

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6064756号
(P6064756)

(45) 発行日 平成29年1月25日 (2017. 1. 25)

(24) 登録日 平成29年1月6日 (2017. 1. 6)

(51) Int. Cl.

F I

G 0 6 F 11/34 (2006.01)

G 0 6 F 11/34 1 0 9

G 0 6 F 11/34 1 3 3

請求項の数 6 (全 18 頁)

(21) 出願番号 特願2013-82397 (P2013-82397)
 (22) 出願日 平成25年4月10日 (2013. 4. 10)
 (65) 公開番号 特開2014-206786 (P2014-206786A)
 (43) 公開日 平成26年10月30日 (2014. 10. 30)
 審査請求日 平成28年1月13日 (2016. 1. 13)

(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 100079049
 弁理士 中島 淳
 (74) 代理人 100084995
 弁理士 加藤 和詳
 (74) 代理人 100099025
 弁理士 福田 浩志
 (72) 発明者 松尾 美由紀
 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

最終頁に続く

(54) 【発明の名称】 性能データ収集プログラム、装置、及び方法

(57) 【特許請求の範囲】

【請求項 1】

演算処理装置と転送制御部とを有する情報処理装置に、

所定の関数で記述された解析対象のプログラムを実行する前記演算処理装置が呼び出した関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納させ、

採取した性能データに基づいて、前記演算処理装置の処理負荷を判定させ、

前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送させる

ことを特徴とする性能データ収集プログラム。

【請求項 2】

前記演算処理装置の処理負荷の判定は、

前記情報処理装置に、

採取した性能データに基づいて、前記演算処理装置がアイドル状態かを判定させることを特徴とする請求項 1 記載の性能データ収集プログラム。

【請求項 3】

前記演算処理装置の処理負荷の判定は、

前記情報処理装置に、

所定期間に採取した複数の性能データの各々に含まれる情報が示す関数におけるアイド

ル関数の割合が所定割合以上の場合に、前記演算処理装置がアイドル状態であると判定させることを特徴とする請求項 2 記載の性能データ収集プログラム。

【請求項 4】

前記主記憶部以外の記憶部への転送は、
前記情報処理装置に、

前記主記憶部に格納した性能データが所定量を超えた場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さず、前記主記憶部以外の記憶部へ転送させることを特徴とする請求項 1 ~ 請求項 3 のいずれか 1 項記載の性能データ収集プログラム。

【請求項 5】

所定の関数で記述された解析対象のプログラムを実行する演算処理装置と、
前記演算処理装置が呼び出した関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納する採取格納部と、
採取された性能データに基づいて、前記演算処理装置の処理負荷を判定する判定部と、
前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納された性能データの少なくとも一部を、前記演算処理装置を介さず、前記主記憶部以外の記憶部へ転送する転送制御部と、
を含む性能データ収集装置。

【請求項 6】

演算処理装置と転送制御部とを有する情報処理装置が、
所定の関数で記述された解析対象のプログラムを実行する前記演算処理装置が呼び出した関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納し、
採取した性能データに基づいて、前記演算処理装置の処理負荷を判定し、
前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さず、前記主記憶部以外の記憶部へ転送することを特徴とする性能データ収集方法。

【発明の詳細な説明】

【技術分野】

【0001】

開示の技術は、性能データ収集プログラム、装置、及び方法に関する。

【背景技術】

【0002】

コンピュータで動作するプログラムにレスポンス低下などの性能劣化の問題が発生した際の原因解明や、プログラムの性能向上のためのチューニング箇所を特定するために、プログラムの性能を解析することが行われている。プログラムの性能解析は、プログラム実行時に収集した性能データなどの各種情報に基づいて行われる。性能データの収集の手法として、解析対象のプログラムの実行中に、動作したプロセスや呼び出された関数等を示す性能データをサンプリングにより採取し、複数の性能データを収集する手法が存在する。

【0003】

例えば、所定のサンプリング間隔で CPU (Central Processing Unit) が実行しているアドレスを取得し、取得したアドレスを取得時刻を関連づけ、アドレスデータとしてメインメモリ上のアドレスデータ記憶部に時系列に格納する技術が提案されている。この技術では、サンプリング期間終了後に、アドレスデータ記憶部に格納されたアドレスデータを、HDD (Hard Disk Drive) 上のファイルに書き出し、一括して解析処理を行っている。

【0004】

また、サンプリングにより採取したデータを、カーネル空間の一時バッファ領域を介し

10

20

30

40

50

てユーザ空間にコピーし、コピーしたデータに対して、デーモン（バックグラウンド処理）で解析処理を行い、解析結果をHDD上のファイルに書き出す技術が存在する。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2007-213205号公報

【非特許文献】

【0006】

【非特許文献1】“openSUSE 12.3 システム分析とチューニングガイド”、[online]、[平成25年3月8日検索]、インターネット<URL: <http://opensuse-man-ja.berlios.de/opensuse-html/cha.tuning.oprofile.html>>

10

【発明の概要】

【発明が解決しようとする課題】

【0007】

1台のコンピュータに複数のプロセッサを搭載しているコンピュータで実行されるプログラムを解析対象のプログラムとする場合には、コンピュータに搭載されているプロセッサ毎に性能データを収集する必要がある。また、近年、プロセッサのマルチコア化に伴い、1つのLSI（Large Scale Integrated circuit）パッケージ内に組み込まれるコア数が増加している。また、1つのコアが複数のスレッドをそれぞれ実行するマルチスレッド化も進んでいる。このようにマルチコア化されたプロセッサで実行されるプログラムや、マルチスレッドを利用したプログラムを解析対象とする場合には、コア毎またはスレッド毎に性能データを収集する必要がある。

20

【0008】

また、HPC（High Performance Computing）などの科学技術計算向けアプリケーションプログラムでは、プログラムの実行時間が長いため、1回に収集すべき性能データのデータ量が増大する。さらに、金融などのミッションクリティカルなシステムで使用されるアプリケーションプログラムでは、高速なレスポンスが要求される。例えば、株式売買システムにおける注文処理で、1/1000秒（1ms）以下のレスポンスタイムを実現するシステムが存在する。なお、注文処理とは、例えば、ユーザ端末からの注文を受け付け、注文を受け付けたことを確認し、注文内容をチェックしてからサーバへ登録し、登録した旨をユーザ端末へ通知する一連の処理である。このような注文処理のシステムを実現するアプリケーションプログラムの性能解析を行う場合には、システムの挙動を捉えるために、より細かいサンプリング間隔で性能データを採取して、サンプリング期間分の性能データを収集する必要がある。

30

【0009】

収集する性能データのデータ量は、例えば下記のように計算することができる。

・収集する性能データのデータ量 = 1性能データのサイズ × 1システム当たりのプロセッサ数 × 1プロセッサ当たりのコア数 × 1コア当たりのスレッド数 × 1スレッド当たりの収集する性能データの個数

・なお、収集する性能データの個数 = サンプリング期間 ÷ サンプリング間隔

40

従って、マルチプロセッサ化、マルチコア化、マルチスレッド化、またはシステムの高速化により、収集する性能データのデータ量は増大する。このように、プログラムの性能解析では、大量の性能データを収集することが要求される。

【0010】

しかし、サンプリング期間終了後に一括して解析処理を行う従来技術では、収集可能なデータ量はメモリサイズに依存するが、メモリサイズを大きくするには限界がある、という問題がある。

【0011】

また、性能データのサンプリング処理を行いながらデーモンで解析処理を行う従来技術では、カーネル空間の一時バッファ領域に保持した解析結果をHDD上のファイルに書き

50

出した後は、カーネル空間の一時バッファ領域を再利用することができる。しかし、同一のコンピュータが性能データのサンプリングを行いながら転送指示を行うことにより、コンピュータの処理負荷が高い場合に転送指示を行っているため、解析対象のプログラムの動作に影響を与えてしまい、収集した性能データの解析時において性能データに影響を与えた要因の切り分けが困難であった、という問題がある。

【 0 0 1 2 】

開示の技術は、一つの側面として、収集した性能データの解析時において、性能データに影響を与えた要因の切り分けを容易にすることが目的である。

【課題を解決するための手段】

【 0 0 1 3 】

10

開示の技術は、演算処理装置と転送制御部とを有する情報処理装置に、所定の関数で記述された解析対象のプログラムを実行する演算処理装置が呼び出した関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納させる。また、開示の技術は、前記情報処理装置に、採取した性能データに基づいて、前記演算処理装置の処理負荷を判定させる。また、開示の技術は、前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送させる。

【発明の効果】

【 0 0 1 4 】

開示の技術は、一つの側面として、演算処理装置の処理負荷が所定値以下の場合に、転送指示を行うため、収集した性能データの解析時において性能データに影響を与えた要因の切り分けが容易になる、という効果を有する。

20

【図面の簡単な説明】

【 0 0 1 5 】

【図 1】本実施形態に係る性能データ収集装置の構成の一例を示すブロック図である。

【図 2】性能データの一例を示す図である。

【図 3】シンボルテーブル (S y s t e m . m a p ファイル) の一部の一例を示す図である。

【図 4】性能データ格納部のアドレス管理を説明するための図である。

【図 5】性能データ格納部のアドレス管理を説明するための図である。

30

【図 6】性能データ格納部のアドレス管理を説明するための図である。

【図 7】性能データ格納部のアドレス管理を説明するための図である。

【図 8】性能データ収集装置として機能するコンピュータの一例を示す概略ブロック図である。

【図 9】解析対象プログラムの実行を示すフローチャートである。

【図 10】性能データ収集処理を示すフローチャートである。

【図 11】アイドル状態判定処理を示すフローチャートである。

【図 12】転送制御処理を示すフローチャートである。

【図 13】性能データ収集装置の他の構成例を示すブロック図である。

【発明を実施するための形態】

40

【 0 0 1 6 】

以下、図面を参照して開示の技術の実施形態の一例を詳細に説明する。

【 0 0 1 7 】

図 1 に、本実施形態の概略を示す。図 1 に示すように、本実施形態に係る性能データ収集装置 10 は、C P U (Central Processing Unit) 12、メモリ 14、及び H C A (Host Channel Adapter) 16 を含んでいる。性能データ収集装置 10 は、コンピュータ等の情報処理装置である。C P U 12 は、解析対象のプログラム及び後述する性能データ収集プログラムを実行する演算処理装置である。メモリ 14 は、例えば R A M (random access memory) などの主記憶部である。H C A 16 は、例えば I n f i n i B a n d (登録商標) 等のインタコネクトを利用した R D M A (Remote Direct Memory Access) を行うた

50

めの通信部としてのインターフェースカードである。

【0018】

また、性能データ収集装置10は、例えばInfiniBand（登録商標）等のインタコネクトを利用した通信を行うことができる通信ケーブル48を介して、性能データ収集装置10とは異なるコンピュータであるリモート装置40と接続されている。リモート装置40は、性能データ収集装置10と同様の構成とすることができ、少なくとも主記憶部であるメモリ44及び通信部であるHCA46を備えている。性能データ収集装置10は、メモリ14に格納されたデータを、CPU12を介することなく、HCA16、通信ケーブル48、及びリモート装置40のHCA46を介して、リモート装置40のメモリ44へ直接転送するRDMA転送を行う。

10

【0019】

また、図1には、性能データ収集装置10の各機能を示す機能ブロックも合わせて示している。性能データ収集装置10は、CPU12の各機能として、採取格納部32、判定部34、及び転送制御部36を備えている。また、メモリ14上に性能データを格納する領域である性能データ格納部38を備えている。

【0020】

採取格納部32は、複数の関数で記述された解析対象のプログラムを実行中のCPU12が呼び出した関数を示す情報を含む性能データを、設定されたサンプリング期間（サンプリングの開始から終了までの期間）において、設定されたサンプリング間隔で採取する。採取格納部32は、例えば図2に示すような性能データを採取する。図2の例では、各性能データは、CPUID、PID、及び実行アドレスを含んでいる。CPUIDは、解析対象のプログラムを実行中のCPUを識別するための識別番号である。PIDは、CPU12が実行中のプロセス（関数を含むプログラムの実行単位）を識別するための識別番号である。実行アドレスは、CPU12が実行中の関数が格納されたメモリ14上のアドレスである。また、性能データに、採取時刻を示すタイムスタンプを付加してもよい。

20

【0021】

例えばLinux（登録商標）の場合には、解析対象のプログラムに関するプロセスを表現するための情報を含む構造体がカーネル空間に生成される。従って、採取格納部32は、カーネル空間に生成された構造体から性能データとして必要な情報を採取することができる。構造体に含まれる情報は、CPU12内のレジスタに格納された値などに基づいて、カーネルにより設定される。より具体的には、解析対象のプログラムを実行中のCPUのCPUIDを返す関数（例えば、`smp_processor_id`）を呼び出して実行することにより、CPUIDを採取することができる。また、CPUで実行中のプロセスを示すシンボル（`current`）を用いて、PIDを返す関数（例えば、`current->pid`）を呼び出して実行することにより、実行中のプロセスのPIDを採取することができる。また、実行中のプロセスの状態に関する情報が格納された構造体`pt_regs`を参照して、実行アドレスを返す関数（例えば、`pt_regs->ip`）を呼び出して実行することにより、実行アドレスを採取することができる。

30

【0022】

なお、マルチプロセッサを利用して解析対象のプログラムが実行される場合には、上記のCPUIDによりそれぞれのCPUを識別可能であるため、CPU毎の性能データを採取することができる。前述のCPUIDを返す関数で取得するCPUIDは論理的なCPUを識別する識別番号であり、マルチコア及びマルチスレッドの場合にも異なるCPUIDとなる。このようにして、コア毎及びスレッド毎の性能データを採取することができる。

40

【0023】

また、採取格納部32は、採取した性能データを格納するための領域である性能データ格納部38をメモリ14上に確保し、性能データ格納部38の先頭アドレスを`addr_RS`、末尾アドレスを`addr_RE`で定められるアドレス範囲の記憶領域に記憶する。採取格納部32は、採取した性能データを、性能データ格納部38の空き領域の先頭から

50

順に格納していく。

【 0 0 2 4 】

判定部 3 4 は、性能データ格納部 3 8 に格納された性能データに基づいて、C P U 1 2 の処理負荷を判定する。例えば、判定部 3 4 は、C P U 1 2 がアイドル状態の場合には、処理負荷が所定値以下であり、C P U が性能に影響を与えるような処理を実行していないと判定することができる。C P U 1 2 がアイドル状態か否かの判定は、例えば、所定期間に採取された性能データに含まれる実行アドレスが示す関数におけるアイドル関数の割合に基づいて判定することができる。

【 0 0 2 5 】

より具体的には、所定期間に採取された性能データに含まれる実行アドレスが示す関数におけるアイドル関数の割合が所定値以上の場合に、C P U 1 2 がアイドル状態であると判定することができる。例えば、サンプリング期間または解析対象のプログラムの動作時間に基づいて、アイドル状態を判定するための所定期間を定めておく。例えば、下記のように所定期間を定めることができる。

例 1 : サンプリング期間が、例えば 1 0 0 秒の場合に、サンプリング期間の 1 % である 1 秒

例 2 : サンプリング期間内に動作時間が例えば 1 m s の間、解析対象のプログラムが複数回実行される場合に、例えば 1 0 回分のプログラムの動作時間 1 0 m s

上記のように定めた所定期間をサンプリング間隔で割り、所定期間内に採取される性能データの数、すなわち所定期間内の関数の数を計算しておく。例えばサンプリング間隔を 1 0 0 μ s とすると、所定期間内の関数の数は、上記例 1 の場合は 1 0 0 0 0 個、上記例 2 の場合は 1 0 0 個である。

【 0 0 2 6 】

この所定期間内の関数の数に対するアイドル関数の数の割合が所定値以上か否かを判定する。例えば所定値を 9 5 % とすると、上記例 1 の場合は 9 5 0 0 個以上、上記例 2 の場合は 9 5 個以上がアイドル関数であれば、C P U 1 2 がアイドル状態であると判定する。また、所定期間に採取された性能データに含まれる実行アドレスが示す関数のうち、アイドル関数以外の関数の割合が所定値以下の場合に、C P U 1 2 がアイドル状態であると判定してもよい。例えば所定期間のうち 5 % は他の関数が含まれていてもよいとすると、所定期間内で他の関数が許容される許容数は、上記例 1 の場合は 5 0 0 個以下、上記例 2 の場合は 5 個以下である。

【 0 0 2 7 】

また、所定期間内の関数に対するアイドル関数の数の割合が、例えば 1 0 0 % か否かを判定するようにしてもよい。アドレス関数の数の割合が 1 0 0 % の場合は、実行アドレスがアイドル関数であることを示す性能データが、所定期間連続して採取されたことを示すこととなる。

【 0 0 2 8 】

なお、実行アドレスが示す関数がアイドル関数か否は、関数名を含むシンボルの名前とメモリ 1 4 上のアドレスとの対応関係を示すシンボルテーブルに基づいて判定することができる。例えば L i n u x (登録商標) の場合、アイドル関数はカーネル関数であり、カーネルが使用するシンボルテーブル S y s t e m . m a p ファイルを参照することにより、カーネル関数の名前とアドレスとの対応関係が得られる。図 3 に、S y s t e m . m a p ファイルの一部の一例を示す。図 3 の例では、先頭の項目は各関数のアドレス範囲の先頭アドレス、次の項目はシンボルの型、及び最後の項目は関数名である。C P U 1 2 がアイドル状態か否かを判定するためのアイドル関数は、例えば「p o l l _ i d l e」など、予め関数名で識別することができる。そこで、S y s t e m . m a p を参照して、対象のアイドル関数のアドレス範囲を取得することができる。上記の p o l l _ i d l e 関数を対象のアイドル関数とする場合には、ffffff810148d0 アイドル関数のアドレス範囲 <ffffff810149480 となる。従って、実行アドレスが対象のアイドル関数のアドレス範囲に含まれる場合には、実行アドレスが示す関数はアイドル関数であると判定することが

10

20

30

40

50

できる。

【0029】

転送制御部36は、性能データ格納部38の空き領域のアドレス及び格納された性能データの数を管理する。また、転送制御部36は、判定部34によりCPU12がアイドル状態であると判定され、かつ性能データ格納部38に格納された性能データが所定量を超えた場合に、性能データの転送制御を行う。以下、詳細に説明する。

【0030】

例えば図4に示すように、転送制御部36は、性能データ格納部38内の空き領域の先頭アドレス`addr_1`、末尾アドレス`addr_2`のアドレス範囲の記憶領域に性能データを記憶する。図4に示すように、性能データ格納部38に性能データが格納されていない状態では、性能データ格納部38の先頭アドレス`addr_RS = addr_1`、性能データ格納部38の末尾アドレス`addr_RE = addr_2`である。また、転送制御部36は、空き領域の先頭アドレス`addr_1`と末尾アドレス`addr_2`とを比較して、空き領域のサイズを判定する。`addr_1 > addr_2`の場合には、アドレス`addr_1`からアドレス`addr_RE`までの領域と、アドレス`addr_RS`からアドレス`addr_2`までの領域が空き領域となる。すなわち、性能データ格納部38を、先頭アドレス`addr_RS`と末尾アドレス`addr_RE`とを繋いだ循環したメモリ領域とみなす。

【0031】

採取格納部32により採取された性能データは、性能データ格納部38の空き領域の先頭から順次格納される。転送制御部36は、性能データ格納部38に格納された性能データの個数 k をカウントする。また、図5に示すように、性能データ格納部38に k 番目に格納された性能データの末尾アドレスを`addr_kE`とする。転送制御部36は、性能データ格納部38に k 番目の性能データが格納された際に、空き領域の先頭アドレス`addr_1`を、アドレス`addr_kE`に更新する。

【0032】

また、転送制御部36は、性能データ格納部38に格納された性能データが所定量を超えたか否かを判定する。この判定は、性能データ格納部38に格納された性能データの個数 k が閾値 K ($1 \leq k \leq K$)を超えたか否かにより判定することができる。例えば、 K は一度に転送する性能データの個数とすることができる。転送制御部36は、性能データ格納部38に格納された性能データが所定量を超えたと判定すると、例えば図6に示すように、空き領域の末尾アドレス`addr_2`を、転送する性能データの先頭アドレス`addr_KS`に指定する。また、 K 番目に格納された性能データの末尾アドレスを、転送する性能データの末尾アドレス`addr_K E`に指定する。

【0033】

ここで、一度に転送する性能データの個数は、性能データ格納部38に格納された性能データの全てであってもよいし、一部であってもよい。一度に全ての性能データを転送する場合には、上記の K 番目に格納された性能データの末尾アドレス`addr_K E`は`addr_1`となる。また、一度に転送する性能データのデータ量を、判定部34でCPU12のアイドル状態を判定する際のアイドル関数の連続時間やデータ転送速度などを考慮して決定してもよい。例えばデータ転送速度が2GBps (1秒間に2GB)で、アイドル関数の連続時間を1秒としている場合には、一度に転送する性能データのデータ量を2GBとすることができる。さらに、性能データ格納部38の性能データの格納以外の利用を考慮して、その半分の1GBを一度に転送するデータ量としてもよい。このように定めたデータ量を1性能データ当たりのデータ量で割って、一度に転送する性能データの個数 K を定めておく。

【0034】

転送制御部36は、転送する性能データの先頭アドレス`addr_KS`及び末尾アドレス`addr_K E`と、転送先のリモート装置40のメモリ44を指定する情報を含む転送情報を通信部であるHCA16に設定することにより、HCA16に性能データの転送を

10

20

30

40

50

指示する。性能データの転送を指示されたH C A 1 6は、設定された転送情報に基づいて、C P U 1 2を介することなく、性能データ格納部3 8内の指定されたアドレスに格納された性能データを、指定された転送先のリモート装置4 0のメモリ4 4へR D M A転送する。

【0035】

なお、上記では、性能データ格納部3 8に格納された性能データが所定量を超えたか否かを、性能データ格納部3 8に格納された性能データの数kが、一度に転送する性能データの数Kを超えたか否かにより判定する場合について説明したが、これに限定されない。一度に転送する性能データの個数Kと、性能データ格納部3 8に格納された性能データが所定量を超えたか否かを判定するための閾値とは、それぞれ異なる値を設定してもよい。

10

【0036】

また、図7に示すように、転送制御部3 6は、空き領域の末尾アドレスa d d r __ 2を、転送する性能データの末尾アドレスa d d r __ K Eに更新する。これにより、格納されていた性能データが転送された領域を、再び空き領域として利用して、新たな性能データを格納することができる。

【0037】

性能データ収集装置1 0は、例えば図8に示すように、C P U 1 2、メモリ1 4、及びH C A 1 6に加え、不揮発性の記憶部1 8、及び入出力インターフェース(I / F) 2 0を備えたコンピュータ8 0で実現することができる。C P U 1 2、メモリ1 4、H C A 1 6、記憶部1 8、及び入出力I / F 2 0は、バス2 2を介して互いに接続されている。入出力I / F 2 0には、マウス、キーボード、ディスプレイ等の入出力装置が接続されている。

20

【0038】

記憶部1 8はH D D、S S D (Solid State Drive)、フラッシュメモリ等によって実現できる。記録媒体としての記憶部1 8には、コンピュータ8 0を性能データ収集装置1 0として機能させるための性能データ収集プログラム5 0が記憶されている。また、記憶部1 8には、解析対象プログラム6 0も記憶されている。C P U 1 2は、解析対象プログラム6 0を記憶部1 8から読み出してメモリ1 4に展開し、解析対象プログラム6 0が有するプロセスを順次実行する。また、C P U 1 2は、解析対象プログラムを実行中に、性能データ収集プログラム5 0を記憶部1 8から読み出してメモリ1 4に展開し、性能データ収集プログラム5 0が有するプロセスを順次実行する。解析対象プログラムは特に限定されないため、ここでは詳細な説明を省略する。

30

【0039】

性能データ収集プログラム5 0は、採取格納プロセス5 2、判定プロセス5 4、及び転送制御プロセス5 6を有する。C P U 1 2は、採取格納プロセス5 2を実行することで、図1に示す採取格納部3 2として動作する。また、C P U 1 2は、判定プロセス5 4を実行することで、図1に示す判定部3 4として動作する。また、C P U 1 2は、転送制御プロセス5 6を実行することで、図1に示す転送制御部3 6として動作する。これにより、性能データ収集プログラム5 0を実行したコンピュータ8 0が、性能データ収集装置1 0として機能することになる。

40

【0040】

なお、C P U 1 2により実現される各機能は、例えば半導体集積回路、より詳しくはA S I C (Application Specific Integrated Circuit) 等で実現することも可能である。

【0041】

次に、本実施形態に係る性能データ収集装置1 0の作用について説明する。まず、C P U 1 2が解析対象プログラム6 0の実行を開始する。具体的には、図9に示すように、ステップ1 2 0で、C P U 1 2が、解析対象プログラム6 0を記憶部1 8から読み出してメモリ1 4に展開する。この際、C P U 1 2内のレジスタに格納された値などに基づいて、解析対象のプログラムに関するプロセスを表現するための情報を含む構造体がカーネル空間に生成される。

50

【 0 0 4 2 】

次に、ステップ 1 2 2 で、CPU 1 2 が、プログラムカウンタに格納されたアドレスから関数を取り出して命令レジスタに取り込み、命令レジスタに取り込んだ関数を実行する。この際、CPU 1 2 は、プログラムカウンタの値を、次に実行すべき関数が格納されたメモリ 1 4 上のアドレスに更新するなど、CPU 1 2 内の各レジスタの値を適宜更新する。これにより、カーネル空間に生成された構造体に設定される情報も更新される。

【 0 0 4 3 】

次に、ステップ 1 2 4 で、CPU 1 2 4 は、プログラムを終了するか否かを判定し、終了しない場合には、ステップ 1 2 2 へ戻って、次の関数を実行し、終了すると判定した場合には、解析対象プログラムの処理を終了する。

10

【 0 0 4 4 】

図 9 に示す解析プログラムの実行中に、CPU 1 2 が、図 1 0 に示す性能データ収集処理を実行する。

【 0 0 4 5 】

図 1 0 に示す性能データ収集処理のステップ 1 0 0 で、採取格納部 3 2 が、メモリ 1 4 上に性能データ格納部 3 8 を確保し、性能データ格納部 3 8 の先頭アドレス `addr_RS`、及び末尾アドレス `addr_RE` を所定の記憶領域に記憶する。また、転送制御部 3 6 が、性能データ格納部 3 8 内の空き領域の先頭アドレス `addr_1` を性能データ格納部 3 8 の先頭アドレス `addr_RS` に設定する。また、性能データ格納部 3 8 内の空き領域の末尾アドレス `addr_2` を性能データ格納部 3 8 の末尾アドレス `addr_RE` に設定する。

20

【 0 0 4 6 】

次に、ステップ 1 0 2 で、採取格納部 3 2 が、カーネル空間に生成された構造体から必要な情報を取得することにより、性能データを採取する。ここでは、例えば図 2 に示すように、CPU 1 2 の識別番号である CPU ID、実行中のプロセスの識別番号である PID、及び実行アドレスを含む性能データを採取するものとする。採取格納部 3 2 は、採取した性能データを、性能データ格納部 3 8 の空き領域の先頭アドレスから格納する。

【 0 0 4 7 】

次に、ステップ 1 0 4 で、判定部 3 4 が図 1 1 に示すアイドル状態判定処理を実行する。

30

【 0 0 4 8 】

図 1 1 に示すアイドル状態判定処理のステップ 1 0 4 0 で、判定部 3 4 が、CPU 1 2 が実行中のプロセスが、カーネルプロセスかユーザプロセスかを判定する。後段の処理で判定されるアイドル関数はカーネル関数であるため、実行中のプロセスがユーザプロセスの場合には、実行アドレスを判定するまでもなく、CPU 1 2 がアイドル状態ではないと判定することができる。例えば Linux（登録商標）の場合、性能データに含まれる PID が 0 であれば、実行中のプロセスはカーネルプロセス、PID が 0 以外であればユーザプロセスであると判定することができる。CPU 1 2 により実行中のプロセスがカーネルプロセスの場合には、ステップ 1 0 4 2 へ移行し、ユーザプロセスの場合には、ステップ 1 0 5 2 へ移行する。

40

【 0 0 4 9 】

ステップ 1 0 4 2 では、判定部 3 4 が、カーネル関数の名前とアドレスとの対応関係を定めたシンボルテーブルを参照して、性能データに含まれる実行アドレスが、アイドル関数のアドレス範囲に含まれるか否かを判定する。実行アドレスがアイドル関数のアドレス範囲に含まれる場合には、実行アドレスがアイドル関数であると判定し、ステップ 1 0 4 4 へ移行する。一方、実行アドレスがアイドル関数のアドレス範囲に含まれない場合には、実行アドレスがアイドル関数ではないと判定し、ステップ 1 0 5 2 へ移行する。

【 0 0 5 0 】

次に、1 0 4 4 では、判定部 3 4 が、アイドル関数の連続数を示す変数 `m` を 1 インクリメントする。次に、ステップ 1 0 4 6 で、判定部 3 4 が、変数 `m` が、アイドル状態を判定

50

するためのアイドル関数の連続数の閾値Mを超えたか否かを判定する。 $m > M$ の場合には、ステップ1048へ移行し、判定部34が、CPU12がアイドル状態であると判定する。次に、ステップ1050で、判定部34が、変数m及び変数n（後述）を0に設定（リセット）し、性能データ収集処理へリターンする。一方、上記ステップ1046で、判定部34が、 $m = M$ と判定した場合には、CPU12がアイドル状態であると判定することなく、アイドル関数の連続数のカウントを継続するため、そのまま性能データ収集処理へリターンする。

【0051】

また、上記ステップ1040またはステップ1042が否定判定の場合、すなわち、実行アドレスが示す関数がアイドル関数ではない場合には、ステップ1052で、判定部34が、変数mが0を超えているか否かを判定する。 $m > 0$ の場合には、判定部34は、アイドル関数の連続数のカウントを継続中であると判定して、ステップ1054へ移行する。

【0052】

ステップ1054では、判定部34が、アイドル関数以外の他の関数の出現回数を示す変数nを1インクリメントする。次に、ステップ1056で、判定部34が、変数nが、変数mが閾値Mに達するまでの期間内で他の関数が許容される許容数Nを超えたか否かを判定する。 $n > N$ の場合には、アイドル関数の連続数のカウントをリセットするため、ステップ1050へ移行して、変数m及び変数nを0に設定する。一方、判定部34が、上記ステップ1052で、 $m = 0$ と判定した場合、または上記ステップ1056で、 $n = N$ と判定した場合には、アイドル関数の連続数のカウントを継続するため、そのまま性能データ収集処理へリターンする。

【0053】

図10に示す性能データ収集処理に戻って、ステップ106で、判定部34が、上記ステップ104において、CPU12がアイドル状態であると判定したか否かを判定する。CPU12がアイドル状態であると判定した場合には、ステップ108へ移行し、転送制御部36が、図12に示す転送制御処理を実行する。CPU12がアイドル状態であると判定しなかった場合には、性能データ収集処理を終了する。

【0054】

図12に示す転送制御処理のステップ1080で、転送制御部36が、性能データ格納部38に格納された性能データの個数を示す変数kを1インクリメントする。次に、ステップ1082で、変数kが、一度に転送する性能データの個数として定めた閾値Kを超えたか否かを判定する。 $k > K$ の場合には、ステップ1084へ移行し、転送制御部36が、例えば図6に示すように、空き領域の末尾アドレスaddr_2を、転送する性能データの先頭アドレスaddr_KSに指定する。また、K番目に格納された性能データの末尾アドレスを、転送する性能データの末尾アドレスaddr_KEに指定する。転送制御部36は、アドレスaddr_KS及びaddr_KEと、転送先のリモート装置40のメモリ44を指定する情報とを含む転送情報をHCA16に設定することにより、HCA16に性能データの転送を指示する。

【0055】

次に、ステップ1086で、転送制御部36が、例えば図7に示すように、空き領域の末尾アドレスaddr_2を、転送する性能データの末尾アドレスaddr_KEに更新する。次に、ステップ1088で、転送制御部36が、変数kを、kから転送した性能データの個数Kを差し引いた値に更新して、図10に示す性能データ収集処理へリターンする。

【0056】

一方、上記ステップ1082で、転送制御部36が、 $k = K$ と判定した場合には、まだ、性能データの転送タイミングではないため、ステップ1090へ移行する。ステップ1090では、転送制御部36が、例えば図5に示すように、空き領域の先頭アドレスaddr_1を、性能データ格納部38にk番目に格納された性能データの末尾アドレスad

10

20

30

40

50

d r _ k E に更新して、図 1 0 に示す性能データ収集処理へリターンする。

【 0 0 5 7 】

図 1 0 に示す性能データ収集処理に戻って、ステップ 1 1 0 で、転送制御部 3 6 が、空き領域の先頭アドレス a d d r _ 1 と末尾アドレス a d d r _ 2 とを比較して、空き領域のサイズを求める。転送制御部 3 6 は、求めた空き領域のサイズが所定サイズ以上か否かを判定する。空き領域が所定サイズ以上の場合には、性能データの採取及び格納を継続可能であると判定して、ステップ 1 1 2 へ移行する。

【 0 0 5 8 】

ステップ 1 1 2 では、採取格納部 3 2 が、予め設定されたサンプリング期間を終了したか、または解析対象プログラム 6 0 の実行が終了したか否かを判定することにより、性能データの採取を終了するか否かを判定する。性能データの採取を継続する場合には、ステップ 1 0 2 へ戻って、サンプリング間隔毎に性能データの採取及び格納を繰り返す。

【 0 0 5 9 】

一方、上記ステップ 1 1 0 で、転送制御部 3 6 が、空き領域が所定サイズに満たないと判定した場合、または上記ステップ 1 1 2 で、採取格納部 3 2 が、性能データの採取を終了すると判定した場合には、性能データ収集処理を終了する。これにより、メモリ 1 4 及びリモート装置 4 0 のメモリ 4 4 に、複数の性能データが収集される。

【 0 0 6 0 】

以上説明したように、本実施形態に係る性能データ収集装置 1 0 によれば、C P U がアイドル状態の場合、すなわち C P U が性能に影響を与えるような処理を実行していない場合に、メモリ上の性能データ格納部に格納された性能データをリモート装置へ R D M A 転送する。これにより、C P U の処理負荷が所定値以下の場合に転送指示を行うため、収集した性能データの解析時において、性能データに影響を与えた要因の切り分けが容易になる。そして、解析対象のプログラムの動作に与える影響を抑制して、採取した性能データを転送しない場合に比べて、大量の性能データを収集することができる。

【 0 0 6 1 】

また、C P U がアイドル状態であり、かつメモリに格納された性能データが所定量を超えた場合に、性能データの転送を行うことで、転送回数を低減させることができるため、より解析対象のプログラムの動作に与える影響を抑制することができる。

【 0 0 6 2 】

収集した性能データは、解析対象プログラムの解析に用いることができる。例えば、収集した性能データを集計し、性能ボトルネックとなる呼び出し回数の多い関数を見つけることができる。また、性能データに採取時刻を示すタイムスタンプを付加した場合には、サンプリング期間全体の統計だけでなく、時系列分析も可能になる。

【 0 0 6 3 】

なお、上記実施形態では、C P U を介さないデータ転送として R D M A 転送を行う場合について説明したが、C P U を介さない、または C P U の負荷が小さい転送方法であれば、例えば D M A (Direct Memory Access) 転送などの他の転送方法を用いてもよい。

【 0 0 6 4 】

図 1 3 に、D M A 転送を行う場合の性能データ収集装置 2 1 0 の一例を示す。性能データ収集装置 2 1 0 は、C P U 1 2、メモリ 1 4、記憶部 1 8、及び D M A C (Direct Memory Access controller) 2 4 を含んでいる。なお、上記実施形態に係る性能データ収集装置 1 0 と同一の部分については同一符号を付している。性能データ収集装置 2 1 0 では、主記憶部であるメモリ 1 4 から、リモート装置 4 0 のメモリ 4 4 ではなく、自装置内の補助記憶部である記憶部 1 8 へデータが転送される。転送制御部 3 6 は、上記と同様に、転送する性能データのアドレス範囲及び転送先を指定する転送情報を D M A C 2 4 に設定することにより、D M A C 2 4 に性能データの転送を指示する。性能データの転送を指示された D M A C 2 4 は、設定された転送情報に基づいて、C P U 1 2 を介することなく、性能データ格納部 3 8 内の指定されたアドレスに格納された性能データを、記憶部 1 8 へ D M A 転送する。

10

20

30

40

50

【 0 0 6 5 】

また、上記では開示の技術における性能データ収集プログラムの一例である性能データ収集プログラム 5 0 が記憶部 1 8 に予め記憶（インストール）されている態様を説明した。しかし、開示の技術における性能データ収集プログラムは、C D - R O M や D V D - R O M 等の記録媒体に記録されている形態で提供することも可能である。

【 0 0 6 6 】

以上の実施形態に関し、更に以下の付記を開示する。

【 0 0 6 7 】

（付記 1）

演算処理装置と転送制御部とを有する情報処理装置に、

所定の関数で記述された解析対象のプログラムを実行する前記演算処理装置が呼び出した関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納させ、

採取した性能データに基づいて、前記演算処理装置の処理負荷を判定させ、

前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送させる

ことを特徴とする性能データ収集プログラム。

【 0 0 6 8 】

（付記 2）

前記演算処理装置の処理負荷の判定は、

前記情報処理装置に、

採取した性能データに基づいて、前記演算処理装置がアイドル状態かを判定させることを特徴とする付記 1 記載の性能データ収集プログラム。

【 0 0 6 9 】

（付記 3）

前記演算処理装置の処理負荷の判定は、

前記情報処理装置に、

所定期間に採取した複数の性能データの各々に含まれる情報が示す関数におけるアイドル関数の割合が所定割合以上の場合に、前記演算処理装置がアイドル状態であると判定させることを特徴とする付記 2 記載の性能データ収集プログラム。

【 0 0 7 0 】

（付記 4）

前記主記憶部以外の記憶部への転送は、

前記情報処理装置に、

前記主記憶部に格納した性能データが所定量を超えた場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送させることを特徴とする付記 1 ～付記 3 のいずれか 1 つに記載の性能データ収集プログラム。

【 0 0 7 1 】

（付記 5）

所定の関数で記述された解析対象のプログラムを実行する演算処理装置と、

前記演算処理装置が呼び出した関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納する採取格納部と、

採取された性能データに基づいて、前記演算処理装置の処理負荷を判定する判定部と、

前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納された性能データの少なくとも一部を、前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送する転送制御部と、

を含む性能データ収集装置。

【 0 0 7 2 】

(付記 6)

前記判定部は、採取された性能データに基づいて、前記演算処理装置がアイドル状態かを判定することを特徴とする付記 5 記載の性能データ収集装置。

【0073】

(付記 7)

前記判定部は、所定期間に採取された複数の性能データの各々に含まれる情報が示す関数におけるアイドル関数の割合が所定割合以上の場合に、前記演算処理装置がアイドル状態であると判定することを特徴とする付記 6 記載の性能データ収集装置。

【0074】

(付記 8)

前記転送制御部は、前記主記憶部に格納された性能データが所定量を超えた場合に、前記主記憶部に格納された性能データの少なくとも一部を、前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送させることを特徴とする付記 5 ～付記 7 のいずれか 1 つに記載の性能データ収集プログラム。

【0075】

(付記 9)

演算処理装置と転送制御部とを有する情報処理装置が、

所定の関数で記述された解析対象のプログラムを実行する前記演算処理装置が呼び出した関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納し、

採取した性能データに基づいて、前記演算処理装置の処理負荷を判定し、

前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送する

ことを特徴とする性能データ収集方法。

【0076】

(付記 10)

前記演算処理装置の処理負荷の判定は、

前記情報処理装置が、

採取した性能データに基づいて、前記演算処理装置がアイドル状態かを判定することを特徴とする付記 9 記載の性能データ収集方法。

【0077】

(付記 11)

前記演算処理装置の処理負荷の判定は、

前記情報処理装置が、

所定期間に採取した複数の性能データの各々に含まれる情報が示す関数におけるアイドル関数の割合が所定割合以上の場合に、前記演算処理装置がアイドル状態であると判定させることを特徴とする付記 10 記載の性能データ収集方法。

【0078】

(付記 12)

前記主記憶部以外の記憶部への転送は、

前記情報処理装置が、

前記主記憶部に格納した性能データが所定量を超えた場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送することを特徴とする付記 9 ～付記 11 のいずれか 1 つに記載の性能データ収集方法。

【0079】

(付記 13)

演算処理装置と転送制御部とを有する情報処理装置に、

所定の関数で記述された解析対象のプログラムを実行する前記演算処理装置が呼び出し

10

20

30

40

50

た関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納させ、

採取した性能データに基づいて、前記演算処理装置の処理負荷を判定させ、

前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送させ、

前記主記憶部及び前記主記憶部以外の記憶部に収集した複数の性能データを集計して、前記解析対象のプログラムを解析させる

ことを特徴とする解析プログラム。

【 0 0 8 0 】

10

(付記 1 4)

所定の関数で記述された解析対象のプログラムを実行する演算処理装置と、

前記演算処理装置が呼び出した関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納する採取格納部と、

採取された性能データに基づいて、前記演算処理装置の処理負荷を判定する判定部と、

前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納された性能データの少なくとも一部を、前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送する転送制御部と、

前記主記憶部及び前記主記憶部以外の記憶部に収集された複数の性能データを集計して、前記解析対象のプログラムを解析する解析部と、

20

を含む解析装置。

【 0 0 8 1 】

(付記 1 5)

演算処理装置と転送制御部とを有する情報処理装置が、

所定の関数で記述された解析対象のプログラムを実行する前記演算処理装置が呼び出した関数を示す情報を含む性能データを、所定のサンプリング間隔で採取して、主記憶部に格納し、

採取した性能データに基づいて、前記演算処理装置の処理負荷を判定し、

前記演算処理装置の処理負荷が所定値以下の場合に、前記主記憶部に格納した性能データの少なくとも一部を、前記転送制御部により前記演算処理装置を介さずに、前記主記憶部以外の記憶部へ転送し、

30

前記主記憶部及び前記主記憶部以外の記憶部に収集した複数の性能データを集計して、前記解析対象のプログラムを解析する

ことを特徴とする解析方法。

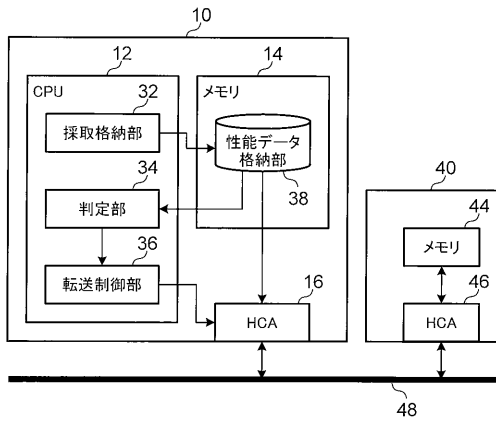
【符号の説明】

【 0 0 8 2 】

- 1 0 性能データ収集装置
- 1 4 メモリ
- 1 6 H C A
- 1 8 記憶部
- 2 4 D M A C
- 3 2 採取格納部
- 3 4 判定部
- 3 6 転送制御部
- 3 8 性能データ格納部
- 4 0 リモート装置
- 4 4 リモート装置のメモリ
- 8 0 コンピュータ

40

【図 1】



【図 2】

CPUID PID 実行アドレス

CPU 1 PID 12980 0x0500

CPU 1 PID 12980 0x1100

CPU 1 PID 12980 0x1300

CPU 1 PID 12980 0x2300

CPU 1 PID 12980 0x1100

CPU 1 PID 12980 0x2300

CPU 1 PID 12980 0x2300

CPU 1 PID 12980 0x2300

⋮

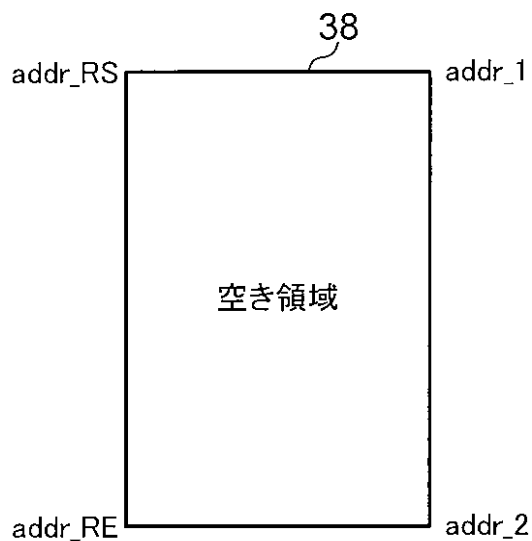
【図 3】

先頭アドレス 型 関数名

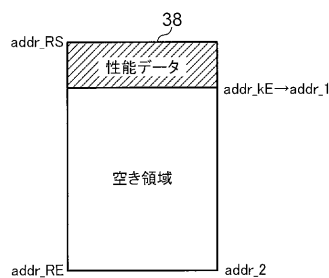
⋮ ⋮ ⋮

poll_idle関数の
アドレス範囲 { ffffffff81014800 t mwait_idle
 fffffff810148d0 t poll_idle
 fffffff81014980 T default_idle
 ⋮

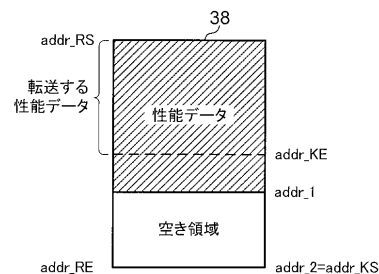
【図 4】



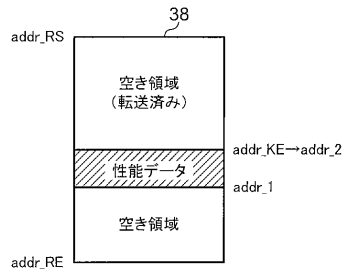
【図 5】



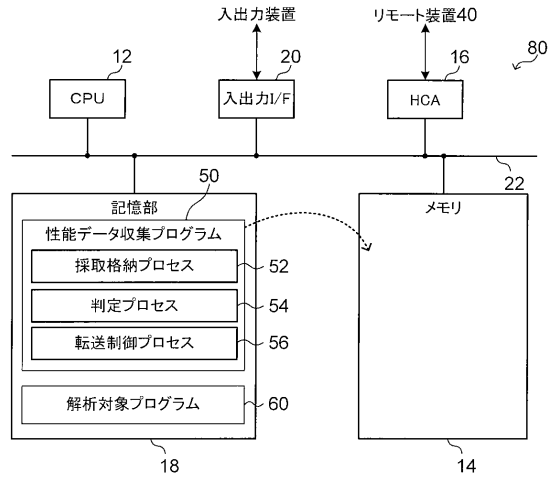
【図 6】



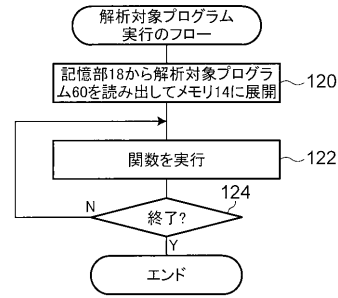
【図 7】



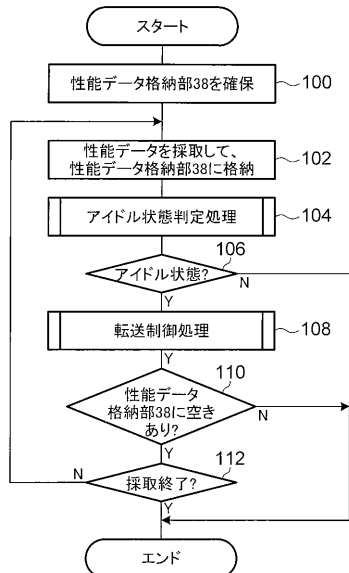
【図 8】



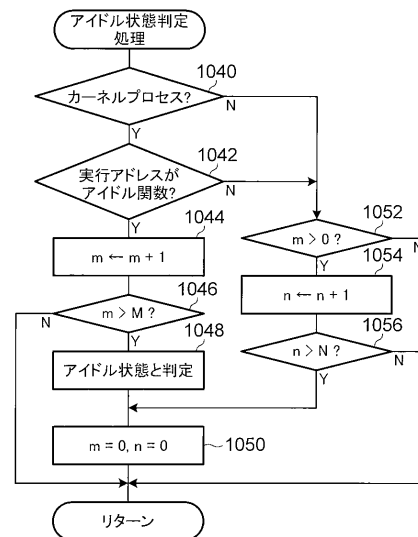
【図 9】



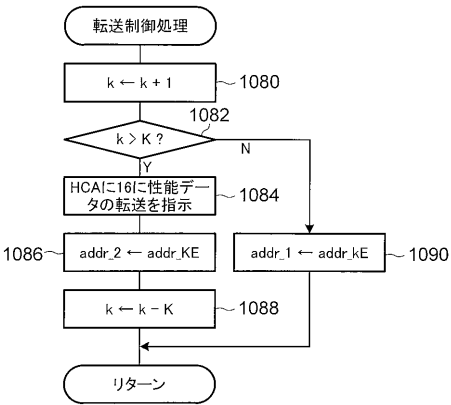
【図 10】



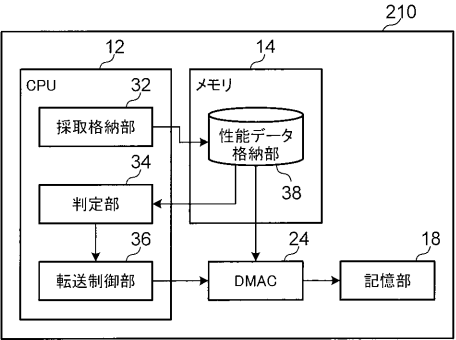
【図 11】



【図 1 2】



【図 1 3】



フロントページの続き

(72)発明者 中島 耕太

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

審査官 多胡 滋

(56)参考文献 特開2005-339107(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 11/34