

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5285690号
(P5285690)

(45) 発行日 平成25年9月11日(2013.9.11)

(24) 登録日 平成25年6月7日(2013.6.7)

(51) Int.Cl. F I
GO6F 15/177 (2006.01) GO6F 15/177 C
GO6F 15/173 (2006.01) GO6F 15/173 640C

請求項の数 17 (全 14 頁)

(21) 出願番号	特願2010-503440 (P2010-503440)	(73) 特許権者	390009531
(86) (22) 出願日	平成20年3月20日 (2008.3.20)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公表番号	特表2010-525433 (P2010-525433A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公表日	平成22年7月22日 (2010.7.22)		アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード
(86) 国際出願番号	PCT/EP2008/053377	(74) 代理人	100108501
(87) 国際公開番号	W02008/128836		弁理士 上野 剛史
(87) 国際公開日	平成20年10月30日 (2008.10.30)	(74) 代理人	100112690
審査請求日	平成22年10月29日 (2010.10.29)		弁理士 太佐 種一
(31) 優先権主張番号	11/736,811	(74) 代理人	100091568
(32) 優先日	平成19年4月18日 (2007.4.18)		弁理士 市位 嘉宏
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 並列コンピュータ・システム、並列コンピュータ・システム上のノード・トラフィックを動的に再経路指定するためのコンピュータ実装方法、コンピュータ可読記録媒体及びコンピュータ・プログラム

(57) 【特許請求の範囲】

【請求項 1】

並列コンピュータ・システムであって、
 論理的な3次元デカルト配列内で、各ノードをその6個の最近隣ノード(X+、X-、Y+、Y-、Z+及びZ-)に接続するトラス・ネットワークによって接続された複数のノードと、

前記並列コンピュータ・システムの前記複数のノード及びネットワーク接続を監視し且つノード及びネットワーク接続の問題リストを作成するノード/ネットワーク監視機構とを備え、

前記複数のノードの少なくともいずれかは、前記問題リスト内のノード及びネットワーク接続を回避するように、各ノードから送信されるデータ・パケット内に6個の最近隣ノード(X+、X-、Y+、Y-、Z+及びZ-)用に設定されたヒント・ビットを使用して、データ・パケットを障害ノード又は障害ネットワーク接続を迂回して1つ以上の前記ネットワーク接続を介して動的に経路指定するものである、並列コンピュータ・システム

10

【請求項 2】

前記ヒント・ビットが、トラス・ネットワーク上でトラフィックを送信すべき好ましい方向を指示する複数の2進値であり、6つのビットをもって、XYZの順序で、X+ X- Y+ Y- Z+ Z-のように定義され、X+又はX-の何れか一方を設定することができるが、両方を設定することができないものである、請求項1に記載の並列コン

20

ピュータ・システム。

【請求項 3】

前記ヒント・ビット(516)が、前記ネットワーク接続を介して送信されるデータ・パケット(510)のヘッダ(512)内に含まれる、請求項1又は請求項2に記載の並列コンピュータ・システム。

【請求項 4】

前記ノードは、更新された前記問題リストを検出し、前記検出に応じて当該ノード上で実行中のアプリケーションを一時的に停止し、更新された前記問題リストに対応する前記ヒント・ビットを使用してデータ・パケットを動的に経路指定するようにし、前記アプリケーションを一時的に停止された点から再開するものである、請求項1ないし請求項3の何れか1項に記載の並列コンピュータ・システム。

10

【請求項 5】

前記並列コンピュータ・システムが、前記複数のノードを3次元トラス・ネットワークで相互接続した超並列コンピュータ・システムである、請求項1ないし請求項4の何れか1項に記載の並列コンピュータ・システム。

【請求項 6】

ヒント・ビットを使用することにより、並列コンピュータ・システム上で実行中のアプリケーションを再始動することなく、ヒント・ビットを使用して、当該並列コンピュータ・システム内の論理的な3次元デカルト配列内で、各ノードをその6個の最近隣ノード(X+、X-、Y+、Y-、Z+及びZ-)に接続するトラス・ネットワークによって接続された複数の計算ノード上のノード・トラフィックを動的に再経路指定するためのコンピュータ実装方法であって、

20

問題について前記複数の計算ノード及びネットワーク接続を監視し且つ障害ノード及び障害ネットワーク接続を問題リスト内で特定するステップと、

更新済みの問題リストを検出するステップと、

アプリケーションを実行中の計算ノードの実行を一時的に停止させるステップと、

前記問題リスト内のノード又はネットワーク接続を分離するように、各ノードから送信されるデータ・パケット内に6個の最近隣ノード(X+、X-、Y+、Y-、Z+及びZ-)用に設定されたヒント・ビットを設定するステップと、

一時的に停止させていた前記アプリケーションの実行を再開するステップとを含む、コンピュータ実装方法。

30

【請求項 7】

前記ヒント・ビットが、トラス・ネットワーク上でトラフィックを送信すべき好ましい方向を指示する複数の2進値であり、6つのビットをもって、XYZの順序で、X+ X- Y+ Y- Z+ Z-のように定義され、X+又はX-の何れか一方を設定することができるが、両方を設定することができないものである、請求項6に記載のコンピュータ実装方法。

【請求項 8】

前記ヒント・ビット(516)が、前記ネットワーク接続を介して送信されるデータ・パケット(510)のヘッダ(512)内に含まれる、請求項6又は請求項7に記載のコンピュータ実装方法。

40

【請求項 9】

前記並列コンピュータ・システムが、前記複数の計算ノードを3次元トラス・ネットワークで相互接続した超並列コンピュータ・システムである、請求項6ないし請求項8の何れか1項に記載のコンピュータ実装方法。

【請求項 10】

少なくとも1つの障害ノードのプロセスを少なくとも1つのバックアップ・ノードに移動させるステップをさらに含む、請求項6ないし請求項9の何れか1項に記載のコンピュータ実装方法。

【請求項 11】

50

計算ノードに対する輻輳ネットワークを検出するステップと、
 トラフィックを当該輻輳ネットワークを迂回して経路指定するように少なくとも1つの
 ヒント・ビットを設定するステップとをさらに含む、請求項6ないし請求項10の何れか
 1項に記載のコンピュータ実装方法。

【請求項12】

論理的な3次元デカルト配列内で、各ノードをその6個の最近隣ノード(X+、X-、
 Y+、Y-、Z+及びZ-)に接続するトーラス・ネットワークを構成するネットワーク
 接続によって接続された複数のノードを有する並列コンピュータ・システムにおいて、
 コンピュータを、

前記複数のノード及びネットワーク接続を監視し且つノード及びネットワーク接続の問
 題リストを作成するノード/ネットワーク監視機構として機能させ、そして 10

前記問題リスト内のノード及びネットワーク接続を回避するように、各ノードから送信
 されるデータ・パケット内に6個の最近隣ノード(X+、X-、Y+、Y-、Z+及びZ-
)用に設定されたヒント・ビットを使用して、データ・パケットを前記ネットワーク接
 続を介して動的に経路指定するノード経路指定機構として機能させるためのコンピュータ
 ・プログラム。

【請求項13】

前記ヒント・ビットが、トーラス・ネットワーク上でトラフィックを送信すべき好まし
 い方向を指示する複数の2進値であり、6つのビットをもって、XYZの順序で、X+
 X- Y+ Y- Z+ Z-のように定義され、X+又はX-の何れか一方を設定する 20
 ことができるが、両方を設定することができないものである、請求項12に記載のコンピ
 ュータ・プログラム。

【請求項14】

前記ヒント・ビット(516)が、前記ネットワーク接続を介して送信されるデータ・
 パケット(510)のヘッダ(512)内に含まれる、請求項12又は請求項13に記載
 のコンピュータ・プログラム。

【請求項15】

前記ノードは、更新された前記問題リストを検出し、前記検出に応じて当該ノード上
 で実行中のアプリケーションを一時的に停止し、更新された前記問題リストに対応する前記
 ヒント・ビットを使用してデータ・パケットを動的に経路指定するようにし、前記アプリ
 ケーションを一時的に停止された点から再開するものである、請求項12ないし請求項1
 4の何れか1項に記載のコンピュータ・プログラム。 30

【請求項16】

前記並列コンピュータ・システムが、前記複数のノードを3次元トーラス・ネットワー
 クで相互接続した超並列コンピュータ・システムである、請求項12ないし請求項15の
 何れか1項に記載のコンピュータ・プログラム。

【請求項17】

請求項6ないし請求項11の何れか1項に記載の方法の各ステップをコンピュータに実
 行させるためのコンピュータ・プログラム。

【発明の詳細な説明】 40

【技術分野】

【0001】

本発明は、並列コンピュータ・システム内の障害回復に係り、さらに詳細に説明すれば
 、超並列コンピュータ・システム上で実行中のアプリケーションを再始動することなく、
 ヒント・ビットを使用して、超並列コンピュータ・システムの計算ノード上のノード・ト
 ラフィックを動的に再経路指定するための装置及び方法に係る。

【背景技術】

【0002】

高性能のコンピュータ・システムについては、故障時間及び修理コストを減少させるた
 めに、効率的な障害回復が重要である。多数の計算ノードを有する並列コンピュータ・シ 50

システムでは、単一コンポーネントの障害が生じる結果として、大きな部分又はコンピュータ全体が、修理のためにオフラインに強制されることがある。アプリケーションの再始動は、かかる障害に先行する相当な量の処理時間を浪費することがある。

【0003】

超並列コンピュータ・システムは、相互接続された多数の計算ノードを有する、1つのタイプの並列コンピュータ・システムである。かかる超並列コンピュータのファミリーは、「Blue Gene」という名前の下で、国際的なビジネス・マシナリー・コーポレーション (IBM) によって開発中である。Blue Gene/Lシステムは、スケラブル・システムであり、計算ノードの現在の最大数は、65,536である。Blue Gene/Lシステムのノードは、2つのCPU及びメモリを有する単一のASIC (特定用途向け集積回路) から成る。コンピュータ全体は、64個のラック内に収容され、各ラック内には32個のノード・ボードがある。

10

【0004】

Blue Gene/Lシステムは、幾つかの通信ネットワークを介して通信する。65,536個の計算ノードは、1つの論理ツリー・ネットワーク及び1つの3次元トラス・ネットワークの両方に配列される。論理ツリー・ネットワークは、各計算ノードが1つの親ノード及び1つ又は2つの子ノードと通信するように、計算ノードをツリー構造状に接続する。トラス・ネットワークは、各計算ノードがコンピュータの1つのセクション内でその6個の最近隣ノードと通信するように、計算ノードを3次元の格子構造状に接続する。計算ノードが、隣接ノードとの通信を必要とするトラス及びツリー・ネットワーク内で配列されるので、単一ノードのハードウェア障害が生じると、障害ハードウェアを修理することができるまで、システムの大きな部分が停止させられることになる。例えば、単一ノードの障害が生じると、トラス・ネットワークの完全なセクションが作動不能になることがある。Blue Gene/Lシステム内のトラス・ネットワークの1つのセクションは、ラックの半分又は512個のノードである。さらに、障害が修正されるまで、障害を有する区画に割り当てられた全てのハードウェアをオフラインにする必要がある。

20

【発明の概要】

【発明が解決しようとする課題】

【0005】

従来技術の超並列コンピュータ・システムでは、実行中に単一ノードの障害が生じると、ソフトウェア・アプリケーションを最初から又は保存されたチェックポイントから再始動することをしばしば必要とする。障害事象が生じる場合、障害ノードの処理を他のノードに移動できるようにすることが有利であろう。そのようにすれば、アプリケーションを、最小の遅延で、バックアップ・ハードウェア上で再開することが可能となり、その結果、システム全体の効率が向上するからである。障害ノード又は障害を生じる可能性がある指示されたノードから一層効果的に回復する方法がなければ、並列コンピュータ・システムは、潜在的なコンピュータ処理時間を浪費し続け、その結果、運用コストが増大するであろう。

30

【課題を解決するための手段】

40

【0006】

超並列コンピュータ・システム上で実行中のアプリケーションを再始動することなく、障害ノード又は輻輳ネットワークを迂回して経路指定するように、ヒント・ビットを使用して、超並列コンピュータ・システムの計算ノード上のノード・トラフィックを動的に再経路指定するための装置及び方法が開示される。1つのノードが障害を有するか、又は当該ノードが障害を生じる可能性があるという指示が存在する場合、超並列コンピュータ・システム上のアプリケーション・ソフトウェアが中断され、その間に、障害ノード上のデータがバックアップ・ノードに移動される。トラス・ネットワークのトラフィックは、障害ノードを迂回して経路指定され、そして当該障害ノード用のトラフィックは、バックアップ・ノードに再経路指定される。同様に、ネットワーク・トラフィックは、輻輳ネッ

50

トワークを迂回して経路指定することができる。

【0007】

本明細書に開示する複数の例及び説明は、Blue Geneアーキテクチャに向けられているが、これらの例及び説明は、ノード・ハードウェアが他のノードからのカットスルー・トラフィックを処理する、ネットワーク構造状に配列された複数のプロセッサを有する任意の並列コンピュータ・システムまで及ぶ。

【発明の効果】

【0008】

本発明は、ノード・トラフィックを動的に再経路指定することにより、故障時間の量を著しく減少させることができ、その結果、コンピュータ・システムの効率を向上させることができるという効果を奏する。

10

【図面の簡単な説明】

【0009】

【図1】超並列コンピュータ・システムのブロック図である。

【図2】超並列コンピュータ・システム内の計算ノードの入出力接続を示すブロック図である。

【図3】超並列コンピュータ・システム内の計算ノードのブロック図である。

【図4】超並列コンピュータ・システム内の計算ノードのトラス・ネットワーク・ハードウェアのブロック図である。

【図5】超並列コンピュータ・システム内のトラス・ネットワーク用データ・パケットのブロック図である。

20

【図6】ノード・トラフィックを動的に再経路指定する1つの例を示すために、超並列コンピュータ・システムの一部を表すブロック図である。

【図7】ノード・トラフィックを動的に再経路指定する他の例を示すために、超並列コンピュータ・システムの一部を表すブロック図である。

【図8】並列コンピュータ・システム内の問題リストを作成するために、ノード及びネットワークを監視するための方法のフローチャートである。

【図9】並列コンピュータ・システム内のヒント・ビットを使用して、ノード・プロセスを動的に再経路指定するための方法のフローチャートである。

【発明を実施するための形態】

30

【0010】

本発明は、超並列コンピュータ・システム上で実行中のアプリケーションを再始動することなく、ヒント・ビットを使用して、超並列コンピュータ・システムの計算ノード上のノード・トラフィックを動的に再経路指定するための装置及び方法に向けられている。1つのノードが障害を有するか、又は当該ノードが障害を生じる可能性があるという指示が存在する場合、超並列コンピュータ・システム上のアプリケーション・ソフトウェアが中断され、その間に、障害ノード上のデータがバックアップ・ノードに移動される。トラス・ネットワークのトラフィックは、障害ノードを迂回して経路指定され、そして当該障害ノード用のトラフィックは、バックアップ・ノードに再経路指定される。以下、インターナショナル・ビジネス・マシーンズ・コーポレーション（IBM）によって開発中のBlue Gene/Lシステムに関連して、複数の例を説明する。

40

【0011】

図1は、Blue Gene/Lシステムのような超並列コンピュータ・システム100を表すブロック図である。Blue Gene/Lシステムは、スケラブル・システムであり、その計算ノードの最大数は65,536である。各計算ノード110は、「Blue Gene/L」計算チップとも呼ばれる、1つの特定用途向け集積回路（ASIC）112を有する。計算チップ112は、2つのプロセッサ又は中央処理装置（CPU）を有し、ノード・ドータ・カード114に取り付けられる。一般に、計算ノード110は、512メガバイトのローカル・メモリ（図示せず）を有する。1つのノード・ボード120は、32個のノード・ドータ・カード114を収容し、各ノード・ドータ・カード1

50

14は、1つの計算ノード110を有する。従って、各ノード・ボード120は、32個の計算ノード110を有し、各計算ノード110は、2つのプロセッサ及び各プロセッサに関連するメモリを有する。ラック130は、32個のノード・ボード120を収容するハウジングである。各ノード・ボード120は、ミッドプレーン・コネクタ134を介して、ミッドプレーン132に接続する。ミッドプレーン132は、ラックの内部にあり、図1には示されていない。Blue Gene/Lシステムの全体は、各々が32個のノード・ボード120を有する、64個のラック130内に収容されるであろう。従って、システムの全体は、65,536個の計算ノード110及び131,072個のCPU(64ラック×32ノード・ボード×32ノード×2CPU)を有するであろう。

【0012】

Blue Gene/Lシステムの構造は、1つのI/Oノード面を有する、1つの計算ノード・コアとして記述することができる。すなわち、1,024個の計算ノード110への通信が、サービス・ノード140に接続されたI/Oプロセッサ170を有する、各I/Oノードによって処理される。これらのI/Oノードは、ローカル・ストレージを有していない。これらのI/Oノードは、論理ツリー・ネットワークを通して計算ノード110に接続され、また機能的なネットワーク(図示せず)を通して機能的な広域ネットワーク能力を有する。この機能的なネットワークは、サービス・ノード140から多数の計算ノード110への通信を処理するための1つのノード・ボード120に位置する、1つのI/Oプロセッサ(又は「Gene/L」リンク・チップ)170に接続される。Blue Gene/Lシステムは、ノード・ボード120に接続された1つのI/Oボー

【0013】

コンピュータ・システム100に含まれるサービス・ノード140は、ノードへのソフトウェアのローディングを処理するとともに、システム全体の動作を制御する。一般に、サービス・ノード140は、リナックス(登録商標)を実行する、制御コンソール(図示せず)付きの「IBM pSeries」サーバのようなミニコンピュータ・システムである。サービス・ノード140は、制御システム・ネットワーク150を介して、計算ノード110のラック130に接続される。制御システム・ネットワーク150は、Blue Gene/Lシステム用の制御、テスト及び立ち上げインフラストラクチャを提供する。制御システム・ネットワーク150は、超並列コンピュータ・システムに必要な通信を提供する、種々のネットワーク・インタフェースを含む。これらのネットワーク・インタフェースについては、後述する。

【0014】

サービス・ノード140は、システム管理専用の制御システム・ネットワーク150を管理する。制御システム・ネットワーク150に含まれる100Mbpsの専用イーサネット(登録商標)・リンクは、サービス・ノード140から多数のノードへの通信を処理するノード・ボード120上に位置する、Idoチップ180に接続される。このネットワークがJTAGネットワークとも呼ばれるのは、これがJTAGプロトコルを使用して通信するためである。ノード・ボード120上にある計算ノード110の制御、テスト及び立ち上げの全ては、サービス・ノード140と通信するJTAGポートを通して管理される。さらに、サービス・ノード140に含まれるノード/ネットワーク監視機構142は、回避すべき障害ノード、障害を生じる可能性があるノード又はネットワーク・リンクを指示する、問題リスト144を維持する。ノード/ネットワーク監視機構142は、サービス・ノード140内のソフトウェアから成るが、システムのノード上で実行中のオペレーティング・システムによって支援されることがある。

【0015】

Blue Gene/Lシステムは、幾つかの通信ネットワークを介して通信する。図

10

20

30

40

50

2は、Blue Gene/Lシステム内の1つの計算ノードのI/O接続を示すブロック図である。65,536個の計算ノード110及び1,024個のI/Oプロセッサ170は、論理ツリー・ネットワーク及び論理的な3次元トラス・ネットワークの両方に配列される。トラス・ネットワークは、各計算ノード110がその6個の最近隣ノードと通信するように、計算ノードを3次元の格子構造状に接続する。図2では、トラス・ネットワークは、当該ノードを6個の隣接ノードにそれぞれ接続するネットワーク接続X+、X-、Y+、Y-、Z+及びZ-によって例示される。一方、ツリー・ネットワークは、図2のTree 0、Tree 1及びTree 2接続によって表される。当該ノードに接続される他の通信ネットワークは、1つのJTAGネットワーク及び1つのグローバル割り込みネットワークを含む。JTAGネットワークは、サービス・ノード140から図1の制御システム・ネットワーク150を介して、テスト及び制御用の通信を提供する。グローバル割り込みネットワークは、或るタスクの完了時に異なる処理段階に移動するように、複数の計算ノード上の同様のプロセスを同期化するためのソフトウェア・バリアを実装するために使用される。従って、このグローバル割り込みネットワークは、ノードの1つの区画で実行中のアプリケーションを始動、停止又は一時停止させるために使用することができる。さらに、各計算ノード110に対するクロック及び電力信号がある。

【0016】

Blue Gene/Lシステムのトラス相互接続は、論理的な3次元デカルト配列内で、各ノードをその6個の最近隣ノード(X+、X-、Y+、Y-、Z+及びZ-)に接続する。6個の最近隣ノードへの接続は、ノード・レベル及びミッドプレーン・レベルで行われる。各ミッドプレーンは、8×8×8ノードの配列である。ミッドプレーン内のノード配列の6個の面(X+、X-、Y+、Y-、Z+及びZ-)の各々は、8H8=64ノードのサイズを有する。6個の面の各々にある64ノードからの各トラス・ネットワーク信号は、当該ミッドプレーンに接続されたリンク・カード(図示せず)を通して、隣接ミッドプレーン内の対応するノードに通信される。当該ミッドプレーンが、任意の次元において1ミッドプレーンの深さを有する区画内で使用される場合、各面の信号は、対向面上にある同じミッドプレーンの入力に戻されることがある。

【0017】

図3は、従来技術に従った、Blue Gene/Lシステム内の計算ノード110のブロック図である。計算ノード110は、ノード計算チップ112を有し、当該チップは、2つのプロセッサ310A、310Bを有する。各プロセッサ310A、310Bは、それぞれ処理コア312を有する。これらのプロセッサは、レベル3のメモリ・キャッシュ(L3キャッシュ)320及びスタティックRAM(SRAM)330に接続される。L3キャッシュ320からのデータは、ダブル・データ・レート(DDR)メモリ・コントローラ350を介して、DDR同期型ダイナミックRAM(SDRAM)340にロードされる。

【0018】

SRAM330は、計算チップ112からのデータをIdoチップ180に通信する、JTAGインタフェース360に接続される。サービス・ノード(図1)は、Idoチップ180及び制御システム・ネットワーク150(図1)の一部であるイーサネット(登録商標)・リンクを通して、計算ノード110と通信する。Blue Gene/Lシステムでは、ノード・ボード120当たり1つのIdoチップが存在し、各ミッドプレーン132(図1)内のボード上に他のものが存在する。Idoチップ180は、信頼のおける100Mbpsの専用イーサネット(登録商標)制御ネットワーク上の生のUDPパケットを使用して、サービス・ノードからコマンドを受信する。Idoチップは、計算ノードとの通信のために、種々のシリアル・プロトコルをサポートする。JTAGプロトコルは、サービス・ノード140(図1)から計算ノード110内のSRAM330の任意のアドレスを対象とする読み取り及び書き込みを行うために使用され、またシステム初期化及びブート・プロセスのために使用される。

【0019】

10

20

30

40

50

また、図3のノード計算チップ112は、ネットワーク・ハードウェア390を含む。ネットワーク・ハードウェア390は、トラス・ネットワーク・ハードウェア392、ツリー・ネットワーク・ハードウェア394及びグローバル割り込みネットワーク・ハードウェア396を含む。前述のように、Blue Gene/Lシステムのこれらのネットワークは、計算ノード110がシステム内の他のノードと通信するために使用される。ネットワーク・ハードウェア390は、計算ノード110がトラス・ネットワークを介してデータ・パケットを受信し且つこれを受け渡すことを可能にする。ネットワーク・ハードウェア390は、ネットワーク・データのトラフィックを独立的に処理し、その結果、計算ノード110のプロセッサ310A、310Bは、トラス・ネットワーク上を流れるデータの量によって負担を負わされることはない。他のノードを宛先として当該ノードをパススルーするネットワーク・データは、「カットスルー」トラフィックと呼ばれる。

10

【0020】

図4は、図3のトラス・ネットワーク・ハードウェア392のブロック図である。トラス・ネットワーク・ハードウェア392は、3つの主要なユニット、すなわちプロセッサ・インタフェース410、トラス・センダ420及びトラス・レシーバ430から成る、プロセッサ・インタフェース410は、先入れ先出し式のキューとして、プロセッサ注入FIFO412及びプロセッサ受信FIFO414から成る。これらのFIFOに対するアクセスは、プロセッサ310A、310B(図3)から2つの浮動小数点ユニット(FPU)レジスタ(図示せず)を介して行われる。すなわち、1対のFPUレジスタから128ビットのメモリ・マップド・ストアを介して、データがこれらのFIFOにロードされる。一方、これらのFIFOからデータが読み取られ、128ビット単位で、1対のFPUレジスタにロードされる。全部で8個あるプロセッサ注入FIFO412は、2グループに編成される。すなわち、一方のグループは、(ノード間オペレーティング・システム・メッセージ用の)2つの高優先順位FIFOから成り、他方のグループは、6個の通常優先順位FIFOから成る(それらのFIFOは最近隣ノード接続性のために十分である)。全てのFIFO内のパケットは、トラス・ネットワーク上の任意の方向に送出することができる。プロセッサ受信FIFO414にも、FIFOの2つのグループがある。各グループは、7個のFIFOを含み、そのうち1つのFIFOは高優先順位用のものであり、6個のFIFOの各々は6入力方向の各々に専用される。具体的には、各レシーバとその対応する受信FIFOの間に、1つの専用バスがある。ストレージについては、全てのトラスFIFOは、エラー検出及び訂正(ECC)によって保護されたSRAMチップを使用し、全ての内部データ経路はパリティについて検査される。

20

30

【0021】

前述のトラス・ネットワーク・ハードウェア392は、種々のトラス・ネットワークを横切って、可変サイズのデータ・パケットを送信する。図5は、トラス・ネットワーク・パケット510の1例を示す。Blue Gene/Lシステムでは、各パケット510は、nH32バイトである。但し、n=1~8個の「チャンク」である。例えば、MPI(Message Passing Interface)規格に準拠するメッセージは、関連する1つ又は両方のプロセッサ310A、310B(図3)上で実行中のソフトウェアによって構築され、送信され且つ受信される、多数のパケットから成ることがある。各パケットの最初の8バイトは、パケット・ヘッダ512である。パケット・ヘッダ512は、リンク・レベルのプロトコル情報(例えば、シーケンス番号)と、宛先を含む経路指定情報と、仮想チャンネル及びサイズと、伝送中のヘッダ・データの破損を検出する、バイト幅の巡回冗長検査(CRC)514とを保持する。また、パケット・ヘッダ512は、後述するヒント・ビット516を保持する。

40

【0022】

パケット・ヘッダ512の後には、多数のデータ・バイト518が保持される。さらに、各パケットには、1バイトの有効性標識520とともに、24ビットのCRCが付加される。有効性標識が必要であるのは、パケットが完全に受信される前にその転送が開始されることがあるためである。このCRCは、各パケットが各リンクを介して送信される際

50

に、当該各パケットの検査を可能にする。タイムアウト機構は、破損パケットを再送するために使用される。8ビットのパケット・ヘッダCRCを使用することは、パケット・ヘッダ・エラーの早期検出を可能にする最適化である。というのは、ヘッダCRCが、完全なパケットCRC内に含まれるからである。

【0023】

前述のように、ヘッダ512は、6つの「ヒント」ビット516を保持する。ヒント・ビット516は、トラス・ネットワークの3次元内でパケットを経路指定することができる方向を指示する。ヒント・ビット516は、XYZの順序で、 $X+$ $X-$ $Y+$ $Y-$ $Z+$ $Z-$ のように定義される。例えば、100100のヒント・ビットは、当該パケットを $X+$ 及び $Y-$ 方向に経路指定することができることを意味する。 $X+$ 又は $X-$ ヒント・ビットの何れか一方を設定することができるが、両方を設定することはできない。というのは、設定された1つのビットは、その次元内でパケットを送信すべき方向を指示するからである。デフォルトは、パケットを任意の方向に送信できることを指示するように、全てのヒント・ビットを設定解除又はゼロに設定することであろう。

10

【0024】

一般に、トラス・ネットワーク内では、ノード間をデータが流れる次元順序が存在する。本明細書で開示する複数の例では、その次元順序がXYZであると仮定するが、他の順序も使用することができる。次元順序がXYZであるという意味は、最初に、データが1つのノードからX次元内で流れ、次に、Y次元内で幾つかのノードを通り、次に、Z次元内で流れるというものである。XYZヒント・ビットは、XYZ次元内の経路指定のためにそれぞれ使用される。

20

【0025】

各ノードは、トラス機能(図示せず)を制御する、ソフトウェアで構成可能な1組のレジスタを維持する。例えば、1組のレジスタは、その隣接ノードの座標を保持する。ヒント・ビットがゼロに設定されるのは、1つのパケットが1つのノードを1つの方向に離れ、そして隣接ノード座標レジスタによって決定されるように、当該パケットがその次元内のその宛先に到着するような場合である。ヒント・ビットは、ヘッダ内で早期に現れ、その結果、調停を効率的にパイプライン化することができる。ヒント・ビットは、ソフトウェア又はハードウェアの何れかによって初期化することができる。ハードウェアによって初期化される場合、適切な方向を決定するために、1次元当たり1組2つのレジスタが使用される。これらのレジスタは、最小ホップ経路指定を提供するように構成することができる。経路指定を完全に行うためには、ヒント・ビット及び仮想チャンネルを調べることが必要である。すなわち、経路指定テーブルは、存在しない。パケットは、動的に又は決定論的な次元順序(XYZ)で、経路指定することができる。すなわち、パケットは、他のトラフィックに基づいて最小輻輳の経路を辿るか、又は固定経路上で経路指定することができる。2地点間パケットのほかに、ヘッダ内の1つのビットを設定することにより、1つのパケットを任意のデカルト次元内でブロードキャストさせ且つこれを各ノードで受信させることができる。後述するように、ソフトウェアは、ヒント・ビットを適切に設定することにより、「デッド」ノード又はリンクを回避することができる。3つまでの非共線型(noncolinear)障害ノードが存在する場合には、完全な接続性を維持することができる。

30

40

【0026】

図6は、ノード・トラフィックを動的に再経路指定する1つの例を示すために、図1の超並列コンピュータ・システムの部分600を表すブロック図である。この部分600は、ノード1610~ノード9612と表記した9個のノードを例示する。図面の内容を簡潔にするために、図6は、X次元及びY次元内のノードのみを示しているが、当該コンピュータ・システムは、Z次元内に位置するノードも有することを理解されたい。X次元及びY次元は、XY軸614によって示される通りである。この例については、アプリケーションがノード1610~ノード8622上で実行中であるものと仮定する。障害又は潜在的な障害がノード5618上で検出される場合、アプリケーションが中断又は一時

50

停止され、そしてF I F O内の全てのネットワーク・トラフィックがクリアされるまで待機することにより、ネットワークが休止される。次に、障害ノード5 6 1 8上のアプリケーションは、スペア・ノード(ノード9 6 1 2)に移動される。次に、障害ノード5 6 1 8を通してデータを移動させる必要がある各ノードは、更新済みの問題リスト(図1の1 4 4)を全てのノード又は少なくとも影響を受けるノードに送信することにより、障害ノード5 6 1 8を回避するように更新される。

【0 0 2 7】

これらのノードは、更新済みの問題リストを使用して、データが障害ノード又は障害ネットワークを迂回して経路指定されることを保証する。次に、データが障害ノード5 6 1 8を迂回して経路指定されるように、各ノードから送信されるデータ・パケット内に適切なヒント・ビットが設定される。図6の例では、ノード2 6 2 0からのデータ・パケットは、X - 用に設定されたヒント・ビットを有するから、X - 方向に進行してノード8 6 2 2に至り、その結果、障害ノード5 6 1 8を回避するであろう。同様に、ノード8 6 2 2からのデータ・パケットは、X + 用に設定されたヒント・ビットを有するから、X + 方向に進行してノード2 6 2 0に至る。また、ノード4 6 2 4からのデータ・パケットは、Y + 用に設定されたヒント・ビットを有するから、Y + 方向に進行してノード6 6 2 6に至り、その結果、障害ノード5 6 1 8を回避するであろう。これに対し、ノード6 6 2 6からのデータ・パケットは、Y - 用に設定されたヒント・ビットを有するから、Y - 方向に進行してノード4 6 2 4に至る。

【0 0 2 8】

図7は、ノード・トラフィックを動的に再経路指定する他の例を示すために、図1の超並列コンピュータ・システムの部分7 0 0を表すブロック図である。この例は、非隣接ノードについてのヒント・ビットの使用法を示す。図6の例と同様に、この部分7 0 0は、ノード1 6 1 0 ~ ノード9 6 1 2と表記した9個のノードを例示する。この例では、障害又は潜在的な障害がノード8 6 2 2上で検出される。アプリケーションが一時停止され、ネットワークが休止され、そして障害ノード8 6 2 2上のアプリケーションがスペア・ノード(ノード9 6 1 2)に移動される。次に、図6の例と同様に、障害ノード8 6 2 2を通してデータを移動させる必要がある各ノードは、更新済みの問題リスト(図1の1 4 4)に影響を受けるノードに送信することにより、障害ノード8 6 2 2を回避するように更新される。次に、データが障害ノード8 6 2 2を迂回して経路指定されることを保証するように、ヒント・ビットが設定される。図7の例では、ノード1 6 1 0からのデータ・パケットは、障害ノード8 6 2 2がX方向に存在しないという理由で、X方向用に設定されたヒント・ビットを有しないであろう。しかし、ノード1 6 1 0は、パケットがY + 方向に進行するように、設定済みのY + ヒント・ビットを有するであろう。ノード1 6 1 0からのパケットがノード7 6 2 8に到着し且つY次元内の進行を開始する場合、当該パケットは、設定済みのY + ヒント・ビットによって指示されるようにY + 方向のノード9 6 1 2に進行し、その結果、障害ノード8 6 2 2を回避するであろう。

【0 0 2 9】

前述のように、ヒント・ビットは、輻輳ネットワークを迂回して動的に経路指定するためにも使用することができる。1例として、図7のノード8 6 2 2とノード5 6 1 8との間のネットワーク7 1 0を検討する。もし、ネットワーク7 1 0が、ネットワーク監視機構(図1の1 4 2)によって輻輳ネットワークとして標識付けされるならば、障害ノード8 6 2 2を迂回して経路指定するためにヒント・ビットを使用する前述の方法と同様に、このネットワークを迂回して動的に経路指定するためにヒント・ビットを使用することができる。代替的に、カットスルー・トラフィックによって、ノードの負担が過度に大きくなることがある。例えば、ネットワーク監視機構が、ノード8 6 2 2を通るカットスルー・トラフィックに起因して、ノード8 6 2 2が過負荷であると決定した場合、ノード8 6 2 2上で実行していたプロセス又はアプリケーションに関するカットスルー・トラフィック負荷を軽減するために、ノード8 6 2 2上のプロセスは、利用可能な交換ノードに動的に再経路指定される。

【 0 0 3 0 】

図 8 は、並列コンピュータ・システム内の障害ノードのプロセスを動的に再経路指定するために、ノードを監視するための方法 8 0 0 を示す。この方法は、サービス・ノード上のソフトウェアによって実施されるが、必要とされる情報を集めるためにソフトウェア及びハードウェアの一方又は両方を必要とすることがある。最初に、ネットワークを監視し（ステップ 8 1 0）、ネットワーク・ホットスポットを問題リスト内に記録する（ステップ 8 2 0）。次に、ノードを監視し（ステップ 8 3 0）、障害ノード又は障害を生じる可能性があるノードを問題リスト内に記録する（ステップ 8 4 0）。その後、この方法が終了する。

【 0 0 3 1 】

図 9 は、並列コンピュータ・システム内の障害ノードのプロセスを動的に再経路指定するための方法 9 0 0 を示す。好ましくは、この方法は、並列コンピュータ・システムの各ノード上のソフトウェア及びハードウェアの一方又は両方によって実施される。最初に、サービス・ノード上のネットワーク監視機構によって送信され且つ回避すべきノード又はネットワークを保持する、更新済みの問題リストを検出する（ステップ 9 1 0）。次に、障害ノードを有する並列システムの区画で実行中のアプリケーションを一時的に停止する（ステップ 9 2 0）。次に、トーラス・ネットワーク・ハードウェアの F I F O がそれらのメッセージの送信を完了するまで待機することにより、ネットワークを休止する（ステップ 9 3 0）。次に、ネットワーク用の交換ノード又は代替経路を見つけた後（ステップ 9 4 0）、障害ノードのプロセスを交換ノードに移動させる（ステップ 9 5 0）。次に、障害ノード又は障害ネットワークを通してネットワーク・トラフィックを送信するであろうノードに対し、ヒント・ビットを使用して、当該障害ノードを分離し且つ当該障害ノード又は輻輳ネットワークを迂回してネットワーク・トラフィックを経路指定するように通知する（ステップ 9 6 0）。次に、アプリケーションをその一時的に停止した点から再開することができる（ステップ 9 7 0）。その後、この方法が終了する。

【 0 0 3 2 】

本明細書は、超並列コンピュータ・システム上で実行中のアプリケーションを再始動することなく、ヒント・ビットを使用して、超並列コンピュータ・システムの計算ノード上のノード・トラフィックを動的に再経路指定するための装置及び方法を開示する。ノード・トラフィックを動的に再経路指定すると、故障時間の量を著しく減少させることができ、その結果、コンピュータ・システムの効率を向上させることができる。本方法は、コンピュータ・ソフトウェアの形態で実施することができる。

【 0 0 3 3 】

当業者には明らかなように、本発明の範囲内で多数の変形が可能である。本発明は特定の実施形態を参照して説明されたが、本発明の精神及び範囲から逸脱することなく、形式及び詳細に関する種々の変更を施し得ることは明らかであろう。

【 符号の説明 】

【 0 0 3 4 】

- 1 0 0 . . . 超並列コンピュータ・システム
- 1 1 0 . . . 計算ノード
- 1 4 0 . . . サービス・ノード
- 1 4 2 . . . ノード/ネットワーク監視機構
- 1 4 4 . . . 問題リスト
- 1 5 0 . . . 制御システム・ネットワーク
 - 3 9 2 . . . トーラス・ネットワーク・ハードウェア
 - 3 9 4 . . . ツリー・ネットワーク・ハードウェア
 - 3 9 6 . . . グローバル割り込みネットワーク・ハードウェア
 - 5 1 0 . . . トーラス・ネットワーク・パケット
 - 5 1 2 . . . パケット・ヘッダ
 - 5 1 6 . . . ヒント・ビット

10

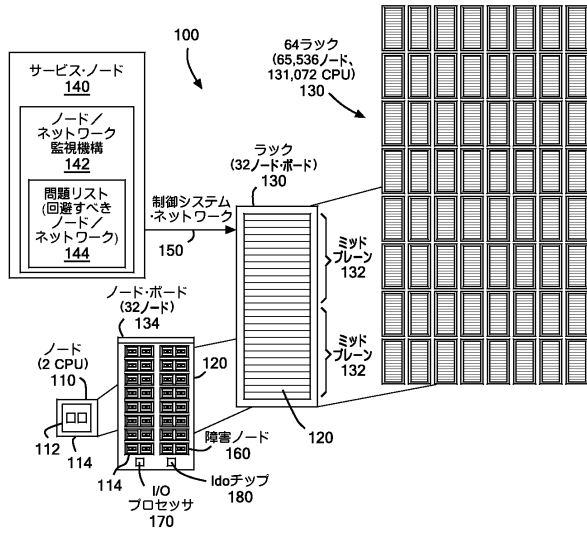
20

30

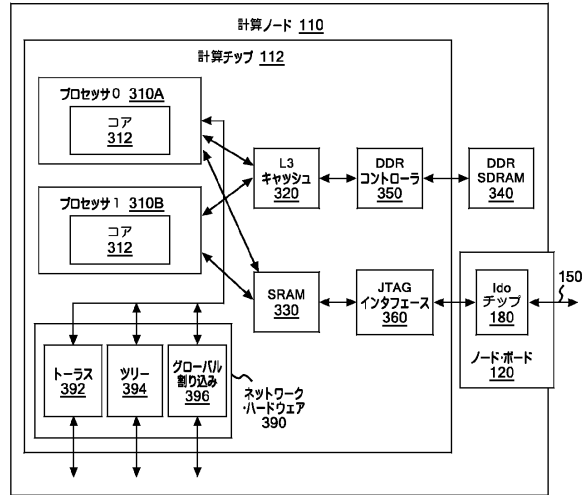
40

50

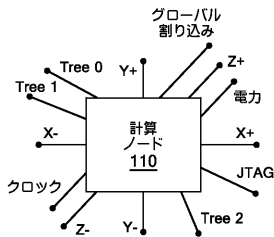
【図1】



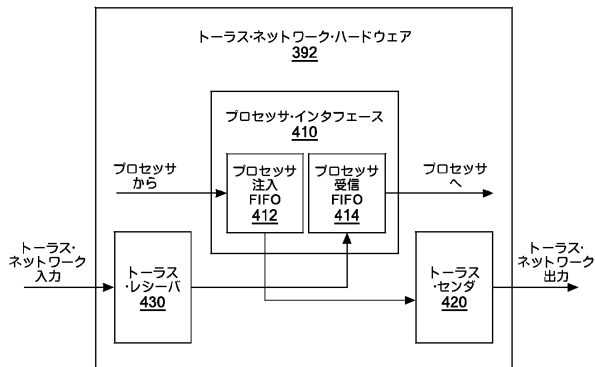
【図3】



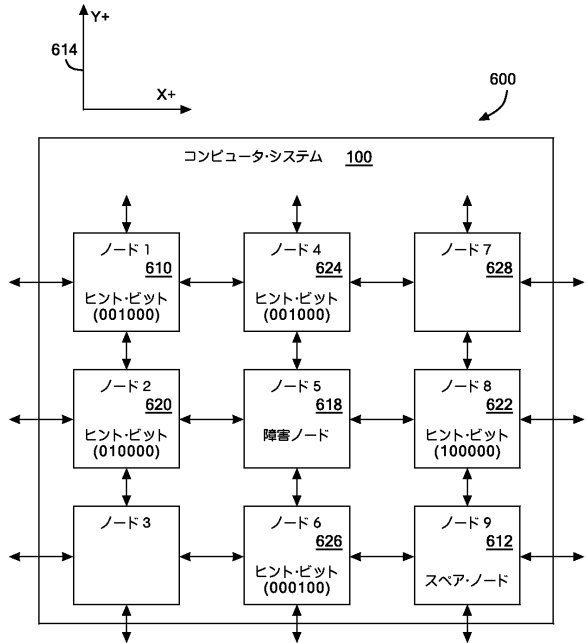
【図2】



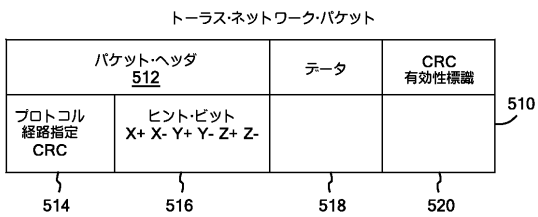
【図4】



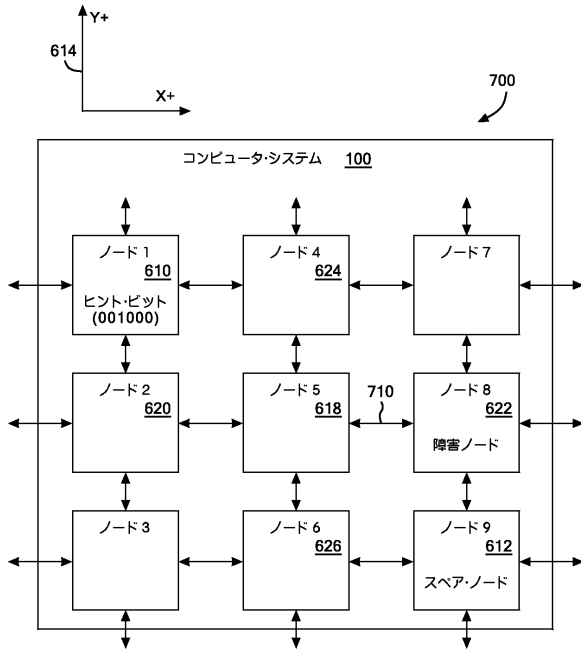
【図6】



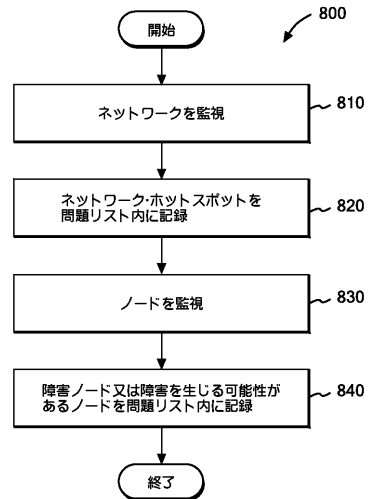
【図5】



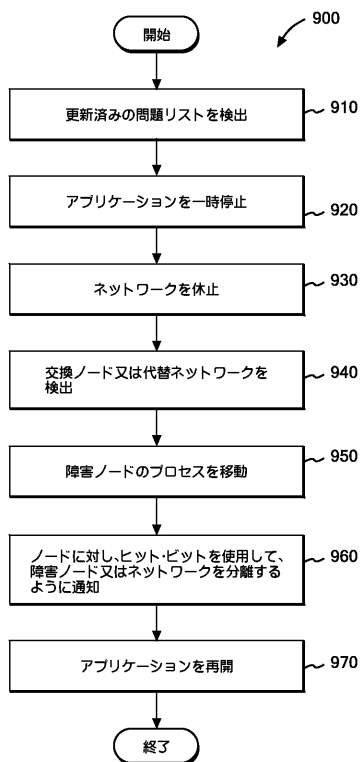
【 図 7 】



【 図 8 】



【 図 9 】



フロントページの続き

- (72)発明者 ピーターズ、アマンダ
アメリカ合衆国 5 5 9 0 1 ミネソタ州 ロチェスター ニクラウス・ドライブ エヌ ダブリュー
ー 5 4 1 9
- (72)発明者 シデルニク、アルバート
アメリカ合衆国 6 1 8 0 1 イリノイ州 アーバナ サウス・リン・ストリート 4 0 2
- (72)発明者 ダーリントン、デーヴィッド
アメリカ合衆国 5 5 9 0 6 ミネソタ州 ロチェスター センチュリー・ヒルズ・ドライブ エヌ
イー 2 0 4 5
- (72)発明者 マッカーシー、パトリック、ジョセフ
アメリカ合衆国 5 5 9 0 1 ミネソタ州 ロチェスター ヴァルキュリヤ・コート エヌ ダブリ
ュー 2 1 2 9
- (72)発明者 シュワルツ、プレント、アラン
アメリカ合衆国 5 4 7 2 9 ウィスコンシン州 チペワ・フォールズ 第 6 5 アヴェニュー 1 9
1 4 1
- (72)発明者 スミス、ブライアン、エドワード
アメリカ合衆国 5 5 9 0 1 ミネソタ州 ロチェスター ノル・レーン エヌ ダブリュー 3 1
2 6

審査官 清木 泰

- (56)参考文献 特開平 0 7 - 2 2 1 7 5 4 (J P , A)
特開平 0 7 - 2 3 9 8 3 5 (J P , A)
N.R.Adiga et al. , Blue Gene/L torus interconnection network , IBM Journal of Research a
nd Development , 米国 , IBM , 2 0 0 5 年 5 月 , Vol:49, No:2/3 , Pages:265-276
Chao Wang, Frank Mueller, Christian Engelmann, Stephen L. Scott , A Job Pause Service u
nder LAM/MPI+BLCR for Transparent Fault Tolerance , Proceedings of IEEE International P
arallel and Distributed Processing Symposium 2007(IPDPS2007) , IEEE , 2 0 0 7 年 3 月 2
6 日 , Pages:1-10
A.Gara et al. , Overview of the Blue Gene/L system architecture , IBM Journal of Researc
h and Development , 米国 , IBM , 2 0 0 5 年 5 月 , Vol:49, No:2/3 , Pages:195-212
Jyothish Varma, Chao Wang, Frank Mueller, Christian Engelmann, Stephen L. Scott , Scala
ble, Fault-Tolerant Membership for MPI Tasks on HPC Systems , Proceedings of the 20th a
nnual International Conference on Supercomputing (ICS2006) , ACM , 2 0 0 6 年 6 月 2 8
日 , Pages:219-228
Yawei Li, Zhiling Lan , Exploit Failure Prediction for Adaptive Fault-Tolerance in Clus
ter Computing , Proceedings of Sixth IEEE International Symposium on Cluster Computing
and the Grid (CCGRID 06) , IEEE , 2 0 0 6 年 5 月 1 6 日 , all 8 pages
Daniel A. Reed, Charng-da Lu, Celso L. Mendes , Reliability challenges in large systems
 , Future Generation Computer Systems , NL , Elsevier Science Publishers B. V. , 2 0 0 6
年 2 月 , Volume:22, Issue:3 , Pages:293-302

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 F 1 5 / 1 6 - 1 5 / 1 7 7

G 0 6 F 1 5 / 8 0

(54)【発明の名称】並列コンピュータ・システム、並列コンピュータ・システム上のノード・トラフィックを動的に再経路指定するためのコンピュータ実装方法、コンピュータ可読記録媒体及びコンピュータ・プログラム