

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6476647号
(P6476647)

(45) 発行日 平成31年3月6日(2019.3.6)

(24) 登録日 平成31年2月15日(2019.2.15)

(51) Int.Cl. F I
H03M 7/40 (2006.01) H03M 7/40

請求項の数 9 (全 29 頁)

<p>(21) 出願番号 特願2014-167895 (P2014-167895) (22) 出願日 平成26年8月20日 (2014. 8. 20) (65) 公開番号 特開2016-46602 (P2016-46602A) (43) 公開日 平成28年4月4日 (2016. 4. 4) 審査請求日 平成29年5月11日 (2017. 5. 11)</p> <p>前置審査</p>	<p>(73) 特許権者 000005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番 1号 (74) 代理人 110002147 特許業務法人酒井国際特許事務所 (72) 発明者 片岡 正弘 神奈川県川崎市中原区上小田中4丁目1番 1号 富士通株式会社内 (72) 発明者 東 秀人 神奈川県川崎市中原区上小田中4丁目1番 1号 富士通株式会社内 (72) 発明者 小澤 崇記 神奈川県川崎市中原区上小田中4丁目1番 1号 富士通株式会社内</p> <p style="text-align: right;">最終頁に続く</p>
--	---

(54) 【発明の名称】 圧縮プログラム、圧縮装置、圧縮方法、伸長プログラム、伸長装置および伸長方法

(57) 【特許請求の範囲】

【請求項1】

コンピュータに、
 数値を示す、少なくとも1つの数字を含む数値文字列と、該数値文字列以外の単語とを含む入力データから、前記数値文字列を抽出し、
 数値文字列と圧縮符号とを対応付けて記憶する辞書を参照し、前記抽出された数値文字列に対応する圧縮符号に変換し、
前記辞書は、所定の複数の数値文字列を予め記憶する第1の辞書と、前記第1の辞書に登録されている数値文字列を記憶する第2の辞書とであり、
前記抽出された数値文字列が、前記第1の辞書に登録されており、かつ、該数値文字列に対応する圧縮符号が前記第1の辞書に登録されていない場合、該数値文字列に圧縮符号を付与して前記第1の辞書における該数値文字列に対応付けて該圧縮符号を登録し、
前記変換する処理は、前記登録された圧縮符号に対応する数値文字列が抽出された場合に、前記抽出された数値文字列を、前記登録された圧縮符号に変換する
 処理を実行させることを特徴とする圧縮プログラム。

【請求項2】

前記第1の辞書は、抽出され得る数値文字列とともに、数値の大きさの順番に応じて割り当てられるコードがさらに対応付けて記憶された辞書であることを特徴とする請求項1に記載の圧縮プログラム。

【請求項3】

前記抽出する処理で数値文字列が抽出された際に、前記数値文字列の検索対象となる複数の圧縮ファイルの各々が、前記数値文字列に対応する圧縮符号を格納するか否かを、前記圧縮ファイル毎に情報ビット列により示すインデックスを参照することにより、前記圧縮符号を格納する圧縮ファイルを特定し、前記情報ビット列において該圧縮ファイルに対応する情報ビットが未格納を示す値の場合、格納を示す値に更新する処理をさらにコンピュータに実行させることを特徴とする請求項 1 または 2 に記載の圧縮プログラム。

【請求項 4】

数値を示す、少なくとも 1 つの数字を含む数値文字列と、該数値文字列以外の単語とを含む入力データから、前記数値文字列を抽出する抽出部と、

数値文字列と圧縮符号とを対応付けて記憶する辞書を参照し、前記抽出された数値文字列に対応する圧縮符号に変換する変換部とを有し、

前記辞書は、所定の複数の数値文字列を予め記憶する第 1 の辞書と、前記第 1 の辞書に登録されている数値文字列を記憶する第 2 の辞書とであり、

前記抽出された数値文字列が、前記第 1 の辞書に登録されており、かつ、該数値文字列に対応する圧縮符号が前記第 1 の辞書に登録されていない場合、該数値文字列に圧縮符号を付与して前記第 1 の辞書における該数値文字列に対応付けて該圧縮符号を登録し、

前記変換する処理は、前記登録された圧縮符号に対応する数値文字列が抽出された場合に、前記抽出された数値文字列を、前記登録された圧縮符号に変換する

ことを特徴とする圧縮装置。

【請求項 5】

コンピュータが、

数値を示す、少なくとも 1 つの数字を含む数値文字列と、該数値文字列以外の単語とを含む入力データから、前記数値文字列を抽出し、

数値文字列と圧縮符号とを対応付けて記憶する辞書を参照し、前記抽出された数値文字列に対応する圧縮符号に変換し、

前記辞書は、所定の複数の数値文字列を予め記憶する第 1 の辞書と、前記第 1 の辞書に登録されている数値文字列を記憶する第 2 の辞書とであり、

前記抽出された数値文字列が、前記第 1 の辞書に登録されており、かつ、該数値文字列に対応する圧縮符号が前記第 1 の辞書に登録されていない場合、該数値文字列に圧縮符号を付与して前記第 1 の辞書における該数値文字列に対応付けて該圧縮符号を登録し、

前記変換する処理は、前記登録された圧縮符号に対応する数値文字列が抽出された場合に、前記抽出された数値文字列を、前記登録された圧縮符号に変換する

処理を実行することを特徴とする圧縮方法。

【請求項 6】

コンピュータに、

数値を示す、少なくとも 1 つの数字を含む数値文字列とともに数値の大きさの順番に応じて割り当てられるコードと圧縮符号とを対応付けて記憶する辞書を用いて、圧縮ファイルから抽出した圧縮符号をコードに変換し、

前記変換したコードと、数値文字列の検索範囲に対応する 1 以上のコードとを比較することで、前記変換したコードが前記検索範囲内に含まれるか否かを判定し、

前記検索範囲内に含まれると判定した場合に、前記変換したコードを数値文字列に伸長して表示する

処理を実行させることを特徴とする伸長プログラム。

【請求項 7】

前記数値文字列の検索対象となる複数の圧縮ファイルの各々が、前記数値文字列に対応する圧縮符号を格納するか否かを、前記圧縮ファイル毎に情報ビット列により示すインデックスを参照することにより、前記検索範囲内の数値文字列を含む圧縮ファイルを特定する処理をさらに実行することを特徴とする請求項 6 に記載の伸長プログラム。

【請求項 8】

数値を示す、少なくとも 1 つの数字を含む数値文字列の大きさの順番に応じて割り当て

10

20

30

40

50

られるコードと圧縮符号とを対応付けて記憶する辞書を用いて、圧縮ファイルから抽出した圧縮符号をコードに変換する変換部と、

前記変換したコードと、数値文字列の検索範囲に対応する1以上のコードとを比較することで、前記変換したコードが前記検索範囲内に含まれるか否かを判定する判定部と、

前記検索範囲内に含まれると判定した場合に、前記変換したコードを数値文字列に伸長して表示する伸長部と

を有することを特徴とする伸長装置。

【請求項9】

コンピュータが、

数値を示す、少なくとも1つの数字を含む数値文字列の大きさの順番に応じて割り当てられるコードと圧縮符号とを対応付けて記憶する辞書を用いて、圧縮ファイルから抽出した圧縮符号をコードに変換し、

前記変換したコードと、数値文字列の検索範囲に対応する1以上のコードとを比較することで、前記変換したコードが前記検索範囲内に含まれるか否かを判定し、

前記検索範囲内に含まれると判定した場合に、前記変換したコードを数値文字列に伸長して表示する処理を実行することを特徴とする伸長方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、圧縮プログラムおよび伸長プログラム等に関する。

【背景技術】

【0002】

圧縮対象のテキストファイルから数値を抽出し、抽出した数値を数字単位で圧縮する技術が存在する。かかる技術では、0～9までの各数字に圧縮符号を割り当て、圧縮対象のテキストファイルから抽出した数値に含まれる各数字を圧縮符号に変換する。さらに、変換された圧縮符号に、例えば全角または半角の指定、カンマの有無、小数点の有無、有効桁数等の数値に関する情報を表す符号を付加することで数値を圧縮符号に変換する。

【先行技術文献】

【特許文献】

【0003】

【特許文献1】国際公開第2008/047432号

【特許文献2】特開2013-150041号公報

【特許文献3】特開平05-174064号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら、数値を数字単位で圧縮すると、桁数分の数字に対応した圧縮符号が必要となり、数値全体の圧縮符号が長くなるため、圧縮率が低下する場合がある。

【0005】

一つの側面では、数値を圧縮する場合の圧縮率を向上させる圧縮プログラムおよび伸長プログラム等を提供することを目的とする。

【課題を解決するための手段】

【0006】

第1の案では、圧縮プログラムは、コンピュータに、入力されたデータから、数値を示す、少なくとも1つの数字を含む数値文字列を抽出する。圧縮プログラムは、コンピュータに、数値文字列と圧縮符号とを対応付けて記憶する辞書を参照し、抽出された数値文字列に対応する圧縮符号に変換する処理を実行させる。

【発明の効果】

【0007】

本発明の1実施態様によれば、圧縮処理時に数値に割り当てる符号長を短くできるとい

10

20

30

40

50

う効果を奏する。

【図面の簡単な説明】

【0008】

【図1】図1は、実施例1の数値の圧縮処理を説明するための図である。

【図2】図2は、実施例1の圧縮処理の全体の流れについて説明するための図である。

【図3A】図3Aは、参考例1の符号長を説明するための図である。

【図3B】図3Bは、実施例1の符号長を説明するための図である。

【図4】図4は、参考例1および実施例1の符号長の比較を説明するための図である。

【図5】図5は、情報処理装置のシステム構成の例を説明するための図である。

【図6】図6は、実施例1の圧縮処理に係るシステム構成の例を示す図である。

10

【図7】図7は、ビットフィルタの数値部の第1の例を示す図である。

【図8】図8は、数値文字列に割り当てられる単語コードのデータ構造の例を示す図である。

【図9】図9は、ビットフィルタの単語部の例を示す図である。

【図10】図10は、動的辞書の一例を示す図である。

【図11】図11は、ビットフィルタの数値部の第2の例を示す図である。

【図12】図12は、実施例1の圧縮処理の流れの例を示す図である。

【図13】図13は、実施例2の圧縮処理に係るシステム構成の例を示す図である。

【図14】図14は、ビットマップ型全文インデックスの一例を示す図である。

【図15】図15は、実施例2の圧縮処理の流れを示す図である。

20

【図16】図16は、実施例3の大小比較処理に係るシステム構成の一例を示す図である。

【図17】図17は、圧縮ファイル選択の処理の流れの例を示す第1の図である。

【図18】図18は、圧縮ファイル選択の処理の流れの例を示す第2の図である。

【図19】図19は、伸長辞書の構造を説明するための図である。

【図20】図20は、数値文字列の大小比較を説明するための図である。

【図21】図21は、実施例3の大小比較処理の全体の流れを示す図である。

【図22】図22は、実施例3の圧縮ファイル選択処理の流れを示す図である。

【図23】図23は、実施例3の数値文字列の単語コード抽出処理の流れを示す図である。

30

【図24】図24は、実施例1～3の情報処理装置のハードウェア構成を示す図である。

【図25】図25は、コンピュータで動作するプログラムの構成例を示す図である。

【図26】図26は、実施形態のシステムにおける装置の構成例を示す図である。

【発明を実施するための形態】

【0009】

以下に、本願の開示する圧縮プログラムの実施例を図面に基づいて詳細に説明する。なお、この実施例によりこの権利範囲が限定されるものではない。各実施例は、処理内容を矛盾させない範囲で適宜組み合わせることが可能である。

【実施例1】

【0010】

40

(実施例1の圧縮処理)

図1を用いて、実施例1の情報処理装置100による圧縮処理について説明する。図1は、実施例1の数値の圧縮処理を説明するための図である。図1の例のように、情報処理装置100は、圧縮処理の対象である対象ファイル10に含まれる「He pays 1,200 yen . . .」を「He」「pays」「1,200」「yen」のように数値または単語ごとに区切り、各数値または各単語を取得する。以下、1以上の数字を有する数値を数値文字列と呼ぶ。数値文字列は、数字以外にプラス・マイナス、カンマおよび小数点等の符号を含んでもよい。情報処理装置100は、取得した単語のうち数値文字列「1,200」を抽出してビットフィルタの数値部121aに出力する。ビットフィルタの数値部121aは、各数値文字列に対して単語コードと、圧縮符号とを対応付ける辞書である。ビットフィルタの数値部

50

1 2 1 a は、それぞれの数値文字列に対応する単語コード 1 1 があらかじめ登録されている。例えば、ビットフィルタの数値部 1 2 1 a には、整数「0」「1」「2」「3」「4」... に対応する単語コード「B0000h」「B00010h」「B00020h」「B00030h」「B00040h」... が順番にあらかじめ登録されている。ビットフィルタの数値部 1 2 1 a のデータ構造に関する詳細は後述する。

【 0 0 1 1 】

なお、数値文字列以外の「He」「pays」「yen」等の単語は、後述するようにビットフィルタの単語部 1 2 1 b に出力される。また、以降、ビットフィルタの数値部 1 2 1 a をビットフィルタ 1 2 1 a と表記する場合があります。ビットフィルタの単語部 1 2 1 b をビットフィルタ 1 2 1 b と表記する。ビットフィルタ 1 2 1 a およびビットフィルタ 1 2 1 b 10
の詳細は後述する。

【 0 0 1 2 】

次いで、情報処理装置 1 0 0 は、ビットフィルタ 1 2 1 a から数値文字列「1,200」に対応する単語コード「B04B01h」を取得する。次いで、情報処理装置 1 0 0 は、取得した単語コード「B04B01h」と、動的辞書 1 2 2 に登録した順に動的に付される圧縮符号「A005h」とを対応づけて動的辞書 1 2 2 に登録する。なお、圧縮符号単語コード「B04B01h」および「A005h」等の末尾に示す「h」は、16進数で表記されていることを示す符号である。

【 0 0 1 3 】

次いで、情報処理装置 1 0 0 は、動的辞書 1 2 2 において動的に付された圧縮符号「A005h」を、数値文字列「1,200」および単語コード「B04B01h」に対応づけてビットフィルタ 1 2 1 a に登録する。そして、情報処理装置 1 0 0 は、ビットフィルタ 1 2 1 a を基にして数値文字列「1,200」に対応する圧縮符号「A005h」を取得し、圧縮ファイル 1 2 20
に出力する。

【 0 0 1 4 】

また、情報処理装置 1 0 0 は、次回、対象ファイル 1 0 に数値文字列「1,200」が出現した場合、既にビットフィルタ 1 2 1 a に登録されている圧縮符号「A005h」を取得し、圧縮ファイル 1 2 に出力する。

【 0 0 1 5 】

図 2 を用いて実施例 1 の圧縮処理全体の流れについて説明する。図 2 は、実施例 1 の圧縮処理の全体の流れについて説明するための図である。図 2 の例のように、情報処理装置 1 0 0 は、対象ファイル 1 0 から単語または数値文字列を抽出する。例えば、情報処理装置 1 0 0 は、対象ファイル 1 0 から単語「pays」を抽出した場合、単語「pays」を符号化部 1 A に格納する。情報処理装置 1 0 0 は、ビットフィルタ（単語部）1 2 1 b から「pays」の圧縮符号を取得し、記憶領域 1 B に格納する。 30

【 0 0 1 6 】

一方、情報処理装置 1 0 0 は、対象ファイルから数値文字列「1,200」を抽出した場合、数値文字列「1,200」を符号化部 2 A に格納する。情報処理装置 1 0 0 は、ビットフィルタ（数値部）1 2 1 a から数値文字列「1,200」の圧縮符号を取得し、記憶領域 2 B に格納する。情報処理装置 1 0 0 は、記憶領域 1 B および記憶領域 2 B に格納した圧縮符号を圧縮ファイル 1 2 に出力する。 40

【 0 0 1 7 】

このように、情報処理装置 1 0 0 は、対象ファイル 1 0 から単語を抽出した場合はビットフィルタ 1 2 1 b を用いて単語を圧縮符号に変換し、数値文字列を抽出した場合はビットフィルタ 1 2 1 a を用いて数値文字列を圧縮符号に変換する。

【 0 0 1 8 】

（参考例 1 および実施例 1 の比較）

上述したように実施例 1 では、数値文字列全体を一つの単位として圧縮する。一方、参考例 1 では、数値文字列に含まれる個々の数字を一つの単位として圧縮するものとする。参考例 1 のように数値文字列を数字単位で圧縮すると、ヘッダ、カンマの有無、小数点の 50

有無等の付加情報を圧縮符号に加えることになる。また、参考例 1 では、数値文字列が大きくなると数値文字列の桁数に比例して圧縮符号の符号長が長くなる。これに対して、実施例 1 のように数値文字列単位で圧縮すれば、数値の桁数に関係なく、固定長の圧縮符号が付与され、参考例 1 と比べて安定して圧縮符号の符号長を短くすることができる。

【 0 0 1 9 】

また、実施例 1 では、情報処理装置 1 0 0 は、カンマの有無や小数点の有無等によって、同じ大きさの数値文字列に異なる圧縮符号を付与してもよい。例えば、情報処理装置 1 0 0 は、カンマが付与されていない「1200」と、カンマが付与された「1,200」とで異なる圧縮符号を付与してもよい。

【 0 0 2 0 】

図 3 A、図 3 B および図 4 を用いて、参考例 1 および実施例 1 で割り当てられる符号長の比較を説明する。図 3 A は、参考例 1 の符号長を説明するための図である。図 3 A の例のように、数値文字列 1 「1,200」を参考例 1 のように数字単位で圧縮する場合、圧縮符号 1 は情報部と数字部とを有する。情報部は、6 ビットのヘッダ情報、1 ビットのカンマの有無情報、6 ビットの有効桁数情報を有する。一方、数字部は、数値文字列の各桁に圧縮符号が割り当てられ、6 ビット×数値文字列の桁数がビット数となる。例えば、数値文字列が 4 桁の場合、数字部のビット数が 2 4 ビットとなる。したがって、圧縮符号 1 の符号長は、数値文字列が 4 桁の場合、情報部と数値部を合わせると 3 7 ビットとなる。

【 0 0 2 1 】

図 3 B は、実施例 1 の符号長を説明するための図である。実施例 1 に係る情報処理装置 1 0 0 は、例えば、対象ファイルから数値文字列を抽出し、抽出した数値文字列に対して 1 6 ビットの圧縮符号を割り当てる。すなわち、情報処理装置 1 0 0 は、対象のファイルから抽出した数値文字列に対し、抽出した順番に 1 6 ビットの圧縮符号「A000h」「A001h」「A002h」「A003h」... を割り当てる。

【 0 0 2 2 】

実施例 1 においては、対象のファイルから抽出した数値文字列が圧縮符号の符号長が 1 6 ビットの固定長となる。このため、実施例 1 に係る情報処理装置 1 0 0 は、参考例 1 と比べて安定して短い符号長の圧縮符号を割り当てることができる。

【 0 0 2 3 】

例えば、図 3 B の例のように、数値文字列 1 「1,200」を実施例 1 のように数値単位で圧縮する場合、情報処理装置 1 0 0 は、数値文字列 1 「1,200」に圧縮符号 2 として 1 6 ビットの固定長の圧縮符号「101000000000101」(A005h) を付与する。

【 0 0 2 4 】

図 4 は、参考例 1 および実施例 1 の符号長の比較を説明するための図である。表 1 は、数値文字列の桁数に対応する符号長を示す図である。例えば、圧縮する数値文字列が 1 桁の場合に、実施例 1 のように数値文字列単位で圧縮すると符号長が 1 6 ビットとなり、参考例 1 のように数字単位で圧縮すると、情報部が 1 3 ビットで数字部が 6 ビットとなり符号長が全体で 1 9 ビットとなる。例えば、圧縮する数値文字列が 3 桁の場合に、実施例 1 のように数値文字列単位で圧縮すると符号長が 1 6 ビットとなり、参考例 1 のように数字単位で圧縮すると、情報部が 1 3 ビットで数字部が 1 8 ビットとなり符号長が全体で 3 1 ビットとなる。また、数値文字列が 5 桁の場合に、数値文字列単位で圧縮すると符号長が 1 6 ビットとなり、数字文字列単位で圧縮すると符号長が 4 3 ビットとなる。なお、参考例 1 において数字 1 桁に割り当てられる符号長は 6 ビットであるものとする。

【 0 0 2 5 】

このように、参考例 1 のように数字単位で圧縮する場合、ヘッダ情報等を有する情報部がある分、圧縮符号の符号長が長くなる。また、数値文字列の桁数が大きくなると、各桁に所定長の圧縮符号が割り当てられ、さらに圧縮符号全体の符号長が長くなる。これに対して、実施例 1 のように数値文字列単位で圧縮する場合、例えば、対象ファイルに出現した順番に昇順に 1 6 ビットの固定長の圧縮符号を割り当てるので、数値の桁数に関係なく数

10

20

30

40

50

値文字列に割当てする符号長を安定して短く設定することができる。

【 0 0 2 6 】

(実施例 1 の圧縮処理に関する処理部の構成)

図 5 を用いて、情報処理装置 1 0 0 の圧縮部と記憶部との関係について説明する。図 5 は、情報処理装置のシステム構成の例を説明するための図である。図 5 の例に示すように、情報処理装置 1 0 0 の記憶部 1 2 0 は、圧縮部 1 1 0 と処理部 1 5 0 とに接続される。記憶部 1 2 0 は例えば、R A M (Random Access Memory)、R O M (Read Only Memory)、フラッシュメモリなどの半導体メモリ素子、ハードディスクや光ディスクなどの記憶装置に対応する。

【 0 0 2 7 】

また、情報処理装置 1 0 0 は、圧縮部 1 1 0 と、処理部 1 5 0 とを有する。圧縮部 1 1 0 および処理部 1 5 0 の機能は例えば、C P U (Central Processing Unit) が所定のプログラムを実行することで実現することができる。また、圧縮部 1 1 0 および処理部 1 5 0 の機能は例えば、A S I C (Application Specific Integrated Circuit) や F P G A (Field Programmable Gate Array) などの集積回路により実現することができる。

【 0 0 2 8 】

圧縮部 1 1 0 は、入力されたデータから、数値を示す、少なくとも 1 つの数字を含む数値文字列を抽出する。圧縮部 1 1 0 は、数値文字列と圧縮符号とを対応付けて記憶する辞書を参照し、抽出された数値文字列に対応する圧縮符号に変換する。

【 0 0 2 9 】

辞書は、所定の複数の数値文字列を予め記憶する第 1 の辞書と、第 1 の辞書に登録されている数値文字列を記憶する第 2 の辞書とである。圧縮部 1 1 0 は、抽出された数値文字列が、第 1 の辞書に登録されており、かつ、該数値文字列に対応する圧縮符号が第 1 の辞書に登録されていない場合、該数値文字列に圧縮符号を付与して第 1 の辞書における該数値文字列に対応付けて該圧縮符号を登録する。なお、第 1 の辞書は、例えば、ビットフィルタであり、第 2 辞書は、例えば、動的辞書である。

【 0 0 3 0 】

第 1 の辞書は、抽出され得る数値文字列とともに、数値の大きさの順番に応じて割り当てられるコードがさらに対応付けて記憶された辞書である。

【 0 0 3 1 】

図 6 を用いて、実施例 1 の圧縮処理に係るシステム構成について説明する。図 6 は、実施例 1 の圧縮処理に係るシステム構成の一例を示す図である。図 6 の例に示されるように、情報処理装置 1 0 0 は、圧縮部 1 1 0 と、記憶部 1 2 0 とを有する。圧縮部 1 1 0 は、ファイルリード部 1 1 1、圧縮符号付与部 1 1 2 およびファイルライト部 1 1 3 を有する。記憶部 1 2 0 は、ビットフィルタ 1 2 1 および動的辞書 1 2 2 を有する。ビットフィルタ 1 2 1 は、数値部 1 2 1 a および単語部 1 2 1 b を有する。以下、実施例 1 の圧縮部 1 1 0 および記憶部 1 2 0 の構成について詳細に説明する。

【 0 0 3 2 】

(記憶部の各構成)

ビットフィルタの数値部 1 2 1 a について説明する。ビットフィルタ 1 2 1 a は、数値文字列に単語コードおよび圧縮符号を対応付けた辞書である。図 7 は、ビットフィルタの数値部の第 1 の例を示す図である。図 7 の例のように、ビットフィルタ 1 2 1 a は、2 グラムと、ビットマップと、数値文字列と、文字列長と、単語コードと、登録番号と、圧縮符号とを対応付ける。「2 グラム」は、各単語に含まれる連数字である。例えば「115」は、「11」「15」に対応する 2 グラムを有する。

【 0 0 3 3 】

「ビットマップ」は、2 グラムの文字列に対応するビットマップを表す。例えば、「115」は、数値文字列へのポイントによって、2 グラム「11」のビットマップ「0 _ 1 _ 0 _ 0 _ 0」と、2 グラム「15」のビットマップ「0 _ 0 _ 1 _ 0 _ 0」とに対応付けられ

10

20

30

40

50

る。例えば、情報処理装置 1 0 0 は、対象ファイルから「115」を取得した場合に、2 グラム「11」のビットマップ「0__1__0__0__0」と、2 グラム「15」のビットマップ「0__0__1__0__0」とを用いて数値文字列「115」にアクセスする。

【 0 0 3 4 】

「数値文字列」は、ビットフィルタ 1 2 1 a に登録された数値文字列である。「数値文字列」は、「0」「1」「2」...「100」「101」...「999」「1000」...のように数値文字列が連番で登録される。また、「数値文字列」は、3 桁ごとにカンマ「,」を有する数値文字列を含んでもよい。例えば、図 4 の例のようにビットフィルタ 1 2 1 a は、カンマを有する文字列として「1,000」「1,001」「1,002」...を有する。なお、「数値文字列」は、負の値や小数点以下の値を有してもよい。「文字列長」は、各数値文字列の桁数である。

10

【 0 0 3 5 】

「単語コード」は、各数値文字列に割り当てるコードである。「単語コード」は、数値の大きさの順に、昇順に割り当てられる。例えば、「単語コード」は、昇順に列挙された数値文字列「0」「1」「2」「3」「4」...に対して、単語コード「B00000h」「B00010h」「B00020h」「B00030h」...を昇順にそれぞれ割り当てられる。なお、6 桁の 1 6 進数で表される単語コードの末尾 1 桁は、半角 / 全角の別、カンマの有無、小数点の有無、正の値 / 負の値の別などを示す情報ビットである。単語コードの下 2 桁以上の領域が数値文字列に対応する。このように、各数値文字列に対して昇順に数値文字列の単語コードを割り当てることで、各数値文字列に対して、数値文字列の大きさに応じた単語コードが割り当てられる。

20

【 0 0 3 6 】

また、「単語コード」は、ビットフィルタ 1 2 1 a がカンマ、負の値、小数点以下の値を有する数値文字列を含む場合、それぞれに別々の単語コードが割り当てられる。例えば、図 7 の例のように、「1000」に対して単語コード「B03E80h」が付され、「1,000」に対して単語コード「B03E81h」が割り当てられる。なお、図 7 の例では、各項目のデータがレコードとして関連づけられて記憶されている例を示したが、上記説明において互いに関連づけられた項目どうしの関係が保たれれば、データは他の記憶のされ方をして構わない。後述する図 9 ~ 1 1 に示されるビットフィルタおよび動的辞書に関する詳細は後述する。

【 0 0 3 7 】

各数値文字列に割り当てられる単語コードの例に関して説明する。図 8 は、数値文字列に割り当てられる単語コードのデータ構造の例を示す図である。ビットフィルタ 1 2 1 a において、あらかじめ各数値文字列に対して、3 バイト、4 バイトまたは 5 バイトの単語コードが割り当てられる。図 8 の例には、3 バイト、4 バイトおよび 5 バイトの単語コードのコード体系が示される。

30

【 0 0 3 8 】

図 8 の例のように、単語コード c 1 は、3 バイトコードで「B0000h ~ B3FFFh」までのコード領域を有する。単語コード c 1 は、c 1 - 1 領域、c 1 - 2 領域および c 1 - 3 領域を有する。c 1 - 1 領域は、3 バイトコードであることを示す「101100」の固定ビットを有する。c 1 - 2 領域は、0 から 1 6 , 3 8 3 までの整数に対応するビットを有する。c 1 - 3 の「****」は、半角 / 全角の別、カンマの有無、小数点の有無、正の値 / 負の値の別などを示す情報ビットである。例えば、c 1 - 3 の最後の 1 ビットが「0」の場合、数値文字列にカンマが含まれないことを示し、c 1 - 3 の最後の 1 ビットが「1」の場合、数値文字列にカンマが含まれることを示す。c 2 - 3 および c 3 - 3 においても同様である。

40

【 0 0 3 9 】

単語コード c 2 は、4 バイトコードで「B400000h ~ B7FFFFFFh」までのコード領域を有する。単語コード c 2 は、c 2 - 1 領域、c 2 - 2 領域および c 2 - 3 領域を有する。c 2 - 1 領域は、4 バイトコードであることを示す「101101」の固定ビットを有する。c 2 - 2 領域は、1 6 , 3 8 4 から 1 , 0 4 8 , 5 7 5 までの整数に対応するビットを有する。

50

c 2 - 3の「****」は、c 1 - 3と同様に情報ビットである。単語コードc 3は、5バイトコードで「B80000h~BBFFFFh」までのコード領域を有する。c 3 - 1領域は、5バイトコードであることを示す「101001」の固定ビットを有する。c 3 - 2領域は、1, 0 4 8, 5 7 6から1, 0 7 3, 7 4 1, 8 2 3までの整数に対応するビットを有する。c 3 - 3の「****」は、c 1 - 3と同様に情報ビットである。なお、3バイトコード、4バイトコード、5バイトコードは、整数以外に小数点を有する数値、負の数値、カンマを有する数値等を含んでもよい。

【0040】

次に、ビットフィルタ1 2 1の単語部1 2 1 bについて説明する。ビットフィルタ1 2 1の単語部1 2 1 bは、基礎単語に単語コードおよび圧縮符号を対応付けた辞書である。図9は、ビットフィルタの単語部の例を示す図である。図9の例のように、ビットフィルタ1 2 1 bは、2グラムと、ビットマップと、基礎単語と、文字列長と、出現頻度と、圧縮符号と、単語コードと、登録番号とを対応付ける。「2グラム」「ビットマップ」に関しては、ビットフィルタ1 2 1 aと同じであるので説明を省略する。「基礎単語」は、ビットフィルタ1 2 1 bに登録された単語である。例えば、「基礎単語」とは、辞典、テキスト群等からあらかじめ抽出された約19万語の単語である。なお、「基礎単語」として登録される単語数は、任意の語数でよい。

10

【0041】

「文字列長」は、基礎単語の文字列の長さである。「出現頻度」は、頻度集計用のテキストファイル群において各基礎単語が出現した回数である。ここで、頻度集計用のテキストファイル群とは、対象ファイルとは別に用意された各基礎単語の出現頻度を集計するための1以上のテキストファイルである。

20

【0042】

基礎単語の出現頻度の集計について説明する。情報処理装置1 0 0は、頻度集計用のテキストファイルを読み込み、頻度集計用のテキストファイル群に存在する単語を適宜抽出してビットフィルタ1 2 1 bに登録する。さらに、情報処理装置1 0 0は、ビットフィルタ1 2 1 bに登録された各々の基礎単語に関し、頻度集計用のテキストファイル群において出現した回数を出現頻度としてカウントする。例えば、図9の例において、ビットフィルタ1 2 1 bは、基礎単語「able」が、頻度集計用のテキストファイル群において「785」回出現したことを示す。

30

【0043】

「圧縮符号」は、各基礎単語に割り当てられた圧縮符号である。情報処理装置1 0 0は、出現頻度が高い基礎単語に対してより短い符号長を有する圧縮符号を割り当てる。「単語コード」は、各単語に割り当てるコードである。「登録番号」は、後述する動的辞書1 2 2に圧縮符号が登録された際に、圧縮符号に一意に付される番号である。「登録番号」は、例えば、動的辞書1 2 2に登録された順番を示す。

【0044】

次に、動的辞書1 2 2について説明する。動的辞書1 2 2は、単語コードと、単語コードに付与された圧縮符号とを対応付ける辞書である。図10は、動的辞書の一例を示す図である。図10の例のように、動的辞書1 2 2は、圧縮符号とコードとを対応付ける。「コード」は、ビットフィルタ1 2 1 aまたはビットフィルタ1 2 1 bで取得された単語コードである。「圧縮符号」は、例えば動的辞書1 2 2に単語コードが登録された順番に単語コードに昇順に付与される固定長の圧縮符号である。

40

【0045】

例えば、情報処理装置1 0 0が、対象ファイルから「order」「box」「1000」「him」...の順番に単語および数値文字列を抽出した場合、各々の単語および数値文字列の単語コードに対して、圧縮符号「A000h」「A001h」「A002h」「A003h」...を割当てる。そして、情報処理装置1 0 0は、各々の単語および数値文字列に対して割り当てた圧縮符号を動的辞書1 2 2に登録する。例えば、情報処理装置1 0 0は、オフセット「0x0000h」の位置に圧縮符号「A000h」と「order」の単語コードとを対応付ける。また、情報処理装置1 0

50

0 は、オフセット「0x0001h」の位置に圧縮符号「A001h」と「box」の単語コードとを対応付ける。また、情報処理装置 1 0 0 は、オフセット「0x0002h」の位置に圧縮符号「A002h」と「1000」の単語コードとを対応付ける。なお、図 1 0 において「c (単語)」の表記は、カッコ内の数値文字列または単語に対応する単語コードを表す。例えば、「c (order)」は、「order」の単語コードを表す。

【 0 0 4 6 】

(圧縮部の各構成)

圧縮部 1 1 0 の各構成について説明する。ファイルリード部 1 1 1 は、対象ファイルを読みだして、対象ファイルから単語および数値文字列を抽出する処理部である。ファイルリード部 1 1 1 は、対象ファイルの文字列中の空白文字によって文字列を単語または数値文字列ごとに区切り、対象ファイルから各単語および各数値文字列を抽出する。ファイルリード部 1 1 1 は、抽出した単語および数値文字列を圧縮符号付与部 1 1 2 に出力する。

10

【 0 0 4 7 】

圧縮符号付与部 1 1 2 は、対象ファイルから抽出された単語および数値文字列に圧縮符号を付与する処理部である。圧縮符号付与部 1 1 2 は、ファイルリード部 1 1 1 から数値文字列を受け付けると、ビットフィルタ 1 2 1 a に登録されている「数値文字列」にアクセスする。

【 0 0 4 8 】

圧縮符号付与部 1 1 2 は、アクセスした「数値文字列」に対応する圧縮符号がビットフィルタ 1 2 1 a に既に登録されていた場合、ビットフィルタ 1 2 1 a から数値文字列に対応する圧縮符号を取得してファイルライト部 1 1 3 に出力する。

20

【 0 0 4 9 】

一方、圧縮符号付与部 1 1 2 は、アクセスした数値文字列に対応する圧縮符号がビットフィルタ 1 2 1 a に登録されていない場合、ビットフィルタ 1 2 1 a から数値文字列に対応する単語コードを取得する。次いで、圧縮符号付与部 1 1 2 は、ビットフィルタ 1 2 1 a から取得した単語コードを、動的辞書 1 2 2 への登録順に付与される圧縮符号に対応付けて動的辞書 1 2 2 に登録する。

【 0 0 5 0 】

図 1 0 を用いて、動的辞書 1 2 2 への圧縮符号の登録の具体例について説明する。図 1 0 の例のように圧縮符号付与部 1 1 2 は、3 番目に登録された「1000」の単語コードに圧縮符号「A002h」を対応づけて動的辞書 1 2 2 に登録する。また、圧縮符号付与部 1 1 2 は、6 番目に登録された「1,200」の単語コードに圧縮符号「A005h」を対応づけて動的辞書 1 2 2 に登録する。このように、圧縮符号付与部 1 1 2 は、各数値文字列の単語コードに対して動的辞書 1 2 2 へ登録された順番に対応する圧縮符号を付与し、数値文字列の単語コードと、付与された圧縮符号とを対応付けて動的辞書 1 2 2 に登録する。

30

【 0 0 5 1 】

次いで、圧縮符号付与部 1 1 2 は、動的辞書 1 2 2 に登録した圧縮符号を、単語コードに対応付けてビットフィルタ 1 2 1 a に登録する。さらに、圧縮符号付与部 1 1 2 は、登録した圧縮符号に登録番号を付与し、登録番号を圧縮符号に対応付けて動的辞書 1 2 2 に登録する。なお、登録番号とは、動的辞書 1 2 2 に登録された順番を表す番号である。

40

【 0 0 5 2 】

図 1 1 を用いて、登録番号および圧縮符号が登録された後のビットフィルタ 1 2 1 a の具体例について説明する。図 1 1 は、ビットフィルタの数値部の第 2 の例を示す図である。図 1 1 の例のように、ビットフィルタ 1 2 1 a に登録番号および圧縮符号が登録される。ビットフィルタ 1 2 1 a において、「登録番号」は、単語コードが動的辞書 1 2 2 に登録された順番を表す番号である。「圧縮符号」は、数値文字列に対応する圧縮符号である。例えば、ビットフィルタ 1 2 1 a において数値文字列「115」に対応する登録番号「15」は、数値文字列「115」に係る圧縮符号「A017」が 1 5 番目に動的辞書 1 2 2 に登録されたことを表す。ビットフィルタ 1 2 1 a には、数値文字列「115」に付与された圧縮符号「A017h」が登録される。また、ビットフィルタ 1 2 1 a において数値文字列「121」に

50

対応する登録番号「12」は、数値文字列「121」に係る圧縮符号「A00E」が12番目に動的辞書122に登録されたことを表す。ビットフィルタ121aには、数値文字列「121」に付与された圧縮符号「A00Eh」が登録される。

【0053】

そして、圧縮符号付与部112は、ビットフィルタ121aに登録した圧縮符号をファイルライト部113に出力する。

【0054】

ファイルライト部113は、圧縮符号付与部122から出力された圧縮符号を基に圧縮ファイルを生成する処理部である。ファイルライト部113は、例えば圧縮符号付与部122から出力された数値文字列または単語の各圧縮符号を、それぞれバッファに格納して圧縮データを生成する。ファイルライト部113は、バッファに生成された圧縮データを基に圧縮ファイルを生成する。

【0055】

(実施例1の圧縮処理の流れ)

次に、実施例1の圧縮処理の流れについて説明する。図12は、実施例1の圧縮処理の流れを説明するための図である。図12の例のように、情報処理装置100は、前処理をおこなう(ステップS10)。例えば、情報処理装置100は、前処理においてビットフィルタ121を保持する領域や、動的辞書122を作成する作業領域を確保する。ファイルリード部111は、対象ファイルを読み出し(ステップS11)、対象ファイルから数値文字列を抽出する(ステップS12)。

【0056】

圧縮符号付与部112は、対象ファイルから抽出された数値文字列に対応する圧縮符号がビットフィルタ121aに登録されているか否かを判定する(ステップS13)。圧縮符号付与部112は、ビットフィルタ121aに圧縮符号が登録されている場合(ステップS13Yes)、ステップS18の処理に移行する。

【0057】

一方、圧縮符号付与部112は、ビットフィルタ121aに圧縮符号が登録されていない場合(ステップS13No)、ビットフィルタ121aから単語コードを取得する(ステップS14)。次いで、圧縮符号付与部112は、ビットフィルタ121aから取得した単語コードと、動的辞書122に単語コードを登録する順に付与される圧縮符号とを対応付けて動的辞書122に登録する(ステップS15)。例えば、圧縮符号付与部112は、単語コードが動的辞書122に登録される順番に、単語コードに圧縮符号「A000h」「A001h」「A002h」「A003h」「A004h」「A005h」・・・を付与する。圧縮符号付与部112は、動的辞書122から登録された圧縮符号を取得する(ステップS16)。次いで、圧縮符号付与部112は、動的辞書122から取得された圧縮符号を、単語コードに対応づけてビットフィルタ121aに登録する(ステップS17)。

【0058】

圧縮符号付与部112は、ビットフィルタ121aから数値文字列に対応する圧縮符号を取得する(ステップS18)。ファイルライト部113は、ビットフィルタ121aから取得された圧縮符号を圧縮ファイルに書き込む(ステップS19)。

【0059】

ファイルリード部111は、ファイルの読み出し位置がファイルの終端であるか否かを判定する(ステップS20)。ファイルリード部111は、読み出し位置がファイルの終端である場合(ステップS20Yes)、処理を終了させる。一方、ファイルリード部111は、読み出し位置がファイルの途中である場合(ステップS20No)、ステップS11の処理に戻る。

【0060】

以上のように実施例1の情報処理装置100は、対象ファイルに含まれる各数値文字列に圧縮符号を割り当てるので圧縮処理時に数値文字列に割り当てる符号長を短くできる。

【0061】

(実施例 1 の効果)

圧縮部 1 1 0 は、入力されたデータから、数値を示す、少なくとも 1 つの数字を含む数値文字列を抽出する。圧縮部 1 1 0 は、数値文字列と圧縮符号とを対応付けて記憶する辞書を参照し、抽出された数値文字列に対応する圧縮符号に変換する。これにより、圧縮処理時に数値文字列に割り当てる符号長を短くできる。

【 0 0 6 2 】

辞書は、所定の複数の数値文字列を予め記憶する第 1 の辞書と、第 1 の辞書に登録されている数値文字列を記憶する第 2 の辞書とである。圧縮部 1 1 0 は、抽出された数値文字列が、第 1 の辞書に登録されており、かつ、該数値文字列に対応する圧縮符号が第 1 の辞書に登録されていない場合、該数値文字列に圧縮符号を付与して第 1 の辞書における該数値文字列に対応付けて該圧縮符号を登録する。これにより、入力された数値にだけ動的に圧縮符号を付与しつつ圧縮することができ、抽出され得るすべての数値に予め圧縮符号を割り当てる場合に比較して、圧縮符号の長さを短くすることができる。

10

【 0 0 6 3 】

第 1 の辞書は、抽出され得る数値文字列とともに、数値の大きさの順番に応じて割り当てられるコードがさらに対応付けて記憶された辞書である。これにより、コードの状態で数値の大きさを比較できる。

【実施例 2】

【 0 0 6 4 】

図 1 3 を用いて、実施例 2 の圧縮処理に係るシステム構成について説明する。図 1 3 は、実施例 2 の圧縮処理に係るシステム構成の例を示す図である。図 1 3 の例に示されるように、情報処理装置 2 0 0 は、圧縮部 2 1 0 と、記憶部 2 2 0 とを有する。圧縮部 2 1 0 は、ファイルリード部 2 1 1、圧縮符号付与部 2 1 2 およびファイルライト部 2 1 3 を有する。記憶部 2 2 0 は、ビットフィルタ 2 2 1、動的辞書 2 2 2 およびビットマップ型全文インデックス 2 2 3 を有する。ビットフィルタ 2 2 1 は、数値部 2 2 1 a および単語部 2 2 1 b を有する。なお、実施例 1 と同じ構成に関しては、番号の下 2 桁を同一にして適宜説明を省略する。

20

【 0 0 6 5 】

圧縮部 2 1 0 は、抽出する処理で数値文字列が抽出された際に、複数の圧縮ファイルのうち数値文字列が含まれる圧縮ファイルを示す情報ビット列を数値文字列ごとに対応付けて保持するインデックスを更新する。以下、実施例 2 の圧縮部 2 1 0 および記憶部 2 2 0 の構成について詳細に説明する。

30

【 0 0 6 6 】

実施例 2 の情報処理装置 2 0 0 は、記憶部 2 2 0 がビットマップ型全文インデックス 2 2 3 を有する点で、実施例 1 の情報処理装置 1 0 0 と異なる。図 1 4 を用いてビットマップ型全文インデックス 2 2 3 のデータ構造について説明する。図 1 4 は、ビットマップ型全文インデックスの一例を示す図である。図 1 4 の例のように、ビットマップ型全文インデックス 2 2 3 は、静的単語および動的単語に係る圧縮付号ごとにビットマップを対応付ける。ビットマップとは、静的単語および動的単語がいずれの圧縮ファイルに含まれるかを表す符号ビット列である。ビットマップの各ビットが、各圧縮ファイルに静的単語または動的単語が含まれているか否かを表す。

40

【 0 0 6 7 】

ビットマップ型全文インデックス 2 2 3 は、例えば、8 0 0 0 種類の静的単語ごと、および 2 4 0 0 0 種類の動的単語ごとにビットマップを対応付ける。静的単語とは、出現頻度集計用のテキストファイル群において各単語の出現頻度を集計した場合に、出現頻度の高い単語を表す。例えば、静的単語は、頻度集計用のテキストファイル群での出現頻度が上位 8 0 0 0 位までの単語である。また、動的単語とは、頻度集計用のテキストファイル群での出現頻度の順位が 8 0 0 0 位未満であって、対象ファイルから抽出された数値文字列または単語である。

【 0 0 6 8 】

50

例えば、ビットマップ型全文インデックス 2 2 3 の有効行 1 行目は、圧縮符号「0001h」に対応する「a」のビットマップが「1011110110...」となっている。ビットマップ型全文インデックスの有効行 1 行目のビットマップは、「a」の圧縮符号が含まれるファイルを表す。ビットマップ「1011110110...」は、1 ビット目に「1」が格納されているのでファイル 1 に「a」が含まれ、2 ビット目に「0」が格納されているのでファイル 2 に「a」が含まれず、3 ビット目に「1」が格納されているのでファイル 3 に「a」が含まれることを表す。また、ビットマップ「1011110110...」は、4 ビット目に「1」が格納されているのでファイル 4 に「a」が含まれ、5 ビット目に「1」が格納されているのでファイル 5 に「a」が含まれていることを表す。なお、ビットマップ「1011110110...」は、ファイル 6 以降の他の各ファイルに「a」が含まれるか否かについても表す。

10

【 0 0 6 9 】

次に、ビットマップ型全文インデックス 2 2 3 の更新について説明する。ファイルライト部 2 1 3 は、圧縮符号付与部 2 1 2 から受け付けた圧縮符号がビットマップ型全文インデックス 2 2 3 に登録されているか否かを判定する。ファイルライト部 2 1 3 は、ビットマップ型全文インデックス 2 2 3 に、受け付けた圧縮符号に対応するビットマップが登録されている場合、受け付けた圧縮符号に対応するビットマップを参照する。ファイルライト部 2 1 3 は、参照したビットマップのうち、対象ファイルに対応するビットが「0」の場合、ビットを「1」に更新する。なお、ファイルライト部 2 1 3 は、対象ファイルに対応するビットが「1」の場合、ビットマップを更新しない。

【 0 0 7 0 】

20

一方、ファイルライト部 2 1 3 は、受け付けた圧縮符号に対応するビットマップがビットマップ型全文インデックス 2 2 3 に、登録されていない場合、ビットマップ型全文インデックス 2 2 3 に新しくビットマップを登録する。

【 0 0 7 1 】

具体的には、ファイルライト部 2 1 3 は、圧縮符号付与部 2 1 2 が対象ファイル中の単語または数値文字列を動的辞書 2 2 2 に登録した場合に、単語または数値文字列に付与された圧縮符号を取得する。かかる場合において、ファイルライト部 2 1 3 は、取得した圧縮符号に係るビットマップをビットマップ型全文インデックス 2 2 3 に登録する。このビットマップには、対象ファイルの数分のビット「0」が含まれる。さらに、ファイルライト部 2 1 3 は、登録したビットマップのビットのうち、対象ファイルに対応するビットを「1」に更新する。すなわち、ファイルライト部 2 1 3 は、対象ファイルにおいて初出の単語または数値文字列を動的辞書 2 2 2 に登録した際に、ビットマップ型全文インデックス 2 2 3 に登録された単語または数値文字列に対応するビットマップを登録する。このようにしてビットマップ型全文インデックス 2 2 3 を生成する。

30

【 0 0 7 2 】

(実施例 2 の圧縮処理の流れ)

次に、実施例 2 の圧縮処理の流れについて説明する。図 1 5 は、実施例 2 の圧縮処理の流れを説明するための図である。図 1 5 の例のように、情報処理装置 2 0 0 は、前処理をおこなう (ステップ S 3 0) 。例えば、情報処理装置 2 0 0 は、前処理としてビットマップ型全文インデックス 2 2 3 を生成するための作業領域を確保する。ファイルリード部 2 1 1 は、対象ファイルを読み出し (ステップ S 3 1) 、対象ファイルから数値文字列を抽出する (ステップ S 3 2) 。

40

【 0 0 7 3 】

圧縮符号付与部 2 1 2 は、ビットフィルタ 2 2 1 a に、対象ファイルから抽出した数値文字列に対応する圧縮符号が登録されている場合 (ステップ S 3 3 Yes) 、ステップ S 3 8 の処理に移行する。

【 0 0 7 4 】

一方、圧縮符号付与部 2 1 2 は、ビットフィルタ 2 2 1 a に、対象ファイルから抽出した数値文字列に対応する圧縮符号が登録されていない場合 (ステップ S 3 3 No) 、ビットフィルタ 2 2 1 a から数値文字列の単語コードを取得する (ステップ S 3 4) 。圧縮符号

50

付与部 2 1 2 は、取得した単語コードと、動的辞書 2 2 2 へ登録する順番に付与される圧縮符号とを対応付けて動的辞書 2 2 2 に登録する（ステップ S 3 5）。圧縮符号付与部 2 1 2 は、動的辞書 2 2 2 に登録された圧縮符号を取得する（ステップ S 3 6）。圧縮符号付与部 2 1 2 は、動的辞書 2 2 2 から取得した圧縮符号をビットフィルタ 2 2 1 a に登録する（ステップ S 3 7）。圧縮符号付与部 2 1 2 は、ビットフィルタ 2 2 1 a から圧縮符号を取得する（ステップ S 3 8）。

【 0 0 7 5 】

ファイルライト部 2 1 3 は、圧縮符号付与部 2 1 2 が取得した圧縮符号を基にしてビットマップ型全文インデックス 2 2 3 を更新する（ステップ S 3 9）。例えば、ファイルライト部 2 1 3 は、ビットマップ型全文インデックス 2 2 3 に圧縮符号に対応するビットマップが登録されている場合、ビットマップに含まれるビットのうち、対象ファイルに対応するビットを「1」に更新する。一方、ファイルライト部 2 1 3 は、ビットマップ型全文インデックス 2 2 3 に圧縮符号に対応するビットマップが登録されていない場合、ビットマップ型全文インデックス 2 2 3 に新しく圧縮符号に対応するビットマップを登録する。

【 0 0 7 6 】

ファイルライト部 2 1 3 は、圧縮符号付与部 2 1 2 が取得した圧縮符号を圧縮ファイルに書き込む（ステップ S 4 0）。

【 0 0 7 7 】

ファイルリード部 2 1 1 は、ファイルの読み出し位置がファイルの終端であるか否かを判定する（ステップ S 4 1）。ファイルリード部 2 1 1 は、読み出し位置がファイルの終端である場合（ステップ S 4 1 Yes）、処理を終了させる。一方、ファイルリード部 2 1 1 は、読み出し位置がファイルの途中である場合（ステップ S 4 1 No）、ステップ S 3 1 の処理に戻る。

【 0 0 7 8 】

以上のように実施例 2 の情報処理装置 2 0 0 は、ファイル圧縮を行う際にビットマップ型全文インデックス 2 2 3 を生成する。これにより、情報処理装置 2 0 0 は、複数の圧縮ファイルを基に数値文字列検索をおこなう際に、検索対象の数値文字列を有するファイルを特定でき、オープンする圧縮ファイルを絞り込むことができるので数値文字列検索を高速化することができる。

【 0 0 7 9 】

（実施例 2 の効果）

圧縮部 2 1 0 は、抽出する処理で数値文字列が抽出された際に、複数の圧縮ファイルのうち数値文字列が含まれる圧縮ファイルを示す情報ビット列を数値文字列ごとに対応付けて保持するインデックスを更新する。これにより、複数の圧縮ファイルを基に数値文字列検索をおこなう際に、検索対象の数値文字列を有するファイルを特定でき、オープンする圧縮ファイルを絞り込むことができるので数値文字列検索を高速化することができる。

【実施例 3】**【 0 0 8 0 】**

数値を数字単位で圧縮された状態で大小比較が可能となるように圧縮すると、出現頻度の高い 0、1 等と、出現頻度の低い 8、9 等を同等の長さの符号長で圧縮する必要があり、出現頻度が低い数字に短い符号を割当てることになるから、圧縮ファイル全体として圧縮率が低下する。

【 0 0 8 1 】

例えば、0 ~ 9 までの各数字に対して 4 ビットの符号長を割当てる場合、出現頻度が「0.0625」の文字又は単語と同等の符号長が割当てられることになり、出現頻度の低い 8、9 等の数字が、出現頻度の高い文字または単語と同等に取り扱われる。このため、他の文字または単語に係る圧縮符号の領域が狭められ、他の文字又は単語に割当てられる符号長が長く補正されることになるので、圧縮ファイル全体として圧縮率が低下するという問題がある。

【 0 0 8 2 】

図16を用いて、実施例3の大小比較処理に係るシステム構成について説明する。図16は、実施例3の大小比較処理に係るシステム構成の一例を示す図である。図16の例に示されるように、情報処理装置300は、処理部330と、記憶部320とを有する。処理部330は、検索範囲受付部331、ファイル選択部332、伸長辞書生成部333、ファイルリード部334、比較部335およびファイルライト部336を有する。記憶部320は、ビットマップ型全文インデックス323および伸長辞書324を有する。

【0083】

記憶部320は例えば、RAM、ROM、フラッシュメモリなどの半導体メモリ素子、ハードディスクや光ディスクなどの記憶装置に対応する。また、処理部330の機能は例えば、CPUが所定のプログラムを実行することで実現することができる。また、処理部330の機能は例えば、ASICやFPGAなどの集積回路により実現することができる。ビットマップ型全文インデックス323のデータ構造に関しては、実施例1または実施例2と同じであるので説明を省略する。

10

【0084】

処理部330は、数値を示す、少なくとも1つの数字を含む数値文字列とともに数値の大きさの順番に応じて割り当てられるコードと圧縮符号とを対応付けて記憶する辞書を用いて、圧縮ファイルから抽出した圧縮符号をコードに変換する。処理部330は、変換したコードと、数値文字列の検索範囲に対応する1以上のコードとを比較することで、変換したコードが検索範囲内に含まれるか否かを判定する。処理部330は、検索範囲内に含まれると判定した場合に、変換したコードを数値文字列に伸長して表示する。

20

【0085】

また、処理部330は、各数値文字列を含む圧縮ファイルに関する情報を有するインデックスを用いて、検索範囲内の数値文字列を含む圧縮ファイルを特定する。以下、実施例3の処理部330および記憶部320の構成について詳細に説明する。

【0086】

検索範囲受付部331は、ユーザによって入力された数値文字列検索の範囲を受け付ける処理部である。検索範囲受付部331は、例えば、入力用フォーマットに入力された検索範囲の最大値および最小値を取得することで数値文字列検索の範囲を受け付ける。検索範囲とは、指定された最大値および最小値に属する数値文字列の範囲である。検索範囲受付部331は、受け付けた最大値および最小値を検索範囲としてファイル選択部332および比較部335に出力する。

30

【0087】

ファイル選択部332は、ビットマップ型全文インデックス323を用いて、検索範囲内の数値文字列を有する圧縮ファイルを選択する処理部である。ファイル選択部332は、ビットマップ型全文インデックス323から、検索範囲内の圧縮符号に対応する1以上のビットマップを抽出する。次いで、ファイル選択部332は、抽出した1以上のビットマップ同士でor演算することで選択結果マップを生成する。選択結果マップとは、検索範囲内の数値文字列を1つ以上含む圧縮ファイルを示すビットマップである。ファイル選択部332は、選択結果マップをファイルリード部334に出力する。

40

【0088】

次に、図17を用いて選択結果マップを生成する処理の流れを説明する。図17は、圧縮ファイル選択の処理の流れの例を示す第1の図である。図17の例には、ビットマップ型全文インデックス323の動的単語に対応する部分が示される。図17の例のように、ファイル選択部332は、例えば検索範囲が110～125であった場合、検索範囲内に属する「115」のビットマップと、「121」のビットマップとを抽出する。次いで、ファイル選択部332は、「115」のビットマップ「1011011000...」と、「121」のビットマップ「1001010101...」とでor演算を行い、選択結果マップ「1011011101...」を生成する。ファイル選択部332によって生成された選択結果マップ「1011011101...」は、数値文字列「115」および「121」のいずれか一方または両方を含むそれぞれの圧縮ファイルを示す。

【0089】

50

なお、ビットマップ型全文インデックス 3 2 3 に、検索範囲内に属するビットマップが 3 以上のある場合、ファイル選択部 3 3 2 は、各ビットマップに対して or 演算を実行して選択結果マップを生成する。例えば、ファイル選択部 3 3 2 は、検索範囲内にビットマップ A、ビットマップ B、ビットマップ C が属する場合、条件式 (A or B) or C を算出することで、選択結果マップを生成する。

【 0 0 9 0 】

なお、ファイル選択部 3 3 2 は、ビットマップ型全文インデックス 3 2 3 のうち検索範囲内に該当するビットマップを特定する際に、後述する大小比較の方法を用いてもよい。

【 0 0 9 1 】

次に、図 1 8 を用いて圧縮ファイルを選択する処理の流れを説明する。図 1 8 は、圧縮ファイル選択の処理の流れの例を示す第 2 の図である。図 1 8 の例のように、帳票フォルダ 5 1 には、複数の圧縮ファイル 5 2 a、5 2 b、5 2 c 等が含まれる。ファイル選択部 3 3 2 は、選択結果マップ 5 0 が「1011011101...」となっているので、帳票フォルダ 5 1 に含まれるファイルのうち、1 番目の圧縮ファイル 5 2 a と 3 番目の圧縮ファイル 5 2 c とをオープンする。なお、ファイル選択部 3 3 2 は、他にも 4 番目、6 ~ 8 番目、1 0 番目のファイルをオープンする。ファイル選択部 3 3 2 によってオープンされた圧縮ファイル 5 2 a には、例えば検索範囲 1 1 0 ~ 1 2 5 に属する数値文字列に係る圧縮符号として「115(円)」に係る圧縮符号が格納されている。

【 0 0 9 2 】

伸長辞書生成部 3 3 3 は、伸長辞書 3 2 4 を生成する処理部である。伸長辞書生成部 3 3 3 は、圧縮ファイルに含まれる圧縮データに基づいて伸長辞書 3 2 4 を生成する。図 1 9 は、伸長辞書の構造を説明するための図である。図 1 9 の例のように、伸長辞書 3 2 4 は根 3 2 4 a と、枝 3 2 4 b (1) ~ 3 2 4 b (4)、葉 3 2 4 c (1) ~ 3 2 4 c (4) を有する。

【 0 0 9 3 】

枝 3 2 4 b (1) ~ 3 2 4 b (4) は、それぞれ 3 2 4 c (1) ~ 3 2 4 c (4) に格納された数値文字列に対応する圧縮符号が含まれる。情報処理装置 3 0 0 は、例えば、圧縮ファイルから圧縮符号を読み込んだ際に、読み込んだ圧縮符号と、枝 3 2 4 b (1) ~ 3 2 4 b (4) に含まれる圧縮符号とを比較することで、読み込んだ圧縮符号に対応する葉 3 2 4 c を特定する。

【 0 0 9 4 】

葉 3 2 4 c (1) ~ 3 2 4 c (4) に対応する葉の構造体は、例えば葉の構造体 3 2 4 C に表される。葉の構造体 3 2 4 C は、葉識別情報、圧縮符号長、文字コードまたは数値テーブルへのポインタ等を有する。「葉識別情報」は、葉を一意に識別する情報である。「圧縮符号長」は、後述するファイルリード部 3 3 4 によって取得された圧縮データのビット列のうち、有効な長さを示す情報である。例えば、各数値文字列には、1 6 ビットの固定長符号が割り当てられるので、数値文字列に対応する圧縮符号長には、「16」ビットが格納される。「文字コード」は、例えばアスキーコード等の文字コードを示す。「数値テーブルへのポインタ」は、複数の数値文字列の単語コードを格納する数値テーブルにおいて、取得された圧縮符号に対応する数値文字列の単語コードが格納されている位置を示すポインタである。なお、葉の構造体 3 2 4 C には、文字コードまたは数値テーブルへのポインタのいずれか一方を有する。

【 0 0 9 5 】

数値テーブル 3 2 4 d について説明する。数値テーブル 3 2 4 d は、ビットフィルタ 1 2 1 においてそれぞれの数値文字列に割り当てられた全ての単語コードを有する。例えば、数値テーブル 3 2 4 d は、「110」の単語コード 3 2 4 d (1) 「B006E0h」、「115」の単語コード 3 2 4 d (2) 「B00730h」、「125」の単語コード 3 2 4 d (3) 「B007D0h」を有する。葉 3 2 4 c (1) ~ 3 2 4 c (4) はそれぞれ、数値テーブルへのポインタによって数値テーブル 3 2 4 d に含まれる数値文字列の単語コードに対応付けられる。例えば、数値文字列「110」に対応する葉の構造体 3 2 4 C に係る数値テーブルのポイン

10

20

30

40

50

タには、領域 3 2 4 d (1) の先頭アドレスに対応するオフセットが格納されている。また、数値文字列「115」に対応する葉の構造体 3 2 4 C に係る数値テーブルのポインタには、領域 3 2 4 d (2) の先頭アドレスに対応するオフセットが格納されている。

【 0 0 9 6 】

ファイルリード部 3 3 4 は、圧縮ファイルを読み込む処理部である。ファイルリード部 3 3 4 は、圧縮ファイルから圧縮符号を取得する。ファイルリード部 3 3 4 は、圧縮ファイルから取得した圧縮符号を伸長辞書 3 2 4 と照らし合わせる。すなわち、ファイルリード部 3 3 4 は、圧縮ファイルから取得した圧縮符号と、葉 3 2 4 c (1) ~ 3 2 4 c (4) とを比較し、圧縮ファイルから取得した圧縮符号に対応する葉の構造体 3 2 4 C を特定する。ファイルリード部 3 3 4 は、圧縮符号に対応する葉の構造体 3 2 4 C にアクセスする。ファイルリード部 3 3 4 は、アクセスした葉の構造体 3 2 4 C に数値テーブルへのポインタが含まれる場合、数値テーブルへのポインタに基づいて数値テーブル 3 2 4 d から単語コードを取得する。

10

【 0 0 9 7 】

例えば、ファイルリード部 3 3 4 は、取得した圧縮符号「101000000010111」が枝 3 2 4 b (1) にヒットする場合、葉 3 2 4 c (1) にアクセスする。ファイルリード部 3 3 4 は、葉 3 2 4 c (1) の葉の構造体 3 2 4 C に数値テーブルへのポインタが格納されていた場合、数値テーブルへのポインタを取得する。ファイルリード部 3 3 4 は、数値テーブルへのポインタに基づいて、数値テーブル 3 2 4 d 内の領域 3 2 4 d (2) の物理アドレスを特定する。ファイルリード部 3 3 4 は、領域 3 2 4 d (2) から数値文字列「115」に対応する単語コード「B00730h」を取得する。そして、ファイルリード部 3 3 4 は、取得した数値文字列「115」に対応する単語コード「B00730h」を比較部 3 3 5 に出力する。

20

【 0 0 9 8 】

比較部 3 3 5 は、検索範囲に対応する単語コードと、ファイルリード部 3 3 4 によって取得された各単語コードとを比較する処理部である。比較部 3 3 5 は、検索範囲の最大値と最小値とに対応する単語コードを取得する。比較部 3 3 5 は、伸長辞書 3 2 4 から検索範囲の最大値と最小値とに対応する単語コードを取得してもよい。次いで、ファイルリード部 3 3 4 は、圧縮ファイルを読み込み、適宜、圧縮ファイルから数値文字列の単語コードを取得し、取得した数値文字列の単語コードを比較部 3 3 5 に出力する。比較部 3 3 5 は、検索範囲の最大値と最小値とに対応する単語コードと、ファイルリード部 3 3 4 によって取得された単語コードとを比較する。比較部 3 3 5 は、ファイルリード部 3 3 4 によって取得された単語コードが検索範囲内に該当するか否かを判定する。そして、比較部 3 3 5 は、比較対象の数値文字列が検索範囲内であると判定した場合に、比較対象の数値文字列に対応する単語コードをファイルライト部 3 3 6 に出力する。

30

【 0 0 9 9 】

ここで、数値文字列に割り当てられる単語コードについて説明する。図 8 の例を用いて説明したように、0 から 1 6 , 3 8 3 までの整数に対して 3 バイトの単語コードが割り当てられ、1 6 , 3 8 4 から 1 , 0 4 8 , 5 7 5 までの整数に対して 4 バイトの単語コードが割り当てられる。さらに、1 , 0 4 8 , 5 7 6 から 1 , 0 7 3 , 7 4 1 , 8 2 3 までの整数に対して 5 バイトの単語コードが割り当てられる。図 8 の例に示す c 1 - 3、c 2 - 3 および c 3 - 3 は、半角 / 全角の別、カンマの有無、小数点の有無、正の値 / 負の値の別などを示す情報ビットである。すなわち、単語コードの末尾 4 ビットが情報ビットに該当する。

40

【 0 1 0 0 】

例えば、整数「0」「1」「2」「3」「4」「5」... に対して、それぞれ 3 バイトの単語コード「B00000h」「B00010h」「B00020h」「B00030h」「B00040h」... が割り当てられる。すなわち、若い整数から順番に単語コードを割り当てる場合、各整数に対して「B00000h」から昇順に単語コードが割り当てられる。これにより、単語コードの状態の数値文字列同士を大小比較することができる。例えば、比較部 3 3 5 は、整数「1」に割り当てられた単語

50

語コード「B00010h」と、整数「3」に割当てられた単語コード「B00030h」とを比較することで、整数「3」の方が大きいと判定することができる。

【 0 1 0 1 】

半角/全角の別、カンマの有無等の表記形式の異なる数値文字列同士の比較について説明する。例えば、数値文字列「1200」には、単語コード「B04B00h」が割当てられ、数値文字列「1,200」には、単語コード「B04B01h」が割当てられ、それぞれ異なる単語コードが割当てられるが、これらは同じ大きさの数値文字列である。比較部 3 3 5 は、表記形式の異なる数値文字列に係る単語コード同士を比較するためにマスク処理を実行する。

【 0 1 0 2 】

マスク処理の具体例について説明する。比較部 3 3 5 は、数値文字列「1200」の単語コード「B04B00h」と数値文字列「1,200」の単語コード「B04B01h」とを比較する場合、それぞれの単語コードにマスク処理用のビット列を乗算する。すなわち、比較部 3 3 5 は、数値文字列「1200」の単語コード(2進数)「1011 0000 0100 1011 0000 0000」に対して、マスク処理用のビット列「1111 1111 1111 1111 1111 0000」を乗算し、ビット列「1011 0000 0100 1011 0000 0000」を取得する。また、比較部 3 3 5 は、数値文字列「1,200」の単語コード(2進数)「1011 0000 0100 1011 0000 0001」に対して、マスク処理用のビット列「1111 1111 1111 1111 1111 0000」を乗算し、ビット列「1011 0000 0100 1011 0000 0000」を取得する。そして、比較部 3 3 5 は、取得したビット列同士を比較し、数値文字列「1200」の単語コード「B04B00h」と数値文字列「1,200」の単語コード「B04B01h」とが等しいと判定する。

10

20

【 0 1 0 3 】

このように、比較部 3 3 5 は、比較対象の各数値文字列に対し、末尾4ビットが「0」となっているマスク処理用のビット列を乗算することで、表記形式の異なる数値文字列に係る単語コード同士を比較することができる。

【 0 1 0 4 】

図 2 0 を用いて、比較部 3 3 5 でなされる大小比較処理の具体例について説明する。図 2 0 は、数値文字列の大小比較を説明するための図である。表 D 1 は、検索範囲の最大値および最小値に対応する単語コードを表す。一方、表 D 2 は、比較対象の数値文字列に対応する単語コードを表す。なお、比較対象の数値文字列とは、ファイルリード部 3 3 4 によって圧縮ファイルから取得された圧縮符号に係る各単語コードである。

30

【 0 1 0 5 】

比較部 3 3 5 は、検索範囲の最小値「110」に対応する単語コード「B006E0h」と、比較対象の数値文字列「115」に対応する単語コード「B00730h」とを比較する。比較部 3 3 5 は、16進数で表記される単語コードを2進数に置き換える。すなわち、比較部 3 3 5 は、検索範囲の最小値「110」に対応する単語コード「B006E0h」を、「1011」「0000」「0000」「0110」「1111」「0000」に置き換える。また、比較部 3 3 5 は、比較対象の数値文字列「115」に対応する単語コード「B00730h」を、「1011」「0000」「0000」「0111」「0011」「0000」に置き換える。そして、比較部 3 3 5 は、2進数に置き換えられた単語コードを4ビット毎に比較し、比較対象の数値文字列「115」が検索範囲の最小値よりも大きいと判定する。

40

【 0 1 0 6 】

比較部 3 3 5 は、検索範囲の最大値「125」に対応する単語コード「B007D0h」と、比較対象の数値文字列「115」に対応する単語コード「B00730h」とを比較する。比較部 3 3 5 は、検索範囲の最大値「125」に対応する単語コード「B007D0h」を、「1011」「0000」「0000」「1000」「0010」「0000」に置き換える。そして、比較部 3 3 5 は、2進数に置き換えられた4ビット毎に比較し、比較対象の数値文字列「115」が検索範囲の最大値よりも小さいと判定する。

【 0 1 0 7 】

比較部 3 3 5 は、比較対象の数値文字列が検索範囲の最小値よりも大きく、検索範囲の最大値よりも小さいので、比較対象の数値文字列が検索範囲内に含まれると判定する。次

50

いで、比較部 335 は、比較対象の数値文字列「115」に対応する単語コード「B00730h」をファイルライト部 336 に出力する。

【0108】

ファイルライト部 336 は、検索範囲内の数値文字列の単語コードを数値文字列に変換し、検索結果として出力する処理部である。ファイルライト部 336 は、比較部 335 から出力された単語コード「B00730h」を数値文字列「115」に変換して、所定の形式でモニタ、プリンタ等の出力媒体に表示する。例えば、ファイルライト部 336 は、数値文字列「115」と共に、数値文字列「115」が含まれるファイル名、数値文字列「115」が含まれるページ数および行数を出力媒体に表示する。

【0109】

(実施例 3 の大小比較処理の全体の流れ)

次に、実施例 3 の大小比較処理の全体の流れについて説明する。図 21 は、実施例 3 の大小比較処理の全体の流れを示す図である。図 21 の例のように、情報処理装置 300 は、前処理をおこなう(ステップ S50)。例えば、情報処理装置 300 は、前処理として大小比較処理を行うための作業領域を確保する。ファイル選択部 332 は、ビットマップ型全文インデックス 323 の情報に基づいてオープンする圧縮ファイルを選択する(ステップ S51)。比較部 335 は、ファイル選択部 332 によって選択された圧縮ファイルから検索範囲内に含まれる単語コードを抽出する(ステップ S52)。ファイルライト部 336 は、検索範囲内に含まれる各単語コードをそれぞれ数値文字列に変換し、変換された各数値文字列を所定の形式で表示媒体に出力することで比較結果を表示する(ステップ S53)。なお、表示媒体は、例えば、モニタ、プリンタ等である。

【0110】

(実施例 3 の圧縮ファイル選択処理の流れ)

次に、実施例 3 の圧縮ファイル選択処理の流れについて説明する。図 22 は、実施例 3 の圧縮ファイル選択処理の流れを示す図である。図 22 に示される処理は、ステップ S51 に対応する。図 22 の例のように、情報処理装置 100 は、前処理をおこなう(ステップ S60)。例えば、情報処理装置 100 は、前処理として選択結果マップを生成するための作業領域を確保する。

【0111】

ファイル選択部 332 は、検索範囲受付部 331 が受け付けた検索範囲の最大値および最小値に対応する単語コードを取得する(ステップ S61)。ファイル選択部 332 は、取得した検索範囲の最大値および最小値に対応する単語コードに基づいて、検索範囲内に属するビットマップをビットマップ型全文インデックス 323 から抽出する(ステップ S62)。ファイル選択部 332 は、ビットマップ型全文インデックス 323 から抽出した複数のビットマップ間で OR 演算をおこなうことで選択結果マップを生成する(ステップ S63)。ファイル選択部 332 は、生成された選択結果マップに基づいてオープンする圧縮ファイルを選択する。

【0112】

(実施例 3 の数値文字列の単語コード抽出処理の流れ)

次に、実施例 3 の数値文字列の単語コード抽出処理の流れについて説明する。図 23 は、実施例 3 の数値文字列の単語コード抽出処理の流れを示す図である。図 23 に示される処理は、ステップ S52 に対応する。図 23 の例のように、情報処理装置 300 は、前処理をおこなう(ステップ S70)。例えば、情報処理装置 300 は、前処理として、伸長辞書 324 を記憶するための領域や大小比較処理を行うための作業領域を確保する。

【0113】

伸長辞書生成部 333 は、伸長辞書を生成する(ステップ S71)。比較部 335 は、検索範囲の最大値および最小値に対応する数値文字列の単語コードを取得する(ステップ S72)。例えば、比較部 335 は、伸長辞書 324 から検索範囲の最大値および最小値に対応する数値文字列の単語コードを取得する。ファイルリード部 334 は、選択結果マップに基づいて選択された圧縮ファイルを読み出す(ステップ S73)。ファイルリード

10

20

30

40

50

部 3 3 4 は、選択された圧縮ファイルから順次 1 6 ビットの数値文字列に係る圧縮符号を取得する（ステップ S 7 4）。比較部 3 3 5 は、圧縮ファイルから取得した圧縮符号に対応する数値文字列の単語コードを伸長辞書 3 2 4 から取得する（ステップ S 7 5）。

【 0 1 1 4 】

比較部 3 3 5 は、数値文字列の単語コードが検索範囲内に該当するか否かを判定する（ステップ S 7 6）。比較部 3 3 5 は、数値文字列の単語コードが検索範囲内に該当する場合（ステップ S 7 6 Yes）、数値文字列の単語コードを伸長文字（数値文字列）に変換し（ステップ S 7 7）、比較結果として伸長文字をモニタ等の表示媒体に出力する（ステップ S 7 8）。一方、比較部 3 3 5 は、数値文字列の単語コードが検索範囲内に該当しない場合（ステップ S 7 6 No）、ステップ S 7 9 の処理に移行する。

10

【 0 1 1 5 】

ファイルリード部 3 3 4 は、ファイルの終端まで至ったか否かを判定する（ステップ S 7 9）。ファイルリード部 3 3 4 は、ファイルの終端まで至った場合（ステップ S 7 9 Yes）、処理を終了させる。一方、ファイルリード部 3 3 4 は、ファイルの途中である場合（ステップ S 7 9 No）、ステップ S 7 4 の処理に戻る。

【 0 1 1 6 】

（実施例 3 の効果）

処理部 3 3 0 は、数値を示す、少なくとも 1 つの数字を含む数値文字列とともに数値の大きさの順番に応じて割り当てられるコードと圧縮符号とを対応付けて記憶する辞書を用いて、圧縮ファイルから抽出した圧縮符号をコードに変換する。処理部 3 3 0 は、変換したコードと、数値文字列の検索範囲に対応する 1 以上のコードとを比較することで、変換したコードが検索範囲内に含まれるか否かを判定する。処理部 3 3 0 は、検索範囲内に含まれると判定した場合に、変換したコードを数値文字列に伸長して表示する。これにより、処理部 3 3 0 は、全ての単語コードを数値文字列に伸長しなくても、検索範囲に含まれる単語コードのみを数値文字列に伸長すればよいので、数値文字列の検索を高速化できる。

20

【 0 1 1 7 】

処理部 3 3 0 は、各数値文字列を含む圧縮ファイルに関する情報を有するインデックスを用いて、検索範囲内の数値文字列を含む圧縮ファイルを特定する。これにより、検索範囲内の数値文字列を有する圧縮ファイルを特定し、オープンする圧縮ファイルを絞り込むので、複数の圧縮ファイルがある場合でも数値文字列検索を高速化できる。

30

【 0 1 1 8 】

（実施例 1 ~ 3 に関連する他の態様）

以下、上述の実施形態における変形例の一部を説明する。下記の変形例のみでなく、本発明の本旨を逸脱しない範囲の設計変更は適宜行われうる。

【 0 1 1 9 】

また、圧縮処理の対象は、ファイル内のデータ以外にも、システムから出力される監視メッセージなどでもよい。例えば、バッファに順次格納される監視メッセージを上述の圧縮処理により圧縮し、ログファイルとして格納するなどの処理が行なわれる。また、例えば、データベース内のページ単位に圧縮が行なわれてもよいし、複数のページをまとめた単位で圧縮が行なわれてもよい。

40

【 0 1 2 0 】

また、実施例 1 に示した処理手順、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。

【 0 1 2 1 】

（情報処理装置のハードウェア構成）

図 2 4 は、実施例 1 ~ 3 の情報処理装置のハードウェア構成を示す図である。図 2 4 の例が示すように、コンピュータ 4 0 0 は、各種演算処理を実行する CPU 4 0 1 と、ユーザからのデータ入力を受け付ける入力装置 4 0 2 と、モニタ 4 0 3 とを有する。また、コンピュータ 4 0 0 は、記憶媒体からプログラム等を読み取る媒体読取装置 4 0 4 と、他の

50

装置と接続するためのインターフェース装置 405 と、他の装置と無線により接続するための無線通信装置 406 とを有する。また、コンピュータ 400 は、各種情報を一時記憶する RAM 407 と、ハードディスク装置 408 とを有する。また、各装置 401 ~ 408 は、バス 409 に接続される。

【0122】

ハードディスク装置 408 には、例えば図 6 に示したファイルリード部 111、圧縮符号付与部 112 およびファイルライト部 113 の各処理部と同様の機能を有する情報処理プログラムが記憶される。さらに、ハードディスク装置 408 には、ファイルリード部 111、圧縮符号付与部 112 およびファイルライト部 113 の各処理部と同様の機能を有する情報処理プログラムが記憶される。また、ハードディスク装置 408 には、情報処理プログラムを実現するための各種データが記憶される。

10

【0123】

CPU 401 は、ハードディスク装置 408 に記憶された各プログラムを読み出して、RAM 407 に展開して実行することで各種の処理を行う。これらのプログラムは、コンピュータ 400 を、例えば図 6 に示したファイルリード部 111、圧縮符号付与部 112 およびファイルライト部 113 として機能させることができる。さらに、これらのプログラムは、コンピュータ 400 を、ファイルリード部 111、圧縮符号付与部 112 およびファイルライト部 113 として機能させることができる。

【0124】

なお、上記の情報処理プログラムは、必ずしもハードディスク装置 408 に記憶されている必要はない。例えば、コンピュータ 400 が読み取り可能な記憶媒体に記憶されたプログラムを、コンピュータ 400 が読み出して実行するようにしてもよい。コンピュータ 400 が読み取り可能な記憶媒体は、例えば、CD-ROM や DVD ディスク、USB (Universal Serial Bus) メモリ等の可搬型記録媒体、フラッシュメモリ等の半導体メモリ、ハードディスクドライブ等が対応する。また、公衆回線、インターネット、LAN (Local Area Network) 等に接続された装置にこのプログラムを記憶させておき、コンピュータ 400 がこれらからプログラムを読み出して実行するようにしてもよい。

20

【0125】

図 25 は、コンピュータで動作するプログラムの構成例を示す図である。コンピュータ 400 において、図 24 に示すハードウェア群 26 (401 ~ 409) の制御を行なう OS (オペレーティング・システム) 27 が動作する。OS 27 に従った手順で CPU 401 が動作して、ハードウェア群 26 の制御・管理が行なわれることにより、アプリケーションプログラム 29 やミドルウェア 28 に従った処理がハードウェア群 26 で実行される。さらに、コンピュータ 400 において、ミドルウェア 28 またはアプリケーションプログラム 29 が、RAM 407 に読み出されて CPU 401 により実行される。

30

【0126】

CPU 401 により圧縮機能が呼び出された場合、ミドルウェア 28 またはアプリケーションプログラム 29 の少なくとも一部に基づく処理を行なうことで、(それらの処理を OS 27 に基づいてハードウェア群 26 を制御して) 圧縮部 110 の機能が実現される。圧縮機能は、それぞれアプリケーションプログラム 29 自体に含まれてもよいし、アプリケーションプログラム 29 に従って呼び出されることで実行されるミドルウェア 28 の一部であってもよい。

40

【0127】

アプリケーションプログラム 29 (またはミドルウェア 28) の圧縮機能により得られる圧縮ファイルは、部分的に伸張することも可能である。圧縮ファイルの途中を伸張する場合には、伸張対象の部分までの圧縮データの伸張処理が抑制されるため、CPU 401 の負荷が抑制される。また、伸張対象の圧縮データを部分的に RAM 407 上に展開するので、ワークエリアも削減される。

【0128】

図 26 は、実施形態のシステムにおける装置の構成例を示す図である。図 26 のシステ

50

△は、コンピュータ400a、コンピュータ400b、基地局30およびネットワーク40を含む。コンピュータ400aは、無線または有線の少なくとも一方により、コンピュータ400bと接続されたネットワーク40に接続している。

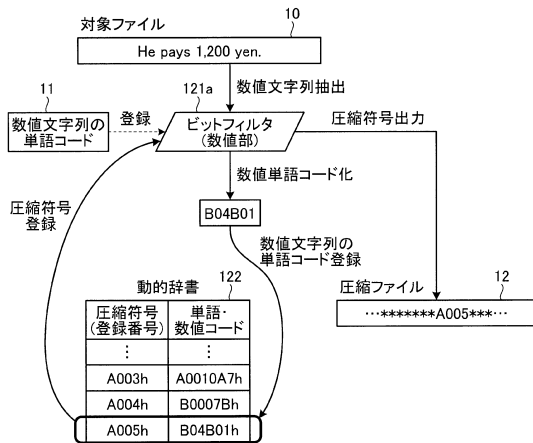
【符号の説明】

【0129】

- 100 情報処理装置
- 110 圧縮部
- 111 ファイルリード部
- 112 圧縮符号付与部
- 120 記憶部
- 121 ビットフィルタ
- 121a 数値部
- 121b 単語部
- 122 動的辞書

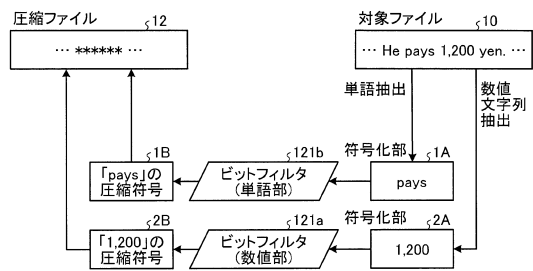
【図1】

実施例1の数値の圧縮処理を説明するための図



【図2】

実施例1の圧縮処理の全体の流れについて説明するための図



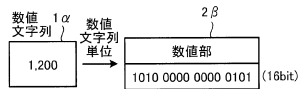
【図3A】

参考例1の符号長を説明するための図



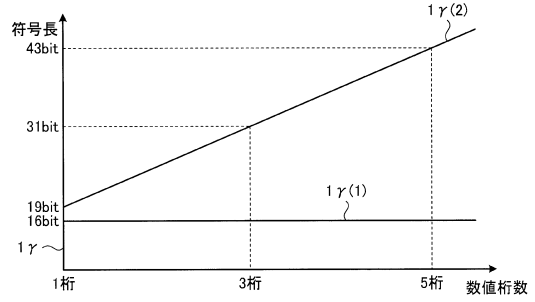
【図3B】

実施例1の符号長を説明するための図



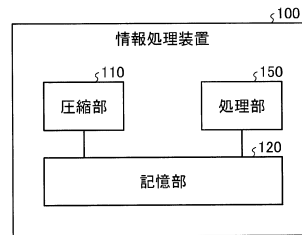
【図4】

参考例1および実施例1の符号長の比較を説明するための図



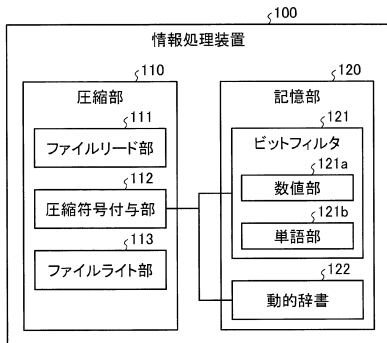
【図5】

情報処理装置のシステム構成の例を説明するための図

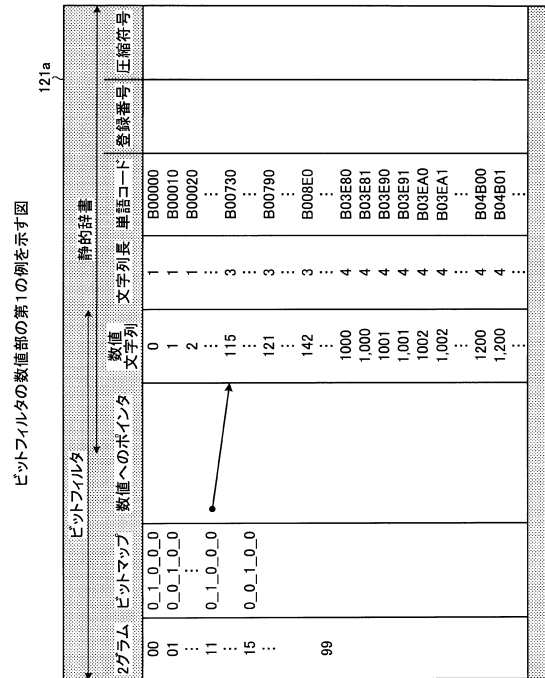


【図6】

実施例1の圧縮処理に係るシステム構成の例を示す図

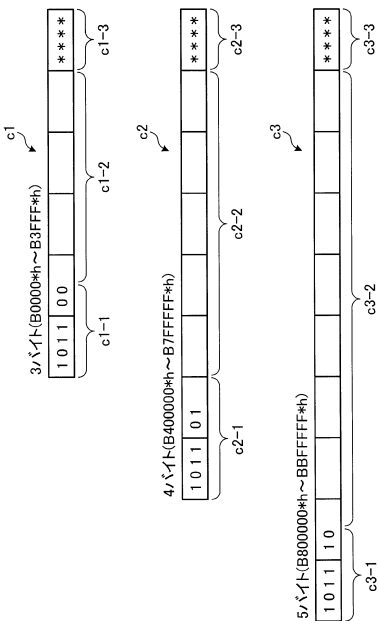


【図7】



【図8】

数値文字列に割り当てられる単語コードのデータ構造の例を示す図



【図9】

ビットフィルタの単語部の例を示す図

動的辞書		静的辞書	
2グラム	ビットマップ	基礎単語へのポイント	登録番号
aa	0_0_0_0_0	aa	1001000...
ab	1_0_0_0_0	able	0101110...
ac	1_0_0_0_0
...
bl	0_1_0_0_0	above	1101110...
bo	0_1_0_0_0	abracadabra	1011101...
ct	0_1_0_0_0
eΔ	0_0_0_1_1	act	1101010...
le	0_0_1_0_0
ou	0_0_1_0_0
ov	0_0_1_0_0
tΔ	0_0_1_0_0
ut	0_0_0_1_0
ve	0_0_0_1_0
...
zΔ	0_0_0_0_0
...

【図10】

動的辞書の一例を示す図

オフセット	圧縮符号	コード
0x0000	A000h	c(order)
0x0004	A001h	c(box)
0x0007	A002h	c(1000)
0x0010	A003h	c(him)
0x0013	A004h	c(pays)
0x0018	A005h	c(1,200)
...
0x002F	A016h	c(121)
...
0x0035	A018h	c(able)
...
0x003F	A025h	c(Kataoka)
...
0x0046	A036h	c(1024)
...

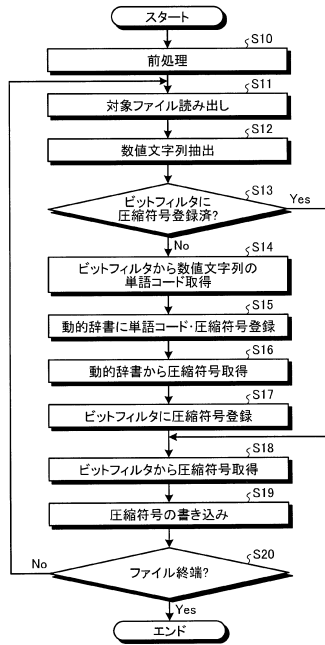
【図11】

ビットフィルタの数値部の第2の例を示す図

動的辞書		静的辞書	
2グラム	ビットマップ	数値文字列	登録番号
00	0_1_0_0_0	0	B00000
01	0_0_1_0_0	1	B00010
...
11	0_1_0_0_0	2	B00020
...
15	0_1_0_0_0	115	B00730
...
99	0_0_1_0_0	121	B00790
...
...	...	142	B008E0
...	...	1000	B03E80
...	...	1000	B03E81
...	...	1001	B03E90
...	...	1001	B03E91
...	...	1002	B03EA0
...	...	1002	B03EA1
...	...	1200	B04B00
...	...	1200	B04B01

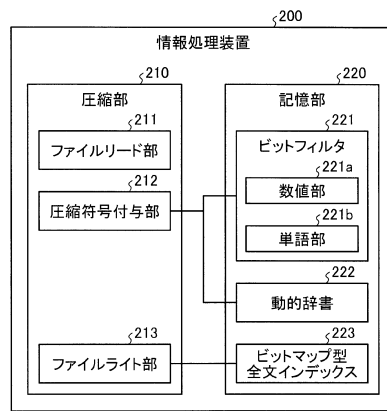
【図12】

実施例1の圧縮処理の流れの例を示す図



【図13】

実施例2の圧縮処理に係るシステム構成の例を示す図



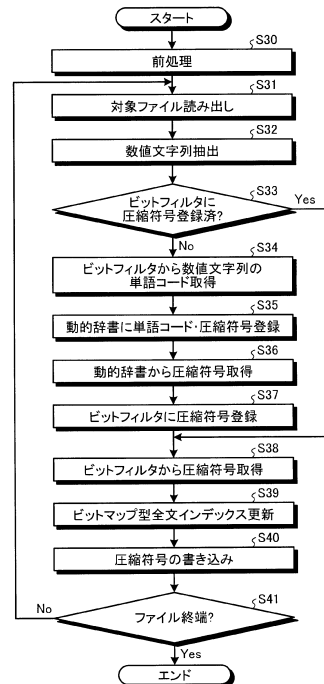
【図14】

ビットマップ型全文インデックスの一例を示す図

圧縮符号	静的単語	ビットマップ	10,000
0000h	「a」	1011110110...	223
...	
0108h	「I」	1011110100...	
021Fh	「zoo」	0001011100...	
動的単語			
A000h	「aardvark」	1011110110...	
...	
A004h	「pays」	1011001000...	
A005h	「1,200」	1000110011...	
...	
A01Eh	「11,800」	0101110001...	
...	
A026h	「115」	1011011000...	
...	
A02Fh	「121」	1001010101...	
...	
FFFFh	「expense」	0001101100...	

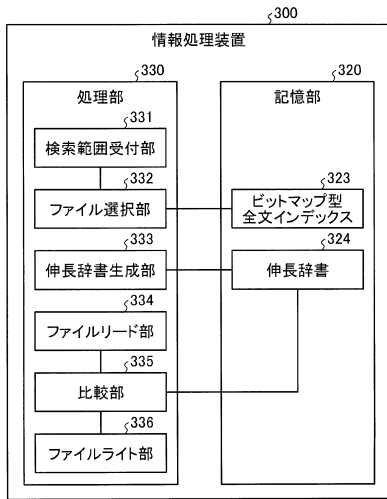
【図15】

実施例2の圧縮処理の流れを示す図



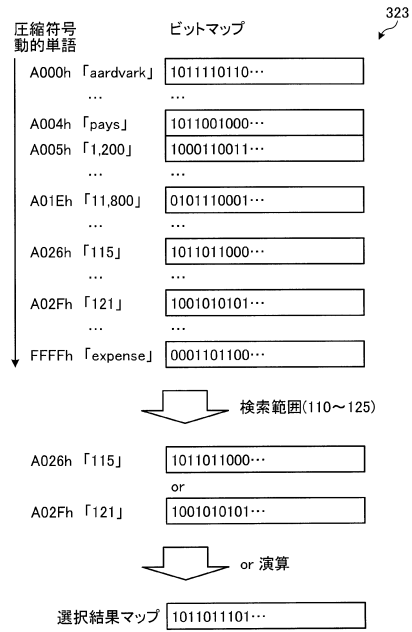
【図16】

実施例3の大小比較処理に係るシステム構成の一例を示す図



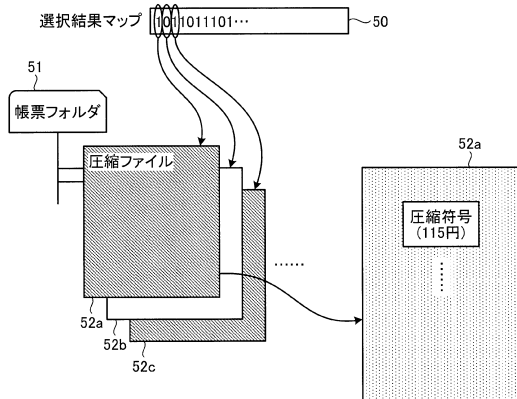
【図17】

圧縮ファイル選択の処理の流れの例を示す第1の図



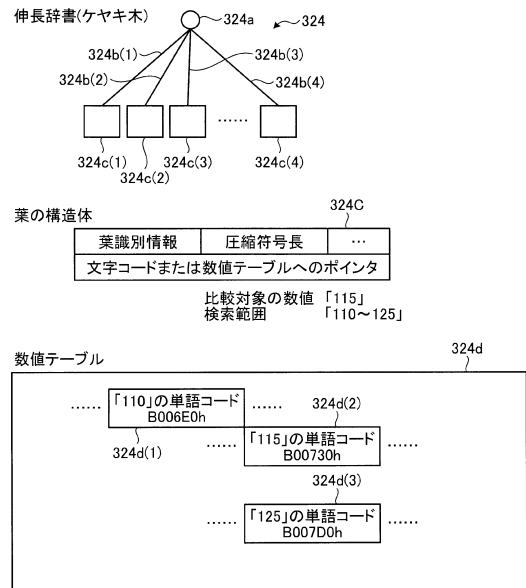
【図18】

圧縮ファイル選択の処理の流れの例を示す第2の図

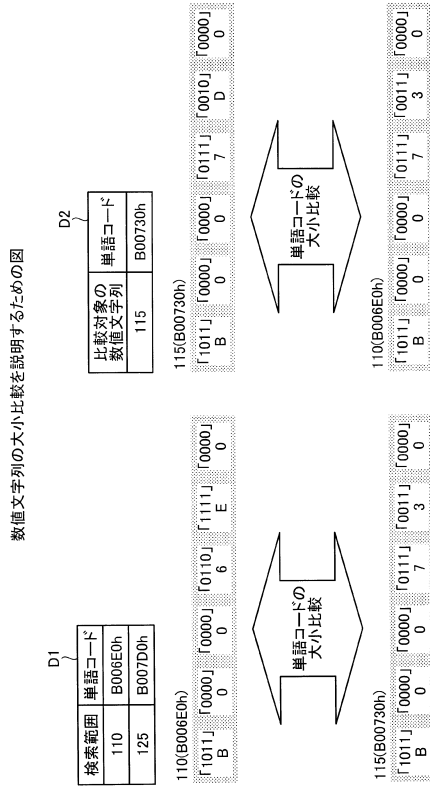


【図19】

伸長辞書の構造を説明するための図

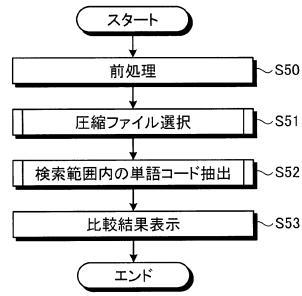


【図20】



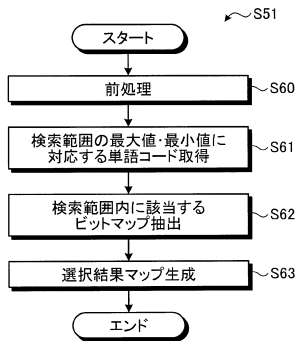
【図21】

実施例3の大小比較処理の全体の流れを示す図



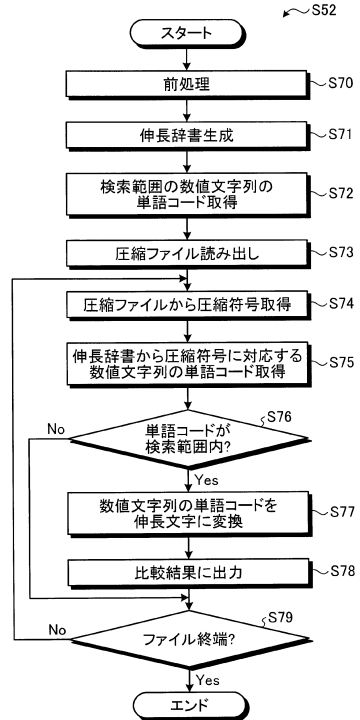
【図22】

実施例3の圧縮ファイル選択処理の流れを示す図



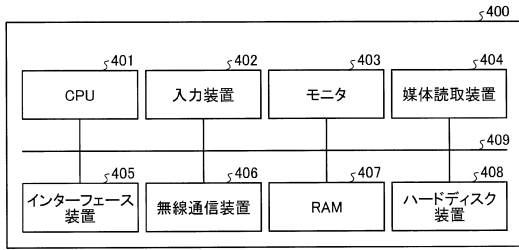
【図23】

実施例3の数値文字列の単語コード抽出処理の流れを示す図



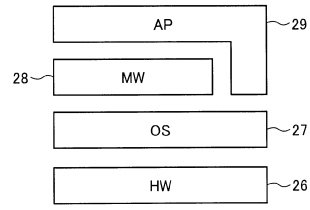
【図24】

実施例1～3の情報処理装置のハードウェア構成を示す図



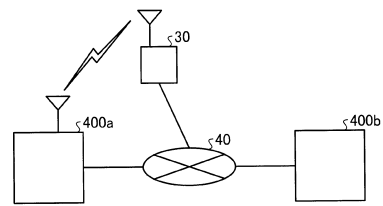
【図25】

コンピュータで動作するプログラムの構成例を示す図



【図26】

実施形態のシステムにおける装置の構成例を示す図



フロントページの続き

審査官 北村 智彦

(56)参考文献 国際公開第2008/047432(WO, A1)

特開平05-120358(JP, A)

特開2011-019011(JP, A)

特開2008-287412(JP, A)

(58)調査した分野(Int.Cl., DB名)

H03M 3/00-9/00

IEEE Xplore