

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7066694号
(P7066694)

(45)発行日 令和4年5月13日(2022.5.13)

(24)登録日 令和4年5月2日(2022.5.2)

(51)国際特許分類		F I		
G 0 6 F	9/50 (2006.01)	G 0 6 F	9/50	1 5 0 A
G 0 6 F	15/167 (2006.01)	G 0 6 F	15/167	
G 0 6 F	15/173 (2006.01)	G 0 6 F	15/173	

請求項の数 11 (全19頁)

(21)出願番号	特願2019-518274(P2019-518274)	(73)特許権者	507402886 パルテック・クラスター・コンペテン ・センター・ゲゼルシャフト・ミット・ ベシュレンクテル・ハフツング PARTEC CLUSTER COMP ETENCE CENTER GMBH ドイツ、8 1 6 7 9 ミュンヘン、ポツ サートシュトラッセ、2 0
(86)(22)出願日	平成29年10月5日(2017.10.5)	(74)代理人	110001195 特許業務法人深見特許事務所
(65)公表番号	特表2019-533249(P2019-533249 A)	(72)発明者	リッペルト, トーマス ドイツ、6 3 7 4 3 アシャッフエンブ ルク、ゾンネンシュトラッセ、4 3
(43)公表日	令和1年11月14日(2019.11.14)	審査官	井上 宏一
(86)国際出願番号	PCT/EP2017/075375		
(87)国際公開番号	WO2018/065530		
(87)国際公開日	平成30年4月12日(2018.4.12)		
審査請求日	令和1年8月7日(2019.8.7)		
(31)優先権主張番号	16192430.3		
(32)優先日	平成28年10月5日(2016.10.5)		
(33)優先権主張国・地域又は機関	欧州特許庁(EP)		
前置審査			

最終頁に続く

(54)【発明の名称】 高性能のコンピューティングシステムおよび方法

(57)【特許請求の範囲】

【請求項 1】

アプリケーションプログラムの計算を行なうためのモジュラーコンピューティングシステムであって、
少なくともクラスタモジュールとブースタモジュールとによって形成される異なるモジュールを備え、各モジュールは複数のノードを含み、
 前記モジュラーコンピューティングシステムは、前記複数のノードに分散されたモジュラーコンピューティング抽象化層をさらに備え、
 前記モジュラーコンピューティング抽象化層は、前記異なるモジュールの前記ノードのために、共有メモリ通信を設定し、前記モジュラーコンピューティング抽象化層によって共有仮想アドレス空間通信を使用するモジュール内およびモジュール間通信ならびに管理機能を提供する、モジュラーコンピューティングシステム。

【請求項 2】

前記異なるモジュールは、少なくとも、ストレージモジュールによって形成されていることを特徴とする、請求項 1 に記載のモジュラーコンピューティングシステム。

【請求項 3】

前記モジュラーコンピューティング抽象化層は、前記ノード内に設けられたノードマネージャによって実現されることを特徴とする、請求項 1 または 2 に記載のモジュラーコンピューティングシステム。

【請求項 4】

前記ノードマネージャは、管理ネットワークを介して互いに通信することを特徴とする、請求項 3 に記載のモジュラーコンピューティングシステム。

【請求項 5】

前記複数のノードは、通信ネットワークを介して通信することを特徴とする、請求項 1 ~ 4 のいずれか 1 項に記載のモジュラーコンピューティングシステム。

【請求項 6】

前記通信ネットワークは、共有メモリ通信を用いて実現されることを特徴とする、請求項 5 に記載のモジュラーコンピューティングシステム。

【請求項 7】

前記ノードマネージャ間の通信は、前記モジュール間の通信から分離されている、請求項 3 に記載のモジュラーコンピューティングシステム。

10

【請求項 8】

前記ノードマネージャは、前記ノードマネージャの各々が各ノードおよび前記ノード間の各インターコネクタのステータスについての情報のセットを維持するように情報をやり取りする、請求項 3 に記載のモジュラーコンピューティングシステム。

【請求項 9】

前記情報のセットに応じて前記ノードにリソースを割り当てることによって前記システムの動作を動的に調整するように前記情報のセットが用いられる、請求項 8 に記載のモジュラーコンピューティングシステム。

【請求項 10】

20

アプリケーションプログラムの計算を行なうためのモジュラーコンピューティングシステムを管理するための方法であって、前記モジュラーコンピューティングシステムは、少なくともクラスタモジュールとブスタモジュールとによって形成される異なるモジュールを備え、各モジュールは複数のノードを含み、前記モジュラーコンピューティングシステムは、前記複数のノードに分散されたモジュラーコンピューティング抽象化層を備え、前記方法は、

前記モジュラーコンピューティング抽象化層を用いて、前記異なるモジュールの前記ノードのために、共有メモリ通信を設定し、前記モジュラーコンピューティング抽象化層によって共有仮想アドレス空間通信を使用するノード内およびノード間通信ならびに管理機能を提供するステップを含む、方法。

30

【請求項 11】

ノードマネージャをさらに備え、前記方法はさらに、前記ノードマネージャが、前記ノードのステータスについての情報を収集し、前記収集されたステータスの情報に応じて、前記ノードにリソースを割り当てるステップをさらに含む、請求項 10 に記載のモジュラーコンピューティングシステムを管理するための方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、高性能のコンピューティング（HPC）の技術分野に関する。特に、本発明はヘテロジニアスコンピューティングシステムに関し、より特定のには、協働して計算タスクを処理するように自由に割り当てられ得る異なるモジュールを含むコンピューティングシステムに関する。モジュールコンピューティング抽象化層（MCAL）と呼ばれる制御エンティティが設けられる。この制御エンティティにより、異なるモジュールによって提供されるさまざまなリソースを動的に割り当てるのが可能になる。本発明は、コンピューティングにおけるさまざまな要求に対する調整において柔軟性があるため、クラウドコンピューティングサービスを提供するための基礎となるシステムとしても適用可能である。クラウドコンピューティングサービスは、共有のコンピュータ処理リソースおよびデータを、要求に応じて（多くの場合、インターネットを通じて）コンピュータや他のデバイスに提供するものである。

40

【背景技術】

50

【 0 0 0 2 】

H P C システムの用途は、幅広い技術分野をカバーする。ほんの数例挙げてみただけでも、人間の脳のシミュレーションを始めとして、天気および気候の予測、コンピュータによる流体工学、地震学的描像、電波天文学、データ解析などがある。高性能のコンピューティングへの要求は、絶えず高まっている。

【 0 0 0 3 】

現在主流の H P C アーキテクチャは、クラスタコンピュータおよび超並列処理 (M P P) システムである。クラスタコンピュータは、標準ネットワーク技術を用いて互いに通信する複数のクラスタノードを含む。個々のクラスタノードは、既製の汎用プロセッサを備える。M P P (超並列処理) システムでは、多数のプロセッサが同時に (したがって、並列して) 計算を行なう。それらもまた、通常は既製の汎用プロセッサを使用する。しかしながら、プロセッサは、通信のために専用のネットワーク技術および個別的に適合されたソフトウェアを使用する傾向にある。

10

【 0 0 0 4 】

より高い計算能力への強い要求を満たすために、コンピュータクラスタ内のクラスタノードの数、および M P P システム内のプロセッサの数が近年大幅に増加した。しかしながら、既製の汎用プロセッサを用いてシステム内での数を増やすだけでは解決策にならないことも分かってきた。

【 0 0 0 5 】

既製の汎用プロセッサにも、もちろんプラスの側面はある。それらの単スレッド性能は高い。マルチコアプロセッサの場合、それらは、プロセッサまたはコア 1 つ当たり大容量のメモリを有する。標準プログラミングを使用することができるため、これらのシステムにアプリケーションを移植することが容易である。しかしながら、この技術は、限界に達しつつある。既製の汎用プロセッサに関して、そのクロック速度 (頻度) は、この 1 0 年で実際的には増加していない。このことは、特にシステムの冷却に関して限界とならないように、チップ 1 つ当たりのエネルギー消費が 1 0 0 ワットを大きく超えるべきではないということに密接に関係している。さらに、それらのプロセッサは比較的高価であり、かつエネルギー効率が低い。

20

【 0 0 0 6 】

解決策は、アクセラレータを採用することに見出された。アクセラレータを用いることによって、非常にスケーラブルなアプリケーションの部分を計算することができる。アプリケーションが費用効果の高い方法で多くのリソースを追加することによって高負荷を処理することができる場合、そのアプリケーションは非常にスケーラブルであると言える。コードにおける並列して計算することのできないシークエンシャル部分が、通常、スケーラビリティへの最も重大な制約要因のうちの 1 つである。

30

【 0 0 0 7 】

2 種類以上のプロセッサを用いたシステムは、ヘテロジニアスコンピュータシステムと呼ばれる。異なるアクセラレータは、汎用プロセッサ、具体的には、メニーコアプロセッサおよび G P U (グラフィックス・プロセッシング・ユニット) と組み合わせて使用され得る。メニーコアプロセッサは、5 0 以上のコアを有する。G P U は、何百個もの単純なコンピューティングコアを有する。これらのタイプのプロセッサは、ともにエネルギー効率が良い。

40

【 0 0 0 8 】

ヘテロジニアスコンピュータシステムでは、各クラスタノードが、1 つ以上のマルチコアプロセッサまたは G P U によって形成されるアクセラレータを備える。アクセラレータは、基本的には、クラスタノードのプライマリプロセッサの機能を補うコプロセッサとして機能する。クラスタノードは、コンポーネント同士がインターコネクトスイッチを通じてデータを伝送し合うファブリックを介して通信する。「ファブリック」という用語は、高性能のコンピューティングの分野では、システムトポロジーを意味するために使用される。また、インターコネクトという用語は、クラスタノード間の通信インフラストラクチ

50

ャについて使用され、別々のデバイスを接続する電氣的または光学的な接続またはケーブルを意味する。各ノードにアクセラレータを取り付けることの大きな問題は、割り当てが静的であることである。汎用プロセッサとGPUとの比は、システムの設計時に固定される。

【0009】

国際出願WO2012-049247A1は、コンピュータクラスタ配列、および、当該紹介されたコンピュータクラスタ配列を動作させるための方法によって、上記欠点に取り組んでいる。このコンピュータクラスタ配列は計算ノードCNを備える。これにより、ブースタBに対する特定の計算タスクを動的に外部調達している。このように、計算ノードCNへのブースタBの割り当て手法が紹介されている。割り当ては、実行時に動的に行なわれる。このアーキテクチャでは、ブースタを形成するクラスタノードおよびアクセラレータノードが1つのファブリックに入る。これは、加速クラスタノードと比較して、より柔軟性がある。なぜなら、プロセスの開始時に、加速処理のためにブースタをクラスタノードに割り当てることができるからである。この配列については、Norbert Eicker等による「メニーコア時代におけるクラスタコンピューティングを追究するDEEPプロジェクト(The DEEP project Pursuing cluster-computing in the many-core era)」、第42回並列処理に関する国際会議(International Conference on Parallel Processing)、885~892頁、2013年、および、「メニーコア時代におけるヘテロジニアスクラスタコンピューティングへの代替アプローチであるDEEPプロジェクト(The DEEP Project An alternative approach to heterogeneous cluster-computing in the many-core era)」、同時実行および計算：実践と経験(Concurrency and Computation: Practice and Experience)、28、2394~2411頁、2015年にさらに記載されている。

10

20

【0010】

1つのアプリケーションがクラスタブースタ型などのヘテロジニアスシステムによって計算されるとき、個々の計算タスクは、複数の異なるクラスタノード上で処理される必要があるだけでなく、2つの異なるアーキテクチャ、すなわち、コンピュータクラスタ上およびアクセラレータ上で処理される必要がある。クラスタコンピュータ内、およびブースタ内、さらには、クラスタコンピュータとブースタとの間において、効率的な通信が確保されなければならない。さらに、個々のパーツの信頼性が高い場合であっても、コンポーネント数の増加によってコンポーネントが故障するリスクが生じる。したがって、システムは、高い回復力、すなわち、個々のコンポーネントの故障への耐性を確保しつつ、機能性も保つ必要がある。

30

【発明の概要】

【課題を解決するための手段】

【0011】

本発明は、モジュラーコンピューティングシステムを提供する。モジュラーコンピューティングシステムは、クラスタモジュール、ブースタモジュール、ストレージモジュール、ニューロモーフィックモジュール、データ解析モジュール、グラフィックスモジュール、量子コンピューティングモジュールなどの異なるモジュールを含む。各モジュールは、複数のノードを含み得る。クラスタモジュールは複数のクラスタノードを含み、ブースタモジュールは複数のブースタノードを含み、ストレージモジュールは複数のストレージノードを含む。その他同様である。

40

【0012】

個々のノードは、すべて1つのファブリックに接続されている。異なるモジュールは異なるインターコネクトを有してもよい。クラスタモジュールは、例えばインフィニバンドインターコネクトを使用する。一方、ブースタモジュールは、特殊インターフェースノードを介してインフィニバンドインターコネクトに接続されたインテルOmni-Pathファブリックを使用する。同じことが、他のタイプのモジュールについても当てはまる。

【0013】

50

異なるモジュールを結合して1つのシステムとするために、モジュラーコンピューティング抽象化層(MCAL)が設けられる。MCALは、通信機能および管理機能を提供する。Par Tec Cluster Competence Center社によるPara Stationを使用すれば、以下で説明する機能および特性を有するMCALを実現するための基盤を形成することができる。

【0014】

MCALは、モジュール内での通信(モジュール内通信と呼ぶ)と、モジュール間の通信(モジュール間通信と呼ぶ)とを可能にする。そのために、MCALは、それぞれのファブリックにインターフェースを提供する。例えば、システム全体に対してファブリックが1つの場合もあるし、または、モジュール毎にファブリックが特化されている場合もある(例えば、クラスタモジュールにはインフィニバンド、ブースタモジュールにはインテル Omni-Pathファブリック、など)。イーサネット(登録商標)、10G、ギガビット、または共有メモリ通信など、他のインターコネクトも可能である。MCALは、如何なる種類の通信を用いてもよく、可用性、アプリケーションプログラムによる指示、性能、またはコストなどの要素に基づいて、特定の通信方法を選択する。

10

【0015】

共有メモリ通信は、2つ以上のプロセスが通信するための非常に高速かつ効率的な方法を提供する。異なるタイプの共有メモリ通信が、MCALによって採用され得る。共有メモリを有する同じ物理ノード上で並列タスクの2つ以上のプロセスが実行される場合、これらのプロセス間の共有メモリ通信のために、この共有メモリが使用され得る。これは、典型的には、CPU1つ当たり1つのプロセスが生成されるSMPノード上で起こり得る。対称型マルチプロセッシング(SMP)は、2つ以上の同一のプロセッサが単一の共有メインメモリに接続しているアーキテクチャを意味する。MCALはプロセスについての知識、および、どのノードでプロセスが行なわれているかについての知識を有しているため、MCALは、そのような状況で共有メモリ通信を確立することができる。

20

【0016】

しかしながら、共有メモリ通信はSMP状況に限定されない。プロセスが1つのモジュール内の異なるノード上で実行される場合、または、異なるモジュールに属するノード上で実行される場合であっても、MCALは、共有の仮想アドレス空間通信を確立する。共有の仮想アドレス空間通信を用いると、メッセージのメタデータおよびデータ自体の仮想アドレスを共有するだけで、サイズの大きなメッセージをノード間で伝達することができる。実際のメッセージは、プロセスにより、その割り当てられた仮想アドレス空間を介してアクセス可能である。したがって、高速通信が可能である。MCALは、通信に関係するプロセスまたはノードによって直接共有されるメモリ領域を設けることによって、仮想アドレス空間通信を確立する。

30

【0017】

上記の例から分かるように、MCALは、通信機能を提供するだけでなく、通信方法を選択し、それによって性能および効率を向上させる。MCALは、例えば、プロセスを生成しながら、すなわち、新たな子プロセスを作成および実行しながら、どのインターコネクトが通信に使用可能かを決定することができる。同時に、MCALは、現在行なわれているプロセスを監視し、子プロセスを見つけて終了させるか、または非同期的に実行し続け得る。

40

【0018】

MCALが提供する第2の主要な機能は、管理である。管理機能は、さまざまな局面、具体的には、ノード管理、モジュール管理、プロセス管理、およびリソース管理をカバーする。

【0019】

ノード管理は、個々のノードを対象とする。例えば、クラスタモジュールにおけるクラスタノード、ブースタモジュールにおけるブースタノード、ストレージモジュールにおけるストレージノードなどであり、上で挙げたようなシステム内に存在する他のタイプのモジュールに関しても同様である。初めに、MCALは使用可能なノードを検出する。使用可

50

能なノードとは、使用できる状態にあるノードのことである。これは、例えば、温度、メモリの空き容量など、および、それらの特性の経時的な変化などのシステムパラメータに基づいて、M C A Lが、正常に機能していないノードまたは故障寸前のノードを特定するように構成されていることを含む。したがって、M C A Lはフォールトトレランスを実現し、それによって高い回復力を確保している。

【 0 0 2 0 】

これは、M C A Lが仮想ノードの概念を採用することによって実現する。仮想ノードは、クラスタノード、ブースタノード、ストレージノードなどの実際のハードウェアを表わすものであるが、M C A Lは、自由に動的にノードをマップすることができる。例えば、アプリケーションにおいてX個のクラスタノード、Y個のブースタノード、およびZ個のストレージノード（X、Y、Zは整数）が必要である場合、M C A Lは、すべての利用可能なリソースの中からそれらのノードを提供することができる。ノードが故障している場合、または、ノードがシステムに追加された場合には、M C A Lは、仮想ノードと物理ノードとの間のマッピングを変更するだけで、迅速に反応することができる。利用可能なノードのリストの中から正常に機能していないノードを取り出すことができ、システム全体が機能不全に陥ることがない。

10

【 0 0 2 1 】

また、M C A Lは、到達可能性、レイテンシ、または物理的配置などの要因に応じて、特定のノードを選択することもできる。これは、実行中のアプリケーションに完全に透過的に行なわれ得る。

20

【 0 0 2 2 】

M C A Lは、システム全体を見渡し、各ノードおよび各インターコネクトの状態を絶えず最新の状態に維持している。M C A Lは、ノード1つ当たり利用可能なC P Uの数、実行中のサービス、利用可能なストレージ、インターコネクト帯域幅、空きメモリ、温度、システムクロック速度などの、各ノードおよび各インターコネクトの特定の特性について、常時通知される。

【 0 0 2 3 】

もう一つの機能は、モジュール管理である。モジュール管理は、ノード管理に類似しており、特定のモジュールについてのノード管理を含み得るが、モジュールに特有の側面を考慮する。モジュール管理は、モジュールが他のモジュールと通信するために用いるインターフェースノード（設けられる場合）の状態など、モジュールに特有の側面に焦点を当てる。ノードと同様に、M C A Lは、モジュールの利用増大を可能にするために、モジュールの状態を監視する。

30

【 0 0 2 4 】

M C A Lは、プロセス管理も担う。実行中のアプリケーションの各々は、複数の個々のプロセスに分割される。複数の個々のプロセスは、開始、分散、および監視する必要がある。同時に、リソースを割り当てる必要がある。これは、プロセスが計算能力およびストレージを必要とすることに対応する。プロセスが特定の処理を必要とする場合、M C A Lは、1つ以上のブースタノード、または、ニューロモーフィックノード、データ解析ノード、グラフィックスノード、または、量子コンピューティングノードなどのさらに特定のノードを、それぞれのモジュールから割り当てる。

40

【 0 0 2 5 】

そのために、M C A Lは、プロセスとその要件との間の依存を認識および制御するように構成されている。これに基づいて、M C A Lは、ブースタモジュールまたはクラスタモジュールの場合にはノード上にプロセスを作成し、または、ストレージモジュールの場合にはストレージを割り当てる。プロセスは監視され、依存プロセスはグループとして扱われる。これにより、並列タスクのうちのいずれかが1つで発生するエラーに迅速に応答することが可能である。したがって、プロセスを速やかに再開させることができ、その結果、実行中のアプリケーションの全体の遅延を最小限に抑えることができる。

【 0 0 2 6 】

50

リソース管理は、M C A L がモジュラーコンピューティングシステムにおいて行なうもう一つのタスクである。リソース管理の一部として、M C A L は、オフロード機構を実現する。すなわち、例えばクラスタモジュール上で実行中のプロセスは、ブースタモジュール中のいくつかのブースタノードを必要とし得る。そのような場合には、M C A L は、タスクを並列して実行するために、選択されたブースタノード上にプロセスを生成することができる。

【 0 0 2 7 】

オフロードの一部として、1つのブースタ上で実行されるように作成されるプロセスは、1つだけでなく複数であってもよい。複数のプロセスは、ホストプロセッサの介入を必要とせず、モジュール通信ファブリックを用いて互いに直接通信する。

10

【 0 0 2 8 】

また、リソース管理は、例えば異なるプロセス間の接続数を減少させることによって、メモリ消費を低減することも保証する。これは、必要な時にのみ実際に確立するオンデマンド接続を用いることによって行なわれる。これにより、接続1つ当たりのメモリ消費が劇的に低減される。したがって、ノードなどのリソースは仮想的に扱われるだけでなく、仮想割り当ては物理ネットワークにも適用される。

【 0 0 2 9 】

リソースを割り当てる際に、M C A L は、1つのプロセスのリクエストを考慮するだけでなく、他のプロセスのリクエストも考慮する。これにより、統合制御に基づいて通信インフラストラクチャの変更に柔軟に応答することができる。

20

【 0 0 3 0 】

統合制御は、M C A L によって提供されるもう一つの機能局面である。統合制御は、M C A L が、管理機能から収集した情報を用いて通信を制御すること（およびその逆）を意味する。このように、通信機能と管理機能とは別個のものではなく、システム性能全体をさらに最適化するために統合されたものである。

【 0 0 3 1 】

すべてのモジュールにおけるすべてのノードのすべての異なる局面、ならびにすべてのインターコネクト、モジュール内通信、およびモジュール間通信のステータスについてのM C A L の知識を用いて、システムを作動させる。これにより、システムの挙動の変更、および、システム上で実行中のアプリケーションの要求の変更に対する動的な調整が可能になる。

30

【 0 0 3 2 】

統合制御は、実際のシステム状態およびアプリケーション要求を考慮するだけでなく、次に起こるモジュールのリクエストまたはインターコネクト使用の予測も行なう。

【 0 0 3 3 】

システム状態の情報は、ノードの動的な割り当てのために用いられる。これは、フィードバックループ登録システムの割り当て決定の広範な影響を用いてリソースの割り当てが行なわれることを含む。その結果、将来の割り当てが相応に調整される。また、システムは、自己学習の方法で、そのような情報を記憶することによって将来のケースの割り当て決定を向上させる。

40

【 0 0 3 4 】

アプリケーションの主要部分が通常はクラスタモジュールで開始されるとしても、M C A L は、クラスタモジュールまたはブースタモジュール上でアプリケーションのどの部分を実行させるかを特定する。決定の根拠は、例えば、実行中のコード性能の監視結果をアプリケーションコード自体で示すものである。これにより、アプリケーションのコンピューティング中に割り当てを調整することができる。

【 0 0 3 5 】

M C A L の実現例としては、M C A L は、システム全体にわたって分散されている。M C A L は、各ノード上に実装されるノードマネージャ（N M ）を含む。ノードマネージャはハードウェアで実現されてもよく、またはソフトウェアで実現されてもよく、例えば、そ

50

それぞれのノード上で実行されるデーモンプロセスとして実現されてもよい。ノードマネージャは、モジュールおよびシステム全体を大域的に見るために、絶えず情報を収集する。システム上で実行中のアプリケーションは、管理機能と通信することによって、または個々のノードマネージャと直接通信することによって、この情報を活用することができる。これは特定のインターフェースによって実現される。この特定のインターフェースは、ハードウェアであってもよいし、API（アプリケーションプログラミングインターフェース）であってもよい。このように、MCALがアプリケーションから入力を得るだけでなく、アプリケーションも、MCALによって収集された情報を個々のノードマネージャを介して活用する。

【0036】

ノードマネージャの通信トラフィックは、モジュール内およびモジュール間通信から厳密に分離されている。これにより、モジュール内およびモジュール間通信は、管理機能によって中断されることなく高速で実行される。さらに、このような分離によって、システムセキュリティを高めることができる。

【0037】

実際の通信ネットワークの使用は、MCALによって制御される。MCALは、この機能をメッセージパッシングインターフェース(MPI)を介してアプリケーションに提供する。MCALの一部は、モジュール毎のMPIスタックであるとともに、システム全体に及ぶグローバルMPIである。個々のMPIは、例えば、コンピュータクラスタモジュール、ブースタモジュール、ストレージモジュールなどについて実現される。

【0038】

MCALの通信は、管理ネットワークを用いる。管理ネットワークは、モジュール内およびモジュール間通信に用いられる物理ネットワークと同じ物理ネットワークであってもよいし、異なる物理ネットワークであってもよい。例えば、モジュール内およびモジュール間通信は、インフィニバンド、ミリネットなどであってもよく、管理ネットワークは、ノードマネージャを接続するイーサネットによって形成されてもよい。

【0039】

システムには、管理タスクまたは制御機能がシステムを作動させることを可能にするフロントエンドを提供するための特殊ノードが含まれてもよい。フロントエンドが通信ネットワークに接続されることは必要ではないが、通信ネットワークへだけ接続されている。例えばワークステーションコンピュータ、パーソナルコンピュータ(PC)、ラップトップコンピュータ、または任意のスマートコンピューティングデバイスによって形成され得るフロントエンドノードを介して、モジュラーコンピューティングシステムは管理され得る。モジュラーコンピューティングシステム上で実行中のアプリケーションは、スケジュール、開始、停止、および監視され得る。これは、バッチ処理によって行なわれ得る。すなわち、手動による介入なく、したがって非対話方式で、モジュラーコンピューティングシステム上で一連のジョブを実行することによって行なわれ得る。

【0040】

手動による介入なくモジュラーコンピューティングシステム上で一連のジョブを実行することに加えて、MCALは、対話型コンピューティングを提供し得る。対話型コンピューティングでは、モジュラーコンピューティングシステム上で実行中のアプリケーションが、その実行中に人からの入力を受け付ける。この入力は、ワークステーションコンピュータ、パーソナルコンピュータ(PC)、ラップトップコンピュータ、または任意のスマートコンピューティングデバイスを介して提供され得る。そのようなシナリオでは、例えばラップトップコンピュータを介してモジュラーコンピューティングシステムにアクセスし、ラップトップコンピュータがモジュラーコンピューティングシステムの計算能力を有しているかのように見せることができる。これは、即時または現在進行中であるとユーザが感じるような時間枠内での応答が必要なリアルタイムのアプリケーションに非常に有用である。また、膨大な量の構造化データ、半構造化データ、および非構造化データを情報のマイニングのためにアクセス可能にするビッグデータアプリケーションにおいても、その

10

20

30

40

50

ような配列を利用することができる。これは、計算能力およびストレージを提供するための基礎となるリソースとしてモジュラーコンピューティングシステムを有するクラウドコンピューティングサービスにも当てはまる。

【発明の効果】

【0041】

上述のモジュラーコンピューティングシステムにおいて、MCALとシステムとが強固に統合されることによって、通信レイテンシが低減され、通信を向上させることができる。

【図面の簡単な説明】

【0042】

【図1】本発明に従うモジュラーコンピューティングシステムの機能図である。

10

【図2】本発明に従うモジュラーコンピューティングシステムの第1の実施形態を示す図である。

【図3】本発明に従うモジュラーコンピューティングシステムの第2の実施形態を示す図である。

【発明を実施するための形態】

【0043】

図1は、本発明に従うモジュラーコンピューティングシステム100の機能図である。アプリケーション層110がモジュラーコンピューティングシステムの最上部を形成する。モジュラーコンピューティングシステム100において、アプリケーション層では、ソフトウェアプログラムなどのアプリケーションプログラムが実行され得る。アプリケーション層110がシステムのリソースにアクセスすることを可能にするために、モジュラーコンピューティング抽象化層120(MCAL)が設けられる。システムのリソースは、クラスタモジュール130、ブースタモジュール140、ストレージモジュール150、データ解析モジュール160、および他のモジュール170、例えば、ニューロモフィックモジュール、グラフィックスモジュール、および量子コンピューティングモジュールからなるグループのうちの1つ以上など、さまざまなモジュールによって形成される。同じ種類のモジュールを複数有することも可能である。

20

【0044】

MCAL120は、モジュール内での通信(モジュール内通信と呼ぶ)と、モジュール間の通信(モジュール間通信と呼ぶ)とを可能にする。MCAL120は、如何なる種類の通信を用いてもよく、可用性、アプリケーションプログラムによる指示、性能、またはコストなどの要素に基づいて、特定の通信方法を選択する。

30

【0045】

さらに、MCAL120は管理機能を提供する。管理機能は、さまざまな局面、具体的には、ノード管理、モジュール管理、プロセス管理、およびリソース管理をカバーする。ノード管理は、個々のノードを対象とするものであり、図2および図3を参照して、さらに説明する。

【0046】

MCAL120は、システム全体を見渡し、モジュール130、140、150、160、および170の各々の状態を絶えず最新の状態に維持している。MCAL120は、利用可能なプロセッサの数、実行中のサービス、利用可能なストレージ、インターコネクタ帯域幅、空きメモリ、温度、システムクロック速度などの、モジュール130、140、150、160、および170の特定の特性について、常時通知される。MCAL120が提供するもう一つの機能は、モジュール管理である。モジュール管理は、モジュールに特有の側面を考慮する。したがって、MCALは、クラスタモジュール130、ブースタモジュール140、ストレージモジュール150、データ解析モジュール160、および他のモジュール170を管理するように特別に適合された機能を含む。MCAL120は、プロセス管理も担う。実行中のアプリケーション110の各々は、複数の個々のプロセスに分割される。複数の個々のプロセスは、開始、分散、および監視する必要がある。同時に、リソースを割り当てる必要がある。これは、プロセスが計算能力およびストレージ

40

50

を必要とすることに対応する。プロセスが特定の処理を必要とする場合、M C A L 1 2 0 は、1つ以上のブースタノード、または、ニューロモーフィックノード、データ解析ノード、グラフィックスノード、または、量子コンピューティングノードなどのさらに特定のなノードを、それぞれのモジュールから割り当てる。そのために、M C A L 1 2 0 は、プロセスとその要件との間の依存を認識および制御するように構成されている。リソース管理は、M C A L 1 2 0 がモジュラーコンピューティングシステムにおいて行なうもう一つのタスクである。リソース管理の一部として、M C A L 1 2 0 は、オフロード機構を実現する。すなわち、例えばクラスタモジュール130上で実行中のプロセスは、ブースタモジュール140中のいくつかのブースタノードを必要とし得る。そのような場合には、M C A L 1 2 0 は、タスクを並列して実行するために、ブースタモジュール140から選択されたブースタノード上にプロセスを生成することができる。

10

【0047】

統合制御は、M C A L 1 2 0 によって提供されるもう一つの機能局面である。統合制御は、M C A L 1 2 0 が、管理機能から収集した情報を用いて通信を制御すること（およびその逆）を意味する。このように、通信機能と管理機能とは別個のものではなく、システム性能全体をさらに最適化するために統合されたものである。これにより、システムの挙動の変更、および、システム上で実行中のアプリケーションの要求の変更に対する動的な調整が可能になる。統合制御は、実際のシステム状態およびアプリケーション要求を考慮するだけでなく、次に起こるモジュールのリクエストまたはインターコネクト使用の予測も行なう。

20

【0048】

図2は、本発明に従うモジュラーコンピューティングシステム200の第1の実施形態を示す。モジュラーコンピューティングシステム200は、クラスタモジュール202と、ブースタモジュール204と、ストレージモジュール206とを含む。

【0049】

クラスタモジュール202は、複数のクラスタノード208（CN）を含む。各クラスタノード208上には、ノードマネージャ209（NM）が設けられる。ノードマネージャ209（NM）は、他のノードマネージャ209と協働して、モジュラーコンピューティングシステム200に通信機能および管理機能を提供する。ブースタモジュール204は、複数のブースタノード210（BN）を含む。各ブースタノード210上には、ノードマネージャ209が設けられる。ノードマネージャ209は、他のノードマネージャ209と協働して、モジュラーコンピューティングシステム200に通信機能および管理機能を提供する。ストレージモジュール206は、複数のストレージノード212（SN）を含む。各ストレージノード212上には、ノードマネージャ209が設けられる。ノードマネージャ209は、他のノードマネージャ209と協働して、モジュラーコンピューティングシステム200に通信機能および管理機能を提供する。

30

【0050】

クラスタノードは、汎用マイクロプロセッサ（例えば、インテルXeonプロセッサ）、メモリ、インターフェース（例えば、ネットワークカード）を有するコンピュータによって形成され得る。ブースタノードは、GPUまたはメニーコアプロセッサ、メモリ、および、インターフェース（例えば、ネットワークカード）を備える処理システムによって形成され得る。ストレージノードは、RAM（ランダムアクセスメモリ）、不揮発性メモリ（例えば、フラッシュメモリ）、SSD（ソリッドステートディスク）、ハードディスクなどであり得るメモリを含む。さらに、ストレージノードは、ノードコントローラ（例えば、マイクロプロセッサ）、およびインターフェース（例えば、ネットワークカード）を含む。

40

【0051】

2つの異なるネットワーキングインフラストラクチャが設けられる。1つ目は、クラスタノード208、ブースタノード210、およびストレージノード212が通信可能になるように、すべてのノード208、210、212を互いに接続する通信ネットワーク22

50

0である。

【0052】

2つ目は、さまざまなノード208、210、および212において設けられたすべてのノードマネージャ209を接続するように提供される管理ネットワーク224である。モジュラーコンピューティングシステム全体の制御を可能にするフロントエンド230が、管理ネットワーク224に取り付けられている。フロントエンド230もまた、ノードマネージャ209を含む。このノードマネージャ209は、他のノードマネージャ209への実際の通信を行なう。

【0053】

異なるモジュールを結合して1つのシステムとするために、図1を参照して説明したモジュラーコンピューティング抽象化層(MCAL)が設けられる。MCALは、通信機能および管理機能を提供する。Par Tec Cluster Competence Center社によるPara Stationを使用すれば、本明細書で説明する機能および特性を有するMCALを実現するための基盤を形成することができる。

10

【0054】

MCALは、個々のノード上のノードマネージャ209によって実現され、モジュール内での通信(モジュール内通信と呼ぶ)と、モジュール間の通信(モジュール間通信と呼ぶ)とを可能にする。そのために、ノードマネージャは、それぞれのファブリックにインターフェースを提供する。例えば、システム全体に対してファブリックが1つの場合もあるし、または、モジュール毎にファブリックが特化されている場合もある(例えば、クラスタモジュール202にはインフィニバンド、ブースタモジュール204にはインテルOmni-Pathファブリック、など)。イーサネット、10G、ギガビット、または共有メモリ通信など、他のインターコネクとも可能である。ノードマネージャ309は、如何なる種類の通信を用いてもよく、可用性、アプリケーションプログラムによる指示、性能、またはコストなどの要素に基づいて、特定の通信方法を選択する。

20

【0055】

共有メモリ通信および共有の仮想アドレス空間通信が採用されてもよい。これらは、2つ以上のプロセスが通信するための非常に高速かつ効率的な方法を提供する。

【0056】

ノードマネージャ209は、一括してノード管理を提供する。ノード管理は、個々のノードを対象とする。例えば、クラスタモジュール202におけるクラスタノード208、ブースタモジュール204におけるブースタノード210、ストレージモジュール206におけるストレージノード212などである。

30

【0057】

さらに、ノードマネージャ209は、到達可能性、レイテンシまたは物理的配置などの要因に応じて、特定のノードを選択する。これは、実行中のアプリケーションに完全に透過的に行なわれ得る。ノードマネージャ209は、システム全体についての情報を収集およびやり取りし、各ノードおよび各インターコネクの状態を絶えず最新の状態に維持している。これにより、ノードマネージャ309は、ノード1つ当たり利用可能なCPUの数、実行中のサービス、利用可能なストレージ、インターコネク帯域幅、空きメモリ、温度、システムクロック速度などの、各ノードおよび各インターコネクの特定の特性について、常時通知される。

40

【0058】

ノードマネージャによって実施されるもう一つの機能は、モジュール管理である。モジュール管理は、ノード管理に類似しており、特定のモジュールについてのノード管理を含み得るが、モジュールに特有の側面を考慮する。モジュール管理は、モジュールが他のモジュールと通信するために用いるインターフェースノード(設けられる場合)の状態など、モジュールに特有の側面に焦点を当てる。ノードと同様に、ノードマネージャ309は、モジュールの利用増大を可能にするために、モジュールの状態を監視する。

【0059】

50

リソース管理は、ノードマネージャ209がモジュラーコンピューティングシステム200において行なうもう一つのタスクである。リソース管理の一部として、ノードマネージャ209は、オフロード機構を実現する。すなわち、例えばクラスタモジュール202上で実行中のプロセスは、ブースタモジュール204中のいくつかのブースタノード210を必要とし得る。そのような場合には、ノードマネージャは、タスクを並列して実行するために、選択されたブースタノード210上にプロセスを生成することができる。

【0060】

統合制御は、ノードマネージャによって実施されるもう一つの機能局面である。統合制御は、ノードマネージャが、管理機能から収集した情報を用いて通信を制御すること（およびその逆）を意味する。このように、通信機能と管理機能とは別個のものではなく、システム性能全体をさらに最適化するために統合されたものである。

10

【0061】

すべてのモジュールにおけるすべてのノードのすべての異なる局面、ならびにすべてのインターコネクト、モジュール内通信、およびモジュール間通信のステータスについてのノードマネージャ209の知識を用いて、システムを作動させる。これにより、システムの挙動の変更、および、システム上で実行中のアプリケーションの要求の変更に対する動的な調整が可能になる。

【0062】

統合制御は、実際のシステム状態およびアプリケーション要求を考慮するだけでなく、次に起こるモジュールのリクエストまたはインターコネクト使用の予測も行なう。

20

【0063】

ノードマネージャ209はハードウェアで実現されてもよく、またはソフトウェアで実現されてもよく、例えば、それぞれのノード上で実行されるデーモンプロセスとして実現されてもよい。ノードマネージャ209は、モジュールおよびシステム全体を大域的に見るために、絶えず情報を収集する。システム上で実行中のアプリケーションは、管理機能と通信することによって、または個々のノードマネージャ209と直接通信することによって、この情報を活用することができる。これは、特定のインターフェースによって実現される。この特定のインターフェースは、ハードウェアであってもよいし、API（アプリケーションプログラミングインターフェース）であってもよい。

【0064】

ノードマネージャ209の通信トラフィックは、モジュール内およびモジュール間通信から厳密に分離されている。これにより、モジュール内およびモジュール間通信は、管理機能によって中断されることなく高速で実行される。さらに、このような分離によって、システムセキュリティを高めることができる。

30

【0065】

実際の通信ネットワーク220の使用は、ノードマネージャ209によって制御される。ノードマネージャ209は、この機能をメッセージパッシングインターフェース(MPI)を介してアプリケーションに提供する。

【0066】

ノードマネージャ209は、管理ネットワーク224を用いて互いに通信する。管理ネットワーク224は、モジュール内およびモジュール間通信に用いられる物理ネットワーク（ここでは、通信ネットワーク220と称する）と同じ物理ネットワークであってもよいし、異なる物理ネットワークであってもよい。例えば、モジュール内およびモジュール間通信は、インフィニバンド、ミリネットなどであってもよく、管理ネットワークは、ノードマネージャ209を接続するイーサネットによって形成されてもよい。

40

【0067】

システムには、管理タスクまたは制御機能がシステムを作動させることを可能にするフロントエンド230を提供するための特殊ノードが含まれてもよい。フロントエンドが通信ネットワーク220に接続されることは必要ではないが、通信ネットワーク220へだけ接続されている。例えばワークステーションコンピュータ、パーソナルコンピュータ(P

50

C)、ラップトップコンピュータ、または任意のスマートコンピューティングデバイスによって形成され得るフロントエンドノードを介して、モジュラーコンピューティングシステム200は管理され得る。モジュラーコンピューティングシステム200上で実行中のアプリケーションは、スケジュール、開始、停止、および監視され得る。これは、バッチ処理によって行なわれ得る。すなわち、手動による介入なく、したがって非対話方式で、モジュラーコンピューティングシステム200上で一連のジョブを実行することによって行なわれ得る。

【0068】

手動による介入なくモジュラーコンピューティングシステム200上で一連のジョブを実行することに加えて、ノードマネージャ309は、対話型コンピューティングを提供し得る。対話型コンピューティングでは、モジュラーコンピューティングシステム200上で実行中のアプリケーションが、その実行中に人からの入力を受け付ける。この入力は、ワークステーションコンピュータ、パーソナルコンピュータ(PC)、ラップトップコンピュータ、または任意のスマートコンピューティングデバイスを介して提供され得る。そのようなシナリオでは、例えばラップトップコンピュータを介してモジュラーコンピューティングシステム200にアクセスし、ラップトップコンピュータがモジュラーコンピューティングシステム200の計算能力を有しているかのように見せることができる。これは、即時または現在進行中であるとユーザが感じるような時間枠内での応答が必要なリアルタイムのアプリケーションに非常に有用である。また、膨大な量の構造化データ、半構造化データ、および非構造化データを情報のマイニングのためにアクセス可能にするビッグデータアプリケーションにおいても、そのような配列を利用することができる。これは、計算能力およびストレージを提供するための基礎となるリソースとしてモジュラーコンピューティングシステム200を有するクラウドコンピューティングサービスにも当てはまる。

【0069】

図3は、本発明に従うモジュラーコンピューティングシステム300の第2の実施形態を示す。モジュラーコンピューティングシステム300は、クラスタモジュール302と、ブースタモジュール304と、ストレージモジュール306とを含む。

【0070】

クラスタモジュール302は、複数のクラスタノード310(CN)を含む。各クラスタノード310上には、ノードマネージャ309(NM)が設けられる。ノードマネージャ309(NM)は、他のノードマネージャ309と協働して、モジュラーコンピューティングシステム300に通信機能および管理機能を提供する。クラスタノード310は、クラスタインターコネクタ312を介して互いに接続されている。

【0071】

ブースタモジュール304は、複数のブースタノード320(BN)を含む。各ブースタノード320上には、ノードマネージャ309が設けられる。ノードマネージャ309は、他のノードマネージャ309と協働して、モジュラーコンピューティングシステム300に通信機能および管理機能を提供する。ブースタノード320は、ブースタインターコネクタ322を介して互いに接続されている。

【0072】

ストレージモジュール306は、複数のストレージノード330(SN)を含む。各ストレージノード330上には、ノードマネージャ309が設けられる。ノードマネージャ309は、他のノードマネージャ309と協働して、モジュラーコンピューティングシステム300に通信機能および管理機能を提供する。ストレージノード330は、ストレージインターコネクタ332を介して互いに接続されている。

【0073】

クラスタノードは、汎用マイクロプロセッサ(例えば、インテルXeonプロセッサ)、メモリ、インターフェース(例えば、ネットワークカード)を有するコンピュータによって形成され得る。ブースタノードは、GPUまたはメニーコアプロセッサ、メモリ、およ

10

20

30

40

50

び、インターフェース（例えば、ネットワークカード）を備える処理システムによって形成され得る。ストレージノードは、RAM（ランダムアクセスメモリ）、不揮発性メモリ（例えば、フラッシュメモリ）、SSD（ソリッドステートディスク）、ハードディスクなどであり得るメモリを含む。さらに、ストレージノードは、ノードコントローラ（例えば、マイクロプロセッサ）、およびインターフェース（例えば、ネットワークカード）を含む。

【0074】

クラスタモジュール302は、クラスタインターコネクタ312とストレージインターコネクタ332とを接続するインターフェースノード340を介して、ストレージモジュール306と通信可能である。ストレージモジュール306は、ストレージインターコネクタ332とブスタインターコネクタ322とを接続するインターフェースノード342を介して、ブスタモジュール304と通信可能である。ブスタモジュール304は、ブスタインターコネクタ322とクラスタインターコネクタ312とを接続するインターフェースノード344を介して、クラスタモジュール302と通信可能である。

10

【0075】

クラスタモジュール302は、例えばインフィニバンドインターコネクタを使用する。一方、ブスタモジュール304は、インターフェースノード344を介してインフィニバンドインターコネクタに接続されたインテルOmni-Pathファブリックを使用する。性能およびスループットを向上させるために、各インターフェースノードは複数のノードによって形成されてもよい。同じことが、他のタイプのモジュールについても当てはまる。

20

【0076】

さまざまなノード310、310、312において設けられたすべてのノードマネージャ309は、同じ通信インフラストラクチャを使用する。モジュラーコンピューティングシステム全体の制御を可能にするフロントエンド350が、クラスタインターコネクタ312に取り付けられている。フロントエンド350もまた、ノードマネージャ309を含む。このノードマネージャ309は、他のノードマネージャ309への実際の通信を行なう。

【0077】

異なるモジュールを結合して1つのシステムとするために、図1を参照して説明したモジュラーコンピューティング抽象化層（MCAL）が設けられる。MCALは、通信機能および管理機能を提供する。Par Tec Cluster Competence Center社によるPara Stationを使用すれば、以下で説明する機能および特性を有するMCALを実現するための基盤を形成することができる。

30

【0078】

MCALは、各ノードにおいてノードマネージャ309によって実現される。プロセス間の通信のために、提供される如何なる種類の通信リンクを用いてもよい。

【0079】

前述のように、ノードマネージャ309はノード管理も提供する。ノード管理は、個々のノードを対象とする。例えば、クラスタモジュール302におけるクラスタノード310、ブスタモジュール304におけるブスタノード320、ストレージモジュール306におけるストレージノード330などである。また、ノードマネージャ309は、到達可能性、レイテンシ、または物理的配置などの要因に応じて、特定のノードを選択することもできる。これは、実行中のアプリケーションに完全に透過的に行なわれ得る。ノードマネージャ309は、常にシステム全体を見渡し、各ノードおよび各インターコネクタの状態を絶えず最新の状態に維持している。ノードマネージャ309は、ノード1つ当たり利用可能なCPUの数、実行中のサービス、利用可能なストレージ、インターコネクタ帯域幅、空きメモリ、温度、システムクロック速度などの、各ノードおよび各インターコネクタの特定の特性について、常時、互いに通知し合う。

40

【0080】

ノードマネージャ309が実施するもう一つの機能は、モジュール管理である。モジュー

50

ル管理は、ノード管理に類似しており、特定のモジュールについてのノード管理を含み得るが、モジュールに特有の側面を考慮する。モジュール管理は、モジュールが他のモジュールと通信するために用いるインターフェースノード（設けられる場合）の状態など、モジュールに特有の側面に焦点を当てる。

【0081】

ノードマネージャ309は、プロセス管理も実施する。実行中のアプリケーションの各々は、複数の個々のプロセスに分割される。複数の個々のプロセスは、開始、分散、および監視する必要がある。同時に、リソースを割り当てる必要がある。これは、プロセスが計算能力およびストレージを必要とすることに対応する。ノードマネージャ309は、プロセスとその要件との間の依存を認識および制御するように構成されている。これに基づいて、ノードマネージャ309は、ブースタモジュール304またはクラスタモジュール302の場合にはノード上にプロセスを作成し、または、ストレージモジュール306の場合にはストレージを割り当てる。プロセスは監視され、依存プロセスはグループとして扱われる。これにより、並列タスクのうちいずれか1つで発生するエラーに迅速に応答することが可能である。したがって、プロセスを速やかに再開させることができ、その結果、実行中のアプリケーションの全体の遅延を最小限に抑えることができる。

10

【0082】

リソース管理は、ノードマネージャ309がモジュラーコンピューティングシステム300において実施するもう一つのタスクである。リソース管理の一部として、ノードマネージャ309は、オフロード機構を実現する。すなわち、例えばクラスタモジュール302上で実行中のプロセスは、ブースタモジュール304中のいくつかのブースタノード320を必要とし得る。そのような場合には、ノードマネージャ309は、タスクを並列して実行するために、選択されたブースタノード320上にプロセスを生成する。

20

【0083】

オフロードの一部として、1つのブースタ上で実行されるように作成されるプロセスは、1つだけでなく複数であってもよい。複数のプロセスは、ホストプロセッサの介入を必要とせずに、モジュール通信ファブリックを用いて互いに直接通信する。

【0084】

また、リソース管理は、例えば異なるプロセス間の接続数を減少させることによって、メモリ消費を低減することも保証する。これは、必要な時にのみ実際に確立するオンデマンド接続を用いることによって行なわれる。これにより、接続1つ当たりのメモリ消費が劇的に低減される。したがって、ノードなどのリソースは仮想的に扱われるだけでなく、仮想割り当ては物理ネットワークにも適用される。

30

【0085】

リソースを割り当てる際に、ノードマネージャ309は、1つのプロセスのリクエストを考慮するだけでなく、他のプロセスのリクエストも考慮する。これにより、統合制御に基づいて通信インフラストラクチャの変更に柔軟に 대응することができる。

【0086】

統合制御は、ノードマネージャ309によって提供されるもう一つの機能局面である。統合制御は、ノードマネージャが、管理機能から収集した情報を用いて通信を制御すること（およびその逆）を意味する。このように、通信機能と管理機能とは別個のものではなく、システム性能全体をさらに最適化するために統合されたものである。

40

【0087】

すべてのモジュールにおけるすべてのノードのすべての異なる局面、ならびにすべてのインターコネクト、モジュール内通信、およびモジュール間通信のステータスについてのノードマネージャの知識を用いて、システムを作動させる。これにより、システムの挙動の変更、および、システム上で実行中のアプリケーションの要求の変更に対する動的な調整が可能になる。

【0088】

統合制御は、実際のシステム状態およびアプリケーション要求を考慮するだけでなく、次

50

に起こるモジュールのリクエストまたはインターコネクト使用の予測も行なう。システム状態の情報は、ノードの動的な割り当てのために用いられる。これは、フィードバックループ登録システムの割り当て決定の広範な影響を用いてリソースの割り当てが行なわれることを含む。その結果、将来の割り当てが相応に調整される。また、システムは、自己学習の方法で、そのような情報を記憶することによって将来のケースの割り当て決定を向上させる。

【 0 0 8 9 】

アプリケーションの主要部分が通常はクラスタモジュール 3 0 2 で開始されるとしても、ノードマネージャは、クラスタモジュール 3 0 2 またはブースタモジュール 3 0 4 上でアプリケーションのどの部分を実行させるかを特定する。決定の根拠は、例えば、実行中のコード性能の監視結果をアプリケーションコード自体で示すものである。これにより、アプリケーションのコンピューティング中に割り当てを調整することができる。

10

【 0 0 9 0 】

ノードマネージャ 3 0 9 はハードウェアで実現されてもよく、またはソフトウェアで実現されてもよく、例えば、それぞれのノード上で実行されるデーモンプロセスとして実現されてもよい。ノードマネージャ 3 0 9 は、モジュールおよびシステム全体を大域的に見るために、絶えず情報を収集する。システム上で実行中のアプリケーションは、管理機能と通信することによって、または個々のノードマネージャ 3 0 9 と直接通信することによって、この情報を活用することができる。これは特定のインターフェースによって実現される。この特定のインターフェースは、ハードウェアであってもよいし、API (アプリケーションプログラミングインターフェース) であってもよい。

20

【 0 0 9 1 】

ノードマネージャの通信トラフィックは、モジュール内およびモジュール間通信から厳密に分離されている。これにより、モジュール内およびモジュール間通信は、管理機能によって中断されることなく高速で実行される。さらに、このような分離によって、システムセキュリティを高めることができる。

【 0 0 9 2 】

実際の通信ネットワークの使用は、ノードマネージャ 3 0 9 によって制御される。ノードマネージャ 3 0 9 は、この機能をメッセージパッシングインターフェース (MPI) を介してアプリケーションに提供する。ノードマネージャ 3 0 9 は、管理ネットワーク (図示せず) を用いる。管理ネットワーク (図示せず) は、モジュール内およびモジュール間通信に用いられる物理ネットワークと同じ物理ネットワークであってもよいし、異なる物理ネットワークであってもよい。例えば、モジュール内およびモジュール間通信は、インフィニバンド、ミリネットなどであってもよく、管理ネットワーク (図示せず) は、ノードマネージャ 3 0 9 を接続するイーサネットによって形成されてもよい。

30

【 0 0 9 3 】

システムには、管理タスクまたは制御機能がシステムを作動させることを可能にするフロントエンドを提供するための特殊ノードが含まれてもよい。フロントエンド 3 5 0 が通信ネットワークに接続されることは必要ではないが、通信ネットワークへだけ接続されている。例えばワークステーションコンピュータ、パーソナルコンピュータ (PC)、ラップトップコンピュータ、または任意のスマートコンピューティングデバイスによって形成され得るフロントエンド 3 5 0 ノードを介して、モジュラーコンピューティングシステム 3 0 0 は管理され得る。モジュラーコンピューティングシステム 3 0 0 上で実行中のアプリケーションは、スケジュール、開始、停止、および監視され得る。これは、バッチ処理によって行なわれ得る。すなわち、手動による介入なく、したがって非対話方式で、モジュラーコンピューティングシステム 3 0 0 上で一連のジョブを実行することによって行なわれ得る。

40

【 0 0 9 4 】

手動による介入なくモジュラーコンピューティングシステム 3 0 0 上で一連のジョブを実行することに加えて、ノードマネージャは、対話型コンピューティングを提供し得る。対

50

話型コンピューティングでは、モジュラーコンピューティングシステム 300 上で実行中のアプリケーションが、その実行中に人からの入力を受け付ける。この入力は、ワークステーションコンピュータ、パーソナルコンピュータ (PC)、ラップトップコンピュータ、または任意のスマートコンピューティングデバイスを介して提供され得る。そのようなシナリオでは、例えばラップトップコンピュータを介してモジュラーコンピューティングシステム 300 にアクセスし、ラップトップコンピュータがモジュラーコンピューティングシステム 300 の計算能力を有しているかのように見せることができる。これは、即時または現在進行中であるとユーザが感じるような時間枠内の応答が必要なリアルタイムのアプリケーションに非常に有用である。また、膨大な量の構造化データ、半構造化データ、および非構造化データを情報のマイニングのためにアクセス可能にするビッグデータアプリケーションにおいても、そのような配列を利用することができる。これは、計算能力およびストレージを提供するための基礎となるリソースとしてモジュラーコンピューティングシステム 300 を有するクラウドコンピューティングサービスにも当てはまる。

10

【図面】

【図 1】

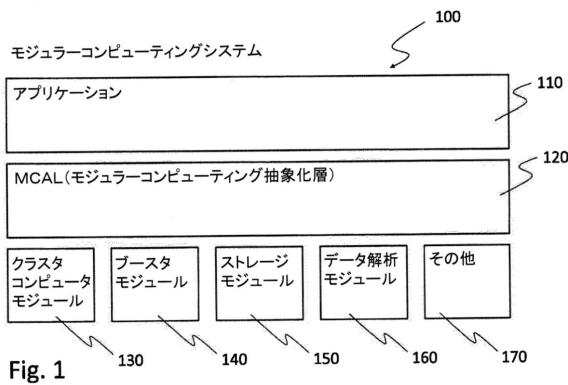
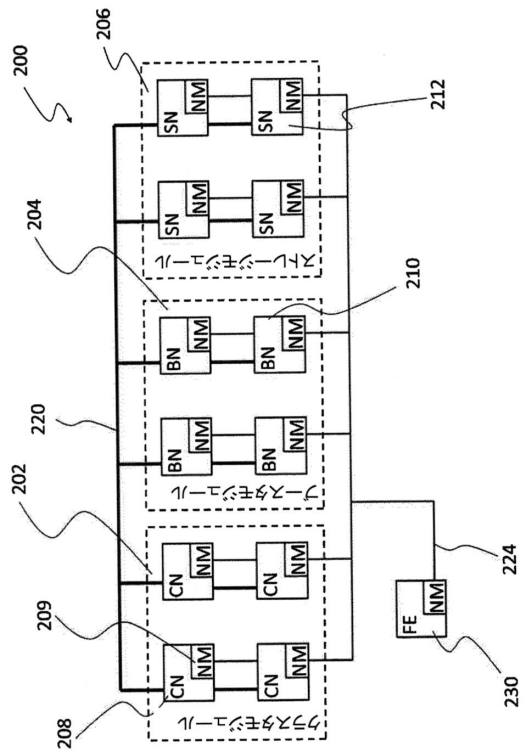


Fig. 1

【図 2】



20

30

Fig. 2

40

50

【 図 3 】

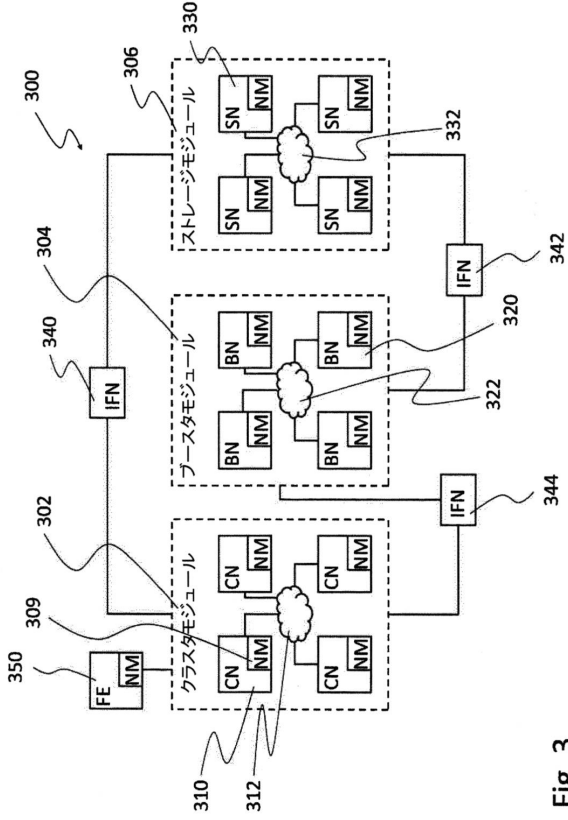


Fig. 3

10

20

30

40

50

フロントページの続き

- (56)参考文献 特表2013-539881(JP,A)
特開2016-012347(JP,A)
特開2012-048330(JP,A)
特表2013-515991(JP,A)

- (58)調査した分野 (Int.Cl., DB名)
G06F 9/50
G06F 15/167
G06F 15/173