。ビデオ説明情報 1 1 8 A は、撮影シーンの抽出されたタイミング情報 1 1 8 B によって示される第 1 の時間間隔において再生することができる。1 又は 2 以上の実施形態では、回路 2 0 2 を、撮影シーンの抽出されたタイミング情報によって示される第 1 の時間間隔においてビデオ説明情報 1 1 8 A の再生（テキスト表現、オーディオ表現、又はテキスト表現及びオーディオ

表現の両方)を制御するように構成することができる。制御は終了に進むことができる。

【0095】

図6は、本開示の実施形態による、シーン説明の再生制御のための例示的な動作を示す第2のフローチャートである。図6の説明は、図1、図2、図3A、図3B、図4及び図5の要素に関連して行う。図6にはフローチャート600を示す。602~610の動作は、例えばメディアレンダリング装置102又は回路202などのいずれかのコンピュータ装置上で実施することができる。動作は602から開始して604に進むことができる。

【0096】

604において、第1の撮影シーン114Aの一連の第2の時間間隔を決定することができる。一連の第2の時間間隔の各々は、一連の撮影シーン114における撮影シーンのオーディオ部分116を再生するための時間間隔を示すことができる。1又は2以上の実施形態では、回路202を、一連の撮影シーンにおける撮影シーンのオーディオ部分116を再生するための時間間隔をそれぞれが示すことができる、撮影シーンの一連の第2の時間間隔を決定するように構成することができる。一連の第2の時間間隔の決定の詳細については、例えば図4で説明している。

【0097】

606において、撮影シーンのビデオ説明情報118Aのオーディオ表現の第3の時間間隔を決定することができる。1又は2以上の実施形態では、回路202を、撮影シーンのビデオ説明情報118Aのオーディオ表現の第3の時間間隔を決定するように構成することができる。第3の時間間隔の決定の詳細については、例えば図4で説明している。

【0098】

608において、決定された一連の第2の時間間隔及び決定された第3の時間間隔に基づいて、ビデオ説明情報118Aのオーディオ表現を再生する速度を決定することができる。1又は2以上の実施形態では、回路202を、決定された一連の第2の時間間隔及び決定された第3の時間間隔に基づいて、ビデオ説明情報118A(すなわち、シーン説明)のオーディオ表現を再生する速度を決定するように構成することができる。ビデオ説明情報の再生速度の決定に関する詳細については、例えば図4で説明している。

【0099】

610において、決定された速度に基づいてビデオ説明情報118Aのオーディオ表現の再生を制御することができる。ビデオ説明情報118Aのオーディオ表現は、一連の第2の時間間隔とは異なることができる第1の時間間隔において再生することができる。1又は2以上の実施形態では、回路202を、決定された速度に基づいてビデオ説明情報118Aのオーディオ表現の再生を第1の時間間隔において制御するように構成することができる。ビデオ説明情報118Aのオーディオ表現の再生を制御する詳細については、例えば図4で説明している。制御は終了に進むことができる。

【0100】

本開示の様々な実施形態は、メディアレンダリング装置402などの機械及び/又はコンピュータが実行できる命令を記憶した非一時的コンピュータ可読媒体及び/又は記憶媒体を提供することができる。これらの命令は、一連の撮影シーンを含むことができるメディアコンテンツの検索を含むことができる動作を機械及び/又はコンピュータに実行させることができる。メディアコンテンツは、ビデオ説明情報及びタイミング情報を含むテキスト情報を含むことができる。ビデオ説明情報は、一連の撮影シーン内の撮影シーンを説明することができる。動作は、撮影シーンのテキスト情報から、ビデオ説明情報を再生するためのタイミング情報を抽出することをさらに含むことができる。動作は、ビデオ説明情報の再生を、抽出された撮影シーンのタイミング情報によって示される第1の時間間隔においてテキスト表現又はテキスト表現及びオーディオ表現のいずれかで制御することをさらに含むことができる。

【0101】

他のいくつかの実施形態では、動作が、撮影シーンの一連の第2の時間間隔を決定することを含むことができる。一連の第2の時間間隔の各々は、一連の撮影シーンにおける撮

10

20

30

40

50

影シーンのオーディオ部分を再生するための時間間隔を示すことができる。動作は、撮影シーンのビデオ説明情報のオーディオ表現の第3の時間間隔を決定することをさらに含むことができる。動作は、決定された一連の第2の時間間隔及び決定された第3の時間間隔に基づいてビデオ説明情報のオーディオ表現を再生する速度を決定することをさらに含むことができる。動作は、決定された速度に基づいてビデオ説明情報のオーディオ表現の再生を第1の時間間隔において制御することをさらに含むことができる。第1の時間間隔は、一連の第2の時間間隔とは異なることができる。

【0102】

本開示の例示的な態様は、(回路202などの)回路を含むことができる(図1のメディアレンダリング装置102などの)メディアレンダリング装置を含むことができる。回路は、(一連の撮影シーン114などの)一連の撮影シーン、(オーディオ部分116などの)オーディオ部分及び(テキスト情報118などの)テキスト情報を含むことができるメディアコンテンツを検索するように構成することができる。テキスト情報は、(ビデオ説明情報118Aなどの)テキストベースのビデオ説明情報及び(タイミング情報118Bなどの)タイミング情報を含むことができる。ビデオ説明情報118Aは、一連の撮影シーンにおける(撮影シーン304などの)撮影シーンを説明することができる。メディアコンテンツは、一連の撮影シーンの各々のオーディオ部分を表すことができるクローズドキャプション情報をさらに含むことができる。一連の撮影シーンの各々を説明するビデオ説明情報は、クローズドキャプション情報と共にメディアコンテンツ内に符号化することができる。ある実施形態では、回路を、撮影シーンのテキスト情報をビデオ説明情報のオーディオ表現に変換するようにさらに構成することができる。

【0103】

ある実施形態では、回路を、撮影シーンのテキスト情報から、ビデオ説明情報を再生するためのタイミング情報を抽出するようにさらに構成することができる。回路は、ビデオ説明情報の再生を、抽出された撮影シーンのタイミング情報によって示される第1の時間間隔においてテキスト表現又はテキスト表現及びオーディオ表現のいずれかで制御するようにさらに構成することができる。

【0104】

別の実施形態では、回路を、撮影シーンのテキスト情報から、ビデオ説明情報を再生するための速度情報を抽出するようにさらに構成することができる。テキスト情報は、速度情報をさらに含むことができる。回路は、抽出された速度情報に基づいて、抽出された撮影シーンのタイミング情報によって示される第1の時間間隔においてビデオ説明情報のオーディオ表現の再生を制御するようにさらに構成することができる。

【0105】

いくつかの実施形態では、回路を、撮影シーンの(一連の第2の時間間隔416などの)一連の第2の時間間隔を決定するように構成することができる。一連の第2の時間間隔の各々は、一連の撮影シーンにおける撮影シーンのオーディオ部分を再生するための時間間隔を示すことができる。回路は、撮影シーンのビデオ説明情報のオーディオ表現の(第3の時間間隔418などの)第3の時間間隔を決定するようにさらに構成することができる。回路は、ビデオ説明情報のオーディオ表現を再生する速度を決定するようにさらに構成することができる。ビデオ説明情報のオーディオ表現を再生する速度は、決定された一連の第2の時間間隔及び決定された第3の時間間隔に基づいて決定することができる。ある実施形態では、決定される速度が、変換されたオーディオ表現の実際の再生速度よりも低いことができる。別の実施形態では、決定される速度が、変換されたオーディオ表現の実際の再生速度よりも高いことができる。

【0106】

いくつかの実施形態では、回路を、ビデオ説明情報のオーディオ表現の再生を(第1の時間間隔422などの)第1の時間間隔において制御するように構成することができる。ビデオ説明情報のオーディオ表現の再生は、決定された速度に基づいて制御することができる。ある実施形態では、回路を、撮影シーンの(一連の第4の時間間隔420A~42

10

20

30

40

50

0 C などの) 一連の第 4 の時間間隔を決定するように構成することができる。一連の第 4 の時間間隔の各々は、一連の第 2 の時間間隔とは異なることができる。回路は、一連の第 4 の時間間隔から、時間間隔閾値よりも高いことができる第 1 の時間間隔を選択するように構成することができる。第 1 の時間間隔は、一連の第 2 の時間間隔とは異なることができる。ある実施形態では、第 1 の時間間隔が、撮影シーンの第 1 のせりふと第 2 のせりふとの間であることができる。第 1 のせりふは、撮影シーンの第 1 のショットの最後の単語であることができ、第 2 のせりふは、撮影シーンの第 2 のショットの最初の単語であることができる。第 1 のショット及び第 2 のショットは、撮影シーンの連続するショットであることができる。別の実施形態では、第 1 の時間間隔が、撮影シーンの開始と撮影シーンの第 1 のせりふとの間であることができる。

10

【0107】

いくつかの実施形態では、回路を、メディアレンダリング装置に関連する規定の速度設定に基づいて、ビデオ説明情報のオーディオ表現を再生する速度を決定するように構成することができる。規定の速度設定は、ビデオ説明情報のオーディオ表現の最大再生速度を示すことができる。回路は、テキスト情報と共に速度情報を受け取り、決定された速度及び規定の速度設定に基づいて撮影シーンの画像部分又はオーディオ部分の一方の再生を制御するようにさらに構成することができる。いくつかの実施形態では、撮影シーンを説明するビデオ説明情報が、撮影シーン内に存在する生物オブジェクト又は無生物オブジェクトに関する認知情報を含むことができる。回路は、撮影シーンのビデオ説明情報に含まれる認知情報の再生を制御するように構成することができる。

20

【0108】

ある実施形態では、回路を、一連の撮影シーンのうちの 1 つの撮影シーンの説明に対応できる第 1 のユーザ入力を受け取るように構成することができる。回路は、受け取られた第 1 のユーザ入力を、一連の撮影シーンの各々に関連するビデオ説明情報内で検索するようにさらに構成することができる。回路は、検索に基づいて、メディアコンテンツを再生するための再生タイミング情報を決定するようにさらに構成することができる。回路は、決定された再生タイミング情報に基づいてメディアコンテンツの再生を制御するようにさらに構成することができる。

【0109】

別の実施形態では、回路を、メディアコンテンツをレンダリングできる対象であるユーザのプロファイル情報を示すことができる第 2 のユーザ入力を受け取るように構成することができる。回路は、受け取られた第 2 のユーザ入力に基づいて、ビデオ説明情報のオーディオ表現を再生する速度設定を決定するように構成することができる。

30

【0110】

いくつかの実施形態では、メディアレンダリング装置に関連する(メモリ 204 などの)メモリを、(訓練済み機械学習(ML)モデル 212 などの)訓練済み ML モデルを記憶するように構成することができる。回路は、撮影シーンの少なくとも 1 つの特性の分析に基づいて撮影シーンのコンテキスト情報を決定するように構成することができる。回路は、撮影シーンの決定されたコンテキスト情報に対する訓練済み ML モデルの適用に基づいて、ビデオ説明情報のオーディオ表現を再生するためのオーディオ特性を決定するようにさらに構成することができる。回路は、決定された速度及び決定されたオーディオ特性に基づいて、ビデオ説明情報のオーディオ表現の再生を第 1 の時間間隔において制御するようにさらに構成することができる。

40

【0111】

ある実施形態では、メディアレンダリング装置が、ビデオ説明情報のテキスト表現を再生する(又は表示する)ように構成されたディスプレイ装置を含むことができる。別の実施形態では、ビデオ説明情報のオーディオ表現の再生に加えてテキスト表現を表示することができる。

【0112】

別の実施形態では、回路を、オーディオレンダリング装置を制御するようにさらに構成

50

することができる。オーディオレンダリング装置は、メディアレンダリング装置に関連することができる。オーディオレンダリング装置は、ビデオ説明情報のオーディオ表現及び撮影シーンのオーディオ部分を再生するように制御することができる。

【 0 1 1 3 】

本開示は、ハードウェアで実現することも、又はハードウェアとソフトウェアとの組み合わせで実現することもできる。本開示は、少なくとも1つのコンピュータシステム内で集中方式で実現することも、又は異なる要素を複数の相互接続されたコンピュータシステムにわたって分散できる分散方式で実現することもできる。本明細書で説明した方法を実行するように適合されたコンピュータシステム又はその他の装置が適することができる。ハードウェアとソフトウェアとの組み合わせは、ロードされて実行された時に本明細書で説明した方法を実行するようにコンピュータシステムを制御することができるコンピュータプログラムを含む汎用コンピュータシステムとすることができる。本開示は、他の機能も実行する集積回路の一部を含むハードウェアで実現することができる。

【 0 1 1 4 】

本開示は、本明細書で説明した方法の実装を可能にする全ての特徴を含み、コンピュータシステムにロードされた時にこれらの方法を実行できるコンピュータプログラム製品に組み込むこともできる。本文脈におけるコンピュータプログラムとは、情報処理能力を有するシステムに特定の機能を直接的に、或いは a) 別の言語、コード又は表記法への変換、b) 異なる内容形態での複製、のいずれか又は両方を行った後に実行させるように意図された命令セットの、あらゆる言語、コード又は表記法におけるあらゆる表現を意味する。

【 0 1 1 5 】

いくつかの実施形態を参照しながら本開示を説明したが、当業者であれば、本開示の範囲から逸脱することなく様々な変更を行うことができ、同等物を代用することもできると理解するであろう。また、本開示の範囲から逸脱することなく、特定の状況又は内容を本開示の教示に適合させるように多くの修正を行うこともできる。従って、本開示は、開示した特定の実施形態に限定されるものではなく、添付の特許請求の範囲内に収まる全ての実施形態を含むように意図される。

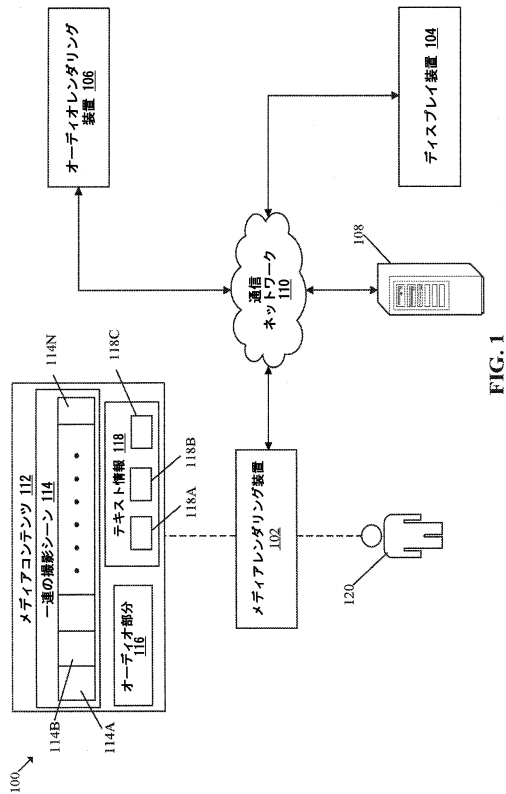
【符号の説明】

【 0 1 1 6 】

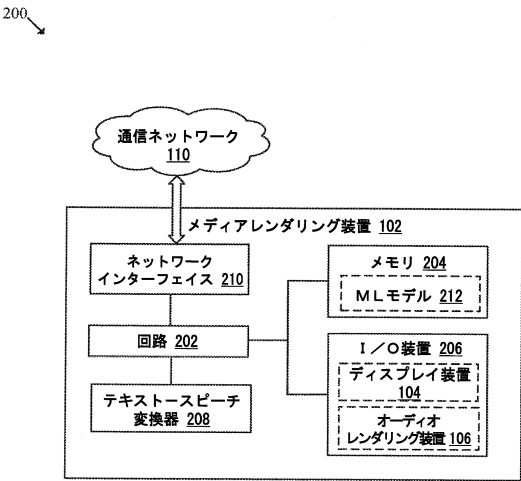
- | | |
|---------|---------------|
| 1 0 0 | ネットワーク環境 |
| 1 0 2 | メディアレンダリング装置 |
| 1 0 4 | ディスプレイ装置 |
| 1 0 6 | オーディオレンダリング装置 |
| 1 0 8 | サーバ |
| 1 1 0 | 通信ネットワーク |
| 1 1 2 | メディアコンテンツ |
| 1 1 4 | 一連の撮影シーン |
| 1 1 4 A | 第 1 の撮影シーン |
| 1 1 4 B | 第 2 の撮影シーン |
| 1 1 4 N | 第 N の撮影シーン |
| 1 1 6 | オーディオ部分 |
| 1 1 8 | テキスト情報 |
| 1 1 8 A | ビデオ説明情報 |
| 1 1 8 B | タイミング情報 |
| 1 1 8 C | 速度情報 |
| 1 2 0 | ユーザ |

【図面】

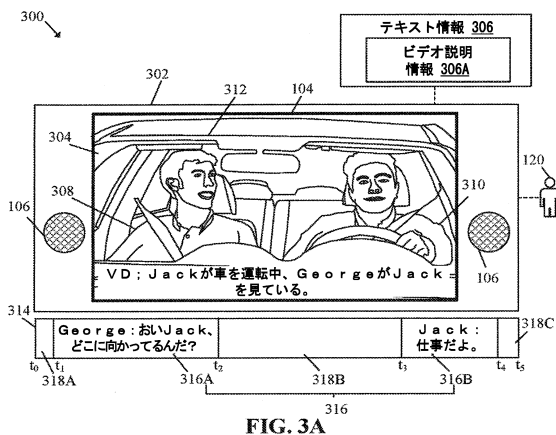
【図 1】



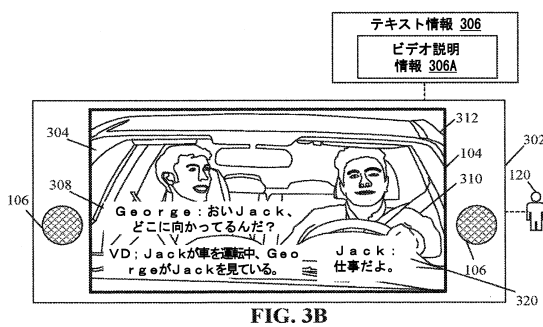
【図 2】



【図 3 A】



【図 3 B】



10

20

30

40

50

【図 4】

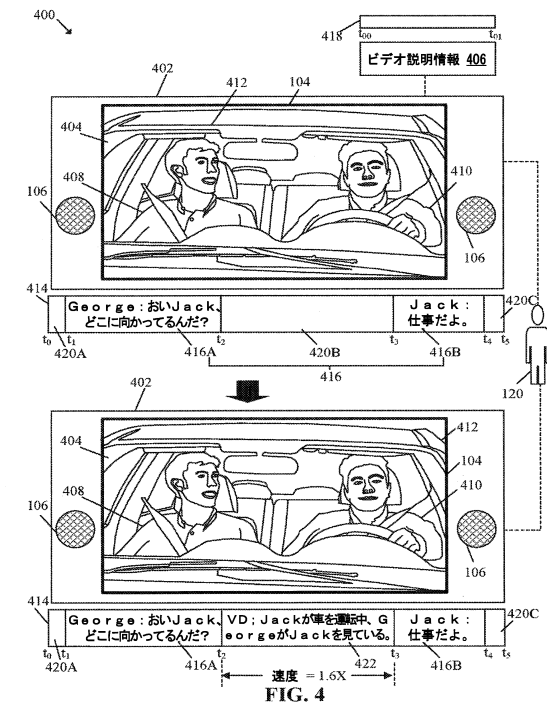


FIG. 4

【図 5】

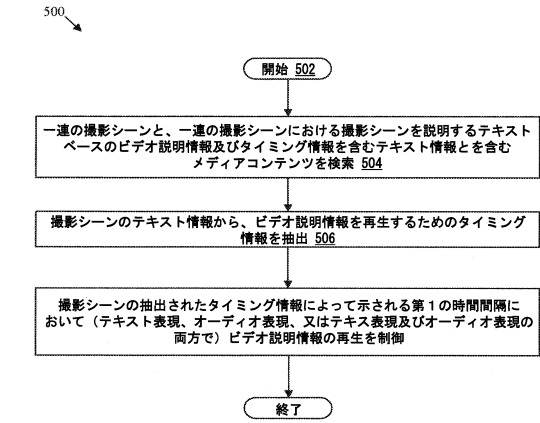


FIG. 5

【図 6】

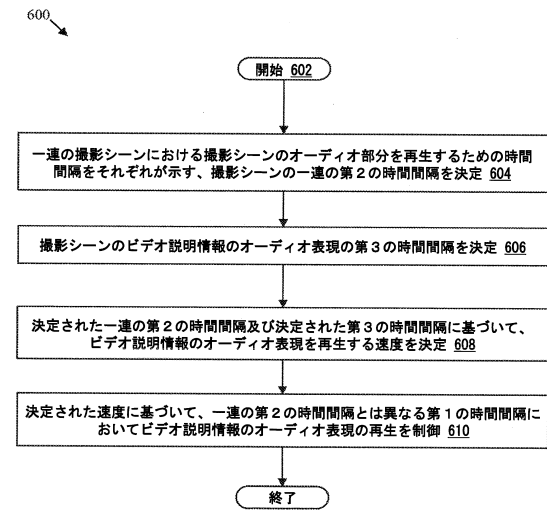


FIG. 6

フロントページの続き

- 1 2 7 サンディエゴ ヴィア エスプリロ 1 6 5 3 5 ソニー コーポレイション オブ アメリカ内
(72)発明者 ネジャット マヤル マイク
アメリカ合衆国 カリフォルニア州 9 2 1 2 7 サンディエゴ ヴィア エスプリロ 1 6 5 3 5 ソ
ニー コーポレイション オブ アメリカ内
(72)発明者 シンタニ ピーター
アメリカ合衆国 カリフォルニア州 9 2 1 2 7 サンディエゴ ヴィア エスプリロ 1 6 5 3 5 ソ
ニー コーポレイション オブ アメリカ内
(72)発明者 ブランチャード ロバート
アメリカ合衆国 カリフォルニア州 9 2 1 2 7 サンディエゴ ヴィア エスプリロ 1 6 5 3 5 ソ
ニー コーポレイション オブ アメリカ内
審査官 大西 宏
(56)参考文献 特開 2 0 0 0 - 2 5 0 5 7 5 (J P , A)
特開 2 0 0 3 - 1 4 3 5 7 5 (J P , A)
特開 2 0 0 3 - 2 5 9 3 2 0 (J P , A)
特開 2 0 0 4 - 0 6 2 7 6 9 (J P , A)
特表 2 0 1 7 - 5 3 1 9 3 6 (J P , A)
国際公開第 2 0 1 8 / 2 1 1 7 4 8 (WO , A 1)
国際公開第 2 0 2 0 / 1 1 2 8 0 8 (WO , A 1)
Ruxandra Tapu et al. , DEEP-HEAR: A Multimodal Subtitle Positioning System Dedicated to
Deaf and Hearing-Impaired People , IEEE Access , 米国 , IEEE , 2019年07月01日 , Volum
e: 7 , pp.88150-88162 , <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=875>
1956 , IEL Online
(58)調査した分野 (Int.Cl. , D B 名)
H 0 4 N 2 1 / 0 0 - 2 1 / 8 5 8