

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2020-52902
(P2020-52902A)

(43) 公開日 令和2年4月2日(2020.4.2)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 40/289 (2020.01)	G06F 17/27 680	5B091
G06F 16/30 (2019.01)	G06F 17/30 170A	
G06F 16/00 (2019.01)	G06F 17/30 220Z	
G06F 40/242 (2020.01)	G06F 17/27 635	

審査請求 未請求 請求項の数 7 O L (全 17 頁)

(21) 出願番号	特願2018-183861 (P2018-183861)	(71) 出願人	000003078 株式会社東芝 東京都港区芝浦一丁目1番1号
(22) 出願日	平成30年9月28日 (2018.9.28)	(71) 出願人	301063496 東芝デジタルソリューションズ株式会社 神奈川県川崎市幸区堀川町72番地34
		(74) 代理人	100108855 弁理士 蔵田 昌俊
		(74) 代理人	100103034 弁理士 野河 信久
		(74) 代理人	100075672 弁理士 峰 隆司
		(74) 代理人	100153051 弁理士 河野 直樹

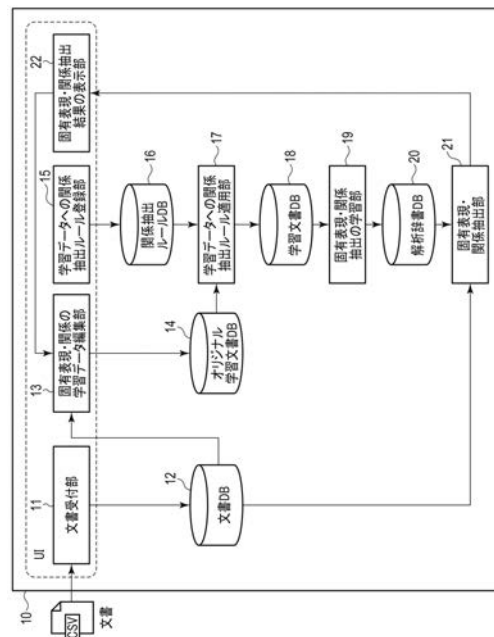
最終頁に続く

(54) 【発明の名称】 固有表現抽出装置、方法およびプログラム

(57) 【要約】

【課題】 文書からの固有表現抽出の精度を向上させる。
 【解決手段】 実施形態における固有表現抽出装置は、文書データから固有表現および固有表現同士の関係を抽出する抽出ルールを定めた抽出用辞書を格納する手段と、抽出対象である抽出用文書データおよび学習用文書データの入力を受け付ける文書受付手段と、抽出用辞書を用いて抽出用文書データから固有表現および関係を抽出する抽出手段と、入力操作に従い、学習用文書データにおける文字列のうち抽出させる固有表現に対応する文字列を指定する指定手段と、抽出用文書データから抽出させる、固有表現の分類同士の関係を定めた関係抽出ルールを格納する手段と、関係抽出ルールを適用することで、指定された固有表現のうち関係抽出ルールの分類に属する固有表現同士の関係を設定した学習文書を生成する生成手段と、学習文書に基づいて抽出用辞書を学習する学習手段とを有する。

【選択図】 図 1



【特許請求の範囲】**【請求項 1】**

文書データから当該文書データの固有表現および固有表現同士の間を抽出する抽出ルールを定めた抽出用辞書を格納する第 1 の格納手段と、

前記固有表現および前記関係の抽出対象である抽出用文書データおよび前記抽出用辞書の学習に用いられる学習用文書データの入力を受け付ける文書受付手段と、

前記抽出用辞書を用いて、前記文書受付手段により受け付けた抽出用文書データから固有表現および固有表現同士の間をそれぞれ抽出する抽出手段と、

入力操作に従い、前記文書受付手段により受け付けた学習用文書データにおける文字列のうち前記抽出手段により抽出させる固有表現に対応する文字列を指定する指定手段と、

前記抽出用文書データから抽出させる、固有表現の分類同士の間を定めた関係抽出ルールを格納する第 2 の格納手段と、

前記第 2 の格納手段に格納された関係抽出ルールを適用することで、前記指定手段により指定された固有表現のうち前記関係抽出ルールで定められた分類に属する固有表現同士の間を設定した学習文書を生成する生成手段と、

前記生成手段により生成された学習文書に基づいて、前記抽出用辞書を学習する学習手段と、

を備えた固有表現抽出装置。

【請求項 2】

前記抽出手段により抽出された固有表現を出力する第 1 の出力手段と、

前記第 1 の出力手段により出力された固有表現のうち入力操作で指定された固有表現の抽出元の文書データを出力する第 2 の出力手段と、

入力操作に従い、前記第 2 の出力手段により出力された抽出元の文書データの固有表現を編集する編集手段と、

をさらに備えた請求項 1 に記載の固有表現抽出装置。

【請求項 3】

前記抽出手段により抽出された固有表現同士の間を出力する第 1 の出力手段と、

前記第 1 の出力手段により表示された固有表現同士の間のうち入力操作で指定された固有表現同士の間を抽出元の文書データを出力する第 2 の出力手段と、

入力操作に従い、前記第 2 の出力手段により出力された抽出元の文書データの固有表現同士の間を編集する編集手段と、

をさらに備えた請求項 1 に記載の固有表現抽出装置。

【請求項 4】

文書データから当該文書データの固有表現および固有表現同士の間を抽出する抽出ルールを定めた抽出用辞書を格納する格納手段と、

前記固有表現および前記関係の抽出対象である抽出用文書データおよび前記抽出用辞書の学習に用いられる学習用文書データの入力を受け付ける文書受付手段と、

前記抽出用辞書を用いて、前記文書受付手段により受け付けた抽出用文書データから固有表現および固有表現同士の間をそれぞれ抽出する抽出手段と、

入力操作に従い、前記文書受付手段により受け付けた学習用文書データにおける文字列のうち前記抽出手段により抽出させる固有表現に対応する文字列および固有表現同士の間をそれぞれ指定する指定手段と、

前記指定手段により指定された固有表現に対応する文字列および固有表現同士の間に基づいて、前記抽出用辞書を学習する学習手段と、

前記指定手段により指定された固有表現に対応する文字列および固有表現同士の間のうち、前記抽出手段により抽出されなかった、固有表現に対応する文字列および固有表現同士の間を出力する出力手段と、

を備えた固有表現抽出装置。

【請求項 5】

文書データから当該文書データの固有表現および固有表現同士の間を抽出する抽出ル

10

20

30

40

50

ールを定めた抽出用辞書を格納する格納手段と、

前記固有表現および前記関係の抽出対象である抽出用文書データおよび前記抽出用辞書の学習に用いられる学習用文書データの入力を受け付ける文書受付手段と、

前記抽出用辞書を用いて、前記文書受付手段により受け付けた抽出用文書データから固有表現および固有表現同士の間をそれぞれ抽出する抽出手段と、

入力操作に従い、前記文書受付手段により受け付けた学習用文書データにおける文字列のうち前記抽出手段により抽出させる固有表現に対応する文字列および固有表現同士の関係をそれぞれ指定する指定手段と、

前記指定手段により指定された固有表現に対応する文字列および固有表現同士の関係に基づいて、前記抽出用辞書を学習する学習手段と、

前記指定手段により指定されない固有表現に対応する文字列および固有表現同士の関係のうち、前記抽出手段により抽出された、固有表現に対応する文字列および固有表現同士の関係を出力する出力手段と、

を備えた固有表現抽出装置。

【請求項6】

固有表現抽出装置に適用する方法であって、

固有表現および固有表現同士の関係の抽出対象である抽出用文書データおよび、前記抽出用文書データから当該抽出用文書データの固有表現および固有表現同士の関係を抽出する抽出ルールを定めた抽出用辞書の学習に用いられる学習用文書データの入力を受け付ける処理と、

文書データから当該文書データの固有表現および固有表現同士の関係を抽出する抽出ルールを定めた抽出用辞書を用いて、前記受け付けた抽出用文書データから固有表現および固有表現同士の関係をそれぞれ抽出する処理と、

入力操作に従い、前記受け付けた学習用文書データにおける文字列のうち前記抽出させる固有表現に対応する文字列を指定する処理と、

前記抽出用文書データから抽出させる、固有表現の分類同士の関係を定めた関係抽出ルールを適用することで、前記指定された固有表現のうち前記関係抽出ルールで定められた分類に属する固有表現同士の関係を設定した学習文書を生成する処理と、

前記生成された学習文書に基づいて、前記抽出用辞書を学習する処理と、
を実行する固有表現抽出方法。

【請求項7】

コンピュータを、

文書データから当該文書データの固有表現および固有表現同士の関係を抽出する抽出ルールを定めた抽出用辞書を格納する第1の格納手段、

前記固有表現および前記関係の抽出対象である抽出用文書データおよび前記抽出用辞書の学習に用いられる学習用文書データの入力を受け付ける文書受付手段、

前記抽出用辞書を用いて、前記文書受付手段により受け付けた抽出用文書データから固有表現および固有表現同士の関係をそれぞれ抽出する抽出手段、

入力操作に従い、前記文書受付手段により受け付けた学習用文書データにおける文字列のうち前記抽出手段により抽出させる固有表現に対応する文字列を指定する指定手段、

前記抽出用文書データから抽出させる、固有表現の分類同士の関係を定めた関係抽出ルールを格納する第2の格納手段、

前記第2の格納手段に格納された関係抽出ルールを適用することで、前記指定手段により指定された固有表現のうち前記関係抽出ルールで定められた分類に属する固有表現同士の関係を設定した学習文書を生成する生成手段、および

前記生成手段により生成された学習文書に基づいて、前記抽出用辞書を学習する学習手段、

として機能させる固有表現抽出処理プログラム。

【発明の詳細な説明】

【技術分野】

10

20

30

40

50

【0001】

本発明の実施形態は、固有表現抽出装置、方法およびプログラムに関する。

【背景技術】

【0002】

従来、人手によるルール又は機械学習などの様々な手法によって、文書データ中に出現する固有表現を抽出する仕組みが提案されてきた。

【0003】

また、文書データから抽出した固有表現について、この固有表現の分類名が出現する度合いから、分類名の重みを算出することで、どの固有表現を出力するかを判定するなどの応用技術も存在する。

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特開2007-148785号公報

【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、他の認識技術と同様に、固有表現抽出にも、理想的には100%の精度が期待されており、さらなる精度向上が求められている。

【0006】

本発明が解決しようとする課題は、文書からの固有表現抽出の精度を向上させることが可能な固有表現抽出装置、方法およびプログラムを提供することである。

【課題を解決するための手段】

【0007】

実施形態に係る固有表現抽出装置は、文書データから当該文書データの固有表現および固有表現同士の間接関係を抽出する抽出ルールを定めた抽出用辞書を格納する第1の格納手段と、前記固有表現および前記関係の抽出対象である抽出用文書データおよび抽出用辞書の学習に用いられる学習用文書データの入力を受け付ける文書受付手段と、前記抽出用辞書を用いて、前記文書受付手段により受け付けた抽出用文書データから固有表現および固有表現同士の間接関係をそれぞれ抽出する抽出手段と、入力操作に従い、前記文書受付手段により受け付けた学習用文書データにおける文字列のうち前記抽出手段により抽出させる固有表現に対応する文字列を指定する指定手段と、抽出用文書データから抽出させる、固有表現の分類同士の間接関係を定めた関係抽出ルールを格納する第2の格納手段と、前記第2の格納手段に格納された関係抽出ルールを適用することで、前記指定手段により指定された固有表現のうち前記関係抽出ルールで定められた分類に属する固有表現同士の間接関係を設定した学習文書を生成する生成手段と、前記生成手段により生成された学習文書に基づいて、前記抽出用辞書を学習する学習手段とを有する。

【発明の効果】

【0008】

本発明によれば、文書からの固有表現抽出の精度を向上させることができる。

【図面の簡単な説明】

【0009】

【図1】実施形態に係る固有表現抽出装置の機能構成例を示すブロック図。

【図2】実施形態に係る固有表現抽出装置の文書DBに格納される文書データの一例を表形式で示す図。

【図3】実施形態に係る固有表現抽出装置のオリジナル学習文書DBに格納される、固有表現に係る学習データの一例を表形式で示す図。

【図4】実施形態に係る固有表現抽出装置のオリジナル学習文書DBに格納される、固有表現同士の間接関係に係る学習データの一例を表形式で示す図。

【図5】実施形態に係る固有表現抽出装置の関係抽出ルールDBに格納される関係抽出ル

10

20

30

40

50

ールの一例を表形式で示す図。

【図 6】実施形態に係る固有表現抽出装置の解析辞書 DB に格納される解析辞書の一例を表形式で示す図。

【図 7】実施形態に係る固有表現抽出装置の第 1 の処理手順の一例を示すフローチャート。

【図 8】実施形態に係る固有表現抽出装置による、文書データの記述へのタグの付与時の表示画面の一例を示す図。

【図 9】実施形態に係る固有表現抽出装置の第 2 の処理手順の一例を示すフローチャート。

【図 10】実施形態に係る固有表現抽出装置による、抽出されたタグ、およびタグ同士の関係の表示画面の一例を示す図。 10

【図 11】実施形態に係る固有表現抽出装置による、抽出されたタグ、およびタグ同士の関係の編集画面の一例を示す図。

【図 12】実施形態に係る固有表現抽出装置の第 3 の処理手順の一例を示すフローチャート。

【図 13】実施形態に係る固有表現抽出装置による、文書データのタグ同士の関係の付与時の表示画面の一例を示す図。

【図 14】実施形態に係る固有表現抽出装置による、抽出されなかったタグ、およびタグ同士の関係の表示画面の一例を示す図。

【図 15】実施形態に係る固有表現抽出装置の第 4 の処理手順の一例を示すフローチャート。 20

【図 16】実施形態に係る固有表現抽出装置による、学習されなかったが抽出されたタグ、およびタグ同士の関係の表示画面の一例を示す図。

【発明を実施するための形態】

【0010】

以下、実施形態について図面を用いて説明する。

図 1 は、本発明の実施形態に係る固有表現抽出装置の機能構成例を示すブロック図である。

図 1 に示すように、実施形態に係る固有表現抽出装置 10 は、文書受付部 11、文書 DB (データベース) 12、固有表現・関係の学習データ編集部 13、オリジナル学習文書 DB 14、学習データへの関係抽出ルール登録部 15、関係抽出ルール DB 16、学習データへの関係抽出ルール適用部 17、学習文書 DB 18、固有表現・関係抽出の学習部 19、解析辞書 DB 20、固有表現・関係抽出部 21、および固有表現・関係抽出結果の表示部 22 を有する。 30

【0011】

また、固有表現抽出装置 10 は、パーソナルコンピュータ (PC) などのコンピュータデバイスを用いたシステムにより実現可能である。例えば、コンピュータデバイスは、CPU (Central Processing Unit) などのプロセッサと、プロセッサに接続されるメモリと、入出力インタフェースとを備える。このうちメモリは、不揮発性メモリなどの記憶媒体を有する記憶装置により構成される。 40

【0012】

文書受付部 11、固有表現・関係の学習データ編集部 13、学習データへの関係抽出ルール登録部 15、学習データへの関係抽出ルール適用部 17、固有表現・関係抽出の学習部 19、固有表現・関係抽出部 21、および固有表現・関係抽出結果の表示部 22 の機能は、例えば、プロセッサがメモリに格納されているプログラムを読み出して実行することにより実現される。なお、これらの機能の一部または全部は、特定用途向け集積回路 (ASIC) などの回路によって実現されてもよい。

【0013】

上記の機能のうち、文書受付部 11、固有表現・関係の学習データ編集部 13、学習データへの関係抽出ルール登録部 15、および固有表現・関係抽出結果の表示部 22 の機能 50

は、ユーザインタフェース（UI）における図示しない入力装置および表示装置と協働した機能として実現することができる。入力装置は、例えばキーボードおよびマウスである。表示装置は、例えば液晶ディスプレイである。

【0014】

文書DB12、オリジナル学習文書DB14、関係抽出ルールDB16、学習文書DB18、解析辞書DB20は、上記メモリのうち随時書込および読み出しが可能な不揮発性メモリに設けられる。

【0015】

固有表現抽出装置10は、文書データにおける固有表現（以下、タグと称することがある）の抽出結果と、固有表現同士の関係（以下、リンクと称することがある）の抽出結果とをあわせて表示装置に表示することができる。また、ユーザに、表示を参照して固有表現および固有表現同士の関係の誤抽出および未抽出を発見させることで、文書データからの固有表現および固有表現同士の関係抽出ルールを定めた解析辞書（抽出用辞書と称することもある）の学習に用いられる学習データを修正することを支援することもできる。

10

【0016】

文書受付部11は、1つ以上の文書データの入力（登録）を受け付けて、この受け付けた文書データを文書DB12に格納する。この格納される文書データは、（1）固有表現、および固有表現同士の関係抽出の対象となる抽出用文書データである場合と、（2）抽出用文書データからの固有表現、および固有表現同士の関係抽出ルールを定めた解析辞書の学習に用いる学習用文書データである場合とがある。

20

【0017】

図2は、実施形態に係る固有表現抽出装置10の文書DB12に格納される文書データの一例を表形式で示す図である。

図2に示した例では、文書DB12に格納される文書データは、（1）文書データに固有のコンテンツID、（2）タイトル、（3）本文などが関連付けられる。

【0018】

固有表現・関係の学習データ編集部13は、ユーザからの入力装置に対する操作に従い、文書DB12に格納される学習用文書データ中の、抽出させる（抽出されるべき）固有表現に対応する文字列と、固有表現同士の関係として抽出させる固有表現の組とを指定（付与）することで、固有表現、固有表現同士の関係の学習データ（オリジナル学習文書）を生成する。この学習データは、オリジナル学習文書DB14に格納される。

30

固有表現・関係の学習データ編集部13は、固有表現として抽出させる文字列と、固有表現同士の関係として抽出させる固有表現の組とを指定する指定手段と呼ぶこともできる。

オリジナル学習文書DB14に格納される学習データは、固有表現に係る学習データと、固有表現同士の関係に係る学習データとに区分される。

【0019】

図3は、実施形態に係る固有表現抽出装置10のオリジナル学習文書DB14に格納される、固有表現に係る学習データの一例を表形式で示す図である。

図3に示した例では、オリジナル学習文書DB14に格納される、固有表現に係る学習データは、（1）固有表現に固有のタグID、（2）固有表現が記述される文書データのコンテンツID、（3）タグの種類、（4）タグの値などが関連付けられる。

40

タグの種類は、固有表現の分類名、例えば「人名」、「地名」などである。タグの値は、具体的な固有表現の記述、例えば具体的な人名、地名などである。

【0020】

図4は、実施形態に係る固有表現抽出装置10のオリジナル学習文書DB14に格納される、固有表現同士の関係に係る学習データの一例を表形式で示す図である。

図4に示した例では、オリジナル学習文書DB14に格納される、固有表現同士の関係に係る学習データは、（1）固有表現同士の関係に固有の関係ID、（2）第1のタグID、（3）第1のタグの役割、（4）第2のタグID、（5）第2のタグの役割などが関

50

連付けられる。

【0021】

図4に示した例では、関係ID「1」に関する関係として、タグIDが「1」である固有表現の役割「住人」と、タグIDが「2」である固有表現の役割「住んでいる地域」との間に関係が存在することが示される。また、この例では、関係ID「2」に関する関係として、タグIDが「3」である固有表現の役割「スポーツ」と、タグIDが「4」である固有表現の役割「順位」との間に関係が存在することが示される。

図4では、2種類の固有表現同士に存在する関係について定義された例について示したが、これに限らず3種類以上の固有表現同士に存在する関係について定義されてもよい。

【0022】

学習データへの関係抽出ルール登録部15は、抽出用文書データから関係が抽出されるべき固有表現の分類名(種類)の組を定める関係抽出ルールを、UIに対するユーザからの入力操作に従って指定(登録)して、関係抽出ルールDB16に格納する。

【0023】

図5は、実施形態に係る固有表現抽出装置10の関係抽出ルールDB16に格納される関係抽出ルールの一例を表形式で示す図である。

図5に示した例では、関係抽出ルールDB16に格納される関係抽出ルールは、(1)関係抽出ルールに固有のルールID、(2)第1のタグの種類、(3)第2のタグの種類、(4)第1のタグの役割、(5)第2のタグの役割などが関連付けられる。関係抽出ルールDB16が設けられる不揮発性メモリは、関係抽出ルールを格納する格納手段と呼ぶこともできる。

【0024】

図5に示した例では、ルールID「1」に関する関係として、第1のタグの種類「人名」と、第2のタグの種類「地名」と、第1のタグの役割「住人」と、第2のタグの役割「住んでいる地域」との間に関係が存在することが示される。また、この例では、ルールID「2」に関する関係として、第1のタグの種類「スポーツ」と、第2のタグの種類「順位」と、第1のタグの役割「競技名」と、第2のタグの役割「競技結果」との間に関係が存在することが示される。

【0025】

学習データへの関係抽出ルール適用部17は、オリジナル学習文書DB14に格納される学習データに、関係抽出ルールDB16に格納される関係抽出ルールを適用することで、当該学習データにおける固有表現同士の関係のうち関係抽出ルールで定められる分類名で示される分類に属する固有表現同士の関係を一括で登録する。これにより、学習データへの関係抽出ルール適用部17は、固有表現同士の関係が登録された学習データである学習文書を生成する。この学習文書は学習文書DB18に格納される。学習データへの関係抽出ルール適用部17は、学習文書を生成する生成手段と呼ぶこともできる。

学習文書DB18に格納される学習文書の各項目は、上記のオリジナル学習文書DB14に格納される各種の学習データ(図3、4参照)と同じである。

【0026】

固有表現・関係抽出の学習部19は、学習文書DB18に格納された学習文書の内容を解析辞書DB20に格納される解析辞書に反映することで、固有表現と固有表現同士の関係との抽出ルールを定めた抽出用辞書を学習する。

【0027】

図6は、実施形態に係る固有表現抽出装置10の解析辞書DB20に格納される解析辞書の一例を表形式で示す図である。

図6に示した例では、解析辞書DB20に格納される解析辞書(抽出用辞書)は、各行に固有の辞書ID、タグの種類、タグの特徴、タグの値、複数種類のタグ同士の関係などが関連付けられる。タグの特徴とは、タグの記述形式、例えばバイナリデータを示す。解析辞書DB20が設けられる不揮発性メモリは、解析辞書を格納する格納手段と呼ぶこともできる。

10

20

30

40

50

【0028】

この解析辞書は、抽出用文書データから固有表現と固有表現同士の関係とを抽出するために照合される辞書である。この解析辞書は、過去の学習用文書データに基づく学習結果が反映されており、また、新たな学習用文書データに基づく学習結果が反映される。この解析辞書は、ニューラルネットワークから構成される学習器であってもよい。

【0029】

固有表現・関係抽出部21は、解析辞書DB20に格納される解析辞書と、文書DB12に格納される抽出用文書データとを照合することで、抽出用文書データから固有表現と固有表現同士の関係をそれぞれ抽出する。

【0030】

固有表現・関係抽出結果の表示部22は、固有表現・関係抽出部21による固有表現および固有表現同士の関係の抽出結果を表示装置に表示する。固有表現・関係抽出結果の表示部22は、固有表現および固有表現同士の関係の抽出結果を出力する出力手段と呼ぶこともできる。

また、固有表現・関係抽出結果の表示部22は、固有表現の抽出結果と固有表現同士の関係の抽出結果とを重ねて表示装置に表示することもできる。これにより、ユーザは、固有表現の誤抽出および未検出を発見しやすくなる。

【0031】

(第1の処理)

次に、固有表現抽出装置10の第1の処理について説明する。

図7は、実施形態に係る固有表現抽出装置10の第1の処理手順の一例を示すフローチャートである。

まず、ユーザからの入力操作に従って、文書受付部11は、学習用文書データの登録を受け付けて、登録した学習用文書データを文書DB12に格納する(S11)。

【0032】

文書DB12に格納された学習用文書データは表示装置に表示される。この表示された状態で、表示画面上の学習用文書データの本文の記述に対するユーザからの入力操作に従って、固有表現・関係の学習データ編集部13は、学習用文書データにおける、ユーザの入力操作により指定された記述にタグであることを示すマーク(下線)を付与する(以下、タグを付与すると称することがある)。タグの付与により生成された、固有表現に係る学習データ(図3参照)はオリジナル学習文書DB14に格納される(S12)。なお、固有表現同士の関係に係る学習データ(図4参照)は、第1の処理では生成されない。

【0033】

図8は、実施形態に係る固有表現抽出装置10による、文書データの記述へのタグの付与時の表示画面G1の一例を示す図である。

図8に示した例では、固有表現抽出装置10の表示装置に表示された画面G1上の文書データの本文中の各記述「搬送異常」、「ボルト」、「ボルトが緩んでいます。」、「アームに付いているネジを締めました。」への画面上のポインタによる指定などにより、各記述にタグをそれぞれ付与することができる。

【0034】

画面G1上の分類名にかかるウインドウに対するポインタによる指定などにより、タグが付与される各記述には当該タグの分類名をあわせて付与できる。図8に示した例では、タグが付与された上記の記述「搬送異常」にはタグの分類名「現象」を、上記の記述「ボルト」には分類名「部位」を、上記の記述「ボルトが緩んでいます。」には分類名「原因」を、上記の記述「アームに付いているネジを締めました。」には分類名「対処」をそれぞれ付与できる。

【0035】

また、画面G1と異なる図示しない設定画面上でのユーザからの入力操作に従って、学習データへの関係抽出ルール登録部15は、タグの任意の第1の分類名とタグの任意の第2の分類名との間に関係(リンク)を付与する。この付与により生成される関係抽出ルー

10

20

30

40

50

ル（図5参照）は関係抽出ルールDB16に格納される（S13）。この設定画面は、画面G1における学習用文書データの表示と並べて表示することができる。第1の処理におけるタグの任意の分類名同士の関係の付与は、上記の学習用文書データの記述によらない付与である。この付与は、1つの分類名と複数の分類名との間で行なうこともできる。

【0036】

ここで、学習データへの関係抽出ルール適用部17は、以下条件の時に、関係抽出ルールDB16に格納される関係抽出ルールの登録内容に合わせて、オリジナル学習文書DB14に格納される学習データで示されるタグのうち、後述のある分類名に係るタグと、別の分類名に係るタグとの関係を追加、編集、または削除する処理を行なう。上述の条件とは、S13において、（1）学習データへの関係抽出ルール登録部15によって、関係抽出ルールに対して、ある分類名に係るタグと、別の分類名に係るタグとの間の関係付与（登録）が完了している場合（S14のNO）、または、（2）完了する前で、当該付与を新たに行なう場合（S15のYES）である。この処理により生成された学習文書は学習文書DB18に格納される（S16）。例えば、関係抽出ルールで、分類名Aと分類名Bとの関係が定義されていれば、学習データにおける、分類名Aに属するタグと分類名Bに属するタグとの間に関係が付与されることになる。

上記の第1の処理によれば、学習データで示される、ある分類名に係るタグと、別の分類名に係るタグとの関係を一括で登録できる。

【0037】

（第2の処理）

次に、固有表現抽出装置10の第2の処理について説明する。

図9は、実施形態に係る固有表現抽出装置10の第2の処理手順の一例を示すフローチャートである。

第2の処理では、まず、固有表現・関係抽出結果の表示部22は、固有表現・関係抽出部21により抽出用文書データから抽出したタグ、およびタグ同士の関係を分類名ごとにグループ化した抽出結果の表示画面G2を表示装置に表示する（S21）。固有表現・関係抽出部21による抽出結果と抽出元文書データとの関係を示す情報は、固有表現・関係抽出結果の表示部22に接続される内部メモリに格納されているとする。

【0038】

図10は、実施形態に係る固有表現抽出装置10による、抽出されたタグ、およびタグ同士の関係の表示画面G2の一例を示す図である。

図10に示した表示画面G2では、分類（分類名）A、B、C、Dなどに属する複数種類のタグが示され、ある分類に属するタグと異なる分類に属するタグとの間の関係を示す。

図10では、分類Aと分類Bとの間、分類Bと分類Cとの間、分類Cと分類Dとの間でのタグ同士の関係がそれぞれ設定される例を示すが、これに限らず、例えば分類Aと分類Cとの間、分類Bと分類Dとの間などでのタグ同士の関係が設定されてもよい。

【0039】

ユーザは、抽出結果の表示画面G2で示される、気になるタグ、またはタグ同士の関係を入力操作により指定することができる（S22）。気になるタグ、またはタグ同士の関係とは、抽出用文書データからの抽出結果として適切でない可能性があるタグ、またはタグ同士の関係である。S22での指定に伴い、固有表現・関係抽出結果の表示部22は、上記の内部メモリに格納される、固有表現・関係抽出部21による抽出結果と抽出元文書データとの関係を示す情報を固有表現・関係の学習データ編集部13に渡す。

【0040】

S22での指定を受けて、固有表現・関係の学習データ編集部13は、指定されたタグ、またはタグ同士の関係の抽出元文書データを上記の渡された情報から検索して、この検索された抽出元文書データの本文などを表示装置に表示する（S23）。

【0041】

この表示を受けて、ユーザからの入力操作により、固有表現・関係の学習データ編集部

10

20

30

40

50

13は、抽出元文書データの記述に付与されているタグ、またはタグとタグの関係を編集する(S24)。

【0042】

図11は、実施形態に係る固有表現抽出装置10による、抽出されたタグ、およびタグ同士の関係の編集画面G3の一例を示す図である。

図11に示した例では、分類Bに属する1つ目のタグと分類Cに属する1つ目のタグ同士の関係が編集対象として指定された例を示す。この画面G3では、ユーザからの入力操作に従って、固有表現・関係の学習データ編集部13は、指定された関係の変更、例えば分類Bに属する1つ目のタグと、分類Cに属する2つ目以降のタグ同士の関係への修正、または関係の削除などを行なうことができる。

10

また、上記のように、付与済みのタグ自体の修正または削除などを行なうこともできる。タグ自体の修正とは、例えば分類名の修正、対象となる記述の変更である。タグ自体の削除とは、対象となる記述に対する固有表現としての指定の解除である。

第2の処理によれば、タグ、およびタグとタグとの関係の抽出結果のうち、指定された抽出結果の抽出元文書を容易に表示することができる。また、タグとタグとの関係の確認、編集を容易に行うことができる。

【0043】

(第3の処理)

次に、固有表現抽出装置10の第3の処理について説明する。

図12は、実施形態に係る固有表現抽出装置10の第3の処理手順の一例を示すフローチャートである。

20

まず、ユーザからの入力操作に従って、文書受付部11は、学習用文書データの登録を受け付けて、登録した学習用文書データを文書DB12に格納する(S31)。ここでは抽出用文書データは文書DB12に格納済みであるとする。

【0044】

文書DB12に格納される学習用文書データは表示装置に表示される。この表示された状態で、表示画面上の学習用文書データの記述に対するユーザからの入力操作に従って、固有表現・関係の学習データ編集部13は、学習用文書データの記述にタグを付与する。タグの付与により生成された、固有表現に係る学習データ(図3参照)はオリジナル学習文書DB14に格納される。

30

ここでは、学習用文書データの記述へのタグの付与時の表示画面は図8に示した表示画面G1であるとする。

【0045】

この表示画面G1に表示される学習用文書データの記述に対するユーザからの入力操作にしたがって、固有表現・関係の学習データ編集部13は、学習用文書データにおける記述に付与された第1のタグと第2のタグとの間に関係(リンク)を付与する。この付与により生成された、固有表現同士の関係に係る学習データ(図4参照)はオリジナル学習文書DB14に格納される(S32)。第3の処理におけるタグ同士の関係の付与は、上記の学習用文書データの本文の記述に対する付与である。ここでの関係の付与は、1つのタグと複数のタグとの間で行なうこともできる。

40

【0046】

図13は、実施形態に係る固有表現抽出装置10による、文書データのタグ同士の関係の付与時の表示画面G4の一例を示す図である。

図13では、表示装置に表示された画面G4上の文書データの本文中の第1の記述「ボルトが緩んでいます。」に付与された、分類名「原因」に係るタグと、本文中の第2の記述「アームに付いているネジを締めました。」に付与された、分類名「対処」に係るタグとの関係を示す線L1が付与された例を示す。

【0047】

第3の処理では、第1の処理で説明した、学習データへの関係抽出ルール登録部15による処理は行なわれず、S32でオリジナル学習文書DB14に格納された各種学習デー

50

タは、学習データへの関係抽出ルール適用部 17 を介して学習文書として学習文書 DB 18 に格納される。

【0048】

次に、固有表現・関係抽出の学習部 19 は、学習文書 DB 18 に格納された学習文書の内容を解析辞書 DB 20 に格納される解析辞書に反映することで、固有表現と固有表現同士の関係との抽出ルールを学習する (S33)。

【0049】

固有表現・関係抽出部 21 は、解析辞書 DB 20 に格納される解析辞書を用いて、文書 DB 12 に格納される抽出用文書データから、タグ、およびタグとタグの関係をそれぞれ抽出する (S34)。

【0050】

固有表現・関係抽出結果の表示部 22 は、S34 で抽出されたタグ、およびタグ同士の関係を分類名ごとにグループ化した抽出結果の表示画面 G2 を表示装置に表示する (S35)。

【0051】

固有表現・関係抽出結果の表示部 22 は、学習文書 DB 18 に格納された学習文書と S34 での抽出結果とを照合する。この照合により、固有表現・関係抽出結果の表示部 22 は、固有表現・関係抽出の学習部 19 により学習文書として生成されたが、S34 で当該抽出用文書データから抽出されなかったタグ、およびタグとタグの関係を特定し、この特定した結果を示す表示画面 G5 を表示装置に表示する (S36)。

上記の、学習文書として生成されたが、抽出用文書データから抽出されなかったタグ、およびタグとタグの関係は、例えば、固有表現・関係抽出の学習部 19 による解析辞書への学習の不具合、ここでは解析辞書に反映させる定義の欠落などに起因して生ずる。

【0052】

図 14 は、実施形態に係る固有表現抽出装置 10 による、抽出されなかったタグ、およびタグ同士の関係の表示画面 G5 の一例を示す図である。

図 14 に示した例では、点線で囲まれる、分類 B に属する 1 つ目のタグ、分類 C に属する 1 つ目および 3 つ目のタグは、学習文書に含まれていたが抽出用文書データから抽出されなかったタグとして示される。

【0053】

また、図 14 に示した例では、点線で示される、(1) 分類 A に属する 2 つ目のタグと分類 B に属する 1 つ目のタグとの間の関係、(2) 分類 B に属する 3 つ目のタグと分類 C に属する 3 つ目のタグとの間の関係、および (3) 分類 C に属する 2 つ目のタグと分類 D に属する 3 つ目のタグとの間の関係は、学習文書に含まれていたが抽出用文書データから抽出されなかった関係として示される。

第 3 の処理により、タグ、およびタグとタグの関係の抽出結果の抽出漏れを容易に確認することができる。

【0054】

(第 4 の処理)

次に、固有表現抽出装置 10 の第 4 の処理について説明する。

図 15 は、実施形態に係る固有表現抽出装置 10 の第 4 の処理手順の一例を示すフローチャートである。

第 4 の処理では、第 3 の処理で説明した S31 ~ S35 までの処理がなされる (S41 ~ S45)。

【0055】

そして、固有表現・関係抽出結果の表示部 22 は、学習文書 DB 18 に格納された学習文書と S44 (S34 と同様) での抽出結果とを照合する。この照合により、固有表現・関係抽出結果の表示部 22 は、固有表現・関係抽出の学習部 19 により学習文書として生成されておらず、直近の学習された解析辞書にも定義されていないが、S44 で当該抽出用文書データから抽出されたタグ、およびタグとタグの関係を特定し、この特定した結果

10

20

30

40

50

の表示画面 G 6 を表示装置に表示する (S 4 6)。

上記の、学習文書として生成されなかったが、抽出用文書データから抽出されたタグ、およびタグとタグの関係は、例えば、固有表現・関係抽出の学習部 1 9 による解析辞書への学習の不具合、ここでは解析辞書に対する不必要な定義の追加などに起因して生ずる。

【 0 0 5 6 】

図 1 6 は、実施形態に係る固有表現抽出装置 1 0 による、学習されなかったが抽出されたタグ、およびタグ同士の関係の表示画面 G 6 の一例を示す図である。

図 1 6 に示した例では、二重線で囲まれる、分類 C に属する 1 つ目のタグは、学習文書には含まれなかったが抽出用文書データから抽出されたタグとして示される。

【 0 0 5 7 】

また、図 1 6 に示した例では、二重線で示される、(1) 分類 A に属する 2 つ目のタグと分類 B に属する 1 つ目のタグとの間の関係、(2) 分類 B に属する 1 つ目のタグと分類 C に属する 1 つ目のタグとの間の関係が示される。これらの関係は、学習文書には含まれなかったが抽出用文書データから抽出された関係として示される。

第 4 の処理により、タグ、およびタグとタグの関係の抽出結果の誤抽出を容易に確認することができる。

【 0 0 5 8 】

以上説明したように、実施形態に係る固有表現抽出装置は、学習データにおけるタグ同士の関係を一括で登録したり、抽出元文書を容易に表示したり、抽出結果の抽出漏れ又は誤抽出を容易に確認できたりするので、文書からの固有表現抽出の精度を向上させることができる。

【 0 0 5 9 】

本発明のいくつかの実施形態を説明したが、これらの実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。これら新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。これら実施形態やその変形は、発明の範囲や要旨に含まれるとともに、特許請求の範囲に記載された発明とその均等の範囲に含まれる。

【 0 0 6 0 】

また、各実施形態に記載した手法は、計算機 (コンピュータ) に実行させることができるプログラム (ソフトウェア手段) として、例えば磁気ディスク (フロッピー (登録商標) ディスク、ハードディスク等)、光ディスク (C D - R O M、D V D、M O 等)、半導体メモリ (R O M、R A M、フラッシュメモリ等) 等の記録媒体に格納し、また通信媒体により伝送して頒布することもできる。なお、媒体側に格納されるプログラムには、計算機に実行させるソフトウェア手段 (実行プログラムのみならずテーブルやデータ構造も含む) を計算機内に構成させる設定プログラムをも含む。本装置を実現する計算機は、記録媒体に記録されたプログラムを読み込み、また場合により設定プログラムによりソフトウェア手段を構築し、このソフトウェア手段によって動作が制御されることにより上述した処理を実行する。なお、本明細書でいう記録媒体は、頒布用に限らず、計算機内部あるいはネットワークを介して接続される機器に設けられた磁気ディスクや半導体メモリ等の記憶媒体を含むものである。

【 符号の説明 】

【 0 0 6 1 】

1 0 ... 固有表現抽出装置、1 1 ... 文書受付部、1 2 ... 文書 D B、1 3 ... 固有表現・関係の学習データ編集部、1 4 ... オリジナル学習文書 D B、1 5 ... 学習データへの関係抽出ルール登録部、1 6 ... 関係抽出ルール D B、1 7 ... 学習データへの関係抽出ルール適用部、1 8 ... 学習文書 D B、1 9 ... 固有表現・関係抽出の学習部、2 0 ... 解析辞書 D B、2 1 ... 固有表現・関係抽出部、2 2 ... 固有表現・関係抽出結果の表示部。

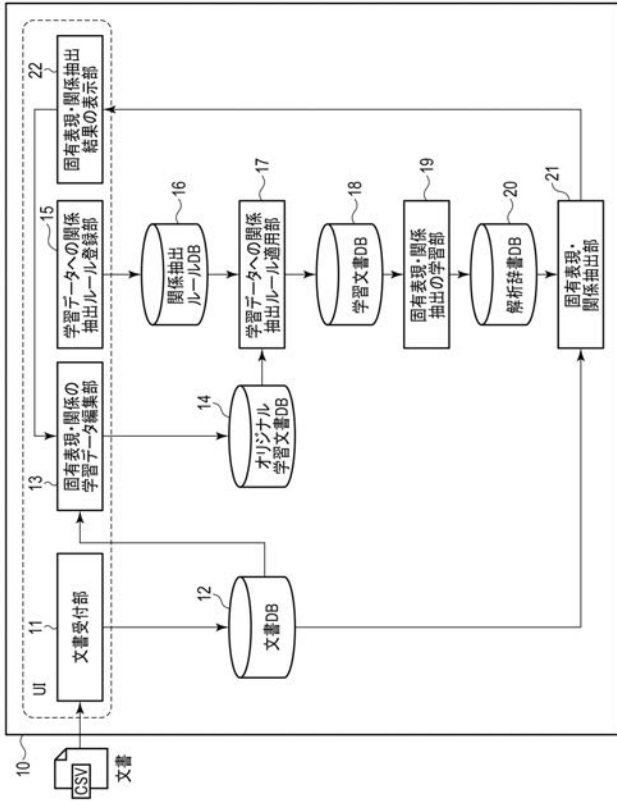
10

20

30

40

【 図 1 】



【 図 2 】

(文書DB)

コンテンツID	タイトル	本文	...
1
2
...

【 図 3 】

(オリジナル学習文書DB(固有表現))

タグID	コンテンツID	タグの種類	タグの値	...
1	1	人名	〇芝太郎	...
2	1	地名	〇芝町	...
...

【 図 4 】

(オリジナル学習文書DB(関係))

関係ID	第1のタグID	第1のタグの役割	第2のタグID	第2のタグの役割	...
1	1	住人	2	住んでいる地域	...
2	3	スポーツ	4	順位	...
...

【 図 6 】

(解析辞書DB)

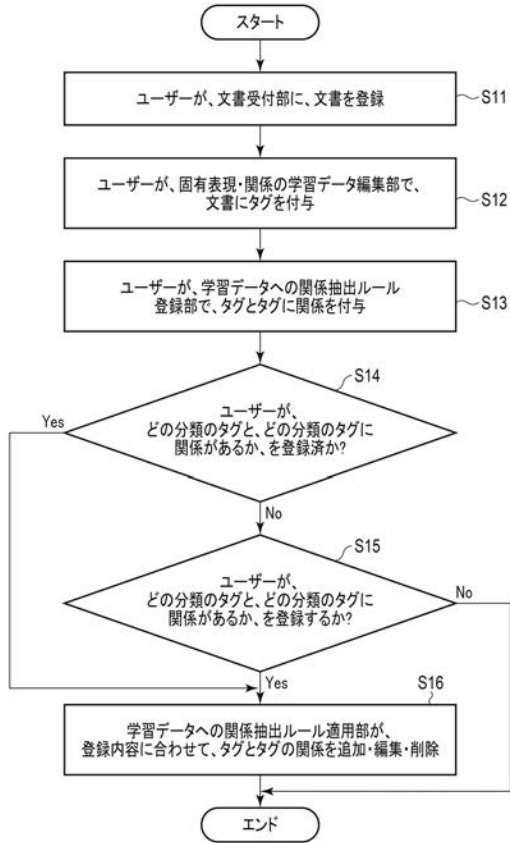
辞書ID	タグの種類	タグの特徴	...
1	人名	バイナリデータ	...
2	地名	バイナリデータ	...
...

【 図 5 】

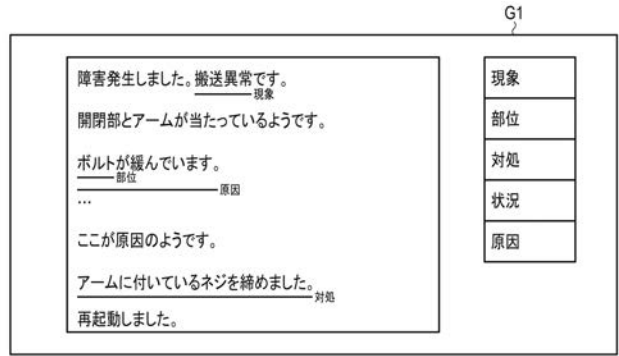
(関係抽出ルールDB(固有表現))

ルールID	第1のタグの種類	第2のタグの種類	第1のタグの役割	第2のタグの役割	...
1	人名	地名	住人	住んでいる地域	...
2	スポーツ	順位	競技名	競技結果	...
...

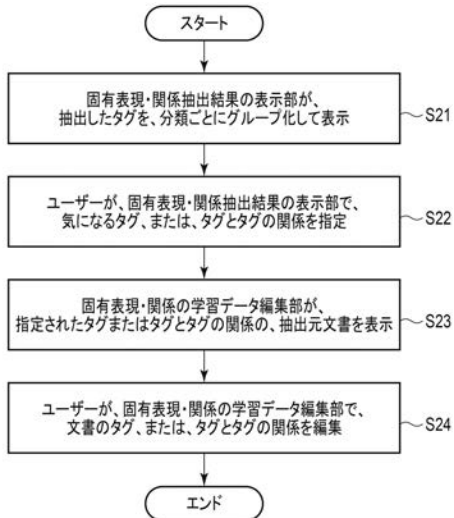
【 図 7 】



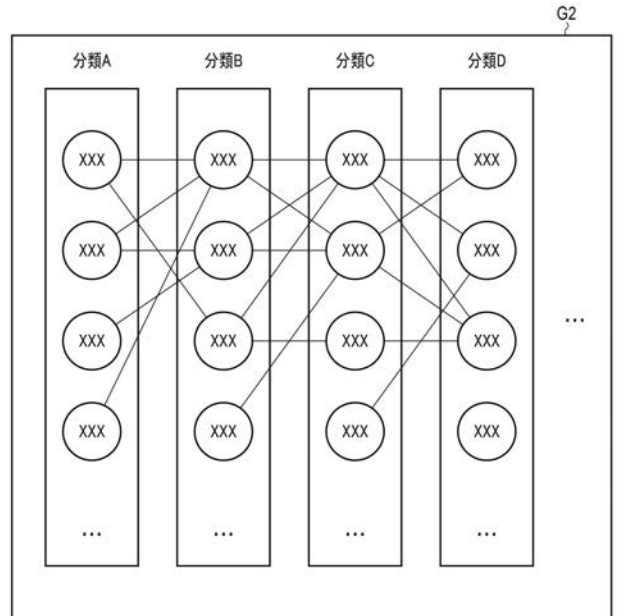
【 図 8 】



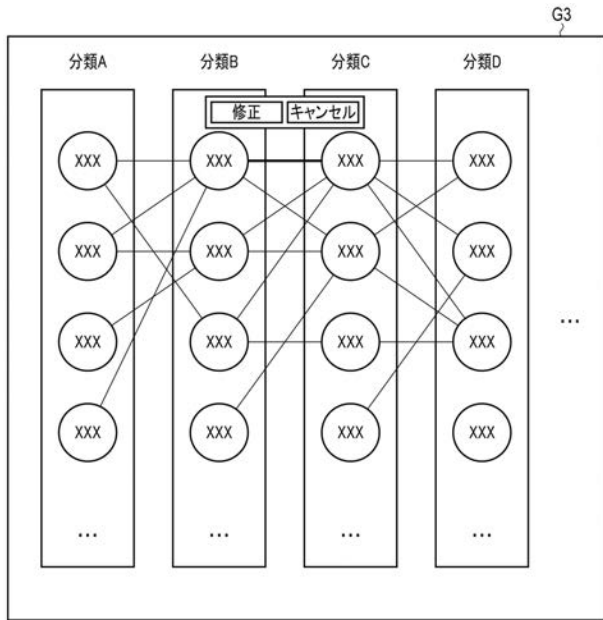
【 図 9 】



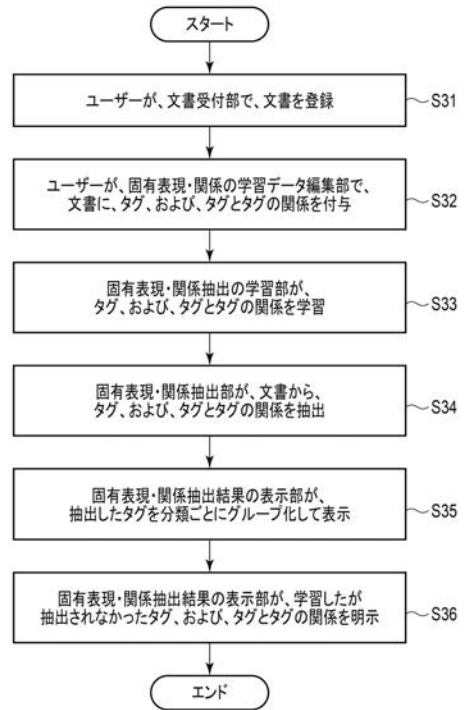
【 図 10 】



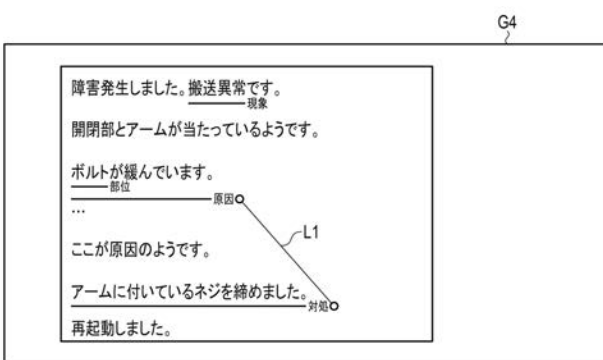
【 図 1 1 】



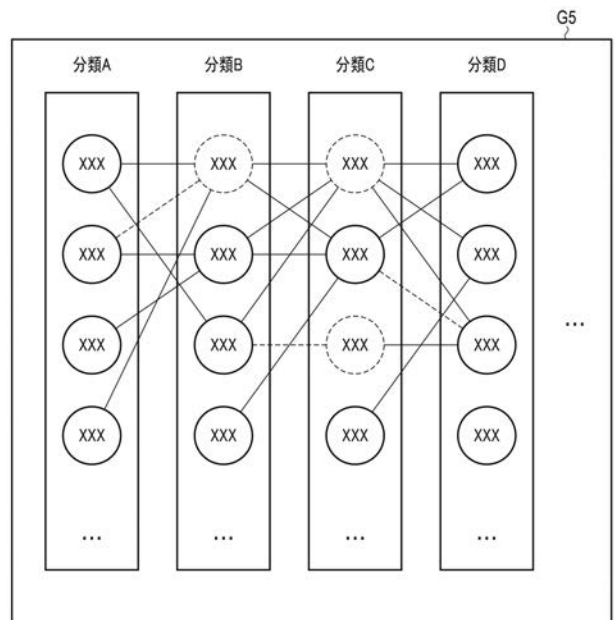
【 図 1 2 】



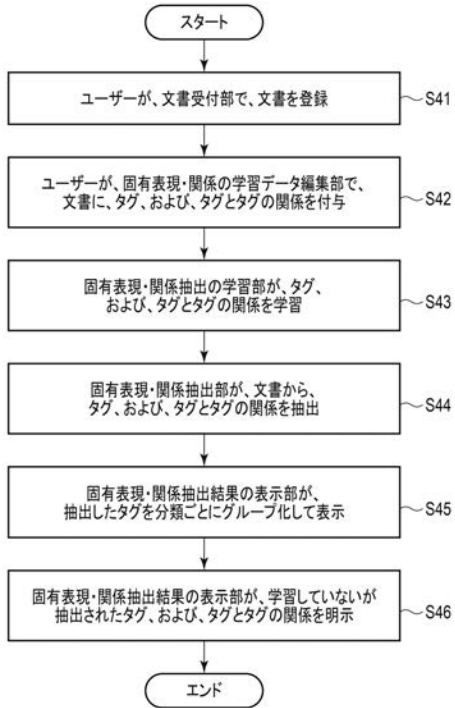
【 図 1 3 】



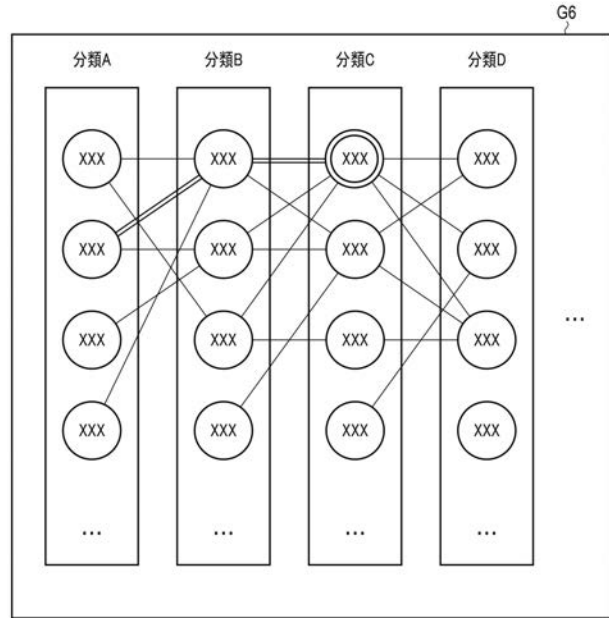
【 図 1 4 】



【 図 1 5 】



【 図 1 6 】



フロントページの続き

(74)代理人 100189913

弁理士 鷗飼 健

(72)発明者 飛田 義賢

神奈川県川崎市幸区堀川町7番地34 東芝デジタルソリューションズ株式会社内

(72)発明者 鈴木 優

神奈川県川崎市幸区堀川町7番地34 東芝デジタルソリューションズ株式会社内

Fターム(参考) 5B091 AA15 AB06 CA01 CC03 EA01