

【特許請求の範囲】

【請求項 1】

検索において識別された情報項目の集合から情報項目のマップを表すデータを受け取り、そのマップは情報項目の相互類似性に基づく配列内の位置関係でその識別された情報項目を提供し、類似する情報項目はその配列内で類似の位置にマッピングされ、

前記マップデータを処理して情報項目の階層クラスタ分割を形成し、その階層クラスタ分割は情報項目の第 1 のレベルのクラスタ分割と、その第 1 のレベルのクラスタ内の情報項目クラスタに対する少なくとももう 1 つのクラスタ分割レベルを提供することを特徴とするマッピングプロセッサよりなる情報検索装置。

【請求項 2】

前記情報項目は複数の特徴付ける情報特徴を含み、各情報項目のその特徴付ける情報特徴は各情報項目に対する特徴ベクトルを形成するのに使われ、その特徴ベクトルは配列内の位置にその情報項目をマッピングすることによりそのマップデータを形成するのに使われることを特徴とする請求項 1 に記載の情報検索装置。

【請求項 3】

前記マッピングプロセッサはその第 1 のクラスタ分割レベルの情報項目にその第 1 のレベルのクラスタのそれぞれと関連する特徴付ける情報特徴を与え、その第 1 のレベルのクラスタ内の情報項目のクラスタに対する特徴付ける情報特徴を他の階層レベルにおいて提供することを特徴とする請求項 1 又は請求項 2 に記載の情報検索装置。

【請求項 4】

各第 1 のレベルのクラスタと関連する特徴付ける情報特徴とその他の情報項目クラスタ分割レベル内の各クラスタと関連する特徴付ける情報特徴とは、それぞれのクラスタと関連する情報項目内に存在する最も共通的な特徴付ける情報特徴から生成されることを特徴とする請求項 1 乃至 3 のいずれか 1 項に記載の情報検索装置。

【請求項 5】

ある下位レベルクラスタ内の情報項目のクラスタ同士は相互に関連するが、一方、第 1 のレベルのクラスタ同士はその下位レベルクラスタ内の情報項目に関する追加情報項目クラスタであることを特徴とする請求項 1 乃至 4 のいずれか 1 項に記載の情報検索装置。

【請求項 6】

各クラスタと関連する特徴付ける情報項目は各クラスタ内の情報項目のそれぞれと関連するテキスト情報の最も共通的なワードであることを特徴とする請求項 3 乃至 5 のいずれか 1 項に記載の情報検索装置。

【請求項 7】

前記情報項目はテキスト情報からなり、その特徴付ける情報特徴はワードであり、情報項目に対する特徴ベクトルはその情報項目内におけるワードグループそれぞれの、出現頻度の集合を表すことを特徴とする請求項 1 乃至 6 のいずれか 1 項に記載の情報検索装置。

【請求項 8】

前記情報項目はテキスト情報を含み、その特徴付ける情報特徴はワードであり、配列内の位置がそのテキスト情報の少なくとも部分的な相互類似性によりマッピングされることを特徴とする請求項 7 に記載の情報検索装置。

【請求項 9】

情報項目の集合中で出現頻度が閾値以上の頻度を有するテキスト情報内のワードを除外するマッピングに対する情報項目の前処理を行うことを特徴とする請求項 7 又は 8 に記載の情報検索装置。

【請求項 10】

情報項目の集合中で出現頻度が閾値以下の頻度を有するテキスト情報内のワードを除外するマッピングに対する情報項目の前処理を行うことを特徴とする請求項 7 乃至 9 のいずれか 1 項に記載の情報検索装置。

【請求項 11】

図形表示装置の表示領域内の n 次元表示配列の表示点として、識別された情報項目に対

10

20

30

40

50

応する配列位置の少なくとも幾つかの表示を表示する、グラフィカルユーザインタフェース（GUI）と組み合わせたディスプレイプロセッサを備えることを特徴とする請求項 1 乃至 10 のいずれか 1 項に記載の情報検索装置。

【請求項 12】

前記表示領域は少なくとも 2 つの領域を含み、その 1 つの領域は第 1 の階層レベルのクラスタの n 次元表示を提供し、その他の領域はその他の階層レベルのクラスタの n 次元表示を提供し、 n は整数であることを特徴とする請求項 11 に記載の情報検索装置。

【請求項 13】

前記情報項目のワード検索を実行する検索プロセッサを備え、前記検索プロセッサとそのグラフィカルユーザインタフェースは識別された情報項目に対応する表示点だけを表示するように協調するよう構成されていることを特徴とする請求項 11 又は 12 に記載の情報検索装置。 10

【請求項 14】

前記ディスプレイプロセッサは、それがグラフィカルユーザインタフェース上に表示されたときに、階層レベルの 1 つにおける第 1 のクラスタを見ているユーザに、その階層レベル内の別のクラスタの位置の n 次元空間の相対方向を提供する指示を表すデータを生成することを特徴とする請求項 11 乃至 13 のいずれか 1 項に記載の情報検索装置。

【請求項 15】

前記ディスプレイプロセッサは、その他のクラスタ内の情報項目の数を表すデータを生成し、前記情報項目の数はその第 1 のクラスタに関係するその他のクラスタの n 次元空間内の相対方向の指示と関連することを特徴とする請求項 14 に記載の情報検索装置。 20

【請求項 16】

前記ディスプレイプロセッサはグラフィカルユーザインタフェースと組合わさってその図形表示装置の第 1 の領域内のその他のクラスタの相対位置の指示を表示し、そのクラスタ内の情報項目の数を表すデータがその指示により表示可能であることを特徴とする請求項 14 又は 15 に記載の情報検索装置。

【請求項 17】

ユーザ制御ポインタを使ってその n 次元空間内の情報項目或いは情報項目のクラスタを選択するユーザ制御装置を備え、その指示上に置かれるそのポインタに応答して、相対方向の指示に関する情報項目の数が表示されることを特徴とする請求項 16 に記載の情報検索装置。 30

【請求項 18】

次元数は 2 であることを特徴とする請求項 12 乃至 17 のいずれか 1 項に記載の情報検索装置。

【請求項 19】

前記情報項目はテキスト情報を有するビデオデータを含むことを特徴とする請求項 1 乃至 18 のいずれか 1 項に記載の情報検索装置を備えるビデオ収集及び / 又は処理装置。

【請求項 20】

情報項目を含む記憶装置と、

前記記憶装置を情報検索装置と接続するデータ通信ネットワークとを備えることを特徴とする請求項 19 に記載のビデオ収集及び / 又は処理装置。 40

【請求項 21】

前記情報項目はその情報項目からの代表画像を提供する代表キースタンプを含むことを特徴とする請求項 19 に記載のビデオ収集及び / 又は処理装置。

【請求項 22】

クラスタと関連する共通特徴付ける情報特徴はそのクラスタに共通する代表キースタンプを含むことを特徴とする請求項 21 に記載のビデオ収集及び / 又は処理装置。

【請求項 23】

検索で識別された情報項目の集合から情報項目のマップを表すデータを受け取るステップと、ここでそのマップはその情報項目の相互類似性に基づく配列内の位置に関してその 50

識別された情報項目を提供し、類似の情報項目はその配列内の類似の位置にマッピングされ、

第1のクラスタ分割レベルの情報項目及びその第1のレベルのクラスタ内の情報項目クラスタに対する少なくとももう1つのクラスタ分割レベルの情報項目を提供する階層分割された情報項目を構成するようにそのマップデータを処理するステップとからなることを特徴とする情報検索及び表示方法。

【請求項24】

前記マップデータ処理ステップは、第1のクラスタ分割レベルの情報項目に、その第1のレベルのクラスタの情報項目のそれぞれと関連する特徴付ける情報特徴を与えるステップと、その第1のレベルのクラスタ内の情報項目のクラスタに対する特徴付ける情報特徴をその他の階層レベルにおいて提供するステップとを含むことを特徴とする請求項23に記載の情報検索及び表示方法。

10

【請求項25】

上記情報項目は複数の特徴付ける情報特徴を含み、各情報項目のその特徴付ける情報特徴は各情報項目の特徴ベクトル形成に使用され、その特徴ベクトルはその情報項目を配列内の位置にマッピングに用いられることを特徴とする請求項24に記載の情報検索及び表示方法。

【請求項26】

図形表示装置上の表示領域内に2次元表示配列の表示点として少なくともその配列の幾つかの位置の表示を表示するステップを含むことを特徴とする請求項23乃至25のいずれか1項に記載の情報検索及び表示方法。

20

【請求項27】

前記表示領域は少なくとも2つの領域を含み、その1つの領域は第1の階層レベルのクラスタのn次元表示を提供し、他方の領域はその他の階層レベルのクラスタのn次元表示を提供することを特徴とする請求項26に記載の情報検索及び表示方法。

【請求項28】

次元nの数は2であることを特徴とする請求項27に記載の情報検索及び表示方法。

【請求項29】

請求項23乃至27のいずれか1項に記載された情報検索及び表示方法を実行するためのプログラムコードを備える計算機ソフトウェア。

30

【請求項30】

請求項29に記載のプログラムコードを提供するための提供媒体。

【請求項31】

前記媒体は記憶媒体であることを特徴とする請求項30に記載の提供媒体。

【請求項32】

前記媒体が伝送媒体であることを特徴とする請求項30に記載の提供媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、大量のコンテンツを扱う情報検索装置及びその方法に関する。

40

【背景技術】

【0002】

キーワードに基づく検索によって情報（例えば、文書、画像、電子メール、特許、音声又は映像コンテンツなどのインターネットコンテンツ又はメディアコンテンツ）を探し出す定着した多数のシステムがある。その具体例として、「グーグル（Google）（登録商標）」や「ヤフー（Yahoo）（登録商標）」を始めとするインターネットサーチエンジンがあり、キーワードで実行される検索により、サーチエンジンによって認識され、関連度の順にランク付けされた結果の一覧が得られる。

【0003】

しかし、大量コンテンツコレクションと称されることが多い大量のコンテンツを包含す

50

るシステムにおいては、効果的な検索クエリを策定し、比較的短い検索「ヒット(hits)」の一覧を得ることは困難である。例えば、本願作成時点に行った、キーワード「大量コンテンツコレクション(massive document collection)」に対するグーグルでの検索では、243000件のヒットが引き出された。インターネットを通じて蓄積されるコンテンツ量は概して経時的に増大するので、検索がこの後繰り返される場合、このヒット数は増大することが予想される。そのようなヒット一覧を精査することには、非常に時間がかかる可能性が大である。

【0004】

概して、大量コンテンツコレクションが良好に利用されない理由には、以下の点が問題である。

- ・利用者が関連するコンテンツの存在を知らない。
- ・利用者は関連するコンテンツの存在を知っているが、そのコンテンツが置かれている場所を知らない。
- ・利用者はコンテンツの存在を知っているが、それが関連性のあるものかどうかを知らない。
- ・利用者は関連コンテンツの存在及びその見出し方を知っているが、コンテンツを見つけるのに長い時間がかかる。

【0005】

論文「大量文書コレクションの自己組織化(Self Organization of a Massive Document Collection)」、コホネン(Kohonen)他、ニューラルネットワークに関するIEEEトランザクション(IEEE Transactions on Neural Networks)、第11巻、第3号、2000年5月、第574～585頁には、所謂「自己組織化マップ(Self Organizing Map(SOM))」を用いた技術が開示されている。これらの自己組織化マップは、各文書の特性を表す「特徴ベクトル(feature vectors)」がSOMのノード上にマッピングされる、所謂非管理型自己学習ニューラルネットワークアルゴリズム(unsupervised self-learning neural network algorithm)を利用する。

【0006】

コホネン等の論文において、第1のステップは、文書テキストを前処理し、次いで、各前処理された文書から特徴ベクトルが導かれることである。1つの形態において、これは、各単語の大辞書での出現頻度を示すヒストグラムであり得る。ヒストグラム中の各データ値(すなわち、各々の辞書単語の各出現頻度)は、 n 値ベクトル中の値となるが、ここで n は辞書中の候補単語の総数である(この論文において記載されている例においては43222)。重み付けが n 値ベクトルに与えられ得、それによって、ある幾つかの単語の増大した関連度又は改善された区別が強調されることになる可能性がある。

【0007】

次いで、 n 値ベクトルは、より大きさが小さいベクトル、すなわち、 n よりも実質的に小さい数 m (この論文における例では500)を有するベクトルにマッピングされる。マッピングは、乱数配列からなる($n \times m$)の「射影行列」でベクトルを乗算することによって達成される。この技術は、いずれか2つの縮小された大きさのベクトルが、2つの各々の入力ベクトルと同等のベクトル内積(dot product)を有する、より小さい大きさのベクトルを生じさせることが示されている。このベクトルマッピングプロセスは、論文「ランダムマッピングによる次元圧縮: クラスタリングのための高速類似性演算(Dimensionality Reduction by Random Mapping: Fast Similarity Computation for Clustering)」、カスキ(Kaski)、Proc. IJCNN、第413～418頁、1998年に記載されている。

【0008】

次いで、次元が圧縮されたベクトルは、各ベクトルを「モデル(model)」(別のベクトル)で乗算するプロセスによってSOM上のノード(ニューロンとも称される)上にマッピングされる。モデルは、SOM上への相互類似性によって自動的にモデルを配列する学習プロセスによって作成され、SOMはノードの二次元グリッドとして通常表される。

10

20

30

40

50

これは簡単な処理ではなく、コホネンらはこれに、700万を丁度下回る数の文書の文書データベースのために、800MBのメモリを有する6プロセッサのコンピュータで6週間かかった。最後に、ユーザがマップの複数の領域にズームしてノードを選択できる状態にSOMを形成するノードのグリッドが表示されるが、これによってユーザインタフェースがそのノードにリンクされた文書を含むインターネットのページへのリンクを提供する。

【発明の開示】

【発明が解決しようとする課題】

【0009】

本発明は情報項目の大規模データベースから情報項目の検索結果を提供する実際的かつ管理可能な方法を見出す技術課題を取り扱う。 10

【課題を解決するための手段】

【0010】

本発明の一態様によれば、検索中に確認された情報項目集合から情報項目のマップを表すデータを受け取り動作可能なマッピングプロセッサからなる情報検索装置が提供される。そのマップは識別された情報項目をその情報項目の相互類似性に応じた配列中における位置として提供する。マッピングデータは類似の情報項目はその配列中で類似の位置をマッピングするように構成される。そのマッピングプロセッサは、第1のクラスタ分割レベルの情報項目とその第1のレベルのクラスタ内の情報項目クラスタに対する少なくとももう1つのクラスタ分割レベルの情報項目とを提供する階層的クラスタ分割情報項目を形成 20
するようにマップデータを処理できる。情報項目を階層構造のクラスタに分割形成することにより、情報項目のナビゲーションや表示が容易になる。

【0011】

更にマッピングプロセッサは第1のクラスタ分割レベルの情報項目に対し第1のレベルのクラスタの情報項目のそれぞれと関連付けられる特徴付ける情報特徴を与えるようにしてもよい。マッピングプロセッサはそれに対応して第1のレベルのクラスタ内の情報項目クラスタに対する特徴付ける情報特徴を別の階層レベルで与えるようにしてもよい。特徴付ける情報特徴は1つのクラスタを別のクラスタから容易に区別する手段を提供する。

【0012】

一般に、情報項目は複数の特徴付ける情報特徴を含んでおり、各情報項目の特徴付ける 30
情報特徴は各情報項目に対する特徴ベクトルを形成するのに用いられ、その特徴ベクトルはその情報項目を配列内の1つの位置にマッピングするのに使われる。このように、各第1のレベルのクラスタと関連付けられた特徴付ける情報特徴とその他の情報項目レベル内の各クラスタと関連付けられた特徴付ける情報特徴とは各クラスタと関係する情報項目内に存在する最も共通的な特徴付ける情報特徴から形成される。

【0013】

例えば、情報項目はテキスト情報であってもよく、特徴付ける情報特徴はワードであり、情報項目に対する特徴ベクトルはその情報項目のワードグループのそれぞれの出現頻度の集合を現している。この例に関しては、各クラスタと関連する特徴付ける情報特徴が各 40
クラスタ内の情報項目のそれぞれと関連するテキスト情報の最も共通的なワードである。

【0014】

情報検索装置は図形表示装置上の表示領域にn次元の表示配列の表示点として少なくともその配列の幾つかの点を表す表示を行うことのできるグラフィカルユーザインタフェースと組み合わせたディスプレイプロセッサからなってもよい。その表示領域は少なくとも2つの領域を含んでもよく、1つの領域は第1の階層レベルのクラスタのn次元表示領域を提供し、もう一方の領域はもう一方の階層レベルのクラスタのn次元表示領域を提供するものである。n次元の次数は整数であればよく、限定するものではないがその次数は2であってもよい。もちろん1でも3でもよいことは当然理解されるべきである。

【発明の効果】

【0015】

本発明の実施の形態によって提供される利点は、1つ以上の表示領域部分を備え、異なる階層レベルの情報項目を表示するのに便利な手段を提供する。例えば、第1のレベルの情報項目を一方の領域に表示しながら第1の領域から選択されたクラスタに現れる情報項目を第2の領域に表示することができる。したがって、検索によって疎母集団の配列であることが分かれば、第1の領域で明らかにされた異なるクラスタ間の相対的なナビゲーションがより簡単に管理できるようになり、より詳細な情報項目の表示が第2の領域内で選択表示されたクラスタで提供できる。

【発明を実施するための最良の形態】

【0016】

添付図面を参照して、本発明の実施の形態を例としてのみ説明する。

10

【0017】

図1は、プログラム及びデータ用のディスク記憶装置30を備えたプロセッサユニット20と、イーサネット（登録商標）ネットワーク又はインターネットなどのネットワーク50に接続されたネットワークインタフェースカード40と、陰極線管装置60などの表示装置と、キーボード70と、マウス80などのユーザ入力装置とを有する汎用コンピュータ10をベースとする情報記憶及び検索システムの概略図である。情報記憶及び検索システムはプログラム制御下で動作し、プログラムは、ディスク記憶装置30上に記憶され、例えば、ネットワーク50、着脱式ディスク（図示せず）又はディスク記憶装置30上へのプリインストールによって与えられる。

【0018】

20

この情報記憶及び検索システムは、2つの一般的な動作モードで動作する。第1のモードにおいては、一組の情報項目（例えば、テキスト情報項目）がディスク記憶装置30又はネットワーク50を介して接続されたネットワークディスクドライブ上で編集され、検索動作に備えて分類及び索引付けされる。第2の動作モードは、索引付け及び分類されたデータに対して実際に検索を行うことである。

【0019】

実施の形態は多くの種類の情報項目に適用可能である。適切な種類の情報を全て網羅するものではないが、この一覧には、特許、映像素材、電子メール、プレゼンテーション、インターネットコンテンツ、放送コンテンツ、商用レポート、音声素材、グラフィック及びクリップアート、写真など、又はこれらのいずれもの組合せ又は合成を含む。この説明においては、テキスト情報項目に言及する。テキスト情報項目は、非テキスト項目と関連付けられても、又はリンクされてもよい。したがって、例えば、音声及び/又は映像素材は、テキスト用語においてその素材を定義するテキスト情報項目である「メタデータ（Metadata）」と関連付けられることが可能である。

30

【0020】

情報項目は、従来の方でディスク記憶装置30にロードされる。好ましくは、これらの情報項目は、項目の検索及び索引付けをより容易にすることを可能にするデータベース構造の一部として記憶されるが、これは絶対的ではない。情報及び項目が一旦このように記憶されると、検索を行うためにこれらを配置するために用いられるプロセスは図2に概略的に示される。

40

【0021】

索引付けされた情報項目は、ディスク記憶装置30上に記憶される必要がないことが理解されるであろう。情報項目は、ネットワーク50を介して情報記憶及び検索システム（汎用コンピュータ）10に接続される外付けのリモートドライブ上に記憶されることが可能である。あるいは、情報は、例えば、インターネット中の様々なサイトに分散されて記憶されてもよい。情報が異なるインターネット又はネットワークサイトに記憶される場合、情報記憶の第2のレベルは、遠隔情報への「リンク（link）」（例えば、ユニバーサルリソースインジケータ：URI）をローカルに記憶するために用いられることが可能であり、そのリンクに関連付けられた関連した概要、要約又はメタデータを有する可能性がある。したがって、ユーザが関連リンクを選択しない（例えば、以下に説明する結果一覧領

50

域 260 から) 限り、遠隔的に保持された情報はアクセスされないが、以下の技術的な説明のために、遠隔的に保持された情報又は要約 / 概要 / メタデータあるいはリンク / URI は、「情報項目 (information item)」として考慮することが可能である。

【0022】

言い換えれば、「情報項目 (information item)」の形式的な定義は、特徴ベクトルが導かれ処理されて (以下を参照)、SOM へのマッピングを提供する項目である。結果一覧領域 260 (以下を参照) に示されるデータは、ユーザが検索する実際の情報項目 (これがローカルに保持され、好都合な表示を行うのに十分短い場合) であっても、又は 1 つ又はそれ以上のメタデータ、URI、要約、一組のキーワード、代表的なキースタンプ画像などの情報項目を表現する及び / 又は指示するデータであってもよい。これは、常にではないが、一組の項目を表現するデータの一覧表示を含むことが多い、動作「一覧 (list)」に固有である。

10

【0023】

別の例において、情報項目は、研究チーム又は法律事務所などのネットワーク化された作業グループを通じて記憶されることが可能である。複合的な手法は、ローカルに記憶された幾つかの情報項目及び / 又はローカルエリアネットワークに亘って記憶された幾つかの情報項目及び / 又は広域ネットワークに亘って記憶された幾つかの情報項目を包含し得る。この場合、情報検索及び検索システムは、例えば、大規模な多国間研究開発組織における、他人による同様の作業の位置指定において有用であることが可能であり、同様な研究作業は、SOM (以下を参照) 中の同様な出力ノードにマッピングされる傾向にある。あるいは、新しいテレビ番組が計画中である場合、この技術は、同様の内容を有する以前のプログラムを検出することによってその独自性をチェックするためにも用いられることが可能である。

20

【0024】

図 1 の情報記憶及び検索システム (汎用コンピュータ) 10 は、索引付けされた情報項目を有することが可能なシステムの一例でしかないことも理解されるであろう。初期 (索引付け) 段階は、適度に強力なコンピュータ、最も可能性が高くは、非ポータブルコンピュータによって実行され、情報へアクセスするというその後の段階は、「パーソナルデジタルアシスタント (personal digital assistant: PDA) (概して片手に入る、表示装置及びユーザ入力装置を有するデータ処理装置)」などのポータブルマシン、ラップト

30

【0025】

プロセスは、特定数の情報項目に限定されない。

【0026】

情報項目の自己組織化マップ (SOM) 表現を生成させるプロセスを、図 2 ~ 図 6 を参照して説明する。図 2 は、所謂「特徴抽出 (feature extraction)」プロセスに次いで SOM マッピングプロセスを図示する概略的なフローチャートである。

40

【0027】

特徴抽出は、生データを抽象表現に変換するプロセスである。次いで、これらの抽象表現は、パターン分類、クラスタリング及び認識などのプロセスに用いられる。このプロセスにおいて、所謂「特徴ベクトル (feature vector)」が生成されるが、これは、文書内で用いられる用語の頻度を表す抽象表現である。

【0028】

特徴ベクトルの作成による視覚化形成プロセスは、以下を含む。

- ・用語の「文書データベース辞書 (document database dictionary)」の作成
- ・「文書データベース辞書」に基づく各個々の文書についての「用語頻度ヒストグラム (term frequency histogram)」作成

50

・ランダムマッピングを用いた「用語頻度ヒストグラム (term frequency histogram) 」の縮小

・情報空間の二次元視覚化の作成

これらのステップをより詳細に検討すると、各文書 (情報項目) 1 0 0 が順に開かれる。ステップ 1 1 0 で、全ての「ストップワード (stop word) 」が文書から除去される。ストップワードとは、「a」、「the」、「however」、「about」、「and」及び「the」などの、前もって作成された一覧にある非常に一般的な単語である。これらの単語は非常に一般的であるので、これらは、概して、十分な長さの全ての文書において同様の頻度で出現する傾向にある。このため、これらの単語は特定の文書の内容を特徴付ける試みにおいてほとんど効果がなく、したがって、除去されるべきである。

10

【 0 0 2 9 】

ストップワードの除去後、ステップ 1 2 0 で残りの単語の語幹分析がされるが、これは単語の変形の共通語幹を見出すことである。例えば、「thrower」、「throws」及び「throwing」は、共通語幹「throw」を有する。

【 0 0 3 0 】

文書中に出現する語幹分析された単語 (「ストップ (stop) 」ワードを除く) の「辞書 (dictionary) 」が維持される。新たな単語に遭遇すると、この単語は辞書に加えられ、文書コレクション全体 (情報項目の集合) においてその単語が出現した回数の実行カウントも記録される。

【 0 0 3 1 】

結果として、集合内の中の全ての文書において用いられる用語をそれらの用語が現れる頻度と共に示した一覧が得られる。余りにも高い又は低い頻度で現れる単語は度外視され、これはすなわち、これらの単語が辞書から除去され、続いて行われる分析には加わらないということである。余りにも低い頻度で現れる単語は綴り間違いであるか、造語であるか、あるいは文書の集合によって表される分野に関連しないかである可能性がある。余りにも高い頻度で現れる単語は、集合の中の文書を区別するためには余り適切ではない。例えば、用語「News」は、放送に関連する文書の試験集合中の総文書の約 3 分の 1 の率で用いられるが、用語「football」は、その試験集合中の文書の約 2 % でしか用いられない。したがって、「football」は「News」よりも文書内容を特徴付けるためにより良い用語であると仮定することができる。逆に、単語「fottball」 (「football」の綴り間違い) は、文書の集合全体において一度しか現れず、したがって、出現が余りにも少ないために度外視される。このような単語は、平均出現頻度から 2 標準偏差を引いた (- 2) 値よりも低い、又は平均出現頻度に 2 標準偏差を足した (+ 2) 値よりも高い出現頻度を有する単語として定義されることが可能である。

20

30

【 0 0 3 2 】

次いで、特徴ベクトルがステップ 1 3 0 で生成される。

【 0 0 3 3 】

これを行うために、集合中の各文書について用語頻度ヒストグラムが作成される。用語頻度ヒストグラムは、辞書 (その文書の集合に属したもの) に存在する単語が個々の文書内で出現する回数をカウントすることによって構成される。辞書中の用語の大半が 1 つの文書中に存在することはないために、これらの用語は頻度ゼロを有する。2 つの異なる文書についての用語頻度ヒストグラムの概略的な例を、図 3 a 及び図 3 b に示す。

40

【 0 0 3 4 】

この例から、ヒストグラムが文書内容をどのように特徴付けるかがわかる。これらの例を検討することによって、文書 1 では文書 2 よりも用語「MPEG」及び「Video」の出現回数が多く、文書 2 自体は用語「MetaData」の出現がより多い。対応する単語が文書中に存在しないので、ヒストグラム中の見出し項目の多くはゼロである。

【 0 0 3 5 】

現実の例においては、実際の用語頻度ヒストグラムは、例におけるよりも大幅に多い数の用語を有する。代表的には、ヒストグラムは 5 0 0 0 0 を超える異なる用語の頻度をブ

50

ロットし得、50000を超える大きさをヒストグラムに与える。このヒストグラムの大きさは、SOM情報空間の構成に用いられる場合には、大幅に縮小される必要がある。

【0036】

用語頻度ヒストグラム中の各見出し項目は、その文書を表す特徴ベクトル中の対応する値として用いられる。このプロセスの結果として、文書コレクション中の各文書についての辞書によって特定される全ての用語の頻度を含む(50000×1)ベクトルが得られる。値の大半は代表的にはゼロであり、その他の値の大半が代表的には1などの非常に小さい数であるために、ベクトルは「スパース(sparse)」と称され得る。

【0037】

特徴ベクトルのサイズ、したがって、用語頻度ヒストグラムの大きさは、ステップ140で縮小される。ヒストグラムの大きさを縮小するプロセスには、2つの方法が提案される。

【0038】

i) ランダムマッピング：ヒストグラムが乱数行列によって乗算される技術である。これは、計算上安価なプロセスである。

【0039】

ii) 潜在意味的索引付け：文書内に同時に出現する可能性が高い用語のグループを探すことによって、ヒストグラムの大きさを縮小する技術である。次いで、これらの単語グループは、単一のパラメータに縮小されることが可能である。これは、計算上高価なプロセスである。

20

【0040】

本実施の態様における用語頻度ヒストグラムの大きさを縮小するために選択された方法は、上記で参照したカスキ(Kaski)の論文において詳細に説明されているような、「ランダムマッピング(random mapping)」である。ランダムマッピングは、乱数行列でヒストグラムを乗算することによって、ヒストグラムの大きさの縮小を達成する。

【0041】

上述のように、「生(raw)」の特徴ベクトル(図4aに概略的に図示)は、代表的には、50000個の値の領域におけるサイズを有するスパースベクトルである。これは約200のサイズ(図4bの概略図を参照)に縮小されることが可能であり、特徴ベクトルの相対的直交特性、すなわち、他の同様に処理された特徴ベクトルとの相対角度(ベクトル内積)などの関係を保持している。特定の直交ベクトル数は限られているが、略直交ベクトルの数は大幅に多いので、これは良好に働く。

30

【0042】

実際に、ベクトルの大きさが増大するに従って、ランダムに生成されたベクトルの任意の集合は互いにほぼ直交する。この特性は、この乱数行列によって乗算されたベクトルの相対方向が保持されることを意味する。これは、それらの内積を調べることによりランダムマッピングの前後のベクトルの類似性を示すことによって表されることが可能である。

【0043】

50000個の値から200個の値にスパースベクトルを縮小することによって、それらの相対的類似性が保持されることを経験的に示すことができる。しかし、このマッピングは完全なものではないが、文書の内容を簡潔に特徴付けるという目的のためには十分である。

40

【0044】

特徴ベクトルが文書コレクションについて生成されて、コレクションの情報空間を規定すると、これらの特徴ベクトルはステップ150で二次元SOMに投影されて、意味マップが作成される。以下の節では、コホネンの自己組織化マップを用いた特徴ベクトルのクラスタリングによる二次元へのマッピングのプロセスを説明する。説明するに当たり、図5も参照される。

【0045】

コホネンの自己組織化マップは、各文書について生成された特徴ベクトルをクラスタリ

50

ング及び組織化するために用いられる。

【0046】

自己組織化マップは、入力ノード170と、二次元平面185として図示されるノードの二次元配列又はグリッド中の出力ノード180とからなる。マップを調整するために用いられる特徴ベクトル中に存在する値と同数の入力ノードが存在する。マップ上の各出力ノードは、重み付けされた結合190（各結合について1つの重み）によって入力ノードに結合されている。

【0047】

初めに、これらの各重みが乱数に設定され、次いで、対話式プロセスによって重みが「調整（trained）」される。マップは、各特徴ベクトルをマップの入力ノードに与えることによって調整される。「最も近接した（closest）」出力ノードが、入力ベクトルと、各出力ノードに関連付けられた重みとの間のユークリッド距離を演算することによって算出される。

10

【0048】

入力ベクトルと、そのノードに関連付けられた重みとの間の最小ユークリッド距離によって識別される最も近接したノードは「勝者（winner）」と称され、このノードの重みは、入力ベクトルに「近接して（closer）」移動するように、重みの値をわずかに変えることによって調整される。勝利ノードに加えて、勝利ノードの近隣にあるノードも調整され、入力ベクトルにわずかに近づいて移動させられる。

【0049】

20

マップが一旦調整されると、ノードの二次元マップ内の入力空間のトポロジーの多くを保持することを可能にするのは、1つのノードの重みのみではなくマップ上のノード領域の重みも調整するこのプロセスである。

【0050】

マップが一旦調整されると、各文書がマップに与えられて、その文書についての入力特徴ベクトルにどの出力ノードが最も近接しているかを見ることが可能になる。重みが特徴ベクトルと同一である可能性は低く、特徴ベクトルとマップ上のその最も近接したノードとの間のユークリッド距離はその「量子化誤差（quantisation error）」として知られている。

【0051】

30

各文書についての特徴ベクトルをマップに与えて、それがどこに存在するかを見ることによって、各文書についてx及びyマップ位置が生じる。これらのx及びy位置は、文書IDと共にルックアップテーブルに入力されると、文書間の関係を視覚化するために用いられることができる。

【0052】

最後に、ディザ成分がステップ160で付加されるが、これを以下で図6を参照して説明する。

【0053】

上述したプロセスに起こる可能性のある問題は、2つの同一又は実質的に同一の情報項目が、SOMのノード配列中の同一ノードにマップされる可能性があることである。これによってデータの取扱いが困難になることはないが、これは表示画面（以下で説明する）上でのデータの視覚化を行う補助とはならない。特に、データが表示画面上で視覚化されると、複数の非常に類似した項目が特定のノードにある1つの項目に対して区別可能になるために有用であることがわかっている。したがって、各情報項目がマップされているノード位置に「ディザ（dither）」成分が付加される。ディザ成分は、ノード分離の±2分の1を無作為に付加することである。したがって、図6を参照すると、それについてマッピングプロセスが出力ノード200を選択する情報項目は、実際には、図6において点線で境界付けられた領域210内のノード200の周囲のいずれものマップ位置にマッピングされてもよいように付加されたディザ成分を有する。

40

【0054】

50

したがって、情報項目は、SOMプロセスの「出力ノード (output node)」以外のノード位置で図6の平面上の位置へマッピングすると考えられ得る。

【0055】

いずれもの時点で、上記で概説したステップ (すなわち、ステップ110から140) に従い、次いで、「前もって調整された (pre-trained)」SOMモデル、すなわち、マップの自己組織化作成の結果として生じるSOMモデルの集合に、結果として得られる縮小された特徴ベクトルを適用することによって、新しい情報項目がSOMに付加されることが可能になる。したがって、新たに付加された情報項目については、マップは概して「再調整 (retrained)」されないが、その代わりに、全てのSOMモデルが修正されていない状態でステップ150及び160が用いられる。新しい情報項目が付加される毎にSOMを再調整するのは計算上高価であり、マップ中の共通してアクセスされる情報項目の相対位置に慣れていく可能性があるユーザに幾分使いにくいものでもある。

10

【0056】

しかし、再調整プロセスが適切である時点も同様にあり得る。例えば、SOMが初めて作成されて以降、新しい用語 (恐らくは、新しいものの新しい項目又は新しい技術分野) が辞書に入力される場合、それらの用語が出力ノードの既存の集合に特に良好にはマップしないこともあり得る。これは、新たに受け取られた情報項目の既存のSOMへのマッピングの間に検出される、所謂「量子化誤差 (quantisation error)」の増加として検出される可能性がある。本実施の形態においては、量子化誤差は閾値誤差量と比較される。量子化誤差が閾値誤差よりも大きい場合、(a)元の情報項目の全て及びSOMが作成されてから付加されたいずれもの項目を用いて、SOMが自動的に再調整されるか、(b)好都合な時間に再調整プロセスを開始するようにユーザが促されるか、のいずれかが行われる。再調整プロセスは、全ての関連する情報項目の特徴ベクトルを用い、ステップ150及び160全体を再適用する。

20

【0057】

図7は、表示画面60上の表示を概略的に図示する。表示は、検索クエリ領域250、結果一覧領域260及びSOM表示領域270を示す。

【0058】

動作中は、始めはSOM表示領域270は空白である。ユーザは、キーワード検索クエリを検索クエリ領域250に入力する。次いで、ユーザは、例えば、キーボード70上の復改キーを押すことによって、又はマウス80を用いて表示画面の「ボタン (button)」を選択して検索を開始する。次いで、標準キーワード検索技術を用いて、検索クエリ領域250内のキーワードがデータベース中の情報項目と比較される。これによって結果一覧が生成され、この各々が結果一覧領域260中の各々の見出し項目280として示される。次いで、SOM表示領域270は、各結果項目に対応する表示点を表示する。

30

【0059】

SOM表現を生成するために用いられる分類プロセスは、SOM中の相互類似情報項目をグループ化する傾向にあるので、検索クエリの結果として、概して、クラスタ290などのクラスタが生じる傾向にある。ここで、SOM表示領域270上の各点が、結果一覧領域260内の結果の1つと関連付けられたSOM内の各々の見出し項目に対応することと、SOM表示領域270内で点が表示されている位置が、ノード配列内のこれらのノードの配列位置に対応することとが特筆される。

40

【0060】

図8は、「ヒット (hits)」(結果一覧中の結果) 数を低減させるための技術を概略的に図示する。ユーザはマウス80を利用して境界線を引き、この境界線は、本例においては、SOM表示領域270中に表示される一組の表示点を囲む矩形ボックス300である。結果一覧領域260において、境界線300内の点に対応する結果のみが表示される。これらの結果が対象となるものではないことがわかると、ユーザは、表示点の異なる集合を囲む別の境界線を引く。

【0061】

50

結果一覧領域 260 は、境界線 300 内に表示される表示点についての結果のための一覧見出し項目を表示し、これは単語の検索クエリ領域 250 中の検索基準を満たしたことが特筆される。境界線 300 は、ノード配列中の母集団化されたノードに対応する他の表示位置を囲むことが可能であるが、これらが検索基準を満たさなかった場合、これらの表示位置は表示されず、したがって、結果一覧領域 260 内に示される結果のサブセットの一部を形成することはない。

【0062】

図 9 は、本発明の 1 つの実施の形態を図示する。

【0063】

図 9 を参照すると、ステップ 920 で自己組織化マップ SOM が作成されるとき、これはラベルを有さない（コホネンの SOM とは異なる）。ユーザは、マップを探求するための誘導を与えるためにラベルを必要とする、本発明の実施の形態においては、ラベルは、ユーザの特定の必要性を満たすために自動的に作成される。ユーザは、図 7 及び / 又は図 8 を参照して説明したように、検索の結果一覧を作成する。ラベルは、結果に従って自動的にかつ動的に作成され、SOM 表示領域 270 内の表示点のクラスをラベル付けするために用いられる。

【0064】

クロスクラスタ関連付け / 補助キーワード検索

本発明の実施の形態の例を図 10、図 11 及び図 12 を参照して説明する。

【0065】

図 10 において、情報項目のデータベースを含むデータ格納装置 400 は、データ通信ネットワーク 410 によって検索プロセッサ 404 及びマッピングプロセッサ 412 に接続される。マッピングプロセッサ 412 は、ユーザ制御装置 414 及びディスプレイプロセッサ 416 に接続される。ディスプレイプロセッサ 416 の出力は、グラフィカルユーザインタフェース 418 に供給され、グラフィカルユーザインタフェース 418 はディスプレイ 420 に接続されている。ディスプレイプロセッサ 416 は、表示画面上で表示を行うために、マッピングプロセッサ 412 からのデータを処理するように動作可能である。

【0066】

データ格納装置 400 は、マッピングプロセッサ 412 とは別に配置されることが可能である。それによって、検索プロセッサ 404 は、データ格納装置 400、マッピングプロセッサ 412、並びにディスプレイプロセッサ 416、グラフィカルユーザインタフェース 418 及びディスプレイ 420 である、情報を表示するために用いられる、図 10 に示される構成要素とは別に配置されることが可能である。あるいは、マッピングプロセッサ 412、検索プロセッサ 404 及びディスプレイプロセッサ 416 は、図 1 に示されるような汎用コンピュータ 10 上で実行するために、ソフトウェアモジュールの形態で実施されてもよい。したがって、マッピングプロセッサ 412、検索プロセッサ 404 及びディスプレイプロセッサ 416 は別々に製造及び配置されることが可能であることが理解されるであろう。

【0067】

図 10 に示される実施の形態は、図 7、図 8 及び図 9 における図と組み合わせられた、図 1 に示されるような情報記憶及び検索システムと実質的に同様に動作する。図 7、図 8 及び図 9 は、検索クエリに対してどのように情報項目が検索されるか及び検索結果がどのように表示されるかの図示例を提供する。したがって、図 10 に示される実施の形態は、検索クエリ、例えば、ユーザ制御装置 414 からキーワードを受け取るように構成される。キーワードが受け取られると、検索プロセッサ 404 によって検索が実行されて、検索結果として識別される情報項目に対応する配列中の x 及び y 位置の組をマッピングプロセッサ 412 との組合せで識別する。例えば、ノードの 40 × 40 の配列については、正方形の二次元配列中に 1600 個の位置が存在する。上記で説明したように、検索プロセッサ 404 は、検索クエリに従って情報項目を検索する。検索プロセッサ 404 による検索

10

20

30

40

50

によって、検索クエリに対応するものとして検索プロセッサ 404 によって識別された情報項目についての x 及び y 位置の組が得られる。検索結果の x 及び y 位置は、マッピングプロセッサ 412 によって受け取られる。

【0068】

マッピングプロセッサ 412 は、k 平均 (k-means) クラスタリングプロセスを行うことによって、第 1 の大域レベル (global level) での情報項目のクラスタを識別するように動作可能である。k 平均クラスタリングプロセスは、クラスタ及び配列内のクラスタの位置を識別する。k 平均クラスタリングプロセスは、クリストファー・エム・ビショップ (Christopher M. Bishop) による「パターン認識のためのニューラルネットワーク (Neural Networks for Pattern Recognition)」と題された書籍、第 187 ~ 188 頁、オックスフォード大学出版 (Oxford University Press) に開示されている。k 平均クラスタリングアルゴリズムの更なる開示は、ウェブアドレス <http://cne.gmu.edu/modules/dau/stat/clustgalgs/clust5bdy.html> に開示されている。

10

【0069】

図 11 に図示されているように、キーワード「show」についての検索の結果によって、それらのメタデータの一部として単語「show」を有する情報項目に対応する配列中の位置が識別されることが可能である。したがって、配列に k 平均クラスタリングアルゴリズムを行った結果、例えば、「quiz」「game」及び「DIY」である情報項目の 3 つのクラスタが識別される。情報項目のこれらのクラスタは、第 1 の階層レベル H レベル 1 を形成する。ディスプレイプロセッサ 416 は、第 1 の階層レベル H レベル 1 の情報項目のクラスタリングに対応するデータをマッピングプロセッサ 412 から受け取る。ディスプレイプロセッサ 416 は、この第 1 の階層レベル H レベル 1 の二次元表示を表すデータを提供するように、データの第 1 の階層レベルを処理する。ディスプレイプロセッサ 416 によって生成されたデータは、図 12 に示されるように、ディスプレイ 420 上の第 1 の表示領域 430 において表示を行うためにグラフィカルユーザインタフェース 418 に与えられる。

20

【0070】

幾つかの実施の形態においては、k 平均クラスタリングアルゴリズムを用いてクラスタの識別を更に精密にするために、マッピングプロセッサ 412 によって更なる動作が行われることが可能である。更なる動作は、「k 平均クラスタリング及び剪定 (k-means clustering and pruning)」と称される。公知の k 平均クラスタリングプロセスは、類似した情報項目を示す検索結果において識別される情報項目について、配列位置のグループを識別する。次いで、結果項目の x 及び y 位置の隣接するサブクラスタが同一のメインクラスタの一部であるかを決定する更なる剪定プロセスが行われる。2 つのサブクラスタの中心間の距離が閾値よりも小さい場合、これらの 2 つのサブクラスタは、同一のメインクラスタの一部であると考えられる。剪定は、クラスタが安定するまで、公知の方法で対話式に行われる。

30

【0071】

マッピングプロセッサ 412 は、第 1 の階層レベル H レベル 1 で識別された情報項目の各クラスタの更なる分析を行うように動作する。情報項目のクラスタを個々に検討し、かつ、それらの情報項目内で更なるクラスタを識別する機能をユーザに提供するために、マッピングプロセッサ 412 は更なる階層レベルを形成する。したがって、情報項目の各クラスタについて、情報項目のその第 1 の階層レベル内の更なるクラスタを識別するために、k 平均クラスタリングアルゴリズムがそのクラスタについて行われる。したがって、例えば、図 11 に図示されるように、k 平均クラスタリングアルゴリズムが「quiz」クラスタに行われると、3 つの更なるクラスタが第 2 の階層レベル H レベル 2 で識別される。

40

【0072】

第 1 の階層レベルについて図示されたように、各クラスタはキーワードに従ってラベル付けされる。キーワードは、クラスタ内の各情報項目が有する、その情報項目と関連付けられたメタデータ内に存在する最も共通する単語を見出すことによって識別される。した

50

がって、例えば、第1の階層レベルにおいて、単語「quiz」、「game」及び「DIY」によって3つのクラスタが識別される。

【0073】

第1の階層レベルHレベル1のクラスタのラベル付けに対応した方法で、第2の階層レベルHレベル2における各クラスタについてキーワードが識別される。したがって、これらの3つのクラスタは、「the chair」「wipeout」及び「enemy within」とラベル付けされる。これらの3つのクラスタの各々が、quiz showの異なるエピソードを含む。

【0074】

理解されるように、各クラスタの分析の更なる反復を行うことができる。これは、第2の階層レベルHレベル2で識別される各クラスタにk平均クラスタリングアルゴリズムを行うことによって達成される。図11に図示されるように「wipeout」情報クラスタは、k平均クラスタリングアルゴリズムを用いて更に分析される。しかし、第3の階層レベルHレベル3では、個別情報項目のみが明らかにされるために、図11に図示されるように、第3の階層レベルHレベル3は、「wipeout」の個々のエピソードを識別する。

【0075】

したがって、マッピングプロセッサ412は、異なる階層レベルで情報項目のクラスタを識別するように動作可能である。各階層レベルを表すデータが、ディスプレイプロセッサ416に与えられる。したがって、グラフィカルユーザインタフェース418と組み合わせられると、例えば、第2の階層レベルHレベル2に対応する可能性がある第2の領域がディスプレイ420上に表示されることが可能である。したがって、ズームコントロールを用いて、ユーザは第1の階層レベルHレベル1で表示されるクラスタにズームし得る。ズームコントロールは、ユーザ制御装置414を用いて動作させられることが可能である。したがって、特定のクラスタへズームすることで、情報項目の第2の階層レベルHレベル2を現す効果を有することができる。あるいは、第1の表示領域430内の「現在の目視」領域を選択するためにユーザ制御装置414を用いてもよい。したがって、第1の表示Hレベル1において示される第1の階層レベルで識別される「quiz」クラスタ内で識別されるクラスタに対して、第2の表示が行われる。

【0076】

本発明の実施の形態によって提供される更なる利点は、第2の又はそれに続く領域において表示される第2の又はそれに続くレベルに、他のクラスタの標識が与えられ得る構成である。標識は、より低い階層レベルで目視されるクラスタと関連付けられたキーワードに対する代替的なクラスタにユーザを導く。したがって、第2の表示領域440内でより低い階層レベルで図示されているクラスタは、目視されているクラスタに対する代替的なクラスタを有する。例えば、図12において、第1の表示領域430内で、第1の階層レベルは、「quiz」、「game」及び「DIY」の3つのクラスタを示す。ズームコントロールは「quiz」クラスタにズームするために用いられるので、第2の表示領域440は、「the chair」、「enemy within」及び「wipeout」である、「quiz」クラスタ内のクラスタの表示を与える。しかし、「quiz」クラスタに対する代替的なキーワードは、第1の表示領域430において図示されるように「DIY」、「horror」及び「game」である。したがって、矢印444、446及び448は、第2の表示領域440において表示されている「quiz」クラスタと同一の階層レベルにある情報項目のクラスタにユーザを導くために与えられる。したがって、次いでユーザが第1の階層レベルから異なるクラスタを閲覧して、第2の階層レベルにおけるクラスタを現すことを望む場合、ユーザは第1の階層レベル内の代替的なクラスタにナビゲートするために矢印を使用することができる。さらに、有利なことに、矢印は、第1の階層レベルで現れるクラスタについてのキーワードラベルでラベル付けされる。他の実施の形態において、クラスタ内の相対数の項目の図示をユーザに与えるために、この数は、方向を指示する矢印と関連付けられたキーワードと並んで示される。ユーザコントロール及びディスプレイは、マウスポインタMPが指示矢印上を通過する、又はその上に位置付けられると、この数を指すように配置されることが可能である。

【 0 0 7 7 】

幾つかの実施の形態の更なる有利な特徴は、付加的なキーワードの一覧、すなわち、第 1 のレベルのクラスタ内の第 2 のレベルのクラスタと関連付けられたキーワードを提供することである。クラスタリングについて図 1 2 において図示されるように、「horror」の更なる第 1 のレベルのクラスタを提供することによって、マウスポインタ M P が「horror」と関連付けられた矢印上に位置付けられると、その第 1 のレベルのクラスタ「horror」内の第 2 のレベルのクラスタに対応する付加的な単語が生じる。その結果、ユーザには、第 1 のレベルのクラスタを第 2 の表示領域 4 4 0 内で目視する必要なく、これらのクラスタと関連付けられた情報項目の内容の非常に有効な図示が与えられる。図 1 2 に図示されるように、表示領域は、第 1 の表示領域 4 3 0 内に出現する情報項目を検覧するため、及びそれらの周囲をナビゲートするための両方に用いられる、概して 4 5 0 で示されるコントロールアイコンを更に含むことが可能である。

10

【 0 0 7 8 】

マルチモード絞込み検索

本発明の別の実施の形態の例を、図 1 3 ~ 図 1 7 と組み合わせて図 1 0 を参照して説明する。図 1 3 は、情報項目と関連付けられて記憶されている特徴付け情報特徴のタイプを図示したものを示す。例えば、情報項目は、テレビ番組からの音声 / 映像データの一部であることが可能である。本例においては、番組はサッカーの試合のハイライトを提供する。したがって、データ項目は、映像データ 4 6 0 及び音声データを含む。音声データと関連付けられているのは、ボックス 4 6 2 内に図示されている音声メタデータである。音声メタデータは、映像データと関連付けられた音声信号の内容及びタイプを示す。本例については、音声データは「音楽 (music)」、「コメンタリ (commentary)」及び「群衆の騒音 (crowd noise)」を含むが、音声信号のタイプを示すメタデータの 1 つ又はそれ以上の他のタイプを含むことが可能である。映像データ及び音声データに加えて、情報項目は、映像及び音声データの内容又は属性を記載する他のメタデータも含むことが可能である。本例については、メタデータは、ボックス 4 6 4 内に図示されており、映像番組の内容の説明を含むことが示されている。S O M が作成される元となる特徴ベクトルを構築するために用いられるのは、このメタデータに含まれる単語である。しかし、本発明の他の実施の形態において、データ格納装置 4 0 0 に含まれる情報項目の集合に、音声メタデータ 4 6 2 である音声データに対する、又は映像データに対する検索が行われることが可能である。この目的のために、映像データ 4 6 0 のフレームから代表キースタンプが生成されることが可能である。

20

30

【 0 0 7 9 】

代表キースタンプ R K S は、映像データの各フレームのカラーヒストグラムを形成することによって生成される。全ての又は選択された映像フレームについてのカラーヒストグラムは組み合わせられ、次いで正規化されて、図 1 3 において棒グラフ 4 6 6 として代表的な形態で図示される、複合カラーヒストグラムが作成される。次いで、複合カラーヒストグラムは、各映像フレームについてのカラーヒストグラムと比較される。各映像フレームについての各列の複合ヒストグラムの対応する列に対する距離を加算することによって、各フレームについてのカラーヒストグラムと複合カラーヒストグラムとの距離が決定される。複合カラーヒストグラムに対して最小距離を有するカラーヒストグラムを有する代表キースタンプ R K S が選択される。次いで、したがって、サッカーの試合を表す番組については、作成された代表キースタンプは、サッカーの競技場の一部の映像画像である可能性が最も高く、これは図 1 3 に示される代表キースタンプによって図示される。

40

【 0 0 8 0 】

他の実施の形態において、R K S は、以下の方法のいずれかによって、各情報項目について映像フレームから作成されることが可能である。

・ユーザは、情報項目の内容全体に対応する最も代表的なフレームであると考えられるフレームを選択することが可能である。情報項目を主観的に表す映像フレームが選択されることをユーザが確実にするので、この方法によって信頼性が改善され得る。しかし、この

50

方法にはより時間がかかる。

- ・ユーザは、情報項目内の第1のフレーム又は無作為のフレームを選択することが可能である。これは、適切なRKSを選択するのには信頼性が低い方法である可能性がある。
- ・画像フレームの内容に基づいて映像フレームを処理し、RKSを選択する他の方法も考えられる。

【0081】

本発明の実施の形態によって、選択された特徴付け情報特徴に基づいて絞込み検索を生じさせる機能が提供され得る。1つの実施の形態において、検索プロセッサ404は、メタデータの項目、映像画像又は音声データのいずれかに関連付けられた一回目の検索において識別されたこれらの情報項目を検索するように動作可能である。代替的な実施の形態においては、検索は、メタデータのみ、映像データのみ、又は音声データのみ、あるいはそれらのいずれもの組合せに対して行われることが可能である。検索クエリの形成を容易にするために、図10に示されるディスプレイ420は、図14に示されるグラフィカルユーザインタフェース418によって与えられる更なるグラフィカルディスプレイを含んでいてもよい。

10

【0082】

図14において、表示領域472内の第1の行470は、メタデータに基づいてクエリ情報を選択する機能をユーザに与える。したがって、情報項目からの画像代表キースタンプがこの行のウィンドウ内に配置される場合、この情報項目と関連付けられたメタデータ（図13に図示されるように）が検索クエリに付加される。したがって、異なる情報項目からの1つ又はそれ以上の代表キースタンプが、タイプメタデータの特徴付け情報特徴についての検索クエリに導入されることが可能である。それに従って、第2の行474において、ユーザによって選択された映像フレームが導入され、検索クエリの一部が形成される。例えば、ユーザは映像データの特定の項目をブラウズし、対象となるフレームを選択することが可能である。次いで、ユーザは行474中にこの画像フレームを配置し、検索クエリの一部を形成させることが可能である。ユーザは、1つ又はそれ以上の映像フレームを導入することが可能である。

20

【0083】

ユーザは、検索される情報項目を、その情報項目内の音声データに従って選択することも可能である。したがって、表示領域476内の第3の行は、その情報項目の代表画像を導入して、検索クエリが、検索クエリにおけるその情報項目に対応する音声データを含むものであることを音声データについての行内で識別する機能をユーザに与える。

30

【0084】

特徴付け情報特徴のタイプに従って検索される情報項目を選択することに加えて、本発明の実施の形態は、選択された情報項目間でブール演算子に従って検索を行う機能も提供する。図14に図示されるように、メタデータ検索について選択される情報項目は、初めの2列478及び480の間に示されるような「AND」演算子に従って検索されるべきである。しかし、検索クエリにおける第1のメタデータと第1の映像画像項目検索クエリとの間の検索クエリは、「OR」演算子によって結合される。映像画像データについて検索される2つの項目は、「AND」演算子によって結合される。音声データに従って検索される情報項目もまた、「NOT」演算子に従って検索クエリにおいて検索されるものである。

40

【0085】

検索クエリを構築した後、検索プロセッサ404は、ユーザによって行われた選択によって構築された、図14に図示される検索クエリに従って、キーワード検索から識別された情報項目を検索するように動作可能である。検索プロセッサは、以下の節で説明されるように、選択された特徴付け情報特徴のタイプに依存して異なった方法で情報項目を検索する。

【0086】

メタデータなどの特徴付け情報特徴についての検索の例については、いずれもの情報項

50

目についても、メタデータから生成されるその情報項目についての特徴ベクトルが、その特徴ベクトルに対応する二次元配列内の点を識別するために使用されることが可能である。したがって、配列内のその識別された位置の所定距離内にある情報項目は、検索クエリの結果として戻されることが可能である。しかし、1つを超える情報項目がメタデータ検索行内で選択された場合、選択されたブール演算子に従ってこれらの項目の両方を検索するように、検索クエリが構築されなければならない。

【0087】

「AND」ブール演算子の例については、各情報項目についての特徴ベクトルが組み合わされて、図15に図示されるような複合特徴ベクトルを形成する。この目的のために、メタデータ内の各単語と関連付けられた値が加算され、正規化されて複合特徴ベクトルが作成される。したがって、図15に図示されるように、行470、列478~480並びにメタデータ検索クエリライン470で図示されるそれらの代表キースタンプを有する、ユーザが選択したメタデータと関連付けられた2つの特徴ベクトルAと特徴ベクトルBとが組み合わされて、特徴ベクトルCが形成される。次いで、検索プロセッサは特徴ベクトルCを取り上げ、これをSOMと比較することが可能である。複合特徴ベクトルCに対応する配列内の最も近い位置を識別した後、配列内のその識別された位置から配列内の所定数の位置内にある情報項目が検索クエリの結果として戻される。

10

【0088】

対応するメタデータ検索のブール「OR」演算子の例については、第1の特徴ベクトルA及び第2の特徴ベクトルBについて、これらの特徴ベクトルについての配列内の対応する位置が識別される。このように、検索クエリの結果として、配列内のこれらの識別された各点の所定数の位置内の全ての情報項目を戻すこととなる。これは図16及び図17に図示される。図17において、二次元配列内の、特徴ベクトルAに対応する位置及び特徴ベクトルBに対応する位置が識別される。図17に示されるように、特徴ベクトルA及びBについての配列位置の所定半径内の配列内の位置は、次いで、検索クエリの結果として識別されたものとして戻されることが可能である。しかし、更なる特徴ベクトルCが検索クエリで識別され、「NOT」ブール演算子がこの更なる特徴ベクトルについて指定される場合、特徴ベクトルCに対応する配列中の位置がここでも識別される。したがって、特徴ベクトルCからの配列位置の所定半径内の情報項目がここでも識別されることが可能である。しかし、「NOT」演算子の結果として、特徴ベクトルC並びに特徴ベクトルA及びBについての配列位置からの半径間で識別されるいずれもの相互的に包括的な配列位置が検索結果から排除される。したがって、検索プロセッサは、特徴ベクトルCからではなく特徴ベクトルA又はBから作成された配列内の位置に対応する情報項目を戻すように構成される。

20

30

【0089】

検索の特徴付け特徴である映像画像データに対応する検索クエリ中の第2列目について、検索プロセッサは、選択されたユーザ映像画像に対応する代表キースタンプについての映像データを検索するように動作可能である。この目的のために、ユーザが選択した映像画像と関連付けられたカラーヒストグラムは、情報項目と関連付けられた各代表キースタンプについてのカラーヒストグラムと比較される。各情報項目の代表キースタンプのカラーヒストグラムと、ユーザ指定の映像画像のカラーヒストグラムとの間の距離が算出される。これは、その画像の色成分を表す各列の間の距離を算出し、各列についてこれらの距離を合算することによって行われる。ユーザ選択映像画像のカラーヒストグラムと、その配列位置に対応する代表キースタンプのカラーヒストグラムとの間の距離が最小である情報項目に対応する配列位置が識別される。ここでもまた、クエリの結果として、識別された配列位置からの所定数の位置内の配列位置を有する情報項目が戻される。

40

【0090】

ブール演算子の場合について、ここでもまた、ブール「AND」演算子について選択及び指定された2つの画像についてのカラーヒストグラムを組み合わせることによって、カラーヒストグラムが形成されることが可能である。複合カラーヒストグラムの形成プロセ

50

スは、図 18 に図示される。図 14 に図示される表示領域内の映像画像検索クエリ行の行 474 並びに列 478 及び 480 において与えられる第 1 及び第 2 のユーザ選択画像についてのカラーヒストグラムは、カラーヒストグラムの各列内の値を平均化することによって組み合わせられる。したがって、図 18 a 及び図 18 b に図示される 2 つのカラーヒストグラムは組み合わせられて、図 18 c において形成されるカラーヒストグラムを形成する。検索される情報項目の代表キースタンプに対して検索されるのは、このカラーヒストグラムである。

【0091】

音声データの例については、検索プロセッサは、選択された情報項目と関連付けられた音声メタデータから特徴ベクトルを形成することが可能である。例えば、音声メタデータは、音声信号中に存在する高調波、スピーチデータ、又は音声メタデータによって表される音声信号内に音楽が存在するかを識別することが可能である。さらに、メタデータは、トニープレーアなどの特定の話し手又はジョンモトソンなどの特定の解説者が音声信号上に存在するかを識別することが可能である。したがって、ここでもまた、特に音声データと関連付けられる他の特徴ベクトルに対して検索されることが可能である選択された音声データから、特徴ベクトルが生成されることが可能である。上記の説明に対応した方法で、ブール演算子が、1 つを超える音声メタデータタイプについての検索を組み合わせるために用いられることが可能である。「AND」演算子の例については、音声メタデータ項目が組み合わせられて、複合メタデータ項目が作成されることが可能である。この複合項目に最も近い特徴ベクトルを有する対応する情報項目を検索することによって、情報項目が識別される。次いで、「OR」演算子が指定されると、検索プロセッサは、両方のメタデータ項目について配列内の所定数の位置の中にある情報項目を回復させることが可能である。ここでもまた、「NOT」ブール演算子は、検索クエリの結果から、一致する音声データを有する戻された情報項目を排除する機能を有する。

【0092】

識別された情報項目からの検索の絞込みについて、本発明の実施の形態が与えられた。しかし、他の実施の形態において、図 14 において図示されるディスプレイによって形成される検索クエリ、並びにメタデータ、映像データ及び音声データに対するその検索クエリの用途は、データ格納装置 400 内の情報の集合全体を検索するために与えられることが可能であることが理解されるであろう。

【0093】

関連検索

本発明の実施の形態の一例に従って上記で説明したように、図 14 に示されるグラフィカルユーザインタフェースを用いて構築された検索クエリによる情報項目は、検索クエリによって識別された特定の配列位置に隣接する項目を識別することによって検索されることが可能である。しかし、他の実施の形態例においては、どのような理由のためであっても、識別された情報項目から関連検索が行われることが可能である。しかし、代表的には、特定のキーワードによる検索によって、識別された情報項目の集合が得られる。これらの情報項目から、ユーザは、これらのうちの 1 つが特に対象となるものであることを決定することが可能である。次いで、関連検索によって、SOM によるこの情報と幾分かの相関を有する項目が与えられることが可能である。これは、例えば、対象とする情報項目に対応する配列位置から所定半径内にある、配列位置に対応する情報項目を識別することによって達成される。

【0094】

本発明の範囲から逸脱することなく、上述の実施の形態に様々な改変を行うことが可能である。本発明の様々な態様及び特徴は、添付の請求項に定義される。

【図面の簡単な説明】

【0095】

【図 1】情報記憶及び検索システムを概略的に示す図である。

【図 2】自己組織化マップ (SOM) の生成を示す概略的なフローチャートである。

10

20

30

40

50

【図 3】(a) 及び (b) は、用語頻度ヒストグラムを概略的に示す図である。

【図 4】(a) は生の特徴ベクトルを概略的に示す図であり、(b) は縮小された特徴ベクトルを概略的に示す図である。

【図 5】SOM を概略的に示す図である。

【図 6】ディザ処理を概略的に示す図である。

【図 7】SOM によって表現される情報へのアクセスのためのユーザインタフェースを提供する表示画面を概略的に示す図である。

【図 8】SOM によって表現される情報へのアクセスのためのユーザインタフェースを提供する表示画面を概略的に示す図である。

【図 9】SOM によって表現される情報へのアクセスのためのユーザインタフェースを提供する表示画面を概略的に示す図である。 10

【図 10】本発明の実施の形態による情報検索装置の概略的なブロック図である。

【図 11】検索において識別された情報項目の階層配置を図示したものである。

【図 12】図 11 において示される階層の異なるレベルを表示する 2 つの領域を提供する表示画面を概略的に示す図である。

【図 13】情報項目の例について 3 つのタイプの特徴付け情報特徴を図示したものである。

【図 14】本発明の実施の形態の例による検索クエリを形成するためのグラフィカルユーザインタフェースを概略的に示す図である。

【図 15】ブール AND 演算による複合特徴ベクトルの形成の概略図である。 20

【図 16】ブール OR 演算子による 2 つの特徴ベクトルと、ブール NOT 演算子による第 3 の特徴ベクトルとの組合せを示す図である。

【図 17】図 16 のブール演算子及び特徴ベクトルによる検索結果を示す識別された情報項目の二次元マップの一部を概略的に示す図である。

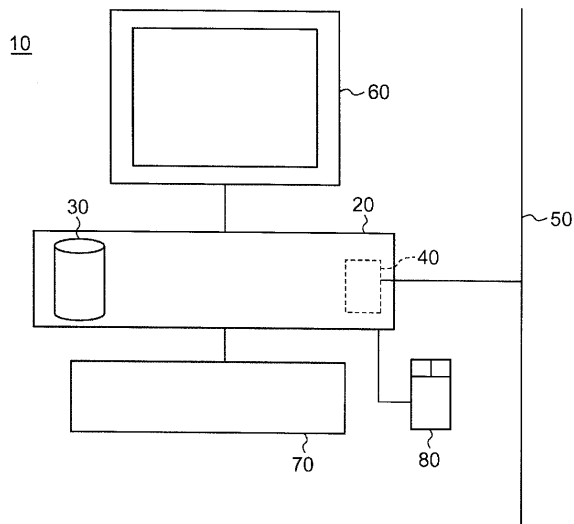
【図 18】(a) 及び (b) は、検索クエリを形成する 2 つの映像画像についてのカラーヒストグラムの 2 つの例を与える例示的な棒グラフであり、(c) は、(a) 及び (b) のカラーヒストグラムを組み合わせることによって作成される例示的な棒グラフである。

【符号の説明】

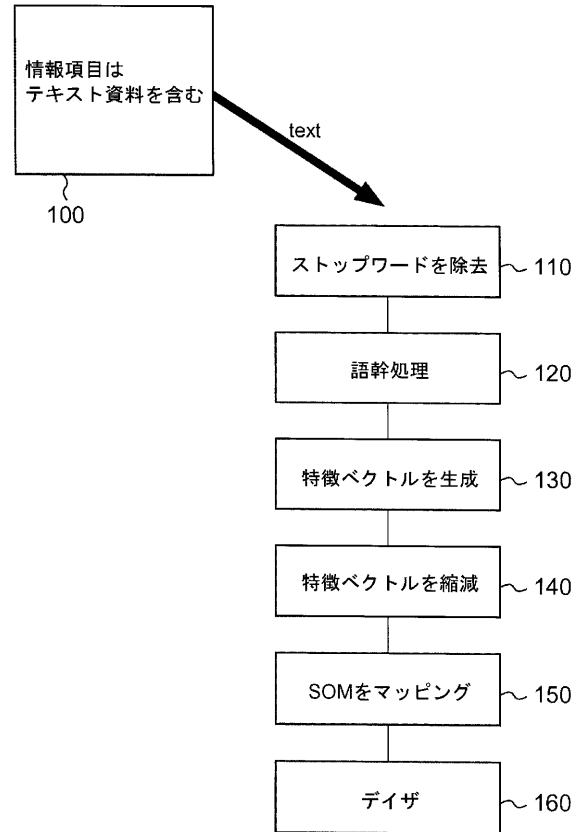
【0096】

10 汎用コンピュータ、20 プロセッサユニット、30 ディスク記憶装置、40 ネットワークインタフェースカード、50 ネットワーク、60 陰極線管表示装置、70 キーボード、80 マウス、400 データ格納装置、404 検索プロセッサ、410 通信ネットワーク、412 マッピングプロセッサ、414 ユーザ制御装置、416 ディスプレイプロセッサ、418 グラフィカルユーザインタフェース (GUI)、420 ディスプレイ 30

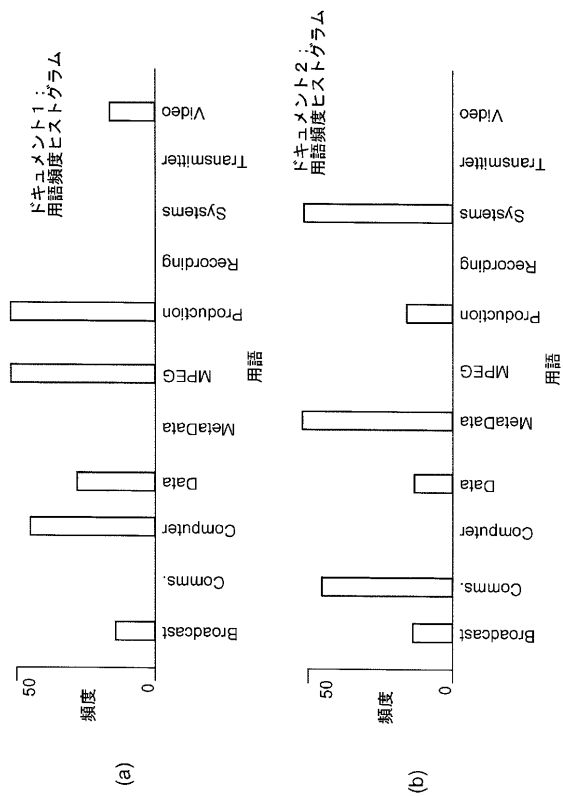
【図 1】



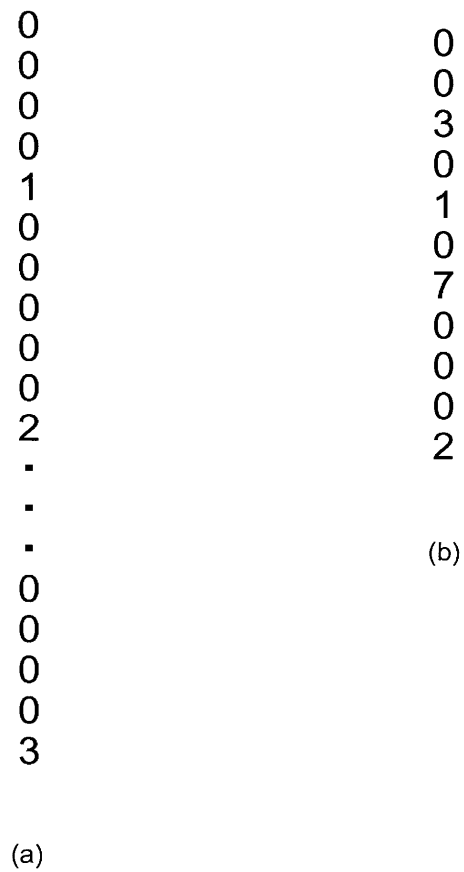
【図 2】



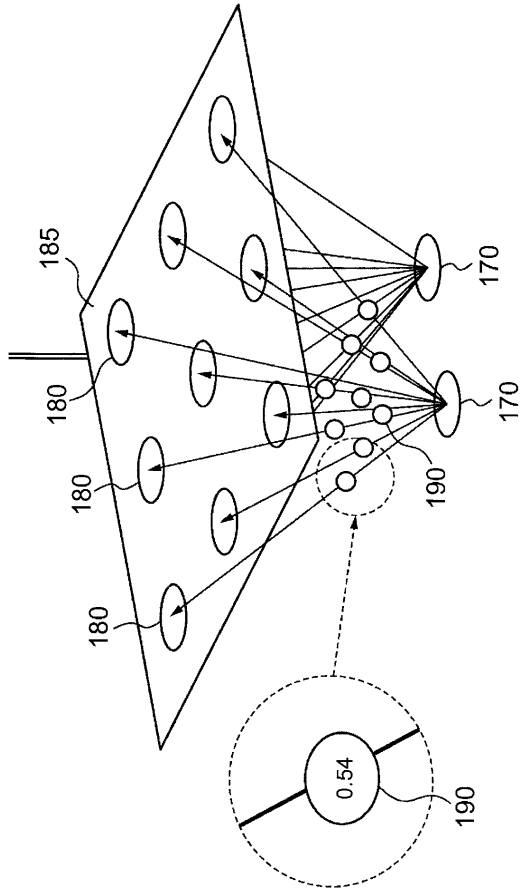
【図 3】



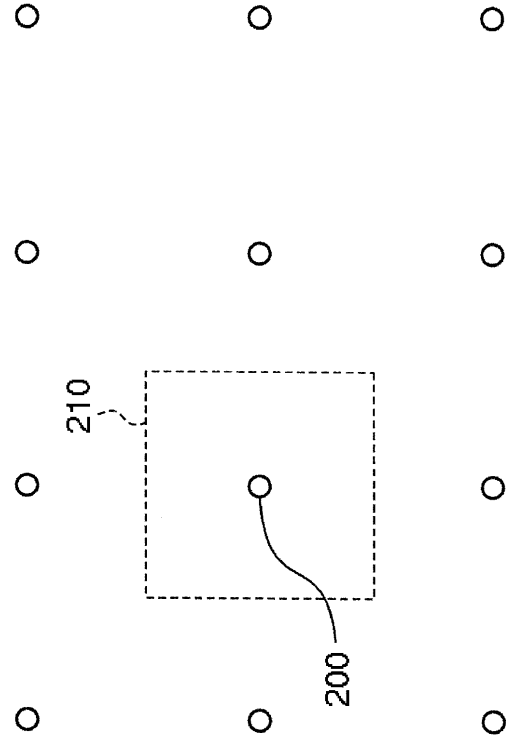
【図 4】



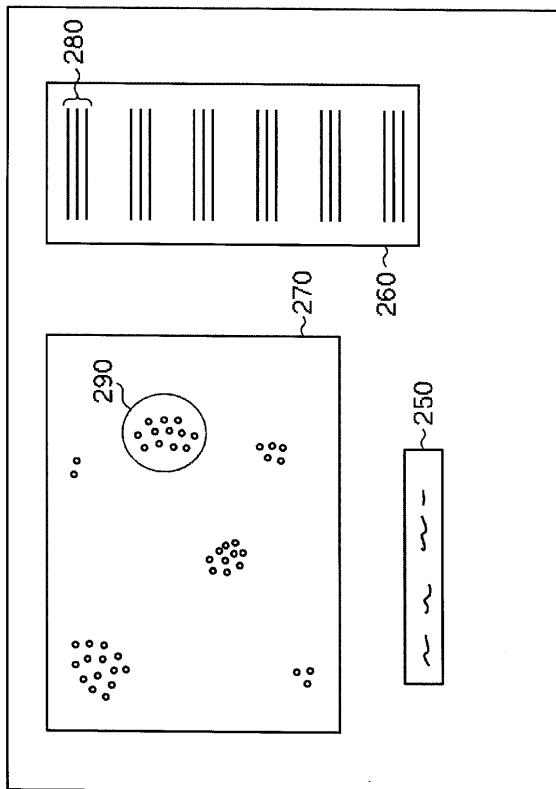
【図 5】



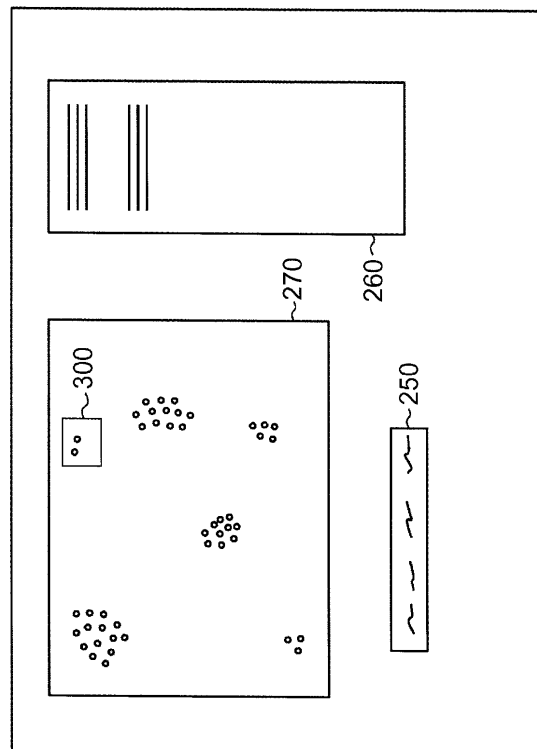
【図 6】



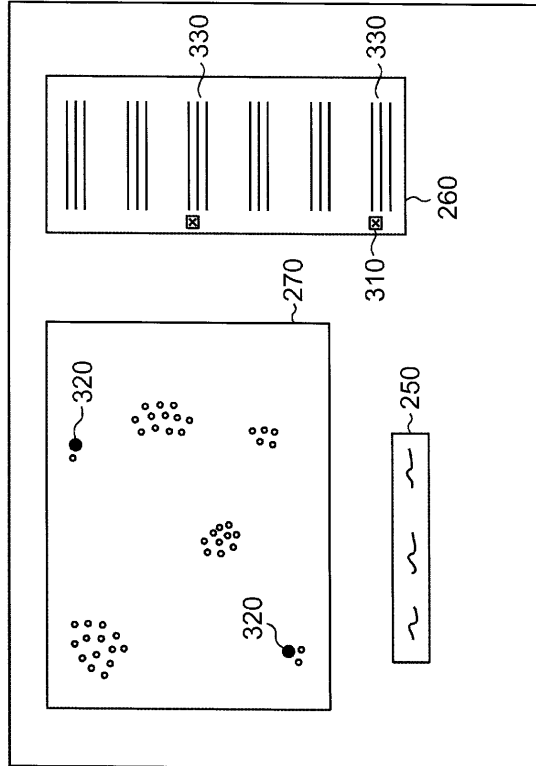
【図 7】



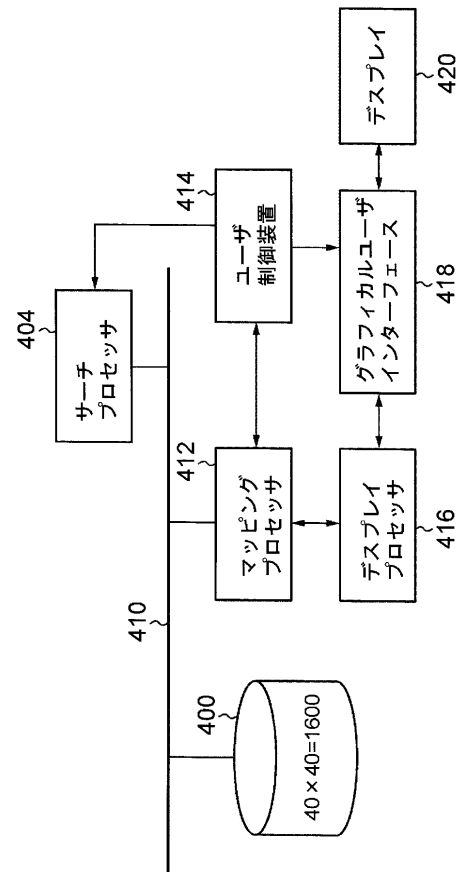
【図 8】



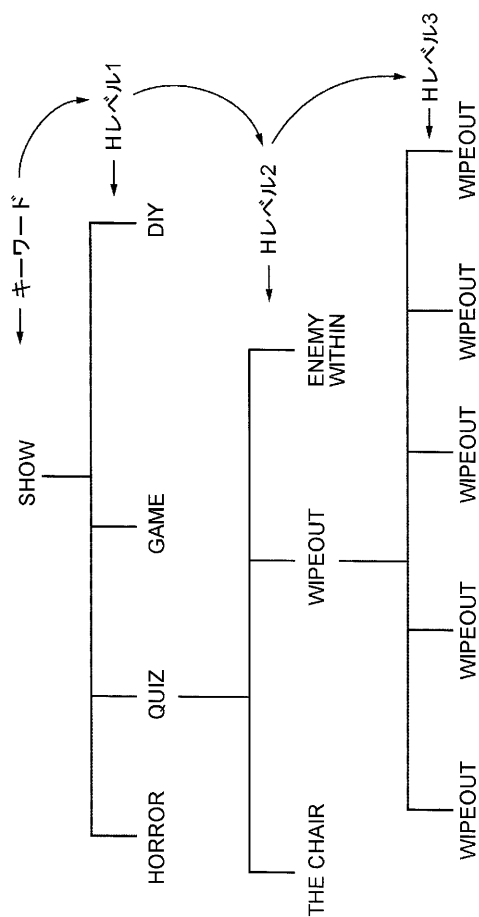
【図 9】



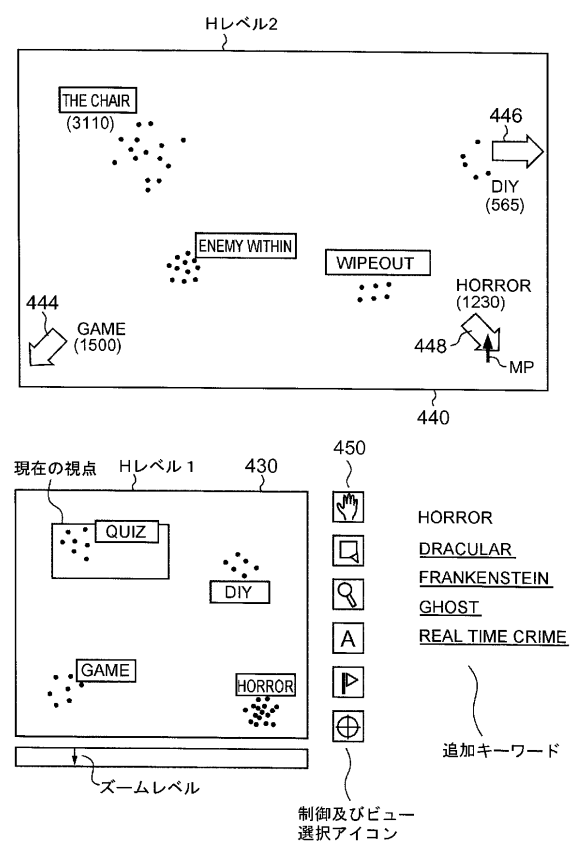
【図 10】



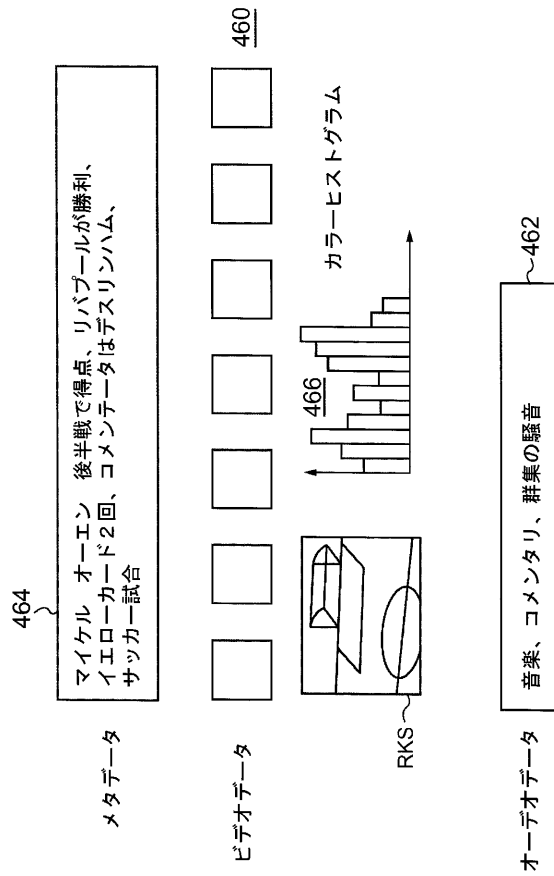
【図 11】



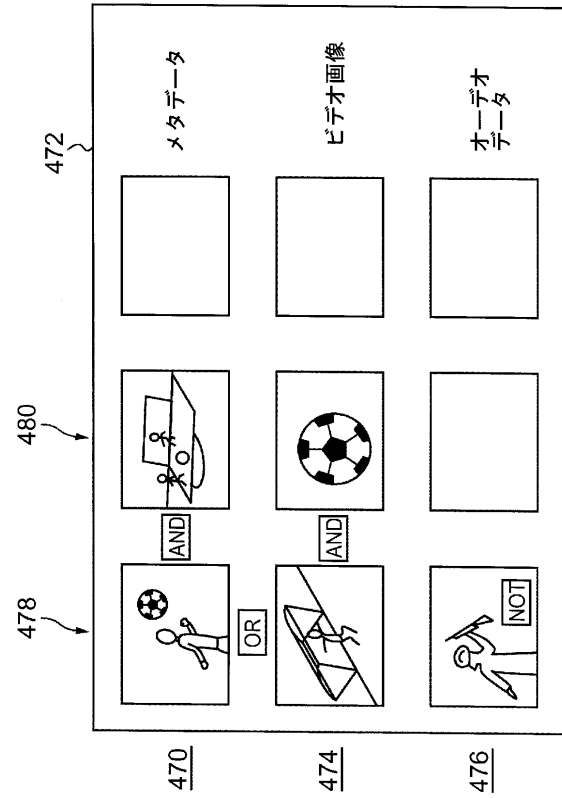
【図 12】



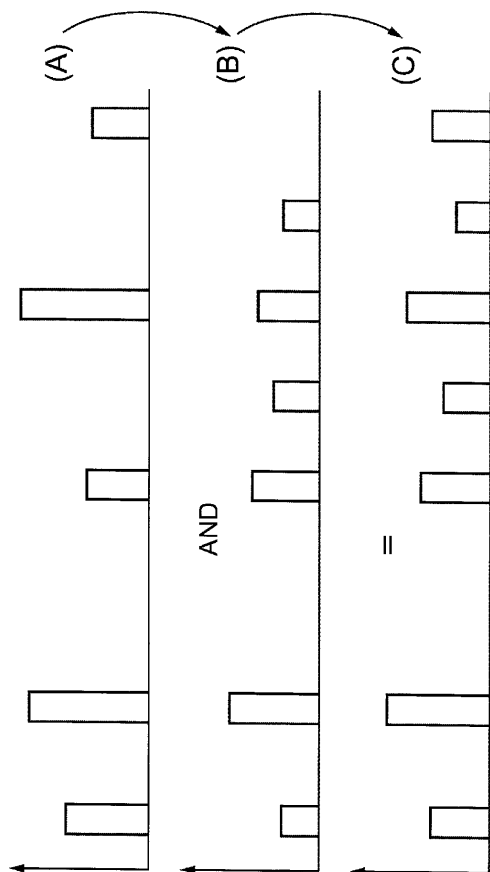
【図 13】



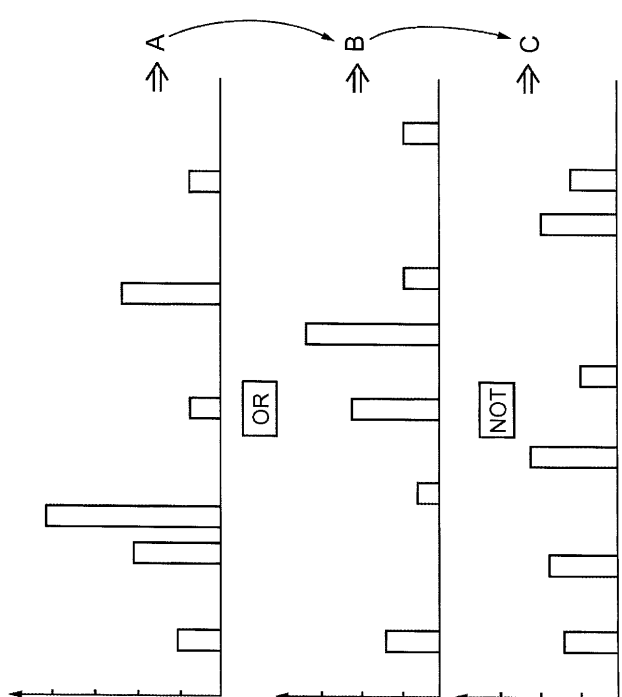
【図 14】



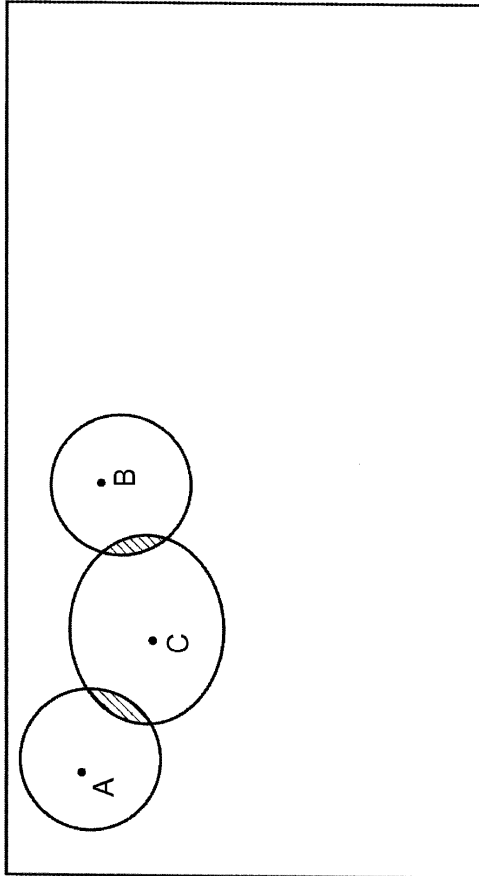
【図 15】



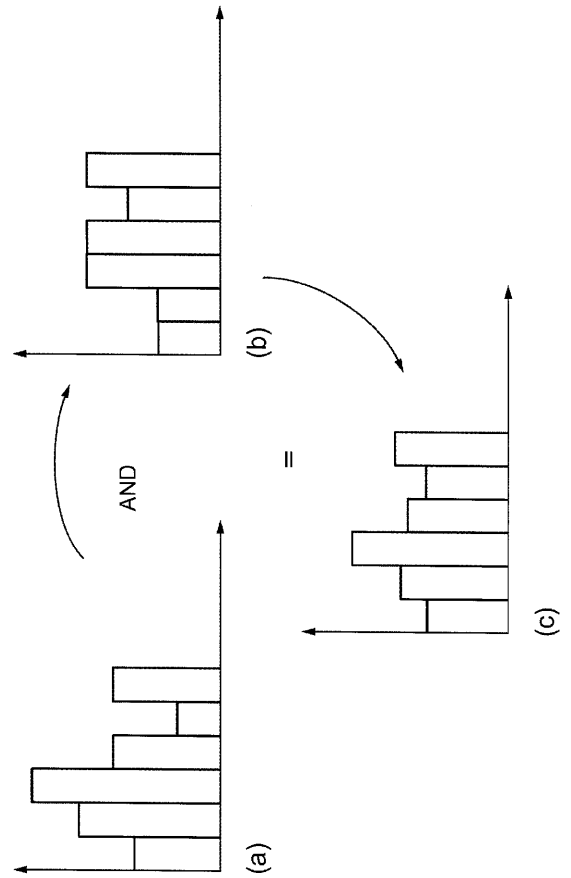
【図 16】



【図 17】



【図 18】



フロントページの続き

(72)発明者 トレペス、デヴィッド ウィリアム
イギリス国 K T 1 3 O X W サリー、ウェイブリッジ、ブルックランズ、ザ ハイツ (番地
無し) ソニー ユナイテッド キングダム リミテッド内

(72)発明者 ソープ、ジョナサン リチャード
イギリス国 K T 1 3 O X W サリー、ウェイブリッジ、ブルックランズ、ザ ハイツ (番地
無し) ソニー ユナイテッド キングダム リミテッド内

F ターム(参考) 5B075 NK43 NR12 PQ02 PQ13