

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2007-508753

(P2007-508753A)

(43) 公表日 平成19年4月5日(2007.4.5)

(51) Int. Cl.	F I	テーマコード (参考)
H03M 7/30 (2006.01)	H03M 7/30 Z	5B082
G06F 12/00 (2006.01)	G06F 12/00 511A	5J064

審査請求 未請求 予備審査請求 未請求 (全 18 頁)

(21) 出願番号 特願2006-534542 (P2006-534542)
 (86) (22) 出願日 平成16年10月15日 (2004.10.15)
 (85) 翻訳文提出日 平成18年4月13日 (2006.4.13)
 (86) 国際出願番号 PCT/AU2004/001406
 (87) 国際公開番号 W02005/039057
 (87) 国際公開日 平成17年4月28日 (2005.4.28)
 (31) 優先権主張番号 2003905688
 (32) 優先日 平成15年10月17日 (2003.10.17)
 (33) 優先権主張国 オーストラリア (AU)

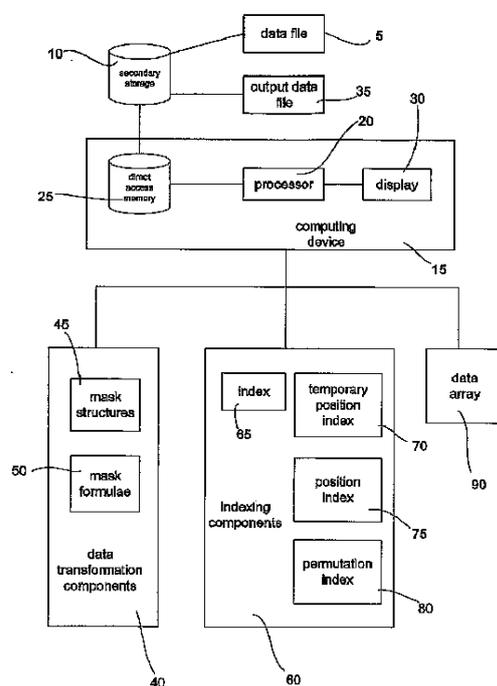
(71) 出願人 506126093
 パクバイト ソフトウェア プロプライエ
 タリー リミティド
 オーストラリア国、ニューサウスウェール
 ズ 2121, エッピング, ラングストン
 プレイス 36エー
 (74) 代理人 100099759
 弁理士 青木 篤
 (74) 代理人 100092624
 弁理士 鶴田 準一
 (74) 代理人 100102819
 弁理士 島田 哲郎
 (74) 代理人 100133835
 弁理士 河野 努

最終頁に続く

(54) 【発明の名称】 データ圧縮システム及び方法

(57) 【要約】

本発明は、所定長以上の長さのバイトシーケンスを有するデータファイルの圧縮方法を提供する。その方法は、二次記憶装置からデータファイルを取り出すステップと、データファイルをダイレクトアクセスメモリに保存するステップと、データファイルの所定長を超えない長さを持つサブシーケンス内で一意なバイト値の頻度を計算するステップと、サブシーケンス内で計算された一意なバイト値の頻度を表すデータ値を含むインデックスを作成するステップと、所定の閾値未満である一意なバイト値の頻度を持つサブシーケンス上で、サブシーケンス内で一意なバイト値の頻度を増加するためにサブシーケンスに対してデータ変換を適用し、データ変換を表すデータ値をインデックスに追加するステップと、所定の閾値を超える一意なバイト値の頻度を持つサブシーケンス上で、サブシーケンス内で1以上の一意な値の位置を表すデータ値をインデックスに追加するステップと、ファイルタイプ識別子を持つ出力データファイルを作成するステップと、出力データファイルにインデックスを追加するステップとを含む。



【特許請求の範囲】**【請求項 1】**

所定長以上の長さのバイトシーケンスを有するデータファイルの圧縮方法であって、
二次記憶装置からデータファイルを取り出すステップと、
前記データファイルをダイレクトアクセスメモリに保存するステップと、
前記データファイルの前記所定長を超えないサブシーケンス内で一意なバイト値の頻度を計算するステップと、
前記サブシーケンス内で計算された前記一意なバイト値の頻度を表すデータ値を含むインデックスを前記サブシーケンスに対して作成するステップと、
所定の閾値未満である一意なバイト値の頻度を持つ前記サブシーケンス上で、該サブシーケンス内で該一意なバイト値の頻度を増加するために該サブシーケンスに対してデータ変換を適用し、該データ変換を表すデータ値を前記インデックスに追加するステップと、
所定の閾値を超える一意なバイト値の頻度を持つ前記サブシーケンス上で、該サブシーケンス内で 1 以上の一意な値の位置を表すデータ値を前記インデックスに追加するステップと、
ファイルタイプ識別子を持つ出力データファイルを作成するステップと、
前記出力データファイルに前記インデックスを追加するステップと、
を含むことを特徴とする方法。

10

【請求項 2】

前記サブシーケンスに対してデータ変換を適用するステップは、
コンピュータメモリ内に複数の変換データセットを保持するステップであって、該変換データセットは一連のバイト値を持ち、且つ変換データセット識別子によって識別されるステップと、
コンピュータメモリから前記変換データセットの一つを読み出すステップであって、該読み出された変換データセットは前記データファイルのサブシーケンスの長さを実質的に等しい長さを持つステップと、
各々のバイトの値に対して、前記読み出された変換データセットにおいて対応する各々のバイトの値に基づいてデータ変換を適用するステップと、
を含む請求項 1 に記載の方法。

20

【請求項 3】

前記読み出された変換データセットの少なくとも一つに基づいてデータ変換された前記サブシーケンスは、該データ変換前の前記サブシーケンスと実質的に同一である、請求項 2 に記載の方法。

30

【請求項 4】

前記変換データセットの少なくとも一つは、ランダムに生成されたバイトレートのシーケンスを含む、請求項 2 に記載の方法。

【請求項 5】

前記変換データセットの少なくとも一つは、バイトレートの所定のシーケンスを含む、請求項 2 に記載の方法。

【請求項 6】

前記変換データセットの少なくとも一つは、前記データファイルのサブシーケンス以外の該データファイルの一部から得られたバイト値のシーケンスを含む、請求項 2 に記載の方法。

40

【請求項 7】

前記サブシーケンスに対して適用された前記データの変換データセットの前記変換データセット識別子を前記インデックスに対して追加するステップをさらに含む、請求項 2 ~ 6 の何れか一項に記載の方法。

【請求項 8】

前記サブシーケンス内で 1 以上の一意な値の前記位置を計算するステップをさらに含む、請求項 1 ~ 7 の何れか一項に記載の方法。

50

【請求項 9】

前記サブシーケンス内で前記 1 以上の一意な値の前記位置を計算するステップは、コンピュータメモリ内にテンポラリポジションインデックスを作成するステップと、前記サブシーケンスから連続したバイトの値を読み出すステップと、各バイト値の読み出しにおいて、該読み出されたバイト値が一意なバイト値か繰り返された値かを決定するステップと、

一意なバイト値の検出において、二つのビットの値の一つを前記テンポラリポジションインデックスに追加し、さもなければ、該二つのビットの値の他方を前記テンポラリポジションインデックスに追加するステップと、

前記テンポラリポジションインデックスから、前記 1 以上の一意な値の前記位置を表すポジションインデックスを作成するステップと、

前記ポジションインデックスから、少なくとも部分的に前記 1 以上の一意な値の前記位置を表す前記データ値を計算するステップと、
を含む請求項 8 に記載の方法。

【請求項 10】

前記サブシーケンス内のバイト数が、前記テンポラリポジションインデックス内のビット数と実質的に等しい、請求項 9 に記載の方法。

【請求項 11】

前記ポジションインデックスのサイズは、前記テンポラリポジションインデックスのサイズよりも小さい、請求項 9 又は 10 に記載の方法。

【請求項 12】

前記サブシーケンス内の一意なバイトの値の順序を表す順列インデックスを作成するステップと、

前記ポジションインデックス及び前記順列インデックスの両方から前記 1 以上の一意な値の前記位置を表す前記データ値を計算するステップと、
をさらに含む請求項 9 ~ 11 の何れか一項に記載の方法。

【請求項 13】

前記 1 以上の一意な値の前記位置を表す前記データ値を形成するために、前記ポジションインデックス及び前記順列インデックスを連結させるステップを含む、請求項 12 に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ圧縮の分野に関連し、特に階乗反復型のロスのない圧縮に基づくデータ圧縮システム及び方法に関する。

【背景技術】

【0002】

電子バイナリファイルには、多くの様々な用途のために多く異なるフォーマットがある。これらのフォーマットには、画像、音、テキスト、データ、実行ファイルなどの保存に適した形式が含まれる。

データを含むバイナリファイルは、暗号化されていなければ、構造化されたフォーマットとなる傾向にある。通常、ヘッダ情報、テキスト、頻度及び他の構成要素間の配置がある。一般に、バイナリファイルの最初の数バイトは、ファイルタイプを表す指標を含むので、どのバイナリファイルを用いるアプリケーションも互換性を有する。実行ファイル若しくは何がしかの機能を実行するために使用されるファイルは、ほとんど構造化されたフォーマットを持たない。しかしながら、これらのファイルは機能を実行するためにオペレーティングシステムと情報をやりとりしなければならないか、オペレーションシステムの一部であるため、構造的要素がある。

【0003】

圧縮され、暗号化されたファイルは、設計によってはファイル内で繰り返される値が除

去されるので、最も構造を持たない。暗号化の場合、置換される値を定義するために鍵が使用される。圧縮については、“省略表現”が繰り返す構造に対して使用される。暗号化されたファイル又は圧縮ファイルの場合、そのファイルは内部構造を持たないばかりか、特に圧縮の場合、ファイルのサイズも変更される。

【0004】

1, 048, 576バイト(1Mb)のサイズのバイナリファイルに対して、数学的に可能なバイトの配列について256¹⁰⁴⁸⁵⁷⁶通りのとり得る構造がある。実際の使用では、この数の一部のみが使用される。実際に使用される数は、異なるファイルタイプの数の推定、実行ファイル又はオペレーショナルファイルの機能性、及び利用可能な圧縮及び暗号化に基づいて近似するしかない。

10

【0005】

データファイル上でデータ圧縮を実行する多くの技術が存在する。幾つかのデータ圧縮アルゴリズムは、インデックス化技術に基づいており、データファイル内の一意の値のインデックス化と計算を含む。最も圧縮されたデータファイルでは、256バイトの各コードセグメント内でデータ値の繰り返しが幾つか存在する。平均的なファイルでは、256バイトのコードセグメントごとに160から170個の繰り返しのない一意な値が存在する。階乗計算に基づいたデータ圧縮技術は、この値の数ではうまく動かない。

【発明の開示】

【0006】

一つの側面において、本発明は、所定長以上の長さのバイトシーケンスを有するデータファイルの圧縮方法を提供する。その方法は、二次記憶装置からデータファイルを取り出すステップと、データファイルをダイレクトアクセスメモリに保存するステップと、データファイルの所定長を超えないサブシーケンス内で一意なバイト値の頻度を計算するステップと、サブシーケンス内で計算された一意なバイト値の頻度を表すデータ値を含むインデックスを作成するステップと、所定の閾値未満である一意なバイト値の頻度を持つサブシーケンス上で、サブシーケンス内で一意なバイト値の頻度を増加するためにサブシーケンスに対してデータ変換を適用し、データ変換を表すデータ値をインデックスに追加するステップと、所定の閾値を超える一意なバイト値の頻度を持つサブシーケンス上で、サブシーケンス内で1以上の一意な値の位置を表すデータ値をインデックスに追加するステップと、ファイルタイプ識別子を持つ出力データファイルを作成するステップと、出力データファイルにインデックスを追加するステップとを含む。

20

30

【発明を実施するための最良の形態】

【0007】

本発明のデータ圧縮システム及び方法の好ましい形式を、ここで添付図面を参照しつつ説明する。

本発明は、データファイル5に対して適用しようとするデータ圧縮システム及び方法を提供する。データファイル5は、BMP、WAV、DOC、XLS、MDB、ZIP、SIT、ARJ、ZOO、TIF、JPG、GIF、MP3、MP4などを含む適当なデータとすることができる。データファイル5を、コンピュータ装置15の一部を形成するか、少なくとも接続される二次記憶装置10に保存することができる。コンピュータ装置15は、ダイレクトアクセスメモリ25及びディスプレイ30と接続されるプロセッサ20を少なくとも含む。コンピュータ装置は、他のコンポーネント、例えばデータ入力装置(図示せず)及び出力装置(図示せず)を含むか、接続されてもよいことを理解されたい。

40

データファイル5は所定長以上の長さのバイトシーケンスを含むことが想定されている。本発明の一つの好ましい実施形式では、この所定長は300バイトである。

【0008】

動作の際、コンピュータ装置15のプロセッサ20は、二次記憶装置10からデータファイル5の全て又は一部を読み出す。読み出されたデータファイル又はその一部は、ダイレクトアクセスメモリ25に保存される。様々な操作がその保存されたデータファイル又はその一部に対して行われる。得られた出力データファイル35は、ダイレクトアクセス

50

メモリ 25 に作成され、二次記憶装置 10 又は他の記憶装置に保存される。多くの場合において、出力データファイル 35 は、データファイル 5 よりも小さいサイズとなることが予想される。

データファイル 5 のサブシーケンスは最初に検査される。サブシーケンスの長さは、所定長の 300 バイトを超えないことが好ましい。特定された一意な値の数が閾値を下回る場合、一連のデータ変換がそのサブシーケンスにおける一意なバイト値の頻度を増加するためにサブシーケンスに対して適用される。

【0009】

複数のデータ変換コンポーネント 40 はダイレクトアクセスメモリ 25 又は二次記憶装置に保存される。データ変換コンポーネント 40 は、ランダムに生成されたバイト値のシーケンス又は所定のバイト値のシーケンスを複数含んでもよい。そのシーケンスはマスク構造 45 として保存される。あるいは、若しくは好ましくは、さらにデータ変換コンポーネントは、追加のマスク構造 45 を生成するために使用することができる複数のマスク式 50 も含む。データ変換コンポーネントのアプリケーションについては、後で説明する。

10

【0010】

システムは、複数のインデックス化コンポーネント 60 も含む。データファイル 5 のサブシーケンスの処理中、出力データファイル 35 に同時に書き込まれるインデックス 65 が作成される。インデックス化コンポーネント 60 はテンポラリポジションインデックス 70、ポジションインデックス 75 及び順列インデックス 80 を含んでもよい。幾つかの場合においては、ポジションインデックス 75 及び順列インデックス 80 の内容がインデックス 65 に追加される。様々なインデックス化コンポーネント 60 の動作については、後で説明する。

20

【0011】

システムは、ダイレクトアクセスメモリ 25 又は二次記憶装置に保存されるデータアレイ 90 を含んでもよい。データアレイ 90 は、データアレイ 90 の内容が出力データファイル 35 に書き込まれる前に、様々なインデックス化コンポーネント 60 及び圧縮されるデータファイル 5 のサブシーケンスの一部を保存するために使用することができる。

【0012】

図 2 から図 4 は、本発明の好ましい形式の動作を示す。バイナリデータファイル 5 は、複数のデータグループに分割されることが好ましい。本発明の一つの好ましい実施形式では、各データグループは 300 バイト以下であることが好ましい。しかし、圧縮されるデータグループのサイズは 5 ビットを超える如何なるサイズであってもよいことを理解されたい。最初に、データファイルは、データファイルの長さが所定長以上か否かを明らかにするためにチェックされる (ステップ 200)。本発明の一つの好ましい実施形式では、所定長の初期値は 300 バイトである。一つの形式では、データファイル全体は二次記憶装置から読み出され、ダイレクトアクセスメモリ 25 内のデータアレイ 90 に保存される。あるいは、データファイル 5 の一部を、データストリームとして二次記憶装置 10 から読み出してもよい。

30

【0013】

データグループは、データグループ内の一意なデータ値の頻度を計算するためにカウントされる (ステップ 205)。そして一意なデータ値の頻度は、所定の閾値と比較される (ステップ 210)。一つの好ましい形式では、所定の閾値は 256 である。300 バイトのサブシーケンス内で一意な値が 256 未満の場合、サブシーケンス内の一意なバイト値の頻度を増加するために、1 以上のデータ変換をそのサブシーケンスに適用してもよい。

40

一意なバイト値の頻度が、300 バイト内で所定の閾値の 256 未満の場合、そのサブシーケンスは、データ変換マスクをサブシーケンスに適用可能か否かを識別するためにテストされる (ステップ 215)。本発明の一つの好ましい形式では、構造ライブラリが、コンピュータメモリ、例えばダイレクトアクセスメモリ 25 に保持される。そのライブラリは、複数のランダムに生成されたデータセットを含むことが好ましい。これらのデータセ

50

ットを、それぞれデータセット識別子によって識別することができ、そのデータセット識別子はコンピュータメモリに保存され、ランダムに生成されたデータセットのそれぞれと関連付けられる。

【0014】

一つの形式では、ランダムに生成されたデータセットの少なくとも一つは、データファイルのサブシーケンスの長さを実質的に等しい長さを持つ。言い換えれば、サブシーケンスのバイト数は、変換データセット又はマスクのバイト数と同じである。そのようなマスクは、サブシーケンス内のそれぞれのバイト値に対して、対応するバイト値及び読み出された変換データセットに基づいて、データ変換を適用することにより、サブシーケンスに適用できる。

10

データ変換の一つの例は、加算剰余(modulus addition)である。サブシーケンスの最初のバイト値とデータセットの最初のバイト値とが加算され、その合計について256を法として剰余が計算される。例えば、サブシーケンスの最初のバイナリ値が168であり、特定されたデータセットの最初のバイナリ値が203の場合、その合計は371である。変換値は、 $371 \text{ MOD } 256$ を計算することにより、115である。その後、シーケンスの2番目のバイトが、データセットの2番目のバイトにより同じ方法で変換される。その後、シーケンスの3番目のバイトが、データセットの3番目のバイトにより同じ方法で変換される、などとなる。

この方法で、マスクがサブシーケンスに適用される(ステップ220)。

【0015】

20

一つの形式において、65, 536個のマスク構造をコンピュータメモリに保存しておくことができ、各マスクは0から65, 536までのインデックスナンバーという形式のデータセット識別子とともに提供される。そのインデックスは、関連するデータセット識別子を示す単なる14ビットセグメントとすることができる。

データ変換コンポーネント40は、マスク式を含むことができる。例えば以下の通りである。

- ・先の300バイト以下のデータファイルのシーケンスの標準偏差。この式は前にサブシーケンスのない、データファイルの最初のシーケンスでは使用できないことを理解されたい。

30

- ・前回のサブシーケンス又は標準偏差に基づいたサブシーケンス内の値の反転
- ・サブシーケンスの構造に基づいて計算される適用可能な構造
- ・関連するサブシーケンスに対して加算又は減算され、ファイル構造に基づいてランダムに生成されるセグメント

【0016】

上記の式は、一連のマスク構造を予め生成するために適用してもよい。あるいは、データ変換中に関連するバイト値を計算してもよい。一つの形式では、512個のランダムに生成された構造又はマスク構造は、ダイレクトアクセスメモリ25に保存される。これらの構造は、300バイトシーケンス内で256以上のヌル値を持つ可能性のあるデータファイルのサブシーケンスに対して適用される。これは、多くのソフトウェアアプリケーションのバイナリファイルのヘッダ構成部分において一般的である。これらのランダムに生成された構造を、高いレベルの反復をもつ他のフォーマットに適用することもできる。

40

【0017】

サブシーケンスでのデータ変換に続けて、サブシーケンスは、300バイトのシーケンス内で256個の一意な値があるか否かを識別するために再度テストされる(ステップ210)。256個の一意な値がなく、そのサブシーケンスに適用可能な別のマスクがなければ、300バイトの閾値は下げられ、プロセスはより小さなサブシーケンスで繰り返される。一つの好ましい実施形態では、閾値は、300未満の8ビット値(バイト)を検査するために、一時的に152個の7ビット値又は77個の6ビット値へ低下させてもよい。それから、閾値は、次のサブシーケンスに対して300バイトに引き上げられる。これについては以下に詳細を述べる。

50

【0018】

ランダムファイルの追加が256バイトセグメント内で256個の一意的な値を生成することは到底無理であり、約10%が可能なところである。一度適切なランダムファイル構造が適用されると、300バイトを超えないデータセグメント内で256個の一意的な値があることが想定される。場合によっては、データ変換の意図はデータグループ内の一意的なデータ値の頻度を増加することにある。

【0019】

本発明は、データグループ内の300個のデータ値のインデックスを計算する。

インデックスは、ダイレクトアクセスメモリ25内のデータアレイ90に保存されることが好ましい。300個のデータ値のインデックスは、最初に2ビットで生成される。256個の一意的なデータ値が300バイトのデータグループ内で特定される場合、そのビット値“01”がインデックスに書き込まれる(ステップ225)。

マスクがサブシーケンスに適用されると、マスク又はデータセット識別子がインデックスに書き込まれる(ステップ230)。このマスク識別子は、0から65,536間のマスク値を識別する16ビット値であることが好ましい。マスク識別子において値が0であるということは、サブシーケンスに対してマスクが適用されなかったか、ヌルマスクが適用されたことを表す。ヌルデータセットがサブシーケンスに適用される場合、データ変換後のサブシーケンスは、データ変換前のサブシーケンスと実質的に同一である。

【0020】

本発明の方法において、次のステップは、テンポラリポジションインデックスを作成することである(ステップ235)。

テンポラリポジションインデックスの作成方法は、データグループの最初のバイトにおいて開始し、256個の一意的な値が300バイトのデータグループから抽出される場合、256個の一意的な値が特定されるまで、データグループ内の次のバイトを検査する。検査される特定の値が、そのデータグループ又は以前のデータグループ内ではじめて現れる値である場合、“1”ビット値がテンポラリインデックスに加算される。一方、検査されるデータ値が既出のデータ値の繰り返しである場合、“0”ビット値がインデックスに書き込まれる。インデックス化の方法は、256個の“1”ビットがインデックスに書き込まれると、直ちに終了する。

【0021】

テンポラリインデックスは、作成される圧縮データストリームにおいて、データグループ内の各データ値の配置と識別を容易にする。インデックスで値“1”の数は、使用されているビット数を示す。例えば、256個の“1”の値が、テンポラリインデックス内の283個のエントリの後のテンポラリインデックスで発生した場合、これはサブシーケンスの283バイト内に256個の一意的なバイト値があることを示す。

300バイトのデータグループ内に256個以上の値がある場合、インデックスの最初の2ビットは既に“01”にセットされている。テンポラリインデックスはメインのインデックスに対して単に加算されるが、一方、この情報を格納する、より効率的な方法がある。サブシーケンスにおいて出現する“1”の値の数は既知である。それらが出現する順番を無視するならば、一意的なバイト値の発生数を記録するだけで十分である。

【0022】

テンポラリインデックスそれ自体を記録するよりも、ポジションインデックスを生成し、このポジションインデックスをメインのインデックスに書き込む方が好ましい。300バイトのサブシーケンスに対して、テンポラリインデックスが、256個の“1”の値とそれに続く44個の“0”の値を含む場合、これにポジションインデックス“0”を割り当ててもよい。44個の“0”の値と256個の“1”の値が300バイトのデータグループ内で取り得る配置の数は nC_r である。これは、256個の“1”の値と44個の“0”の値がある場合、300個の値内に 1.34×10^{53} と等しい $300! / (256! \cdot 44!)$ 通りの組み合わせがあるということの意味する。

1.34×10^{53} 個のこの最大ポジションインデックス値は、 2^{177} よりも小さく

10

20

30

40

50

、値を表現するために177ビットを必要とする。

これは、300ビットの実際のテンポラリインデックスを保存するよりも、テンポラリインデックス内に少なくとも256個の“1”の値があるという事実を利用することによって、代わりにポジションインデックスは177ビット若しくは22.125バイトで記録できるということを意味する。

【0023】

圧縮だけでなく解凍も可能にするために、データグループ内のデータ値の順序も記録することが重要である。これは順列インデックスを生成し、この順列インデックスをメインのインデックスに書き込むことにより達成される(ステップ245)。

順列インデックスの計算は、256個の一意な値を並べることができる方法の数、すなわち繰り返しのない256個の値の順列に基づく。最初の値に対しては、256通り有り、2番目の値に対しては255通り有り、3番目の値に対しては254通り有る、などである。これは256!として表され、“256の階乗”として参照される。そのため256個の一意な値の可能な順列の数は、 8.57×10^{506} 通りである。この値は、 2^{1684} が 8.57×10^{506} よりも大きい 8.6×10^{506} と等しいので、1684ビットで表すことができる。1684ビットは210.5バイトと等価である。

シーケンス0, 1, 2, 3, 4, . . . , 254, 255は順列番号1で表し、シーケンス255, 254, 253, . . . , 3, 2, 1, 0は順列番号 8.57×10^{506} で表す。

【0024】

順列インデックスは、メインのインデックスに書き込まれる。メインのインデックスは、これまでにサブシーケンス内の一意なビット値の計算された頻度を表すデータ値を含む。これは、ビット値“01”、続いて適用されたマスクを表す16ビット、その次にポジションインデックスを表す177ビット、その次に順列インデックスを表す1684ビットが続く。

【0025】

十分な長さのサブシーケンスを得るにはデータファイル内の残りのビットが十分でないか、一意な値の残りが十分でないところに到達すると、インデックスは出力ファイルに書き込まれる(ステップ250)。

出力ファイルは、ファイルタイプを識別するために先頭3バイトを含むことが好ましい。さらにファイルタイプ識別子に続く2バイトは、最大65,536回に対する特定のデータファイル全体にわたって実行された本発明の方法の試行数を示す。

これらの5バイトに続いて、データアレイ90に保存されるインデックスが出力ファイルに加えらる。インデックスに続いて、インデックスで未処理な、すなわちデータファイル内に残るビットの値又は一意な値が十分でない値が加えられる。

たいていの場合、ヘッダの5バイトと、それに続く本体と、出力ファイルの終わりに未圧縮の形式でフルに書き込まれる63以下のビットの値とが存在することが予想される。出力ファイルの本体は、ストリーミングの手法で抽出することを容易にするために連続して書き込まれるインデックスの集合であることが好ましい。

【0026】

図2について上述してきたように、本発明の方法を複数回繰り返した後にデータファイルに残るバイトが300もないか、256個の一意な値及び適用可能な別のマスクがない300バイトのシーケンスがある場合も発生する。ステップ260で示されるように、一つの好ましい形式では、データファイルから読み出されるサブシーケンスのサイズを減らすことができる。

【0027】

図3を参照すると、データファイルに少なくとも133バイトが残っているか否かを識別するために、データファイルはチェックされる(ステップ305)。

データファイルに残る152個の7ビット値を含む少なくとも133バイトが有る場合、152個の7ビット値内で一意な値の数をカウントする(ステップ310)。その後、

10

20

30

40

50

一意な値の数が、閾値（例えば 128）に対してチェックされる（ステップ 315）。133 バイトのサブシーケンスに十分な一意の数がなければ、図 2 のステップ 215 及び 220 と同様の方法で適用可能なマスクが特定され（ステップ 340）、適用される（ステップ 345）。

一度、一意な値の数の閾値がデータファイルの 152 個の 7 ビット値に対して確認されると、ビットシーケンス “10” がインデックスに書き込まれ（ステップ 350）、本方法は図 2 において 230 で示されるステップへ進む。

【0028】

データファイル内に処理されるべき残りが 152 個の 7 ビット値もないか、152 個の 7 ビットサブシーケンス内に 128 個の一意な値を見つけられず、適用可能な別のマスクもない場合、ステップ 355 に示されるように、本方法は図 4 に示されるものへ渡される。図 4 に示すように、検査下にあるデータファイルのビットグループの数は、77 個の 6 ビット値に減らされる。データファイルに 77 個の 6 ビット値が残っている場合（ステップ 405）、その 77 個の 6 ビット値において一意な値の数がカウントされる（ステップ 410）。

10

【0029】

一意な値の数が閾値 64 に対してチェックされる（ステップ 415）。77 個の 6 ビット値において一意な値が 64 よりも少なければ、本方法はマスクが適用可能か否かを明らかにする（ステップ 420）。マスクが適用可能であれば、マスクが適用される（ステップ 425）。これら最後の 2 ステップ（425、430）は、図 2 のステップ 215 と 220、及び図 3 のステップ 340 と 345 と同様である。

20

77 個の 6 ビットサブシーケンスにおいて 64 個の一意な値があれば、値 “11” がインデックスに書き込まれる（ステップ 430）。そして制御は図 2 のステップ 230 へ戻される。

データファイル内に処理されるべき残りが 77 個の 6 ビット値もないか、77 個の 6 ビットサブシーケンス内に 64 個の一意な値を見つけられなければ、ビット値 “00” がインデックスに書き込まれ（ステップ 435）、インデックスは図 2 に示されるステップ 250 と同様の手法で出力ファイルに書き込まれ、且つデータファイルの残りのバイトは出力ファイルに書き込まれる。

【0030】

検査下にあるバイト数に依存して、図 2 のステップ 245 で示される順列インデックスに対しても少々変更を必要とすることを理解されたい。152 個の 7 ビットグループ内で 128 個の一意なデータ値が有る場合、ポジションインデックスは、 $152! / (128! \cdot 24!)$ 通りとなり、これは 5.48×10^{27} と等しい。これは、 $2^{93} = 9.9 \times 10^{27}$ なので、93 ビットで表すことができる。

30

77 個の 6 ビットグループ全体にわたって 64 個の一意な値が有る場合、インデックスは $77! / (64! \cdot 13!)$ 通りとなる。これは、 1.84×10^{14} となる先の値よりも大きい $2.81 \times 10^{14} = 2^{48}$ なので、48 ビットで表すことができる。

【0031】

同様に、検査下にあるバイト数に依存して、図 2 のステップ 245 で示される順列インデックスに対しても少々変更が必要となる。128 個の値に対する順列は、 $128!$ すなわち 3.86×10^{215} である。これは、 $2^{717} = 6.89 \times 10^{215}$ なので、表現するために 717 ビットを必要とする。

40

64 個の値に対する順列は、 $64!$ すなわち 1.27×10^{89} である。これは、 $2^{296} = 1.27 \times 10^{89}$ なので、296 ビットで表現することができる。

【0032】

図 5 は、377 バイト、350 バイト、320 バイト、300 バイト（8 ビットグループ）、152（7 ビットグループ）及び 77（6 ビットグループ）のデータグループサイズにおける想定される結果の表を示す。この表に示されるのは、バリエーションに包含される効果の指標である。これについては後で説明する。

50

【0033】

解凍は、上記の手順を単に反転するものである。インデックスの値は、最初から最後(256番目)までの各値の範囲を示す。範囲を備えることは関連する値を提供する。インデックスは、ヘッダとともに再構成に使用することができる。全てのコンポーネントは一緒にバックされるので、ストリーミングが使用されることが予想される。

【0034】

繰り返される値の配置のインデックス化は、より効率的な方法がある場合、セグメントに対する“0”と“1”の値のストリングから変更してもよい。例えば、繰り返される値が1個か2個のみの場合、バイト数は257か258となる。257番目のビット又は258番目のビットを使用するよりも、最初と最後のバイトがそのセグメントに対して一意となることが知られている。そのため、257値の場合には、8ビットが単一の繰り返し値の位置を提供し、16ビットが258バイトセグメントの場合における両方の繰り返し値の位置を提供する。

10

【0035】

本方法は、全てのファイルタイプ及び構造に適用することができる。PKWareのZIP製品のようなツールでかなりの量まで圧縮されたファイルタイプ又は構造に対して、本発明の方法は1回の試行では同じレベルに到達しないであろう。しかし、本方法は、同じファイルに対して繰り返し適用し、各回でそのサイズを減らすことができる。試行数すなわち繰り返す数は、処理するハードウェア及び/又はユーザが要求する時間に依存する。

20

【0036】

全てのコンポーネントが既知であるため、解凍は極めて高速である。圧縮はランダムデータ構造とのマッチングを必要とするので、解凍は圧縮よりもさらに高速となり得る。

全てのインデックス化は実際のデータそれ自体の中に包含されるので、複数回の解凍ルーチンを同時に実行してもよい。

他のアプリケーションは、ソフトウェア圧縮、データ圧縮、ソニープレイステーション2(登録商標)、マイクロソフト(登録商標)X-Boxなどのような対戦型オンラインゲーム、ボイスオーバーIP、ビデオオンデマンドを含む。本発明は、データ又はバイナリ情報が保存され、変換され、若しくは如何なるフォーマットで使用される如何なるアプリケーションも含む。

30

【0037】

上記の記述は、300バイトのコードセグメント内の256個の一意的な値又はそれよりも小さいものに基づいている。この選択された値は単に説明を目的とするものであることを理解されたい。5ビット又はそれ以上のデータグループ、又は0から31の間の値はこの方法を用いて再構成することができる。ランダムに生成されたデータセット又はオーバーレイファイルの数を減らすことは、3及び4ビット値がよく使用されるということの意味する。

説明してきた8ビット(256値)よりも大きなビット値を用いて、より大きな削減を行うことができる。例えば、9ビット値で圧縮された場合、8ビットの圧縮により達成される圧縮をさらに超える圧縮ゲインがある。

40

【0038】

削減若しくは圧縮は、値に対して使用されるビット数とともに増大する。256値(300バイトセグメント)は、512値(600バイトセグメント)ほど圧縮しない。同様に、512値(600バイトセグメント)は、1024値ほど圧縮しない。計算はファイルサイズに基づくはずなので、上限レベルはない。

上記の300バイトの方法を用いて、これを377バイトグループに拡張することができる。これは、図5に示され、且つこの明細書で説明される好ましい実施形態に対して最適レベルである300バイトを用いて、効果のある範囲が256から377バイトのグループであることを意味する。

【0039】

50

300個の8ビットグループ(バイト)のバリエーション、152個の7ビットグループ及び77個の6ビットグループを、圧縮されたファイルのヘッダに示してもよい。そのバリエーションは、二つのパートから構成されてもよい。それらは、以下のものである。

1. セグメントサイズに対する関連するビットグループの数の表示。8ビットグループに対するサイズの範囲は256から377であり、7ビットで表され得る。7ビットグループに対する範囲は5ビットで表され、6ビットグループに対する範囲は4ビットで表され得る。

2. ビットグループのそれぞれ内で変更があるか否かを示すために上記のそれぞれの終端に追加され得る追加ビット。“0”は変更無しを示し、“1”は変更有りを示すことができる。

ヘッダは、上記の値を示す追加の19ビットを含んでもよい。

ヘッダごとに、変更値が許されるならば、グループごとを基礎として変更値をインデックスに書き込んでもよい。

例えば、8ビットグループの一つのグループのデフォルトは300の値としてもよいが、各セグメントは、含まれる変更値によって示されるように、256から377の間で変化してもよい。

【0040】

本発明の別の実施形態は、複数の繰り返しバイト圧縮エンハンスメントを含む。これは図6及び図7を参照して説明される。

機能的電子ファイルは、複数の異なるカテゴリのバイト構造に分類される。これらは単純な2色のビットマップから、これまでに利用可能なロスのない圧縮アルゴリズムを用いて圧縮されたファイルまで、様々である。

【0041】

ヘッダ情報に続く、2色のビットマップについては、ビット値1は黒を意味し、他は白を意味する。多くの繰り返しがあるので、ロスのない手法でのこれらのファイルの圧縮は単純である。

24ビットのビットマップに移ると、パターンの識別はより困難となり、そのためロスのない圧縮の圧縮率は、現在のアルゴリズムを用いて、より単純なビットマップ構造上での圧縮率と同程度に大きくはならない。

【0042】

ここに説明するプロセスは、24ビットのビットマップよりも単純なパターンを導入し、それは、標準的な写真タイプの画像に対して、圧縮量を飛躍的に増大させるように、現在利用可能なロスのない圧縮アルゴリズムを用いてロスのない圧縮を可能にする。

【0043】

これを達成するために、図6の610に示されるように、オリジナルのイメージは3個の成分に分解され、結合されたもののサイズはオリジナルのイメージよりも非常に大きくなる。それから、620に示されるように、全ての3バイト(24ビット)グループは、10進数の昇順に並べられる。例えば、236, 217, 67は、67, 217, 236に並べ替えられる。バイト配置の変更は、ハフマン構造を用いてインデックスに記録される。

オリジナルの構造が6通りしかないので、これらは以下のビットインデックスを用いて記録される。

00 = 123

01 = 132

100 = 213

101 = 231

110 = 312

111 = 321

上記の数のそれぞれは、バイトの並び替えられた位置と比較したときのそのバイトのオリジナルの位置を表す。

10

20

30

40

50

【 0 0 4 4 】

6 2 5 に示されるように、イメージが完全に走査されると、このインデックスはファイル（ファイル A ）に書き込まれる。

6 3 0 に示されるように、最小のもの全て、又は各グループから現在の最初のバイトの値が別個のファイル（ファイル B ）に書き込まれる。

6 3 5 に示されるように、バイトの値は順番に並んでいるので、2 番目のバイトの値から最初のバイトの値を引いた値がファイル（ファイル C ）に書き込まれ、直ぐに 3 番目のバイトから 2 番目のバイトの値を引いた値が続く。

これで、3 個のファイル、ファイル A、ファイル B、ファイル C が生成される。ファイル B とファイル C を合わせた合計は、オリジナルの 2 4 ビットのビットマップと同じである。ファイル A はバイトのインデックスを表すので、ファイル A はサイズにおいて余分のオーバーヘッドである。

【 0 0 4 5 】

その後、全ての 3 個のファイル（ A、B 及び C ）がロスのないアルゴリズム又は W I N Z I P 6 4 0 のような製品を用いて一つのファイル（ 6 5 0 ）に圧縮されると、得られたファイルは、未修正のイメージファイルにこれらのツールを単に適用することによって得られるものよりも平均で 2 5 % 小さい。

テストでは、最悪の場合のシナリオで 2 . 5 % の低下を示し、最高の場合、画質が 2 4 ビットのトゥルーカラービットマップの 8 2 % であった。同じゲインは J P E G のロスなしの圧縮モードを用いて得ることができる。

このプロセスは、データを保持する 3 バイトのグループ分けを用いて如何なるファイル構造にも適用できる。それはまた、ロスのない圧縮のレベルをより大きくしていくために、4、5、6、7、8 等のバイト構造をカバーするように拡張してもよい。

【 0 0 4 6 】

ビットマップファイルが、イメージを表示するために使用されるように、W a v e (. w a v) ファイルは、音を出すために使用される。圧縮のエンハンスメントプロセスの別の例を、図 7 を参照しつつ W a v e フォーマットファイルに関してここで説明する。より多い色又は質を提供するそれぞれのビットマップファイル（ 2 ビット、4 ビット、8 ビット、1 0 ビット、1 2 ビット、1 6 ビット、2 4 ビット、3 0 ビット）の異なるレベルがあるように、同じことが W a v e ファイルにも生じる。

W a v e ファイルは複数の成分を用いて作成され、その成分は平均サンプリングレート、サンプリングレート、オーディオサンプルサイズ及びチャンネル数である。

サンプリングレートが低いほど、ファイルサイズも小さく、質も劣化することを意味する。また、モノラルファイルはステレオファイルよりも小さい。

【 0 0 4 7 】

ここで取り上げる W a v e フォーマットは、商用 C D に最高品質のステレオ音楽を保存するときに使用されるフォーマットである。このフォーマットは W a v e フォーマットから C D フォーマットに変換される。

平均データレートが 1 7 6 . 4 K b / 秒、サンプリングレート 4 4 . 1 k H z、オーディオサンプルサイズ 1 6 ビット、2 チャンネル（ステレオ）の W a v e ファイルに対して、次のように適用される。

【 0 0 4 8 】

ファイル内の全てのバイトの値が、n をファイルの最後のバイトとして（通常のオーディオファイルに対して、これは 5 0 , 0 0 0 , 0 0 0 のオーダとなるであろう）、1 から n の番号で表される場合、全ての偶数の位置のバイトの値が一つのファイル（ファイル 1 ）に書き込まれ（ 7 2 5 ）、全ての奇数の位置のバイトの値が別のファイル（ファイル 2 ）に書き込まれる（ 7 3 0 ）。例えば、以下の表の通りとなる。

10

20

30

40

【表 1】

バイト値	255	167	33	0	0	24	24	167	167
順序値	1	2	3	4	5	6	7	8	9
ファイル1 (odd)	255	33	0	24	167				
ファイル2 (even)	167	0	24	167					

それから、両方のファイル（ファイル1及びファイル2）が、再度640で示されるように、ロスのないアルゴリズム又はWINZIPのような製品を用いて一つのファイルに圧縮されると、得られたファイル650は、未修正のイメージファイルにこれらのツールを単に適用することによって得られるものよりも平均で20%小さい。

10

テストでは、圧縮ファイルのサイズにおいて最悪の場合のシナリオで10%のさらなる低下を示し、最高の場合、サイズにおいて43%の低下であった。

抽出/解凍は単純であり、二つのファイルが関連するロスのないツールを用いて解凍された後に、ファイル2からのバイトがファイル1のバイトのそれぞれの間に挿入される。

【0049】

上記は、その好ましい形式を含む本発明を説明する。当業者にとって明らかな変更や修正は、添付の特許請求の範囲で規定されるように、本発明の範囲に含まれることが意図されている。

20

【図面の簡単な説明】

【0050】

【図1】本発明のシステムの好ましい形式を示すである。

【図2】本発明の好ましい形式のデータ圧縮プロセスのフローチャートである。

【図3】本発明の好ましい形式のデータ圧縮プロセスのフローチャートである。

【図4】本発明の好ましい形式のデータ圧縮プロセスのフローチャートである。

【図5】本発明の好ましい実施形態について想定されるデータ圧縮結果の表を示す図である。

【図6】複数の繰り返しバイト圧縮エンハンスメントに関する本発明のさらなる側面を表す図である。

30

【図7】複数の繰り返しバイト圧縮エンハンスメントに関する本発明のさらなる側面を表す図である。

【 図 1 】

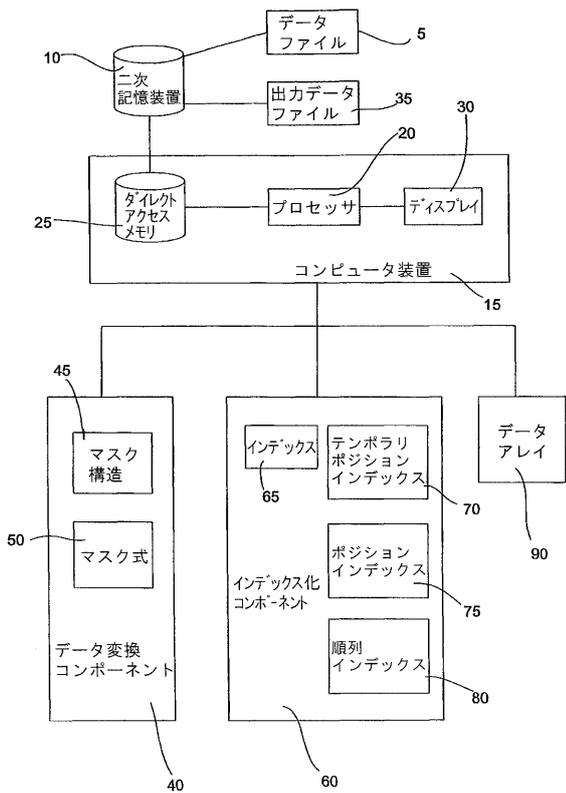


FIGURE 1

【 図 2 】

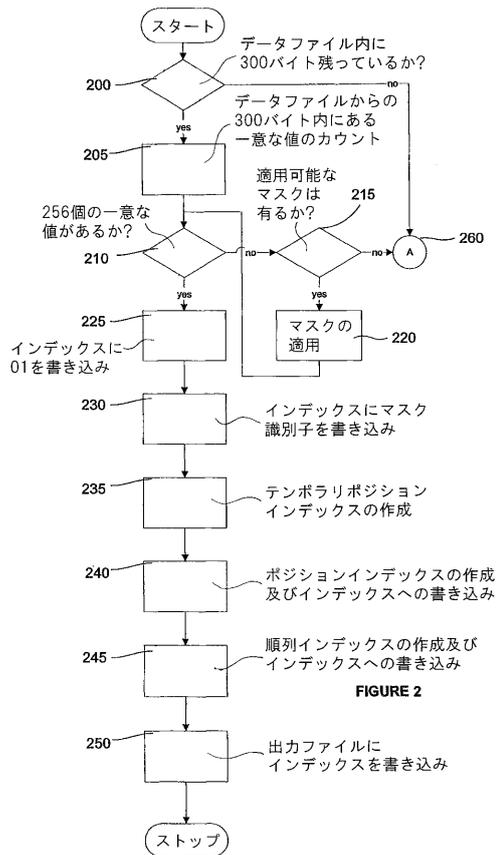


FIGURE 2

【 図 3 】

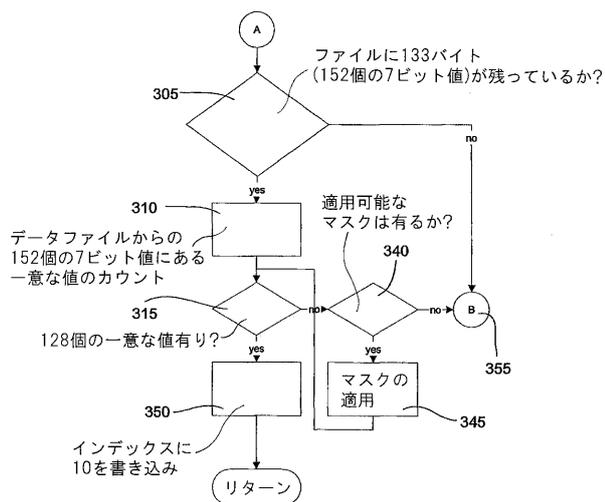


FIGURE 3

【 図 4 】

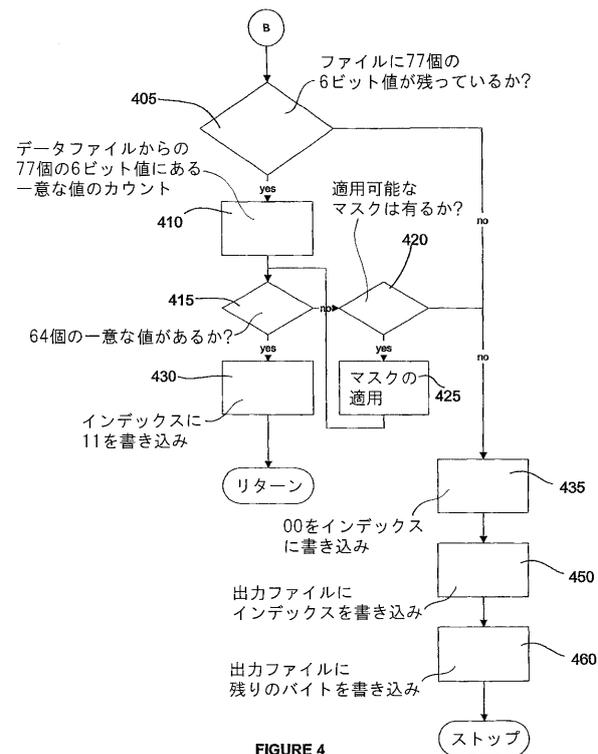


FIGURE 4

【 図 5 】

シーケンス サイズ	バイト (ビット)数	シーケンス サイズを示す インデックス ビット	変更の必要が あるか否かを 示すインデックス ビット	適用される マスキングを 示すビット	階層値を基 づく必要と されるビット	組み合わせ の順序を 示すビット	合計	削減率
377	256 (2048)	2	7	16	1684	337	2046	0.10%
377	256 (2048)	2	7	16	1684	337	2039	0.44%
350	256 (2048)	2	7	16	1684	290	1999	2.39%
350	256 (2048)	2	7	16	1684	290	1992	2.73%
320	256 (2048)	2	7	16	1684	227	1936	5.47%
320	256 (2048)	2	7	16	1684	227	1929	5.81%
300	256 (2048)	2	7	16	1684	178	1886	7.91%
300	256 (2048)	2	7	16	1684	178	1879	8.25%
152	(892)	2	7	16	717	93	828	7.17%
77	(384)	2	7	16	296	48	362	5.73%

FIGURE 5

【 図 6 】

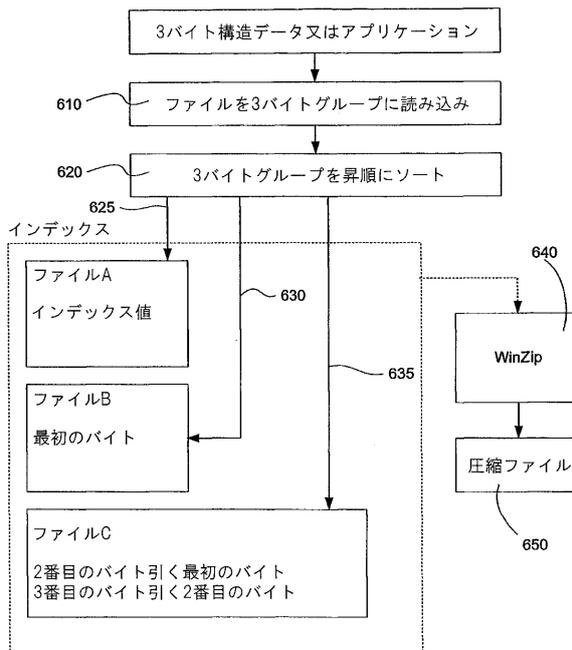


FIGURE 6

【 図 7 】

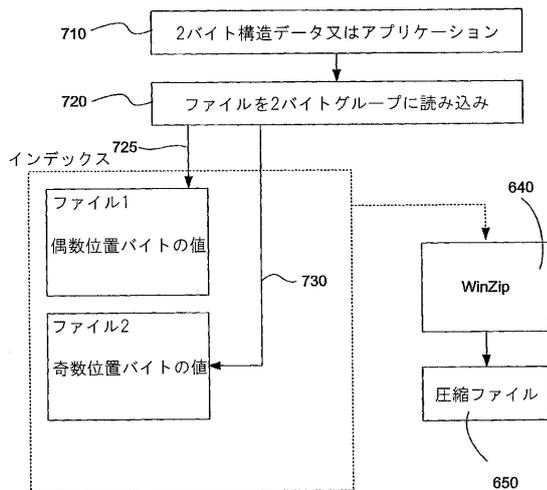


FIGURE 7

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No.

PCT/AU2004/001406

A. CLASSIFICATION OF SUBJECT MATTER		
Int. Cl. ⁷ : H03M 7/30		
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)		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) WPAT, ESP@CE, USPTO, JPO and INTERNET: Keywords (data, compression, unique, byte, frequency, index) and similar terms.		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	JP 2002-325252 A (NORITZ CORP) 8 November 2002 See whole document	1 - 13
A	Goebel G., DATA COMPRESSION [1.0] INTRODUCTION / LOSSLESS DATA COMPRESSION [online], 1 May 2003, pages 1 - 15, [retrieved 6 December 2004], Retrieved from The Internet:< http://www.vectorsite.net/tdomp1.html	1 - 13
A	JP 06-319047 A (SEIKO EPSON CORP) 15 November 1994 See whole document	1 - 13
(English translations for both JP documents sourced from the JPO website for searching the Patent Abstracts of Japan, 6 December 2004, URL:> http://www19.ipdl.ncipi.go.jp/PA1/cgi-bin/PAINDEX)		
<input type="checkbox"/> Further documents are listed in the continuation of Box C <input checked="" type="checkbox"/> See patent family annex		
* Special categories of cited documents:		
"A"	document defining the general state of the art which is not considered to be of particular relevance	"T"
"E"	earlier application or patent but published on or after the international filing date	"X"
"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)	"Y"
"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 6 December 2004		Date of mailing of the international search report 21 DEC 2004
Name and mailing address of the ISA/AU AUSTRALIAN PATENT OFFICE PO BOX 200, WODEN ACT 2606, AUSTRALIA E-mail address: pct@ipaaustralia.gov.au Facsimile No. (02) 6285 3929		Authorized officer BEN TUOHY Telephone No : (02) 6283 7918

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No.

PCT/AU2004/001406

This Annex lists the known "A" publication level patent family members relating to the patent documents cited in the above-mentioned international search report. The Australian Patent Office is in no way liable for these particulars which are merely given for the purpose of information.

Patent Document Cited in Search Report	Patent Family Member
JP 2002325252	
JP 6319047	

Due to data integration issues this family listing may not include 10 digit Australian applications filed since May 2001.

END OF ANNEX

フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), EP(AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PL, PT, RO, SE, SI, SK, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NA, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW

(74)代理人 100108383

弁理士 下道 晶久

(72)発明者 パーカー, ブルース

オーストラリア国, ニューサウスウエールズ 2079, マウント コラー, ベリル アベニュー
23

Fターム(参考) 5B082 GA01

5J064 AA02 BA18 BC01 BC14 BC22 BD02 BD03 BD04