

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5351958号

(P5351958)

(45) 発行日 平成25年11月27日 (2013.11.27)

(24) 登録日 平成25年8月30日 (2013.8.30)

(51) Int. Cl.

F I

G 0 6 T 1/00 (2006.01)

G 0 6 T 1/00 2 0 0 D

G 0 6 T 7/00 (2006.01)

G 0 6 T 7/00 3 0 0 F

G 0 6 F 17/30 (2006.01)

G 0 6 F 17/30 1 7 0 B

H 0 4 N 5/76 (2006.01)

G 0 6 F 17/30 2 1 0 A

H 0 4 N 5/76 B

請求項の数 2 (全 17 頁)

(21) 出願番号 特願2011-512451 (P2011-512451)  
 (86) (22) 出願日 平成21年5月22日 (2009.5.22)  
 (65) 公表番号 特表2011-525012 (P2011-525012A)  
 (43) 公表日 平成23年9月8日 (2011.9.8)  
 (86) 国際出願番号 PCT/US2009/003160  
 (87) 国際公開番号 W02009/148518  
 (87) 国際公開日 平成21年12月10日 (2009.12.10)  
 審査請求日 平成24年3月22日 (2012.3.22)  
 (31) 優先権主張番号 61/058, 201  
 (32) 優先日 平成20年6月2日 (2008.6.2)  
 (33) 優先権主張国 米国 (US)  
 (31) 優先権主張番号 12/331, 927  
 (32) 優先日 平成20年12月10日 (2008.12.10)  
 (33) 優先権主張国 米国 (US)

(73) 特許権者 513077243  
 インテレクチュアル ベンチャーズ ファ  
 ンド 83 エルエルシー  
 アメリカ合衆国、89128 ネバダ州、  
 ラスベガス、ウエスト レイク ミード  
 ブールバード 7251、スイート 30  
 O  
 (74) 代理人 100107766  
 弁理士 伊東 忠重  
 (74) 代理人 100070150  
 弁理士 伊東 忠彦  
 (74) 代理人 100091214  
 弁理士 大貫 進介

最終頁に続く

(54) 【発明の名称】 デジタルコンテンツ記録のための意味論的イベント検出

(57) 【特許請求の範囲】

【請求項 1】

イベントに関連する画像記録のグループの意味論的イベント分類を容易にする方法であって、

前記画像記録の各々から複数の視覚的特徴を抽出するステップと、

前記視覚的特徴を使用して前記画像記録の各々に対する複数の概念スコアを生成するステップと、

前記画像記録の前記概念スコアに基づいて、各イベントを記述するための B O F 特徴ベクトルを、意味論的イベントに対応する予め定められたコードブックに前記イベントの前記画像記録の前記概念スコアをマッピングすることにより生成するステップと、

前記イベントに意味論的イベントが現れる確率の指標である検出スコアを生成する意味論的イベント分類器に前記マッピングされた特徴ベクトルを供給するステップと、  
 を包含し、

各前記概念スコアは、視覚的概念に対応し、前記画像記録が前記視覚的概念を含む確率の指標であることを特徴とする方法。

【請求項 2】

請求項 1 に記載のイベントに関連する画像記録のグループの意味論的イベント分類を容易にする方法であって、

前記画像記録の対の間のペアワイズ類似性を決定するステップと、

スペクトルクラスタ化を適用して、前記決定されたペアワイズ類似性に基づいて、前記

10

20

意味論的イベントの訓練画像記録を各クラスが一つのコードワードに対応する異なるクラスにグループ化することによって各前記意味論的イベントの前記コードブックを生成するステップと、

前記訓練イベントの前記画像記録の前記概念スコアを意味論的イベントに対応する前記コードブックにマッピングして、各前記訓練イベントを記述するためのB O F特徴ベクトルを生成するステップと、

前記イベント分類器を前記訓練イベントに対応する前記B O F特徴ベクトルに基づいて訓練するステップと、

を包含する訓練プロセスを有することを特徴とする方法。

【発明の詳細な説明】

10

【技術分野】

【0001】

本発明は、デジタル静止画像又はビデオのようなデジタルコンテンツ記録のカテゴリ化に関する。特に、本発明は、意味論的イベント(semantic events)の検出に基づいたデジタルコンテンツ記録のカテゴリ化に関する。

【背景技術】

【0002】

低コストの電子消費者撮像技術の出現は、平均的な消費者によって獲得されるデジタル画像の数の顕著な増加をもたらす結果となっている。実際、様々な形態の電子メモリが時間とともにますます安価になっているので、消費者は、より一層多くのデジタル静止画像及びビデオを撮影するとともに、以前には廃棄したであろうデジタル静止画像及びビデオも保持する傾向にある。結果として、平均的な消費者は、記憶及び後の検索のためにデジタル画像を適切に識別及びカタログ化するにあたって、ますます困難な問題に直面している。一般的に、そのような識別及びカタログ化は通常は手作業で実行され、これは消費者にとって極端に時間を消費するプロセスになることがある。

20

【0003】

単なる一つの描写として、消費者は1回の休暇の間にいくつもの異なる場所に旅行するかもしれない。消費者は、特定の場所の各々で、ならびに他の主題カテゴリ又はイベントに関係している場所の各々で、画像を撮影し得る。例えば、消費者は、それらの場所の各々で家族メンバの画像を撮影し、それらの場所の各々で特定のイベントの画像を撮影し、それらの場所の各々で歴史的な建造物の画像を撮影し得る。旅行から戻ると、消費者は、デジタル画像を人物、誕生日、博物館などの様々なグループ分けに基づいて分類し、デジタル画像をそのグループ分けに基づいて電子アルバムに記憶したいと思うかもしれない。消費者は、現在のところ、何百というデジタル静止画像及びビデオセグメントを特定のイベントで識別するために、それらを手作業で分類するということに直面している。

30

【0004】

上記のことを考慮して、最近、消費者の写真及びビデオの自動アルバム化が大きな関心を集めている。自動アルバム化に対する一つの人気のあるアプローチは、デジタル画像及びビデオを日付順及び画像コンテンツ内の視覚的な類似性によるイベントに従って組織化することである。例えば、非特許文献1には、デジタル画像のグループがどのようにして自動的にイベントにクラスタ化されることができかが記載されている。

40

【先行技術文献】

【非特許文献】

【0005】

【非特許文献1】A.C.Loui及びA.Savakis,「Automated event clustering and quality screening of consumer pictures for digital albuming(デジタルアルバム化のための消費者写真の自動イベントクラスタ化及び質のスクリーニング)」, IEEE Trans. on Multimedia, 2003年, Vol.5, No.3, p.390 - 402

【発明の概要】

【発明が解決しようとする課題】

50

## 【 0 0 0 6 】

画像の基本的なクラスタ化は単一のイベントに関連しているように見える画像をグループ化することができるが、自動アルバム化プロセスを改善するために、クラスタ化されたイベントに意味論的意味 ( semantic meanings ) のタグ付けをすることができることが望ましいであろう。しかし、意味論的イベントの検出は、基本的な問題を提示する。第 1 に、実用的なシステムは、デジタル静止画像及びビデオを同時に処理することができる必要がある。これは、しばしば両方が、実際の消費者画像コレクションに存在するからである。第 2 に、実用的なシステムは実際の消費者コレクション内の様々な意味論的コンテンツを収容し、それによって、各々の特定の意味論的イベントを検出する特定の個別の方法の代わりに、異なる意味論的イベントを検出する包括的な方法を組み込んだシステムを提供することを望ましくする必要がある。最後に、実用的なシステムは、識別及び分類における誤りを防ぐために、ロバストである必要がある。

10

## 【課題を解決するための手段】

## 【 0 0 0 7 】

本発明は、デジタル画像コンテンツ記録における意味論的イベント検出のためのシステム及び方法を提供する。特に、イベントレベルの「 Bag - o f - F e a t u r e s ( 特徴のバッグ ) 」 ( B O F ) 表現がイベントをモデル化するために使用され、包括的な意味論的イベントが、 B O F 表現に基づいて、元の低レベルの視覚的特徴空間の代わりに概念空間で検出される。

## 【 0 0 0 8 】

20

好適な実施形態では、イベントレベル表現が開発され、そこでは各イベントが B O F 特徴ベクトルによってモデル化され、 B O F 特徴ベクトルに基づいて意味論的イベント検出器が直接的に構築される。分類器の訓練のために画像レベル特徴ベクトルが使用される単純なアプローチに比べて、本発明は、イベント内の難しい画像又は誤って組織化された画像に対して、よりロバストである。例えば、任意の所与のイベントにおいて、いくつかの画像は分類が難しいことがある。これらの困難な画像は、通常、決定境界を複雑にし、モデル化を困難にする。イベントレベル特徴表現を適用することによって、イベントレベルの測定における困難な又は誤ったデジタル静止画像及びビデオセグメントの影響を減らすことによって、感度の問題を避けることができる。後述のように、良好な検出性能が、サポートベクトルマシン ( Support Vector Machine , SVM ) 分類器に対する少数のサポートベクトルで達成され得る。すなわち、分類の問題が、イベントレベル表現によって顕著に単純化され得る。

30

## 【 0 0 0 9 】

好適な実施形態では、あるイベントに関連した画像記録のグループの意味論的イベント分類を容易にする方法が提供され、その方法は、画像記録の各々から複数の視覚的特徴を抽出するステップと、それらの視覚的特徴を使用して画像記録の各々に対して複数の概念スコアを生成するステップであって、各概念スコアが視覚的概念に対応し、且つ各概念スコアは画像記録がその視覚的概念を含む確率を示す、ステップと、画像記録の概念スコアに基づいてそのイベントに対応する特徴ベクトルを生成するステップと、イベント分類器に特徴ベクトルを供給するステップであって、当該イベント分類器は、当該イベントに対応する少なくとも一つの意味論的イベント分類器を特定する、ステップと、を含む。

40

## 【 0 0 1 0 】

画像記録は、少なくとも一つのデジタル静止画像及び少なくとも一つのビデオセグメントを含み得る。したがって、このシステムは、通常はデジタル静止画像及びビデオセグメントの両方を含む実生活の消費者画像データセットを取り扱うことができる。

## 【 0 0 1 1 】

複数の視覚的特徴の抽出は、ビデオセグメントからのキーフレームの抽出、ならびにキーフレーム及びデジタル静止画像の両方からの複数の視覚的特徴の抽出を含む。それから初期概念スコアが、抽出された視覚的特徴の各々に対応する各キーフレーム及び各デジタル静止画像に対して生成される。それから好ましくは、アンサンブル概念スコアが、初期

50

概念スコアに基づいて各キーフレーム及び各デジタル静止画像に対して生成される。

【 0 0 1 2 】

アンサンブル概念スコアは好ましくは、所与のキーフレーム又は所与のデジタル静止画像に対する各々の抽出された視覚的特徴に対する個別の概念スコアを融合することによって生成される。

【 0 0 1 3 】

意味論的イベント分類器がひとたび特定されると、デジタル静止画像及びビデオセグメントは、画像及びビデオセグメントの適切な分類、記憶、及び検索を容易にするためにタグ付けされることができる。

【図面の簡単な説明】

10

【 0 0 1 4 】

【図 1】本発明に従った意味論的イベント検出システムの模式的ブロック図である。

【図 2】図 1 に描かれた意味論的イベント検出システムによって利用される処理モジュールを描いた流れ図である。

【図 3】意味論的イベント検出のために図 1 に描かれたシステムを訓練するために利用される処理モジュールを描いた流れ図である。

【図 4】図 1 に描かれたシステムで使用される概念検出器を訓練するために利用される処理モジュールを描いた流れ図である。

【図 5】テストプロセスで使用される異なる意味論的イベントを、それらの詳細な定義を含めて描いた表である。

20

【図 6】図 1 に描かれたシステムの結果と、B O F 特徴ベクトルが元の低レベルの視覚的特徴に基づいて構築される従来のアプローチとの比較を描いたグラフである。

【図 7】本発明の結果を、ベースラインイベント検出器の結果、及び画像レベル概念スコア表現を直接的に使用する S V M 検出器 ( S V M - D i r e c t ) の結果と比較するグラフである。

【図 8】S V M - D i r e c t 方法と比較して本発明によって要求されるサポートベクトルの数を比較するグラフである。

【発明を実施するための形態】

【 0 0 1 5 】

本発明は、ある好適な実施形態及び添付の図面を参照して記述される。

30

【 0 0 1 6 】

複雑な意味論的イベントは通常は初歩的な視覚的概念の同時発生によって生成される。例えば、「結婚」は、「人々」「花」「公園」などのような視覚的概念に関連した意味論的イベントであり、あるパターンで進展する。一般的に、視覚的概念は、画像の画像コンテンツ特性として定義されることができ、通常は、特定のイベントを識別するために使用されるワード ( w o r d ) よりも広いワードによって意味論的に表現される。したがって、視覚的概念は、特定のイベントに与えられることができる画像コンテンツ特性のサブセットを形成する。

【 0 0 1 7 】

本発明では、初歩的な視覚的概念が最初に画像から検出され、意味論的イベント検出器が、元の低レベルの特徴空間の代わりに概念空間に構築される。そのようなアプローチからの恩恵は、少なくとも 2 つの局面を含む。第 1 に、視覚的概念は、元の低レベルの特徴よりも高レベルであり、且つより直観的である。S.Ebadollahi らの IEEE ICME ( 2006 ) 「 V i s u a l e v e n t d e t e c t i o n u s i n g m u l t i d i m e n s i o n a l c o n c e p t d y n a m i c s ( 多次元概念ダイナミクスを使用する視覚的イベント検出 ) 」に記述されているように、概念スコアは、意味論的イベントをモデル化するために強力である。第 2 に、本発明における概念空間は、好ましくは、例えば S.F.Chang らの ACM MIR ( 2007 ) 「 M u l t i m o d a l s e m a n t i c c o n c e p t d e t e c t i o n f o r c o n s u m e r v i d e o b e n c h m a r k ( 消費者ビデオベンチマークのための多モード意味論的概念検出 ) 」に記述されているように、意味論的概念検出器によって形成され、例えば A.C.Loui らの ACM MIR ( 2007 ) 「 K o d a k c o n s u m e r v i d e o b e n c h m a r k d a t a s e t : c o n c e p t d e

40

50

definition and annotation (コダック消費者ビデオベンチマークデータセット：概念の定義と注釈付け)」に記述されたタイプの既知の消費者電子画像データセットで訓練される。これらの意味論的概念検出器は、以前の画像データセットから付加的な情報を組み込む重要な役割を果たし、現在の画像データセットにおける意味論的イベントを検出する手助けをする。

【 0 0 1 8 】

例えば、上述のデータセットが実際の消費者からのデジタル静止画像及びビデオセグメントを含むと仮定すると、データセット全体が最初にマクロイベントのセットに区分され、各マクロイベントがさらにイベントのセットに区分されることが望ましい。この区分は、好ましくは、上述の以前に開発されたイベントクラスタ化アルゴリズムを使用することによって、ビデオセグメントの各デジタル静止画像の撮影時間及びそれらの間の色の類似性に基づく。例えば、 $E_t$  が  $t$  番目のイベントを指し、 $m_p^t$  個の写真及び  $m_v^t$  個のビデオを含むとする。 $I_i^t$  及び  $V_j^t$  は、 $E_t$  における  $i$  番目の写真及び  $j$  番目のビデオを指す。画像はこのアルゴリズムを使用してイベントにグループ化又はクラスタ化されることができ、イベント自身は、意味論的な意味で識別されたり関連付けられたりしない。したがって、本発明のゴールは、特定の意味論的意味、すなわち、特定のイベント  $E_t$  及びそのイベントに対応する画像記録に対して、「結婚」及び「誕生日」のような意味論的イベント  $S_E$  のタグ付けをすることである。

【 0 0 1 9 】

「人々」「公園」及び「花」のような、同時に発生する視覚的概念によって意味論的イベントが生成されることが仮定されるであろう。ここで  $C_1, \dots, C_N$  は  $N$  個の視覚的概念を示す。上述の視覚的概念検出器を使用して、 $21$  個 ( $N = 21$ ) の SVM ベースの概念検出器が、適用されたデータセットに対して、好ましくは低レベルの色、テクスチャ、及びエッジの視覚的特徴を使用して形成される。これらの意味論的概念検出器は、各画像  $I_i^t$  に対する  $21$  個の個別の概念スコア  $p(C_1, I_i^t), \dots, p(C_N, I_i^t)$  を生成するために適用されることができる。これらの概念スコアはそれから、以下により詳細に記述されるように、概念空間における画像  $I_i^t$  を

【数 1】

$$f(I_i^t) = [p(C_1, I_i^t), \dots, p(C_N, I_i^t)]^T$$

と表現するための特徴ベクトルを形成するために使用される。

【 0 0 2 0 】

実際の消費者からのビデオセグメントは、通常は一つの長い撮影映像（ロングショット）からの様々な視覚的コンテンツを有するので、各ビデオ  $V_j^t$  は、好ましくはセグメント  $V_{j,1}^t, \dots, V_{j,m_j}^t$  のセットに区分される。各セグメントは所与の長さ（例えば 5 秒）を有する。それからキーフレームが、ビデオセグメントから一様に周期的にサンプリングされる（例えば 0.5 秒ごとに）。例えば、 $I_{j,k,l}^t$  が  $k$  番目のセグメント  $V_{j,k}^t$  の  $l$  番目のキーフレームであるとする、そのときには、 $I_{j,k,l}^t$  はまた、デジタル静止画像と同じように概念空間内の特徴ベクトル

【数 2】

$$f(I_{j,k,l}^t)$$

によっても表現されることができる。上述されたものとは異なるサンプリングレートが容易に使用され得ることが理解されるであろう。

【 0 0 2 1 】

デジタル静止画像及びビデオセグメントの両方が、 $x$  によって表されるデータポイントとして定義されることができる。例えば、イベント  $E_t$  は合計で、

【数 3】

$$m' = m'_p + \tilde{m}'_v$$

のデータ点を含み、

【数 4】

$$\tilde{m}'_v$$

は  $E_i$  における  $m'_v$  個のビデオクリップからのビデオセグメントの全数である。意味論的イベント検出器がそれから、これらのデータ点、及び概念スコアから開発された対応する特徴ベクトルに基づいて、実行される。

10

【0022】

B O F 表現は、画像に対する包括的概念を検出するために有効であることが証明されてきている。例えば、J.Sivic及びA.Zisserman, 「Video google: a text retrieval approach to object matching in videos (ビデオグーグル: ビデオにおけるオブジェクトマッチングに対するテキスト検索アプローチ)」, ICCV, pp.1470-1477 (2003) を参照のこと。B O F では、画像は、順序なしのローカル記述子のセットによって表現される。クラスタ化技術を通して、中レベルの視覚的語彙が構築され、そこでは各々の視覚的ワードがローカル記述子のグループによって形成される。各々の視覚的ワードは、画像を記述するためのロバストで且つノイズが除去された視覚的用語であるとみなされる。

20

【0023】

例えば、 $S_E$  が意味論的イベント、例えば「結婚」を指し、 $E_1, \dots, E_M$  がこの意味論的イベントを含む  $M$  個のイベントを指すとする。各  $E_i$  は、 $m'_p$  個の写真及び

【数 5】

$$\tilde{m}'_v$$

個のビデオセグメントによって形成されている。視覚的語彙と同様に、概念語彙が、これら

【数 6】

$$\sum_{i=1}^M m'_i$$

30

のデータ点（ここで、

【数 7】

$$m' = m'_p + \tilde{m}'_v$$

) を  $n$  個の概念ワードにクラスタ化することによって構築されることができる。各概念ワードは、 $S_E$  を含む全てのイベントを記述するための共通の特性である概念同時発生のパターンとして取り扱われることができる。特に、静止ビデオ画像及びビデオデータ点の両方を収納するために、スペクトルクラスタ化アルゴリズム（例えばA.Y.Ng, M.Jordan及びY.Weiss, 「On spectral clustering: analysis and an algorithm (スペクトルクラスタ化について: 分析及びアルゴリズム)」, Advances in NIPS (2001) を参照のこと) が適用されて、アースムーバーの距離 (Earth Mover's Distance, EMD) によって測定されたペアワイズ類似性 (pairwise similarity) に基づいて概念語彙を構築する。EMD は、Y.Rubner, C.Tomasi及びL.Guibas, 「The earth mover's distance as a metric for image retrieval (画像検索のための指標としてのアースムーバーの距離)」, IJCV (2000) に記述されている。

40

【0024】

各データ点は画像のセットとして取り扱われる。すなわち、静止ビデオ画像に対して一つの画像、及びビデオセグメントに対して複数の画像である。それから、2つのデータ点

50

(画像セット)の間の類似性を測定するためにEMDが使用される。2つの画像セットの間の距離を計算するためには多くの方法があり、例えば、これら2つのセットにおける画像の間の最大/最小/平均距離がある。これらの方法は雑音の多い異常画像(outlier images)によって容易に影響されるが、EMDは、よりロバストな距離指標を提供する。EMDは、重み正規化された制約の対象となる2つの画像セットの間の対距離(pairwise distance)の全てにおいて最小の重み付け距離を見出し、データ点の間の部分的なマッチングを許容し、異常画像の影響を低減することができる。

【0025】

2つのデータ点の間のEMDは、以下のようにして計算される。データ点 $x_1$ 及び $x_2$ にそれぞれ $n_1$ 個及び $n_2$ 個の画像があるとする。 $x_1$ 及び $x_2$ の間のEMDは、任意の2つの画像 $I_p^1$   $x_1$ 及び $I_q^2$   $x_2$ の間のフロー $f(I_p^1, I_q^2)$ によって重み付けされたグラウンド距離 $d(I_p^1, I_q^2)$ の線形組み合わせである。

【数8】

$$D(x_1, x_2) = \frac{\sum_{p=1}^{n_1} \sum_{q=1}^{n_2} d(I_p^1, I_q^2) f(I_p^1, I_q^2)}{\sum_{p=1}^{n_1} \sum_{q=1}^{n_2} f(I_p^1, I_q^2)} \quad (1)$$

ここで、最適フローマトリクス $f(I_p^1, I_q^2)$ は、以下の線形プログラムから得られる。

【数9】

$$\begin{aligned} & \min \sum_{p=1}^{n_1} \sum_{q=1}^{n_2} d(I_p^1, I_q^2) f(I_p^1, I_q^2) \\ & \text{w.r.t } f(I_p^1, I_q^2), 1 \leq p \leq n_1, 1 \leq q \leq n_2 \\ & \text{s.t. } f(I_p^1, I_q^2) \geq 0, \sum_{q=1}^{n_2} f(I_p^1, I_q^2) \leq w_p^1, \sum_{p=1}^{n_1} f(I_p^1, I_q^2) \leq w_q^2 \\ & \sum_{p=1}^{n_1} \sum_{q=1}^{n_2} f(I_p^1, I_q^2) = \min \left\{ \sum_{p=1}^{n_1} w_p^1, \sum_{q=1}^{n_2} w_q^2 \right\} \end{aligned}$$

ここで、 $w_p^1$ 及び $w_q^2$ はそれぞれデータ点 $x_1$ 及び $x_2$ における画像 $I_p^1$ 及び $I_q^2$ の重みである。ここで、等しい重み $w_p^1 = 1/n_1$ 及び $w_q^2 = 1/n_2$ を取る。概念スコア特徴に対するユークリッド距離が、距離 $d(I_p^1, I_q^2)$ として使用される。式(1)より、EMDは、2つのデータ点における最もマッチする画像の対を見出す。重み正規化制約は、各画像が他のセットで十分なマッチングを有することを確実にする。 $x_1$ 及び $x_2$ の両方が写真であるときには、EMDは単にユークリッド距離である。この対のEMDはそれから、ガウス関数 $S(x_1, x_2) = \exp(-D(x_1, x_2)/r)$ によって、ペアワイズ類似性に変換される。ここで、 $r$ は全訓練データ点の間の全ての対の距離の平均である。

【0026】

上記で言及したスペクトルクラスタ化は、データ点の対の間の類似性からなるデータセットにおけるグループを見出すための技術である。ここで、エング(Eng)らにより開発されたアルゴリズムが適用され、以下のように記述されることができる。類似性マトリクス $S(x_i, x_j)$ が与えられると、

- ・アフィンマトリクス $A_{ij} = S(x_i, x_j)$  ( $i = j$ の場合)、かつ $A_{ii} = 0$ を得る。
- ・対角線マトリクス $D_{ij} = \sum_j A_{ij}$ を定義する。 $L = D^{-1/2} A D^{-1/2}$ を得る。
- ・最大のもことから順に $n$ 個の固有値に対応する $L$ の固有ベクトル

【数10】

$$u_1, \dots, u_n$$

を見出し、

10

20

30

40

50

【数 1 1】

$$U = [u_1, \dots, u_n]$$

を得る。ここで、 $n$  は、保持すべき固有値のエネルギー比によって決定される。

・  $U$  の行が単位長さを有するように再正規化することによって、 $U$  からマトリクス  $V$  を得る。

・  $V$  における各行を  $R^n$  (元の  $i$  番目のデータ点に対応する  $i$  番目の行) における点として取扱い、 $K$ -means アルゴリズムを介して全ての点を  $n$  個のクラスにクラスタ化する。

【0027】

10

スペクトルクラスタ化アルゴリズムによって得られる各データクラスは概念ワードと呼ばれ、全てのクラスは、意味論的イベントを表し且つ検出するための概念語彙を形成する。 $W_j^i$  が意味論的イベント  $S_{E_i}$  に対して学習された  $j$  番目のワードを表し、 $S(x, W_j^i)$  が、 $x$  と  $W_j^i$  におけるメンバ・データ点 (member data points) との間の最大類似性として計算されたワード  $W_j^i$  に対するデータ  $x$  の類似性を指すとする。

【数 1 2】

$$S(x, W_j^i) = \max_{x_k \in W_j^i} S(x_k, x)$$

20

であり、ここで、 $S(x_k, x)$  は上記と同じように定義される。各データ  $x$  に対して、ベクトル  $[S(x, W_1^i), \dots, S(x, W_n^i)]^T$  が、 $x$  に対する B O F 特徴ベクトルとして取り扱われることができる。イベント  $E_t$  が  $m^t$  個のデータ点を含むとし、上記の B O F 特徴ベクトルに基づいて、イベント  $E_t$  はまた、B O F 特徴ベクトル

【数 1 3】

$$f_{\text{bof}}(E_t)$$

によって、

【数 1 4】

30

$$f_{\text{bof}}(E_t) = [\max_{x \in E_t} S(x, W_1^i), \dots, \max_{x \in E_t} S(x, W_n^i)]^T$$

として表されることもできる。最後に、B O F 特徴

【数 1 5】

$$f_{\text{bof}}$$

を使用すると、二値化された一対全 S V M 分類器が、意味論的イベント  $S_{E_i}$  を検出するために学習されることができる。

【0028】

40

ここで図 1 を参照すると、本発明のある実施形態に従ったデジタルコンテンツ記録のための意味論的イベント検出のためのシステム 100 が描かれている。このシステム 100 は、データ処理ユニット 110、周辺ユニット 120、ユーザインターフェースユニット 130、及びメモリユニット 140 を含む。メモリユニット 140、周辺ユニット 120、及びユーザインターフェースユニット 130 は、データ処理システム 110 に通信的に接続されている。

【0029】

データ処理システム 110 は一以上のデータ処理装置を含み、このデータ処理装置が、ここで記述される図 2 ~ 4 の例示的なプロセスを含む本発明の様々な実施形態のプロセスを実現する。「データ処理装置」又は「データプロセッサ」という表現は、中央処理装置

50



(CPU)、デスクトップコンピュータ、ラップトップコンピュータ、メインフレームコンピュータ、個人デジタル端末、ブラックベリー<sup>TM</sup>、デジタルカメラ、携帯電話、あるいは、電氣的、磁氣的、光學的、生物學的構成要素とともに実現されるか又はその他の方法で実現された任意の他のデータ処理、データ管理、又はデータ取扱いのための装置のような任意のタイプのデータ処理装置を含むことが意図されている。

【0030】

メモリユニット140は情報を記憶するように構成された一以上のメモリ装置を含み、記憶される情報は、ここで記述される図2～4の例示的なプロセスを含む本発明の様々な実施形態のプロセスを実行するために必要とされる情報を含む。メモリユニット140は分散プロセッサアクセス可能メモリシステムであってもよく、これは、複数のコンピュータ及び/又は装置を介してデータ処理システム110に通信的に接続された複数プロセッサアクセス可能メモリを含む。一方、メモリユニット140は、分散プロセッサアクセス可能メモリシステムである必要は無く、したがって単一のデータプロセッサ又は装置内に配置された一以上のプロセッサがアクセス可能なメモリを含み得る。さらに、「メモリユニット」という表現は、揮発性又は不揮発性、電子的、磁氣的、光學的、又はその他のものであってもよい任意のプロセッサアクセス可能データ記憶装置を含むことが意図されており、これは、フロッピー（登録商標）ディスク、ハードディスク、コンパクトディスク、DVD、フラッシュメモリ、ROM、及びRAM、又は任意の他のデジタル記憶媒体を含むが、これらに限定されるものではない。

【0031】

「通信的に接続された」という表現は、有線であっても無線であっても、データが通信され得る装置、データプロセッサ、又はプログラムの間の任意のタイプの接続を含むことが意図されている。さらに、「通信的に接続された」という表現は、単一のデータプロセッサ内の装置又はプログラム間の接続、異なるデータプロセッサに配置された装置又はプログラム間の接続、及びデータプロセッサには全く配置されていない装置間の接続を含むことが意図されている。この点に関して、メモリユニット140はデータ処理システム110から離れて示されているが、当業者は、メモリユニット140がデータ処理システム110内に完全に又は部分的に実現され得ることを理解するであろう。さらに、この点に関して、周辺システム120及びユーザインターフェースシステム130がデータ処理システム110から離れて示されているが、当業者は、それらのシステム的一方又は両方がデータ処理システム110内に完全に又は部分的に実現され得ることを理解するであろう。

【0032】

周辺システム120は、データ処理システム110にデジタルコンテンツ記録を提供するように構成された一以上の装置を含み得る。例えば、周辺システム120はデジタルビデオカメラ、携帯電話、通常のデジタルカメラ、又はその他のデータプロセッサを含み得る。加えて、周辺システム120は、データ処理システム110を離れたデータ源に接続するために必要な機器、装置、回路などを含み得る。例えば、システム100は、インターネットを介して、データセットが記憶されるサーバーにリンクされ得る。データセットは、システム100を訓練するために使用されるデジタルコンテンツ記録のデータセット、あるいは、システム100によって分析されるべきデジタルコンテンツ記録を含むデータセットを含み得る。データ処理システム110は、周辺システム120内の装置からデジタルコンテンツ記録を受領すると、そのようなデジタルコンテンツ記録を、さらなる処理のためにプロセッサアクセス可能メモリシステム140に記憶し得て、あるいは、十分な処理パワーが利用可能であれば、受領したデータストリームとしてリアルタイムでデジタルコンテンツ記録を分析し得る。

【0033】

ユーザインターフェースシステム130は、マウス、キーボード、他のコンピュータ、又はデータがそこからデータ処理システム110に入力される任意の装置又は複数の装置の組み合わせを含み得る。これに関して、周辺システム120がユーザインターフェース

システム 130 から離れて示されているが、当業者は、周辺システム 120 がユーザインターフェースシステム 130 の一部として含まれ得ることを理解するであろう。

【0034】

ユーザインターフェースシステム 130 はまた、ディスプレイ装置、プロセッサアクセス可能メモリ、又はデータ処理システム 110 によってデータがそこに出力される任意の装置又は複数の装置の組み合わせを含み得る。これに関して、ユーザインターフェースシステム 130 がプロセッサアクセス可能メモリを含むならば、そのようなメモリは、ユーザインターフェースシステム 130 及びメモリユニット 140 が図 1 では離れて示されているが、メモリユニット 140 の一部であり得る。

【0035】

システムの基本的な動作がここで図 2 を参照して記述される。これは図 1 に描かれたユニットの一以上によって実現される処理モジュールを描く流れ図である。処理モジュールがシステム 100 に設けられたユニットの一以上によって実行される指示を含むことが理解されるべきである。図示されている例では新しいイベント ( $E_0$ ) がデータエントリモジュール 200 を介してシステム 100 に与えられる。ここでは  $E_0$  が特定の意味論的イベントに属する確率が決定されることが望ましい。例えば、ユーザインターフェースユニット 130 を介して受領されたオペレータ指示に基づいて、データ処理ユニット 110 は、 $E_0$  に対応するデータをメモリユニット 140 にダウンロードするように周辺ユニット 120 の動作を制御する。各イベントは、複数のデジタルコンテンツ記録を含み、図示されている例ではデジタル静止画像  $m_{0,p}$  及びビデオコンテンツ  $m_{0,v}$  を含む。これらのデジタルコンテンツ記録は、撮影時刻及び色の類似性に基いて、先に記述されたクラスタ化方法を利用して一緒にグループ化される。クラスタ化方法は、システム 100 への提出に先立って、静止デジタル画像及びビデオセグメントのデータセットに適用されることができる。あるいは、データセットがシステム 100 に与えられて、データエントリモジュール 200 が、 $E_0$  を生成するためにデータ処理ユニット 110 の一つの動作要素としてクラスタ化動作を実行してもよい。

【0036】

例えば、消費者は電子カメラを使用して、複数の異なるイベントの 100 個のデジタル静止画像及びビデオからなるデータセットを撮影し得る。電子カメラからのメモリカードが、周辺ユニット 120 の一部としてのカードリーダーユニットに提供される。ユーザインターフェースユニット 130 を介してユーザによって入力された制御指示に反応して、データ処理ユニット 110 は、データセットをメモリカードからメモリユニット 140 にダウンロードするように、周辺ユニット 120 の動作を制御する。データ処理ユニット 110 はそれから先に進んで、デジタル静止画像及びビデオを複数のイベントに対応する複数のクラスタにグループ化するために、データセットに対してクラスタ化アルゴリズムを実行する。これにより、データエントリモジュール 200 内に提供された指示の機能が完了し、ある数のデジタル静止画像及びビデオ (例えば元の 100 個のうちの 10 個) が、 $E_0$  に関連しているとして識別される。この時点で、10 個のデジタル静止画像及びビデオが  $E_0$  に関連付けられるが、 $E_0$  は、「結婚」のような特定の意味論的イベントにはまだ関連付けられていない。

【0037】

視覚的特徴抽出モジュール 210 がそれから使用されて、 $E_0$  内のビデオセグメントからキーフレームを獲得し、キーフレーム及び  $E_0$  内に含まれるデジタル静止画像の両方から視覚的特徴が抽出される。図示されている例では、視覚的特徴抽出モジュール 210 は、格子ベースの色モーメント、ガボール・テクスチャ (Gabor texture)、及びエッジ方向性ヒストグラムを、デジタル静止画像及びビデオの各々に対して決定する。しかし、図示されている例で使用されているもの以外の視覚的特徴が容易に利用され得ることが、理解されるであろう。

【0038】

データ処理ユニット 110 は、視覚的特徴抽出モジュール 210 内に提供された指示に

10

20

30

40

50

したがって、 $E_0$ とともに含まれるデジタル静止画像及びビデオの各々に対して、従来の技術を利用して必要なキーフレーム及び視覚的特徴の抽出を実行する。したがって、 $E_0$ に対応する10個のデジタル静止画像及びビデオの各々についての3つの視覚的特徴表現が、さらなる分析のためにここでは利用可能である。

#### 【0039】

特徴抽出モジュール210によって抽出された3つの視覚的特徴は、概念検出モジュール220によって使用されて、特定のキーフレーム又は静止デジタル画像が特定の意味論的イベントに関係している確率を反映した概念スコアを生成する。概念検出モジュール220は、好ましくは、2ステップのプロセスを使用して概念スコアを決定する。第1に、概念スコア検出モジュール222が設けられ、これは、21個の上記のSVM意味論的概念決定子(データ処理ユニット110によって実現される)を利用して、各デジタル静止画像及びキーフレームに対する各視覚的特徴空間における各々の個別の分類器に基づいて、概念スコアを生成する。第2に、個々の概念スコアがそれから融合モジュール224(データ処理ユニット110によって実現される)によって融合され、特定のデジタル静止画像及びキーフレームに対するアンサンブル概念検出スコアを生成し、それによって、さらに処理されるべきデータ量を低減する。

#### 【0040】

好適な実施形態では、融合モジュール224は最初に、異なる特徴からの異なる分類出力を、シグモイド関数  $1 / (1 + \exp(-D))$  によって正規化する。ここで、 $D$ は、決定境界までの距離を表すSVM分類器の出力である。融合は、21個の概念の各々に対する異なる視覚的特徴からの分類出力の平均を取ることによって完了されて、アンサンブル概念検出スコアを生成する。

#### 【0041】

単純化された例では、3つの概念「人々」「公園」及び「花」が議論される。「人々」「公園」及び「花」に対する概念スコアが、 $E_0$ の10個の画像の各々の3つの視覚的特徴表現の各々について生成される。例えば、10個の画像のグループの最初の画像の色の特徴表現は、人々を含む確率が90%、公園を含む確率が5%、及び花を含む確率が5%であり得て、最初の画像のテクスチャの特徴表現は、人々を含む確率が5%、公園を含む確率が80%、及び花を含む確率が15%であり得て、最初の画像のエッジ方向の特徴表現は、人々を含む確率が10%、公園を含む確率が50%、及び花を含む確率が40%であり得る。

#### 【0042】

10個の画像の3つの視覚的特徴表現が与えられると、30セットの概念スコアが生成され(各視覚的特徴表現に対して一つ)、各セットは3つの個別の概念スコア(「人々」に対して一つ、「公園」に対して一つ、及び「花」に対して一つ)を含む。最初の画像に対するアンサンブル概念スコアを生成するために、視覚的表現の各々に対する概念の各々についての確率が平均され、第1の画像のアンサンブル概念スコアは、人々を含む確率が35%(人々の確率として色90%、テクスチャ5%、エッジ5%の平均)、公園を含む確率が30%(公園の確率として色5%、テクスチャ80%、エッジ5%の平均)、及び花を含む確率が20%(花の確率として色5%、テクスチャ15%、エッジ40%の平均)となる。

#### 【0043】

アンサンブル概念スコアは引き続いてBOFモジュール230に与えられ、これが $E_0$ に対するBOFベクトルを決定する。 $E_0$ に対するBOF特徴ベクトルは、最初に、各々の各デジタル静止画像及びビデオセグメントのアンサンブル概念スコアを使用して $E_0$ 内に含まれるデジタル静止画像及びビデオセグメントの各々に対する個別の特徴ベクトルを決定することによって得られる。好適な実施形態では、各デジタル静止画像又はビデオセグメントはデータ点として扱われて、 $E_0$ 内の各データ点のアンサンブル概念スコアの間のペアワイズ類似性、及び所与の意味論的イベント(SE)、例えば「結婚」に対する各々の予め定められた正の訓練データ点のアンサンブル概念スコアが、それからEMDを使

用して類似性検出モジュール 2 3 2 によって計算される。効果的には、個別の特徴ベクトルは、 $E_0$ 内に含まれるデジタル静止画像及びビデオセグメントの各々に対して得られる。マッピングモジュール 2 3 4 がそれから、 $E_0$ の個別の特徴ベクトルの各々を意味論的イベントのコードブック（以下により詳細に記述される訓練プロセスの間に先に開発されている）にマッピングするために使用され、 $E_0$ に対するイベント特徴ベクトルが、マッピングされた類似性に基づいて生成される。

#### 【 0 0 4 4 】

イベント特徴ベクトルが、ここで分類器モジュール 2 4 0 に供給されることができる。描かれている例では、分類器モジュール 2 4 0 は SVM 分類器を使用して、 $E_0$ に対するイベント検出スコアを生成する。イベント検出スコアは、新しいイベント  $E_0$  が「結婚」のような所与の意味論的イベントに対応する最終的な確率を表す。イベント検出スコアはそれから、好ましくは予め定められた閾値と比較され、 $E_0$ が結婚イベントとしてカテゴリ化されるべきかどうか決定される。予め定められた閾値は、所与のアプリケーションにてシステム 1 0 0 によって要求される正確さのレベルに依存して変化してもよい。

#### 【 0 0 4 5 】

ひとたび  $E_0$ が適切にカテゴリ化されると、 $E_0$ に対応する静止デジタル画像及びビデオセグメントが、適切な意味論的イベント分類器でタグ付けされ、適切なアルバムフォルダ又はファイルに分類されて、後の検索のためにメモリユニット 1 4 0 内に記憶されることができる。あるいは、タグ付けされた静止デジタル画像及びビデオセグメントは、周辺ユニット 1 2 0 を介して画像記憶媒体に書き込まれることができる。静止デジタル画像及びビデオセグメントの意味論的イベント分類器によるタグ付けは、画像及びビデオセグメントがサーチエンジンによって容易に検索されることを可能にするという付加的な効果を提供する。

#### 【 0 0 4 6 】

システム 1 0 0 の訓練が、ここで図 3 を参照して記述される。最初に、 $T$ 個の正の訓練イベント  $E_1, \dots, E_T$  が、データエントリモジュール 2 0 0 を使用して入力される。各イベント  $E_t$  は、 $m_{t,p}$  個の写真及び  $m_{t,v}$  個のビデオを含み、これらは、先に記述されたクラスタ化方法によって、撮影時間及び色の類似性にしたがってグループ化されている。視覚的抽出モジュール 2 1 0 がそれから使用されて、ビデオセグメントからキーフレームを抽出し、キーフレーム及びデジタル静止画像の両方から視覚的特徴が抽出される。上述の動作の場合と同様に、視覚的特徴は、格子ベースの色モーメント、ガボール・テクスチャ、及びエッジ方向性ヒストグラムを含む。概念検出モジュール 2 2 0 がそれから使用されて、上述のようにキーフレーム及びデジタル画像に対するアンサンブル概念スコアを生成する。

#### 【 0 0 4 7 】

BOF 学習モジュール 2 5 0 がそれから使用されて、システム 1 0 0 を訓練する。最初に、各デジタル画像又はビデオセグメントがデータ点として取り扱われて、データ点の各々の対の間のペアワイズ類似性が、先に記述された類似性検出モジュール 2 3 2 を使用して EMD によって計算される。ペアワイズ類似性マトリクスに基づいて、スペクトルクラスタ化モジュール 2 5 2 が使用されてスペクトルクラスタ化を適用し、データ点を異なるクラスタにグループ化する。ここで、各クラスタは一つのコードワードに対応する。意味論的イベント SE を検出するために分類器を訓練するために、全ての訓練イベント  $E_i$  ( $E_i$  は SE に対する正の訓練イベント及び負の訓練イベントの両方を含む) が上述のコードブックにマッピングされて、マッピングモジュール 2 5 4 によって各訓練イベントに対する BOF 特徴ベクトルが生成される。BOF 特徴ベクトルに基づいて、分類器訓練モジュール 2 6 0 が使用されて、特定の意味論的イベント SE を検出するように二値化 SVM 分類器を訓練する。

#### 【 0 0 4 8 】

図 4 は、概念スコア検出モジュール 2 2 2 で使用されるビデオ概念検出器のための訓練プロセスの詳細を描いている。概念 C に対して、ベンチマーク消費者ビデオデータセット

10

20

30

40

50

からN個の正の訓練ビデオがデータエントリモジュール200を介して提供される。キーフレームがビデオから得られて、視覚的特徴が、先の例におけるように視覚的特徴抽出モジュール210を使用してキーフレームから得られる。視覚的特徴は、格子ベースの色モーメント、ガボール・テクスチャ、及びエッジ方向性ヒストグラムを含む。概念訓練モジュール270がそれから使用されて、概念決定子を訓練する。すなわち、視覚的特徴の各タイプに基づいて、各キーフレームは特徴ベクトルとして表されて、二値化SVM分類器が概念Cを検出するように訓練される。特徴の個々のタイプに対するこれらの分類器の判別機能が一緒に平均されて、概念Cに対するアンサンブル概念検出器を生成する。

#### 【0049】

上述された意味論的決定システム及び方法のテストが、コダックの消費者データセットから1972個の消費者イベントを評価することによって実行された。イベントは10個の異なる意味論的イベントにラベルされ、その詳細な定義が、図5に与えられた表に示されている。合計1261個のイベントが訓練のためにランダムに選択され、残りはテストのために使用された。訓練及びテストデータはマクロイベントレベルで区分された。すなわち、同じマクロイベントからのイベントが、訓練又はテストデータとして一緒に扱われた。これは、同じマクロイベントからの類似のイベントが分離されることを避けるためであり、分類問題を単純化する。

#### 【0050】

平均精度 (average precision, AP) が性能の指標として使用されたが、これは、ビデオ概念検出のための公式指標として使用されている。例えば、ニスト (Nist) の「ツリービデオ検索評価 (Tree video retrieval evaluation (treevid))」, 2001 - 2006, <http://www-nlpir.nist.gov/projects/treevid>を参照のこと。これは、精度 - 再生曲線における異なる再生点での精度値の平均を計算し、これにより特定の意味論的イベントを検出する際の分類器の有効性を評価する。複数の意味論的イベントを考慮するときには、APの平均 (mean of APs, MAP) が使用される。

#### 【0051】

意味論的イベント検出アルゴリズムにおける概念スコア表現の有効性を示すために、本発明の方法とBOF特徴ベクトルが元の低レベルの視覚的特徴に基づいて構築されるアプローチとを比較する実験が行われた。具体的には、上述のS.F.Changらに記述されたものと同じ低レベルの視覚的特徴が使用された。図6は、性能の比較を与える。図6に示されるように、両方の方法は、一致してランダム推測よりも性能がよい。しかし、概念スコアを伴うSE検出は、大抵の概念について、APに関して低レベル特徴を伴うSE検出よりも性能がよく、MAPに関しては20.7%よい。この結果は、意味論的イベントの検出手助けするために先の概念検出モデルを使用するパワーを確かめるものである。

#### 【0052】

第2の実験が行われ、イベントレベル表現対画像レベル表現の比較を行った。この実験では、本発明の意味論的検出方法 (SE検出) と2つの他の検出器、すなわち (1) ベースラインイベント検出器 (ベースライン) 及び (2) 画像レベル概念スコア表現を直接的に使用するSVM検出器 (SVMダイレクト) との間の比較が行われた。図7は、異なる方法のAP比較を与える。示されるように、提案されるSE検出は意味論的イベントの大抵に対して最も良く機能し、「結婚」「クリスマス」及び「学校活動」のような多くの意味論的イベントに対して、2番目に良い方法に比べて20%より多くの顕著な性能の改善を得ている。この結果は、イベントレベルBOF表現の成功を確かめるものである。加えて図8は、異なるアルゴリズムからのサポートベクトルの数の比較を与える。一般的に、サポートベクトルが少ないほど、決定境界が単純になる。図より、決定境界は、イベントレベル表現によって顕著に単純化され、SVM分類器は意味論的イベントを非常によく分離することができている。さらに、ベースライン検出器及びSE検出による「動物」に対するトップ5個の検出イベントの比較は、SE検出方法が100%の精度を達成することができるのに対して、画像ベースのSVMダイレクト方法は20%の精度しか得ることが

10

20

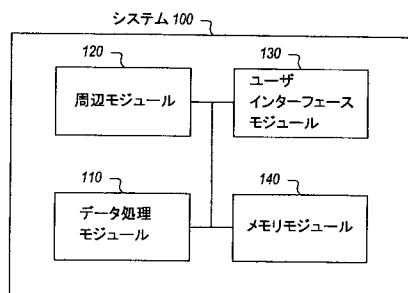
30

40

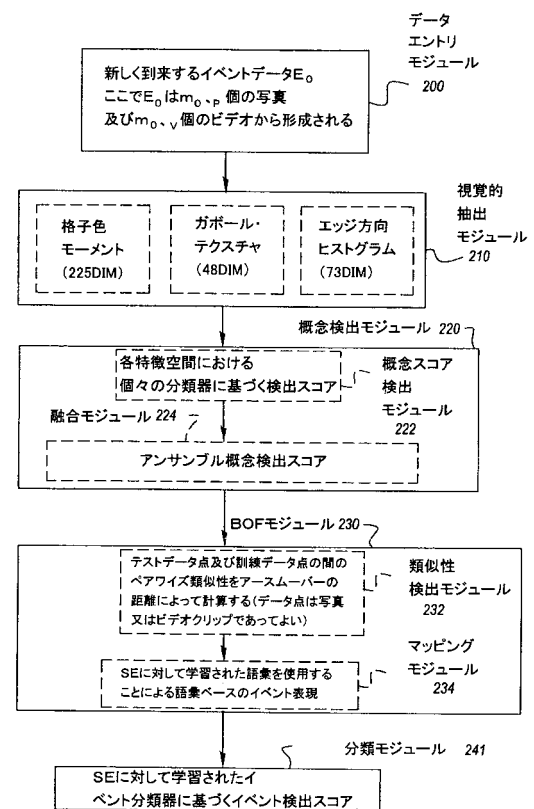
50

できないことを示した。

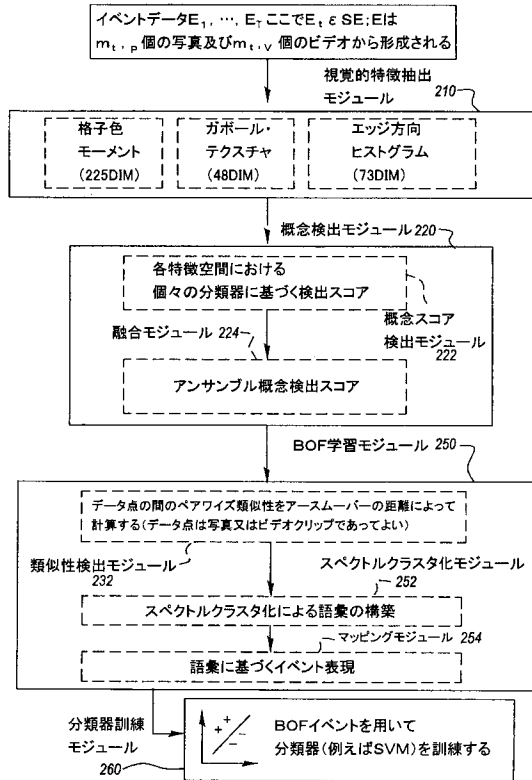
【図 1】



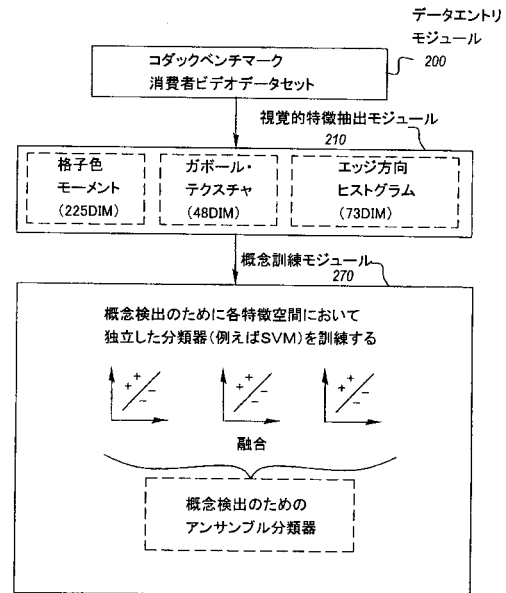
【図 2】



【図 3】



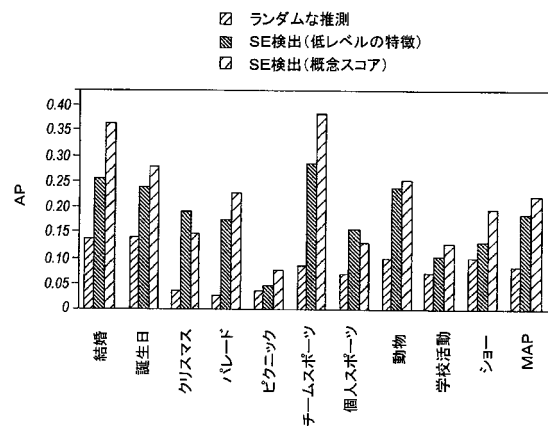
【図 4】



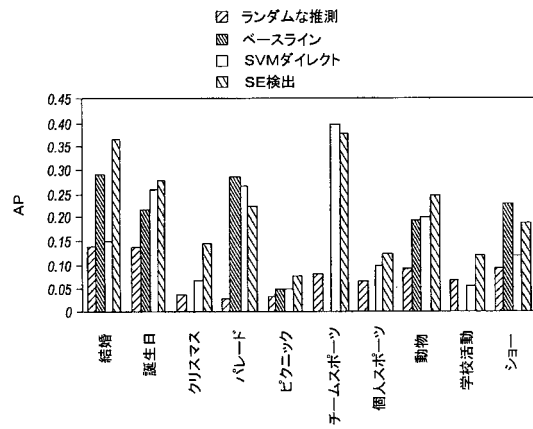
【図 5】

意味論的イベント	定義
結婚	花嫁及び花婿、ケーキ、装飾された車、レセプション、結婚パーティー、あるいは結婚の日に関連した何か
誕生日	誕生日ケーキ、風船、包装されたプレゼント、及び誕生日の帽子 通常は有名な誕生日の歌を伴う
クリスマス	クリスマスツリー及び通常のクリスマス装飾、必ずしもクリスマスの日に撮影されない
パレード	公共の場を動いていく人々又は車両の進行
ピクニック	屋外、ピクニックテーブルがある又はない、日よけがある又はない、視野の中の人々及び食べ物
チームスポーツ	バスケットボール、野球、フットボール、ホッケー、及び他のチームスポーツ
個人スポーツ	テニス、水泳、ボートリング、及び他の個人スポーツ
動物	ペット (例えば、犬、猫、馬、鳥、ハムスター)、野生の動物、動物園、及び動物ショー
学校活動	学校の卒業、学校の初日又は最終日、及び学校に関する他のイベント
ショー	ショー及びコンサート、リサイタル、演劇、及び他のイベント

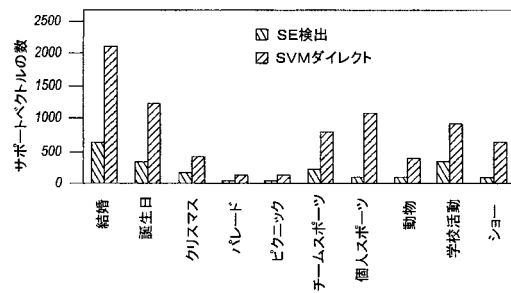
【図 6】



【図 7】



【図 8】





---

フロントページの続き

(72)発明者 ルイ アレクサンダー シー

アメリカ合衆国 ニューヨーク ロチェスター ステート ストリート 343

(72)発明者 ジアン ウェイ

アメリカ合衆国 ニューヨーク ニューヨーク ウェスト 120 ティーエイチ ストリート

マッド 500 エス ダブリュー 1300

審査官 真木 健彦

(56)参考文献 特開2007-317077(JP,A)

特開2003-153007(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06T 1/00

G06T 7/00

G06F 17/30