

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関  
国際事務局

(43) 国際公開日  
2015年2月5日(05.02.2015)



(10) 国際公開番号  
WO 2015/016133 A1

- (51) 国際特許分類:  
G06F 17/30 (2006.01)
- (21) 国際出願番号: PCT/JP2014/069571
- (22) 国際出願日: 2014年7月24日(24.07.2014)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (30) 優先権データ:  
特願 2013-158200 2013年7月30日(30.07.2013) JP
- (71) 出願人: 日本電信電話株式会社(NIPPON TELEGRAPH AND TELEPHONE CORPORATION) [JP/JP]; 〒1008116 東京都千代田区大手町一丁目5番1号 Tokyo (JP).
- (72) 発明者: 岡野 靖(OKANO, Yasushi); 〒1808585 東京都武蔵野市緑町3丁目9-11 NTT 知的財産センター内 Tokyo (JP). 折原 慎吾(ORIHARA, Shingo); 〒1808585 東京都武蔵野市緑町3丁目9-11 NTT 知的財産センター内 Tokyo (JP). 佐藤 徹(SATO, Tohru); 〒1808585 東京

都武蔵野市緑町3丁目9-11 NTT 知的財産センター内 Tokyo (JP). 朝倉 浩志(ASAKURA, Hiroshi); 〒1808585 東京都武蔵野市緑町3丁目9-11 NTT 知的財産センター内 Tokyo (JP).

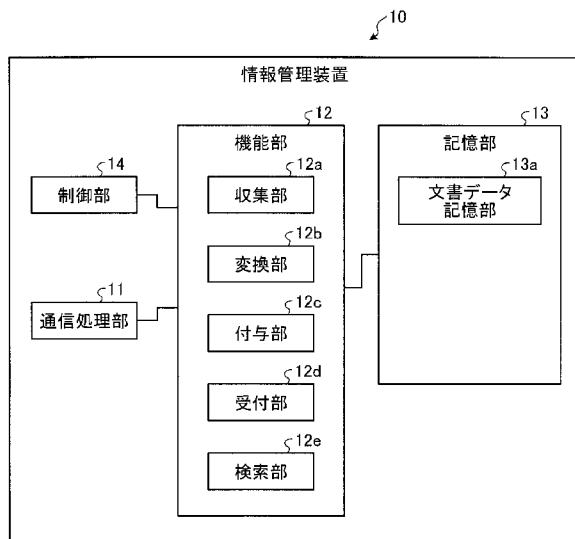
(74) 代理人: 酒井 宏明, 外(SAKAI, Hiroaki et al.); 〒1006020 東京都千代田区霞が関三丁目2番5号霞が関ビルディング 酒井国際特許事務所 Tokyo (JP).

(81) 指定国 (表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

[続葉有]

(54) Title: INFORMATION MANAGEMENT DEVICE, AND INFORMATION MANAGEMENT METHOD

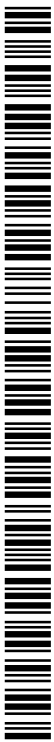
(54) 発明の名称: 情報管理装置及び情報管理方法



- 10 Information management device
- 11 Communication processing unit
- 12 Functional unit
- 12a Collection unit
- 12b Conversion unit
- 12c Assigning unit
- 12d Reception unit
- 12e Search unit
- 13 Storage unit
- 13a Document data storage unit
- 14 Control unit

(57) Abstract: An information management device (10) collects a plurality of pieces of document data on a network. The information management device (10) then classifies the respective pieces of document data according to predetermined fields, using words included in the collected document data, and assigns tag information corresponding to the fields to each piece of the document data. The information management device (10) further receives designation of a document data field as a search object. Thereafter, the information management device (10) searches for the document data to which the tag information corresponding to the received field is assigned.

(57) 要約: 情報管理装置(10)では、ネットワーク上における複数の文書データを収集する。続いて、情報管理装置(10)では、収集された各文書データに含まれる単語を用いて、各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与する。そして、情報管理装置(10)では、検索対象とする文書データの分野の指定を受け付ける。続いて、情報管理装置(10)では、受け付けられた分野に対応するタグ情報が付与された文書データを検索する。



WO 2015/016133 A1



(84) 指定国 (表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR),

OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

添付公開書類:

— 国際調査報告 (条約第 21 条(3))

## 明 細 書

**発明の名称**：情報管理装置及び情報管理方法

### 技術分野

[0001] 本発明は、情報管理装置及び情報管理方法に関する。

### 背景技術

[0002] 従来、インターネット上における不特定多数の文書を対象とし、特定の単語（キーワード）に関連する文書データを検索する手段として、様々なサーチエンジンが提供されている。例えば、サーチエンジンでは、キーワードの入力を受け付けることで、キーワードに関連する文書データを検索し、該文書データを出力する。

[0003] このようにサーチエンジンで検索された文書データから主要コンテンツのみを自動で抽出する技術が知られている（例えば、特許文献1参照）。また、検索対象となるインターネット上の複数の文書データを、類似する内容同士に分類する技術が知られている（例えば、特許文献2参照）。

### 先行技術文献

#### 特許文献

[0004] 特許文献1：特開2010-117941号公報

特許文献2：特許第4125951号

### 発明の概要

#### 発明が解決しようとする課題

[0005] しかしながら、従来の技術では、利用者が望むジャンルの文書データが全体の文書データの量に比べて少ない場合には、適切に文書データを検索できない場合があるという問題があった。例えば、セキュリティに関する記事のように、もともと話題が少ない記事を検索しようとした場合に、類似する記事や関連する記事を適切に検索することが困難であった。

[0006] そこで、この発明は、利用者が望むジャンルの文書データが全体の文書データの量に比べて少ない場合であっても、適切に文書データを検索すること

を目的とする。

### 課題を解決するための手段

[0007] 上述した課題を解決し、目的を達成するため、情報管理装置は、ネットワーク上における複数の文書データを収集する収集部と、前記収集部によって収集された各文書データに含まれる単語を用いて、前記各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与する付与部と、検索対象とする文書データの分野の指定を受け付ける受付部と、前記受付部によって受け付けられた分野に対応するタグ情報が付与された文書データを検索する検索部と、を有することを特徴とする。

[0008] また、情報管理方法は、情報管理装置によって実行される情報管理方法であって、ネットワーク上における複数の文書データを収集する収集工程と、前記収集工程によって収集された各文書データに含まれる単語を用いて、前記各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与する付与工程と、検索対象とする文書データの分野の指定を受け付ける受付工程と、前記受付工程によって受け付けられた分野に対応するタグ情報が付与された文書データを検索する検索工程と、を含んだことを特徴とする。

### 発明の効果

[0009] 本願に開示する情報管理装置及び情報管理方法は、利用者が望むジャンルの文書データが全体の文書データの量に比べて少ない場合であっても、例えば、検索漏れや関係のない文書データの混入を減らし、適切に文書データを検索することが可能である。

### 図面の簡単な説明

[0010] [図1]図1は、第一の実施形態に係る情報管理装置の構成の一例を示す図である。

[図2]図2は、第一の実施形態に係る文書データ記憶部によって記憶される情報の一例を示す図である。

[図3]図3は、第一の実施形態に係る情報管理装置において、収集した記事に

タグを付与し、タグに基づく検索を行う一連の処理について説明する図である。

[図4]図4は、第一の実施形態に係る情報管理装置におけるタグ付与処理の流れを説明するためのフローチャートである。

[図5]図5は、第一の実施形態に係る情報管理装置における情報検索処理の流れを説明するためのフローチャートである。

[図6]図6は、第二の実施形態に係る情報管理装置の構成の一例を示す図である。

[図7]図7は、第二の実施形態に係る分野別単語リスト記憶部によって記憶される情報の一例を示す図である。

[図8]図8は、第二の実施形態に係る情報管理装置において、分野別単語リストを参照して記事にタグを付与し、タグに基づく検索を行う一連の処理について説明する図である。

[図9]図9は、第二の実施形態に係る情報管理装置におけるタグ付与処理の流れを説明するためのフローチャートである。

[図10]図10は、第三の実施形態に係る情報管理装置において、記事にタグを付与し、キーワードに関連する記事の検索を行う一連の処理について説明する図である。

[図11]図11は、第三の実施形態に係る情報管理装置における情報検索処理の流れを説明するためのフローチャートである。

[図12]図12は、情報管理プログラムを実行するコンピュータを示す図である。

### 発明を実施するための形態

[0011] 以下に添付図面を参照して、この発明に係る情報管理装置及び情報管理方法の実施形態を詳細に説明する。なお、この実施形態によりこの発明が限定されるものではない。

[0012] [第一の実施形態]

以下の実施形態では、第一の実施形態に係る情報管理装置及び情報管理方

法による処理の流れを順に説明し、最後に第一の実施形態による効果を説明する。

[0013] [情報管理装置の構成]

図1に示した情報管理装置10の構成を説明する。図1は、第一の実施形態に係る情報管理装置10の構成を説明するための図である。図1に示すように、情報管理装置10は、通信処理部11、機能部12、記憶部13および制御部14を有する。また、情報管理装置10は、インターネットに接続されている。

[0014] 通信処理部11は、インターネットにおける装置との間でやり取りする各種情報に関する通信を制御する。例えば、通信処理部11は、インターネットにおけるサーバに対して記事等を含む文書データを要求し、文書データを受信する。

[0015] 記憶部13は、図1に示すように、文書データ記憶部13aを有する。記憶部13は、例えば、RAM (Random Access Memory)、フラッシュメモリ (Flash Memory) 等の半導体メモリ素子、又は、ハードディスク、光ディスク等の記憶装置などである。

[0016] 文書データ記憶部13aは、インターネット上のニュースサイト、BBS (Bulletin Board System)、Twitter (登録商標) などから収集された記事や投稿の文書データを記憶する。また、文書データ記憶部13aは、文書データに対応付けて、該文書データのジャンル (分野) を示すタグ情報を記憶する。なお、文書データ記憶部13aは、一般のデータベース (MySQLやPostgreSQL等) を用いてもよいし、表形式やテキスト形式での格納など、その蓄積方法の種類は問わない。

[0017] 例えば、文書データ記憶部13aは、図2に例示するように、文書データの内容を示す「記事本文」と、記事本文のジャンルを示す「タグ」とを対応付けて記憶する。ここで、「タグ」は、一つの記事に対して、一つであってもよいし、複数であってもよい。具体的な例を挙げて説明すると、図2に示すように、記事本文「スマホに充電機器経由でウィルス感染する脆弱性が発

見される」と、タグ「セキュリティ、携帯」とが対応付けて記憶されている。

[0018] 図1の説明に戻って、機能部12は、収集部12a、変換部12b、付与部12c、受付部12dおよび検索部12eを有する。ここで、機能部12は、各処理を受け持つところであり、実際にはソフトウェア（の1コンポーネント）またはミドルウェアとして実現される。また、制御部14は、通信処理部11、機能部12、記憶部13の動作を制御し、情報管理装置10の動作を司るもので、実際にはCPU（Central Processing Unit）やMPU（Micro Processing Unit）等の集積回路等で実現される。

[0019] 収集部12aは、ネットワーク上における複数の文書データを収集する。例えば、収集部12aは、インターネット上のニュースサイト、BBS、Twitterなどから記事を収集する。ここで、ニュースサイト、BBSについては、収集部12aは、事前にユーザが定めた収集先リストに基づいて、サイトへアクセスし、記事を収集する。

[0020] また、Twitterについては、収集部12aは、例えばStreaming APIやSearch APIを用いて、全Tweetから一部を取得したり、ユーザが事前に定めたキーワードやTwitterユーザIDに基づき、条件に当てはまるTweetを取得する。

[0021] さらに、収集部12aは、収集した記事を、分析に活用できるように整形する。具体的には、ニュースやBBSについては、不必要なHTMLタグやスクリプト、あるいは記事と関係ない広告を取り除いたりする。

[0022] 変換部12bは、収集部12aによって収集された各文書データに含まれる単語に基づいて、該文書データを特徴ベクトルに変換する。具体的には、変換部12bは、収集した記事データについて、不要文字の除去および文字種の統一を行った後、記事データを機械学習エンジンにかけるための特徴ベクトル変換を行う。

[0023] ここで、変換部12bは、不要文字の除去として、例えば、記事データに対して、余計な空白や言語処理の障害となるURL等の削除を行う。また、

例えば、変換部12bは、文字種の統一として、記事データに使用されている文字について、英大文字小文字やいわゆる半角全角の統一を行う。

[0024] また、変換部12bは、特徴ベクトルへの変換について、例えば、形態素解析によるもの、n-gramによるもの、区切り文字によるもの、のいずれかを利用することができる。変換部12bは、形態素解析によるものを利用した場合には、記事データを品詞によって分割し、それらの特徴ベクトルへ変換する。このような形態素解析には、例えばオープンソースのMeCab等のライブラリを利用することができる。例えば、変換部12bは、記事データが「Twitterの使い方が、まだ、よくわからん。」という文章だった場合に、形態素解析を利用し、「Twitter／の／使い方／が／、／まだ／、／よく／わから／ん／。」と分割する。

[0025] また、変換部12bは、n-gramによるものを利用した場合には、記事データを先頭から1文字ずつずらしながらn文字の組を作り、それらの特徴ベクトルへ変換する。例えば、変換部12bは、記事データが「Twitterの使い方が、まだ、よくわからん。」という文章だった場合に、n-gram (n=3) を利用し、「Twi／wit／itt／tte／ter／erの／rの使／の使い／・・・」と分割する。

[0026] また、変換部12bは、区切り文字によるものを利用した場合には、記事データを別途定めた区切り文字（空白やカンマ“,” など）によって分割し、それらの特徴ベクトルへ変換する。一般に、形態素解析は日本語の文章に、空白区切りは英語に適用されることが多い。例えば、変換部12bは、記事データが「Twitterの使い方が、まだ、よくわからん。」という文章で区切り文字にカンマ“,”を指定した場合に、区切り文字を利用し、「Twitterの使い方が／まだ／よくわからん。」と分割する。

[0027] そして、変換部12bは、このようにして要素に分割された記事データの特徴ベクトルに変換する。特徴ベクトルの変換手法としては、例えば、各要素の出現回数をそのまま特徴ベクトルとする方式、回数によらず出現するかどうかを1または0に対応させる方式、文章全体の出現回数を考慮した重みづ

けを行う方式などがある。これらは、使用する機械学習ライブラリの具備する手法であれば、どのようなものを用いても構わない。

[0028] 付与部12cは、収集部12aによって収集された各文書データに含まれる単語を用いて、各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与する。具体的には、付与部12cは、変換部12bによって変換された特徴ベクトルを用いて、各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与する。

[0029] 例えば、付与部12cは、変換部12bによって変換された特徴ベクトルを機械学習のエンジンに与え、事前に与えられたカテゴリに分類する。ここで事前に与えられたカテゴリとしては、例えば、「セキュリティ」、「プログラミング」、「携帯」などのユーザにとって関心がある分野とする。ここで使用する機械学習のエンジンは種類を問わないが、例えばオープンソースのJubatus等を用いることができる。そして、付与部12cは、記事データに対して、機械学習で分類されたカテゴリのタグを付与し、記事とタグを合わせて文書データ記憶部13aに格納する。

[0030] 受付部12dは、検索対象とする文書データの分野の指定を受け付ける。具体的には、受付部12dは、検索対象となる文書データの分野の候補を表示し、表示した分野の候補のなかから分野の指定を受け付ける。

[0031] 例えば、受付部12dは、検索を開始する検索指示を受け付けると、検索対象となる文書データの分野の候補として、例えば、「セキュリティ」、「プログラミング」、「携帯」などの分野を示す単語を表示するとともに、各単語の近傍にチェックボックスを表示する。そして、受付部12dは、チェックボックスにチェック記号が入れられた単語を、指定された分野として受け付ける。なお、受付部12dが指定を受け付ける分野の数は1つでもよいし、複数であってもよい。

[0032] 検索部12eは、受付部12dによって受け付けられた分野に対応するタグ情報が付与された文書データを検索する。例えば、検索部12eは、分野

「セキュリティ」について検索の指示を受け付けた場合には、「セキュリティ」のタグが付与された文書データを文書データ記憶部 13 a から検索する。そして、検索部 12 e は、検索した文書データを表示する。

[0033] なお、検索部 12 e は、複数の分野について検索の指示を受け付けた場合には、全ての分野に対応するタグが付与された文書データを文書データ記憶部 13 a から検索してもよいし、複数の分野のうちのいずれかの分野に対応するタグが付与された文書データを全て検索するようにしてもよい。

[0034] ここで、図 3 を用いて、情報管理装置 10 が、収集した記事にタグを付与し、タグに基づく検索を行う一連の処理について説明する。図 3 は、第一の実施形態に係る情報管理装置において、収集した記事にタグを付与し、タグに基づく検索を行う一連の処理について説明する図である。図 3 に示すように、情報管理装置 10 の収集部 12 a は、インターネット上のニュースサイト、Twitter、BBS 等から記事等の情報を収集する（図 3 の（1）参照）。

[0035] そして、変換部 12 b が収集部 12 a によって収集された各記事に含まれる単語に基づいて、該記事を特徴ベクトルに変換する。その後、付与部 12 c は、変換部 12 b によって変換された特徴ベクトルを機械学習のエンジンに与え、事前に与えられたカテゴリに分類し、カテゴリに対応するタグを記事等に付与する（図 3 の（2）参照）。そして、検索部 12 e は、ユーザに指定された分野に対応するタグ情報が付与された文書データを検索する（図 3 の（3）参照）。

[0036] [情報管理装置による処理]

次に、図 4、5 を用いて、第一の実施形態に係る情報管理装置 10 による処理を説明する。図 4 は、第一の実施形態に係る情報管理装置におけるタグ付与処理の流れを説明するためのフローチャートである。図 5 は、第一の実施形態に係る情報管理装置における情報検索処理の流れを説明するためのフローチャートである。

[0037] まず、図 4 を用いて、第一の実施形態に係る情報管理装置 10 におけるタ

グ付与処理の流れを説明する。図4に示すように、情報管理装置10の収集部12aは、インターネット上のWebサイト（ニュースサイト、BBS、Twitter、ブログ等）から記事を収集する（ステップS101）。

[0038] そして、変換部12bは、収集した記事について、不要文字の除去を行う（ステップS102）。例えば、変換部12bは、不要文字の除去として、記事データに対して、余計な空白や言語処理の障害となるURL等の削除を行う。

[0039] 続いて、変換部12bは、収集した記事について、文字種の統一を行う（ステップS103）。例えば、変換部12bは、文字種の統一として、記事データに使用されている文字について、英大文字小文字やいわゆる半角全角の統一を行う。

[0040] 変換部12bは、収集した記事について、不要文字の除去および文字種の統一を行った後、機械学習エンジンにかけるための特徴ベクトル変換を行う（ステップS104）。例えば、変換部12bは、特徴ベクトルへの変換について、形態素解析によるもの、n-gramによるもの、区切り文字によるもの、のいずれかを利用して記事を分割し、特徴ベクトルの変換を行う。

[0041] 続いて、付与部12cは、変換部12bによって変換された特徴ベクトルを機械学習のエンジンに与え、事前に与えられたカテゴリに分類する（ステップS105）。そして、付与部12cは、記事データに対して、機械学習で分類されたカテゴリのタグを付与する（ステップS106）。その後、付与部12cは、収集した記事と付与されたカテゴリを文書データ記憶部13aに格納する（ステップS107）。

[0042] 次に、図5を用いて、第一の実施形態に係る情報管理装置10における情報検索処理の流れを説明する。図5に示すように、情報管理装置10の受付部12dは、検索を開始する検索指示を受け付けると（ステップS201肯定）、検索対象となり得る複数の分野の候補を表示する（ステップS202）。

[0043] 例えば、受付部12dは、検索を開始する検索指示を受け付けると、検索

対象となる文書データの分野の候補として、例えば、「セキュリティ」、「プログラミング」、「携帯」などの分野を示す単語を表示するとともに、各単語の近傍にチェックボックスを表示する。そして、受付部12dは、チェックボックスにチェック記号が入れられた単語を、指定された分野として受け付ける。なお、受付部12dが指定を受け付ける分野の数は1つでもよいし、複数であってもよい。

[0044] そして、受付部12dは、表示した分野の候補のなかから分野の指定を受け付けた否かを判定する（ステップS203）。この結果、受付部12dが表示した分野の候補のなかから分野の指定を受け付けたと判定した場合に（ステップS203肯定）、選択された分野に対応するタグを有する記事を検索する（ステップS204）。例えば、検索部12eは、分野「セキュリティ」について検索の指示を受け付けた場合には、「セキュリティ」のタグが付与された文書データを文書データ記憶部13aから検索する。そして、検索部12eは、検索された記事を出力する（ステップS205）。

[0045] [第一の実施形態の効果]

上述してきたように、第一の実施形態にかかる情報管理装置10では、ネットワーク上における複数の文書データを収集し、収集された各文書データに含まれる単語を用いて、前記各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与する。そして、情報管理装置10では、検索対象とする文書データの分野の指定を受け付け、受け付けられた分野に対応するタグ情報が付与された文書データを検索する。利用者が望むジャンルの文書データが全体の文書データの量に比べて少ない場合であっても、適切に文書データを検索することが可能である。

[0046] 例えば、情報管理装置10では、インターネット上のWebサイト（ニュース、Twitter、BBS、ブログ等）から記事を収集し、機械学習によってこれらの記事の分類・タグ付けを行い、記事とタグを格納する。そして、記事に付与されたタグを基に、利用者が望む記事を検索することができるため、利用者が望むジャンルの記事の記事全体の量に比べて少ない場合で

も、タグ情報を基に、利用者が望む記事を多くの記事から探し出すことができる。

[0047] また、情報管理装置 10 では、収集された各文書データに含まれる単語に基づいて、該文書データを特徴ベクトルに変換する。そして、情報管理装置 10 では、変換された特徴ベクトルを用いて、各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与する。このため、文書データに対して適切にタグを付与することが可能である。

[0048] また、情報管理装置 10 では、検索対象となる文書データの分野の候補を表示し、表示した分野の候補のなかから分野の指定を受け付ける。このため、ユーザが検索キーワードを知らないような場合、例えば新しい話題に関する記事であっても、検索を行うことが可能である。

[0049] [第二の実施形態]

第二の実施形態において、情報管理装置が、所定の分野に関する単語リストを記憶し、単語リストを参照して、分野に対応するタグ情報を文書データごとに付与するようにしてもよい。そこで、以下では、図 6～図 9 を用いて、所定の分野に関する単語リストを記憶し、単語リストを参照して、各文書データから単語を抽出し、該単語に基づいて、該文書データを特徴ベクトルに変換する場合について説明する。なお、第一の実施形態と共通する構成および処理については、説明を省略する。

[0050] まず、図 6 を用いて、第二の実施の形態に係る情報管理装置 10A の構成を説明する。図 6 は、第二の実施形態に係る情報管理装置の構成の一例を示す図である。第二の実施の形態に係る情報管理装置 10A は、図 1 に示した第一の実施の形態に係る情報管理装置 10 と比較して、分野別単語リスト記憶部 13b を新たに備える点が相違する。

[0051] 分野別単語リスト記憶部 13b は、所定の分野に関する分野別単語リストを記憶する。例えば、分野別単語リスト記憶部 13b は、図 7 に例示するように、分野に対応付けて、各分野に関する単語リストを記憶する。図 7 の例を挙げて説明すると、例えば、分野別単語リスト記憶部 13b は、分野「セ

セキュリティ」に対応付けて単語「脆弱性、ウィルス・・・」を記憶する。図7は、第二の実施形態に係る分野別単語リスト記憶部13bによって記憶される情報の一例を示す図である。ここでは、分野別単語リスト記憶部13bは、ユーザが望むジャンルの単語を分野別単語リストとして記憶しているものとする。

[0052] 変換部12bは、分野別単語リスト記憶部13bに記憶された単語リストを参照して、各文書データから単語を抽出し、該単語に基づいて、該文書データを特徴ベクトルに変換する。

[0053] 例えば、変換部12bは、収集した記事データについて、不要文字の除去および文字種の統一を行った後、分野別単語抽出処理として、あらかじめ与えられた分野別単語リストを基に、各分野のリストに含まれる単語を記事本文から抽出し、抽出した結果である単語を特徴ベクトルに変換する。

[0054] 具体的な例を挙げて説明すると、変換部12bは、記事本文が「スマホに充電機器経由でウィルス感染する脆弱性が発見される」である場合に、図7に例示した分野別単語リストを参照して、リストに含まれる単語を記事本文から抽出処理を行うと、その結果として、分野「セキュリティ」の単語「ウィルス」、「脆弱性」と、分野「携帯」の単語「スマホ」を抽出することとなる。そして、変換部12bは、「ウィルス」、「脆弱性」および「スマホ」を特徴ベクトルに変換する。

[0055] その後、第一の実施形態と同様に、付与部12cは、変換部12bによって変換された特徴ベクトルを機械学習のエンジンに与え、事前に与えられたカテゴリに分類する。そして、付与部12cは、記事データに対して、機械学習で分類されたカテゴリにタグを付与し、記事とタグを合わせて文書データ記憶部13aに格納する。

[0056] なお、上記の処理において、特徴ベクトルに変換する処理を省略し、抽出した単語に対応する分野をタグとして記事データに付与してもよい。つまり、上記の例を用いて説明すると、付与部12cは、例えば、リストに含まれる単語を記事本文から抽出処理が行われた結果、「ウィルス」、「脆弱性」

および「スマホ」が抽出された場合には、ウィルスおよび脆弱性に対応する「セキュリティ」と、スマホに対応する「携帯」とを、タグとして付与してもよいし、単語数が最も多い単語に対応する「セキュリティ」のみをタグとして付与してもよい。

[0057] ここで、図8を用いて、情報管理装置10Aが、収集した記事にタグを付与し、タグに基づく検索を行う一連の処理について説明する。図8は、第二の実施形態に係る情報管理装置において、分野別単語リストを参照して記事にタグを付与し、タグに基づく検索を行う一連の処理について説明する図である。図8に示すように、情報管理装置10Aの収集部12aは、インターネット上のニュースサイト、Twitter、BBS等から記事等の情報を収集する(図8の(1)参照)。

[0058] そして、変換部12bが分野別単語リスト記憶部13bに記憶された単語リストを参照して、各文書データから単語を抽出し、該単語に基づいて、該文書データを特徴ベクトルに変換する。その後、付与部12cは、変換部12bによって変換された特徴ベクトルを機械学習のエンジンに与え、事前に与えられたカテゴリに分類し、カテゴリに対応するタグを記事等に付与する(図8の(2)参照)。そして、検索部12eは、ユーザに指定された分野に対応するタグ情報が付与された文書データを検索する(図8の(3)参照)。

[0059] 次に、図9を用いて、第二の実施形態に係る情報管理装置10Aによる処理を説明する。図9は、第二の実施形態に係る情報管理装置におけるタグ付与処理の流れを説明するためのフローチャートである。

[0060] 図9に示すように、情報管理装置10Aの収集部12aは、インターネット上のWebサイト(ニュースサイト、BBS、Twitter、ブログ等)から記事を収集する(ステップS301)。そして、変換部12bは、収集した記事について、不要文字の除去を行う(ステップS302)。例えば、変換部12bは、不要文字の除去として、記事データに対して、余計な空白や言語処理の障害となるURL等の削除を行う。

- [0061] 続いて、変換部12bは、収集した記事について、文字種の統一を行う（ステップS303）。例えば、変換部12bは、文字種の統一として、記事データに使用されている文字について、英大文字小文字やいわゆる半角全角の統一を行う。
- [0062] 変換部12bは、分野別単語リストを参照し、各分野のリストに含まれる単語を記事本文から抽出する（ステップS304）。そして、機械学習エンジンにかけるための特徴ベクトル変換を行う（ステップS305）。例えば、変換部12bは、特徴ベクトルへの変換について、抽出した単語をそのまま用いるもの、形態素解析によるもの、*n-gram*によるもの、区切り文字によるもの、のいずれかを利用して記事を分割し、特徴ベクトルの変換を行う。
- [0063] 続いて、付与部12cは、変換部12bによって変換された特徴ベクトルを機械学習のエンジンに与え、事前に与えられたカテゴリに分類する（ステップS306）。そして、付与部12cは、記事データに対して、機械学習で分類されたカテゴリのタグを付与する（ステップS307）。その後、付与部12cは、収集した記事と付与されたカテゴリを文書データ記憶部13aに格納する（ステップS308）。
- [0064] このように、第二の実施形態に係る情報管理装置10Aでは、所定の分野に関する単語リストを記憶する。そして、情報管理装置10Aは、単語リストを参照して、各文書データから単語を抽出し、該単語に基づいて、該文書データを特徴ベクトルに変換する。このため、特徴ベクトルへの変換の際に、単語リストの単語を使用することで、より分野に特化した分類が可能である。
- [0065] [第三の実施形態]
- 第三の実施形態では、情報管理装置が、検索対象とする文書データの分野の指定として、分野に関するキーワードの入力を受け付け、受け付けられたキーワードに対応するタグ情報が付与された文書データを検索するようにしてもよい。そこで、以下では、図10および図11を用いて、キーワードに

関連する記事の検索を行い、類似した記事を出力する場合について説明する。なお、第一の実施形態と共通する構成および処理については、説明を省略する。

[0066] まず、図10を用いて、第三の実施形態に係る情報管理装置10Bにおいて、記事にタグを付与し、キーワードに関連する記事の検索を行う一連の処理について説明する。図10は、第三の実施形態に係る情報管理装置において、記事にタグを付与し、キーワードに関連する記事の検索を行う一連の処理について説明する図である。

[0067] 図10に示すように、情報管理装置10Bの収集部12aは、インターネット上のニュースサイト、Twitter、BBS等から記事等の情報を収集する（図10の（1）参照）。

[0068] そして、付与部12cは、変換部12bによって変換された特徴ベクトルを機械学習のエンジンに与え、事前に与えられたカテゴリに分類し、カテゴリに対応するタグを記事等に付与する（図10の（2）参照）。そして、受付部12dは、ユーザからキーワードの入力を受け付ける（図10の（3）参照）。

[0069] 続いて、検索部12eは、キーワードに対応するタグが付与された記事を検索する（図10の（4）参照）。例えば、キーワードとして「脆弱性」が付与された場合には、「脆弱性」に対応するタグ「セキュリティ」が付与された記事を検索する。そして、検索部12eは、検索した結果を推薦結果として、キーワードに関連する記事をユーザへ出力する（図10の（5）参照）。

[0070] 次に、図11を用いて、第三の実施形態に係る情報管理装置10Bによる処理を説明する。図11は、第三の実施形態に係る情報管理装置における情報検索処理の流れを説明するためのフローチャートである。

[0071] 図11に示すように、情報管理装置10Bの受付部12dは、検索を開始する検索指示を受け付けると（ステップS401肯定）、キーワードの入力を受け付けたか否かを判定する（ステップS402）。そして、受付部12

dは、キーワードの入力を受け付け場合には（ステップS402肯定）、キーワードに対応するタグを有する記事を検索する（ステップS403）。例えば、検索部12eは、キーワードとして「脆弱性」が付与された場合には、「脆弱性」に対応するタグ「セキュリティ」が付与された記事を検索する。そして、検索部12eは、検索された記事を出力する（ステップS404）。

[0072] このように、第三の実施形態に係る情報管理装置10Bでは、検索対象とする文書データの分野の指定として、分野に関するキーワードの入力を受け付け、受け付けられたキーワードに対応するタグ情報が付与された文書データを検索する。このため、情報管理装置10Bでは、ユーザが入力したキーワードを基に、適切に文書データを検索することが可能である。

[0073] [システム構成等]

また、図示した各装置の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散・統合の具体的形態は図示のものに限られず、その全部または一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的または物理的に分散・統合して構成することができる。例えば、変換部12bと付与部12cとを統合してもよい。さらに、各装置にて行なわれる各処理機能は、その全部または任意の一部が、CPUおよび当該CPUにて解析実行されるプログラムにて実現され、あるいは、ワイヤードロジックによるハードウェアとして実現され得る。

[0074] また、本実施例において説明した各処理のうち、自動的におこなわれるものとして説明した処理の全部または一部を手動的に行うこともでき、あるいは、手動的におこなわれるものとして説明した処理の全部または一部を公知の方法で自動的におこなうこともできる。この他、上記文書中や図面中で示した処理手順、制御手順、具体的名称、各種のデータやパラメータを含む情報については、特記する場合を除いて任意に変更することができる。

[0075] [プログラム]

また、上記実施形態において説明した情報管理装置 10 が実行する処理をコンピュータが実行可能な言語で記述したプログラムを作成することもできる。例えば、第一の実施形態に係る情報管理装置 10 が実行する処理をコンピュータが実行可能な言語で記述した情報管理プログラムを作成することもできる。この場合、コンピュータが情報管理プログラムを実行することにより、上記実施形態と同様の効果を得ることができる。さらに、かかる情報管理プログラムをコンピュータ読み取り可能な記録媒体に記録して、この記録媒体に記録された情報管理プログラムをコンピュータに読み込ませて実行することにより上記第一の実施形態と同様の処理を実現してもよい。以下に、図 1 に示した情報管理装置 10 と同様の機能を実現する情報管理プログラムを実行するコンピュータの一例を説明する。

[0076] 図 12 は、情報管理プログラムを実行するコンピュータ 1000 を示す図である。図 12 に例示するように、コンピュータ 1000 は、例えば、メモリ 1010 と、CPU 1020 と、ハードディスクドライブインタフェース 1030 と、ディスクドライブインタフェース 1040 と、シリアルポートインタフェース 1050 と、ビデオアダプタ 1060 と、ネットワークインタフェース 1070 とを有し、これらの各部はバス 1080 によって接続される。

[0077] メモリ 1010 は、図 12 に例示するように、ROM (Read Only Memory) 1011 及び RAM 1012 を含む。ROM 1011 は、例えば、BIOS (Basic Input Output System) 等のブートプログラムを記憶する。ハードディスクドライブインタフェース 1030 は、図 12 に例示するように、ハードディスクドライブ 1031 に接続される。ディスクドライブインタフェース 1040 は、図 12 に例示するように、ディスクドライブ 1041 に接続される。例えば磁気ディスクや光ディスク等の着脱可能な記憶媒体が、ディスクドライブ 1041 に挿入される。シリアルポートインタフェース 1050 は、図 12 に例示するように、例えばマウス 1051、キーボード 1052 に接続される。ビデオアダプタ 1060 は、図 12 に例示するよう

に、例えばディスプレイ 1061 に接続される。

[0078] ここで、図 12 に例示するように、ハードディスクドライブ 1031 は、例えば、OS 1091、アプリケーションプログラム 1092、プログラムモジュール 1093、プログラムデータ 1094 を記憶する。すなわち、上記の情報管理プログラムは、コンピュータ 1000 によって実行される指令が記述されたプログラムモジュールとして、例えばハードディスクドライブ 1031 に記憶される。

[0079] また、上記実施形態で説明した各種データは、プログラムデータとして、例えばメモリ 1010 やハードディスクドライブ 1031 に記憶される。そして、CPU 1020 が、メモリ 1010 やハードディスクドライブ 1031 に記憶されたプログラムモジュール 1093 やプログラムデータ 1094 を必要に応じて RAM 1012 に読み出し、各種処理手順を実行する。

[0080] なお、情報管理プログラムに係るプログラムモジュール 1093 やプログラムデータ 1094 は、ハードディスクドライブ 1031 に記憶される場合に限られず、例えば着脱可能な記憶媒体に記憶され、ディスクドライブ等を介して CPU 1020 によって読み出されてもよい。あるいは、情報管理プログラムに係るプログラムモジュール 1093 やプログラムデータ 1094 は、ネットワーク (LAN (Local Area Network)、WAN (Wide Area Network) 等) を介して接続された他のコンピュータに記憶され、ネットワークインタフェース 1070 を介して CPU 1020 によって読み出されてもよい。

## 符号の説明

- [0081] 10、10A、10B 情報管理装置
- 11 通信処理部
  - 12 機能部
    - 12a 収集部
    - 12b 変換部
    - 12c 付与部

1 2 d 受付部

1 2 e 検索部

1 3 記憶部

1 3 a 文書データ記憶部

1 3 b 分野別単語リスト記憶部

1 4 制御部

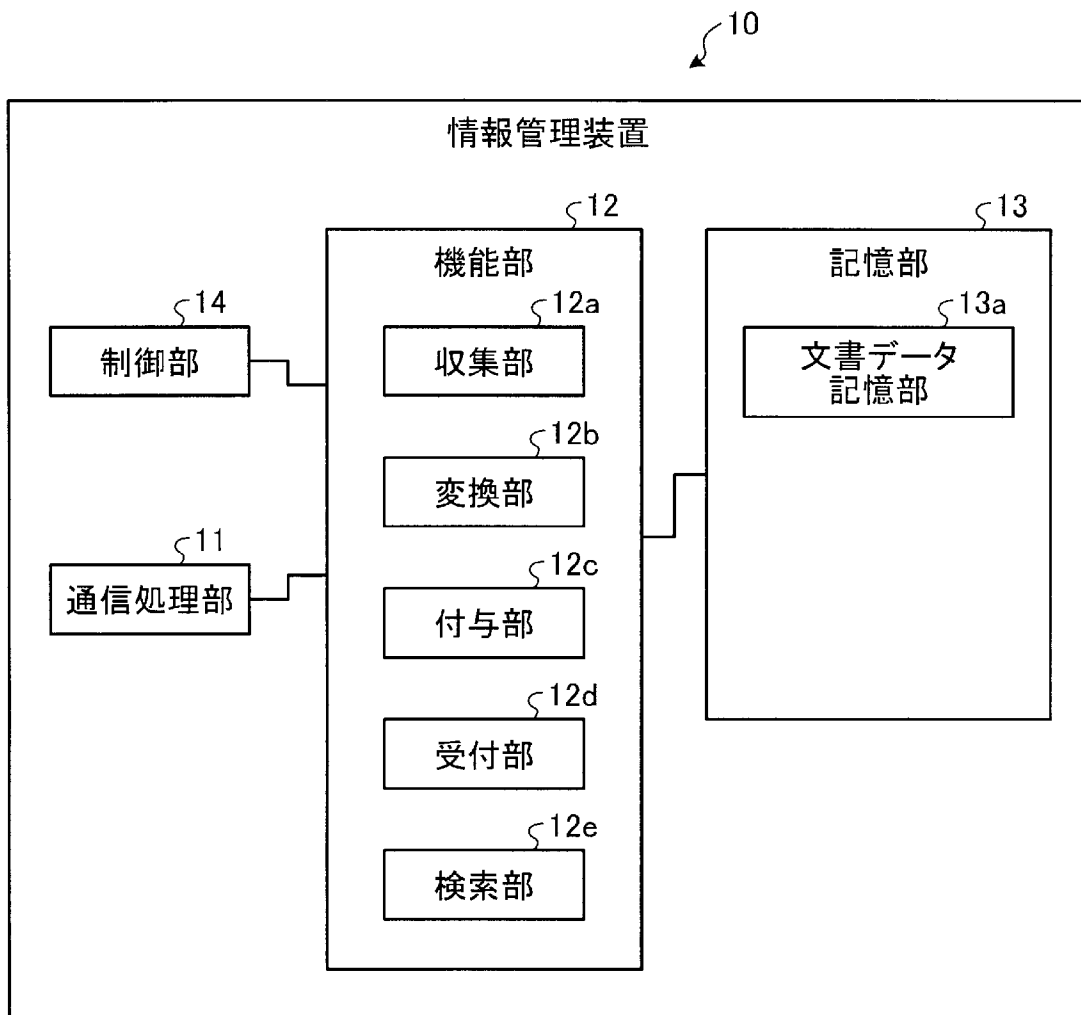
## 請求の範囲

- [請求項1] ネットワーク上における複数の文書データを収集する収集部と、  
前記収集部によって収集された各文書データに含まれる単語を用いて、前記各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与する付与部と、  
検索対象とする文書データの分野の指定を受け付ける受付部と、  
前記受付部によって受け付けられた分野に対応するタグ情報が付与された文書データを検索する検索部と、  
を有することを特徴とする情報管理装置。
- [請求項2] 前記収集部によって収集された各文書データに含まれる単語に基づいて、該文書データを特徴ベクトルに変換する変換部をさらに有し、  
前記付与部は、前記変換部によって変換された特徴ベクトルを用いて、前記各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与することを特徴とする請求項1に記載の情報管理装置。
- [請求項3] 所定の分野に関する単語リストを記憶する記憶部をさらに有し、  
前記変換部は、前記記憶部に記憶された単語リストを参照して、前記各文書データから単語を抽出し、該単語に基づいて、該文書データを特徴ベクトルに変換することを特徴とする請求項2に記載の情報管理装置。
- [請求項4] 前記受付部は、検索対象となる文書データの分野の候補を表示し、表示した分野の候補のなかから分野の指定を受け付けることを特徴とする請求項1～3のいずれか一つに記載の情報管理装置。
- [請求項5] 前記受付部は、検索対象とする文書データの分野の指定として、分野に関するキーワードの入力を受け付け、  
前記検索部は、前記受付部によって受け付けられたキーワードに対応するタグ情報が付与された文書データを検索することを特徴とする請求項1～3のいずれか一つに記載の情報管理装置。

## [請求項6]

情報管理装置によって実行される情報管理方法であって、  
ネットワーク上における複数の文書データを収集する収集工程と、  
前記収集工程によって収集された各文書データに含まれる単語を用いて、前記各文書データを所定の分野ごとに分類し、該分野に対応するタグ情報を文書データごとに付与する付与工程と、  
検索対象とする文書データの分野の指定を受け付ける受付工程と、  
前記受付工程によって受け付けられた分野に対応するタグ情報が付与された文書データを検索する検索工程と、  
を含んだことを特徴とする情報管理方法。

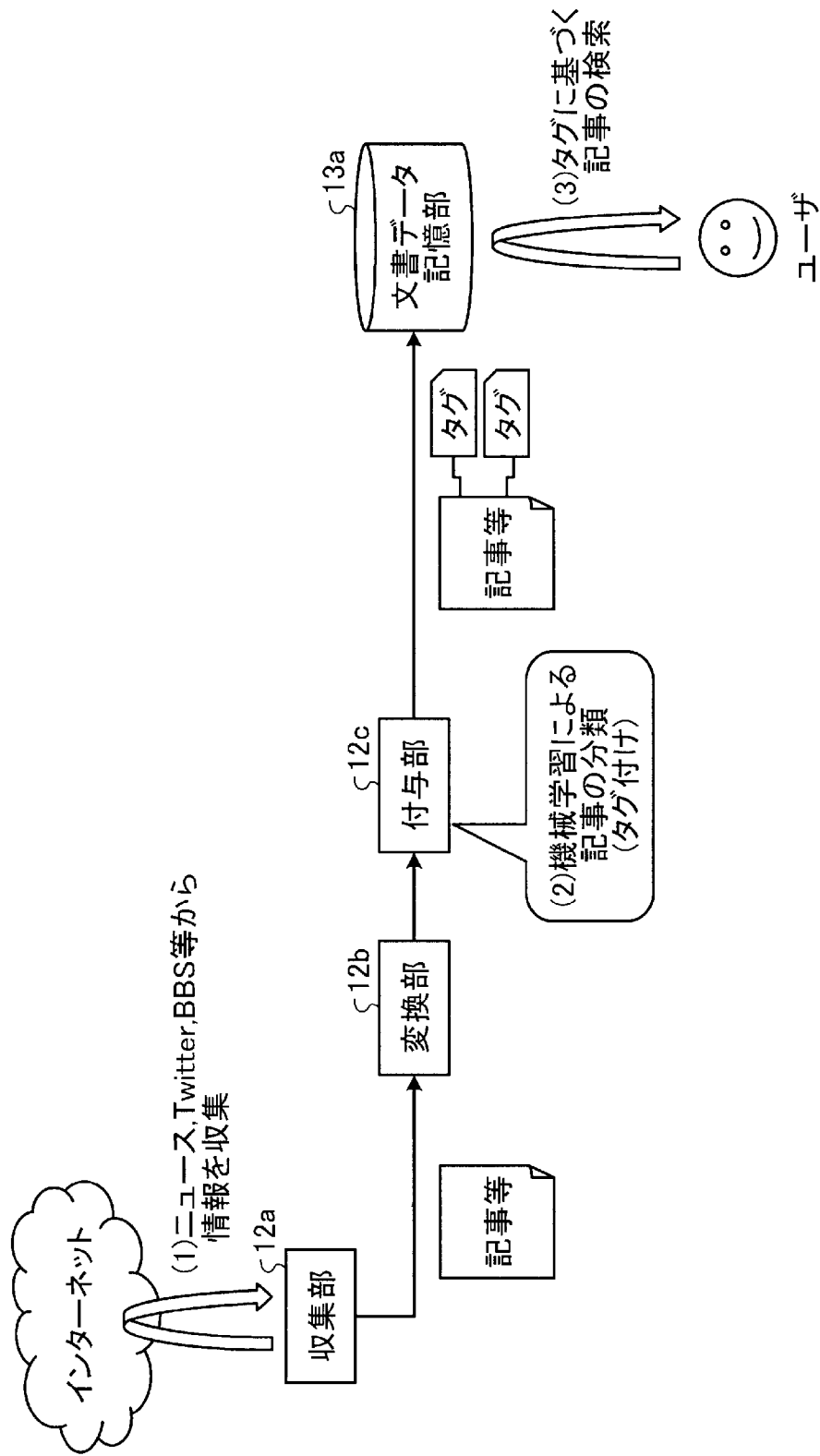
[図1]



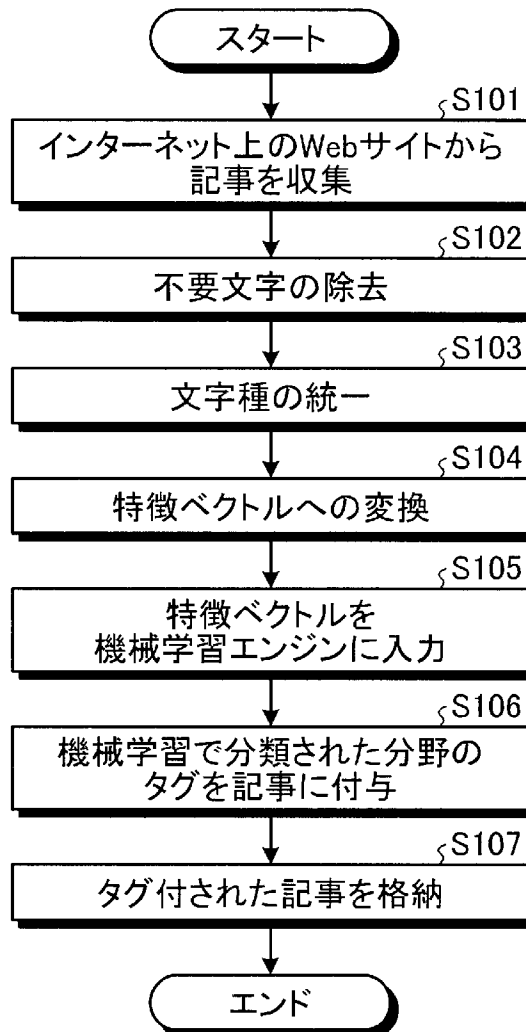
[図2]

記事本文	タグ
スマホに充電機器経由でウイルス感染する脆弱性が発見される	セキュリティ、携帯
⋮	⋮

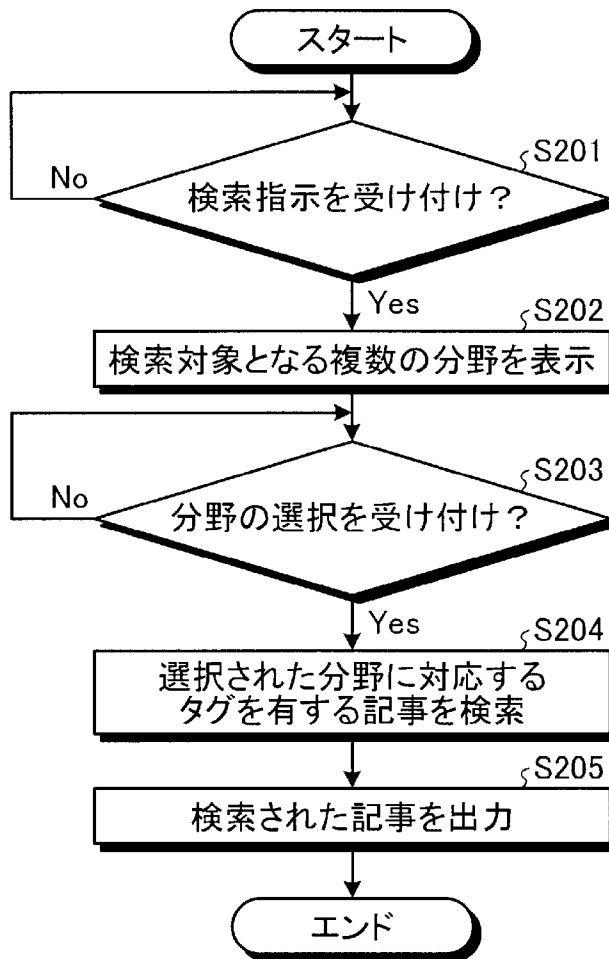
[図3]



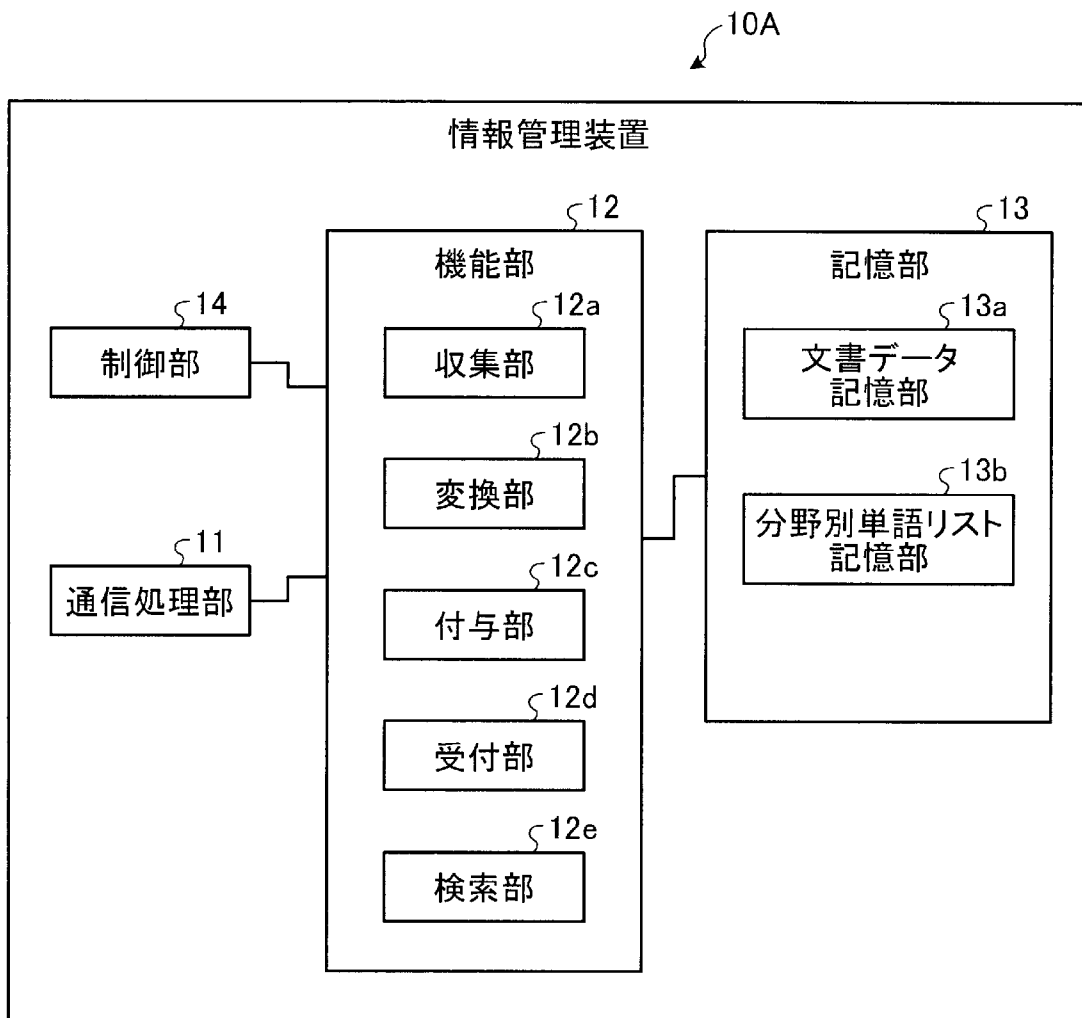
[図4]



[図5]



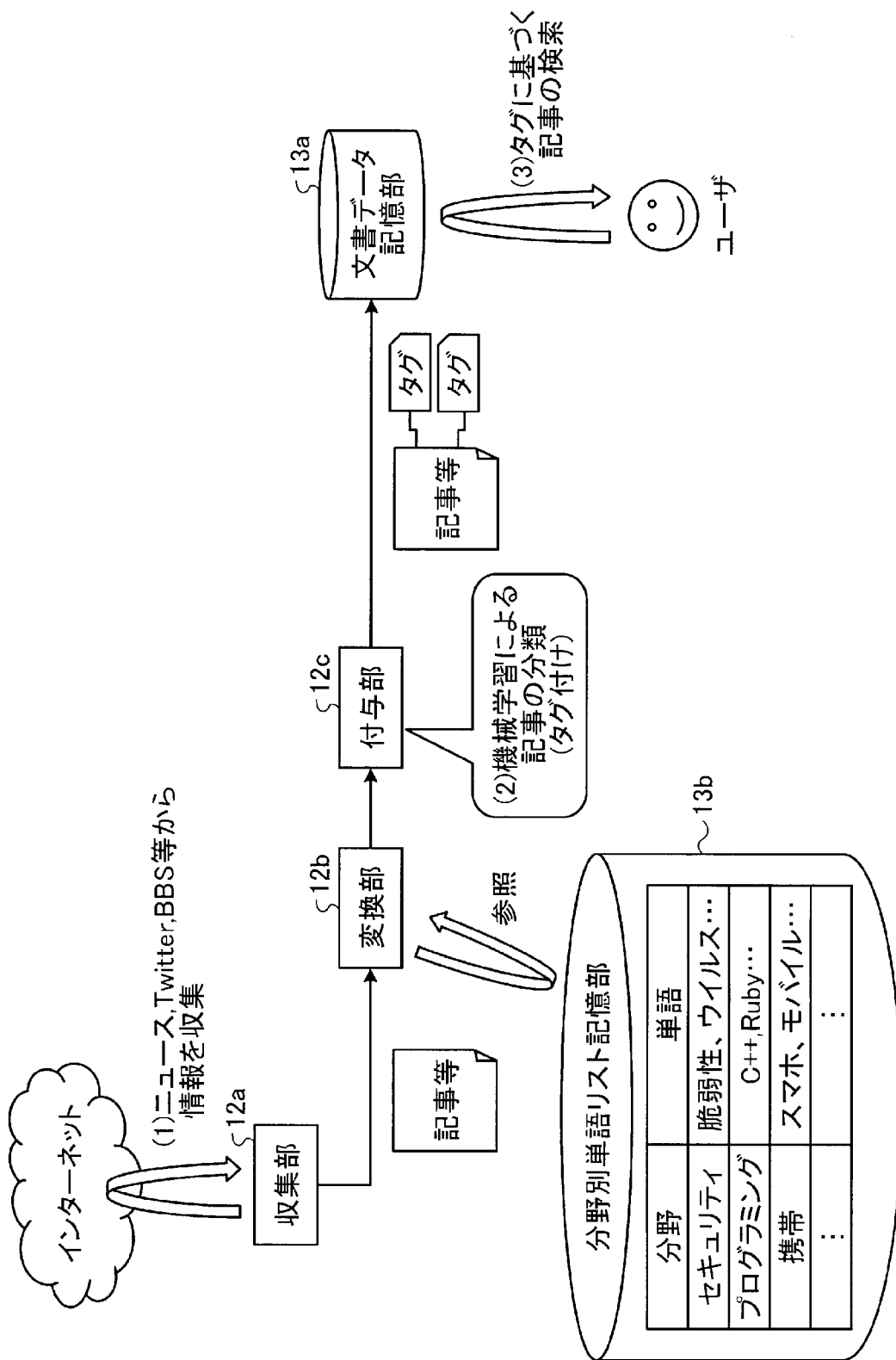
[図6]



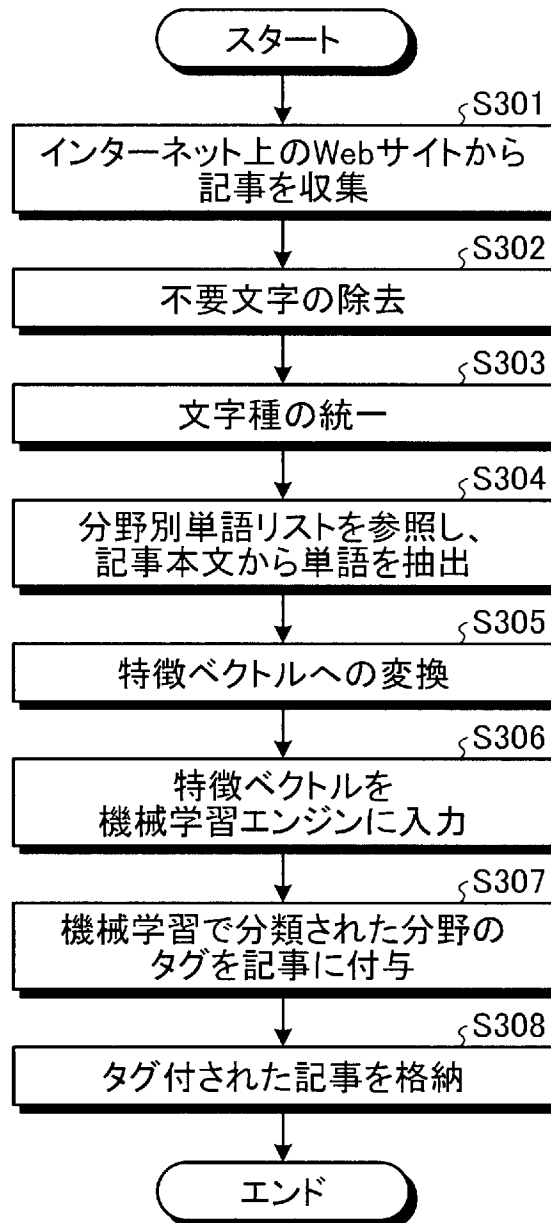
[図7]

分野	単語
セキュリティ	脆弱性、ウイルス…
プログラミング	C++,Ruby…
携帯	スマホ、モバイル…
⋮	⋮

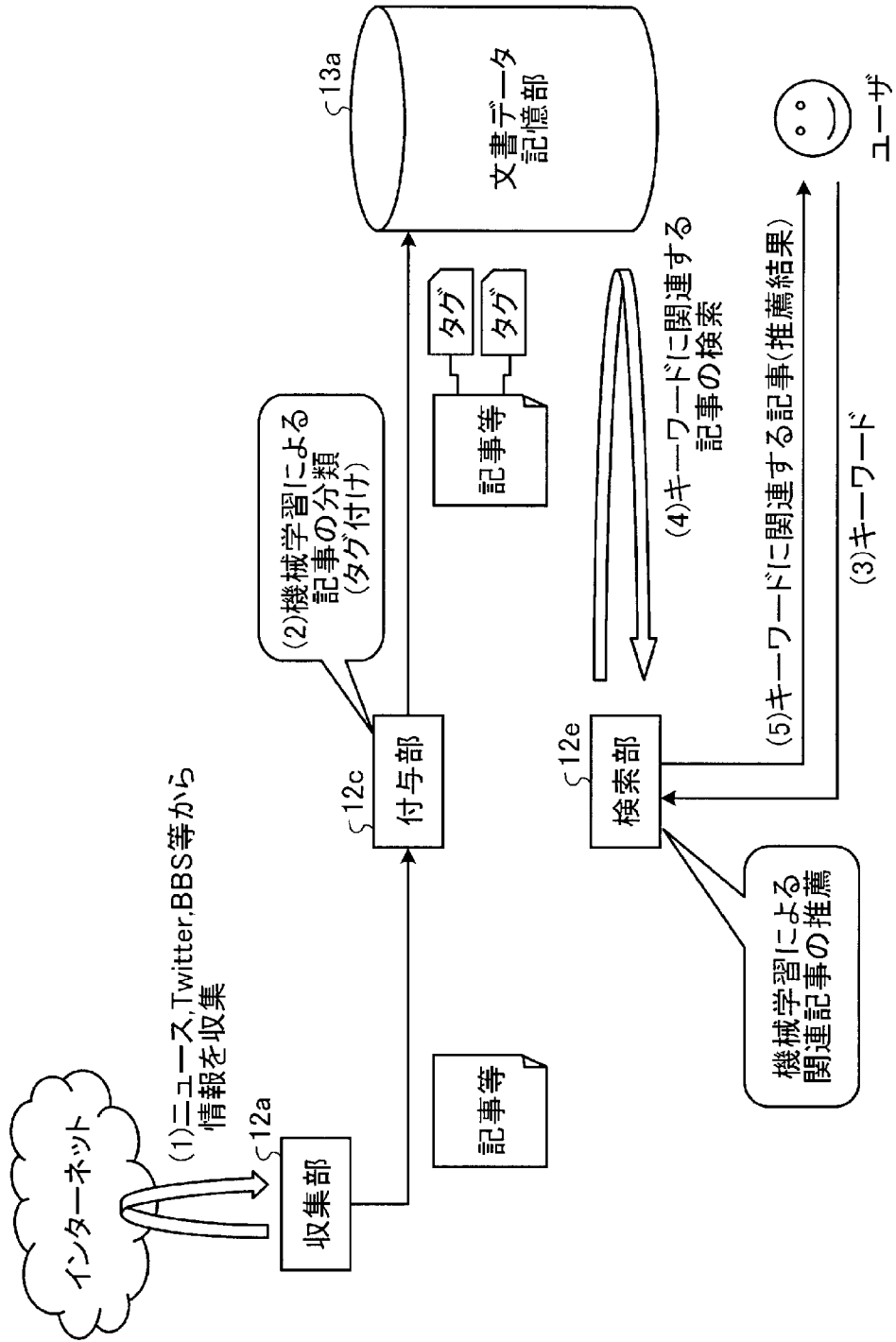
[図8]



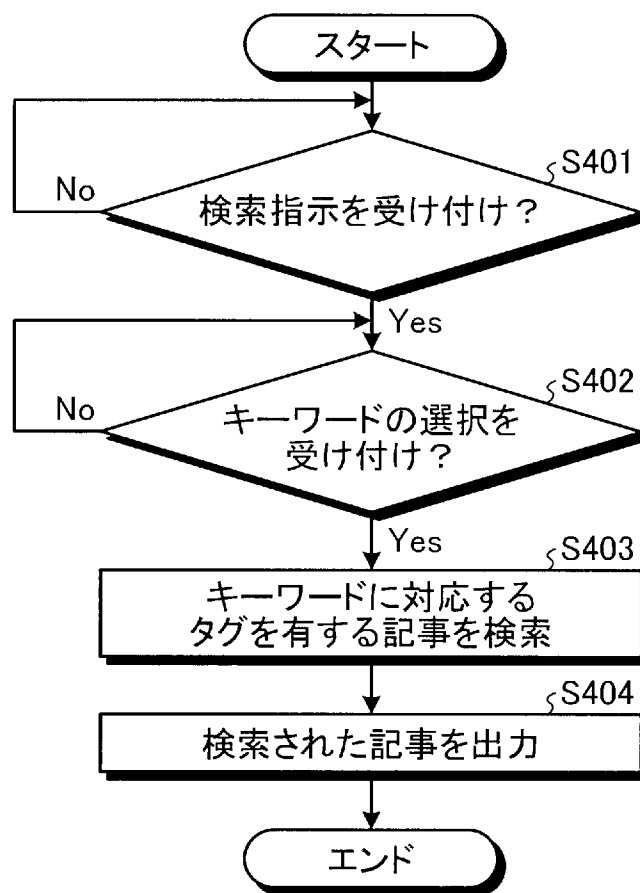
[図9]



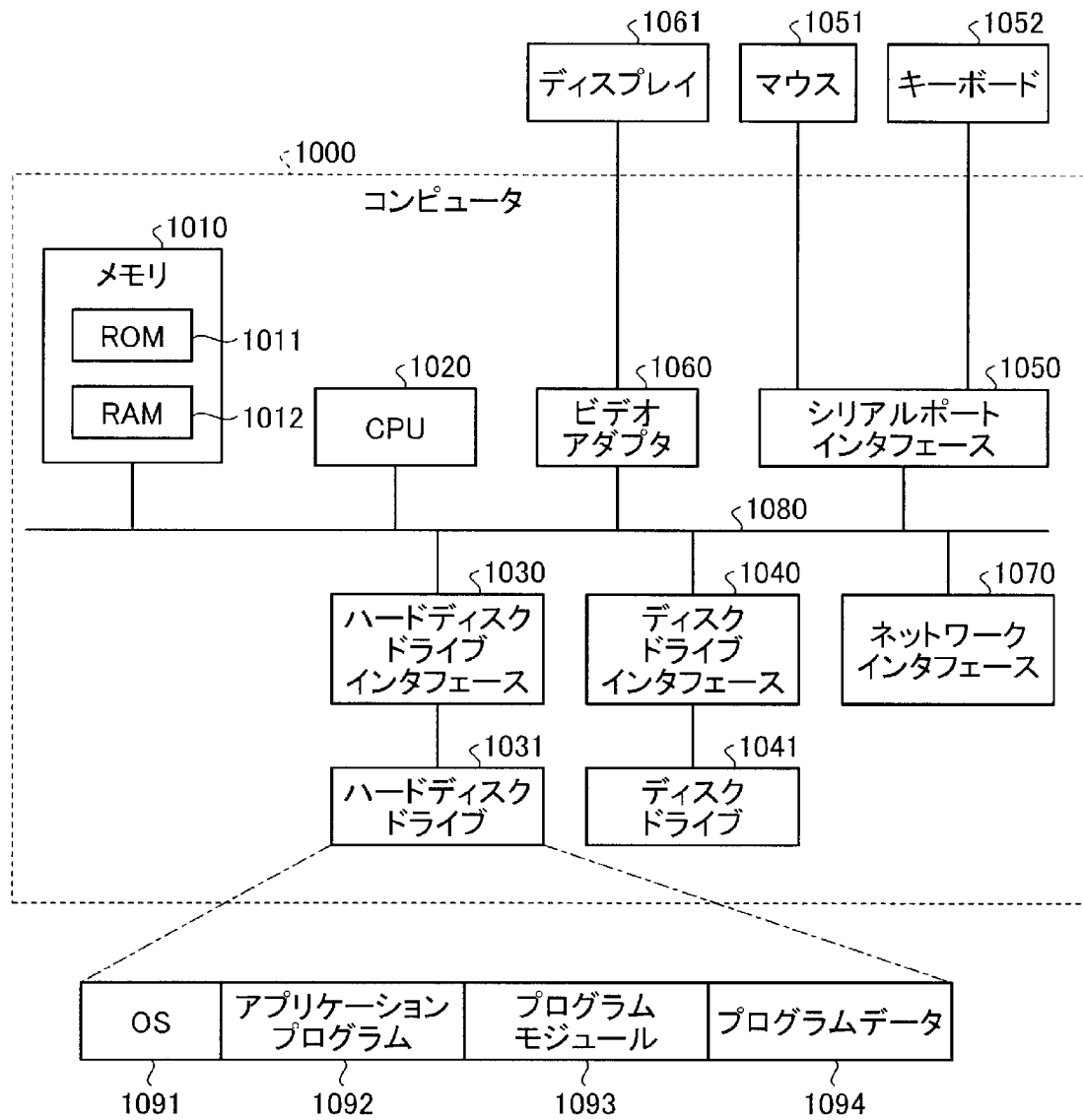
[図10]



[図11]



[図12]



**INTERNATIONAL SEARCH REPORT**

International application No.  
PCT/JP2014/069571

**A. CLASSIFICATION OF SUBJECT MATTER**  
G06F17/30(2006.01) i

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

**B. FIELDS SEARCHED**

Minimum documentation searched (classification system followed by classification symbols)  
G06F17/30

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Jitsuyo Shinan Koho	1922-1996	Jitsuyo Shinan Toroku Koho	1996-2014
Kokai Jitsuyo Shinan Koho	1971-2014	Toroku Jitsuyo Shinan Koho	1994-2014

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

**C. DOCUMENTS CONSIDERED TO BE RELEVANT**

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	JP 2008-217157 A (Nippon Telegraph and Telephone Corp.), 18 September 2008 (18.09.2008), paragraphs [0031], [0036] to [0041]; fig. 5 to 6 (Family: none)	1-6
Y	JP 2009-259248 A (NHN Corp.), 05 November 2009 (05.11.2009), paragraph [0006] & KR 10-2009-0108486 A	1-6
Y	JP 10-143537 A (Ricoh Co., Ltd.), 29 May 1998 (29.05.1998), paragraphs [0015], [0017] (Family: none)	1-6

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

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

Date of the actual completion of the international search 16 October, 2014 (16.10.14)	Date of mailing of the international search report 28 October, 2014 (28.10.14)
--	---

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

**INTERNATIONAL SEARCH REPORT**

International application No.

PCT/JP2014/069571

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	JP 2012-164018 A (NIFTY Corp.), 30 August 2012 (30.08.2012), entire text; all drawings (Family: none)	1-6
A	JP 2010-26923 A (Omron Corp.), 04 February 2010 (04.02.2010), entire text; all drawings (Family: none)	1-6
A	JP 2008-276344 A (Justsystem Corp.), 13 November 2008 (13.11.2008), entire text; all drawings (Family: none)	1-6

A. 発明の属する分野の分類 (国際特許分類 (IPC)) Int.Cl. G06F17/30(2006.01)i		
B. 調査を行った分野 調査を行った最小限資料 (国際特許分類 (IPC)) Int.Cl. G06F17/30		
最小限資料以外の資料で調査を行った分野に含まれるもの 日本国実用新案公報 1922-1996年 日本国公開実用新案公報 1971-2014年 日本国実用新案登録公報 1996-2014年 日本国登録実用新案公報 1994-2014年		
国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)		
C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
Y	JP 2008-217157 A (日本電信電話株式会社) 2008.09.18, 【0031】, 【0036】 - 【0041】, 第5-6図 (ファミリーなし)	1-6
Y	JP 2009-259248 A (エヌエイチエヌ コーポレーション) 2009.11.05, 【0006】 & KR 10-2009-0108486 A	1-6
Y	JP 10-143537 A (株式会社リコー) 1998.05.29, 【0015】, 【0017】 (ファミリーなし)	1-6
<input checked="" type="checkbox"/> C欄の続きにも文献が列挙されている。 <span style="margin-left: 200px;"><input type="checkbox"/> パテントファミリーに関する別紙を参照。</span>		
* 引用文献のカテゴリー 「A」 特に関連のある文献ではなく、一般的技術水準を示すもの 「E」 国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの 「L」 優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す) 「O」 口頭による開示、使用、展示等に言及する文献 「P」 国際出願日前で、かつ優先権の主張の基礎となる出願日の後に公表された文献 「T」 国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの 「X」 特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの 「Y」 特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの 「&」 同一パテントファミリー文献		
国際調査を完了した日 16.10.2014	国際調査報告の発送日 28.10.2014	
国際調査機関の名称及びあて先 日本国特許庁 (ISA/J P) 郵便番号100-8915 東京都千代田区霞が関三丁目4番3号	特許庁審査官 (権限のある職員) 吉田 誠 電話番号 03-3581-1101 内線 3599	5M 3659

C (続き) . 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	JP 2012-164018 A (ニフティ株式会社) 2012.08.30, 全文, 全図 (ファミリーなし)	1-6
A	JP 2010-26923 A (オムロン株式会社) 2010.02.04, 全文, 全図 (ファミリーなし)	1-6
A	JP 2008-276344 A (株式会社ジャストシステム) 2008.11.13, 全文, 全図 (ファミリーなし)	1-6