

(19) 世界知的所有権機関  
国際事務局



(43) 国際公開日  
2005年12月29日 (29.12.2005)

PCT

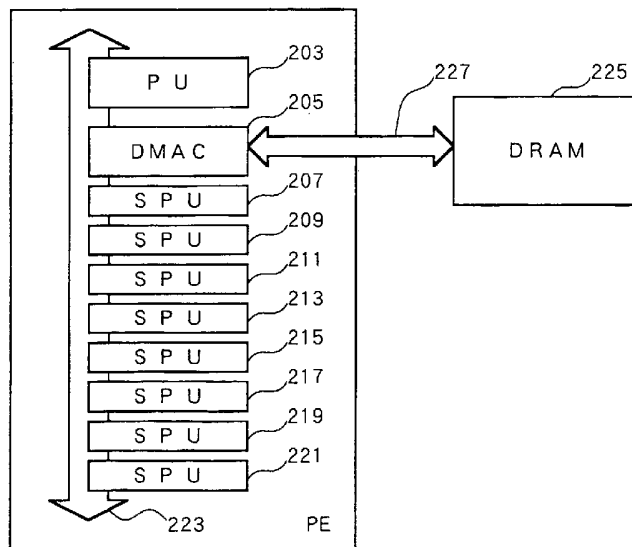
(10) 国際公開番号  
WO 2005/124548 A1

- (51) 国際特許分類: **G06F 9/50**
- (21) 国際出願番号: PCT/JP2005/010986
- (22) 国際出願日: 2005年6月9日 (09.06.2005)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (30) 優先権データ:  
特願2004-177460 2004年6月15日 (15.06.2004) JP
- (71) 出願人 (米国を除く全ての指定国について): 株式会社ソニー・コンピュータエンタテインメント (SONY COMPUTER ENTERTAINMENT INC.) [JP/JP]; 〒1070062 東京都港区南青山二丁目6番21号 Tokyo (JP).
- (72) 発明者; および
- (75) 発明者/出願人 (米国についてのみ): 加藤 裕樹 (KATO, Hiroki) [JP/JP]; 〒1070062 東京都港区南青山二丁目6番21号 株式会社ソニー・コンピュータエンタテインメント内 Tokyo (JP). 前川 博俊 (MAEGAWA, Hiroto) [JP/JP]; 〒1070062 東京都港区南青山二丁目6番21号 株式会社ソニー・コンピュータエンタテインメント内 Tokyo (JP). 石田 隆行 (ISHIDA, Takayuki) [JP/JP]; 〒1070062 東京都港区南青山二丁目6番21号 株式会社ソニー・コンピュータエンタテインメント内 Tokyo (JP).
- (74) 代理人: 鈴木 正剛 (SUZUKI, Seigoh); 〒1050014 東京都港区芝三丁目2番7号 芝NKビル4階 Tokyo (JP).

[続葉有]

(54) Title: PROCESSING MANAGEMENT DEVICE, COMPUTER SYSTEM, DISTRIBUTED PROCESSING METHOD, AND COMPUTER PROGRAM

(54) 発明の名称: 処理管理装置、コンピュータ・システム、分散処理方法及びコンピュータプログラム



(57) Abstract: When performing distributed processing in processing devices connected to a network and a processing managing device for managing the processing devices, it is possible to eliminate the overhead of the processing management device. The processing management device (PU) (203) managing processing devices (SPU) (207) under its control lists up the network address of the SPU (207) and other SPU connected to the network and resource information indicating the current task execution ability of the SPU in a resource list. When one of the SPU transmits a task request to the PU (203), the PU (203) specifies one or more SPU capable of performing the task request in the resource list and requests the specified SPU to execute a task including the execution result specification destination, thereby enabling execution result transmission/reception not using the PU (203).

(57) 要約: ネットワークに接続されている複数の処理装置とこれらの処理装置を管理する処理管理装置とで分散処理を行う際の処理管理装置のオーバー

[続葉有]



WO 2005/124548 A1



(81) 指定国 (表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, KE, KG, KM, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NA, NG, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SM, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW.

(84) 指定国 (表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ,

BY, KG, KZ, MD, RU, TJ, TM), ヨーロッパ (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IS, IT, LT, LU, MC, NL, PL, PT, RO, SE, SI, SK, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

添付公開書類:  
— 国際調査報告書

2文字コード及び他の略語については、定期発行される各PCTガゼットの巻頭に掲載されている「コードと略語のガイダンスノート」を参照。

---

ヘッドを回避する。傘下の複数の処理装置 (SPU) 207等を管理する処理管理装置 (PU) 203が、SPU 207等と、ネットワークに接続されている他のSPUに対するネットワークアドレス及び当該SPUの現在のタスク実行能力を表すリソース情報をリソースリストにリストアップしておく。いずれかのSPUがPU 203に対してタスク要求を送信すると、PU 203は、タスク要求を遂行可能な1又は複数のSPUをリソースリストより特定し、特定したSPUに、その実行結果の指定先を含むタスクの実行を依頼することにより、PU 203を介在させない実行結果の受け渡しを可能にする。

## 明 細 書

処理管理装置、コンピュータ・システム、分散処理方法及びコンピュータプログラム

5

## 技術分野

本発明は、広帯域環境のネットワークに接続される複数のコンピュータによって効率的なタスクの実行を可能にするためのコンピュータ・システム、分散処理方法及びこのような分散処理を可能にするためのコンピュータ装置のアーキテク

10 チャに関する。

## 発明の背景

今日、ネットワークに接続された複数のコンピュータ装置が協働して一つのタスクを分散処理することが一般的になっている。タスクを分散処理する場合、従来は、どのコンピュータ装置にどのタスクを割り当てるかを定めるために、予め

15 ネットワークに接続可能なすべてのコンピュータ装置の処理能力を把握しているサーバ装置の存在が不可欠となる。サーバ装置は、タスクの負荷の大きさと、分散処理を行おうとするときにネットワークに接続されている各コンピュータ装置の余剰処理能力（計算資源）とを特定し、負荷に応じた余剰処理能力を有するコ

20 ンピュータ装置を逐次割り当てていき、その割り当てたコンピュータ装置からタスクの実行結果を受け取る。

サーバ装置を必要とする従来の分散処理方法では、任意の時点でネットワークに接続されたり、非接続になったりするコンピュータ装置の余剰処理能力を、サーバ装置において、迅速に把握することは、非常に困難である。また、サーバ装置が、タスクの分散処理を依頼したコンピュータ装置からその実行結果を受け取

25 って、タスクの依頼元に転送しなければならないため、サーバ装置のオーバーヘッドが大きくなる。そのため、タスクの実行に必要な時間と、ネットワークを介するデータ伝送に必要な時間とが実質的に増加してしまうという問題がしばしば生じていた。

本発明は、このような従来の問題点を解決する分散処理の仕組みを提供することを、主たる課題とする。

#### 発明の開示

5 本発明は、特徴的な処理管理装置、コンピュータ・システム、分散処理管理用のデバイス及びコンピュータプログラムによって上記の課題を解決する。

本発明の処理管理装置は、それぞれ依頼されたタスクを実行し、その実行結果を依頼元の指定先へ送信する機能を備えた複数の処理装置が接続可能なネットワークに接続される処理管理装置であって；前記ネットワークに接続されている処理装置のタスク実行能力を表すリソース情報と当該処理装置とのネットワーク通信を可能にするための通信設定情報とがリストアップされている所定のメモリへのアクセスを可能にする第1管理手段と；ネットワーク通信によりいずれかの処理装置からタスクの実行要求を受け付けたときに当該タスクを実行可能な処理装置を前記メモリにリストアップされているリソース情報より特定し、特定した処理装置についての前記通信設定情報を前記メモリより取得するとともに、前記特定した処理装置と前記タスク要求を行った処理装置の少なくとも一方の処理装置に他方の処理装置の前記通信設定情報を伝達することにより、当該処理装置間のネットワーク通信によるダイレクトな前記実行結果の受け渡しを可能にする第2管理手段と；を有するものである。

20 上記のメモリと、このメモリへの前記リソース情報および通信設定情報の記録を行うリソース情報等保持手段を備えて処理管理装置を構成することができる。

また、タスクは、当該タスクの実行に続く後続タスクの実行依頼とその後続タスクの実行結果の指定先の情報を含むものとすることができる。

前記第2管理手段は、例えば、前記特定した処理装置に対して前記メモリに記録されている当該処理装置の通信設定情報をもとにその実行結果の指定先を含む前記タスクの実行を依頼するように動作する。あるいは、前記第2管理手段は、例えば、前記タスク要求を行った処理装置に、前記特定した処理装置へのダイレクトなタスクの実行を依頼させるように動作する。

本発明の処理管理装置は、複数の処理装置が随時参加することができ且つ脱退

## 3

できる共有空間を前記ネットワーク上に形成する共有空間モジュールを有するものとすることができる。この場合、前記第2管理手段は、前記共有空間モジュールにより形成される共有空間に参加している処理装置の通信設定情報及び現在のリソース情報を当該処理装置より取得し、取得したこれらの情報を前記メモリに  
5 リストアップして使用可能状態にするとともに、前記共有空間から脱退した処理装置について前記メモリにリストアップされている情報を使用不能状態に変更するように動作する。使用不能状態にすることの最も単純な例は、削除することであるが、使用不能のフラグをたてるようにすることもできる。

本発明のコンピュータ・システムは、それぞれ依頼されたタスクを実行し、その実行結果を依頼元の指定先へ送信する機能を備えた処理装置と、この処理装置  
10 と内部バスを介して接続された処理管理装置とを含むものである。前記処理装置及び前記処理管理装置は、それぞれ前記内部バスを介してネットワークに接続されている。

前記処理管理装置は；前記ネットワークに接続されている処理装置のタスク実行能力を表すリソース情報と当該処理装置とのネットワーク通信を可能にするための通信設定情報とがリストアップされている所定のメモリへのアクセスを可能にする第1管理手段と；ネットワーク通信によりいずれかの処理装置からタスク  
15 の実行要求を受け付けたときに当該タスクを実行可能な処理装置を前記メモリにリストアップされているリソース情報より特定し、特定した処理装置についての前記通信設定情報を前記メモリより取得するとともに、前記特定した処理装置と前記タスク要求を行った処理装置の少なくとも一方の処理装置に他方の処理装置の前記通信設定情報を伝達することにより、当該処理装置間のネットワーク通信  
20 によるダイレクトな前記実行結果の受け渡しを可能にする第2管理手段と；を有するものである。

本発明の分散処理管理用のデバイスは、それぞれ依頼されたタスクを実行し、その実行結果を依頼元の指定先へ送信する機能を備えた複数の処理装置が接続可能なネットワークに接続されるコンピュータ・システムに搭載されるデバイスであって；所定のコンピュータプログラムを実行することにより、前記コンピュータ・システムを；前記ネットワークに接続されている処理装置のタスク実行能力  
25

## 4

を表すリソース情報と当該処理装置とのネットワーク通信を可能にするための通信設定情報とがリストアップされている所定のメモリへのアクセスを可能にする第1管理手段；ネットワーク通信によりいずれかの処理装置からタスクの実行要求を受け付けたときに当該タスクを実行可能な処理装置を前記メモリにリストアップされているリソース情報より特定し、特定した処理装置についての前記通信設定情報を前記メモリより取得するとともに、前記特定した処理装置と前記タスク要求を行った処理装置の少なくとも一方の処理装置に他方の処理装置の前記通信設定情報を伝達することにより、当該処理装置間のネットワーク通信によるダイレクトな前記実行結果の受け渡しを可能にする第2管理手段；として動作させるものである。

本発明の分散処理方法は、それぞれ依頼されたタスクを実行し、実行結果を依頼元の指定先へ送信する機能を備えた複数の処理装置と、各処理装置との間でネットワークを通じて通信を行う処理管理装置との協働で行う分散処理方法であって、前記処理管理装置が、前記ネットワークに接続されている処理装置のタスク実行能力を表すリソース情報および当該処理装置へのアクセスを可能にするための通信設定情報を当該処理装置より取得し、取得したこれらの情報を所定のメモリにリストアップする段階と、いずれかの前記処理装置が前記処理管理装置に対してタスク要求を送信する段階と、前記タスク要求を受信した前記処理管理装置が、受信したタスク要求を遂行可能な1又は複数の前記処理装置を前記メモリの記録情報より特定し、特定した処理装置に、その実行結果の指定先を含むタスクの実行を依頼する段階とを有し、前記処理管理装置を介在させない前記実行結果の受け渡しを可能にする方法である。

本発明が提供するコンピュータプログラムは、コンピュータを、それぞれ依頼されたタスクを実行し、その実行結果を依頼元の指定先へ送信する機能を備えた複数の処理装置が接続可能なネットワークに接続される処理管理装置として動作させるためのプログラムであって；前記コンピュータを；前記ネットワークに接続されている処理装置のタスク実行能力を表すリソース情報と当該処理装置とのネットワーク通信を可能にするための通信設定情報とがリストアップされている所定のメモリへのアクセスを可能にする第1管理手段；ネットワーク通信により

いずれかの処理装置からタスクの実行要求を受け付けたときに当該タスクを実行可能な処理装置を前記メモリにリストアップされているリソース情報より特定し、特定した処理装置についての前記通信設定情報を前記メモリより取得するとともに、前記特定した処理装置と前記タスク要求を行った処理装置の少なくとも一方の処理装置に他方の処理装置の前記通信設定情報を伝達することにより、当該処理装置間のネットワーク通信によるダイレクトな前記実行結果の受け渡しを可能にする第2管理手段；として機能させるためのコンピュータプログラムである。

#### 図面の簡単な説明

- 10 図1は、本発明が適用されるコンピュータ・システムの全体図である。  
図2は、すべてのコンピュータ装置に共通となるPEの構造説明図である。  
図3は、複数のPEを有するBEの構造説明図である。  
図4は、SPUの構造説明図である。  
図5は、ネットワークを伝送するパケットの構造説明図である。
- 15 図6は、パケットに含まれるDMAコマンドの例を示した図であり、(a)はSPUのローカル・メモリへのプログラム等のロードを指示するためのコマンドを含むDMAコマンド、(b)はプログラムをキックさせるためのキックコマンドを含むDMAコマンドの例を示す。  
図7は、多数のSPUから構成される情報統合体の説明図である。
- 20 図8は、PUに形成される機能及びテーブル等の説明図である。  
図9は、クラスタテーブル及びリソースリストの利用形態の説明図である。  
図10は、リソースリスト、SPUステータステーブル及びクラスタリストの関係を示した説明図である。  
図11は、PUに形成されるタスクマネージャにおける処理手順図である。
- 25 図12は、PUに形成されるリソースアロケータにおける処理手順図である。  
図13は、タスクマネージャを用いた全体的な処理概要の説明図である。  
図14は、タスクマネージャを用いた処理概要の他の説明図。

発明を実施するための最良の形態

まず、本発明が適用されるコンピュータ・システムの概要を説明する。

#### <ネットワーク型コンピュータ・システムの概要>

図1は、本発明が適用されるコンピュータ・システム101の全体図が示されている。このコンピュータ・システム101は、ネットワーク104を含んでい  
5 る。ネットワーク104の例としては、ローカル・エリア・ネットワーク(LAN)、インターネットのようなグローバルネットワーク、あるいは他のコンピュータ・ネットワークが挙げられる。

ネットワーク104には、複数のコンピュータ装置の各々がそれぞれ任意のタイ  
10 ミングで接続して、他のコンピュータ装置と双方向の通信を行うことができる。コンピュータ装置の例としては、パーソナルコンピュータ106、サーバコンピュータ108、通信機能付のゲームコンソール110、PDA112及びその他の有線または無線コンピュータとコンピューティング・デバイスなどが含まれる。

各コンピュータ装置は、それぞれ、共通の構造を持つプロセッサ・エレメント  
15 (以下、「PE」)を有している。これらのPEは、命令セット・アーキテクチャ(ISA)がすべて同じで、同じ命令セットに従って所要の処理を実行できるものである。個々のコンピュータ装置に含まれるPEの数は、そのコンピュータ装置がタスクを実行する上で必要な処理能力によって決められる。

コンピュータ装置のPEがこのように均質な構成を有することから、コンピュータ・システム101におけるアダプタビリティを改善することができる。また、  
20 各コンピュータ装置が、各々のPEのうち1つまたはそれ以上、またはPEの一部を用いて、他から依頼されたタスクを実行できるようにすることにより、タスクをどのコンピュータ装置で実行するかはさほど重要ではなくなる。依頼されたタスクの実行結果を、タスクを依頼したコンピュータ装置あるいは後続タスクを実行するコンピュータ装置を指定先として伝達するだけで足りる。そのため、個々の  
25 タスクは、ネットワーク104に接続されている複数のコンピュータ装置の間で分散実行することが容易になる。

各コンピュータ装置が、共通の構造のPEと共通のISAとを有するので、コンピュータ装置間の互換性を達成するためのソフトウェアの追加層の計算上の負担が回避される。また、異質なネットワークの混在という問題の多くを防ぐこと

ができる。したがって、このシステム101は、広帯域処理の実現が可能となる。

### <コンピュータ装置のアーキテクチャ>

次に、コンピュータ装置のアーキテクチャを明らかにする。まず、各コンピュータ装置が有するPEの構造例を図2を参照して説明する。

- 5 図2に示されるように、PE201は、プロセッサ・ユニット(PU)203、ダイレクト・メモリ・アクセス・コントローラ(DMAC)205、PUとは異なる複数のプロセッサ・ユニット(SPU)、すなわち、SPU207、SPU209、SPU211、SPU213、SPU215、SPU217、SPU219、SPU221を含んでいる。PEバス223は、SPU、DMAC205、
- 10 PU203を相互に接続する。このPEバス223は、従来型のアーキテクチャなどを備えていても良いし、あるいは、パケット交換式ネットワークとして実現されても良い。パケット交換式ネットワークとして実現される場合、より多くのハードウェアが必要となり、その一方で、利用可能な帯域幅が増加する。

- PE201は、デジタル論理回路を実現する様々な方法を用いて構成可能であるが、単一の集積回路として構成されることが望ましい。PE201は、高帯域メモリ接続部227を介して、ダイナミック・ランダム・アクセス・メモリ(DRAM)225に接続される。DRAM225は、PE201のメイン・メモリとして機能する。なお、メインメモリは、DRAM225であることが望ましいが、
- 15 10 15 20 25
- 他の種類のメモリ、例えばスタティック・ランダム・アクセス・メモリ(SRAM)としての磁気ランダム・アクセス・メモリ(MRAM)、光メモリまたはホログラフィ・メモリなどをメインメモリとして用いることもできる。DMAC205は、PU203及びSPU207等によるDRAM225へのダイレクト・メモリ・アクセス(DMA)を実現する。DMAの実現手法の例としては、クロスバー・スイッチをその前段に設けることなどが挙げられる。

- 25 PU203は、DMAC205に対してDMAコマンドを出すことによりSPU207等の制御を行う。その際、PUは、独立したプロセッサとしてSPU207等を扱う。したがって、SPU207等による処理を制御するために、PUは、遠隔手順呼出しに類似したコマンドを使用する。これらのコマンドは“遠隔手順呼出し(RPC)”と呼ばれる。PU203は、一連のDMAコマンドをDM

AC205へ出すことによりRPCを実行する。DMAC205は、依頼されたタスクの実行に対してSPUが必要とするプログラム（以下、「SPUプログラム」）とそれに関連するデータとをSPUのローカル・ストレージへロードする。PU203は、その後、SPUへ最初のキック・コマンドを出し、SPUプログラムを実行させる。

PU203は、必ずしもネットワーク型のプロセッサである必要はなく、スタンド・アローン型の処理が可能な標準的なプロセッサを主たるエレメントとして含むものであって良い。このプロセッサがDRAM225あるいは図示しない読み出し専用メモリに記録されている本発明のコンピュータプログラムを読み込んで実行することにより、通信モジュールを形成するとともに、自己が管理可能なSPU207, 209, 211, 213, 215, 217, 219, 221、他のPEあるいは後述するBEのPEが管理するSPUの管理を行うための各種機能を形成する。これらの機能については、後で詳しく述べる。

図2に示したPE201のようないくつかのPEを結合して処理能力を向上させることもできる。例えば、図3に示されるように、1以上のチップ・パッケージなどの中に複数のPEを結合してパッケージ化し、これにより単一のコンピュータ装置を形成しても良い。このような構成のコンピュータ装置は、広帯域エンジン(BE)と呼ばれる。

図3の例では、BE301には、4つのPE(PE303、PE305、PE307、PE309)が含まれる。これらのPE間の通信は、BEバス311を介して行われる。また、これらのPEと共用DRAM315とが、広帯域メモリ接続部313によって接続されている。BEバス311の代わりに、BE301のPE間の通信は、共用DRAM315とこの広帯域メモリ接続部313とを介して行うことができる。BE301において、1つのPEが有する複数のSPUのすべては、それぞれ独立に、但し、一定の制限下で、共用DRAM315内のデータへアクセスすることができる。

入力/出力インターフェース(I/O)317と外部バス319とは、BE301と、ネットワーク104に接続されている他のコンピュータ装置との間の通信インターフェースとなるものである。I/O317は、例えばプロセッサ等の能

動素子を含んで構成され、ネットワーク104とBE301内の各PE303、305、307、309との間の通信を制御するほか、ネットワーク104からの種々の割込をも受け付け、これを該当するPEに伝達する。そこで、以下の説明では、これらを総称して「ネットワークインタフェース」と称する場合がある。

- 5 このようなネットワークインタフェースは、必ずしもBE301に内蔵されていなくても良く、ネットワーク104上に設けられていても良い。

なお、図2及び図3では図示していないが、画像処理を行うためのピクセル・エンジン、画像用キャッシュ、ディスプレイへの表示コントローラ等を含んで、PE201あるいはBE301を構成することもできる。このような付加的な機能は、後述するリソース管理において重要となる。

PE又はBEに含まれるSPUは、単一命令、複数データ(SIMD)プロセッサであることが望ましい。PU203の制御によって、SPUは、並列的かつ独立に、PEを通じて依頼されたタスクを実行し、その実行結果を依頼元が指定するSPUに向けて出力する機能を実現する。図2の例では、1つのPE201が  
15 8個のSPUを有するが、SPUの数は、必要とする処理能力に応じて任意であって良い。

次に、SPUの構造を詳しく説明する。SPUは、図4に例示される構造を有する。SPU402には、SPU内で処理されるデータ及びSPU内で実行される各種スレッドなどを格納するためのローカル・ストレージとなるローカル・メモリ406、レジスタ410、複数の浮動小数点演算ユニット412及び複数の整数演算ユニット414が含まれる。浮動小数点演算ユニット412と整数演算ユニット414の数は、予め必要となる処理能力に応じて、任意に決めることができる。

ローカル・メモリ406は、好適にはSRAMとして構成される。SPU402には、さらに、各種スレッドの実行結果、タスク依頼あるいはタスクの実行結果の受け渡しを行うためのバス404が含まれる。SPU402にはさらに内部バス408、420、418が含まれる。内部バス408は、ローカル・メモリ406とレジスタ410とを接続する。内部バス420、418は、それぞれ、レジスタ410と浮動小数点演算ユニット412との間、及び、レジスタ410

と整数演算ユニット414と間を接続する。ある好ましい実施の態様では、レジスタ410から浮動小数点演算ユニット412または整数演算ユニット414へのバス418と420の幅は、浮動小数点演算ユニット412または整数演算ユニット414からレジスタ410へのバス418と420の幅よりも広い。レジスタ410から浮動小数点演算ユニットまたは整数演算ユニットへの上記バスの広い幅によって、レジスタ410からのより広いデータ・フローが可能になる。

なお、図示を省略したが、コンピュータ装置は、絶対タイマを使用する。絶対タイマはSPUとPEの他のすべてのエレメントへクロック信号を出力する。このクロック信号はこれらのエレメントを駆動するクロック信号に依存せず、かつ、このクロック信号より高速である。この絶対タイマによって、SPUによるタスク・パフォーマンスのためのタイム・バジェット（割り当て時間）が決定される。このタイム・バジェットによって、これらのタスクの完了時間が設定されるが、この時間は、SPUによるタスク処理に必要な時間よりも長い時間になる。その結果、個々のタスクについて、タイム・バジェットの範囲内に、ビジーな時間とスタンバイ時間とが存在することになる。SPUプログラムは、SPUの実際の処理時間にかかわらず、このタイム・バジェットに基づいて処理を行うように作成される。

#### <通信データ>

以上のように構成されるコンピュータ装置は、ネットワーク104に接続されている任意のタイミングで、ネットワーク104上の他のコンピュータ装置に対してタスクを依頼したり、あるいは、依頼されて自己が行ったタスクの実行結果等を伝達するための通信データを生成する。この通信データには、タスクの依頼内容、タスクを実行するための処理プログラム、データ、実行結果あるいは後続タスクの内容を伝達する指定先の通信設定情報など、様々な種類の情報が含まれる。この実施形態では、これらの様々な種類の情報を所定のデータ構造のもとでパケット化し、このパケットをコンピュータ装置相互で受け渡しするようにする。

パケットの受け渡しは、通常は、DRAMを介して行われる。例えば、コンピュータ装置が図2のようなPEを一つだけ有する場合はDRAM225、コンピュータ装置が図3のようなBEを有する場合は共用DRAM315から、それぞ

れ、図4に示した各SPUのローカル・メモリ406に読み出され、そのSPUによって直接処理される。

なお、上記の場合、SPUのローカル・メモリ406に、プログラム・カウンタと、スタックと、処理プログラムを実行するために必要な他のソフトウェア・  
5 エレメントとが含まれるようにしても良い。

パケットの受け渡しに際しては、常にDRAMを介さなければならないという  
ものでもない。すなわち、コンピュータ装置が図2のようなPEを一つだけ有す  
る場合はSPUのローカル・メモリ406に直接ロードされるようにし、他方、  
コンピュータ装置が図3のようなBEを有する場合は、I/O317から各SP  
10 Uのローカル・メモリ406に直接ロードされるようにして、それぞれのSPU  
によって処理されるようにしても良い。

図5は、パケットの構造例を示している。パケット2302の中には、ルート  
選定情報セクション2304と本体部分2306とが含まれる。ルート選定情報  
セクション2304に含まれる情報は、ネットワーク104のプロトコルに依っ  
15 て決められる。ルート選定情報セクション2304の中には、ヘッダ2308、  
宛先ID2310、ソースID2312及び応答ID2314が含まれる。宛先  
ID2310には、タスクの実行結果又は後続タスクを実行させるための当該パ  
ケットの伝送先となるSPUの識別子が含まれる。

SPUの識別子は、それが図2に示される一つのPEに所属する場合にはPE  
20 及びSPUのセットから成るネットワーク・アドレス、それが図3に示されるB  
Eの一つのPEに所属する場合にはBE、PE及びSPUのセットから成るネッ  
トワーク・アドレスが含まれる。TCP/IPプロトコルの下でのネットワーク・  
アドレスはインターネット・プロトコル(IP)アドレスである。ソースID23  
12には、当該パケットの伝送元となるSPUの識別子が含まれる。このソース  
25 ID2314により識別されるSPUとからパケットが生成され、ネットワーク  
104に向けて発信される。応答ID2314には、当該パケットに関するクエ  
リとタスク実行結果又は後続タスクを伝送する指定先のSPUの識別子が含ま  
れる。

パケットの本体部分2306には、ネットワークのプロトコルとは無関係の情

報が含まれる。図5の破線で示した分解部分は、本体部分2306の細部を示している。本体部分2306のヘッダ2318によってパケット本体の開始領域が識別される。パケットインターフェース2322には、当該パケットを利用する上で必要となる各種情報が含まれる。これらの情報の中には、グローバルな一意

5 的ID2324と、タスク実行のために指定されるSPU2326と、専有メモリサイズ2328と、前回処理されたSPUで付されたパケットID2330とが含まれる。SPUは、依頼されたタスクを実行した際に、このパケットID2330を更新してネットワーク104に発信することになる。グローバルな一意

10 的ID2324は、ネットワーク104全体を通じて当該パケット2302を一意的に識別するための情報である。このグローバルな一意ID2324は、ソースID2312(ソースID2312内のSPUの一意的識別子など)と、パケット2302の作成または伝送の時刻と日付とに基づいて作成される。

専有メモリサイズ2328によって、タスクの実行に必要なDRAMと関連する必要なSPU内に、他の処理から保護されたメモリサイズが確保される。なお、

15 専有メモリサイズ2328は、必要なデータ構造としては不可欠なものではなく、予めメモリサイズが指定されている場合には、不要となる。前回のパケットID2330によって、シーケンシャルな実行を要求する1グループのタスクにおける前回のタスクを実行したSPUの識別子を事後的に把握可能になる。

実行セクション2332の中には、当該パケットのコア情報が含まれる。この

20 コア情報の中には、DMAコマンド・リスト2334と、タスクの実行に必要な処理プログラム2336と、データ2338とが含まれる。処理プログラム2336には、SPUプログラム2360、2362などの、このパケット2302を受け取ったSPUによって実行されるプログラムが含まれ、データ2338にはこれらのSPUプログラムにより処理される各種データが含まれる。

25 なお、タスク依頼先のSPUにタスク実行に必要な処理プログラムが存在する場合には、処理プログラム2336は不要となる。

DMAコマンド・リスト2334には、プログラムの起動に必要な一連のDMAコマンドが含まれる。これらのDMAコマンドにはDMAコマンド2340、2350、2355、2358が含まれる。

## 1 3

これらのDMAコマンド2340、2350、2355、2358は、例えば、  
図6に示されるようなものである。

すなわち、図6(a)に示されるように、DMAコマンド2340には、VID  
コマンド2342が含まれる。VIDコマンド2342は、DMAコマンドが  
5 出されたときに、DRAM225の物理IDに対して対応づけられるSPUのバ  
ーチャルIDである。DMAコマンド2340にはロード・コマンド2344と  
アドレス2346も含まれる。ロード・コマンド2344は、SPUにDRAM  
225から特定の情報を読み出してローカル・メモリ406へ記録するように命  
令するためのものである。アドレス2346によってこの特定情報を含むDRA  
10 M225内のバーチャル・アドレスが与えられる。この特定情報は、処理プログ  
ラム2336、データ2338、あるいはその他のデータなどであって良い。D  
MAコマンド2340には、また、ローカル・メモリ406のアドレス2348  
が含まれる。このアドレスによって、パケットに記録されているすべての情報を  
ロードできそうなローカル・メモリ406のアドレスが識別される。DMAコマ  
15 ンド2350についても同様となる。

図6(b)に例示されるDMAコマンド2355は、キック・コマンドを含む  
ものである。「キック・コマンド」とは、上述したように、PUから指定先となる  
SPUに出される当該パケット内のSPUプログラムによる実行を開始させるた  
めのコマンドである。このDMAコマンド2355には、VIDコマンド235  
20 2と、上述したキック・コマンド2354と、プログラムカウンタ2356とが  
含まれる。ここにいうVIDコマンド2352は、キックすべき対象SPUを識  
別するためのものであり、キック・コマンド2354は関連するキック・コマン  
ドを与えるものであり、プログラムカウンタ2356は、SPUプログラムの実  
行用プログラムカウンタのためのアドレスを与えるものである。DMAコマンド  
25 2358についても同様の内容となる。

### <運用形態>

次に、図1に示したコンピュータ・システム101の運用形態の一例を説明す  
る。

上述したように、ネットワーク104に接続される各コンピュータ装置は、共

通構造のPEと共通のISAとを有する。そのため、ネットワーク104に接続されているコンピュータ装置間のアーキテクチャの相違は吸収され、ネットワーク104には、図7に示されるように、各PEに含まれる多数のSPUの各々があたかも情報処理の細胞（Cell）のように機能する大規模情報処理統合体WOが形成される。

大規模情報処理統合体WOにおける個々のSPU512は、物理的には自己が所属するPU501, 502, 503, 504, 50n, 50mによって管理され、単独のSPUとして動作したり、そのPU内の他のSPUとグループ化されて協働で動作したりするが、論理的には、PUによる壁はなく、異なるPUのもとで管理される他のSPUとの間でグループ化が可能である。このような形態でグループ化される場合、一つのタスクをグループに属する複数のSPUで分散処理することにより実行することができる。

分散処理の一形態として、グループ間で共通にアクセス可能な共有空間を構築し、あるコンピュータ装置を操作するユーザが、そのコンピュータ装置のいくつかのSPUのリソースを通じて、共有空間にアクセスすることができる。図7において、G01～G04は、それぞれ、大規模情報処理統合体WOにおける複数のSPUのグループを示している。

このような大規模情報処理統合体WOにおける効率的な分散処理を可能にするため、PUには、種々の機能が形成される。図8は、本発明のプログラムとの協働によりPUに形成される機能及び各種リスト並びにテーブルの説明図である。

PUには、タスクマネージャ601と、リソースアロケータ603と、クラスタマネージャ605と、スケジューラ607と、共有空間モジュールの一例となる共有空間マネージャ609と、通信モジュール611とが形成される。通信モジュール611は、内部のSPUとの間でPEバス233（図2）を介して双方向通信を行うとともに、PEバス233及びネットワーク104を介して他のPU（PE/BE）のSPUとの間での双方向通信の手順を制御する。通信モジュール611は、ルーティングデータ及びプロトコルデータを記録したテーブル711を参照することにより、上記制御を行う。

クラスタマネージャ605は、通信モジュール611及びタスクマネージャ6

01を介して通信可能なすべてのSPUをクラスタ化するための処理を行う。この処理の内容を具体的に説明すると、まず、各SPUが処理可能なタスクの種類と現在の処理能力値（処理能力を数値化したデータ）を、例えばSPUにおいて実行されるモニタリング・スレッドと、タスクマネージャ601により実行されるモニタリング・スレッドとを通じて取得する。そして、各SPUにおいて実行可能なタスクの種類と各々のSPUの処理能力値の組み合わせを、その実行が想定されるタスクの種類又はサイズ毎にクラスタ化してクラスタリスト705にリストアップする。

クラスタリスト705の内容は、例えば図9の上段に示されたものとなる。クラスタの最低単位は、1つのSPU7050となる。図9の例では、それぞれ1つの特定のSPUからなる複数種類のSPUクラスタ7051と、複数のSPUの組み合わせからなるSPUクラスタ7053、7055、7057、7059、7061、7063がリスト化されている。それぞれ、1クラスタ毎に、SPUを識別するためのSPU\_idとそのSPUが所属するPUの識別子（PU\_ref）と、割当の可否を表す識別情報（valid/invalid）とが記録される。「valid」は割当可、「invalid」は割当不可を表す。

図10の上段を参照すると、4つのSPUの組み合わせからなるSPUクラスタの具体例が示されている。SPUクラスタは、SPUの数（size）、タスクサイズ（length）に応じて予め定められているテーブル7031に基づいて作成される。図10上段の例では、SPUが4つなので、テーブル7031の「SIZE」が「4」のデータに基づいてSPUクラスタが生成されている。クラスタリスト705の記録内容は、随時、更新される。テーブル7031におけるSPUの数（size）は、任意の数値ではなく、全体のクラスタ数に対する割合を表す関数等で表現されていても良い。

SPUステータステーブル701は、例えば図10後段に示される内容のものである。すなわち、あるクラスタの「PU\_ref」によって識別されるPU3302におけるSPU\_id毎にそのステータス情報（status）がSPUステータステーブル7011に記録される。ステータス情報は、そのSPUが現在クラスタ化されているかどうか、稼働状態にあるかなどを表す情報である。「busy」は

稼働状態であり「none」はそれ故にバックポインタが不要であることが示されている。

タスクマネージャ601は、図11の手順で処理を行う。すなわち、SPUからタスク依頼（タスク要求）があったときに（TS1001:Yes）、依頼されたタスクを実行可能なSPUクラスタの割当をリソースアロケータ603に依頼する（TS1002）。リソースアロケータ603から、割り当てられたSPUクラスタに属するSPUの情報が通知されると（TS1003:Yes）、割り当てられたすべてのSPUのSPU\_idをスケジューラ607に通知して、そのSPUをリザーブさせる（TS1004）。これは、SPUの動作を他の処理用と競合させないようにするためである。

タスクに記述された1個以上の処理プログラム間の呼び出し関数をつなぐデータパスと割当されたSPUの位置情報から生成された関数テーブル707を参照して、SPUに実行させるタスクの種類に応じた実行結果の指定先となるデータパスを検索する（TS1005）。そして、最初のタスクを実行するSPUに向けて、当該タスクの依頼内容とデータパスとを送信する（TS1006）。ここにいうタスクの記述例としては、以下のようなものがある。

<映像処理>

```

    <program="ビデオ用プログラム" name="video" />
    <program="オーディオ用プログラム" name="audio" />
20  <program="パケット化プログラム" name="packet" />
    <sub from="video.Video_OUT" to="packet.Video_IN" />
    <sub from="audio.Audio_OUT" to="packet.Audio_IN" />
    <data file="video_data" to="video.Data_IN" />
    <data file="audio_data" to="audio.DATA_IN" />
25  <data file="packet_data" from="packet.DATA_OUT" />

```

</映像処理>

この際のタスクの依頼内容とデータパスの送信は、PU、すなわちタスクマネージャ601からダイレクトに行ってもよく、データパス等をタスク要求を行ったSPUのローカル・メモリに記録して、そのSPUにてパケットを作成し、こ

のパケットをPEバスを介してネットワークに送信するようにしても良い。

また、変形例として、割り当てられたSPUのSPU\_idより特定されるSPUに向けて、タスク要求を行ったSPUのデータバスとタスク依頼の内容とを伝達し、当該特定されるSPUからタスク要求を行ったSPUにダイレクトに通  
5 信路を確立するようにしても良い。

タスクマネージャ601は、また、通信モジュール611を通じて、自己に属するSPUについてのリソース情報、ステータス情報、スケジュール情報を他のBE、PEあるいはPEに伝達するとともに、他のBE、PEあるいはPEに所属するSPUについてのこれらの情報を取得する。取得したステータス情報はリ  
10 ソースアロケータ603を介して、あるいは、タスクマネージャがダイレクトに、SPUステータステーブル701に記録する。リソース情報も同様にしてリソースリスト703に記録する。

リソース情報は、個々のSPUが現在どのような種類の処理を行うことができ、且つ、その処理能力値がどの程度かを表す情報である。リソースリスト703に  
15 は、このリソース情報が、SPU\_id毎に記録される。上述したように、SPUは、すべてのコンピュータ装置において共通構造であるが、SPUが所属するPEに付加されている処理エンジンの有無、及び、そのような処理エンジンの種類により、SPUが可能となる処理の種類が異なる場合がある。SPUプログラムさえロードすれば、どのような種類の処理あっても可能な場合、処理の種類は、  
20 図9に示されるように「all」となる。「all」以外の特定の処理の種類がある場合は、画像処理のみ、音声処理のみ、パッケージングのみ・・・のように記述される。処理能力値がすべてのSPUにおいて共通で、その値が既知の場合には、処理能力値については、リストアップを省略することができる。

スケジュール情報は、上述したように、SPUをリザーブするための情報である。  
25 る。スケジューラ607は、このスケジュール情報をSPU毎に図示しないタイムテーブル等に記録し、必要に応じてタスクマネージャ601に記録情報を通知する。

リソースアロケータ603は、図12に示した手順で処理を行う。すなわち、タスクマネージャ601を通じて割当依頼があると(AL1001:Yes)、リ

ソースアロケータ603は、タスクサイズ及びタスクの種類を解析する（AL1002）。そして、クラスタリスト705にリストアップされているSPUクラスタから、タスクの実行に最適なSPUクラスタを検索する（AL1003）。

最適なSPUクラスタの検索基準としては、例えば、以下のようなものがある。

- 5 (1) 検索の引数がタスクサイズのみとなる場合、すなわち、処理能力値がすべてのSPUで共通の場合、ごく簡易な例としては、タスクサイズが一致するクラスタリストを選択し、クラスタリスト内で識別情報が valid である先頭のクラスタを割り当てる。ネットワーク距離に基づいて最適かどうかを判定することもできる。この場合は、クラスタを構成するSPU間のすべての組み合わせについて、
- 10 ネットワーク距離（相手のSPUに到達するまでに経由しなければならないホストの数）を算出し、ネットワーク距離の最大値が小さい順にクラスタリストを作成し、クラスタリスト内で識別情報が valid である先頭のクラスタを割り当てる。
- (2) 検索の引数がタスクサイズとタスクの種類になる場合は、例えば、タスクサイズに基づいて、クラスタリストを選択し、次いで、識別情報が valid である
- 15 各クラスタについてクラスタリストの先頭から、必要なタスク処理を備えたSPUが存在するかを順にチェックする。そして、最初に必要なタスク処理を実行できるSPU群を有するクラスタを返す。

- 以上のようにして最適なSPUクラスタを特定すると、リソースアロケータ603は、SPUステータステーブル701を参照して、特定したSPUクラスタ
- 20 に属するSPUのステータスを確認する（AL1004）。ステータスが「busy」等のために割当可能でなかった場合は（AL1005:No）、クラスタリスト705から他のSPUクラスタを特定し（AL1006）、AL1004の処理に戻る。割当可能であった場合は（AL1005:Yes）、割当可能なすべてのSPUをタスクマネージャ601に通知する（AL1007）。

- 25 共有空間マネージャ609は、共有空間管理データをもとに、複数のSPUが任意に参加し、任意に脱退することができる共有空間を、図7に示した大規模情報処理統合体WOの一部の領域に形成する。共有空間管理データは、例えば、複数のSPUが同時期にアクセス可能な領域のアドレス等である。共有空間を形成したときは、上記のステータス情報、リソース情報、スケジュール情報を共有す

る必要がある。そこで、共有空間マネージャ609は、SPUステータステーブル701及びリソースリスト703を、参加しているSPUを管理しているPUに伝達し、記録内容の更新時期を同期させるようにする。クラスタリスト705をもPU間で相互に受け渡すようにしても良い。

- 5 PUは、共有空間に参加しているSPUの識別情報、各SPUの通信設定情報、各SPUの現在のリソース情報等を所定のメモリ領域にSPUの識別情報毎にリストアップして使用可能状態 (valid) とする。他方、共有空間から脱退したSPUについては、リストアップされている情報を使用不能状態 (invalid) に変更する。
- 10 図13は、上記の全体的な処理の例を示した図である。この図において、網掛部分は、クラスタ化されたSPUを表している。リソースアロケータ5007は、本来、各PUに備えられるものであるが、複数のPUに跨る共有空間での処理なので、便宜上、共通の1つのリソースアロケータとして、便宜上、PUの外部において説明する。
- 15 この例では、大規模情報処理統合体WOを構成する、ある1つのSPU4001から、PU(#B)5003に、ビデオデータとオーディオデータとを合成してパッケージングして、自己の元に返して欲しいという内容のタスク要求が、パケットにより発信された場合の流れの例が示されている。
- このタスク要求がPU(#B)5003に渡されると("1")、PU(#B)5003は、リソースアロケータ5007にそのタスクの実行に最適なSPUクラスタの割当を依頼する("2")。リソースアロケータ5007は、リソースリスト703のリソース情報に照らせば、ビデオデータについてはSPU4003、オーディオデータについてはSPU4005、パケット化プログラムはSPU4007が適していると判定し、その旨をPU(#B)5003に通知する("3")。
- 25 PU(#B)5003は、そのSPUを管理するPUのネットワークアドレス宛に、最初のタスクの依頼内容とタスクの実行結果の伝達すべき指定先であるSPU4007の通信設定情報を含むパケットを各SPU4003、4005に送信するための準備を行う。

すなわち、ビデオ用プログラムはPU(#A)のSPU4003、オーディオ

用プログラムはPU（#C）5005のSPU4005、パケット化プログラムはPU#AのSPU4007に配置する。その後、データパスを、各プログラム内に設定する。

- データパスは、タスク中に記述されたデータフローに基づき、各SPUに割り  
5 当てられた処理プログラムをキーとして、関数テーブル707の検索により取得する。例えば、ビデオ用プログラムであれば、以下のようになる。

「video.VIDEO\_OUT → packet.VIDEO\_IN」

SPUの割当情報に基づきこの結果は、以下のように変換される。

「PU#A:SPU#●:VIDEO\_OUT → PU#B:SPU#○:VIDEO\_IN」

- 10 その後、指定先の識別子を付加する。指定先の識別子の形態としては、例えば、TCP、UDPなどの既存のプロトコルを利用するのであれば、PUとIPアドレスを対応付け、さらに、SPUと通信ポート番号とを対応付ける。あるいは、各SPU又はPUに個別のIPアドレスを割り当てる。あるいは、IPアドレス  
15 によって指定先のPE、BEあるいはコンピュータ装置を指定し、これをローカル・ストレージに記録することで、既存のプロトコルを利用する。なお、これに限らず、IP以外のアドレス指定方法、例えばSPU毎のプロセッサ識別子でデータ送信が可能なネットワークのアドレスを指定するようにしても良い。以上の処理によって、例えばビデオデータの指定先を含むパケットのデータ部分は、以下のようなものとなる。

- 20 「192.168:0.1:10001:VIDEO\_OUT → 192.168.0.2:10002:VIDEO\_IN」

このようなデータ部と、それがSPU4003、4005に到達するようなヘッダと利用するプロトコルに応じたパケットを生成し、これを上述した（図3）ネットワークインタフェースを通じてネットワーク104に送信する（"4"）。

- SPU4003を管理するPU（#A）5001及びSPU4005を管理す  
25 るPU（#C）5003は、それぞれ、ネットワークインタフェースからの割込を受けて、パケットを受け取り、SPU4003、4005のローカル・メモリに、そのパケットの内容を書き込む。SPU4003、4005は、これにより、自己に依頼されたタスクを実行する。SPU4003、4005によるタスクの実行結果はパケット用プログラムと共にパケット化され、指定先であるSPU4

## 21

007に伝達されるようにする(”5”)。SPU4007は、パッケージ用プログラムで、ビデオデータとオーディオデータとをパッケージ化し、その結果を各SPU4003、4005が指定した指定先であるSPU4001に向けて送信する(”6”)。これにより、SPU4001は、タスク要求に対応した実行結果  
5 を得ることができる。

なお、それぞれのPUは、パケットのヘッダ情報だけを解析することでパケットを受け取るべきSPUを特定し、PUは各SPUに対して必要なパケットデータをネットワークインタフェース或いはDRAM上から取得するように依頼する実施の形態を採用することもできる。

10 また、上記の例では、パケット化プログラムが、SPU4003あるいはSPU4005の実行結果と共にパケット化され、そのパケットがSPU4007に送信されることでプログラムが実行される。そして、SPU4003、4005のうち先に処理が終わった方のプログラムをSPU4007で実行し、一方(オーディオまたはビデオ)からの実行結果を待つという形態例であるが、このよう  
15 な形態以外にも以下の形態が可能である。

すなわち、ビデオプログラム、オーディオプログラムと同様に”4”のメッセージによってパケット化プログラムをSPU4007で実行し、SPU4003、SPU4005からの実行結果パケットの待ち受け状態にする。SPU4003、SPU4005は、実行結果をパケット化し、これをSPU4007に送信する。  
20 このように、予めすべてのプログラムをSPU4003、4005、4007に配置して、データのフローをつなぐタイプと、処理の進展に応じて先行するSPUで実行結果と後続のプログラムを共にパケット化するタイプの2種類の実施が可能である。

図14は、図13の例において、SPU4005が、SPU4003からのパ  
25 ケットによって動作可能になる場合の例を示している(”1”)~(”5”)。また、SPU4005からSPU4007に送られた実行結果とSPUプログラムは、SPU4007に送られる(”6”)。そして、SPU4007でパッケージ化した後、タスク要求を行ったSPU4001ではなく、そのSPU4001と同じグループにクラスタ化された他のSPU4009に送信する(”7”)。

これらの例から判るように、PUは、リソースアロケータ5007に依頼して、最適なSPUを特定し、特定したSPUにタスクの依頼内容と指定先とを送信することにより、自己を介在させることなく、SPUだけで、以後のタスクの依頼、及びその実行結果の受け渡しが行われるので、タスク要求に際してのPUのオーバーヘッドが回避される。

また、予め想定されるタスクサイズ毎にクラスタ化しておき、タスク要求時に最適なSPUクラスタを割り当てるようにしたので、SPU割当に要する時間も短縮化される。このような効果は、リアルタイム性が要求される用途では、極めて意義の大きなものとなる。

10 なお、図13および図14の例において、PUの機能を当該PUの外部、例えば上述したネットワークインタフェースの一機能として内蔵させることで、PUを介さずに直接SPUがネットワークインタフェースからの割り込みを受けて処理を行う実施形態をとることも可能である。

15 また、PU相当のインタフェースネットワークあるいは処理装置をネットワーク104上に配備し、あるSPUが、目的のSPUにダイレクトにアクセスするための通信設定情報等を、PUを介さずに、ネットワークインタフェースあるいは上記の処理装置に、能動的に取りに行くという構成も可能である。この場合、ネットワークインタフェースないし処理装置には、PUとほぼ同様のアーキテクチャ（プロセッサ等）と所定のプログラムとの協働によって、以下のような機能

20 の全部又は一部の機能を形成するようにする。

(a) リソースアロケータ603、5007、クラスタマネージャ605又はこれらと同様の機能。SPUステータステーブル701、リソースリスト703、クラスタリスト705については、自ら備えていても良いし、これらを管理するサーバないしPUにアクセスするようにしても良い。

25 (b) リソースアロケータ603又は同様の機能実現体に、ネットワーク通信により、タスクの実行に最適なSPUクラスタの割当を依頼し、あるいは自ら割り当てる機能。

(c) クラスタマネージャ605又は同様の機能実現体に、ネットワーク通信により、いずれかのSPUからタスク要求を受け付けたときに当該タスクのサイズ

## 23

又は種類に適合するクラスタをクラスタリストを参照することにより特定し、特定したクラスタに対して当該タスクの割当を依頼し、あるいは自ら割り当てる機能。

(d) SPU又はSPUを管理するPUのネットワークアドレスを管理（自ら記録／更新するか、あるいは、記録されているサーバ等へのアクセスを行う）機能。

(e) 最初のタスクの依頼内容とタスクの実行結果の伝達すべき指定先であるSPUの通信設定情報を含むパケットを、目的のSPUへ送信するための準備を行う機能、あるいは、実際にパケットをSPUへ送信する機能。

本発明によれば、処理管理装置が、タスク要求を受け取ったときに、そのタスク要求に対するタスクを遂行可能な処理装置を特定し、特定した処理装置とタスク要求元との間のダイレクトなタスクの実行結果の受け渡しを可能にするので、処理管理装置におけるオーバーヘッドを回避することができる。これにより、複数の処理装置による分散処理を効率的に行うことができる。

## 請求の範囲

1. それぞれ依頼されたタスクを実行し、その実行結果を依頼元の指定先へ送信する機能を備えた複数の処理装置が接続可能なネットワークに接続される処理管理装置であって；  
5 前記ネットワークに接続されている処理装置のタスク実行能力を表すリソース情報と当該処理装置とのネットワーク通信を可能にするための通信設定情報とがリストアップされている所定のメモリへのアクセスを可能にする第1管理手段と；
- 10 ネットワーク通信によりいずれかの処理装置からタスクの実行要求を受け付けたときに当該タスクを実行可能な処理装置を前記メモリにリストアップされているリソース情報より特定し、特定した処理装置についての前記通信設定情報を前記メモリより取得するとともに、前記特定した処理装置と前記タスク要求を行った処理装置の少なくとも一方の処理装置に他方の処理装置の前記通信設定情報を  
15 伝達することにより、当該処理装置間のネットワーク通信によるダイレクトな前記実行結果の受け渡しを可能にする第2管理手段と；  
を有する処理管理装置。
2. 前記メモリと、このメモリへの前記リソース情報および通信設定情報の記録を行うリソース情報等保持手段を有する、  
20 請求の範囲第1項記載の処理管理装置。
3. 前記タスクは、当該タスクの実行に続く後続タスクの実行依頼とその後続タスクの実行結果の指定先の情報を含む、  
請求の範囲第1項記載の処理管理装置。
4. 前記第2管理手段は、前記特定した処理装置に対して前記受け付けたタスクの実行結果の指定先を含む当該タスクの実行を依頼する、  
25 請求の範囲第3項記載の処理管理装置。
5. 前記第2管理手段は、前記タスク要求を行った処理装置に、前記特定した処理装置へのダイレクトなタスクの実行を依頼させる、  
請求の範囲第3項記載の処理管理装置。

6. 複数の処理装置が随時参加することができ且つ脱退できる共有空間を前記ネットワーク上に形成する共有空間モジュールを有し、

前記第2管理手段は、前記共有空間モジュールにより形成される共有空間に参加している処理装置の通信設定情報及び現在のリソース情報を当該処理装置より  
5 取得し、取得したこれらの情報を前記メモリにリストアップして使用可能状態にするとともに、前記共有空間から脱退した処理装置について前記メモリにリストアップされている情報を使用不能状態に変更する、

請求の範囲第1項記載の処理管理装置。

7. それぞれ依頼されたタスクを実行し、その実行結果を依頼元の指定先へ送信する機能を備えた処理装置と、この処理装置と内部バスを介して接続された処理管理装置とを含み；  
10

前記処理装置及び前記処理管理装置は、それぞれ前記内部バスを介してネットワークに接続されており；

前記処理管理装置は；

15 前記ネットワークに接続されている処理装置のタスク実行能力を表すリソース情報と当該処理装置とのネットワーク通信を可能にするための通信設定情報とがリストアップされている所定のメモリへのアクセスを可能にする第1管理手段と；

ネットワーク通信によりいずれかの処理装置からタスクの実行要求を受け付け  
20 たときに当該タスクを実行可能な処理装置を前記メモリにリストアップされているリソース情報より特定し、特定した処理装置についての前記通信設定情報を前記メモリより取得するとともに、前記特定した処理装置と前記タスク要求を行った処理装置の少なくとも一方の処理装置に他方の処理装置の前記通信設定情報を伝達することにより、当該処理装置間のネットワーク通信によるダイレクトな前  
25 記実行結果の受け渡しを可能にする第2管理手段と；

を有するものである、コンピュータ・システム。

8. それぞれ依頼されたタスクを実行し、その実行結果を依頼元の指定先へ送信する機能を備えた複数の処理装置が接続可能なネットワークに接続されるコンピュータ・システムに搭載されるデバイスであって；

## 26

所定のコンピュータプログラムを実行することにより、前記コンピュータ・システムを；前記ネットワークに接続されている処理装置のタスク実行能力を表すリソース情報と当該処理装置とのネットワーク通信を可能にするための通信設定情報とがリストアップされている所定のメモリへのアクセスを可能にする第1管理手段；

ネットワーク通信によりいずれかの処理装置からタスクの実行要求を受け付けたときに当該タスクを実行可能な処理装置を前記メモリにリストアップされているリソース情報より特定し、特定した処理装置についての前記通信設定情報を前記メモリより取得するとともに、前記特定した処理装置と前記タスク要求を行った処理装置の少なくとも一方の処理装置に他方の処理装置の前記通信設定情報を伝達することにより、当該処理装置間のネットワーク通信によるダイレクトな前記実行結果の受け渡しを可能にする第2管理手段；として動作させる、

分散処理管理用のデバイス。

9. それぞれ依頼されたタスクを実行し、実行結果を依頼元の指定先へ送信する機能を備えた複数の処理装置と、各処理装置との間でネットワークを通じて通信を行う処理管理装置との協働で行う分散処理方法であって、

前記処理管理装置が、前記ネットワークに接続されている処理装置のタスク実行能力を表すリソース情報および当該処理装置へのアクセスを可能にするための通信設定情報を当該処理装置より取得し、取得したこれらの情報を所定のメモリにリストアップする段階と、

いずれかの前記処理装置が前記処理管理装置に対してタスク要求を送信する段階と、

前記タスク要求を受信した前記処理管理装置が、受信したタスク要求を遂行可能な1又は複数の前記処理装置を前記メモリの記録情報より特定し、特定した処理装置に、その実行結果の指定先を含むタスクの実行を依頼する段階とを有し、

前記処理管理装置を介在させない前記実行結果の受け渡しを可能にする、分散処理方法。

10. 前記複数の処理装置のすべてが共通の構造のプロセッサ・エレメントを有する、

請求の範囲第9項記載の分散処理方法。

11. 前記複数の処理装置が有するプロセッサエレメントは、命令セット・アーキテクチャ(I S A)がすべて同じで、同じ命令セットに従って所要の処理を実行できるものである、

5 請求の範囲第10項記載の分散処理方法。

12. コンピュータを、それぞれ依頼されたタスクを実行し、その実行結果を依頼元の指定先へ送信する機能を備えた複数の処理装置が接続可能なネットワークに接続される処理管理装置として動作させるためのプログラムであって；

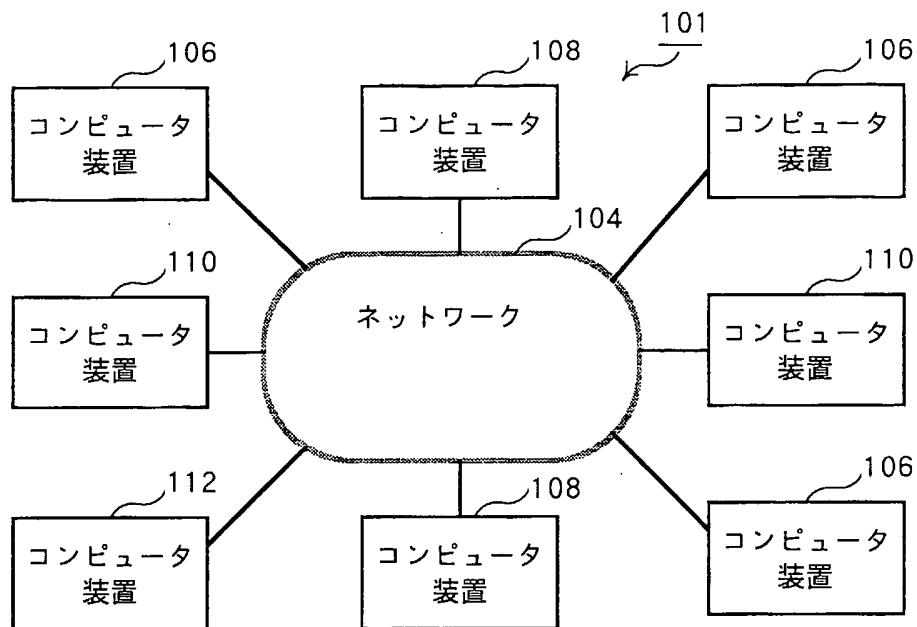
前記コンピュータを；

10 前記ネットワークに接続されている処理装置のタスク実行能力を表すリソース情報と当該処理装置とのネットワーク通信を可能にするための通信設定情報とがリストアップされている所定のメモリへのアクセスを可能にする第1管理手段；

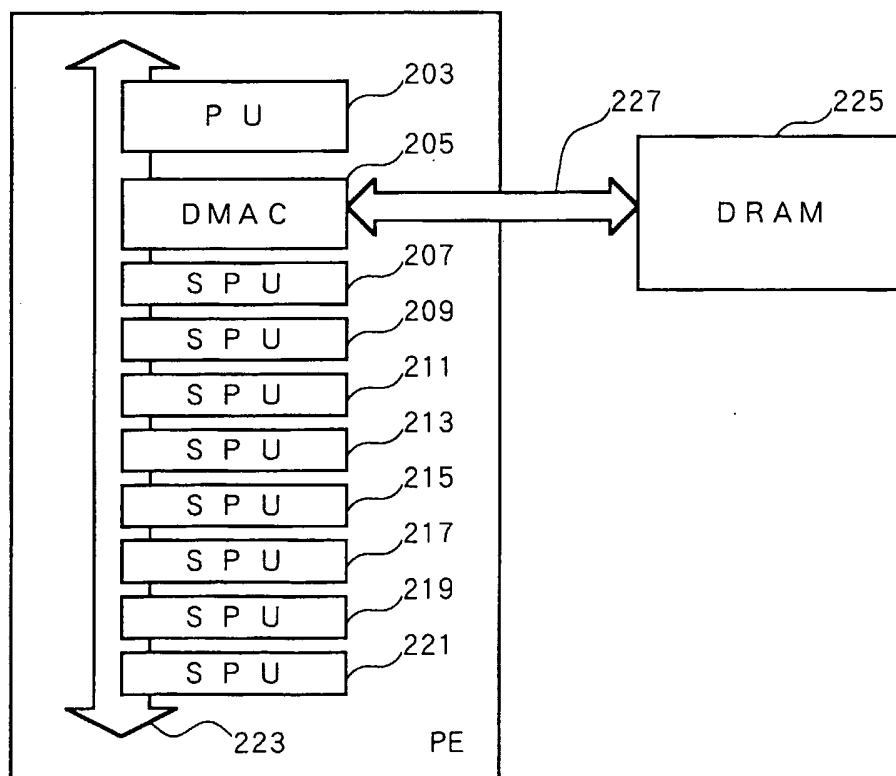
ネットワーク通信によりいずれかの処理装置からタスクの実行要求を受け付けたときに当該タスクを実行可能な処理装置を前記メモリにリストアップされているリソース情報より特定し、特定した処理装置についての前記通信設定情報を前記メモリより取得するとともに、前記特定した処理装置と前記タスク要求を行った処理装置の少なくとも一方の処理装置に他方の処理装置の前記通信設定情報を

15 伝達することにより、当該処理装置間のネットワーク通信によるダイレクトな前記実行結果の受け渡しを可能にする第2管理手段；

20 として機能させるためのコンピュータプログラム。

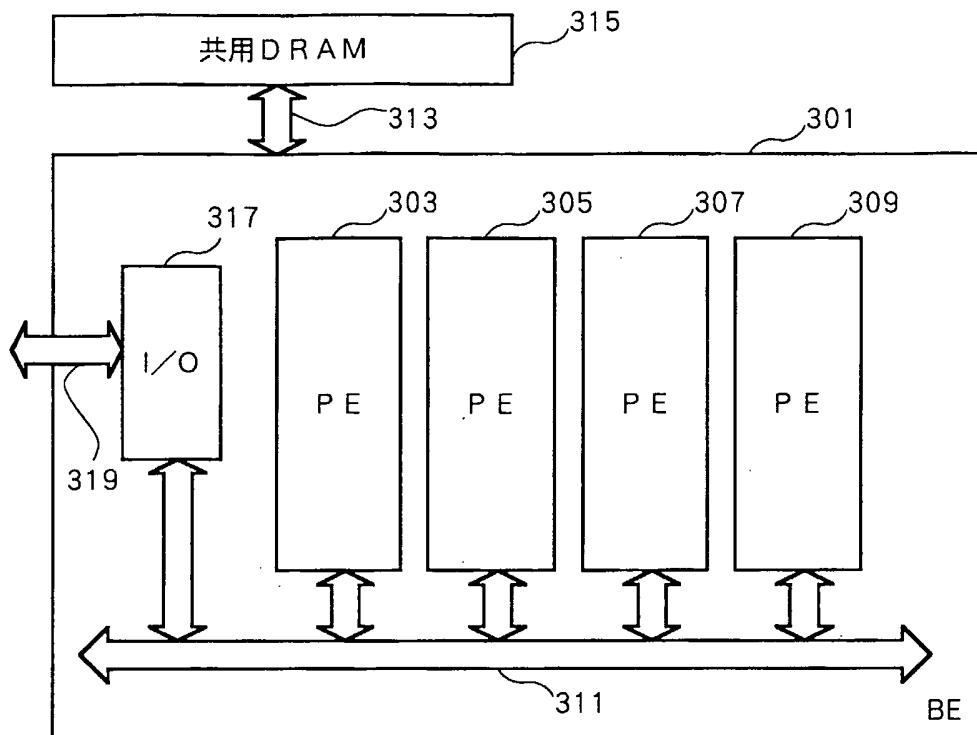


第1図

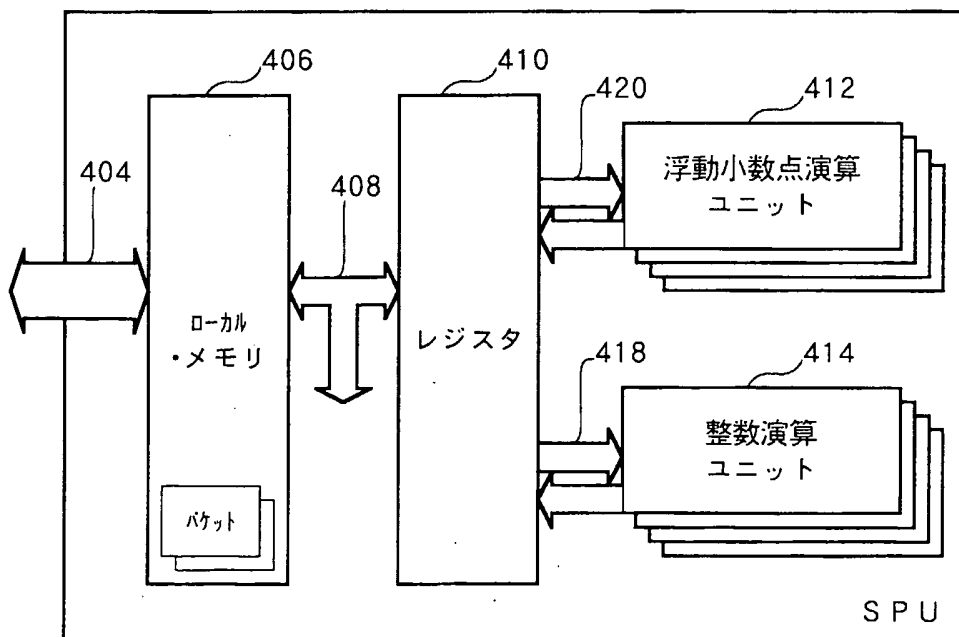


第2図

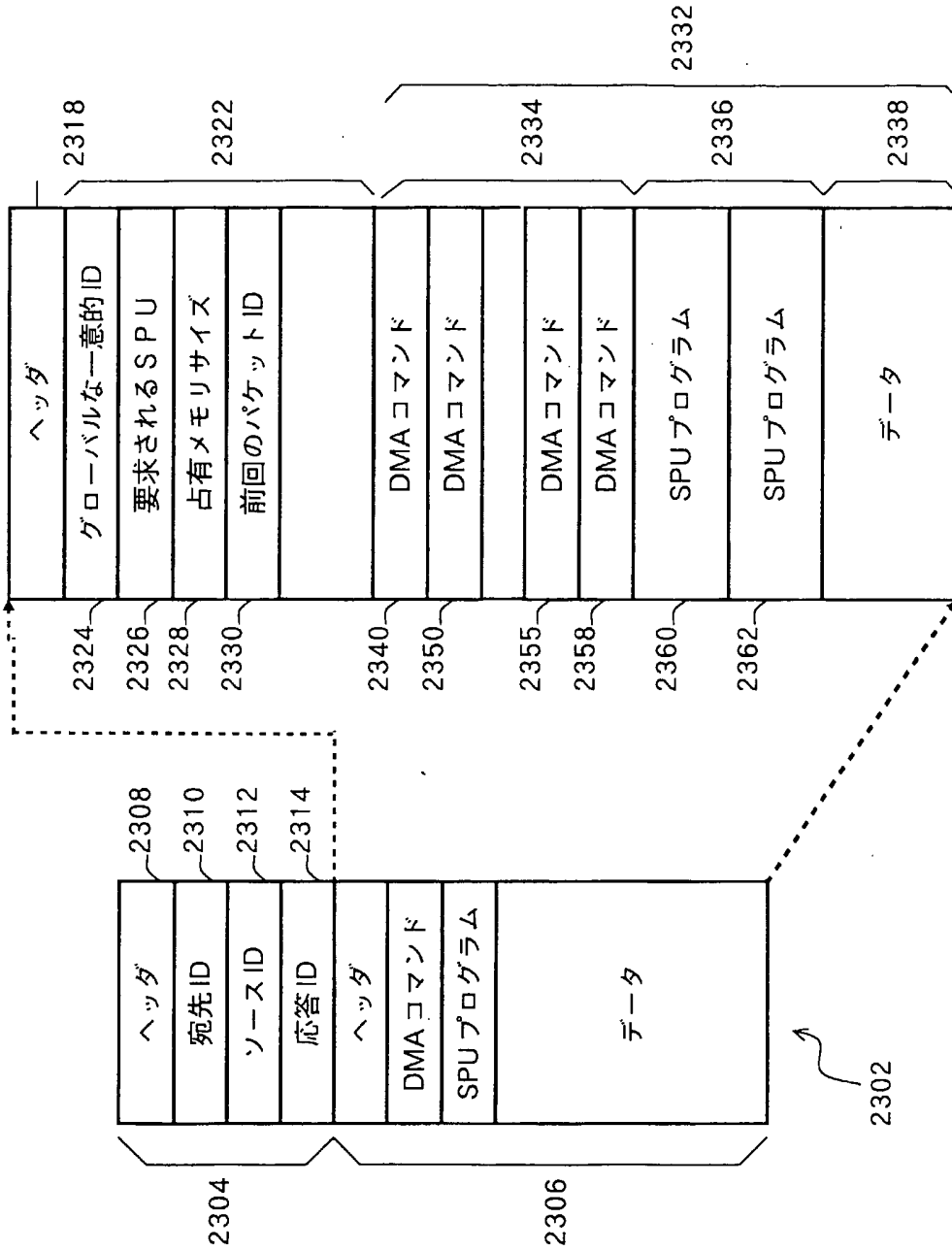
2/12



第3図



第4図



第5図

4/12

(a)

	2342	2344	2346	2348
2340	VID コマンド	ロード コマンド	アドレス	ローカル・メモリ のアドレス
2350	VID コマンド	ロード コマンド	アドレス	ローカル・メモリ のアドレス

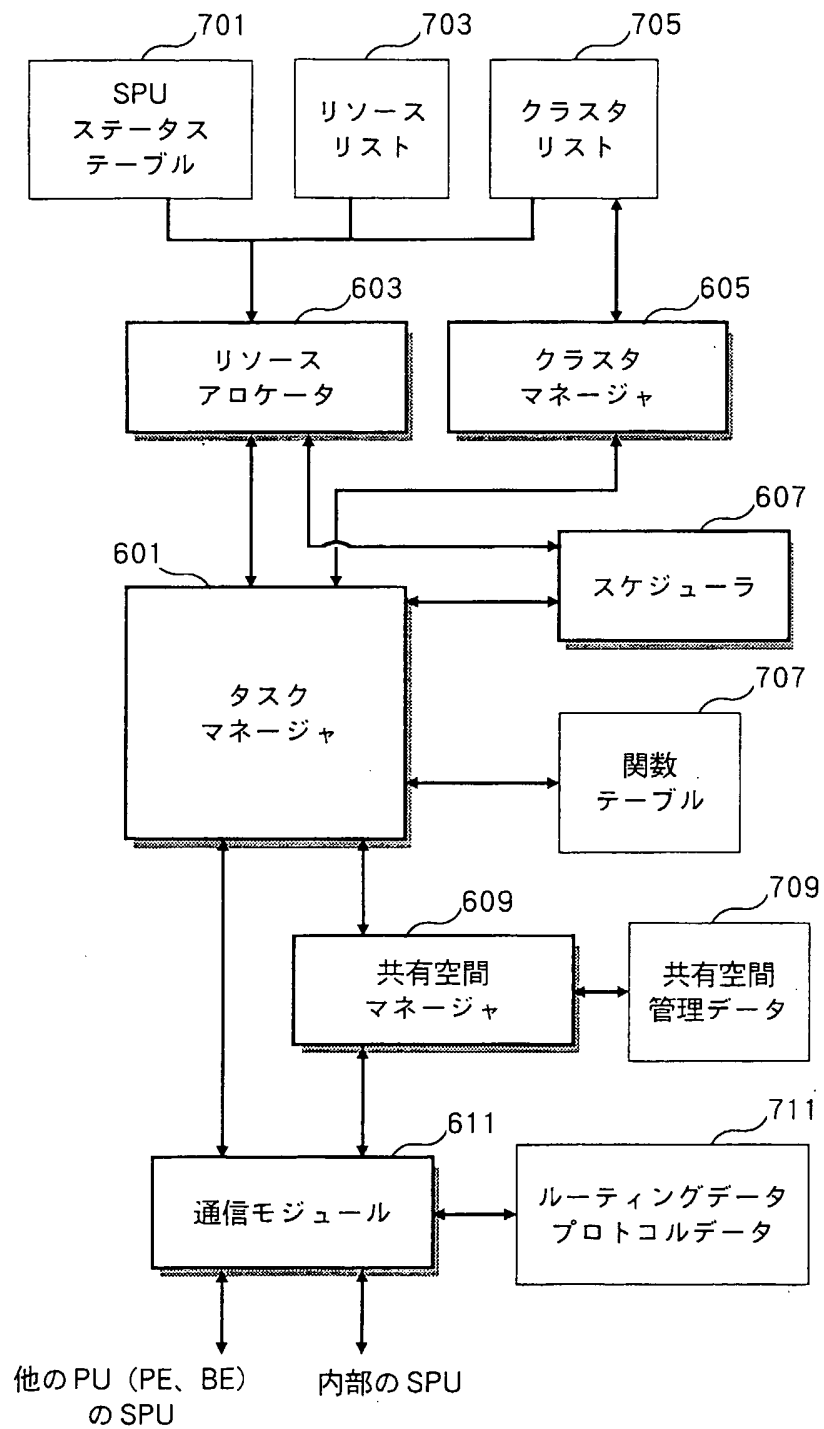
(b)

	2352	2354	2356
2355	VID コマンド	キック コマンド	プログラムカウンタ
2358	VID コマンド	キック コマンド	プログラムカウンタ

第6図

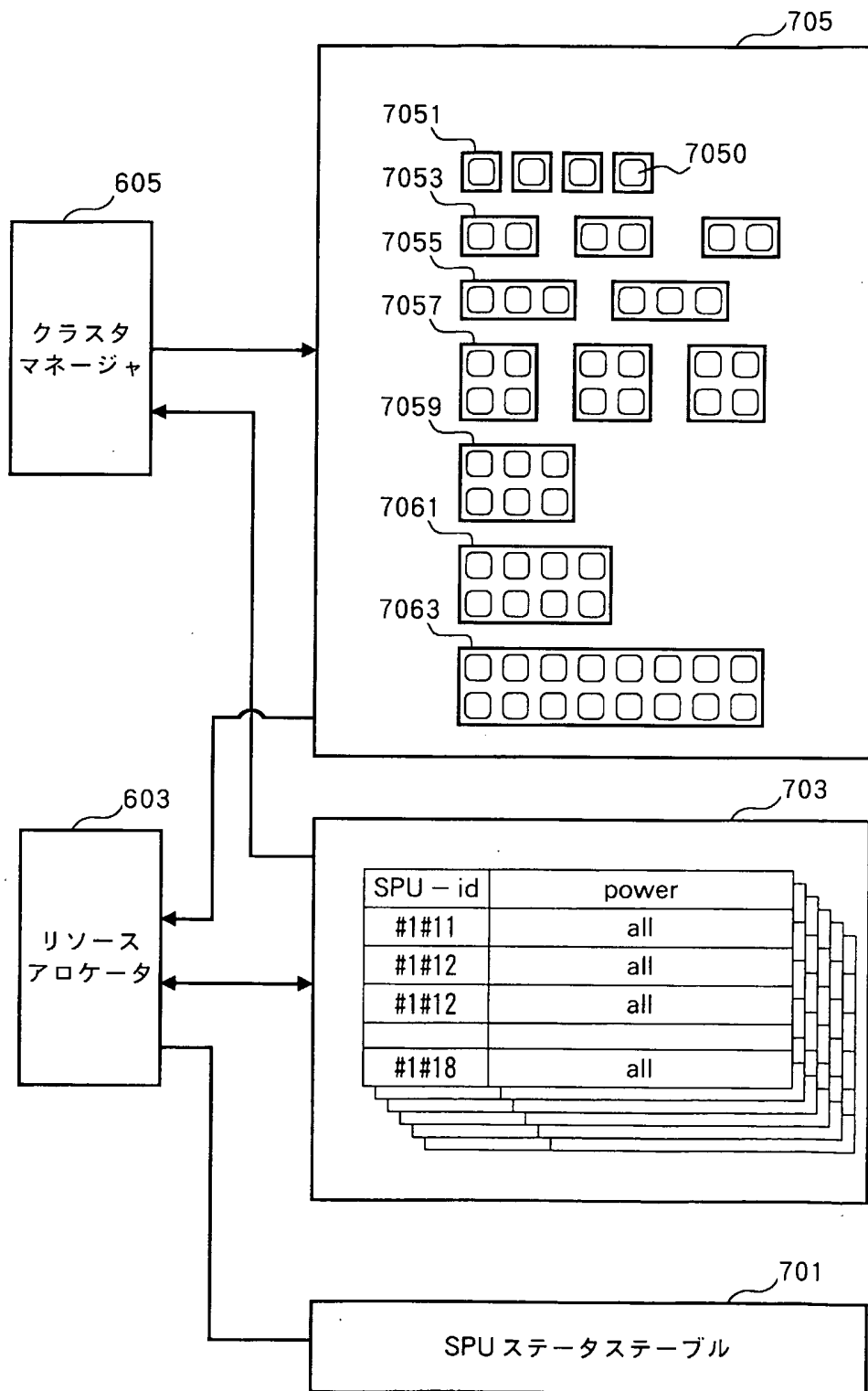


6/12

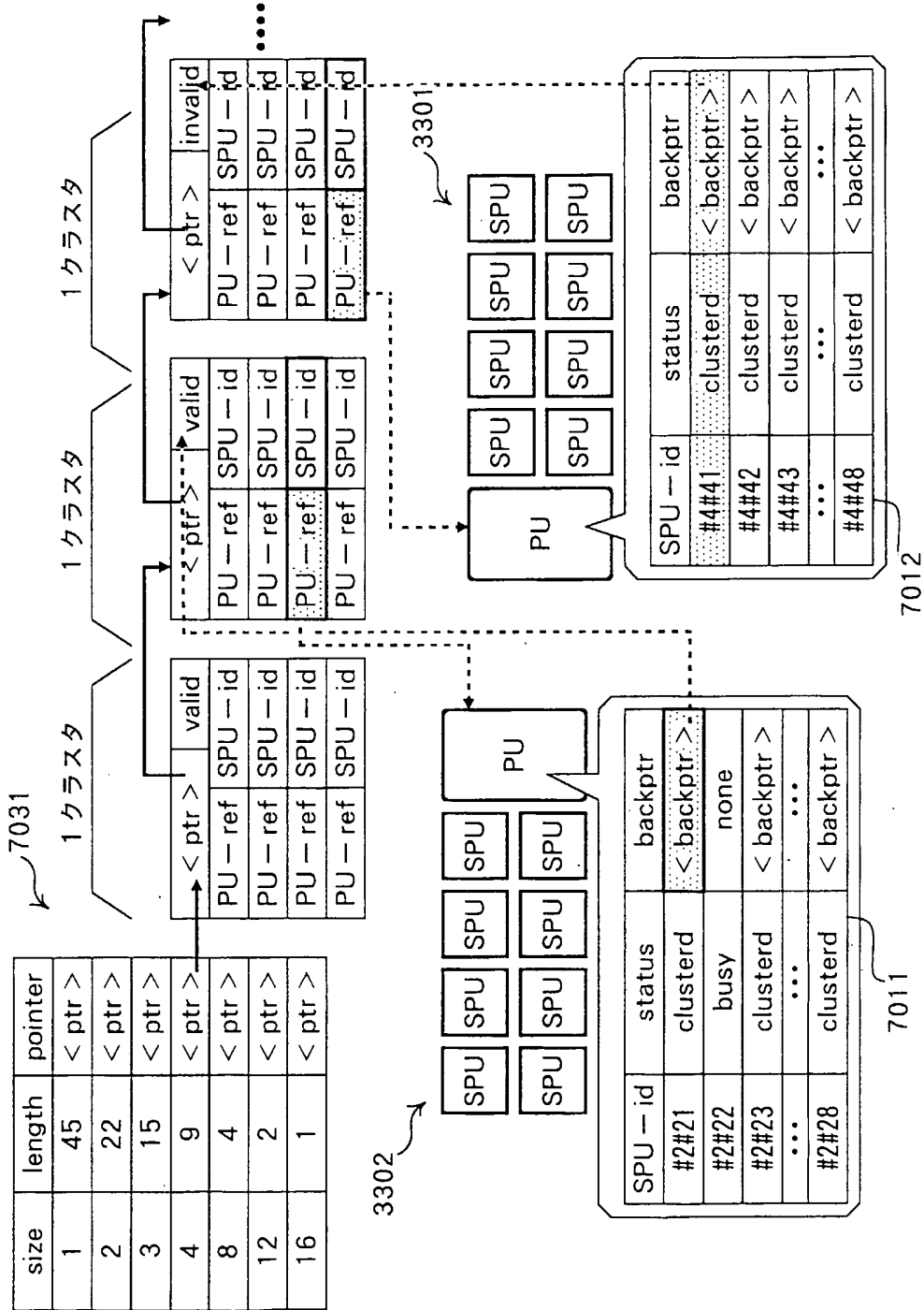


第8図

7/12

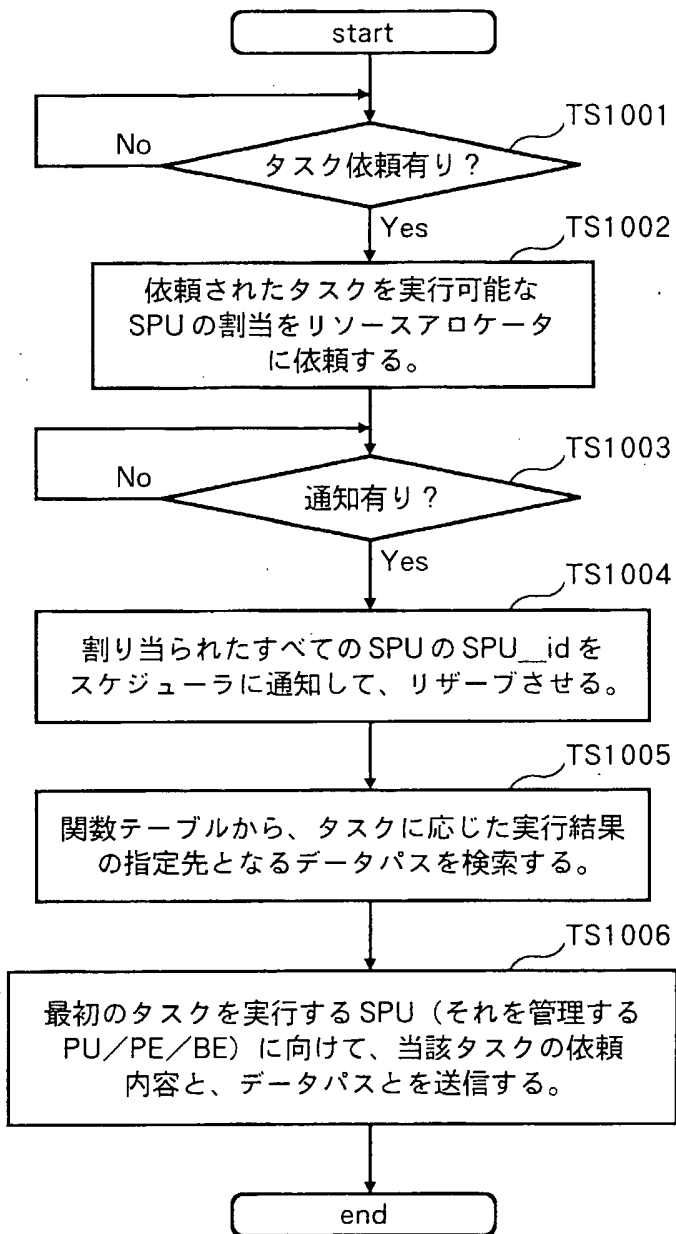


第9図



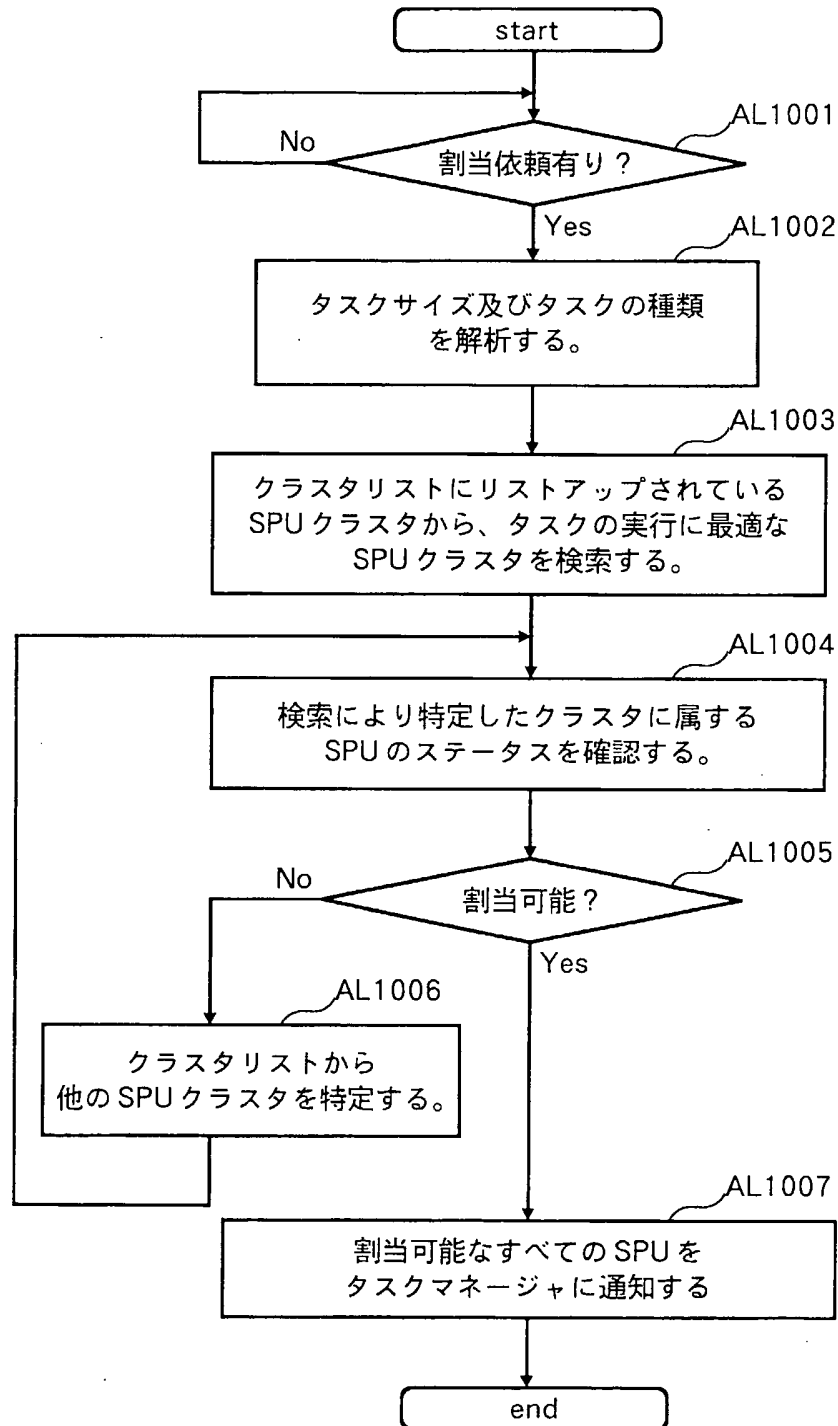
第10図

9/12

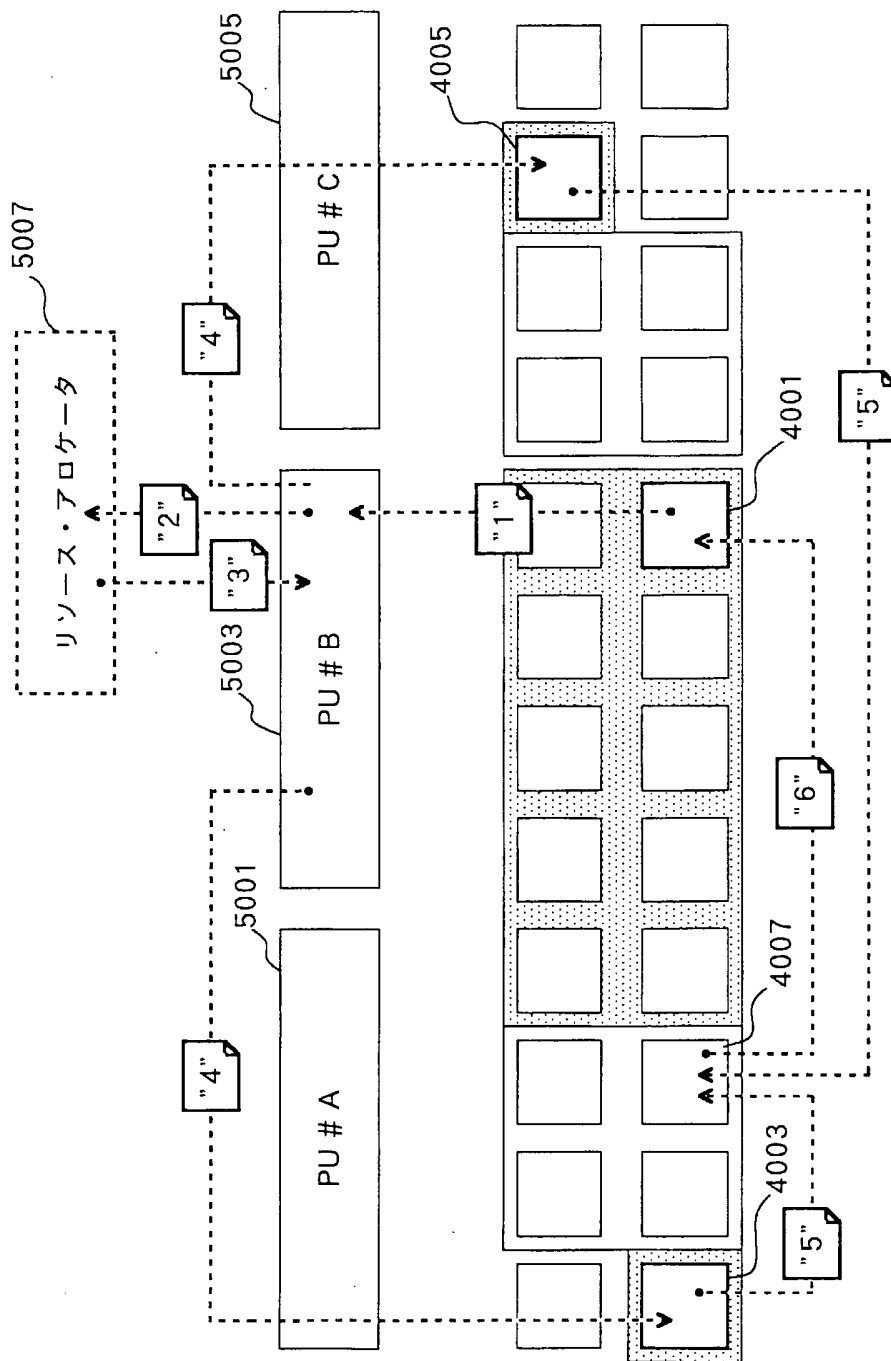


第11図

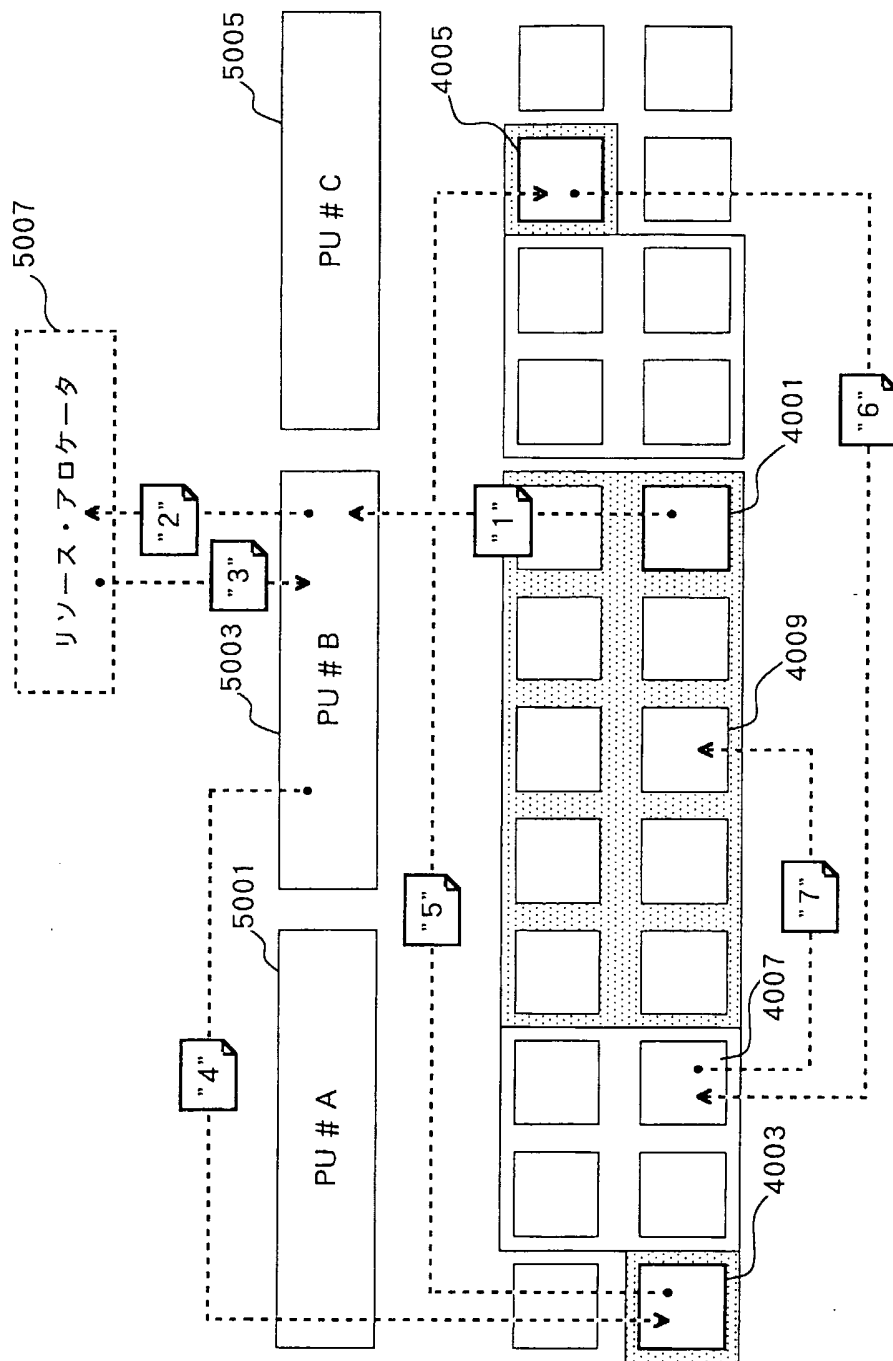
10/12



第12図



第13図



第14図

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2005/010986

A. CLASSIFICATION OF SUBJECT MATTER Int. Cl. <sup>7</sup> G06F9/50		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) Int. Cl. <sup>7</sup> G06F9/50		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Jitsuyo Shinan Koho 1922-1996 Jitsuyo Shinan Toroku Koho 1996-2005 Kokai Jitsuyo Shinan Koho 1971-2005 Toroku Jitsuyo Shinan Koho 1994-2005		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
P, X	JP 2004-287801 A (Sony Computer Entertainment Inc.), 14 October, 2004 (14.10.04), Par. Nos. [0009], [0031] to [0033], [0047] to [0048], [0060]	1-5, 7-12
Y	JP 7-056754 A (International Business Machines Corp.), 03 March, 1995 (03.03.95), Par. Nos. [0053] to [0061]	1-12
Y	JP 7-141302 A (Director General, Agency of Industrial Science and Technology), 02 June, 1995 (02.06.95), Par. Nos. [0017] to [0019]	1-12
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 12 September, 2005 (12.09.05)		Date of mailing of the international search report 27 September, 2005 (27.09.05)
Name and mailing address of the ISA/ Japanese Patent Office		Authorized officer
Facsimile No.		Telephone No.

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2005/010986

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	JP 6-348662 A (Fuji Xerox Co., Ltd.), 22 December, 1994 (22.12.94), Par. Nos. [0013] to [0014], [0020]	1-12
Y	JP 8-083253 A (Toshiba Corp.), 26 March, 1996 (26.03.96), Par. Nos. [0021], [0027]	6

**INTERNATIONAL SEARCH REPORT**  
Information on patent family members

International application No.  
PCT/JP2005/010986

JP 2004-287801 A	2004.10.14	US 2004/205759 A1	2004.10.14
		WO 2004/084069 A2	2004.09.30
JP 7-056754 A	1995.03.03	US 5574911 A1	1996.11.12
JP 7-141302 A	1995.06.02	(Family: none)	
JP 6-348662 A	1994.12.22	(Family: none)	
JP 8-083253 A	1996.03.26	US 5829041 A1	1998.10.27

A. 発明の属する分野の分類 (国際特許分類 (IPC))  
 Int.Cl.<sup>7</sup> G06F9/50

B. 調査を行った分野  
 調査を行った最小限資料 (国際特許分類 (IPC))  
 Int.Cl.<sup>7</sup> G06F9/50

最小限資料以外の資料で調査を行った分野に含まれるもの

日本国実用新案公報	1922-1996年
日本国公開実用新案公報	1971-2005年
日本国実用新案登録公報	1996-2005年
日本国登録実用新案公報	1994-2005年

国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)

C. 関連すると認められる文献

引用文献の カテゴリー*	引用文献名、及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
P, X	JP 2004-287801 A (株式会社ソニー・コンピュータエンタテインメント) 2004.10.14, 段落【0009】、【0031】-【0033】、 【0047】-【0048】、【0060】	1-5, 7-12
Y	JP 7-056754 A (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) 1995.03.03, 段落【0053】-【0061】	1-12

C欄の続きにも文献が列挙されている。  パテントファミリーに関する別紙を参照。

* 引用文献のカテゴリー	の日の後に公表された文献
「A」特に関連のある文献ではなく、一般的技術水準を示すもの	「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの
「E」国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの	「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの
「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)	「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの
「O」口頭による開示、使用、展示等に言及する文献	「&」同一パテントファミリー文献
「P」国際出願日前で、かつ優先権の主張の基礎となる出願	

国際調査を完了した日  
 12.09.2005

国際調査報告の発送日  
 27.9.2005

国際調査機関の名称及びあて先  
 日本国特許庁 (ISA/JP)  
 郵便番号100-8915  
 東京都千代田区霞が関三丁目4番3号

特許庁審査官 (権限のある職員)  
 殿川 雅也  
 5B 9646  
 電話番号 03-3581-1101 内線 3544

C (続き) . 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
Y	JP 7-141302 A (工業技術院長) 1995. 06. 02, 段落【0017】 - 【0019】	1 - 12
Y	JP 6-348662 A (富士ゼロックス株式会社) 1994. 12. 22, 段落【0013】 - 【0014】、【0020】	1 - 12
Y	JP 8-083253 A (株式会社東芝) 1996. 03. 26, 段落【0021】、 【0027】	6

国際調査報告  
パテントファミリーに関する情報

国際出願番号 PCT/JP2005/010986

JP 2004-287801 A	2004. 10. 14	US 2004/205759 A1 WO 2004/084069 A2	2004. 10. 14 2004. 09. 30
JP 7-056754 A	1995. 03. 03	US 5574911 A1	1996. 11. 12
JP 7-141302 A	1995. 06. 02	ファミリーなし	
JP 6-348662 A	1994. 12. 22	ファミリーなし	
JP 8-083253 A	1996. 03. 26	US 5829041 A1	1998. 10. 27