

- 【書類名】 特許願
- 【整理番号】 CP00744
- 【あて先】 特許庁長官殿
- 【国際特許分類】 G06F 3/06
- 【発明者】
- 【住所又は居所】 大韓民国 790-751 ギョンサンブクト ポハンシ ナム
グ ジゴクドン ギョス・アパート 9-1503
- 【氏名】 パク チャンイク
- 【発明者】
- 【住所又は居所】 大韓民国 702-788 デグ ブクグ テジョンドン ヒョ
ップファ・アパート 103-1506
- 【氏名】 パク セジン
- 【特許出願人】
- 【識別番号】 511258411
- 【氏名又は名称】 ポハン工科大学校産学協力団
- 【氏名又は名称原語表記】 POSTECH Academy - Industry
Foundation
- 【代理人】
- 【識別番号】 100121728
- 【弁理士】
- 【氏名又は名称】 井関 勝守
- 【電話番号】 06-6136-6951
- 【ファクシミリ番号】 06-6136-6952
- 【連絡先】 担当
- 【パリ条約による優先権等の主張】
- 【国名】 大韓民国
- 【出願日】 2013年 3月 7日
- 【出願番号】 10-2013-0024356
- 【手数料の表示】
- 【振替番号】 00016713
- 【納付金額】 15,000円
- 【提出物件の目録】
- 【物件名】 明細書 1
- 【物件名】 特許請求の範囲 1
- 【物件名】 要約書 1
- 【物件名】 図面 1
- 【包括委任状番号】 1112740

【書類名】 明細書

【発明の名称】 データ重複除去方法及び装置

【技術分野】

【0001】

本発明は、データ重複除去技術に関し、より詳細には、低い入出力レイテンシ(l a t e n c y)を提供するためのデータ重複除去方法及び装置に関する。

【背景技術】

【0002】

データ重複除去技術は、データ保存装置内に重複されたデータを除去してより多い保存空間を確保するための技術を意味する。現在、多い企業、公共機関などでデータの安全な保管のためにデータのバックアップ(b a c k u p)を周期的に実行している。バックアップデータは、その特性上多い重複的な要素を有し、これによって、バックアップデータの保存空間の効率を向上させるためにデータ重複除去技術が使われている。このようなデータ重複除去技術は、バックアップデータ保存装置の特性上、低い入出力レイテンシ(l a t e n c y)を必要としないので、重複除去率を高める技術を中心として発展している。

【0003】

しかし、このようなデータ重複除去技術は、重複除去のための複雑なアルゴリズム(a l g o r i t h m)を基礎で実行されるため、ノート型P C(n o t e b o o k)、スマートホン(s m a r t p h o n e)、タブレット型(t a b l e t)P Cなどのような携帯用端末に適用しにくい問題点がある。すなわち、このような携帯用端末にデータ重複除去技術を適用する場合、順次的に保存されたデータの物理的手順が変わるようになるのでデータの入出力速度が深刻に遅くなる問題点がある。

【発明の概要】

【発明が解決しようとする課題】

【0004】

したがって、前記のような従来の諸問題点を解消するために提案されたものであって、本発明の目的は、データの入出力特性に応じて適応的に決定された重複除去率を基盤としてデータの重複を除去するためのデータ重複除去方法を提供することにある。

【0005】

本発明の他の目的は、データの入出力特性に応じて適応的に決定された重複除去率を基盤としてデータの重複を除去するためのデータ重複除去装置を提供することにある。

【課題を解決するための手段】

【0006】

前記目的を達成するための本発明の一実施例によるデータの重複除去方法は、データの入力要請または出力要請を基盤として前記データに対する接近特性を獲得する段階と、前記接近特性を基盤として前記データの重複除去単位を決定する段階と、前記重複除去単位を基盤として前記データに対する重複除去を実行する段階と、を含む。

【0007】

ここで、前記接近特性は、前記データに対する接近時間、前記データに対する変更時間、前記データに対する順次的接近回数及び前記データに対する任意的接近回数の中で少なくとも一つを含むことができる。

【0008】

ここで、前記接近特性を獲得する段階は、前記データの入力要請を受信した場合、前記データの入力要請に対する時間情報を基盤として前記接近時間及び前記変更時間を獲得する段階と、前記データの入力要請の連続性を基盤として前記順次的接近回数または前記任意的接近回数を獲得する段階と、を含むことができる。

【0009】

ここで、前記接近特性を獲得する段階は、前記データの出力要請を受信した場合、前記データの出力要請に対する時間情報を基盤として前記接近時間を獲得する段階と、前記データの出力要請の連続性を基盤として前記順次的接近回数または前記任意的接近回数を獲

得する段階と、を含むことができる。

【0010】

ここで、前記重複除去単位を決定する段階は、前記データに対する現在接近時間と前記データに対する以前変更時間に対する第1の差を算出する段階と、前記第1の差が予め定義された第1の時間以下の場合、前記重複除去単位をデータに対する重複除去可能性が一番低い第4の重複除去単位で決定する段階と、前記第1の差が予め定義された第1の時間を超過する場合、前記データに対する現在接近時間と前記データに対する以前接近時間に対する第2の差を算出する段階と、前記第2の差が予め定義された第2の時間を超過する場合、前記重複除去単位をデータに対する重複除去可能性が一番高い第1の重複除去単位で決定する段階と、前記第2の差が予め定義された第2の時間以下の場合、前記データに対する任意的接近回数が順次的接近回数以上であれば、前記重複除去単位を第1の重複除去単位より重複除去可能性が低い第2の重複除去単位で決定する段階と、前記第2の差が予め定義された第2の時間以下の場合、前記データに対する任意的接近回数が順次的接近回数未満であれば、前記重複除去単位を第2の重複除去単位より重複除去可能性が低い第3の重複除去単位で決定する段階と、を含むことができる。

【0011】

ここで、前記第4の重複除去単位は、前記データに対する重複除去を実行しないことを意味する。

【0012】

ここで、前記データに対する重複除去を実行する段階は、前記重複除去単位を基盤として前記データに対する少なくとも一つのデータブロックを生成する段階と、前記データブロックに対する固有の識別子を生成する段階と、前記固有の識別子がインデックステーブル内に存在するか否かを判断する段階と、前記固有の識別子が前記インデックステーブル内に存在する場合、前記固有の識別子に対応したデータブロックを除去する段階と、前記固有の識別子が前記インデックステーブル内に存在しない場合、前記固有の識別子と前記固有の識別子に対応したデータブロックを保存する段階と、を含むことができる。

【0013】

ここで、前記固有の識別子を生成する段階は、ハッシュアルゴリズムを使用して前記データブロックに対する固有の識別子を生成することができる。

【0014】

ここで、前記重複除去単位は、データに対する重複除去可能性を基盤として少なくとも一つの重複除去単位で分類されることができる。

【0015】

本発明の他の目的を達成するための本発明の一実施例によるデータ重複除去装置は、データの入力要請または出力要請を基盤として前記データに対する接近特性を獲得し、前記接近特性を基盤として前記データの重複除去単位を決定し、前記重複除去単位を基盤として前記データに対する重複除去を実行する処理部と、前記処理部で処理される情報及び処理された情報を保存する保存部と、を含む。

【0016】

ここで、前記接近特性は、前記データに対する接近時間、前記データに対する変更時間、前記データに対する順次的接近回数及び前記データに対する任意的接近回数の中で少なくとも一つを含むことができる。

【0017】

ここで、前記処理部は、前記データの入力要請を受信した場合、前記データの入力要請に対する時間情報を基盤として前記接近時間及び前記変更時間を獲得し、前記データの入力要請の連続性を基盤として前記順次的接近回数または前記任意的接近回数を獲得することができる。

【0018】

ここで、前記処理部は、前記データの出力要請を受信した場合、前記データの出力要請に対する時間情報を基盤として前記接近時間を獲得し、前記データの出力要請の連続性を

【0019】

ここで、前記処理部は、前記重複除去単位を決定する場合、前記データに対する現在接近時間と前記データに対する以前変更時間に対する第1の差を算出し、前記第1の差が予め定義された第1の時間以下の場合、前記重複除去単位をデータに対する重複除去可能性が一番低い第4の重複除去単位で決定し、前記第1の差が予め定義された第1の時間を超過する場合、前記データに対する現在接近時間と前記データに対する以前接近時間に対する第2の差を算出し、前記第2の差が予め定義された第2の時間を超過する場合、前記重複除去単位をデータに対する重複除去可能性が一番高い第1の重複除去単位で決定し、前記第2の差が予め定義された第2の時間以下の場合、前記データに対する任意的接近回数が順次的接近回数以上であれば、前記重複除去単位を第1の重複除去単位より重複除去可能性が低い第2の重複除去単位で決定し、前記第2の差が予め定義された第2の時間以下の場合、前記データに対する任意的接近回数が順次的接近回数未満であれば、前記重複除去単位を第2の重複除去単位より重複除去可能性が低い第3の重複除去単位で決定することができる。

【0020】

ここで、前記第4の重複除去単位は、前記データに対する重複除去を実行しないことを意味する。

【0021】

ここで、前記処理部は、前記データに対する重複除去を実行する場合、前記重複除去単位を基盤として前記データに対する少なくとも一つのデータブロックを生成し、前記データブロックに対する固有の識別子を生成し、前記固有の識別子がインデックステーブル内に存在するか否かを判断し、前記固有の識別子が前記インデックステーブル内に存在する場合、前記固有の識別子に対応したデータブロックを除去し、前記固有の識別子が前記インデックステーブル内に存在しない場合、前記固有の識別子と前記固有の識別子に対応したデータブロックを保存することができる。

【0022】

ここで、前記処理部は、前記固有の識別子を生成する場合、ハッシュアルゴリズムを使用して前記データブロックに対する固有の識別子を生成することができる。

【0023】

ここで、前記重複除去単位は、データに対する重複除去可能性を基盤として少なくとも一つの重複除去単位で分類されることができる。

【発明の効果】

【0024】

本発明によれば、データの入出力特性を基盤として重複除去率(すなわち、チャンク(chunk)のサイズ)を適応的に決定することができ、適応的に決定された重複除去率を基盤として重複を除去することができるので、低い入出力レイテンシ(latency)を提供することができる。

【図面の簡単な説明】

【0025】

【図1】 図1は、本発明の一実施例によるデータ重複除去方法を示したフローチャートである。

【図2】 図2は、データに対する接近特性を示した表である。

【図3】 図3は、本発明の一実施例によるデータ重複除去方法において接近特性獲得段階を示した図である。

【図4】 図4は、本発明の一実施例によるデータ重複除去方法において重複除去単位決定段階を示したフローチャートである。

【図5】 図5は、重複除去単位に対する特性を示した表である。

【図6】 図6は、本発明の一実施例によるデータ重複除去方法において重複除去実行段階を示したフローチャートである。

【図7】 図7は、データブロックに対する固有の識別子を生成する過程を示した概念図である。

【図8】 図8は、本発明の一実施例によるデータ重複除去装置を示したブロック図である。

【発明を実施するための形態】

【0026】

本発明は、多様に変更可能であり、さまざまな実施形態を有することができる。ここでは、特定の実施形態を図面に例示して詳細に説明する。

【0027】

しかし、これは本発明の好ましい実施態様に過ぎず、本発明の実施の範囲を限定するものではなく、本発明の明細書及び図面内容に基づいてなされた均等な変更および付加は、いずれも本発明の特許請求の範囲内に含まれるものとする。

【0028】

第1、第2などの用語は、多様な構成要素を説明するために使用することができるが、前記構成要素は前記用語により限定されるものではない。前記用語は、一つの構成要素を他の構成要素から区別するための目的のみで使用される。例えば、本発明の権利範囲を脱しない範囲で、第1の構成要素は第2の構成要素と命名することができ、類似に第2の構成要素も第1の構成要素と命名することができる。及び/またはこの用語は、複数の関連された記載された項目の組み合わせまたは複数の関連された記載された項目の中のいずれか一項目を含む。

【0029】

いかなる構成要素が他の構成要素に「連結されて」いるかあるいは「接続されて」いるとの用語は、ある構成要素が他の構成要素に直接的に連結されるかあるいは接続されることもできるが、中間に他の構成要素が介在することもできることを意味する。一方に、ある構成要素が他の構成要素に「直接連結されて」いるかあるいは「直接接続されて」いるとの用語は、中間に他の構成要素が存在しないことを意味する。

【0030】

本明細書で使用した用語は、但し、特定の実施形態を説明するために使用されたもので、本発明はこれに限定されるものではない。単数の表現は、文脈上明白に相違に記載しない限り複数の表現を含む。本出願において、「含む」または「有する」などの用語は、明細書上に記載された特徴、数字、ステップ、動作、構成要素、部品またはこれらを組み合わせたものが存在することを指定するのであって、一つまたはその以上の他の特徴や数字、ステップ、動作、構成要素、部品またはこれらを組み合わせたものなどの存在または付加可能性をあらかじめ排除することではない。

【0031】

特定しない限り、技術的や科学的な用語を含んでここで使用されるすべての用語は、本発明が属する技術分野で通常知識を有する者により一般的に理解される意味と同一な意味を有する。一般的に使用される辞典に定義された用語は、関連技術の文脈上有する意味と一致することと解でき、本出願において明白に定義しない限り、理想的や過度に形式的な意味を有することで解釈できない。

【0032】

以下、添付図面を参照して実施例を詳しく説明するが、図面符号に関係なく同一であるか対応する構成要素には同一な参照番号を付与し、その重複説明は省略する。

【0033】

図1は、本発明の一実施例によるデータ重複除去方法を示したフローチャートである。

【0034】

図1を参照すれば、本発明の一実施例によるデータ重複除去方法は、データの入力要請または出力要請を基盤としてデータに対する接近特性を獲得する段階(ステップS100)と、接近特性を基盤としてデータの重複除去単位を決定する段階(ステップS200)と、重複除去単位を基盤としてデータに対する重複除去を実行する段階(ステップS300)と

、を含む。

【0035】

ここで、図1に示した各段階は、図8に示したデータ重複除去装置で実行することができ、データ重複除去装置の具体的な構成とその機能に対しては後述する。

【0036】

図2は、データに対する接近特性を示した表である。

【0037】

図2を参照すれば、データに対する接近特性は、データに対する接近時間(a_time)、データに対する変更時間(m_time)、データに対する順次的接近回数(seqCount)、データに対する任意的接近回数(randCount)などを含むことができる。ここで、データに対する接近特性は、前記情報などに限定されないで、データの入力または出力に対する特性を示すことができる情報であれば、いずれも含むことができる。

【0038】

データに対する接近時間(a_time)は、データの入力要請(すなわち、データの書き込み要請)またはデータの出力要請(すなわち、データの読み取り要請)を受信した時間を意味する。データの入力要請を受けた場合、データ重複除去装置は、入力要請時間を該当データに対する接近時間で獲得することができ、獲得した接近時間を保存(すなわち、現在データに対するa_timeフィールドを現在時間で記録)することができる。一方、既に保存された接近時間がある場合、データ重複除去装置は、既に保存された接近時間を最近獲得した接近時間に更新することができる。

【0039】

一方、データの出力要請を受けた場合、データ重複除去装置は、出力要請時間を該当データに対する接近時間で獲得することができ、獲得した接近時間を保存(すなわち、現在データに対するa_timeフィールドを現在時間で記録)することができる。一方、既に保存された接近時間がある場合、データ重複除去装置は、既に保存された接近時間を最近獲得した接近時間に更新することができる。

【0040】

データに対する変更時間(m_time)は、データの入力要請を受信した時間を意味する。データの入力要請を受けた場合、データ重複除去装置は、入力要請時間を該当データに対する変更時間で獲得することができ、獲得した変更時間を保存(すなわち、現在データに対するm_timeフィールドを現在時間で記録)できる。一方、既に保存された変更時間がある場合、データ重複除去装置は、既に保存された変更時間を最近獲得した変更時間に更新することができる。

【0041】

データに対する順次的接近回数(seqCount)は、現在のデータ要請と以前のデータ要請が連続された(すなわち、現在のデータ要請と以前のデータ要請の物理的または論理的ブロック番号が連続的な場合、またはその要請が連続的な形態を有する場合)回数を意味し、データに対する任意的接近回数(randCount)は、現在のデータ要請と以前データ要請が連続されない回数を意味する。ここで、順次的接近回数と任意的接近回数は、現在まで累積された回数を意味する。

【0042】

現在のデータ要請と以前のデータ要請が連続される場合、データ重複除去装置は、現在データに対するseqCountフィールドの値を1増加させる。一方、現在のデータ要請と以前のデータ要請が連続されない場合、データ重複除去装置は、現在データに対するrandCountフィールドの値を1増加させる。

【0043】

図3は、本発明の一実施例によるデータ重複除去方法において接近特性獲得段階を示したフローチャートである。

【0044】

図3を参照すれば、データに対する接近特性を獲得する段階(ステップS100)は、受

受信された要請がデータの入力要請に該当するかデータの出力要請に該当するかを判断する段階(ステップS 1 1 0)と、データの入力要請を受信した場合、データの入力要請に対する時間情報を基盤として接近時間及び変更時間を獲得する段階(ステップS 1 2 0)と、データの入力要請の連続性を基盤として順次的接近回数または任意的接近回数を獲得する段階(ステップS 1 3 0)と、を含む。

【0045】

また、データに対する接近特性を獲得する段階(ステップS 1 0 0)は、データの出力要請を受信した場合、データの出力要請に対する時間情報を基盤として接近時間を獲得する段階(ステップS 1 4 0)と、データの出力要請の連続性を基盤として順次接近回数または任意的接近回数を獲得する段階(ステップS 1 5 0)と、を含む。

【0046】

ステップS 1 1 0において、データ重複除去装置は、受信された要請がデータの入力要請に該当するかデータの出力要請に該当するかを判断する。受信された要請がデータの入力要請に該当する場合、データ重複除去装置は、次の段階として、ステップS 1 2 0、ステップS 1 3 0を実行する。一方、受信された要請がデータの出力要請に該当する場合、データ重複除去装置は、次の段階として、ステップS 1 4 0、ステップS 1 5 0を実行する。

【0047】

ステップS 1 2 0において、データ重複除去装置は、データの入力要請に対する時間情報を基盤として接近時間及び変更時間を獲得する。すなわち、データ重複除去装置は、データの入力要請を受信した時間を接近時間及び変更時間で獲得することができ、獲得した接近時間及び獲得した変更時間をデータベース(d a t a b a s e)に保存する。この時、既に保存された接近時間が存在する場合、データ重複除去装置は、既に保存された接近時間を前記獲得した接近時間に更新することができ、既に保存された変更時間が存在する場合、データ重複除去装置は、既に保存された変更時間を前記獲得した変更時間に更新することができる。

【0048】

ステップS 1 3 0において、データ重複除去装置は、データの入力要請と該当データに対する以前要請の連続性を基盤として順次的接近回数または任意的接近回数を獲得することができる。すなわち、データ重複除去装置は、データの入力要請と該当データに対する以前要請が連続される場合、順次的接近回数を1増加させることができ、データの入力要請と該当データに対する以前要請が連続されない場合、任意的接近回数を1増加させることができる。

【0049】

例えば、データベースに保存された接近特性の中で順次的接近回数が7であり、任意的接近回数が5である場合、データ重複除去装置は、データの入力要請と該当データに対する以前要請が連続される場合、順次接近回数を8に更新し、データの入力要請と該当データに対する以前要請が連続されない場合、任意的接近回数を6に更新する。

【0050】

ここで、ステップS 1 2 0を先に実行した後、ステップS 1 3 0を実行することで説明したが、ステップS 1 2 0とステップS 1 3 0の実行順序は、これに限定されない。すなわち、ステップS 1 3 0は、ステップS 1 2 0と同時に実行するか、またはステップS 1 2 0より先に実行してもよい。

【0051】

ステップS 1 4 0において、データ重複除去装置は、データの出力要請に対する時間情報を基盤として接近時間を獲得することができる。すなわち、データ重複除去装置は、データの出力要請を受信した時間を接近時間で獲得することができ、獲得した接近時間をデータベースに保存することができる。この時、既に保存された接近時間が存在する場合、データ重複除去装置は、既に保存された接近時間を前記獲得した接近時間に更新することができる。

【0052】

ステップS150において、データ重複除去装置は、データの出力要請と該当データに対する以前要請の連続性を基盤として順次的接近回数または任意的接近回数を獲得することができる。すなわち、データ重複除去装置は、データの出力要請と該当データに対する以前要請が連続される場合、順次接近回数を1増加させることができ、データの出力要請と該当データに対する以前要請が連続されない場合、任意的接近回数を1増加させることができる。

【0053】

例えば、データベースに保存された接近特性の中で順次的接近回数が7であり、任意的接近回数が5である場合、データ重複除去装置は、データの出力要請と該当データに対する以前要請が連続される場合、順次接近回数を8に更新し、データの出力要請と該当データに対する以前要請が連続されない場合、任意的接近回数を6に更新する。

【0054】

ここで、ステップS140を先に実行した後、ステップS150を実行することで説明したが、ステップS140とステップS150の実行順序は、これに限定されない。すなわち、ステップS150は、ステップS140と同時に実行するか、またはステップS140より先に実行してもよい。

【0055】

図4は、本発明の一実施例によるデータ重複除去方法において重複除去単位決定段階を示したフローチャートである。

【0056】

図4を参照すれば、重複除去単位を決定する段階(ステップS200)は、データに対する現在接近時間とデータに対する以前変更時間に対する第1の差を算出する段階(ステップS210)と、第1の差が予め定義された第1の時間を超過するか否かを判断する段階(ステップS220)と、第1の差が予め定義された第1の時間以下の場合、重複除去単位をデータに対する重複除去可能性が一番低い第4の重複除去単位で決定する段階(ステップS230)と、第1の差が予め定義された第1の時間を超過する場合、データに対する現在接近時間とデータに対する以前接近時間に対する第2の差を算出する段階(ステップS240)と、第2の差が予め定義された第2の時間以下であるか否かを判断する段階(ステップS250)と、第2の差が予め定義された第2の時間を超過する場合、重複除去単位をデータに対する重複除去可能性が一番高い第1の重複除去単位で決定する段階(ステップS260)と、第2の差が予め定義された第2の時間以下の場合、データに対する任意的接近回数がデータに対する順次的接近回数以下であるか否かを判断する段階(ステップS270)と、第2の差が予め定義された第2の時間以下の場合、データに対する任意的接近回数が順次的接近回数以上であれば、重複除去単位を第1の重複除去単位より重複除去可能性が低い第2の重複除去単位で決定する段階(ステップS280)と、第2の差が予め定義された第2の時間以下の場合、データに対する任意的接近回数が順次的接近回数未満であれば、重複除去単位を第2の重複除去単位より重複除去可能性が低い第3の重複除去単位で決定する段階(ステップS290)と、を含むことができる。

【0057】

以下、ステップS200で決定される重複除去単位に対して図5を参照して詳細に説明する。

【0058】

図5は、重複除去単位に対する特性を示した表である。

【0059】

図5を参照すれば、重複除去単位は、第1の重複除去単位、第2の重複除去単位、第3の重複除去単位、第4の重複除去単位で分類することができる。第1の重複除去単位は、よく使われないデータに適用することができ、よく使われないデータは、データに対する現在接近時間と以前接近時間の差が予め定義された臨界値より大きいデータを意味する。第1の重複除去単位は、全ての重複除去単位の中で一番小さいチャンク(chunk)を使

用する。これによって、全ての重複除去単位の中で一番高いデータ重複除去率を提供することができる。すなわち、第1の重複除去単位は、低いレイテンシ(latency)より高いデータ重複除去率を提供するために使われることができる。

【0060】

第2の重複除去単位は、順次的接近より任意的接近がよく発生するデータに適用することができる。第2の重複除去単位は、第1の重複除去単位より大きくて第3の重複除去単位より小さいサイズのチャンクを使用する。これによって、全ての重複除去単位の中で相対的に高いデータ重複除去率(すなわち、第1の重複除去単位より低く第3の重複除去単位より高い重複除去率)を提供することができる。すなわち、任意的接近がよく発生するデータの場合、物理的な歪みが発生しても任意的接近性能に問題が発生しないので、第2の重複除去単位は、高い重複除去率を提供するために相対的に小さいサイズのチャンクを使用することができる。

【0061】

第3の重複除去単位は、任意的接近より順次的接近がよく発生するデータに適用することができる。第3の重複除去単位は、第2の重複除去単位より大きいサイズのチャンクを使用する。これによって、全ての重複除去単位の中で相対的に小さい重複除去率(すなわち、第2の重複除去単位より低い重複除去率)を提供することができる。すなわち、順次的な接近がよく発生するデータに対して低い入出力レイテンシを提供するため、第3の重複除去単位は、相対的に大きいサイズのチャンクを使用することができる。

【0062】

第4の重複除去単位は、入力がよく発生するデータに適用することができる。第4の重複除去単位は、第3の重複除去単位より大きいサイズのチャンクを使用する。これによって、全ての重複除去単位の中で一番小さい重複除去率(すなわち、第3の重複除去単位より低い重複除去率)を提供することができる。一方、第4の重複除去単位は、データに対する重複除去を実行しないことを意味する。すなわち、データの重複除去は、出力を主とするデータに有利なので、入力を主とするデータの場合、重複除去を実行しないこともある。

【0063】

ここで、重複除去単位の分類は、上述の説明に限定されないで、多様に構成することができる。例えば、重複除去単位を3個分類または5個分類で構成することができる。重複除去単位が3個分類で構成される場合、第1の重複除去単位は、一番高い重複除去率を提供することができ、第2の重複除去単位は、第1の重複除去単位より低い重複除去率を提供することができ、第3の重複除去単位は、第2の重複除去単位より低い重複除去率(すなわち、一番低い重複除去率)を提供することができる。

【0064】

さらに、図4を参照すれば、ステップS210において、データ重複除去装置は、データに対する現在接近時間とデータに対する以前変更時間に対する第1の差を算出することができる。すなわち、データ重複除去装置は、データの入力要請または出力要請から獲得した現在接近時間と、同一なデータに対する以前変更時間(すなわち、データの以前入力要請から獲得した変更時間)の差である第1の差を算出することができる。ここで、第1の差は、該当データがどれくらいの頻度で変更されるか(すなわち、データの入力要請がどれくらいの頻度で発生するか)を意味する。

【0065】

ステップS220において、データ重複除去装置は、第1の差が予め定義された第1の時間を超過するか否かを判断することができる。ここで、予め定義された第1の時間は、データ重複除去を実行するために基準になる時間を意味し、ユーザの設定によって異なる値を有することができる。例えば、予め定義された第1の時間は、1時間、2時間、3時間などで設定することができる。第1の差が予め定義された第1の時間以下の場合、データ重複除去装置は、次の段階として、ステップS230を実行することができ、第1の差が予め定義された第1の時間を超過する場合、データ重複除去装置は、次の段階として、

【0066】

ステップS 2 3 0において、データ重複除去装置は、重複除去単位をデータに対する重複除去可能性が一番低い第4の重複除去単位で決定することができる。すなわち、第1の差が予め定義された第1の時間以下の場合、これは入力がよく発生するデータを意味するので、データ重複除去装置は、重複除去単位の中でデータ重複除去率が一番低い(または、データ重複除去を実行しない)第4の重複除去単位を選択することができる。

【0067】

ステップS 2 4 0において、データ重複除去装置は、データに対する現在接近時間とデータに対する以前接近時間に対する第2の差を算出することができる。すなわち、データ重複除去装置は、データの入力要請または出力要請から獲得した現在接近時間と、同一なデータに対する以前接近時間(すなわち、データの以前入力要請または以前出力要請から獲得した接近時間)の差である第2の差を算出することができる。ここで、第2の差は、該当データに対する接近がどれくらいの頻度で発生するかを意味する。

【0068】

ステップS 2 5 0において、データ重複除去装置は、第2の差が予め定義された第2の時間以下であるか否かを判断することができる。ここで、予め定義された第2の時間は、データに対する接近が発生する可能性が低いデータを区別するために基準になる時間を意味し、ユーザの設定によって異なる値を有することができる。例えば、予め定義された第2の時間は、1時間、2時間、3時間などで設定することができる。第2の差が予め定義された第2の時間を超過する場合、データ重複除去装置は、次の段階として、ステップS 2 6 0を実行することができ、第2の差が予め定義された第2の時間以下の場合、データ重複除去装置は、次の段階として、ステップS 2 7 0を実行することができる。

【0069】

ステップS 2 6 0において、データ重複除去装置は、重複除去単位をデータに対する重複除去可能性が一番高い第1の重複除去単位で決定することができる。すなわち、第2の差が予め定義された第2の時間を超過する場合、これはよく使われないデータ(すなわち、接近可能性が低いデータ)を意味するので、データ重複除去装置は、重複除去単位の中でデータ重複除去率が一番高い第1の重複除去単位を選択することができる。

【0070】

ステップS 2 7 0において、データ重複除去装置は、データに対する任意的接近回数がデータに対する順次的接近回数未満であるか否かを判断する。任意的接近回数が順次的接近回数以上の場合、データ重複除去装置は、次の段階として、ステップS 2 8 0を実行することができ、任意的接近回数が順次的接近回数未満の場合、データ重複除去装置は、次の段階として、ステップS 2 9 0を実行することができる。

【0071】

ステップS 2 8 0において、データ重複除去装置は、重複除去単位を第1の重複除去単位より重複除去可能性が低い第2の重複除去単位で決定することができる。すなわち、データに対する任意的接近回数が順次的接近回数以上の場合、これは任意的接近がよく発生するデータを意味するので、データ重複除去装置は、重複除去単位の中でデータ重複除去率が相対的に高い第2の重複除去単位を選択することができる。

【0072】

ステップS 2 9 0において、データ重複除去装置は、重複除去単位を第2の重複除去単位より重複除去可能性が低い第3の重複除去単位で決定することができる。すなわち、データに対する任意的接近回数が順次的接近回数未満の場合は、順次的接近がよく発生するデータを意味するので、データ重複除去装置は、重複除去単位の中でデータ重複除去率が相対的に低い第3の重複除去単位を選択することができる。

【0073】

図6は、本発明の一実施例によるデータ重複除去方法において重複除去実行段階を示したフローチャートであり、図7は、データブロックに対する固有の識別子を生成する過程

【0074】

以下、図6及び図7を参照して、重複除去を実行する段階(ステップS300)に対して詳細に説明する。

【0075】

本発明の一実施例によるデータ重複除去方法において重複除去を実行する段階(ステップS300)は、重複除去単位を基盤としてデータに対する少なくとも一つのデータブロックを生成する段階(ステップS310)と、データブロックに対する固有の識別子(identifier)を生成する段階(ステップS320)と、固有の識別子がインデックステーブル(index table)内に存在するか否かを判断する段階(ステップS330)と、固有の識別子がインデックステーブル内に存在しない場合、固有の識別子と固有の識別子に対応したデータブロックを保存する段階(ステップS340)と、固有の識別子がインデックステーブル内に存在する場合、固有の識別子に対応したデータブロックを除去する段階(ステップS350)と、を含む。

【0076】

ステップS310において、データ重複除去装置は、重複除去単位を基盤としてデータに対する少なくとも一つのデータブロックを生成することができ、重複除去単位は、チャンクのサイズを意味する。すなわち、データ重複除去装置は、第1の重複除去単位に対応したチャンクのサイズを基礎としてデータブロックを生成することができ、データ重複除去装置は、第2の重複除去単位に対応したチャンクのサイズを基礎としてデータブロックを生成することができ、データ重複除去装置は、第3の重複除去単位に対応したチャンクのサイズを基礎としてデータブロックを生成することができ、データ重複除去装置は、第4の重複除去単位に対応したチャンクのサイズを基礎としてデータブロックを生成することができる。一方、第4の重複除去単位がデータ重複除去を実行しないことを意味する場合、データ重複除去装置は、データに対するデータブロックを生成しないこともある。

【0077】

上述したステップS310を基礎として、データ重複除去装置は、データ30(図7)から複数のデータブロック31(図7)を生成することができる。

【0078】

ステップS320において、データ重複除去装置は、データブロックに対する固有の識別子を生成することができる。ここで、データ重複除去装置は、ハッシュアルゴリズム(例えば、SHA-1、SHA-2、SHA-3など)を使用してデータブロックに対する固有の識別子を生成することができる。データブロックに対する固有の識別子を生成する方法は、上述した方法に限定されなるものではなくて、公知の多様な方法を使用してデータブロックに対する固有の識別子を生成することができる。

【0079】

上述のステップS320を基礎として、データ重複除去装置は、各々のデータブロック31(図7)から固有の識別子32(図7)を生成することができる。

【0080】

ステップS330において、データ重複除去装置は、固有の識別子がインデックステーブル内に存在するか否かを判断する。インデックステーブルは、固有の識別子と固有の識別子に対応したデータブロックを含むことができる。ここで、固有の識別子がインデックステーブル内に存在する場合、これは固有の識別子に対応したデータブロックが既に保存されていることを示し、固有の識別子がインデックステーブル内に存在しない場合、これは固有の識別子に対応したデータブロックが保存されていないことを示す。固有の識別子がインデックステーブル内に存在しない場合、データ重複除去装置は、次の段階として、ステップS340を実行することができ、固有の識別子がインデックステーブル内に存在する場合、データ重複除去装置は、次の段階として、ステップS350を実行することができる。

【0081】

ステップS 3 4 0において、データ重複除去装置は、固有の識別子と固有の識別子に対応したデータブロックを保存することができる。すなわち、固有の識別子に対応したデータブロックが保存されていない状態であるので、データ重複除去装置は、重複除去を実行しないで固有の識別子とデータブロックをデータベース(または、インデックステーブル)に保存する。

【0082】

ステップS 3 5 0において、データ重複除去装置は、固有の識別子に対応したデータブロックを除去することができる。すなわち、固有の識別子に対応したデータブロックが既に保存されている状態であるので、データ重複除去装置は、重複除去(すなわち、固有の識別子に対応したデータブロック削除)を実行することができる。

【0083】

本発明の一実施例によるデータ重複除去方法は、多様なコンピュータ手段を通じて実行できるプログラム命令の形態で具現してコンピュータ判読可能媒体に記録することができる。コンピュータ判読可能媒体は、プログラム命令、データファイル、データ構造などを単独でまたは組み合わせて含むことができる。コンピュータ判読可能媒体に記録されるプログラム命令は、本発明のために特別に設計されて構成されたものであるか、あるいはコンピュータソフトウェア分野の当業者に公知されて使用可能なものであってもよい。

【0084】

コンピュータ判読可能媒体の例には、ROM(Read Only Memory)、RAM、フラッシュメモリー(flash memory)などのようにプログラム命令を保存して実行するように特別に構成されたハードウェア装置が含まれる。プログラム命令の例には、コンパイラ(compiler)により作られるもののような機械語コードだけではなくインタープリター(interpreter)などを使用してコンピュータにより実行できる高級言語コードを含む。上述したハードウェア装置は、本発明の動作を実行するために少なくとも一つのソフトウェアモジュールで作動するように構成されることができ、その逆も同様である。

【0085】

図8は、本発明の一実施例によるデータ重複除去装置を示したブロック図である。

【0086】

図8を参照すれば、本発明の一実施例によるデータ重複除去装置は、処理部10及び保存部20を含む。

【0087】

処理部10は、データの入力要請または出力要請を基盤としてデータに対する接近特性を獲得することができ、接近特性を基盤としてデータの重複除去単位を決定することができ、重複除去単位を基盤としてデータに対する重複除去を実行することができる。

【0088】

ここで、処理部10は、論理的構成である接近特性獲得部11と、重複除去単位決定部12と、重複除去実行部13と、インデックステーブル管理部14と、を含むことができる。一方、処理部10は、物理的構成であるプロセッサ(processor)及びメモリー(memory)を含む。プロセッサは、汎用のプロセッサ(例えば、CPU(Central Processing Unit)及び/またはGPU(Graphics Processing Unit)など)またはデータ重複除去方法の実行のための専用のプロセッサを意味する。メモリーには、データ重複除去方法の実行のためのプログラムコード(program code)が保存される。すなわち、プロセッサは、メモリーに保存されたプログラムコードを読み出すことができ、読み出されたプログラムコードを基盤としてデータ重複除去方法の各段階を実行することができる。

【0089】

ここで、データに対する接近特性は、データに対する接近時間(a_time、図2参照)、データに対する変更時間(m_time、図2参照)、データに対する順次的接近回数(seqCount、図2参照)、データに対する任意的接近回数(randCount、図

2参照)などを含むことができる。ここで、データに対する接近特性は、前記情報などに限定されないで、データの入力または出力に対する特性を示すことができる情報であれば、いずれも含むことができる。

【0090】

データに対する接近時間(a_time)は、データの入力要請(すなわち、データの書き込み要請)またはデータの出力要請(すなわち、データの読み取り要請)を受信した時間を意味する。データに対する変更時間(m_time)は、データの入力要請を受信した時間を意味する。データに対する順次的接近回数(seqCount)は現在のデータ要請と以前のデータ要請が連続される(すなわち、現在のデータ要請と以前のデータ要請が同一である)回数を意味し、データに対する任意的接近回数(randCount)は、現在のデータ要請と以前のデータ要請が連続されない(すなわち、現在のデータ要請と以前のデータ要請が同一ではない)回数を意味する。

【0091】

接近特性を獲得する場合、処理部10は、データの入力要請に対する時間情報を基盤として接近時間及び変更時間を獲得することができ、データの入力要請の連続性を基盤として順次的接近回数または任意的接近回数を獲得することができる。また、処理部10は、データの出力要請に対する時間情報を基盤として接近時間を獲得することができ、データの出力要請の連続性を基盤として順次的接近回数または任意的接近回数を獲得することができる。ここで、接近特性を獲得する過程は、処理部10内の接近特性獲得部11で実行することができる。

【0092】

具体的に、処理部10は、受信された要請がデータの入力要請に該当するかデータの出力要請に該当するかを判断することができる。受信された要請がデータの入力要請に該当する場合、処理部10は、データの入力要請を受信した時間を接近時間及び変更時間で獲得することができ、獲得した接近時間及び獲得した変更時間を保存部20に保存することができる。さらに、処理部10は、データの入力要請と該当データに対する以前要請が連続される場合、順次的接近回数を1増加させることができ、データの入力要請と該当データに対する以前要請が連続されない場合、任意的接近回数を1増加させることができる。

【0093】

一方、受信された要請がデータの出力要請に該当する場合、処理部10は、データの出力要請を受信した時間を接近時間で獲得することができ、獲得した接近時間をデータベースに保存することができる。さらに、処理部10は、データの出力要請と該当データに対する以前要請が連続される場合、順次的接近回数を1増加させることができ、データの入力要請と該当データに対する以前要請が連続されない場合、任意的接近回数を1増加させることができる。

【0094】

重複除去単位を決定する場合、処理部10は、データに対する現在接近時間とデータに対する以前変更時間に対する第1の差を算出することができ、第1の差が予め定義された第1の時間以下の場合、重複除去単位をデータに対する重複除去可能性が一番低い第4の重複除去単位で決定することができる。一方、第1の差が予め定義された第1の時間を超過する場合、処理部10は、データに対する現在接近時間とデータに対する以前接近時間に対する第2の差を算出することができる。

【0095】

ここで、処理部10は、第2の差が予め定義された第2の時間を超過する場合、重複除去単位をデータに対する重複除去可能性が一番高い第1の重複除去単位で決定することができ、第2の差が予め定義された第2の時間以下の場合、データに対する任意的接近回数が順次的接近回数以上であるか否かを判断することができる。

【0096】

ここで、処理部10は、データに対する任意的接近回数が順次的接近回数以上の場合、重複除去単位を第1の重複除去単位より重複除去可能性が低い第2の重複除去単位で決定

することができ、データに対する任意的接近回数が順次的接近回数未満の場合、重複除去単位を第2の重複除去単位より重複除去可能性が低い第3の重複除去単位で決定することができる。

【0097】

上述した重複除去単位を決定する過程は、処理部10内の重複除去単位決定部12で実行することができる。

【0098】

ここで、重複除去単位は、第1の重複除去単位、第2の重複除去単位、第3の重複除去単位、第4の重複除去単位で分類することができる。第1の重複除去単位は、よく使われないデータに適用することができ、第2の重複除去単位は、順次的接近より任意的接近がよく発生するデータに適用することができ、第3の重複除去単位は、任意的接近より順次的接近がよく発生するデータに適用することができ、第4の重複除去単位は、入力がよく発生するデータに適用することができる。

【0099】

具体的に、処理部10は、データの入力要請または出力要請から獲得した現在接近時間と、同一なデータに対する以前変更時間(すなわち、データの以前入力要請から獲得した変更時間)の差である第1の差を算出することができる。ここで、第1の差は、該当データがどれくらいの頻度で変更されるか(すなわち、データの入力要請がどれくらいの頻度で発生するか)を意味する。

【0100】

処理部10は、第1の差が予め定義された第1の時間を超過するか否かを判断することができる。ここで、予め定義された第1の時間は、データ重複除去を実行するために基準になる時間を意味し、ユーザの設定によって異なる値を有することができる。

【0101】

第1の差が予め定義された第1の時間以下の場合、これは入力がよく発生するデータを意味するので、処理部10は、重複除去単位の中でデータ重複除去率が一番低い(または、データ重複除去を実行しない)第4の重複除去単位を選択することができる。

【0102】

一方、第1の差が予め定義された第1の時間を超過する場合、処理部10は、データの入力要請または出力要請から獲得した現在接近時間と、同一なデータに対する以前接近時間(すなわち、データの以前入力要請または以前出力要請から獲得した接近時間)の差である第2の差を算出することができる。ここで、第2の差は、該当データに対する接近がどれくらいの頻度で発生するかを意味する。

【0103】

処理部10は、第2の差が予め定義された第2の時間以下であるか否かを判断することができる。ここで、予め定義された第2の時間は、データに対する接近が発生する可能性が低いデータを区別するために基準になる時間を意味し、ユーザの設定によって異なる値を有することができる。

【0104】

第2の差が予め定義された第2の時間を超過する場合、これはよく使われないデータ(すなわち、接近可能性が低いデータ)を意味するので、処理部10は、重複除去単位の中でデータ重複除去率が一番高い第1の重複除去単位を選択することができる。一方、第2の差が予め定義された第2の時間以下の場合、処理部10は、データに対する任意的接近回数がデータに対する順次接近回数未満であるか否かを判断することができる。

【0105】

データに対する任意的接近回数が順次的接近回数以上の場合、これは任意的接近がよく発生するデータを意味するので、処理部10は、重複除去単位の中でデータ重複除去率が相対的に高い第2の重複除去単位を選択することができる。一方、データに対する任意的接近回数が順次的接近回数未満の場合は、順次的接近がよく発生するデータを意味するので、処理部10は、重複除去単位の中でデータ重複除去率が相対的に低い第3の重複除去

【0106】

データの重複除去を実行する場合、処理部10は、重複除去単位を基盤としてデータに対する少なくとも一つのデータブロックを生成することができ、データブロックに対する固有の識別子を生成することができ、固有の識別子がインデックステーブル内に存在するか否かを判断することができ、固有の識別子がインデックステーブル内に存在する場合、固有の識別子に対応したデータブロックを除去することができ、固有の識別子がインデックステーブル内に存在しない場合、固有の識別子と固有の識別子に対応したデータブロックを保存することができる。

【0107】

ここで、データの重複を除去する過程は、処理部10内の重複除去実行部13で実行することができ、インデックステーブル内に情報を保存、削除、更新する過程は、処理部10内のインデックステーブル管理部14で実行することができる。

【0108】

具体的に、処理部10は、第1の重複除去単位に対応したチャンクのサイズを基礎としてデータブロックを生成することができ、第2の重複除去単位に対応したチャンクのサイズを基礎としてデータブロックを生成することができ、第3の重複除去単位に対応したチャンクのサイズを基礎としてデータブロックを生成することができ、第4の重複除去単位に対応したチャンクのサイズを基礎としてデータブロックを生成することができる。一方、第4の重複除去単位がデータ重複除去を実行しないことを意味する場合、処理部10は、データに対するデータブロックを生成しないこともある。

【0109】

処理部10は、ハッシュアルゴリズム(例えば、SHA-1、SHA-2、SHA-3など)を使用してデータブロックに対する固有の識別子を生成することができる。データブロックに対する固有の識別子を生成する方法は、上述の説明に限定されるものではなくて、公知の多様な方法を使用してデータブロックに対する固有の識別子を生成することができる。

【0110】

処理部10は、固有の識別子がインデックステーブル内に存在するか否かを判断することができる。固有の識別子がインデックステーブル内に存在しない場合、これは固有の識別子に対応したデータブロックが保存されていない状態であるので、処理部10は、重複除去を実行しないで固有の識別子とデータブロックを保存部20に保存することができる。一方、固有の識別子がインデックステーブル内に存在する場合、これは固有の識別子に対応したデータブロックが既に保存されている状態であるので、処理部10は、重複除去(すなわち、固有の識別子に対応したデータブロック削除)を実行することができる。

【0111】

保存部20は、処理部10で処理される情報及び処理された情報を保存することができる。例えば、保存部20は、データの入力要請、データの出力要請、データに対する接近特性、第1の差、予め定義された第1の時間、第2の差、予め定義された第2の時間、インデックステーブル、重複除去単位情報などを保存することができる。

【0112】

以上、添付した図面を参照して本発明の実施形態について説明したが、本発明が属する技術の分野における通常の知識を有する者であれば、本発明の技術的思想を逸脱しない範囲内で、様々な置換、変形及び変更が可能であるので、上述した実施例及び添付された図面に限定されるものではない。

【符号の説明】

【0113】

- 10：処理部
- 11：接近特性獲得部
- 12：重複除去単位決定部

1 3 : 重複除去実行部

1 4 : インデックステーブル管理部

2 0 : 保存部

【書類名】 特許請求の範囲

【請求項 1】

データ重複除去装置で実行するデータ重複除去方法であって、
データの入力要請または出力要請を基盤として前記データに対する接近特性を獲得する段階と、

前記接近特性を基盤として前記データの重複除去単位を決定する段階と、
前記重複除去単位を基盤として前記データに対する重複除去を実行する段階と、
を含むことを特徴とするデータ重複除去方法。

【請求項 2】

前記接近特性は、前記データに対する接近時間、前記データに対する変更時間、前記データに対する順次的接近回数及び前記データに対する任意的接近回数の中で少なくとも一つを含むことを特徴とする請求項 1 に記載のデータ重複除去方法。

【請求項 3】

前記接近特性を獲得する段階は、
前記データの入力要請を受信した場合、前記データの入力要請に対する時間情報を基盤として前記接近時間及び前記変更時間を獲得する段階と、
前記データの入力要請の連続性を基盤として前記順次的接近回数または前記任意的接近回数を獲得する段階と、を含むことを特徴とする請求項 2 に記載のデータ重複除去方法。

【請求項 4】

前記接近特性を獲得する段階は、
前記データの出力要請を受信した場合、前記データの出力要請に対する時間情報を基盤として前記接近時間を獲得する段階と、
前記データの出力要請の連続性を基盤として前記順次接近回数または前記任意的接近回数を獲得する段階と、を含むことを特徴とする請求項 2 に記載のデータ重複除去方法。

【請求項 5】

前記重複除去単位を決定する段階は、
前記データに対する現在接近時間と前記データに対する以前変更時間に対する第 1 の差を算出する段階と、
前記第 1 の差が予め定義された第 1 の時間以下の場合、前記重複除去単位をデータに対する重複除去可能性が一番低い第 4 の重複除去単位で決定する段階と、
前記第 1 の差が予め定義された第 1 の時間を超過する場合、前記データに対する現在接近時間と前記データに対する以前接近時間に対する第 2 の差を算出する段階と、
前記第 2 の差が予め定義された第 2 の時間を超過する場合、前記重複除去単位をデータに対する重複除去可能性が一番高い第 1 の重複除去単位で決定する段階と、
前記第 2 の差が予め定義された第 2 の時間以下の場合、前記データに対する任意的接近回数が順次的接近回数以上であれば、前記重複除去単位を第 1 の重複除去単位より重複除去可能性が低い第 2 の重複除去単位で決定する段階と、
前記第 2 の差が予め定義された第 2 の時間以下の場合、前記データに対する任意的接近回数が順次的接近回数未満であれば、前記重複除去単位を第 2 の重複除去単位より重複除去可能性が低い第 3 の重複除去単位で決定する段階と、を含むことを特徴とする請求項 2 に記載のデータ重複除去方法。

【請求項 6】

前記第 4 の重複除去単位は、前記データに対する重複除去を実行しないことを意味することを特徴とする請求項 5 に記載のデータ重複除去方法。

【請求項 7】

前記データに対する重複除去を実行する段階は、
前記重複除去単位を基盤として前記データに対する少なくとも一つのデータブロック (block) を生成する段階と、
前記データブロックに対する固有の識別子を生成する段階と、
前記固有の識別子がインデックステーブル (index table) 内に存在するか否

かを判断する段階と、

前記固有の識別子が前記インデックステーブル内に存在する場合、前記固有の識別子に対応したデータブロックを除去する段階と、

前記固有の識別子が前記インデックステーブル内に存在しない場合、前記固有の識別子と前記固有の識別子に対応したデータブロックを保存する段階と、を含むことを特徴とする請求項 1 に記載のデータ重複除去方法。

【請求項 8】

前記固有の識別子を生成する段階は、ハッシュアルゴリズム(h a s h a l g o r i t h m)を使用して前記データブロックに対する固有の識別子を生成することを特徴とする請求項 7 に記載のデータ重複除去方法。

【請求項 9】

前記重複除去単位は、データに対する重複除去可能性を基盤として少なくとも一つの重複除去単位で分類されることを特徴とする請求項 1 に記載のデータ重複除去方法。

【請求項 10】

データの入力要請または出力要請を基盤として前記データに対する接近特性を獲得し、前記接近特性を基盤として前記データの重複除去単位を決定し、前記重複除去単位を基盤として前記データに対する重複除去を実行する処理部と、

前記処理部で処理される情報及び処理された情報を保存する保存部と、を含むことを特徴とするデータ重複除去装置。

【請求項 11】

前記接近特性は、前記データに対する接近時間、前記データに対する変更時間、前記データに対する順次的接近回数及び前記データに対する任意的接近回数の中で少なくとも一つを含むことを特徴とする請求項 10 に記載のデータ重複除去装置。

【請求項 12】

前記処理部は、前記データの入力要請を受信した場合、前記データの入力要請に対する時間情報を基盤として前記接近時間及び前記変更時間を獲得し、前記データの入力要請の連続性を基盤として前記順次的接近回数または前記任意的接近回数を獲得することを特徴とする請求項 11 に記載のデータ重複除去装置。

【請求項 13】

前記処理部は、前記データの出力要請を受信した場合、前記データの出力要請に対する時間情報を基盤として前記接近時間を獲得し、前記データの出力要請の連続性を基盤として前記順次的接近回数または前記任意的接近回数を獲得することを特徴とする請求項 11 に記載のデータ重複除去装置。

【請求項 14】

前記処理部は、

前記重複除去単位を決定する場合、前記データに対する現在接近時間と前記データに対する以前変更時間に対する第 1 の差を算出し、前記第 1 の差が予め定義された第 1 の時間以下の場合、前記重複除去単位をデータに対する重複除去可能性が一番低い第 4 の重複除去単位で決定し、前記第 1 の差が予め定義された第 1 の時間を超過する場合、前記データに対する現在接近時間と前記データに対する以前接近時間に対する第 2 の差を算出し、前記第 2 の差が予め定義された第 2 の時間を超過する場合、前記重複除去単位をデータに対する重複除去可能性が一番高い第 1 の重複除去単位で決定し、

前記第 2 の差が予め定義された第 2 の時間以下の場合、前記データに対する任意的接近回数が順次的接近回数以上であれば、前記重複除去単位を第 1 の重複除去単位より重複除去可能性が低い第 2 の重複除去単位で決定し、

前記第 2 の差が予め定義された第 2 の時間以下の場合、前記データに対する任意的接近回数が順次的接近回数未満であれば、前記重複除去単位を第 2 の重複除去単位より重複除去可能性が低い第 3 の重複除去単位で決定することを特徴とする請求項 11 に記載のデータ重複除去装置。

【請求項 15】

前記第4の重複除去単位は、前記データに対する重複除去を実行しないことを意味することを特徴とする請求項14に記載のデータ重複除去装置。

【請求項16】

前記処理部は、前記データに対する重複除去を実行する場合、前記重複除去単位を基盤として前記データに対する少なくとも一つのデータブロック(block)を生成し、前記データブロックに対する固有の識別子を生成し、前記固有の識別子がインデックステーブル(index table)内に存在するか否かを判断し、前記固有の識別子が前記インデックステーブル内に存在する場合、前記固有の識別子に対応したデータブロックを除去し、前記固有の識別子が前記インデックステーブル内に存在しない場合、前記固有の識別子と前記固有の識別子に対応したデータブロックを保存することを特徴とする請求項10に記載のデータ重複除去装置。

【請求項17】

前記処理部は、前記固有の識別子を生成する場合、ハッシュアルゴリズム(hash algorithm)を使用して前記データブロックに対する固有の識別子を生成することを特徴とする請求項16に記載のデータ重複除去装置。

【請求項18】

前記重複除去単位は、データに対する重複除去可能性を基盤として少なくとも一つの重複除去単位で分類されることを特徴とする請求項10に記載のデータ重複除去装置。

【書類名】 要約書

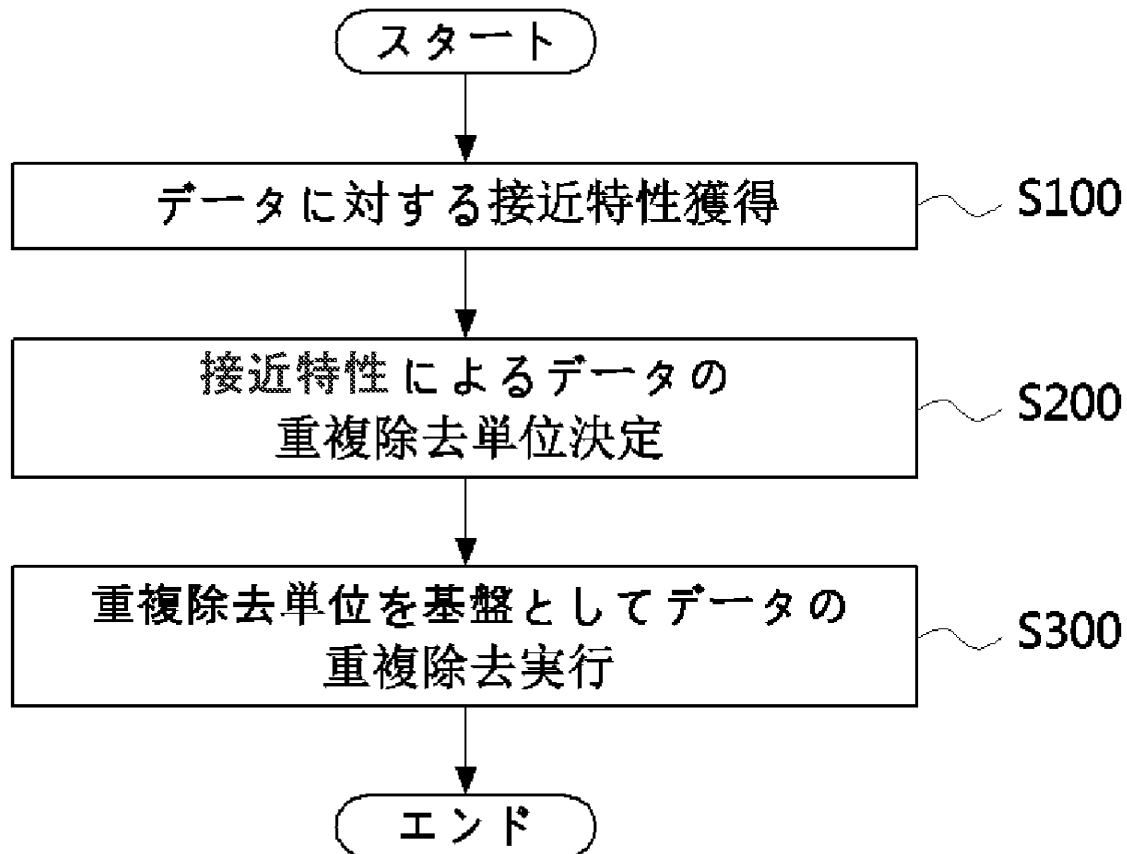
【要約】

本発明は、データ重複除去方法及び装置に関する。データの重複除去方法は、データの入力要請または出力要請を基盤としてデータに対する接近特性を獲得する段階と、接近特性を基盤としてデータの重複除去単位を決定する段階と、重複除去単位を基盤としてデータに対する重複除去を実行する段階と、を含む。したがって、低い入出力レイテンシを提供することができる。

【選択図】 図1

【書類名】 図面

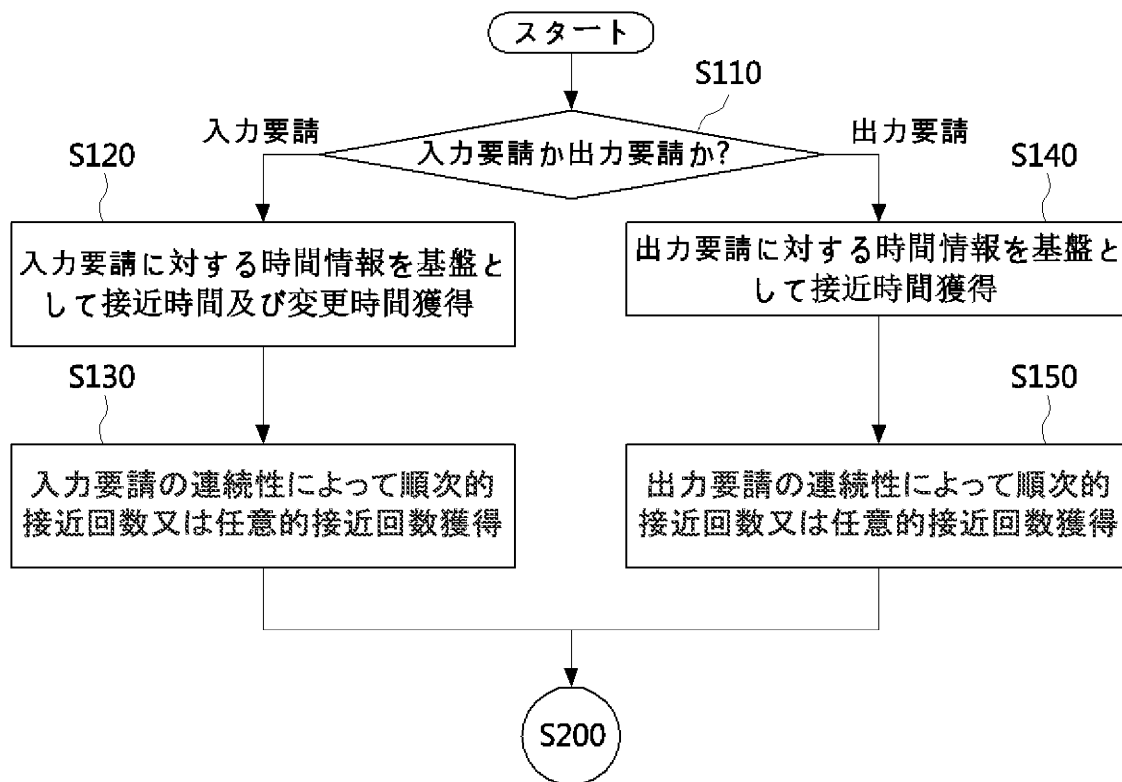
【図 1】



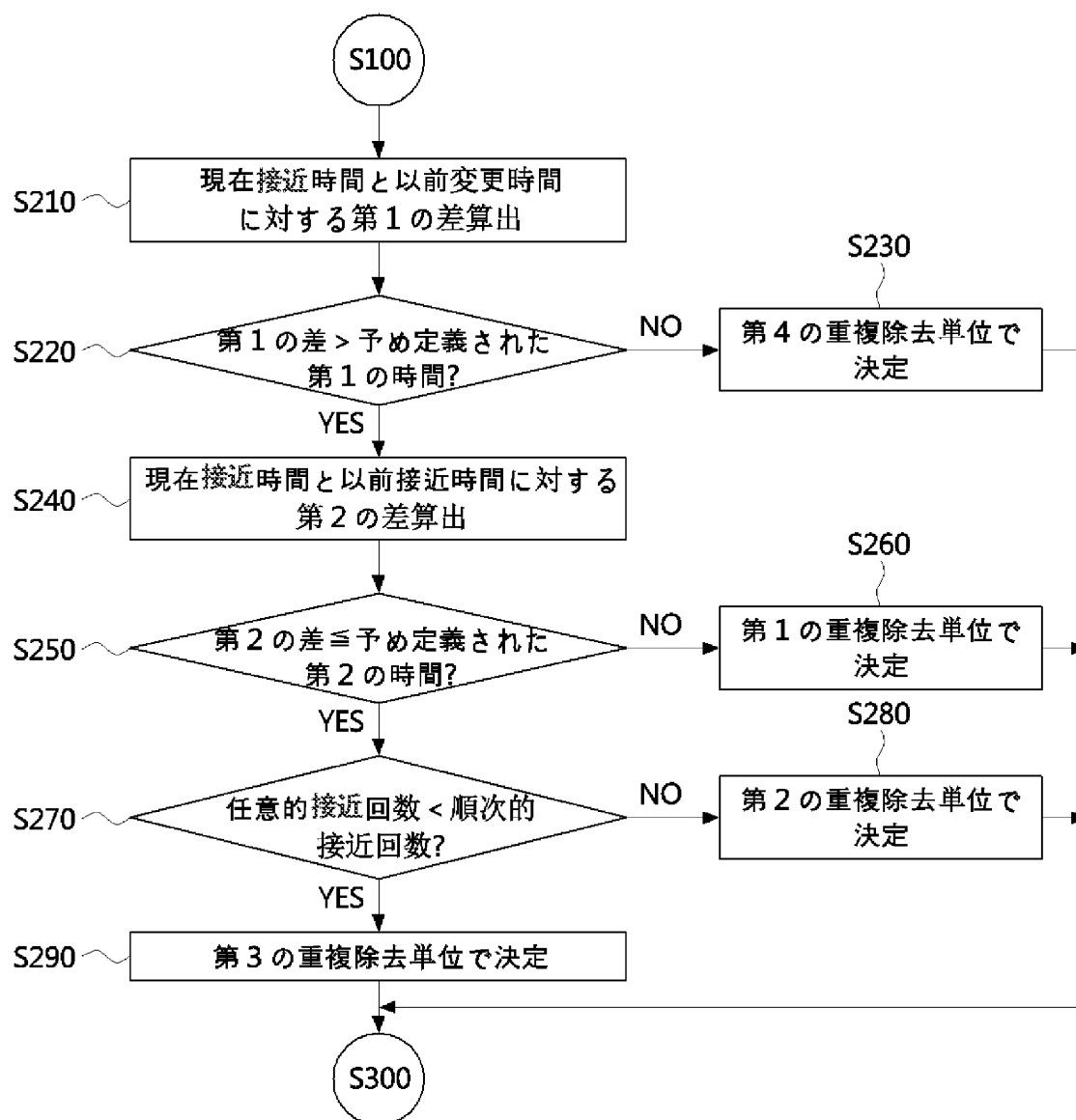
【図 2】

タイプ	フィールド	備考
時間型	m_time	最近変更時間
時間型	a_time	最近接近時間
正数型	seqCount	順次的接近回数
正数型	randCount	任意的接近回数

【図3】



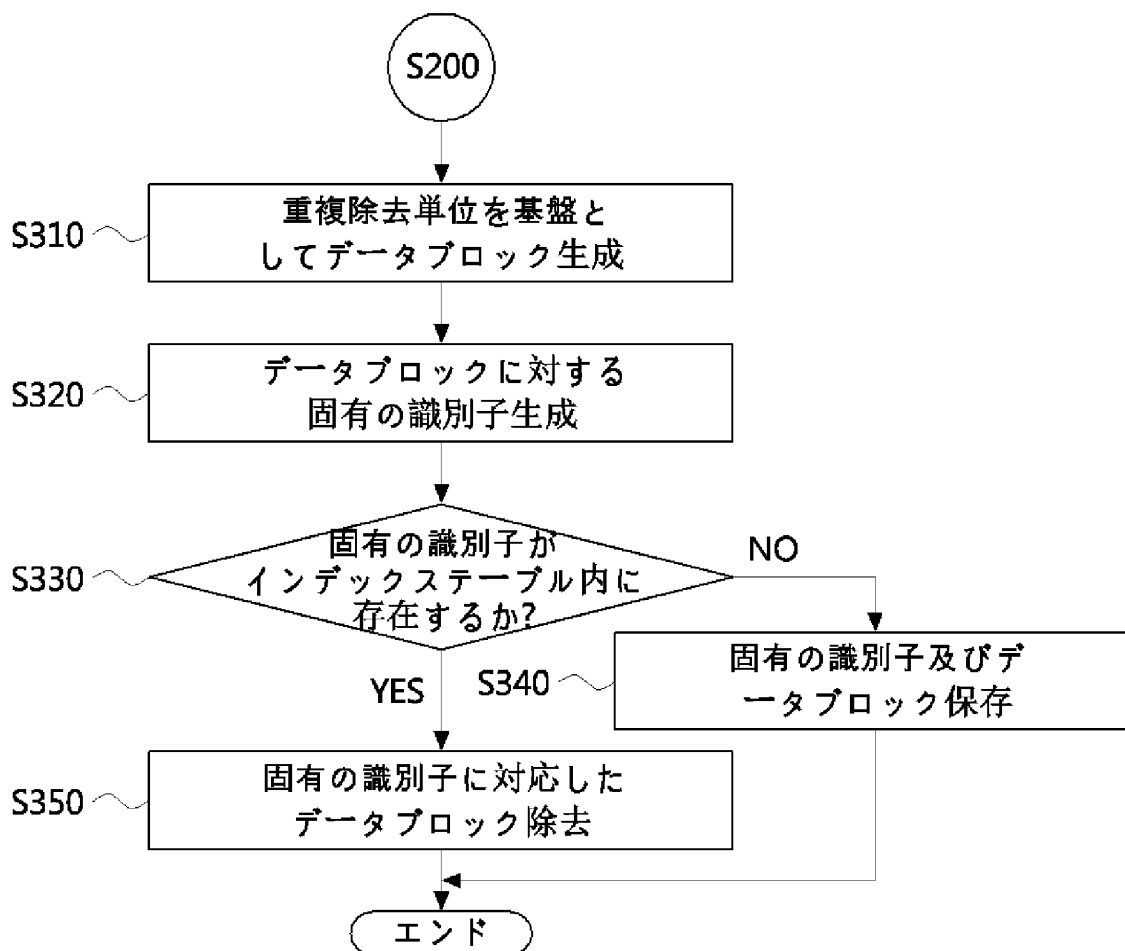
【図 4】



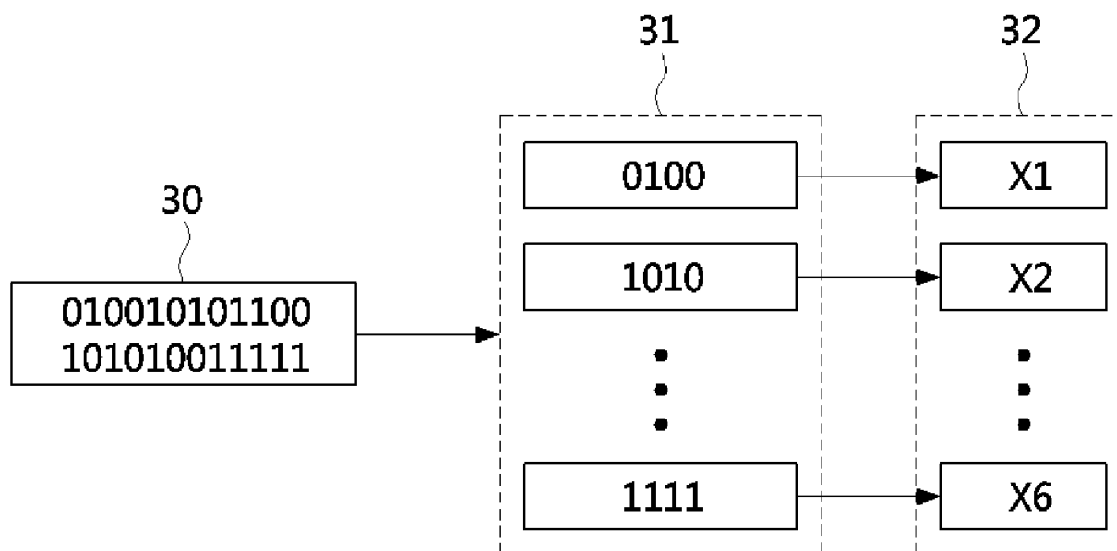
【図 5】

分類	重複除去単位	備考
よく使用しないデータ	第1の重複除去単位	- 一番小さいチャンク使用 - 一番大きい重複除去率
任意的接近がよく発生するデータ	第2の重複除去単位	- 相対的に小さいチャンク使用 - 相対的に大きい重複除去率
順次的接近がよく発生するデータ	第3の重複除去単位	- 相対的に大きいチャンク使用 - 相対的に小さい重複除去率
入力がよく発生するデータ	第4の重複除去単位	- 重複除去を実行しない

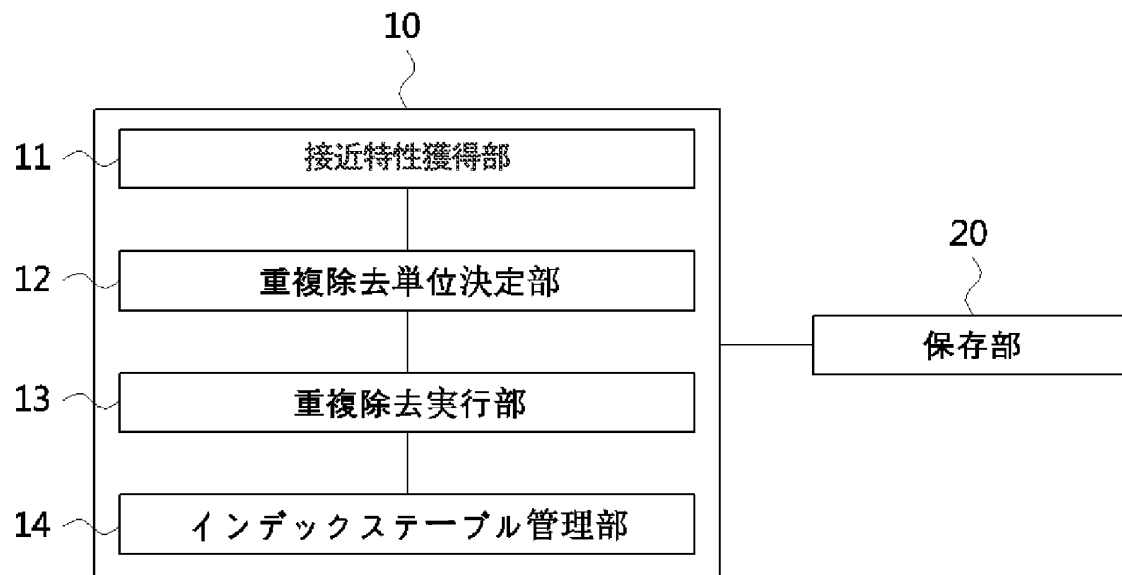
【図 6】



【図 7】



【図8】



【書類名】 出願審査請求書
【整理番号】 CP00744
【あて先】 特許庁長官 殿
【出願の表示】
【出願番号】 特願2014- 45770
【請求項の数】 18
【請求人】
【識別番号】 511258411
【氏名又は名称】 ポハン工科大学校産学協力団
【代理人】
【識別番号】 100121728
【弁理士】
【氏名又は名称】 井関 勝守
【手数料の表示】
【振替番号】 00016713
【納付金額】 190,000円

受領書

平成26年 3月 7日
特許庁長官識別番号 100121728
氏名(名称) 井関 勝守 様

以下の書類を受領しました。

項番	書類名	整理番号	受付番号	提出日	出願番号通知(事件の表示)
1	特許願	CP00744	51400497951	平26. 3. 7	特願2014- 45770 以上

受領書

平成26年 3月 7日
特許庁長官

識別番号 100121728
氏名(名称) 井関 勝守 様

以下の書類を受領しました。

項番	書類名	整理番号	受付番号	提出日	出願番号通知(事件の表示)
1	出願審査請求	CP00744	51400497964	平26. 3. 7	特願2014- 45770 以上