

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5646644号
(P5646644)

(45) 発行日 平成26年12月24日 (2014. 12. 24)

(24) 登録日 平成26年11月14日 (2014. 11. 14)

(51) Int. Cl.	F I
G06F 9/50 (2006.01)	G06F 9/46 465E
G06F 12/00 (2006.01)	G06F 12/00 550A
G06F 15/177 (2006.01)	G06F 15/177 A

請求項の数 17 (全 16 頁)

(21) 出願番号	特願2012-541131 (P2012-541131)	(73) 特許権者	591016172
(86) (22) 出願日	平成22年11月22日 (2010. 11. 22)		アドバンスト・マイクロ・ディバイス・
(65) 公表番号	特表2013-512509 (P2013-512509A)		インコーポレイテッド
(43) 公表日	平成25年4月11日 (2013. 4. 11)		ADVANCED MICRO DEVI
(86) 国際出願番号	PCT/US2010/057561		CES INCORPORATED
(87) 国際公開番号	W02011/066202		アメリカ合衆国、94088-3453
(87) 国際公開日	平成23年6月3日 (2011. 6. 3)		カリフォルニア州、サニibel、ピィ・
審査請求日	平成25年11月22日 (2013. 11. 22)		オウ・ボックス・3453、ワン・エイ・
(31) 優先権主張番号	12/624, 626		エム・ディ・プレイス、メイル・ストップ
(32) 優先日	平成21年11月24日 (2009. 11. 24)		・68 (番地なし)
(33) 優先権主張国	米国 (US)	(74) 代理人	100108833
早期審査対象出願			弁理士 早川 裕司
		(74) 代理人	100111615
			弁理士 佐野 良太

最終頁に続く

(54) 【発明の名称】 分散型多重コアメモリ初期化

(57) 【特許請求の範囲】

【請求項 1】

メモリ初期化タスクを含むブート処理タスクを複数の処理ノードに行わせる方法であって、

制御処理ノードでメモリ初期化タスクを複数のメモリ初期化サブタスクに分割することと、

対応する処理ノードを各メモリ初期化サブタスクが有するように、前記複数のメモリ初期化サブタスクを前記複数の処理ノードの間で分散させることと、

前記対応する処理ノードで各メモリ初期化サブタスクを実行してサブタスク結果を生成することと、

前記複数の処理ノードからのサブタスク結果を前記制御処理ノードで結合することと、を備え、

前記複数のメモリ初期化サブタスクは、前記複数の処理ノードで並列に又はシーケンスに実行され得るものであって、

前記分散させることは、

シリアルプレゼンスディテクト (SPD) 値をデュアルインラインメモリモジュールから読み出すメモリ初期化サブタスクを、前記複数の処理ノードの1つに割り当てることと、

前記SPD値に依存しないメモリ初期化サブタスクを、前記複数の処理ノードの他の1つに割り当てることとを含む、

方法。

【請求項 2】

前記複数のメモリ初期化サブタスクを分散させるのに先立ち前記複数の処理ノードの間の通信リンクを初期化することを更に備える請求項 1 の方法。

【請求項 3】

前記対応する処理ノードで各メモリ初期化サブタスクを実行することは前記対応する処理ノードでの実行のために各メモリ初期化サブタスクをスケジューリングすることを含む請求項 1 の方法。

【請求項 4】

前記対応する処理ノードで各メモリ初期化サブタスクを実行することはソフトウェアルーチンを実行することを含む請求項 1 の方法。

【請求項 5】

前記対応する処理ノードでの前記メモリ初期化サブタスクの実行に関する状況報告を前記制御処理ノードで受信することを更に備える請求項 1 の方法。

【請求項 6】

前記複数の処理ノードの 1 つは、そのメモリ初期化サブタスクを、前記複数の処理ノードの他の 1 つであって当該他の 1 つに対応するメモリ初期化サブタスクを実行する他の 1 つと並列に実行する請求項 1 の方法。

【請求項 7】

システムメモリと、
複数の処理ノードとを含むコンピュータシステムであって、
前記複数の処理ノードの各々は、
プロセッサコアと、前記複数の処理ノードの他の少なくとも 1 つに接続するための通信インタフェースと、前記システムメモリへの及び前記システムメモリからのデータフローを管理するためのメモリ制御器とを含み、

前記複数の処理ノードは、

マスター処理ノードと、複数の実行処理ノードとを含み、

前記マスター処理ノードは、

メモリ初期化タスクを複数のメモリ初期化サブタスクに分割することと、割り当てられたメモリ初期化サブタスクの各々が対応する実行処理ノードを有するように、前記複数のメモリ初期化サブタスクを前記複数の実行処理ノードに割り当てることとによって、前記メモリ初期化タスクを前記システムメモリに対して実行するように構成されており、

前記複数のメモリ初期化サブタスクは、

前記複数の実行処理ノードの 1 つに割り当てられたメモリ初期化サブタスクが、前記複数の実行処理ノードの他の 1 つに割り当てられた他のメモリ初期化サブタスクに依存することなく順序を問わずに実行可能となるように、前記複数の実行処理ノードで並列に又はシーケンスに実行され得るものであり、

前記割り当てることは、

シリアルプレゼンスディテクト (S P D) 値を前記システムメモリのデュアルインラインメモリモジュールから読み出すメモリ初期化サブタスクを、前記複数の処理ノードの 1 つに割り当てることと、

前記 S P D 値に依存しないメモリ初期化サブタスクを、前記複数の処理ノードの他の 1 つに割り当てることとを含む、

コンピュータシステム。

【請求項 8】

各実行処理ノードは、前記マスター処理ノードによって割り当てられた各メモリ初期化サブタスクを実行するとともに、サブタスク結果を生成するように構成されている請求項 7 のコンピュータシステム。

【請求項 9】

前記マスター処理ノードは、前記複数の実行処理ノードからのサブタスク結果を結合す

10

20

30

40

50

るように構成されている請求項 7 のコンピュータシステム。

【請求項 1 0】

前記システムメモリは、複数のデュアルインラインメモリモジュールを含む請求項 7 のコンピュータシステム。

【請求項 1 1】

前記マスター処理ノードは、前記メモリ初期化サブタスクの状況を前記複数の実行処理ノードから受信するように構成されている請求項 7 のコンピュータシステム。

【請求項 1 2】

前記マスター処理ノードは、メモリ初期化サブタスクを実行する複数の実行処理ノードの 1 つである請求項 7 のコンピュータシステム。

10

【請求項 1 3】

前記複数の実行処理ノードの 1 つは、そのメモリ初期化サブタスクを、前記複数の処理ノードの他の 1 つであって当該他の 1 つに対応するメモリ初期化サブタスクを実行する他の 1 つと並列に実行する請求項 7 のコンピュータシステム。

【請求項 1 4】

前記マスター処理ノードは、前記複数の実行処理ノードで前記複数のメモリ初期化サブタスクを実行した結果に基づいて、前記マスター処理ノードの前記メモリ制御器内の 1 つ以上のレジスタをプログラムすることによって、前記メモリ初期化タスクを実行する請求項 7 のコンピュータシステム。

【請求項 1 5】

20

メモリアレイと、
複数の処理ノードとを含むコンピュータシステムであって、
前記複数の処理ノードの各々は、
プロセッサコアと、前記メモリアレイと通信するためのメモリ制御器と、前記複数の処理ノードの他の少なくとも 1 つに接続するための通信インタフェースとを含み、
前記複数の処理ノードのうち第 1 の処理ノードは、
メモリ初期化タスクを複数のサブタスクに分割し、前記複数のサブタスクを前記複数の処理ノードのうち他の処理ノードに割り当て、
前記複数の処理ノードの各々は、

割り当てられたサブタスクを実行しながら前記メモリアレイの任意の部分にアクセスするとともに、サブタスク実行結果を前記複数の処理ノードのうち前記第 1 の処理ノードに戻すことができ、これにより、前記複数の処理ノードが前記メモリアレイのメモリ初期化を実行することを可能にし、

30

前記複数のサブタスクを前記複数の処理ノードのうち他の処理ノードに割り当てることは、

シリアルプレゼンスディテクト (S P D) 値をデュアルインラインメモリモジュールから読み出すメモリ初期化サブタスクを、前記複数の処理ノードの 1 つに割り当てることと、

前記 S P D 値に依存しないメモリ初期化サブタスクを、前記複数の処理ノードの他の 1 つに割り当てることとを含む、
コンピュータシステム。

40

【請求項 1 6】

前記メモリアレイは、ダブルデータレートバスを介して前記複数の処理ノードに接続された複数のデュアルインラインメモリモジュールを含む請求項 1 5 のコンピュータシステム。

【請求項 1 7】

前記複数の処理ノードのうち前記第 1 の処理ノードは、前記複数の処理ノードで前記複数のサブタスクを実行した結果に基づいて、前記第 1 の処理ノードの前記メモリ制御器内の 1 つ以上のレジスタをプログラムすることによって、前記メモリ初期化タスクを実行し、これにより、前記複数の処理ノードの 1 つのみで前記メモリ初期化タスクを実行するこ

50

とによってブートを行うのに必要であろう時間と比べてブート時間要件を低減する請求項15のコンピュータシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は概してデータ処理システムに関する。1つの態様においては、本発明はシステムブート(system boot)に際してのメモリ初期化のための方法及び装置に関する。

【背景技術】

【0002】

データ処理システム又はコンピューティングシステムは、独立したコンピューティング能力を1人以上のユーザに与えるように設計されており、例えばメインフレーム、ミニコンピュータ、ワークステーション、サーバ、パーソナルコンピュータ、インターネット端末、ノートブック、及び埋め込み型システムを含む多くの形態において見出すことができる。一般的にコンピュータシステムアーキテクチャは、選択されたシステムコンポーネント、関連するメモリ及び制御論理(典型的にはシステムボード上の)並びに入力及び/又は出力(I/O)をシステムに提供する複数の周辺デバイスへの高速高帯域幅(high speed, high bandwidth)のアクセスを1つ以上のマイクロプロセッサコアにもたらすことを含むように設計される。例えば図1は、従来のコンピュータシステム100のための例示的なアーキテクチャを示している。コンピュータシステム100は、「ノース(north)」ブリッジ104を介してシステムメモリ108に接続される1つ以上のプロセッサ102を含む。典型的には、メモリアレイ108は1つ以上のメモリモジュールを含み、またメモリモジュールの追加又は交換のためのメモリスロットを含むこともある。ノースブリッジ回路104は種々のメモリモジュールとインタフェースするようにプログラム可能であり、また図示されるように、ノースブリッジ回路104は複数のメモリモジュール108の間で共有される。その結果、異なる複数のメモリモジュールにデータ投入される(populated)場合、ノースブリッジ回路104は、各メモリモジュールが正確に動作することを可能にするパラメータでプログラムされる必要がある。図示されるノースブリッジ回路104は、高速高帯域幅バス(例えばメモリバス107)を介してメモリ108に接続され、また高速高帯域幅バス(例えばエーリンク(Alink)又はPCIバス)を介して「サウス(south)」ブリッジ112にも接続される。「サウス」ブリッジ112は、1つ以上のI/Oデバイス、例えば、周辺コンポーネント相互接続(Peripheral Component Interconnect)(PCI)バス110(PCIバス110は次いでネットワークインタフェースカード(NIC)120に接続される)、シリアルATAタッチメント(serial ATA Attachment)(SATA)インタフェース114、ユニバーサルシリアルバス(USB)インタフェース116、及びローピンカウント(Low Pin Count)(LPC)バス118(LPCバス118は次いでスーパー入力/出力制御器チップ(Super I/O)122及びBIOSメモリ124に接続される)に接続される。理解されるであろうように、他のバス、デバイス、及び/又はサブシステム、例えば、キャッシュ、モデム、並列又はシリアルインタフェース、SCSIインタフェース等がコンピュータシステム100内に含まれ得ることが分かるはずである。また、ノースブリッジ104及びサウスブリッジ112は、単一チップ又は複数のチップで実装することができ、総称して「チップセット」と称される。代替的には、ブリッジチップの機能の全て又は一部はプロセッサ102内にある場合がある。

【0003】

コンピュータシステムは典型的には基本入力/出力システム(BIOS)と称される組み込み型ソフトウェアルーチンを含み、BIOSは、プログラマ又はユーザがシステムハードウェアとやりとりすることができるように、システムハードウェアとオペレーティングシステムの間でのソフトウェアインタフェースを提供する。BIOS命令は、不揮発性メモリ124(例えばROM(リードオンリメモリ)、PROM(プログラム可能なROM)、EPROM(消去可能なROM)、EEPROM(電氣的に消去可能なROM)、

10

20

30

40

50

フラッシュＲＡＭ（ランダムアクセスメモリ）等）内に記憶されており、そしてメモリを試験すること及び初期化すること、システムをインベントリすること(inventorying)及び初期化すること、並びにシステムを試験することを含む電源投入時の重要なコンピュータシステム機能を制御するために用いられる。電源投入時のこれらの機能は、「システムブート(system boot)」又は「システムをブートすること(booting the system)」と称され、システム電源投入のたびに又はシステムがリセットされるたびに生じ得る。図１に示される従来のコンピュータシステム１００においては、ＢＩＯＳは、コンピューティングシステム１００をブートするために、ブートストラッププロセッサ(boot strap processor)と称されるプロセッサ１０２の１つ上で実行する。動作においてブートストラッププロセッサ１０２は、ノースブリッジ１０４を介してメモリアレイ１０８と通信する。ノースブリッジ回路１０４は、ダブルデータレート（ＤＤＲ）バス等のメモリバス１０７とインタフェースする１つ以上のチャネル制御器に接続されるメモリ制御器（ＭＣ）１０６を含む。ノースブリッジ回路１０４はまた、標準的なバス１０９、例えば周辺コンポーネント相互接続（ＰＣＩ）バス上で１つ以上のサウスブリッジ１１２とも通信する。サウスブリッジ１１２は１つ以上の入力／出力（Ｉ／Ｏ）デバイス１２０，１２２，１２４と通信するが、追加的な又はより少数のデバイス（図示せず）がサウスブリッジ１１２に接続されているかもしれない。

10

【発明の概要】

【発明が解決しようとする課題】

【０００４】

20

システム初期化に際して、ブートストラッププロセッサ１０２は、相当な時間を必要とし得るシステムブートの間に、全てのメモリ試験及びクリーニングを行う。大型サーバシステムは、多くの場合にブートするために数分を必要とし、その間、他のプロセッサはアイドルである。図示される例では、ブートストラッププロセッサ上で実行中のＢＩＯＳは、ＤＤＲメモリバス１０７上でのＤＩＭＭメモリ１０８からの情報に基づいてストラッププロセッサ１０２からレジスタをプログラムすることによって、ノースブリッジ回路１０４内のメモリ制御器１０６を初期化する。メモリ初期化は、メモリモジュールのデータ投入状況(population)を検証することと、メモリの適切な動作（スタックビット無し）を検証することと、メモリを既知の値に初期化し又はクリーニングすることと、を含み得る。大きなメモリ１０８（例えば８、１６、又は３２ギガバイトのメモリ）を伴う従来のシステムは、メモリを初期化するのに数分を必要とすることがあり、特に「ＤＤＲトレーニング」処理は、ブートストラッププロセッサ１０２を用いて、メモリ初期化タスク間の依存を問わずタスクの順序付けられたシーケンスを連続的に実行することによって、メモリを初期化する。この遅延は、大量のメモリ１０８が各制御器１０６に接続される場合に深刻である。

30

【０００５】

そこで、本願発明者によって発見されてきた当該分野における種々の問題に対処する改良されたメモリ初期化デバイス、方法論、及びシステムに対する要求が存在し、ここで、従来の解決法及び技術の種々の限界及び不利益は、後述の図面及び詳細な説明を参照して本願の残りを精査した後に当業者にとって明らかになるはずであるが、背景技術欄のこの説明は、説明される主題が従来技術であるとの了解を意図するものではないことが理解されるべきである。

40

【課題を解決するための手段】

【０００６】

概して、本発明の実施形態は、先行するタスクに依存せずに順序不問で実行が達成され得るようにメモリ初期化タスクをコアに分散させることによって多重コアコンピュータシステム内でメモリを初期化するためのシステム、方法、及び装置を提供する。選択された実施形態においては、複数の処理ノード又はコア及び複数のシステムメモリコンポーネントを含むシステムが提供される。選択された実施形態においては、処理ノード又はコアの１つは、メモリ初期化タスクを多重ＣＰＵノード／コアによって処理される個々のサブタ

50

スクに分割するための制御ノード/コア機能を提供するように選択される。個々のサブタスクは、実行のための他の「スレーブ(slave)」コアに送られ又は分散させられる。各スレーブコアでは、単一又は複数のサブタスクを処理する資源をスレーブコアが有している場合に、サブタスクは実行のためにスケジューリングされる。結果として、所与のスレーブコアでの所与のサブタスクの実行の時間は、他のスレーブコアでの他のサブタスクの実行と重複してよいし又はしなくてもよい。その結果、スレーブコアは、サブタスクを並列に又はシリアルに実行することが可能である。サブタスクが処理された後に、それらは制御コアが元のタスクを完了するために再結合される。これにより、スレーブコアはサブタスクを並列に実行することができるので、同時タスク実行が可能になり、結果として全体の実行時間が減少すると共に最適な性能がもたらされる。

10

【0007】

本発明の種々の実施形態によると、メモリ初期化タスク等のブート処理タスクを複数の処理ノードに行わせるための方法が提供される。方法の例示的な実施形態においては、制御処理ノードは、ブート処理タスク(例えばメモリ初期化タスク)を複数のブート処理サブタスク(例えばメモリ初期化サブタスク)に分割し、次いで対応する処理ノードを各サブタスクが有するように、制御処理ノードを含んでいてよい複数の処理ノードの間でサブタスクを分散させる。サブタスクを分散させるのに先立ち、制御処理ノードは、複数の処理ノードの間の通信リンクを初期化する。サブタスク分散の例を提供すると、処理ノードの1つは、シリアルプレゼンスディテクト(serial presence detect)(SPD)値をDIMMメモリから読み出すメモリ初期化サブタスクを割り当てられる一方で、処理ノードの他の1つは、SPD値に依存しない複雑な初期化タスクを実行するメモリ初期化サブタスクを割り当てられる。各メモリ初期化サブタスクは、実行のためにスケジューリングされ、次いで対応する処理ノードで実行されてサブタスク結果を生成する。選択された実施形態においては、対応する処理ノードでメモリ初期化サブタスクを実行するために、ソフトウェアルーチンが用いられる。制御処理ノードは、対応する処理ノードでのメモリ初期化サブタスクの実行に関する状況報告(status reports)を受信し、そして複数の処理ノードからのサブタスク結果が次いで制御処理ノードで結合され、ここで、複数のメモリ初期化サブタスクは、複数の処理ノードで並列に又はシーケンスに実行されてよい。例えば、処理ノードの1つは、それが割り当てられたメモリ初期化サブタスクを、処理ノードの他の1つであってその対応するメモリ初期化サブタスクを実行する他の1つと並列に実行することができる。

20

30

【0008】

他の実施形態においては、システムメモリ(例えばDIMMベースのレイ)と複数の処理ノードとを備えるコンピュータシステムが提供され、ここで、各処理ノードは、プロセッサコアと、他の単一又は複数のノードとの通信インタフェースと、システムメモリへの及びシステムメモリからのデータフローを管理するためのメモリ制御器と、を含む。処理ノードの1つはマスター処理ノードであり、マスター処理ノードは、メモリ初期化タスクを複数のメモリ初期化サブタスクに分割することと、対応する実行処理ノードを各割り当てられたメモリ初期化サブタスクが有するようにメモリ初期化サブタスクを1つ以上の実行処理ノード(マスター処理ノードを含んでいてよい)に割り当てることと、によってメモリ初期化タスクをシステムメモリ上で行うように構成される。各実行処理ノードは、マスター処理ノードによって割り当てられる各メモリ初期化サブタスクを実行すると共にサブタスク結果を生成するように構成される。マスター処理ノードは、実行処理ノードからメモリ初期化サブタスクの状況を受信するように、また、単一又は複数の実行処理ノードからのサブタスク結果を結合して、単一又は複数の実行処理ノードで複数のメモリ初期化サブタスクを実行した結果に基づいてマスター処理ノードのメモリ制御器内で1つ以上のレジスタをプログラムすることによってメモリ初期化タスクを行うように、構成される。選択された実施形態においては、マスター処理ノードは、実行処理ノードの1つに割り当てられたメモリ初期化サブタスクが、実行処理ノードの他の1つに割り当てられた他のメモリ初期化サブタスクに依存せずに順序不問で実行され得るように、メモリ初期化サブ

40

50

タスクを分割すると共に割り当てる。このようにして、複数のメモリ初期化サブタスクが実行処理ノードで並列に又はシーケンスに実行され得るように、第1の実行処理ノードは、そのメモリ初期化サブタスクを、単一又は複数の実行処理ノードの他の1つであってその対応するメモリ初期化サブタスクを実行している他の1つと並列に実行することができる。

【0009】

更に他の実施形態においては、メモリアレイと複数の処理ノードとを備えるコンピュータシステムが提供される。メモリアレイは、ダブルデータレート（DDR）バスを介して複数の処理ノードに接続される複数のデュアルインラインメモリモジュールとして実装されてよい。各処理ノードは、プロセッサ、メモリアレイと通信するためのメモリ制御器、及び複数の処理ノードの他の少なくとも1つとの通信インタフェース、として実装されてよい。動作において、第1の処理ノードは、メモリ初期化タスクを複数のサブタスクに分割すると共に複数のサブタスクを処理ノードに割り当てる。各処理ノードは、割り当てられたサブタスクを実行しながらメモリアレイの任意の部分からの情報を得ることができる。加えて、各処理ノードは、サブタスク実行結果を第1の処理ノードに戻し、それにより複数の処理ノードがメモリアレイのメモリ初期化を効率的に且つ迅速に行うことを可能にする。このようにして、第1の処理ノードは、そのメモリ制御器内で1つ以上のレジスタを、処理ノードで複数のサブタスクを実行した結果に基づいてプログラムすることによってメモリ初期化タスクを行い、それにより、処理ノードの1つのみでメモリ初期化タスクを実行することによってブートするのに必要であろう時間と比べてブート時間要求を低減する。

【図面の簡単な説明】

【0010】

添付の図面を参照することによって、本発明はより良く理解されるであろうし、またその多くの目的、特徴、及び利益が当業者にとって明らかになるであろう。種々の図面を通して同じ参照番号の使用は同様の又は類似の要素を指定する。

【0011】

【図1】図1は多重コアコンピュータシステムの単純化されたアーキテクチャのブロック図である。

【0012】

【図2】図2は本発明の選択された実施形態に従う分散型多重コアメモリ初期化を伴うコンピューティングシステムアーキテクチャを示す図である。

【0013】

【図3】図3は本発明の選択された実施形態に従う例示的な処理ノードを示す図である。

【0014】

【図4】図4は例示的なDIMM（デュアルインラインメモリモジュール）を示す図である。

【0015】

【図5】図5は本発明の選択された実施形態に従う最適化されたメモリ初期化及び試験のフロー図である。

【発明を実施するための形態】

【0016】

マスターコアの制御の下でメモリ初期化タスクをサブタスクに分割するために、分散型多重コアメモリ試験及び初期化の方法及び装置が提供される。マスターコアは次いで、多重コアによる別個の実行のためにサブタスクを送り又は分散し、完了した場合、結果は、元のタスクを完了させるための再結合のためにマスターコアに送り戻される。メモリ初期化タスクを別個のコアに分散させることによって、サブタスクの実行は、順序を問わず且つ先行するタスクに依存せずに達成され得る。この分散させられた実行は、先行するサブタスクが完了することを待たずにシーケンスを問わず幾つかのタスクが完了することを可能にし、それにより、メモリ初期化時間に対して要求される時間を短くしてデバイス性能

を向上させることができる。

【 0 0 1 7 】

本発明の種々の例示的な実施形態が添付図面を参照して以下に詳細に説明される。以下の説明において種々の詳細が記載されるが、これら特定の詳細なしに本発明が実施され得ること、及びデバイス設計者の特定の目標、例えば実装毎に変わるであろうプロセス技術又は設計関連の制約の遵守、を達成するために種々の実装特有の決定がここで説明される本発明に対してなされ得ること、が理解されるはずである。そのような開発努力は得てして複雑で且つ時間を要するものであろうが、この開示の利益を享受する当業者にとっては日常的な取り組みであるはずである。例えば、本発明を限定し又は曖昧にすることを避けるために、選択された態様が詳細によりはむしろブロック図形態で示される。ここに提供される詳細な説明の幾つかの部分は、コンピュータメモリ内に記憶されるデータに対して動作するアルゴリズム及び命令に関して提示される。そのような説明及び表現は、当業者が彼らの取り組みの内容を他の当業者に説明しまた伝えるために用いられる。概して、アルゴリズムは、所望の結果を結果としてもたらすステップの自己矛盾のないシーケンスを参照し、ここで「ステップ」は、必ずしもそうであるとは限らないが、記憶され、転送され、結合され、比較され、そして他の方法で操作されることが可能な電子的又は磁気的な信号の形態をとり得る物理量の操作を参照する。これらの信号をビット、値、要素、記号、文字、用語(terms)、数字、等として参照することは一般的な慣習である。これら及び類似の用語は、適切な物理量に関連してよく、またこれらの量に適用される便利な指標にすぎない。特に断りのない限り、以下の議論から明らかであるように、「処理すること」又は「コンピューティング」又は「計算すること」又は「決定すること」又は「表示すること」等の用語を用いる議論は、コンピュータシステムのレジスタ及びメモリ内で物理的な(電子的な)量として表されるデータを、コンピュータシステムのメモリ若しくはレジスタ又は他のそのような情報記憶デバイス、伝送デバイス若しくは表示デバイス内で物理量として同様に表される他のデータへと操作し及び変換するコンピュータシステム又は同様の電子的コンピューティングデバイスの動作及び処理を参照することが理解される。

【 0 0 1 8 】

図2を参照すると、本発明の選択された実施形態に従う分散型多重コアメモリ初期化を伴うコンピューティングシステムアーキテクチャ200が示されている。図示されるコンピューティングシステム200においては、リンク203を介して互いに通信する複数の処理ノード202[0:3]がある。各処理ノード202は、例えば、プロセッサコア、メモリ制御器、及びリンクインタフェース回路を含む。リンク203は、例えばハイパートランスポート(HyperTransport)(HT)プロトコル等の分割トランザクション(split-transaction)バスプロトコルに従うデュアル点对点(point to point)リンクであってよい。リンク203は、上流へのデータフロー及び下流へのデータフローを含んでいてよい。リンク信号は、典型的には、クロック、制御、コマンド、アドレス、及びデータの情報等のトラフィック、並びにデバイス間で流れるトラフィックを制限し(qualify)また同期させるリンクサイドバンド信号を含む。処理ノード202の各メモリ制御器は、対応するメモリアレイ206[0:3]と通信する。処理ノード202及びメモリアレイ206はシステムの「コヒーレントな(coherent)」部分内にあり、ここでは全てのメモリトランザクションがコヒーレントである。処理ノード202の1つ以上の内部にはノースブリッジが含まれていてよく、あるいはノースブリッジは、別のHTリンクを介して処理ノード202の1つに結合される別個のノースブリッジデバイス208として設けられていてよい。いずれの場合にも、ノースブリッジデバイスは別のHTリンクを介してサウスブリッジ210に結合されていてよい。加えて、1つ以上のI/Oデバイス212がサウスブリッジ210に結合されていてよい。BIOS_ROM214はサウスブリッジ210に結合されていてよい。ノースブリッジ208、サウスブリッジ210、及びI/Oデバイス212は、システムの「非コヒーレントな」部分内にある。

【 0 0 1 9 】

各メモリアレイ 206 は、メモリモジュールの追加又は交換のための、データ投入された又は未投入の(populated or unpopulated)幾つかのメモリスロットから構成され得る。例えばサーバシステム 200 においては、各メモリスロットは、全体的なサーバシステム 200 が大きなメモリ、例えば 32 ギガバイト (G B y t e s) の記憶容量を有するように、512 メガバイト (M b y t e s) の記憶容量を提供することができる。各処理ノード 202 のメモリ制御器は、異なってプログラムされ得るが、関連する処理ノード 202 に結合されるメモリモジュール 206 のローカル種とインタフェースするようにプログラムされる必要がある。システム 200 は図示されるよりも複雑であり得ることが理解されるはずである。例えば、システムのコヒーレントな部分には追加的な処理ノード 202 があってよい。また、処理ノード 202 は梯子型アーキテクチャで図示されているが、処理ノード 202 は種々の方法で相互接続されてよく、また追加的な単一又は複数の H T リンクを介するより複雑な相互間の結合を有していてもよい。

10

【0020】

図 3 は本発明の選択された実施形態に従う例示的な処理ノード 202 を示している。図示されるように、処理ノード 202 は、プロセッサコア 302、多重 H T リンクインタフェース 304、及びメモリ制御器 306 を含む。プロセッサコア 302 は、プロセッサノードによって行われるメモリ初期化及び試験タスクのためのコード命令を実行する。例えばプロセッサノードが制御ノード又はマスターノードとして機能する場合には、プロセッサコア 302 は、メモリ初期化及び試験タスクを、他のプロセッサノードに割り当てられまた分散せられる複数のサブタスクへと分割するためのコード命令を実行する。一方、プロセッサノードが実行ノード又はスレーブノードの 1 つとして機能する場合には、プロセッサコア 302 は、割り当てられたメモリ初期化及び試験サブタスクを獲得し、スケジューリングし、また実行し、そして結果を制御ノード/マスターノードへ戻すためのコード命令を実行する。クロスバー 308 は、要求、応答、及びブロードキャストメッセージをプロセッサ 302 及びノード又は適切な単一若しくは複数の H T リンクインタフェース 304 へ転送する。要求、応答、及びブロードキャストメッセージの転送は、B I O S により構成される必要のある各処理ノード 202 内の多重構成ルーティングテーブル(multiple configuration routing tables)によって直接的に指令される(directed)。また、メモリ制御器 306 は、B I O S によってプログラムされる必要のある動作パラメータのための多重構成レジスタを含む。

20

30

【0021】

図 4 は例示的な D I M M (デュアルインラインメモリモジュール) 400 を示している。理解されるであろうように、メモリアレイ 206 は、幾つかの、典型的には 8 つの D I M M から構成されてよく、ここで各 D I M M 400 は、D D R (ダブルデータレート)メモリチップ等の多重ランダムアクセスメモリ(R A M)集積回路又はチップ 402 [1:N]を含む。また、D I M M 400 は E C C (誤り訂正符号)回路 404 を有していてもよい。E C C 回路 404 は、メモリ誤りを発見して且つ訂正することを可能にする誤り訂正符号を記憶する。更に、D I M M 400 は S P D (シリアルプレゼンスディテクト)回路 406 を有していてもよい。S P D 回路 406 は、D I M M 400 の動作範囲を指定するリードオンリ情報及びデータシートにおいて見出されるであろう他の同様の情報を含む。例えば S P D 回路 406 は、D I M M 400 のメモリ記憶容量、最小サイクル時間等の動作パラメータ、C A S 待ち時間、等を識別する。

40

【0022】

システム初期化に際しては、各メモリモジュール 206 は初期化され且つ試験される必要がある。このことは、メモリモジュールのデータ投入状況(population)を検証することと、メモリの適切な動作(スタックビット無し)を検証することと、メモリを既知の値に初期化し又はクリーニングすることと、を含んでいてよい。各メモリモジュールは、任意のメモリ誤りを訂正するために E C C を周期的に利用することに基づいて洗われ(scrubbed)得る。

【0023】

50

図5は本発明の選択された実施形態に従う最適化されたメモリ初期化及び試験シーケンス500のフロー図を示している。シーケンスは、ステップ502で、各処理ノードが他の処理ノード及びメモリと通信することができるようにシステムリンクを初期化することによって開始する。このことは、典型的には、ブートストラッププロセッサ上で動作しているBIOSによって行われる。代替的には、ハードワイヤードのシステム(hard wired system)においては、処理ノードは、システム電源投入に際して互いに及びメモリと通信するように自動的に構成されてよい。

【0024】

ステップ504では、ブートストラッププロセッサ(BSP)として指定された処理ノードは、メモリ初期化及び/又は試験プロセスが、メモリ初期化及び/又は試験タスクを複数のサブタスクに分割し、次いで各タスクを実行処理ノードに割り当てると共に分散させることによって開始するようにさせる。例えばBSPは、割り当てられたサブタスクと共に開始メッセージを各処理ノードへ送ることができる。追加的に又は代替的に、割り当てられたサブタスクがBSPによって送り届けられ得るように又は割り当てられたサブタスクが単一若しくは複数の実行処理ノードによってレトリブされ(retrieved)得るように、メモリ初期化及び試験プロセスが開始すべきであることを示すビットが各処理ノード内に書き込まれてよい。代替的には、任意のプロセッサ上でのコード実行が1つ以上のレジスタに書き込む場合等に、指令された割り込みがプロセスを開始させることができる。書き込まれた値は、割り当てられたサブタスクに対して1つ以上のプロセッサが特定のアクションをとるようにそれらに指令することができる。例えば書き込まれた値は、指令された割り込みの目標を識別するノード識別子と、どこでコード実行を開始するのかを直接的に又は間接的に指定することを介して実行するための一連の命令及びサブタスクを示すベクターと、を含むことができる。必要であれば、目標とされているプロセッサがベクターにより指定される割り当てられたサブタスクの実行を開始することを促す開始メッセージを、目標とされているプロセッサにハードウェアメカニズムが送り届けることができる。

【0025】

ステップ506では、各処理ノード202は、その割り当てられたメモリ初期化及び試験サブタスクを行い、それらは関連するメモリ206に対応していてもしていなくてもよい。図2に示される例示的な実装においては、コンピューティングシステム200は、4つのコア、1つの制御コア202[0]、及び3つのスレーブコア202[1:3]を含む。メモリ初期化/試験タスクは、制御コア202[0]によって複数のサブタスク(例えばサブタスク1、サブタスク2、サブタスク3等)に分割される。制御コア202[0]は、サブタスク1を第1のスレーブコア202[1]へ送り、サブタスク2を第2のスレーブコア202[2]へ送り、そしてサブタスク3を第3のスレーブコア202[3]へ送る。ステップ506では、各スレーブコア202[1:3]はそのサブタスクを受け取り、そして実行のためにそれをスケジューリングする。また、制御コア202[0]は、異なるサブタスクの実行を継続してよく、あるいはスレーブコア上での実行が完了するのを待ってよく又はそれ自身の別個のサブタスクを実行してよい。各スレーブコア202[1:3]は、それがタスクを処理する資源を有しているタイミングでサブタスクを実行し、その結果、実行の時間は他のスレーブコアと重なるかもしれないし重ならないかもしれない。結果として、割り当てられたメモリ初期化/試験サブタスクは、スレーブコアの資源の能力に応じて、コア202によってシリアルに又は並列に実行されることが可能である。サブタスク実行は、処理ノードによって又は各処理ノードのノースブリッジ内の回路によって行われ得る。

【0026】

理解されるであろうように、制御コア及びスレーブコアは、独立した複数のプラットフォームシステム上にあってよい。そのような実装においては、制御コアは、各参加しているプラットフォームシステム上の1つ以上のコアが、その割り当てられている単一又は複数のサブタスクを各々処理し得るように、サブタスクを当該システムへ送る。このシナリオにおいては、実行はシリアルに又は並列になされてよい。メモリ初期化タスク実行の分

10

20

30

40

50

散させられた性質の結果として、実行が並列に又はシリアルになされるという事実は、タスクの完了とは関係がない。タスクが並列に実行される場合、タスクはより速く完了するであろうが、データを転送する待ち時間は、制御コアでのシリアル化を結果としてもたらず可能性がある。

【 0 0 2 7 】

ステップ 5 0 8 では、各処理ノードは状況を B P S に報告する。状況報告することは、例えば、メモリ初期化及び試験の間に連続的であってよく、あるいは周期的であってよく、あるいはメモリ初期化及び試験が完了したときであってよい。また、状況報告は種々の方法で報告され得る。例えば、ブートストラッププロセッサは、各処理ノードへ周期的にクエリーを送り、割り当てられたサブタスクが実行され終わったかどうかを問い合わせることができる。代替的には、各処理ノードは、その状況をブートストラッププロセッサに示すメッセージを送ることができる。別の代替案としては、各処理ノードは、ローカルレジスタに又はブートストラッププロセッサのローカルメモリに書き込んで状況を報告することもできる。

【 0 0 2 8 】

全ての処理ノードがそれらの割り当てられたメモリ初期化 / 試験を完了した後に、ステップ 5 1 0 では、ブートストラッププロセッサがシステムブートを継続し、そして他の処理ノードは停止する。メモリ初期化 / 試験タスクが実行されるべきである場合（決定 5 1 2 が肯定）には、処理はステップ 5 0 4 で再び始まる。そうでない場合には、処理はステップ 5 1 4 で終了する。

【 0 0 2 9 】

メモリ初期化 / 試験サブタスクを異なる処理ノードに割り当てることによって、システムブートは極めて速く完了することができる。幾つかのサブタスクは先行するサブタスクの完了に依存し、また他のタスクは先行するサブタスクが完了するのを待たずにシーケンスからは自由に完了され得るのではあるけれど、タスクの依存性の有無を問わずにサブタスクを分割し、割り当て、そして分散させることによって、全体的な性能及び起動時間を改善することができる。このことは、サブタスクの幾つかは、タスクが実行順序になる前にそれらが現れ且つそれらの前に依存性のある実行を有していないことを待つ必要性なしに第 1 の処理ノードによって実行されるという理由で「システムブート(system boot)」又は「システムをブートすること(booting the system)」の時間が低減されるという事実

【 0 0 3 0 】

ここに説明されるように、分散させられたメモリ初期化ルーチンは、分散及び別個のノードによる実行のためにメモリタスクを分割することによって多重プロセッサシステム内の多重ノード間で分散させられるメモリを構成し、初期化し、そして試験することが可能な B I O S コードと共に実装されてよく、結果として改良された実行時間をもたらすことができる。例えば、1つのコアは S M バス(SMBus)（低速バス）から S P D を読み出すタスクを課せられ得る一方で、別のコアは S P D 依存でない複雑な初期化タスクを実行することができる。開示される分散型処理手法はまた、特定の B I O S アルゴリズムの精度を改善することができる。例えば、より速い実行のために B I O S が最適化されるようにアルゴリズムをトレーニングすることが開発されてきており、それにより遅延設定の最適配置を決定するために用いられ得るデータアイが生成される。しかし、生成されるデータアイは、トレーニングアルゴリズムを実行するコアにより処理され得るデータの量によって制限される。ここに開示される分散型処理アプローチは、主コアを用いてデータを集めそして情報を処理しまた結果を主コアに戻すシステム内の他のコアへの分散のために当該デ

ータを分割することによって精度を改善するために用いられ得る。開示される分割及び分散型スキームはまた結果として、割り当てられたサブタスクが順序を問わず実行される場合に、より速いブート時間をもたらす。更に別の可能な利益は、ブートストラッププロセッサにBIOS初期化タスクの全てを行わせることに代えて、サブタスクが各コアに分散させられるようにすることによって、処理コアがBIOS初期化の間に効率的に使用され得ることである。

【0031】

更なる例、及びメモリ初期化タスクがどのようにして分割されそして複数の処理ノードにわたって分散させられ得るのかの説明を提供するために、以下のメモリ初期化ルーチン又はアルゴリズムが次に説明され、ここでサブタスク s_t の総数はシステム内のコアの総数以下である。

10

【0032】

$T = \{s_t[1] + s_t[2] + s_t[3] + \dots s_t[X]\}$ 、ここで s_t はメモリ初期化サブタスクである；

【0033】

$C_0 =$ 「 s_t 」を結合して「 T 」を生成することに関与する主初期化コア；

【0034】

$c[n] =$ サブタスクコア、但し $n = \{1, \dots, N\}$ ；

【0035】

$N =$ N X となるコアの最大数；

20

【0036】

$X =$ 任意の T に対するサブタスクの最大数；

【0037】

$n =$ 現在のコア数；

【0038】

$x =$ 現在のサブタスク数；

【0039】

ステップ1 - サブタスクをコアに分散させる；

【0040】

ステップ2 - $c[n]$ に対して；

30

【0041】

ステップ3 - $s_t[x]$ に対して；

【0042】

ステップ4 - C_0 から $c[n]$ まで $s_t[x]$ 初期化情報及びコードを送る；

【0043】

ステップ5 - $c[n]$ は $s_t[x]$ の実行を開始；

【0044】

ステップ6 - C_0 は $c[n]$ から $s_t[x]$ 結果を受け取り実行を開始；

【0045】

ステップ7 - C_0 は n 及び x をインクリメント；

40

【0046】

ステップ8 - C_0 は全ての X に対してステップ2～7を繰り返す；

【0047】

ステップ9 - $n = N$ の場合にステップ10へ移動；

【0048】

ステップ10 - C_0 は分散させられたサブタスクデータを収集；

【0049】

ステップ11 - $c[n]$ に対して；

【0050】

ステップ12 - C_0 は $c[n]$ から $s_t[x]$ データを要求；

50

【 0 0 5 1 】

ステップ 1 3 - $c[n]$ は実行を完了して結果を c_0 へ送る；

【 0 0 5 2 】

ステップ 1 4 - c_0 は n 及び x をインクリメント；

【 0 0 5 3 】

ステップ 1 5 - c_0 は全ての $n = N$ になるまで全ての X に対してステップ 1 1 ~ 1 4 を繰り返す；

【 0 0 5 4 】

ステップ 1 6 - c_0 は全ての $s t[1 - X]$ からのデータを再結合することによって T を生成；

【 0 0 5 5 】

ステップ 1 7 - 全てのステップが次のタスクのために繰り返される。

【 0 0 5 6 】

理解されるであろうように、前述のルーチンは、疑似コード形態で表現され得るし且つ / 又は特定の命令セットアーキテクチャのためのそれらが使用している対応するアセンブリ言語若しくは高レベル言語コードへとトランスレートされ得る。また、ここに説明される動作は、コンピュータシステムユーザによって直接エントリされたコマンド及び / 又はソフトウェアモジュールによって実行されるステップを含み得る。ここに参照される任意のステップの機能性は、モジュール又はモジュールの一部の機能性に対応してよい。ソフトウェアモジュールに加えて、上述のフロー又はフローの一部は、アプリケーション命令として実装され得る。ここに参照される動作は、モジュール又はモジュールの一部（例えばソフトウェア、ファームウェア、又はハードウェアのモジュール）であってよい。例えば、ここで論じられるソフトウェアモジュールは、スクリプト、バッチ若しくは他の実行可能なファイル又はそのようなファイルの組み合わせ及び / 若しくは一部を含んでいてよい。ソフトウェアモジュールは、コンピュータ可読媒体上でエンコードされるコンピュータプログラム又はそのサブルーチンを含んでいてよい。

【 0 0 5 7 】

加えて、当業者であれば、モジュール間の境界は単に例示的なものであり、また代替的な実施形態はモジュールを合併し又はその代わりにモジュールの機能性の分解を課する場合があることを理解するはずである。例えば、ここで論じられるモジュールは、多重コンピュータ処理として実行されるべきサブモジュールへと分解されてよい。また、代替的な実施形態は特定のモジュール又はサブモジュールの多重インスタンスを結合してよい。更に、当業者であれば、例示的な実施形態において説明される動作は例示のみを目的としていることを認識するはずである。本発明に従う追加的な動作においては、動作は結合されてよく、あるいは動作の機能性は分散させられてよい。従って、ここに説明されるフロー、その動作、及びそのためのモジュールは、フローの動作を実行するように構成されるコンピュータシステム上で実行されてよく且つ / 又はコンピュータ可読媒体から実行されてよい。フローは、フローを実行するコンピュータシステムを構成するための機械可読媒体及び / 又はコンピュータ可読媒体において具現化されてよい。従って、ソフトウェアモジュールは、コンピュータシステムメモリ内に記憶され且つ / 又はコンピュータシステムメモリへ伝送されて、モジュールの機能を行うためのコンピュータシステムを構成してよい。

【 0 0 5 8 】

上に開示される特定の実施形態は例示のみを目的としており、また本発明は異なる方法であるがここでの教示の利益を有する当業者に明らかな均等な方法において修正されそして実施され得るので、上に開示される特定の実施形態は本発明を限定するものとして解釈されるべきではない。従って、前述の説明は、記載されている特定の形態に本発明を限定することを意図するものではないが、一方で、本発明の精神及び範囲の最も広い形態から逸脱することなしに当業者が種々の変更、置換及び代替をなし得ると当業者が理解するはずであるように添付の特許請求の範囲によって定義される本発明の精神及び範囲内に含ま

10

20

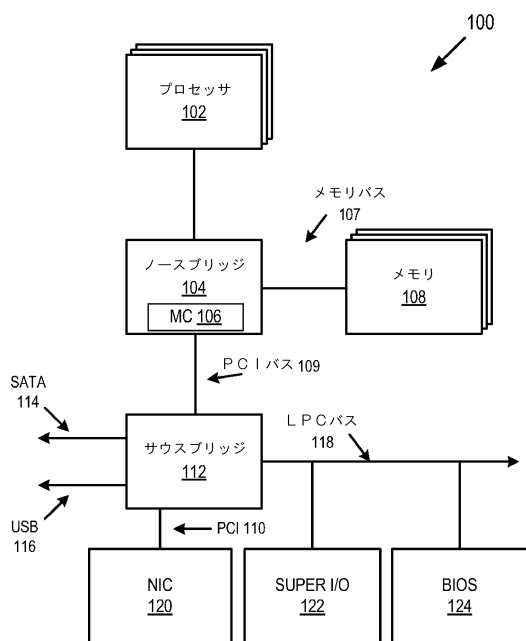
30

40

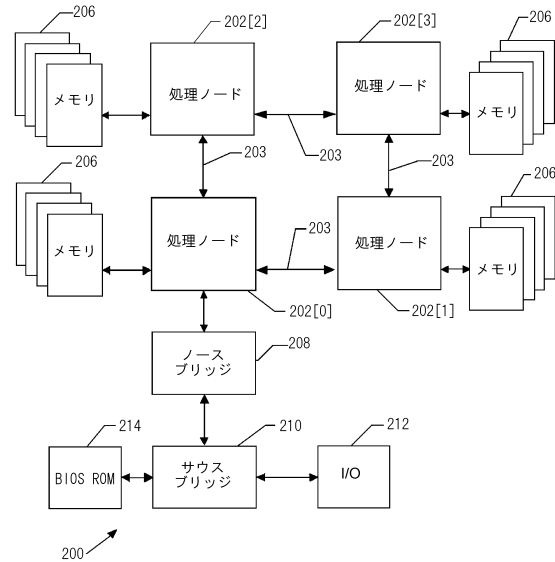
50

れるであろう代替、修正、及び均等なものに及ぶことが意図されている。

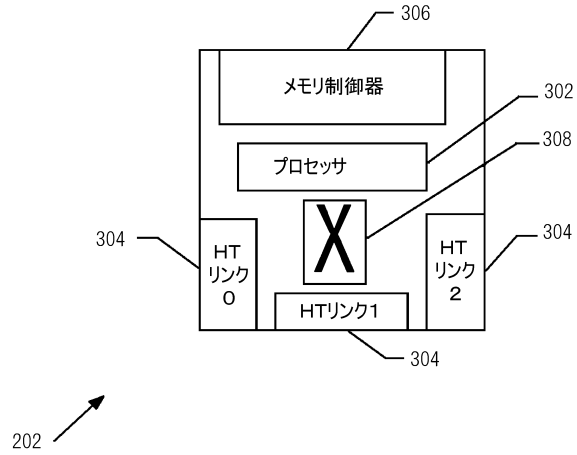
【図 1】



【図 2】



【図3】



【図4】

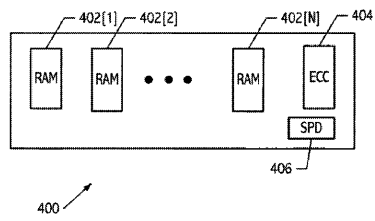
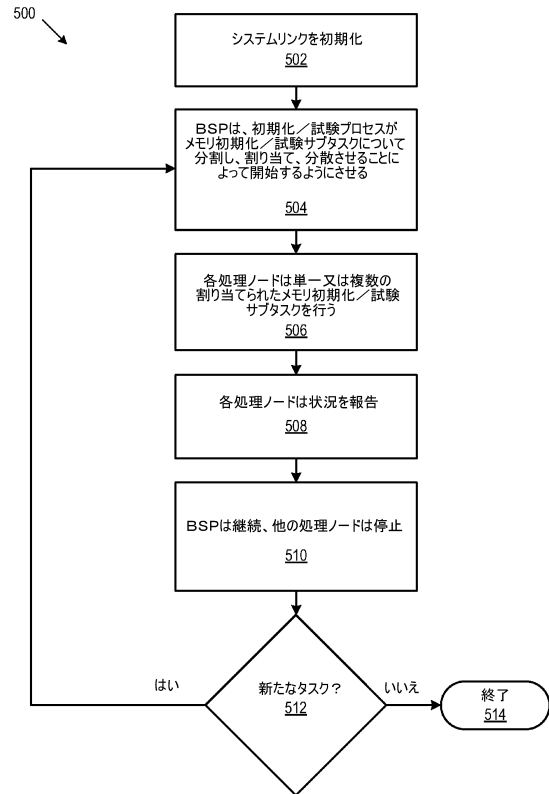


Figure 4

【図5】



フロントページの続き

(74)代理人 100162156

弁理士 村雨 圭介

(72)発明者 オズウィン ハウスティ

アメリカ合衆国 78749 テキサス州、オースティン、センドラ メサ ドライブ 8811

審査官 三坂 敏夫

(56)参考文献 国際公開第00/017750(WO, A1)

米国特許出願公開第2006/0149959(US, A1)

米国特許出願公開第2008/0162878(US, A1)

米国特許第05673388(US, A)

米国特許第06158000(US, A)

特開2009-217297(JP, A)

特開2006-309332(JP, A)

特開2004-145733(JP, A)

特開2009-110429(JP, A)

特表2002-525745(JP, A)

特表2009-506436(JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 9/50

G06F 12/00

G06F 15/177