

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2020-77075
(P2020-77075A)

(43) 公開日 令和2年5月21日(2020.5.21)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 12/06 (2006.01)	G06F 12/06 522D	5B160
G06F 12/1027 (2016.01)	G06F 12/1027 120	5B205

審査請求 未請求 請求項の数 14 O L (全 21 頁)

<p>(21) 出願番号 特願2018-208545 (P2018-208545)</p> <p>(22) 出願日 平成30年11月6日 (2018.11.6)</p> <p>(出願人による申告) 平成28年度 国立研究開発法人新エネルギー・産業技術総合開発機構「IoT推進のための横断技術開発プロジェクト/高速大容量ストレージデバイス・システムの研究開発」委託研究、産業技術力強化法第19条の適用を受ける特許出願</p>	<p>(71) 出願人 000003078 株式会社東芝 東京都港区芝浦一丁目1番1号</p> <p>(74) 代理人 110002147 特許業務法人酒井国際特許事務所</p> <p>(72) 発明者 城田 祐介 東京都港区芝浦一丁目1番1号 株式会社東芝内</p> <p>(72) 発明者 金井 達徳 東京都港区芝浦一丁目1番1号 株式会社東芝内</p> <p>Fターム(参考) 5B160 MM03 5B205 MM51 VV03</p>
---	--

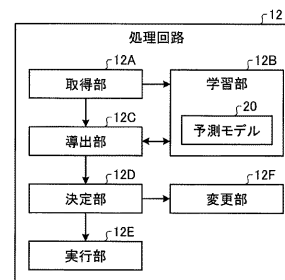
(54) 【発明の名称】 情報処理装置、情報処理方法、およびプログラム

(57) 【要約】

【課題】複数種類のメモリの使い分けに用いる情報を効率よく提供する。

【解決手段】情報処理装置10は、取得部12Aと、導出部12Cと、決定部12Dと、を備える。取得部12Aは、処理回路12の動作統計情報42Aを取得する。導出部12Cは、動作統計情報42Aから処理回路12のメモリアクセス特性42Bを導出するための予測モデル20に基づいて、取得した動作統計情報42Aからメモリアクセス特性42Bを導出する。決定部12Dは、導出したメモリアクセス特性42Bに基づいて、第1記憶部14Aより処理回路12によるアクセス速度が遅い第2記憶部14Bのデータを第1記憶部14Aへ転送し、第1記憶部14A内の該データにアクセスする第1アクセス方式、または、第2記憶部14B内のデータにアクセスする第2アクセス方式、の何れかのアクセス方式を決定する。

【選択図】 図3



【特許請求の範囲】**【請求項 1】**

処理回路の動作統計情報を取得する取得部と、
前記動作統計情報から前記処理回路のメモリアクセス特性を導出するための予測モデルに基づいて、取得した前記動作統計情報から前記メモリアクセス特性を導出する導出部と、
導出した前記メモリアクセス特性に基づいて、第 1 記憶部より前記処理回路によるアクセス速度が遅い第 2 記憶部のデータを前記第 1 記憶部へ転送し、前記第 1 記憶部内の該データにアクセスする第 1 アクセス方式、または、前記第 2 記憶部内のデータにアクセスする第 2 アクセス方式、の何れかのアクセス方式を決定する決定部と、
を備える情報処理装置。

10

【請求項 2】

決定された前記アクセス方式に応じて、前記データの前記第 2 記憶部から前記第 1 記憶部への転送および前記第 1 記憶部内の該データへのアクセス、または、前記第 2 記憶部内のデータへのアクセス、を実行する実行部、
を備える請求項 1 に記載の情報処理装置。

【請求項 3】

前記動作統計情報は、
前記処理回路が実行中のアプリケーションに割当てられた物理メモリサイズ、および、TLB (Translation Lookaside Buffer) ミスに関する動作統計情報、の少なくとも一方を含む、
請求項 1 または請求項 2 に記載の情報処理装置。

20

【請求項 4】

前記メモリアクセス特性は、
前記処理回路が単位期間あたりに使用したメモリサイズを示す、
請求項 1 ~ 請求項 3 の何れか 1 項に記載の情報処理装置。

【請求項 5】

前記決定部は、
導出した前記メモリアクセス特性が第 1 閾値より大きい場合、前記第 2 アクセス方式を決定し、該メモリアクセス特性が前記第 1 閾値以下の場合、前記第 1 アクセス方式を決定する、
請求項 4 に記載の情報処理装置。

30

【請求項 6】

第 1 閾値は、前記処理回路が利用可能な、前記第 1 記憶部のサイズ以上の値である、請求項 5 に記載の情報処理装置。

【請求項 7】

前記決定部は、
取得した前記動作統計情報に関する 1 または複数のアプリケーションの各々に割当てられた物理メモリサイズの合計値に対する、前記メモリアクセス特性の比率が、第 2 閾値より大きい場合、前記第 2 アクセス方式を決定し、前記第 2 閾値以下の場合、前記第 1 アクセス方式を決定する、
請求項 1 ~ 請求項 3 の何れか 1 項に記載の情報処理装置。

40

【請求項 8】

前記第 2 閾値は、前記合計値の N 分の 1 である (N は、2 以上の整数)、請求項 7 に記載の情報処理装置。

【請求項 9】

前記第 1 アクセス方式を決定した場合、
前記第 1 記憶部の利用可能なメモリサイズを変更する変更部、
を備える、請求項 1 ~ 請求項 8 の何れか 1 項に記載の情報処理装置。

【請求項 10】

50

前記動作統計情報と前記メモリアクセス特性との対応を示す教師データを複数含む教師データセットを用いて、前記予測モデルを学習する学習部、
を備える、請求項 1 ~ 請求項 9 の何れか 1 項に記載の情報処理装置。

【請求項 1 1】

前記教師データは、アプリケーションの命令単位ごとの、前記動作統計情報と前記メモリアクセス特性との対応を示す、
請求項 1 0 に記載の情報処理装置。

【請求項 1 2】

前記学習部は、学習用のアプリケーションを少なくとも 2 回実行し、一方の回の該アプリケーションの実行によって前記動作統計情報を取得し、他方の回の該アプリケーションの実行によって前記メモリアクセス特性を取得し、該アプリケーションの命令単位ごとに、取得した前記動作統計情報と取得した前記メモリアクセス特性との対応を示す前記教師データを生成する、
請求項 1 1 に記載の情報処理装置。

【請求項 1 3】

処理回路の動作統計情報を取得するステップと、
前記動作統計情報から前記処理回路のメモリアクセス特性を導出するための予測モデルに基づいて、取得した前記動作統計情報から前記メモリアクセス特性を導出するステップと、
導出した前記メモリアクセス特性に基づいて、第 1 記憶部より前記処理回路によるアクセス速度が遅い第 2 記憶部のデータを前記第 1 記憶部へ転送し、前記第 1 記憶部内の該データにアクセスする第 1 アクセス方式、または、前記第 2 記憶部内のデータにアクセスする第 2 アクセス方式、の何れかのアクセス方式を決定するステップと、
を含む情報処理方法。

【請求項 1 4】

処理回路の動作統計情報を取得するステップと、
前記動作統計情報から前記処理回路のメモリアクセス特性を導出するための予測モデルに基づいて、取得した前記動作統計情報から前記メモリアクセス特性を導出するステップと、
導出した前記メモリアクセス特性に基づいて、第 1 記憶部より前記処理回路によるアクセス速度が遅い第 2 記憶部のデータを前記第 1 記憶部へ転送し、前記第 1 記憶部内の該データにアクセスする第 1 アクセス方式、または、前記第 2 記憶部内のデータにアクセスする第 2 アクセス方式、の何れかのアクセス方式を決定するステップと、
をコンピュータに実行させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、情報処理装置、情報処理方法、およびプログラムに関する。

【背景技術】

【0002】

MRAM (Magnetoresistive Random Access Memory)、ReRAM (Resistive RAM)、PCM (Phase-Change Memory) などの各種のストレージクラスメモリ (SCM) が開発されている。SCM は、DRAM (Dynamic Random Access Memory) に比べてアクセス速度は遅いが、集積度が高い。一方、DRAM は、SCM に比べて集積度は低いが、アクセス速度が速い。このため、複数種類のメモリを搭載したシステムの場合、これらのメモリを使い分けて用いる必要がある。

【0003】

しかし、従来では、複数の種類のメモリの使い分けに用いる情報が管理されておらず、また、この情報を収集する手段を有していなかった。このため、従来では、複数種類のメ

10

20

30

40

50

メモリの使い分けに用いる情報を効率よく提供することは困難であった。

【先行技術文献】

【非特許文献】

【0004】

【非特許文献1】R. F. Freitas and W. W. Wilcke, "Storage-class Memory: The Next Storage System Technology", IBM Journal of Research and Development Vol. 52 No. 4, pp. 439-447, 2008.

【発明の概要】

10

【発明が解決しようとする課題】

【0005】

本発明が解決しようとする課題は、複数種類のメモリの使い分けに用いる情報を効率よく提供することができる、情報処理装置、情報処理方法、およびプログラムを提供することである。

【課題を解決するための手段】

【0006】

実施形態の情報処理装置は、取得部と、導出部と、決定部と、を備える。取得部は、処理回路の動作統計情報を取得する。導出部は、前記動作統計情報から前記処理回路のメモリアクセス特性を導出するための予測モデルに基づいて、取得した前記動作統計情報から前記メモリアクセス特性を導出する。決定部は、導出した前記メモリアクセス特性に基づいて、第1記憶部より前記処理回路によるアクセス速度が遅い第2記憶部のデータを前記第1記憶部へ転送し、前記第1記憶部内の該データにアクセスする第1アクセス方式、または、前記第2記憶部内のデータにアクセスする第2アクセス方式、の何れかのアクセス方式を決定する。

20

【図面の簡単な説明】

【0007】

【図1】情報処理装置の構成の一例を示す模式図。

【図2】物理アドレス空間の模式図。

【図3】処理回路の機能ブロック図。

30

【図4】予測モデルの学習の説明図。

【図5】動作統計情報とメモリアクセス特性との関係の説明図。

【図6A】動作統計情報の模式図。

【図6B】メモリアクセス特性の模式図。

【図7】導出部および決定部の処理の説明図。

【図8】アクセス方式の決定の説明図。

【図9】情報処理の手順のフローチャート。

【図10】情報処理の手順のフローチャート。

【発明を実施するための形態】

【0008】

40

以下に添付図面を参照して、本実施の形態の詳細を説明する。

【0009】

図1は、本実施の形態の情報処理装置10の構成の一例を示す模式図である。情報処理装置10は、処理回路12と、キャッシュメモリ16と、管理装置18と、を備える。情報処理装置10のメモリバスには、記憶部14が接続されている。

【0010】

処理回路12とキャッシュメモリ16、処理回路12と管理装置18、および、キャッシュメモリ16と管理装置18、の各々は、データや信号を授受可能に接続されている。処理回路12および管理装置18と記憶部14とは、データや信号を授受可能に接続されている。

50

【0011】

処理回路12は、1または複数のプロセッサを有する。プロセッサは、例えば、CPU (Central Processing Unit) である。プロセッサは、1または複数のCPUコアを含んでいてもよい。処理回路12は、1または複数のアプリケーションプログラムの実行に応じて、管理装置18を介して、記憶部14からのデータ読出しや、記憶部14へのデータ書込みを行う。

【0012】

なお、以下では、アプリケーションプログラムを、単に、アプリケーション、と称して説明する場合がある。また、記憶部14からのデータ読出および記憶部14へのデータ書込みを総称して説明する場合には、記憶部14へのアクセス、と称して説明する場合がある。

10

【0013】

処理回路12および管理装置18は、記憶部14に記憶されているデータをキャッシュメモリ16に一時的に記憶し、処理に用いる。

【0014】

記憶部14は、処理回路12による作業領域として用いられるメインメモリである。本実施の形態の情報処理装置10は、複数種類の記憶部14を備える。すなわち、本実施の形態の情報処理装置10は、複数種類の記憶部14を、メインメモリとして用いる。

【0015】

複数種類の記憶部14は、処理回路12によるアクセス速度が互いに異なる。なお、以下では、処理回路12によるアクセス速度を、単にアクセス速度と称して説明する場合がある。また、アクセス速度は、アクセス遅延とも呼ばれることもある。アクセス速度が速いというのは、アクセス遅延時間が短いことである。

20

【0016】

本実施の形態では、情報処理装置10は、アクセス速度の異なる複数種類の記憶部14として、第1記憶部14Aと、第2記憶部14Bと、を備える。なお、情報処理装置10は、3種類以上の記憶部14を備えた構成であってもよい。

【0017】

第1記憶部14Aは、第2記憶部14Bに比べてアクセス速度が速い。また、本実施の形態では、第1記憶部14Aは、第2記憶部14Bより集積度が低い。

30

【0018】

第1記憶部14Aは、例えば、揮発性メモリである。具体的には、第1記憶部14Aは、DRAM (Dynamic Random Access Memory) である。なお、第1記憶部14Aは、DRAMと同程度に高速アクセスが可能な、MRAM (Magnetoresistive Random Access Memory) 等の不揮発メモリであってもよい。

【0019】

一方、第2記憶部14Bは、第1記憶部14Aに比べてアクセス速度が遅い。また、本実施の形態では、第2記憶部14Bは、第1記憶部14Aより集積度が高い。すなわち、第2記憶部14Bは、第1記憶部14Aより容量が大きい。

40

【0020】

第2記憶部14Bは、例えば、不揮発性メモリである。具体的には、第2記憶部14Bは、DRAMより大容量な大容量高速不揮発メモリ (Non-volatile Memory) である。

【0021】

更に具体的には、第2記憶部14Bは、MRAM、PCM (Phase Change Memory)、PRAM (Phase Random Access Memory)、PCRAM (Phase Change Random Access Memory)、ReRAM (Resistance Change Random Access Memory)、FeRAM (Ferroelectric Random Access Memory) である。

50

ss Memory)、3DXPointまたはMemristorなどである。

【0022】

また、第2記憶部14Bは、いわゆるストレージクラスメモリ(SCM)と呼ばれるメモリであってもよい。また、第2記憶部14Bは、複数の半導体装置を1つの基板または筐体等に設けたモジュールであってもよい。

【0023】

本実施の形態では、第1記憶部14AがDRAMであり、第2記憶部14BがSCMである場合を、一例として説明する。なお、第1記憶部14Aのアクセス速度が第2記憶部14Bより速ければよく、これらの組合せは、第1記憶部14AがDRAMであり第2記憶部14BがSCMである形態に限定されない。例えば、第1記憶部14AがMRAMであり、第2記憶部14BがReRAMであってもよい。

10

【0024】

なお、第1記憶部14Aおよび第2記憶部14Bを総称して説明する場合には、単に、記憶部14と称して説明する。

【0025】

記憶部14は、複数の第1領域を含む。第1領域は、複数の第2領域を含む。言い換えると、本実施の形態では、処理回路12および管理装置18が、第1記憶部14Aおよび第2記憶部14Bを、第1領域毎に管理すると共に、第1領域内の第2領域ごとに管理する。

【0026】

図2は、処理回路12から見た物理アドレス空間を示す模式図である。

20

【0027】

図2に示すように、第1記憶部14Aおよび第2記憶部14Bは、それぞれ、複数の第1領域を含む。

【0028】

第1領域は、例えば、処理回路12によるデータの管理単位、あるいは、処理回路12上で動作するオペレーティングシステムによるデータの管理単位(例えば、ページ)である。ページは例えば4KBなどである。言い換えると、第1領域は、第1記憶部14Aと第2記憶部14Bとの間でデータを転送するときの、転送単位である。なお、第1領域は、処理回路12によるデータの管理単位の所定数倍の単位などであってもよい。本実施の形態では、第1領域が、ページに相当する場合を、一例として説明する。

30

【0029】

第2領域は、第1領域より小さい領域である。例えば、第2領域は、処理回路12による記憶部14に対するアクセスに伴うデータの書き換え単位である。すなわち、第2領域は、処理回路12によるデータのアクセス単位である。具体的には、第2領域は、キャッシュラインと呼ばれる単位である。キャッシュラインは、キャッシュメモリ16に対するデータの書き換え単位に対応する。すなわち、処理回路12からのメモリアクセス要求を受けた管理装置18は、キャッシュラインの単位で、第1記憶部14Aまたは第2記憶部14Bにアクセスする。

【0030】

キャッシュラインは、例えば、64バイトである。なお、第2領域は、キャッシュラインよりも小さい単位(例えば、バイト単位)であってもよい。また、第2領域は、キャッシュラインのサイズの所定数倍の単位などであってもよい。

40

【0031】

本実施の形態では、処理回路12および管理装置18は、図2に示す物理アドレス空間15にマッピングされた第1記憶部14Aおよび第2記憶部14B内の領域を、第1領域のサイズ(例えば、ページサイズ)に区切って管理する。そして、処理回路12および管理装置18は、ページテーブルを用いて論理アドレスから物理アドレスに変換することで、仮想記憶を実現する。

【0032】

50

図 1 に戻り、説明を続ける。管理装置 1 8 は、処理回路 1 2 による、複数種類の記憶部 1 4 (第 1 記憶部 1 4 A、第 2 記憶部 1 4 B) に対するアクセスを管理する。管理装置 1 8 は、メモリ管理ユニット (MMU: Memory Management Unit) などと称される場合がある。管理装置 1 8 はメモリコントローラなどであってもよい。

【0033】

管理装置 1 8 は、処理回路 1 2 から受付けたメモリアクセス要求を処理する。メモリアクセス要求は、処理回路 1 2 から記憶部 1 4 に対するアクセス要求である。メモリアクセス要求は、記憶部 1 4 へのデータ書込み、または記憶部 1 4 からのデータ読出しを示す。メモリアクセス要求は、アクセス対象の記憶部 1 4 の第 1 領域のアドレス情報および第 2 領域のアドレス情報を含む。これらのアドレス情報は、論理アドレスによって表される。

【0034】

管理装置 1 8 は、処理回路 1 2 から受付けたメモリアクセス要求によって示されるアクセス対象のデータがキャッシュメモリ 1 6 に格納されていない場合、記憶部 1 4 へアクセスする。この場合、管理装置 1 8 は、処理回路 1 2 から受付けたメモリアクセス要求によって示される、アクセス対象の記憶部 1 4 における第 1 領域内の第 2 領域にアクセスする。そして、管理装置 1 8 は、アクセスした該第 2 領域に対して、メモリアクセス要求によって示される処理 (書込みや読出し) を実行する。

【0035】

具体的には、処理回路 1 2 から受付けたメモリアクセス要求が特定の第 2 領域への書込みを示す場合がある。この場合、管理装置 1 8 は、メモリアクセス要求に示される、アクセス対象の記憶部 1 4 における、アクセス対象の第 1 領域内の第 2 領域に、メモリアクセス要求に示されるデータを書込む。また、処理回路 1 2 から受付けたメモリアクセス要求が特定の第 2 領域からのデータ読出しを示す場合がある。この場合、管理装置 1 8 は、メモリアクセス要求に示される、アクセス対象の記憶部 1 4 における、アクセス対象の第 1 領域内の第 2 領域からデータを読出し、キャッシュメモリ 1 6 に格納するとともに、処理回路 1 2 へ出力する。

【0036】

次に、処理回路 1 2 について詳細を説明する。上述したように、処理回路 1 2 は、1 または複数のアプリケーションの実行に応じて、管理装置 1 8 を介して、記憶部 1 4 へのアクセスを行う。

【0037】

図 3 は、処理回路 1 2 の機能ブロック図の一例である。処理回路 1 2 は、取得部 1 2 A と、学習部 1 2 B と、導出部 1 2 C と、決定部 1 2 D と、実行部 1 2 E と、変更部 1 2 F と、を備える。

【0038】

取得部 1 2 A、学習部 1 2 B、導出部 1 2 C、決定部 1 2 D、実行部 1 2 E、および変更部 1 2 F の少なくとも 1 つは、CPU などのプロセッサにプログラムを実行させること、すなわちソフトウェアにより実現してもよい。また、取得部 1 2 A、学習部 1 2 B、導出部 1 2 C、決定部 1 2 D、実行部 1 2 E、および変更部 1 2 F の少なくとも 1 つは、専用の IC (Integrated Circuit) などのハードウェアにより実現してもよい。また、取得部 1 2 A、学習部 1 2 B、導出部 1 2 C、決定部 1 2 D、実行部 1 2 E、および変更部 1 2 F の少なくとも 1 つは、ソフトウェアおよびハードウェアを併用して実現してもよい。複数のプロセッサを用いる場合、各プロセッサは、これらの取得部 1 2 A、学習部 1 2 B、導出部 1 2 C、決定部 1 2 D、実行部 1 2 E、および変更部 1 2 F のうち 1 つを実現してもよいし、各部のうち 2 以上を実現してもよい。

【0039】

取得部 1 2 A は、処理回路 1 2 の動作統計情報を取得する。

【0040】

動作統計情報とは、処理回路 1 2 の動作に関する情報の統計値である。詳細には、動作統計情報は、処理回路 1 2 が 1 または複数のアプリケーションを実行時の、動作に関する

10

20

30

40

50

情報の統計値である。動作に関する情報の統計値とは、単位期間Tあたりの、動作に関する情報を示す。単位期間Tは、予め設定すればよい。動作統計情報は、管理装置18やキャッシュメモリ16や情報処理装置10の動作に関する情報の統計値であってもよい。動作統計情報は、例えば、プロセッサが備える、ハードウェアイベントを測定する性能カウンタにより収集される。動作統計情報は、例えば、OSが管理する、情報処理装置の状態やOS内部の状態を示す情報(例えば、OS内部のイベント発生回数の統計情報)全般であってもよい。

【0041】

具体的には、動作統計情報は、性能カウンタなどで収集される、単位期間Tあたりの、TLB(Translation Lookaside Buffer)ミスの回数、TLBミスに関する動作統計情報、キャッシュメモリの各階層(L1キャッシュ、L2キャッシュ、L3キャッシュ、LLC(Last Level Cache)など)のキャッシュミスの回数、キャッシュミスのミスに関する動作統計情報、記憶部14への書込回数、記憶部14からの読出回数、STLB(Secondary level TLB)ミスの回数、STLBミスに関する動作統計情報のなどのハードウェアイベントのうち、1または複数によって表される。動作統計情報には、さらに、OSが管理する、当該実行期間に実行されているアプリケーションに当該実行期間中に割当てられている物理メモリサイズ(つまりアプリケーションの実行中のいずれかのタイミングでアクセスされる可能性があるメモリのサイズ)なども含まれていてもよい。なお、動作統計情報は、これらに限定されない。

【0042】

取得部12Aは、公知の方法で、処理回路12の動作統計情報を取得すればよい。例えば、取得部12Aは、処理回路12に設けられた性能カウンタから、単位期間Tごとの動作統計情報を順次取得すればよい。性能カウンタは、例えば、Intelプロセッサのパフォーマンスモニタリングカウンタ(Performance Monitoring Counter)などであるが、これに限定されない。取得部12Aは、性能カウンタと一体的に構成してもよい。また、取得部12Aと性能カウンタとを、別体として構成してもよい。本実施の形態では、取得部12Aと性能カウンタとを、別体として構成する形態を一例として説明する。

【0043】

学習部12Bは、教師データを複数含む教師データセットを用いて、予測モデル20を学習する。

【0044】

予測モデル20は、動作統計情報からメモリアクセス特性を導出するためのモデルである。予測モデル20は、学習によって生成される、学習モデルである。

【0045】

メモリアクセス特性とは、第1記憶部14Aおよび第2記憶部14Bに対する、処理回路12によるアクセスの特性を示す。

【0046】

例えば、メモリアクセス特性は、処理回路12がアプリケーション実行中に単位期間Tあたりに使用したメモリサイズによって表される。詳細には、メモリアクセス特性は、例えば、処理回路12がアプリケーション実行中に単位期間Tあたりに記憶部14やキャッシュメモリ16上のデータに対してロード命令とストア命令を発行した場合、そのデータの合計サイズである。具体的には、処理回路12がアプリケーション実行中において、単位期間TあたりNページに対してアクセスしたと想定する(Nは1以上の整数)。この場合、1ページの容量が4Kバイトとした場合、メモリアクセス特性は、"N"に、1ページの容量である"4Kバイト"を乗算した結果($N \times 4K$)によって表される。これは、一般に、ワーキングセットサイズなどとも呼ばれることもある。なお、メモリサイズは、処理回路12が記憶部14あるいは記憶部14のデータをキャッシュするキャッシュメモリ16に対してアクセスしたデータを格納するページのページ数(第1領域の数)で表し

10

20

30

40

50

てもよい。

【 0 0 4 7 】

学習部 1 2 B は、動作統計情報を入力としメモリアクセス特性を出力とする予測モデル 2 0 を、教師データセットを用いて学習する。

【 0 0 4 8 】

図 4 は、予測モデル 2 0 の学習の説明図である。教師データセット 4 0 は、複数の教師データ 4 2 を含む。教師データ 4 2 は、単位期間 T ごとに生成される。教師データ 4 2 は、動作統計情報 4 2 A と、メモリアクセス特性 4 2 B と、を含む、なお、1 つの教師データ 4 2 には、1 または複数の動作統計情報 4 2 A が含まれる。また、1 つの教師データ 4 2 には、1 または複数の動作統計情報 4 2 A に対応する 1 つの正解情報として、1 つのメモリアクセス特性 4 2 B が含まれる。

10

【 0 0 4 9 】

処理回路 1 2 は、予め、教師データセット 4 0 を用意する。例えば、処理回路 1 2 は、1 または複数の学習用のアプリケーション 3 0 (例えば、アプリケーション 3 0 A、アプリケーション 3 0 B、アプリケーション 3 0 C) を実行する。そして、処理回路 1 2 は、アプリケーション 3 0 を実行中における 1 または複数の動作統計情報 4 2 A とメモリアクセス特性 4 2 B との組からなる教師データ 4 2 を、単位期間 T ごとに生成する。この処理により、処理回路 1 2 は、予め、複数の教師データ 4 2 を含む教師データセット 4 0 を用意する。

【 0 0 5 0 】

図 5 は、学習用のアプリケーション 3 0 の実行期間 T A における、動作統計情報 4 2 A とメモリアクセス特性 4 2 B との関係の一例を示す説明図である。図 5 に示すグラフの縦軸は、動作統計情報またはメモリアクセス特性を示す。図 5 に示すグラフの横軸は、時間を示す。なお、図 5 には、動作統計情報 4 2 A が実行中のアプリケーションに割り当てられた物理メモリサイズを示し、メモリアクセス特性 4 2 B が、処理回路 1 2 がアプリケーション実行中に単位期間 T あたりに使用したメモリサイズを示す場合を、一例として示した。

20

【 0 0 5 1 】

処理回路 1 2 が、学習用のアプリケーション 3 0 を実行したときの、動作統計情報 4 2 A およびメモリアクセス特性 4 2 B の時間経過に伴う推移が、図 5 に示す推移を示したと想定する。この場合、学習部 1 2 B は、単位期間 T ごとの、動作統計情報 4 2 A とメモリアクセス特性 4 2 B との対応を、単位期間 T ごとの教師データ 4 2 として生成すればよい。なお、隣接するタイミングの教師データ 4 2 の単位期間 T は、一部が重複するタイミングであってもよく、また、互いに非重複のタイミングであってもよい。

30

【 0 0 5 2 】

学習部 1 2 B は、性能カウンタから取得部 1 2 A を介して、単位期間 T ごとの動作統計情報 4 2 A を取得することで、単位期間 T の教師データ 4 2 の生成に用いればよい。

【 0 0 5 3 】

また、学習部 1 2 B は、教師データ 4 2 に用いる単位期間 T ごとのメモリアクセス特性 4 2 B を、以下の方法により取得すればよい。

40

【 0 0 5 4 】

詳細には、学習部 1 2 B は、単位期間 T の最初のタイミング(例えば、 t_1)で、処理回路 1 2 に予めインストールされた OS (オペレーティングシステム) が管理するページテーブルの一部をリセットする。詳細には、学習部 1 2 B は、ページテーブルの全てのページのアクセス済フラグを " 0 " として " アクセス未 " とすることで、ページテーブルをリセットする。次に、学習部 1 2 B は、該単位期間 T の区間の終了タイミング(例えば、 t_2)に、該ページテーブルにおけるアクセス済フラグ(" 1 " となっているフラグ)を計数する。この計数処理により、学習部 1 2 B は、単位期間 T に処理回路 1 2 によってアクセスされたページ数を求める。そして、学習部 1 2 B は、該ページ数 (N) に、1 ページの容量である " 4 K バイト " を乗算した結果 ($N \times 4 K$) を、該単位期間 T のメモリア

50

クセス特性 4 2 B として取得する。

【 0 0 5 5 】

そして、学習部 1 2 B は、単位期間 T ごとに、上記処理を実行する。そして、学習部 1 2 B は、単位期間 T ごとに取得した、動作統計情報 4 2 A とメモリアクセス特性 4 2 B との対応を示す教師データ 4 2 を生成する。

【 0 0 5 6 】

ここで、上述したように、学習部 1 2 B は、性能カウンタから取得部 1 2 A を介して動作統計情報 4 2 A を取得する。このため、学習部 1 2 B は、学習用のアプリケーション 3 0 の実行中に、リアルタイムで動作統計情報 4 2 A を取得可能である。一方、メモリアクセス特性 4 2 B の取得には、学習部 1 2 B は、ページテーブルのリセット、アクセス済フラグの計数、メモリアクセス特性 4 2 B の計算、といった処理を、単位期間 T ごとに実行する必要がある。このため、学習部 1 2 B は、学習用のアプリケーション 3 0 の実行中に、リアルタイムでメモリアクセス特性 4 2 B を取得することは困難な場合がある。また、ページテーブルのリセット、アクセス済フラグの計数、メモリアクセス特性 4 2 B の計算、といった処理自体が、動作統計情報 4 2 A に影響を与えてしまう（アプリケーションのみを実行した場合の動作統計情報 4 2 A に対して大きく変化してしまう）可能性もあるため、それを回避するのが望ましい。

【 0 0 5 7 】

そこで、本実施の形態では、処理回路 1 2 は、学習用のアプリケーション 3 0 を 2 回実行する。そして、学習部 1 2 B は、1 回目および 2 回目の内の一方の回のアプリケーション 3 0 の実行時に動作統計情報 4 2 A を取得し、他方の回のアプリケーション 3 0 の実行時にメモリアクセス特性 4 2 B を取得する。そして、学習部 1 2 B は、単位期間 T に対応する、アプリケーション 3 0 の命令単位ごとに、動作統計情報 4 2 A とメモリアクセス特性 4 2 B との対応を示す教師データ 4 2 を生成する。

【 0 0 5 8 】

図 6 A は、学習用のアプリケーション 3 0 を実行したときの、動作統計情報 4 2 A の一例を示す模式図である。図 6 A には、動作統計情報 4 2 A として、TLB ミスに関する動作統計情報 4 2 A 1 と、実行中のアプリケーションに割当てられた物理メモリサイズ 4 2 A 2 と、を一例として示した。なお、図 6 A に示すグラフの縦軸は、動作統計情報を示し、横軸は、時間を示す。

【 0 0 5 9 】

図 6 A に示すように、処理回路 1 2 が、学習用のアプリケーション 3 0 を実行したときの、時間経過に伴う動作統計情報 4 2 A の推移が、図 6 A に示す推移を示したと想定する。そして、学習部 1 2 B は、単位期間 T を、該アプリケーション 3 0 が命令単位 S を実行する期間と定める。命令単位 S は、例えば、1 0 万回である。なお、命令単位 S の命令回数は、予め設定すればよく、1 0 万回に限定されない。

【 0 0 6 0 】

この場合、学習部 1 2 B は、アプリケーション 3 0 の命令単位 S ごとの動作統計情報 4 2 A を、性能カウンタから取得することで、命令単位 S ごと（すなわち、単位期間 T ごと）の動作統計情報 4 2 A を取得する。例えば、学習用のアプリケーション 3 0 を実行した期間 T A を、過去から未来に向かって、命令単位 S ごとに、フェーズ P 1、フェーズ P 2、フェーズ P 3、フェーズ P 4、フェーズ P 5 に分割して管理すると想定する。この場合、学習部 1 2 B は、各フェーズ（フェーズ P 1 ~ フェーズ P 5）の各々ごとの、動作統計情報 4 2 A を取得する。

【 0 0 6 1 】

次に、学習部 1 2 B は、同じ学習用のアプリケーション 3 0 を再度実行し、命令単位 S ごとのメモリアクセス特性 4 2 B を取得する。

【 0 0 6 2 】

図 6 B は、教師データ 4 2 用の動作統計情報 4 2 A の取得時と同じ学習用のアプリケーション 3 0 を実行したときの、メモリアクセス特性 4 2 B の一例を示す模式図である。図

10

20

30

40

50

6 Bには、メモリアクセス特性4 2 Bとして、処理回路1 2が単位期間Tあたりに使用したメモリサイズを一例として示した。なお、図6 Bに示すグラフの縦軸は、メモリアクセス特性を示し、横軸は、時間を示す。

【0063】

学習部1 2 Bは、単位期間Tに相当する命令単位Sごとに、ページテーブルのリセット、アクセス済フラグの計数、メモリアクセス特性4 2 Bの計算、を実行する。この処理により、学習部1 2 Bは、命令単位Sごとの、メモリアクセス特性4 2 Bを取得する。このため、学習部1 2 Bは、フェーズ(フェーズP 1~フェーズP 5)の各々の命令単位Sごとの、メモリアクセス特性4 2 Bを取得する。

【0064】

そして、学習部1 2 Bは、学習用のアプリケーション3 0の命令単位Sごとの、動作統計情報4 2 Aとメモリアクセス特性4 2 Bとの対応を示す、教師データ4 2を生成すればよい。

【0065】

ここで、学習部1 2 Bが、命令単位Sごとに、ページテーブルのリセット、アクセス済フラグの計数、メモリアクセス特性4 2 Bの計算、を実行するため、これらの実行に要する時間は、単位期間Tより長い期間T'となる場合がある。

【0066】

しかし、本実施の形態では、学習部1 2 Bが、動作統計情報4 2 Aをリアルタイムで取得するときの単位期間Tに相当する命令単位Sを基準として、該命令単位Sごとのメモリアクセス特性4 2 Bを取得する。このため、学習部1 2 Bは、処理回路1 2が実際にアプリケーション3 0を実行したときの、単位期間Tごとの動作統計情報4 2 Aとメモリアクセス特性4 2 Bとの対応を示す教師データ4 2を、精度良く生成することができる。

【0067】

図4に戻り説明を続ける。そして、学習部1 2 Bは、複数の教師データ4 2を含む教師データセット4 0を用いて、動作統計情報4 2 Aからメモリアクセス特性4 2 Bを導出するための予測モデル2 0を学習する。

【0068】

学習部1 2 Bは、公知の学習アルゴリズムを用いて、予測モデル2 0を学習すればよい。学習アルゴリズムは、例えば、線形回帰、k近傍法(KNN: K - Nearest Neighbor algorithm)、サポートベクターマシン、Random Forestなどであるが、これらに限定されない。

【0069】

学習部1 2 Bは、予め定められたタイミング毎に、教師データセット4 0に含まれる複数の教師データ4 2を用いて、予測モデル2 0を学習すればよい。また、例えば、学習部1 2 Bは、新たな教師データ4 2が教師データセット4 0に登録されるごとに、予測モデル2 0を学習してもよい。教師データ4 2登録は、任意のタイミングで実行すればよい。

【0070】

また、学習部1 2 Bは、新たな教師データ4 2を含む教師データセット4 0を用いて新たな予測モデル2 0を学習した場合には、学習部1 2 Bに登録されている予測モデル2 0を、新たに学習した予測モデル2 0に更新すればよい。すなわち、学習部1 2 Bには、1つの予測モデル2 0が格納された状態となる。

【0071】

図3に戻り説明を続ける。次に、導出部1 2 Cについて説明する。

【0072】

導出部1 2 Cは、学習部1 2 Bが学習した予測モデル2 0に基づいて、取得部1 2 Aで取得した動作統計情報4 2 Aからメモリアクセス特性4 2 Bを導出する。導出部1 2 Cがメモリアクセス特性4 2 Bの導出に用いる動作統計情報4 2 Aは、学習部1 2 Bによる予測モデル2 0の学習時とは異なり、処理回路1 2が学習用以外の実際のアプリケーション3 0を実行したときの動作統計情報4 2 Aである。導出部1 2 Cは、この動作統計情報4

10

20

30

40

50

2 A と、予測モデル 2 0 と、を用いて、メモリアクセス特性 4 2 B を導出する。

【 0 0 7 3 】

図 7 は、導出部 1 2 C および決定部 1 2 D による処理の一例の説明図である。

【 0 0 7 4 】

例えば、導出部 1 2 C は、最適化対象のアプリケーション 3 2 を処理回路 1 2 が実行しているときの動作統計情報 4 2 A を、取得部 1 2 A から取得する。そして、導出部 1 2 C は、取得した動作統計情報 4 2 A を予測モデル 2 0 へ入力することで、予測モデル 2 0 からの出力として、メモリアクセス特性 4 2 B を得る。

【 0 0 7 5 】

すなわち、導出部 1 2 C は、予測モデル 2 0 を用いて、取得部 1 2 A で取得した動作統計情報 4 2 A に対する、メモリアクセス特性 4 2 B の予測値を得る。

10

【 0 0 7 6 】

決定部 1 2 D は、導出部 1 2 C が導出したメモリアクセス特性 4 2 B に基づいて、アクセス方式を決定する。

【 0 0 7 7 】

アクセス方式は、処理回路 1 2 による記憶部 1 4 に対するアクセス方式を示す。本実施の形態では、アクセス方式は、第 1 アクセス方式、または、第 2 アクセス方式を示す。

【 0 0 7 8 】

第 1 アクセス方式は、第 2 記憶部 1 4 B のデータを第 1 記憶部 1 4 A へ転送し、該第 1 記憶部 1 4 A 内の該データにアクセスする、アクセス方式である。転送およびアクセスするデータは、該アクセス方式の決定に用いた動作統計情報 4 2 A の取得時に処理回路 1 2 がアクセスしていたデータである。言い換えると、転送およびアクセスするデータは、該アクセス方式の決定に用いた動作統計情報 4 2 A によって示される動作実行時に処理回路 1 2 がアクセスしていたデータである。

20

【 0 0 7 9 】

本実施の形態では、転送とは、コピーを意味する。上述したように、処理回路 1 2 は、ページ単位（第 1 領域の単位）でデータ転送を行う。また、処理回路 1 2 は、キャッシュライン（第 2 領域）の単位で、記憶部 1 4 へアクセスする。

【 0 0 8 0 】

このため、第 1 アクセス方式は、該アクセス方式の決定に用いた動作統計情報 4 2 A によって示される動作実行時に処理回路 1 2 がアクセスしていたデータを含むページ（第 1 領域）内のデータを、第 2 記憶部 1 4 B から第 1 記憶部 1 4 A へ転送し、転送後の第 1 記憶部 1 4 A の該データへアクセスすることを示す。

30

【 0 0 8 1 】

第 2 アクセス方式は、第 2 記憶部 1 4 B 内のデータにアクセスする方式を示す。本実施の形態では、処理回路 1 2 は、原則、第 2 記憶部 1 4 B 内のデータにアクセスする。そして、処理回路 1 2 は、特定の条件を満たした場合にのみ、第 2 記憶部 1 4 B から第 1 記憶部 1 4 A へデータを転送し、第 1 記憶部 1 4 A 内のデータへアクセスする。このため、第 2 アクセス方式は、データを第 2 記憶部 1 4 B に配置したまま、該第 2 記憶部 1 4 B にダイレクトにアクセスすることを示す。

40

【 0 0 8 2 】

例えば、決定部 1 2 D は、導出部 1 2 C が導出したメモリアクセス特性 4 2 B が第 1 閾値より大きい場合、第 2 アクセス方式を決定する。また、決定部 1 2 D は、該メモリアクセス特性 4 2 B が第 1 閾値以下の場合、第 1 アクセス方式を決定する。

【 0 0 8 3 】

図 8 は、アクセス方式の決定の説明図である。図 8 の横軸は時間を示し、縦軸は動作統計情報 4 2 A を示す。

【 0 0 8 4 】

例えば、導出部 1 2 C が、アプリケーション 3 2 の実行中に取得部 1 2 A で取得した動作統計情報 4 2 A を予測モデル 2 0 へ入力することで、図 8 に示すメモリアクセス特性 4

50

2 Bを導出したと想定する。

【0085】

なお、動作統計情報42Aが、処理回路12で実行中のアプリケーションに割り当てられた物理メモリサイズを示すと想定する。また、メモリアクセス特性42Bが、アプリケーション32実行中、単位期間Tあたりに処理回路12（アプリケーション32）が使用したメモリサイズを示すと想定する。

【0086】

この場合、決定部12Dは、導出したメモリアクセス特性42Bが第1閾値より大きい場合、第2アクセス方式を決定する。また、決定部12Dは、導出したメモリアクセス特性42Bが第1閾値以下である場合、第1アクセス方式を決定する。

10

【0087】

図8に示すように、決定部12Dは、アプリケーション32の実行期間TAの内の前半の期間Aのように、メモリアクセス特性42Bが第1閾値以下の場合、第1アクセス方式を決定する。

【0088】

すなわち、決定部12Dは、メモリアクセス特性42Bが第1閾値以下である状態を、処理回路12によるメモリアクセスが集中し、アクセスのローカリティの高い状態であると推定する。そして、決定部12Dは、メモリアクセス特性42Bが第1閾値以下の場合、第1アクセス方式を決定する。

【0089】

このため、決定部12Dは、処理回路12がアクセスのローカリティの高いデータにアクセスする場合、すなわち、記憶部14におけるメモリアクセスされる場所が集中している場合、データを第2記憶部14Bから第1記憶部14Aへ転送し、処理回路12が第1記憶部14A上のデータをキャッシュライン単位でアクセスするように、アクセス方式を決定することができる。

20

【0090】

一方、決定部12Dは、アプリケーション32の実行期間TAの後半の期間Bのように、メモリアクセス特性42Bが第1閾値を超える場合、第2アクセス方式を決定する。

【0091】

図8に示すように、決定部12Dは、アプリケーション32の実行期間TAの内の後半の期間Bのように、メモリアクセス特性42Bが第1閾値を超える場合、第2アクセス方式を決定する。

30

【0092】

すなわち、決定部12Dは、メモリアクセス特性42Bが第1閾値を超える状態を、処理回路12によるメモリアクセスが分散し、アクセスのローカリティが低く、使用中のメモリサイズが大きい状態であると推定する。そして、決定部12Dは、メモリアクセス特性42Bが第1閾値を超える場合、第2アクセス方式を決定する。

【0093】

メモリアクセス特性42Bが第1閾値を超える場合、第2アクセス方式を決定することで、決定部12Dは、処理回路12によるメモリアクセスの高速化を図ることができる。

40

【0094】

ここで、メモリアクセス特性42Bが第1閾値を超える場合には、第2記憶部14Bのデータを第1記憶部14Aへ転送しても、第1記憶部14Aの空き容量不足により、すぐに第1記憶部14Aから第2記憶部14Bへデータが転送されてしまう。すなわち、メモリアクセス特性42Bが第1閾値を超える場合に、第2記憶部14Bのデータを第1記憶部14Aへ転送すると、第1記憶部14Aと第2記憶部14Bとの間のページ単位のデータ入れ替えが頻発してしまい、処理回路12の性能低下を引き起こす可能性がある。

【0095】

そこで、本実施の形態では、決定部12Dは、メモリアクセス特性42Bが第1閾値を超える場合には、第2アクセス方式を決定する。このため、決定部12Dは、処理回路1

50

2によるメモリアクセスの高速化を図ることができる。

【0096】

なお、第1閾値には、任意の値を予め定めればよい。例えば、第1閾値には、処理回路12が利用可能な、第1記憶部14Aのサイズやこれに近い値であればよい。

【0097】

処理回路12が利用可能なサイズとは、具体的には、処理回路12が実行するアプリケーション32で利用可能な、第1記憶部14Aのサイズや、情報処理装置10で利用可能な記憶部14のサイズなどである。なお、第1閾値は、これらの利用可能なサイズに対して所定の割合大きい値であってもよい。また、第1閾値は、これらの利用可能なサイズに対して所定の割合小さい値であってもよい。

10

【0098】

なお、第1閾値を、処理回路12が実行するアプリケーション32で利用可能な第1記憶部14Aのサイズより大きい値、または、情報処理装置10で利用可能な第1記憶部14Aのサイズより大きい値とすると、以下の効果が得られる。例えば、第1記憶部は高速なSCMを想定しているため、第1アクセス方式が決定された場合における、第1記憶部14Aと第2記憶部14Bとの間のデータの転送は、高速に行う事が可能である。つまり、第2記憶部14Bから第1記憶部14Aへのデータ転送あるいは第1記憶部14Aから第2記憶部14Bへのデータ転送を積極的に行って頻繁にデータを入れ替えても速度低下は僅かで済み高速に処理することができる。このため、利用可能な第1記憶部14Aを最大限活用するという観点では、第1閾値を、利用可能な第1記憶部14Aのサイズより大きい値とすることで、利用可能な第1記憶部14Aを最大限活用して、より大きなメモリサイズを利用するアプリケーションが実行可能になる。また、利用可能な第1記憶部14Aを削減するという観点では、第1閾値を、利用可能な第1記憶部14Aのサイズより大きい値とすることで、実際に利用可能な第1記憶部14Aのサイズより小さいサイズで、使用メモリサイズが小さく(つまり消費電力が低く)速度低下も抑えた効率の良い処理を実行することができる。

20

【0099】

図3に戻り説明を続ける。次に、実行部12Eについて説明する。実行部12Eは、決定部12Dで決定されたアクセス方式に応じて、データの第2記憶部14Bから第1記憶部14Aへの転送および第1記憶部14A内の該データへのアクセス、または、第2記憶部14B内のデータへのアクセス、を実行する。

30

【0100】

すなわち、決定部12Dが第1アクセス方式を決定した場合、実行部12Eは、該アクセス方式の決定に用いた動作統計情報42Aによって示される動作実行時に処理回路12がアクセスしていたデータを含むページ(第1領域)を、第2記憶部14Bから第1記憶部14Aへ転送し、第1記憶部14Aにおける転送した該ページ内の該データへのアクセスを実行する。

【0101】

なお、実行部12Eによる、第1記憶部14Aから第2記憶部14Bへのデータ転送のタイミングは、限定されない。例えば、実行部12Eは、決定部12Dが第1アクセス方式を決定した直後、実行部12Eが次回該データへアクセスする時、予め定めた条件を満たしたタイミング、の何れであってもよい。予め定めた条件を満たすタイミングは、例えば、処理回路12による記憶部14へのメモリアクセスが所定値以下の期間などである。

40

【0102】

一方、決定部12Dが第2アクセス方式を決定した場合、実行部12Eは、該アクセス方式の決定に用いた動作統計情報42Aによって示される動作実行時に処理回路12がアクセスしていた、第2記憶部14B内のデータへのアクセスを継続して実行する。

【0103】

次に、変更部12Fについて説明する。決定部12Dが第1アクセス方式を決定した場合、変更部12Fは、第1記憶部14Aの、利用可能なメモリサイズを変更する。

50

【 0 1 0 4 】

詳細には、変更部 1 2 F は、第 1 記憶部 1 4 A の利用可能なメモリサイズを、第 1 アクセス方式の決定に用いたメモリアクセス特性 4 2 B によって示される、処理回路 1 2 がアプリケーション 3 2 の実行中に単位期間 T あたりに使用したメモリサイズに変更する。また、変更部 1 2 F は、第 1 記憶部 1 4 A の利用可能なメモリサイズを、該使用した該メモリサイズより一定の割合大きいサイズ、または、該使用した該メモリサイズより一定の割合小さいサイズに変更する。そして、変更部 1 2 F は、利用可能なメモリサイズを変更した後の第 1 記憶部 1 4 A における、利用不可能な領域を、パワーオフまたはセルフリフレッシュモードなどの低消費電力モードに設定すればよい。

【 0 1 0 5 】

なお、処理回路 1 2 は、変更部 1 2 F を備えない構成であってもよい。すなわち、処理回路 1 2 は、第 1 記憶部 1 4 A における利用可能なメモリサイズの変更を行わない形態であってもよい。

【 0 1 0 6 】

次に、本実施の形態の処理回路 1 2 が実行する情報処理の手順の一例を説明する。

【 0 1 0 7 】

図 9 は、本実施の形態の処理回路 1 2 が実行する情報処理の手順の一例を示す、フローチャートである。なお、図 9 に示す情報処理の手順の実行前に、学習部 1 2 B が、予測モデル 2 0 を学習済であるものとして説明する。また、図 9 には、変更部 1 2 F が第 1 記憶部 1 4 A の利用可能なメモリサイズの変更処理を実行しない形態を、一例として示した。

【 0 1 0 8 】

まず、取得部 1 2 A が、単位期間 T における、処理回路 1 2 の動作統計情報 4 2 A を取得する (ステップ S 1 0 0)。

【 0 1 0 9 】

次に、導出部 1 2 C が、学習部 1 2 B が学習した予測モデル 2 0 に、ステップ S 1 0 0 で取得した動作統計情報 4 2 A を入力することで、メモリアクセス特性 4 2 B を導出する (ステップ S 1 0 2)。

【 0 1 1 0 】

次に、決定部 1 2 D が、ステップ S 1 0 2 で導出したメモリアクセス特性 4 2 B が第 1 閾値より大きいかなかを判断する (ステップ S 1 0 4)。

【 0 1 1 1 】

メモリアクセス特性 4 2 B が第 1 閾値より大きいと判断した場合 (ステップ S 1 0 4 : Yes)、ステップ S 1 0 6 へ進む。

【 0 1 1 2 】

ステップ S 1 0 6 では、決定部 1 2 D は、第 2 アクセス方式を決定する (ステップ S 1 0 6)。次に、実行部 1 2 E は、ステップ S 1 0 0 で取得した動作統計情報 4 2 A を示す処理を実行中の処理回路 1 2 がアクセス中の、第 2 記憶部 1 4 B のデータを、該第 2 記憶部 1 4 B に配置したまま、該第 2 記憶部 1 4 B にアクセスする (ステップ S 1 0 8)。

【 0 1 1 3 】

次に、処理回路 1 2 は、情報処理を終了するか否かを判断する (ステップ S 1 1 0)。例えば、処理回路 1 2 は、ステップ S 1 0 0 で取得した動作統計情報 4 2 A を示す処理を実行中のアプリケーション 3 2 の終了指示を受付けたかなかを判別することで、ステップ S 1 1 0 の判断を行う。ステップ S 1 1 0 で肯定判断すると (ステップ S 1 1 0 : Yes)、本ルーチンを終了する。一方、ステップ S 1 1 0 で否定判断すると (ステップ S 1 1 0 : No)、上記ステップ S 1 0 0 へ戻る。

【 0 1 1 4 】

一方、上記ステップ S 1 0 4 で、ステップ S 1 0 2 で導出したメモリアクセス特性 4 2 B が第 1 閾値以下であると判断すると (ステップ S 1 0 4 : No)、ステップ S 1 1 2 へ進む。

【 0 1 1 5 】

10

20

30

40

50

ステップ S 1 1 2 では、決定部 1 2 D は、第 1 アクセス方式を決定する（ステップ S 1 1 2）。次に、実行部 1 2 E は、ステップ S 1 0 0 で取得した動作統計情報 4 2 A によって示される動作実行時に処理回路 1 2 がアクセスしていた、第 2 記憶部 1 4 B の第 1 領域内（ページ内）のデータを、該第 2 記憶部 1 4 B から第 1 記憶部 1 4 A へ転送する（ステップ S 1 1 4）。

【 0 1 1 6 】

次に、実行部 1 2 E は、ページテーブルにおける、ステップ S 1 1 4 で転送した第 1 領域の論理アドレスに対応する物理アドレスを、ステップ S 1 1 4 で転送した転送先の第 1 記憶部 1 4 A の格納先を示す物理アドレスに更新する（ステップ S 1 1 6）。このため、処理回路 1 2 は、該データにアクセスする場合には、第 1 記憶部 1 4 A にアクセスすることで、該データにアクセスすることが可能となる。

10

【 0 1 1 7 】

そして、実行部 1 2 E は、ステップ S 1 1 4 で第 1 記憶部 1 4 A へ転送されたデータにアクセスする（ステップ S 1 1 8）。そして、上記ステップ S 1 1 0 へ進む。

【 0 1 1 8 】

以上説明したように、本実施の形態の情報処理装置 1 0 は、取得部 1 2 A と、導出部 1 2 C と、決定部 1 2 D と、を備える。取得部 1 2 A は、処理回路 1 2 の動作統計情報 4 2 A を取得する。導出部 1 2 C は、動作統計情報 4 2 A から処理回路 1 2 のメモリアクセス特性 4 2 B を導出するための予測モデル 2 0 に基づいて、取得した動作統計情報 4 2 A からメモリアクセス特性 4 2 B を導出する。決定部 1 2 D は、導出したメモリアクセス特性 4 2 B に基づいて、第 1 記憶部 1 4 A より処理回路 1 2 によるアクセス速度が遅い第 2 記憶部 1 4 B のデータを第 1 記憶部 1 4 A へ転送し、第 1 記憶部 1 4 A 内の該データにアクセスする第 1 アクセス方式、または、第 2 記憶部 1 4 B 内のデータにアクセスする第 2 アクセス方式、の何れかのアクセス方式を決定する。

20

【 0 1 1 9 】

このように、本実施の形態の情報処理装置 1 0 は、予測モデル 2 0 に基づいて、第 1 アクセス方式または第 1 アクセス方式の何れかのアクセス方式を決定する。

【 0 1 2 0 】

ここで、従来では、複数種類のメモリ（記憶部 1 4）の使い分けに用いる情報を効率よく提供することは困難であった。

30

【 0 1 2 1 】

S C M などの第 2 記憶部 1 4 B は、D R A M などの第 1 記憶部 1 4 A より大容量であるがアクセス速度が遅い。このため、処理対象のデータの特性に合わせて、データを第 1 記憶部 1 4 A と第 2 記憶部 1 4 B とに分散して格納してアクセスすると、処理回路 1 2 は効率よくデータ処理を実行することができる。すなわち、処理回路 1 2 によるアクセスのローカリティが低く、メモリアクセスされる場所が広域に分散されることでアクセスするデータサイズが大きい場合には、データを第 2 記憶部 1 4 B に配置したままとし、処理回路 1 2 が第 2 記憶部 1 4 B にダイレクトにアクセスすることが好ましい。また、メモリアクセスされる場所が集中することで、アクセスのローカリティの高いデータに処理回路 1 2 がアクセスする場合には、データを第 2 記憶部 1 4 B から第 1 記憶部 1 4 A へ転送（コピー）し、処理回路 1 2 は第 1 記憶部 1 4 A 上のデータをキャッシュライン単位でアクセスすることが好ましい。

40

【 0 1 2 2 】

しかし、従来では、複数種類の記憶部 1 4 の使い分けに用いる情報、すなわち、処理回路 1 2 のメモリアクセス特性 4 2 B を、効率よく得ることが困難であった。

【 0 1 2 3 】

一方、本実施の形態の情報処理装置 1 0 では、予測モデル 2 0 に基づいて、取得した動作統計情報 4 2 A からメモリアクセス特性 4 2 B を導出し、第 1 アクセス方式または第 1 アクセス方式の何れかのアクセス方式を決定する。

【 0 1 2 4 】

50

従って、本実施の形態の情報処理装置 10 は、複数種類のメモリの使い分けに用いる情報を効率よく提供することができる。

【0125】

(変形例 1)

なお、上記実施の形態では、情報処理の手順の説明時に、変更部 12F が変更処理を実行しない形態を、一例として示した。

【0126】

しかし、情報処理の手順の実行時に、変更部 12F による変更処理を実行してもよい。

【0127】

この場合、例えば、図 9 に示すステップ S 112 の第アクセス方式を決定した後に、変更部 12F が、第 1 記憶部 14A の利用可能なメモリサイズを変更する変更処理を実行すればよい。そして、その次に、上記ステップ S 114 ~ ステップ S 118 の処理を実行すればよい。

10

【0128】

なお、変更部 12F による変更処理のタイミングは、このタイミングに限定されない。例えば、変更部 12F は、第 1 アクセス方式が決定され、第 2 記憶部 14B の第 1 領域内のデータが第 1 記憶部 14A へ転送された後に、変更処理を実行してもよい。また、変更部 12F は、第 1 アクセス方式が決定され、第 2 記憶部 14B の第 1 領域内のデータが第 1 記憶部 14A へ転送され、更に、実行部 12E が第 1 記憶部 14A のデータにアクセスした後に、変更処理を実行してもよい。

20

【0129】

(変形例 2)

なお、上記実施の形態では、決定部 12D は、導出部 12C が導出したメモリアクセス特性 42B が第 1 閾値より大きい場合、第 2 アクセス方式を決定する形態を説明した。また、決定部 12D は、該メモリアクセス特性 42B が第 1 閾値以下の場合、第 1 アクセス方式を決定する形態を説明した。

【0130】

しかし、決定部 12D は、他の方法により、第 1 アクセス方式または第 2 アクセス方式を決定してもよい。

【0131】

例えば、決定部 12D は、メモリアクセス特性 42B の比率が、第 2 閾値より大きい場合、第 2 アクセス方式を決定する。メモリアクセス特性 42B の比率とは、取得部 12A が取得した動作統計情報 42A によって示される動作実行時に実行中の、1 または複数のアプリケーション 32 の各々に割当てられた、物理メモリサイズの合計値に対する、メモリアクセス特性 42B の比率（割合）を示す。

30

【0132】

具体的には、決定部 12D は、実行中のアプリケーション 32 の各々に割当てられた物理メモリサイズの合計値に対する、導出部 12C で導出されたメモリアクセス特性 42B としての使用中のメモリサイズの比率が、第 2 閾値以下の場合、第 1 アクセス方式を決定する。

40

【0133】

該比率が第 2 閾値以下の状態とは、アプリケーション 32 で利用する可能性のあるメモリの一部の領域に対して、メモリアクセスが集中している状態を示す。このため、この場合、決定部 12D は、処理回路 12 がアクセスのローカリティの高いデータにアクセスしていると判断し、第 2 記憶部 14B 上のデータをページ単位（第 1 領域単位）で第 1 記憶部 14A へ転送し、アクセスする第 1 アクセス方式を決定する。

【0134】

一方、決定部 12D は、該比率が第 2 閾値を超える場合、第 2 アクセス方式を決定する。

【0135】

50

該比率が第2閾値を超える状態とは、アプリケーション32で利用する可能性のあるメモリ領域全体に対する処理回路12によるアクセスが、該メモリ領域全体に分散されている状況であることを示す。このため、該比率が第2閾値を超える状態の場合、アクセスのローカリティが低く、使用中のメモリサイズが大きい状態である。このため、この場合、決定部12Dは、第2アクセス方式を決定する。

【0136】

なお、第2閾値は、予め任意の値を定めればよい。例えば、第2閾値は、上記1または複数のアプリケーション32の各々に割当てられた物理メモリサイズの合計値のN分の1（Nは、2以上の整数）。例えば、第2閾値は、上記合計値の1/3、1/5、1/7、1/10などの値などである。基本的な考え方としては、利用するDRAMに対してそれに見合う高速化が得られる場合は第1アクセス方式を選択したい。あるアプリケーションに対して、そのアプリケーションに割り当てられた物理メモリサイズを10としたときに、DRAMを1あるいは2あるいは3程度利用することでDRAM利用による消費電力増加やコストを抑えつつDRAM利用による局所性の高い処理の高速化が得られるのが好ましいためである。

10

【0137】

次に、決定部12Dが比率を用いてアクセス方式を決定する場合に、処理回路12が実行する情報処理の手順の一例を説明する。

【0138】

図10は、本変形例の処理回路12が実行する情報処理の手順の一例を示す、フローチャートである。

20

【0139】

まず、処理回路12は、上記実施の形態のステップS100～ステップS102（図9参照）と同様にして、ステップS200～ステップS202の処理を実行する。

【0140】

詳細には、取得部12Aが、単位期間Tにおける、処理回路12の動作統計情報42Aを取得する（ステップS200）。次に、導出部12Cが、学習部12Bが学習した予測モデル20に、ステップS200で取得した動作統計情報42Aを入力することで、メモリアクセス特性42Bを導出する（ステップS202）。

【0141】

次に、決定部12Dが、ステップS202で導出したメモリアクセス特性42Bの比率が、第2閾値より大きいか否かを判断する（ステップS204）。

30

【0142】

メモリアクセス特性42Bの比率が第2閾値より大きいと判断した場合（ステップS204：Yes）、上記実施の形態のステップS106～ステップS110（図9参照）と同様にして、ステップS206～ステップS208～ステップS210の処理を実行する。

【0143】

一方、メモリアクセス特性42Bの比率が第2閾値以下と判断した場合（ステップS204：No）、上記実施の形態のステップS112～ステップS118（図9参照）と同様にして、ステップS212～ステップS218の処理を実行する。そして、ステップS210で肯定判断すると（ステップS210：Yes）、本ルーチンを終了する。

40

【0144】

以上説明したように、決定部12Dは、導出部12Cで導出したメモリアクセス特性42Bの比率が、第2閾値より大きいか否かを判断することで、アクセス方式を決定してもよい。この場合についても、上記実施の形態と同様の効果が得られる。

【0145】

なお、上記実施の形態および変形例では、情報処理装置10が、処理回路12と、キャッシュメモリ16と、管理装置18と、を備える形態を一例として説明した（図1参照）。しかし、情報処理装置10は、処理回路12と、キャッシュメモリ16と、管理装置1

50

8と、記憶部14と、を備えた構成であってもよい。また、処理回路12が、キャッシュメモリ16および管理装置18の少なくとも一方を含む構成であってもよい。また、管理装置18が、記憶部14およびキャッシュメモリ16を備えた構成であってもよい。

【0146】

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

10

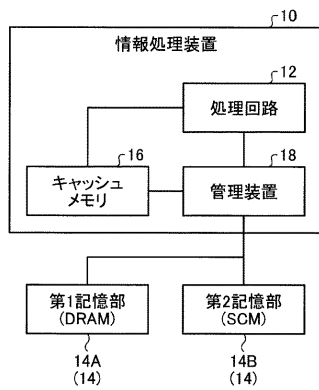
【符号の説明】

【0147】

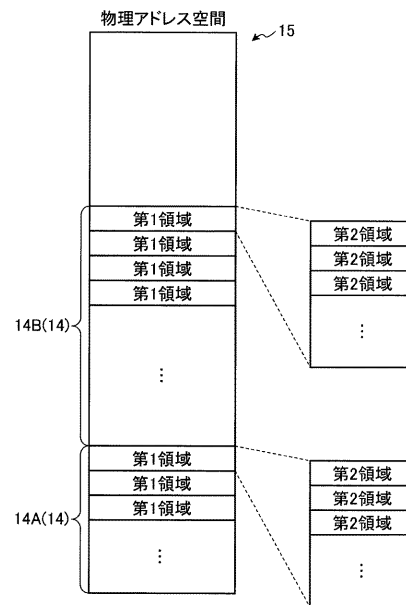
- 10 情報処理装置
- 12 処理回路
- 12A 取得部
- 12B 学習部
- 12C 導出部
- 12D 決定部
- 12E 実行部
- 12F 変更部
- 14 記憶部
- 14A 第1記憶部
- 14B 第2記憶部
- 20 予測モデル

20

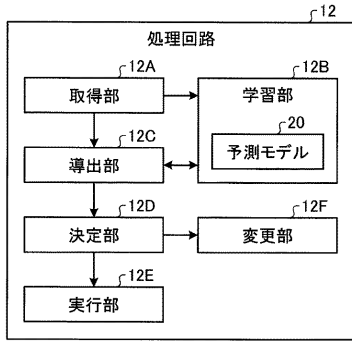
【図1】



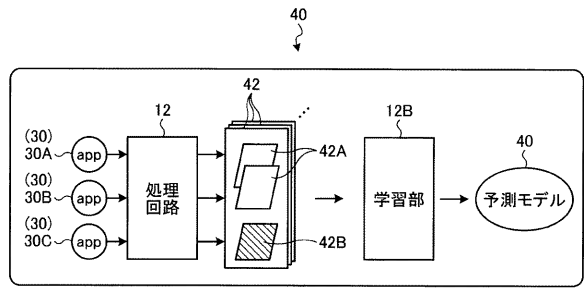
【図2】



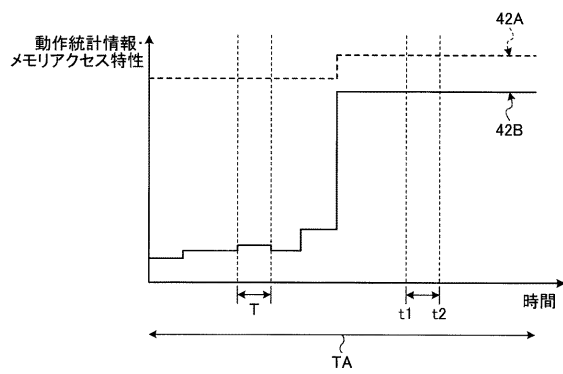
【図 3】



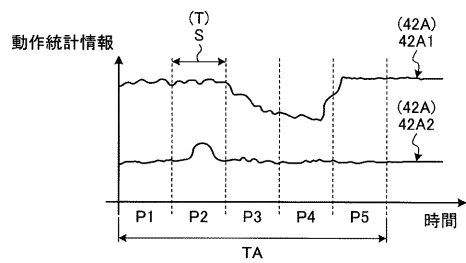
【図 4】



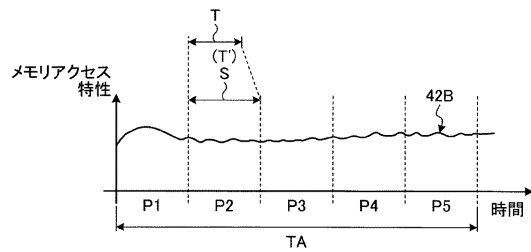
【図 5】



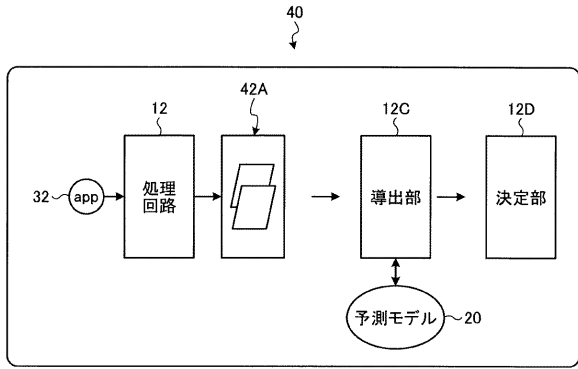
【図 6 A】



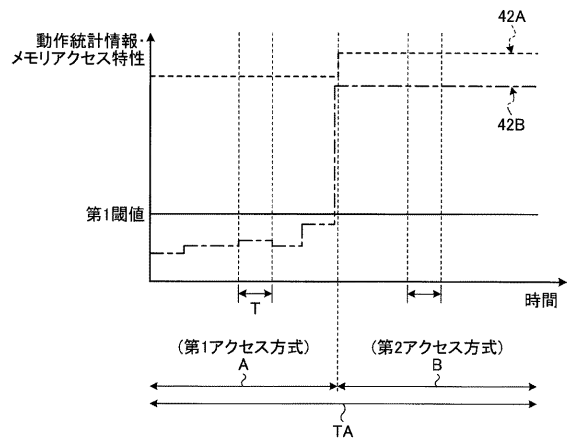
【図 6 B】



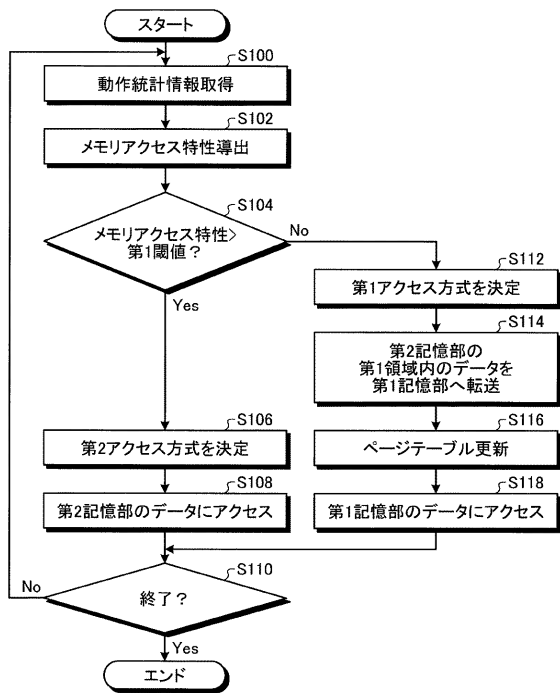
【図7】



【図8】



【図9】



【図10】

