

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2004-220581  
(P2004-220581A)

(43) 公開日 平成16年8月5日(2004.8.5)

(51) Int. Cl. <sup>7</sup>	F I	テーマコード (参考)
G06F 13/14	G06F 13/14 320H	5B005
G06F 12/02	G06F 12/02 570M	5B014
G06F 12/08	G06F 12/08 531Z	5B060
G06F 12/10	G06F 12/08 551C	
	G06F 12/08 553Z	
	審査請求 有 請求項の数 23 O L (全 29 頁) 最終頁に続く	

(21) 出願番号 特願2003-421698 (P2003-421698)  
 (22) 出願日 平成15年12月18日 (2003.12.18)  
 (31) 優先権主張番号 10/339724  
 (32) 優先日 平成15年1月9日 (2003.1.9)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 390009531  
 インターナショナル・ビジネス・マシー  
 ズ・コーポレーション  
 INTERNATIONAL BUSIN  
 ESS MASCHINES CORPO  
 RATION  
 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  
 (74) 代理人 100086243  
 弁理士 坂口 博  
 (74) 代理人 100091568  
 弁理士 市位 嘉宏  
 (74) 代理人 100108501  
 弁理士 上野 剛史

最終頁に続く

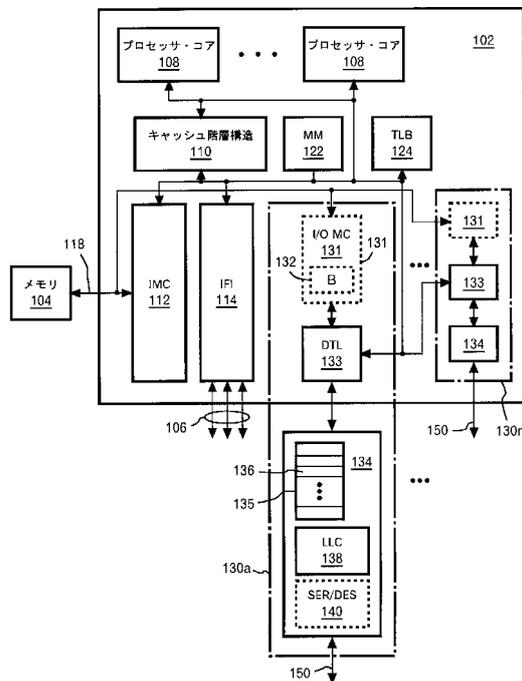
(54) 【発明の名称】 入出力 (I/O) 通信のハードウェア・アクセラレーションを実現するデータ処理システム

(57) 【要約】

【課題】 改善された処理装置、データ処理システム、およびデータ処理方法を提供すること。

【解決手段】 処理装置などの集積回路は、基板、および基板中に形成された集積回路要素を含む。集積回路要素は、命令を実行するプロセッサ・コアと、プロセッサ・コアに結合され、プロセッサ・コアと集積回路の外部のシステム相互接続の間での通信をサポートする相互接続インターフェイスと、プロセッサ・コアに結合され、入出力通信リンクを介して入出力通信をサポートする外部通信アダプタの少なくとも一部分を含む。

【選択図】 図4



**【特許請求の範囲】****【請求項 1】**

基板と、前記基板中に形成された集積回路要素とを備える集積回路であって、前記集積回路要素が、

命令を実行するプロセッサ・コアと、

前記プロセッサ・コアに結合され、前記プロセッサ・コアと前記集積回路の外部のシステム相互接続との間の通信をサポートする相互接続インターフェースと、

前記プロセッサ・コアに結合され、入出力通信リンクを介して入出力通信をサポートする外部通信アダプタの少なくとも一部分と

を含む集積回路。

10

**【請求項 2】**

前記集積回路要素がさらに、

前記プロセッサ・コアに結合されたキャッシュ階層構造と、

前記キャッシュ階層構造に結合された集積化メモリ制御装置と

を備える、請求項 1 に記載の集積回路。

**【請求項 3】**

前記外部通信アダプタを含めて複数の外部通信アダプタの少なくとも一部分を含み、前記複数の外部通信アダプタのうちの少なくとも 2 つがそれぞれ異なる入出力通信プロトコルを実施する、請求項 1 に記載の集積回路。

**【請求項 4】**

前記集積回路要素が、通信要求で指定されるアドレスに基づいて、前記集積回路要素内の通信コマンドを前記相互接続インターフェースおよび前記外部通信アダプタに経路指定するためのメモリ・マップを含む、請求項 1 に記載の集積回路。

20

**【請求項 5】**

前記外部通信アダプタに結合された変換索引バッファをさらに備え、前記変換索引バッファが、入出力 (I/O) コマンドで指定される実効アドレスを実アドレスに変換する、請求項 1 に記載の集積回路。

**【請求項 6】**

前記外部通信アダプタの前記少なくとも一部分が、I/O データにアクセスするためにメモリにアクセスする入出力 (I/O) メモリ制御装置を含む、請求項 1 に記載の集積回路。

30

**【請求項 7】**

前記外部通信アダプタの前記少なくとも一部分が、

前記プロセッサ・コアおよびリンク・レイヤ制御装置に結合され、前記プロセッサ・コアによる I/O コマンドにตอบสนองして、入出力 (I/O) データ転送を制御するデータ転送ロジックを備える、請求項 1 に記載の集積回路。

**【請求項 8】**

前記外部通信アダプタの前記少なくとも一部分が、リンク・レイヤ制御装置をさらに備える、請求項 1 に記載の集積回路。

**【請求項 9】**

前記外部通信アダプタが、I/O データ転送の完了を示すソフトウェア・アクセス可能な標識をセットする手段を含む、請求項 1 に記載の集積回路。

40

**【請求項 10】**

前記集積回路内の前記外部通信アダプタの前記部分が、着信 I/O データと発信 I/O データの少なくとも一方をバッファするバッファを含む、請求項 1 に記載の集積回路。

**【請求項 11】**

請求項 1 に記載の少なくとも 1 つの集積回路と、

前記相互接続インターフェースに結合されたシステム相互接続と、

前記少なくとも 1 つの集積回路に結合されたメモリ・システムと

を備えるデータ処理システム。

50

## 【請求項 1 2】

前記集積回路が、第 1 の集積回路を備え、  
前記外部通信アダプタの前記少なくとも一部分が、第 1 の部分を備え、  
前記データ処理システムがさらに、前記第 1 の集積回路に接続された第 2 の集積回路を含み、前記外部通信アダプタが、前記第 2 の集積回路内で実施される第 2 の部分を含む、請求項 1 1 に記載のデータ処理システム。

## 【請求項 1 3】

前記第 2 の部分が、リンク・レイヤ制御装置を備える、請求項 1 2 に記載のデータ処理システム。

## 【請求項 1 4】

基板と、前記基板中に形成された集積回路要素とを含む第 1 の集積回路であって、前記集積回路要素が

命令を実行するプロセッサ・コアと、

前記プロセッサ・コアに結合された相互接続インターフェースであって、前記プロセッサ・コアが前記少なくとも 1 つのプロセッサ・コアと前記第 1 の集積回路の外部のシステム相互接続との間での通信をサポートする相互接続インターフェースと、

前記プロセッサ・コアに結合され、入出力通信リンクを介して入出力通信をサポートする外部通信アダプタの第 1 の部分とを含む、第 1 の集積回路と、

前記第 1 の集積回路のピンに接続された第 2 の集積回路とを備え、前記外部通信アダプタが前記第 2 の集積回路内に実施される第 2 の部分を含むシステム。

## 【請求項 1 5】

前記外部通信アダプタの前記第 1 の部分が、前記プロセッサ・コアに結合され、プロセッサ・コアによる I / O コマンドにตอบสนองして入出力 ( I / O ) データ転送を制御するデータ転送ロジックを備え、

前記外部通信アダプタの前記第 2 の部分が、リンク・レイヤ制御装置を備える、請求項 1 4 に記載のシステム。

## 【請求項 1 6】

プロセッサ・コアを含む集積回路を動作させる方法であって、

前記集積回路内の相互接続インターフェースを利用して、プロセッサ・コアと前記集積回路の外部のシステム相互接続との間での通信をサポートするステップと、

前記集積回路内の外部通信アダプタを利用して、入出力 ( I / O ) 通信リンクを介して入出力通信をサポートするステップであって、集積回路内の通信だけを利用して、前記プロセッサ・コアから前記外部通信アダプタに I / O 通信コマンドを伝送することを含むステップと

を含む方法。

## 【請求項 1 7】

前記プロセッサ・コアがアクセスする可能性が高いデータおよび命令を前記集積回路内のキャッシュ階層構造にキャッシュするステップと、

前記集積回路内の集積化メモリ制御装置を利用して外部メモリにアクセスするステップと

をさらに含む、請求項 1 6 に記載の方法。

## 【請求項 1 8】

前記集積回路が、第 1 および第 2 の外部通信アダプタの少なくとも一部分を含み、前記 I / O 通信をサポートするステップが、

前記集積回路内の前記第 1 の外部通信アダプタを利用して、第 1 のリンク・レイヤ・プロトコルを使用する I / O 通信をサポートするステップと、

前記集積回路内の前記第 2 の外部通信アダプタを利用して、異なる第 2 のリンク・レイヤ・プロトコルを使用する I / O 通信をサポートするステップと

を含む、請求項 1 6 に記載の方法。

## 【請求項 1 9】

10

20

30

40

50

前記通信要求で指定されるアドレスに基づいて、前記集積回路内の通信コマンドを前記相互接続インターフェースおよび前記外部通信アダプタに経路指定するステップをさらに含む、請求項 16 に記載の方法。

【請求項 20】

I/O 通信をサポートする前記ステップが、前記 I/O コマンド内で指定される実効アドレスを、メモリ位置を識別するために利用される実アドレスに変換するステップをさらに含む、請求項 16 に記載の方法。

【請求項 21】

I/O 通信をサポートする前記ステップが、前記外部通信アダプタ内の入出力 (I/O) メモリ制御装置を利用して、I/O データにアクセスするためにメモリにアクセスするステップを含む、請求項 16 に記載の方法。 10

【請求項 22】

I/O 通信をサポートする前記ステップが、I/O データ転送の完了を示すソフトウェア・アクセス可能な標識をセットするステップを含む、請求項 16 に記載の方法。

【請求項 23】

前記 I/O 通信をサポートするステップが、着信 I/O データと発信 I/O データの少なくとも一方を前記集積回路内にバッファするステップを含む、請求項 16 に記載の方法。

【発明の詳細な説明】

【技術分野】

20

【0001】

本発明は、一般にデータ処理に関し、少なくとも一態様では、データ処理システムによる入出力 (I/O) 通信に関する。

【背景技術】

【0002】

従来データ処理システムにおいては、入出力 (I/O) 通信は、一般に、1本または複数の内部バスによってデータ処理システムの1つまたは複数の処理装置に結合されるメモリ・マップ I/O アダプタによって容易にされている。たとえば、図 1 は、イーサネット (R) 通信リンク 52 を介してリモート・コンピュータ 60 との I/O 通信をサポートする P C I (Peripheral Component Interconnect、周辺装置相互接続) I/O アダプタ 50 を含む従来技術の S M P (Symmetric Multiprocessor、対称型マルチプロセッサ) データ処理システム 8 を示している。 30

【0003】

図に示すように、従来技術の S M P データ処理システム 8 は、S M P システム・バス 11 により通信できるように結合される複数の処理装置 10 を含む。S M P システム・バスは、たとえば 8 バイト幅のアドレス・バスおよび 16 バイト幅のデータ・バスを含むことができ、500 M H z で動作することができる。各処理装置 10 は、プロセッサ・コア 14 およびキャッシュ階層構造 16 を含み、高速 (たとえば 533 M H z の) 専用メモリ・バス 20 を介して、外部システム・メモリ 12 用の関連メモリ制御装置 (M C) 18 と通信する。処理装置 10 は、一般に先端のカスタム集積回路 (I C) 技術を利用して製造され、2 G H z 以上のプロセッサ・クロック周波数で動作することができる。 40

【0004】

処理装置 10 間の通信は、完全にキャッシュ・コヒーレントである。すなわち、各処理装置 10 内のキャッシュ階層構造 16 は、従来 M E S I (Modified, Exclusive, Shared, Invalid) プロトコル、またはその変形を用いて、その処理装置 10 によってアクセスされる現行のキャッシュ・メモリの各グラニュール (granule) が、他の処理装置 10 またはシステム・メモリ 12 あるいはその両方の中の対応するメモリのグラニュールに対してどういう関係にあるかを追跡する。

【0005】

S M P システム・バス 11 には、メザニン (mezzanine) I/O バス制御装置 30 が、 50

またオプションで、1つまたは複数の追加のメザニン・バス制御装置32が結合される。メザニンI/Oバス制御装置30(および他のメザニン・バス制御装置32のそれぞれ)は、通信のために、それぞれのメザニン・バス40をSMPシステム・バス11にインターフェースする。典型的な実装では、メザニン・バス40は、SMPシステム・バス11よりも幅がずっと狭く、より低い周波数で動作する。たとえば、メザニン・バス40は、(アドレスとデータを多重化した)8バイト幅とすることができ、200MHzで動作することができる。

#### 【0006】

図に示すように、メザニン・バス40は、MCA(MicrochannelArchitecture、マイクロチャンネル・アーキテクチャ)IOCC42、PCIEクスプレス(3GIO)IOCC44、およびPCI IOCC46を含めて、いくつかのIOCC(I/Ochannel controllers、I/Oチャンネル制御装置)の接続をサポートする。各IOCC42~46は、固定した最大数の装置の接続をサポートするスロットを提供するバス47~49のそれぞれに結合される。PCI IOCC46の場合には、接続される装置には、I/O通信リンク52を介してネットワーク54およびリモート・コンピュータ60との通信をサポートするPCI I/Oアダプタ50が含まれる。

10

#### 【0007】

データ処理システム8内のI/Oデータと「ローカル」データが、異なるコヒーレンシ・ドメインに属することに留意されたい。すなわち、処理装置10のキャッシュ階層構造16は、コヒーレンシを維持するために、従来のMESIプロトコルまたはその変形を利用するが、リモート・コンピュータ60に転送するためにメザニンI/Oバス制御装置30にキャッシュされるデータ・グラニュールは、通常、共用状態として、またはデータ・グラニュールが後でデータ処理システム8内で修正される場合には無効状態として記憶される。ほとんどのシステムでは、非排他的状態、修正された状態、または同様の排他的状態が、データ処理システム8内でI/Oデータに対してサポートされる。さらに、着信I/Oデータ転送はすべて、データを修正するために処理装置10によって使用されるような、リード・ピフォア・ライト(たとえば、RWITM(read-with-intent-to modify、修正のためのリード)およびDCLAIM)動作ではなくてストア・スルー動作である。

20

#### 【0008】

前述の一般的なハードウェア実装では、SMPデータ処理システム8がI/O通信リンク52上でデータを伝送するために用いる典型的な方法は、アプリケーション・プロセス、動作環境ソフトウェア(たとえばOSおよび関連デバイス・ドライバ)、およびI/Oアダプタ(および他のハードウェア)のそれぞれが1部分を実施する、3つの部分からなる動作として説明することができる。

30

#### 【0009】

所与の任意の時点で、SMPデータ処理システム8の処理装置10は、一般に多数のアプリケーション・プロセスを同時に実行する。これらのプロセスの1つがシステム・メモリ12からリモート・コンピュータ60へとI/Oチャンネル52を介してデータを伝送する必要があるという最も簡単な場合には、このプロセスは、まずI/Oアダプタ50用のロックを取得するために他のプロセスと競合しなければならない。使用する予定の伝送プロトコルの信頼性およびその他の要因に応じて、データ・グラニュールが伝送に先立って他のプロセスによって修正されないことを保証するために、このプロセスがまた、1つまたは複数の伝送すべきデータ・グラニュールに対して1つまたは複数のロックを取得しなければならないこともある。

40

#### 【0010】

このプロセスは、I/Oアダプタ50用のロック(および伝送すべきデータ・グラニュール用の多分1つまたは複数のロック)を取得した後に、OSソケット・インターフェースを介して、そのオペレーティング・システム(OS)に1つまたは複数のコールを行う。これらのソケット・インターフェース・コールは、ソケットを初期化し、ポート・アドレスにソケットをバインドし、接続受入れ準備状態を示し、データを送信または受信しあ

50

るいはその両方を行い、ソケットをクローズすることをオペレーティング・システムに求める要求を含む。これらのソケット・コールでは、呼出し (calling) プロセスは、一般的に利用すべきプロトコル (たとえば TCP、UDP など)、アドレッシングの方法、伝送すべきデータ・グラニュールのベース EA (effective address、実効アドレス)、データ・サイズ、およびリモート・コンピュータ 60 内の宛先メモリ位置を示すフォーリン・アドレス (foreign address) を指定する。

#### 【0011】

次に動作環境ソフトウェアを参照すると、ブートに続いて、OS は、処理装置 10 によって内部で使用される仮想 (または実効) アドレス空間とは別の I/O アドレス空間を割り付けること、およびシステム・メモリ 12 中に TCE (Translation Control Entry、変換制御エントリ) テーブル 24 を作成することを含めて、I/O 通信用のリソースを作成するための様々な動作を実施する。TCE テーブル 24 は、I/O 装置によって生成される I/O アドレスとシステム・メモリ 12 中の RA の間で変換を行う TCE を提供することにより、I/O 通信を実施するのに利用される DMA (Direct Memory Access、ダイレクト・メモリ・アクセス) サービスをサポートする。

10

#### 【0012】

これらおよびその他のリソースの作成に続いて、OS は、I/O 通信をサポートするサービスを提供することにより様々なプロセスのソケット・インターフェース・コールに回答する。たとえば、OS は、まずソケット・インターフェース・コールに含まれる EA を RA (real address、実アドレス) に変換し、次いで、たとえば RA をハッシングすることにより、RA にマッピングすべき PCI I/O アドレス空間のページを決定する。さらに、OS は、システム・メモリ 12 中の TCE テーブル 24 を、要求された I/O 通信を実施するために利用される DMA サービスがサポートされるように動的に更新する。もちろん、TCE テーブル 24 内の TCE が現在使用可能でない場合、OS は、TCE テーブル 24 中のある TCE を犠牲にし、影響を受けるプロセスにその DMA が終了したことを知らせるか、あるいは必要な TCE を解放するプロセスを要求する必要がある。

20

#### 【0013】

次いで、ほとんどのデータ処理システムでは、OS は、I/O アダプタ 50 によるデータ転送のパラメータを指定する CCB (Command Control Block、コマンド制御ブロック) 22 をメモリ 12 内に作成する。たとえば、CCB 22 は、システム・メモリ 12 内で位置を指定する 1 つまたは複数の PCI アドレス空間アドレス、かかる各アドレスに関連するデータ・サイズ、およびリモート・コンピュータ 60 内の CCB のフォーリン・アドレスを含むことができる。データ転送のための TCE および CCB 22 を確立した後、OS は、CCB 22 のベース・アドレスをその呼出しプロセスに返す。使用されるプロトコルによっては、OS はまた、(たとえば、ヘッダでデータをカプセル化することや、フロー制御を提供することなどにより) 追加のデータ処理サービスを提供することもできる。

30

#### 【0014】

CCB 22 のベース・アドレスを受け取るとそれに回答して、このプロセスは、CCB 22 のベース・アドレスを PCI I/O アダプタ 50 内のレジスタに書き込むことにより、システム・メモリ 12 からリモート・コンピュータ 60 へのデータ転送を開始する。この呼出しに回答して、PCI I/O アダプタ 50 は、そのレジスタ中にその呼出しプロセスによって書き込まれるベース・アドレスを利用して、CCB 22 の DMA 読出しを実施する。(一部の簡単なシステムでは、CCB 22 が非変換アドレス領域内にあるので、アドレス変換は、CCB 22 の DMA 読出しにとって必要ではない。ただし、よりハイエンドのサーバー・クラス・システムでは、一般に CCB 22 の DMA 読出しのためにアドレス変換が実施される)。次いで、アダプタ 50 は CCB 22 を読み出し、リモート・コンピュータ 60 に伝送すべき第 1 のデータ・グラニュールの (CCB 22 から読み出された) ベース PCI アドレス空間アドレスを目標とする DMA 読出し動作を発行する。

40

#### 【0015】

PCI アダプタ 50 からの DMA 読出し動作の受取りに回答して、PCI IOCC4

50

6はその内部TCEキャッシュにアクセスして、指定された目標アドレスに対する変換を見つけ出す。TCEキャッシュ・ミスに回答して、PCI IOCC46は、TCEテーブル24の読出しを実施して、関連するTCEを得る。PCI IOCC46が必要なTCEを得た後に、PCI IOCC46は、TCEを参照してDMA読出し動作内で指定されるPCIアドレス空間アドレスをRAに変換し、システム・メモリ12のDMA読出しを実施し、要求されたI/OデータをPCI I/Oアダプタ50に返す。PCI I/Oアダプタ50による(たとえばリンク・レイヤ・プロトコルの要件を満たすための)さらなる可能な処理の後、PCI I/Oアダプタ50は、データ・グラニュールを、リモート・コンピュータ60のシステム・メモリ内のデータ・グラニュールの記憶を制御する、リモート・コンピュータ60内のCCBのフォーリン・アドレスと一緒に、I/O通信リンク52およびネットワーク54を介してリモート・コンピュータ60に伝送する。

10

20

30

40

50

**【0016】**

DMA読出し動作およびデータ伝送の前述のプロセスは、PCI I/Oアダプタ50が、CCB22内で指定されるデータをすべて伝送してしまうまで続行する。その後、PCI I/Oアダプタ50は割り込みをアサートして、データ転送が完了したことを知らせる。当業者には理解されるように、PCI I/Oアダプタ50による割り込みのアサートによって、処理装置10の1つによるコンテキスト切換えおよびFLIH(first-level interrupt handler、第1レベル割り込みハンドラ)の実行がトリガされる。次いで、FLIHは、(たとえばメザニンI/Oバス制御装置30内の)システム割り込み制御レジスタを読み取って、その割り込みがPCI IOCC46に端を發したと判断し、PCI IOCC46の割り込み制御レジスタを読み取って、その割り込みがPCI I/Oアダプタ50によって生成されたと判断し、その後、PCI I/Oアダプタ50のSLIH(second-level interrupt handler、第2レベル割り込みハンドラ)を呼び出して、PCI I/Oアダプタ50の割り込み制御レジスタを読み取って、多分複数のDMAのうちどれが完了したかを判定する。次いで、FLIHは、I/Oデータ転送が完了したことを呼出しプロセスに指示するために、ポーリング・フラグをセットする。

**【発明の開示】****【発明が解決しようとする課題】****【0017】**

本発明は、以上で概要を説明した従来のI/O通信が非効率的であることを認識したものである。先に指摘したように、OSがメモリ中にTCEテーブルを提供し、それによってIOCCがシステム・メモリ中でI/Oドメインからのアドレスを実アドレスに変換することが可能になる。システム・メモリ中でのTCEテーブルの作成と管理に関連するオーバーヘッドは、オペレーティング・システムの性能を低下させ、IOCCによるI/Oアドレスの変換は、各I/Oデータ転送の待ち時間を増加させる。複数のプロセスによるI/Oアダプタおよびシステム・メモリへのアクセスを同期させるためにロックを使用すること、ならびにそれらへのアクセスについて調停すること、ならびにI/O(たとえば、PCI)バス、メザニン・バス、およびSMPシステム・バスによって実施されるプロトコル間で変換を行うことによってさらなる待ち時間が生じる。さらに、SMPシステム・バス上でのI/Oデータ転送の伝送は、そうでなければ処理装置間の(たとえば読出し要求および同期動作の)多分性能にとってクリティカルな通信に利用することができたはずの帯域幅を消費する。

**【0018】**

従来のデータ処理システムの性能は、I/Oアダプタと呼出しプロセスの間の通信を可能にする割り込みハンドラの使用によってさらに低下する。先に指摘したように、従来の実装においては、データ転送が完了したときI/Oアダプタが割り込みをアサートし、割り込みハンドラは、データ転送が完了したことを呼出しプロセスに知らせるためにシステム・メモリ中でポーリング・フラグをセットする。I/Oアダプタと呼出しプロセスの間の通信を容易にするための割り込みの使用は、それによって各データ転送ごとに2つのコンテ

キスト切換えが必要となり、有用な仕事を実施するのではなく割込みハンドラを実行してプロセッサ・サイクルが消費されるので非効率的である。

【0019】

本発明はさらに、データ処理システム内の他のデータと異なるコヒーレンシ・ドメイン内でI/Oデータを管理することが多くの場合に望ましくないことを認識したものである。

【0020】

本発明はまた、データ処理システム性能は、たとえばI/O通信を実施するために用いられる不必要な命令をバイパスすることによりさらに改善できることも認識したものである。たとえば、複数レイヤのプロトコル(たとえば、TCP/IP)を使用するI/O通信では、コンピュータ間でのデータグラムの伝送は、送信側コンピュータでも、受信側コンピュータでも、プロトコル・スタックをトラバースするデータグラムを必要とする。多くのデータ転送では、しばしば結果として得られるアドレス・ポインタ、データ値、またはその他の実行結果が変化することなく、プロトコル・レイヤのうち少なくとも一部の命令が繰り返して実行される。したがって、本発明は、I/O性能、より一般的にはデータ処理システム性能が、かかる反復的コード・シーケンス内の命令をバイパスすることによって、かなり改善できることを認識したものである。

10

【課題を解決するための手段】

【0021】

本発明では、改善された処理装置、データ処理システム、およびデータ処理方法を提供することにより、当技術分野の前記その他の欠点に対処する。本発明の少なくとも一実施形態では、処理装置などの集積回路は、基板および基板中に形成された集積回路要素を含む。この集積回路要素は、命令を実行するプロセッサ・コアと、このプロセッサ・コアに結合され、プロセッサ・コアと集積回路の外部のシステム相互接続との間の通信をサポートする相互接続インターフェースと、プロセッサ・コアに結合され、入出力通信リンクを介して入出力通信をサポートする外部通信アダプタの少なくとも一部分とを含む。

20

【0022】

本発明の目的、特徴、および利点は、すべて以下の詳細に記述された説明で明らかとなる。

【0023】

本発明の特徴と考えられる新しい特徴は、添付の特許請求の範囲中に記載されている。しかしながら、本発明自体、ならびに好ましい使用形態、そのさらなる目的および利点は、実例としての実施形態についての以下の詳細な説明を添付図面と併せ参照することによって最も良く理解されよう。

30

【発明を実施するための最良の形態】

【0024】

再び図面を、具体的には図2を参照すると、本発明を有利に利用することができるネットワーク・システム70の一例が示されている。図に示すように、ネットワーク・システム70は、データ通信ができるようにネットワーク74によって結合された少なくとも2つのコンピュータ・システム(すなわち、ワークステーション・コンピュータ・システム72およびサーバ・コンピュータ・システム100)を含む。ネットワーク74は、任意の数の通信プロトコルを使用する、1つまたは複数の有線、無線、または光によるローカル・エリア・ネットワーク(たとえば、企業のイントラネット)または広域ネットワーク(たとえば、インターネット)を含むことができる。さらに、ネットワーク74はパケット交換サブネットワークおよび回路交換サブネットワークの一方または両方を含むことができる。以下に詳細に説明するように、本発明によれば、データは、入出力(I/O)データ通信のための革新的な方法、システム、および機器を利用して、ワークステーション72およびサーバ100により、あるいはワークステーション72とサーバ100の間でネットワーク74を介して転送することができる。

40

【0025】

50

次に、図3を参照すると、本発明による、改善されたI/O通信を含めて改善されたデータ処理をサポートする、マルチプロセッサ(MP)サーバ・コンピュータ・システム100の実施形態の一例が示されている。図に示すように、サーバ・コンピュータ・システム100は、複数の処理装置102を含み、これらの複数の処理装置は、それぞれメモリ104の1つに結合されている。各処理装置102はさらに、諸処理装置102間でのデータ、命令、および制御情報の通信をサポートする、統合された分散スイッチング機構106に結合される。各処理装置102は、集積回路要素がその上に形成された半導体基板を備える単一の集積回路として実施することが好ましい。複数の処理装置102、およびスイッチング機構106の少なくとも一部分は、有利なことに、共通のバックプレーンまたはチップ・キャリア上に一緒にパッケージ化することができる。

10

**【0026】**

図3にさらに示すように、本発明によれば、1つまたは複数の処理装置102は、スイッチング機構106とは独立したI/O通信を行うためにI/O通信リンク150に結合される。さらに以下で説明するように、通信リンク150に処理装置102を結合することにより、I/O通信をかなり簡単にし、性能を改善することができる。

**【0027】**

データ処理システム100が図示されていない追加の多くの構成要素を含むことができることが当業者には理解されよう。かかる追加の構成要素は本発明を理解するためには必要がないので、これらは図3に示さず、また本明細書でさらに論じてはいない。しかしながら、本発明によって提供されるI/O通信に対する拡張は、任意のシステム・アーキテクチャのデータ処理システムに適用可能であり、一般化したMPアーキテクチャまたは図3に示すSMPシステム構造だけに決して限定されないことも理解すべきである。

20

**【0028】**

次に図4を参照すると、サーバ・コンピュータ・システム100内における処理装置102の実施形態の一例のより詳細なブロック図が示されている。図に示すように、処理装置102内の集積回路要素は、1つまたは複数の命令スレッドをそれぞれ独立かつ同時に実行できる1つまたは複数のプロセッサ・コア108を含む。処理装置102はさらに、プロセッサ・コア108に結合され、プロセッサ・コア108によってアクセスされる可能性の高いデータおよび命令用の待ち時間の小さな記憶域を提供するキャッシュ階層構造110をさらに含む。キャッシュ階層構造110は、たとえば、別々の分岐した、各プロセッサ・コア108ごとのレベル1(L1)の命令キャッシュおよびデータ・キャッシュ、および複数のプロセッサ・コア108によって共用されるレベル2(L2)の大きなキャッシュを含むことができる。かかる各キャッシュは、従来の(または従来と異なった)キャッシュ・アレイ、キャッシュ・ディレクトリ、およびキャッシュ制御装置を含むことができる。キャッシュ階層構造110は、キャッシュされたデータおよび命令のコヒーレンシ状態を追跡するために、そのキャッシュ・ディレクトリ内で、周知のMESI(Modified, Exclusive, Shared, Invalid)キャッシュ・コヒーレンシ・プロトコルまたはそれらの変形プロトコルを実施することが好ましい。

30

**【0029】**

キャッシュ階層構造110はさらに、高周波数の広い帯域幅のメモリ・バス118によって処理装置102に結合されたメモリ104へのアクセスを制御するIMC(integrated memory controller、集積化メモリ制御装置)112に結合される。すべての処理装置102のメモリ104がまとまって、サーバ・コンピュータ・システム100内の最下レベルの揮発性メモリ(しばしば「システム・メモリ」と呼ばれる)を形成し、このメモリは、一般にすべての処理装置102からアクセス可能である。

40

**【0030】**

処理装置102はさらに、スイッチング機構106用のIFI(integrated fabric interface、集積化機構インターフェース)114を含む。IFI114は、IMC112にもキャッシュ階層構造110にも結合されており、スイッチング機構106に対してプロセッサ・コア108が要求する動作を支配するマスタ回路、ならびに(たとえば、コヒ

50

ーレンシを維持するためにキャッシュ階層構造 110 に対する動作をスヌープすることにより、または関連するメモリ 104 から要求されたデータを検索することにより)スイッチング機構 106 から受け取る動作に応答するスヌープ回路を含む。

#### 【0031】

処理装置 102 は、プロセッサ・コア 108 およびメモリ・バス 118 に結合された 1 つまたは複数の ECA (external communication adapters、外部通信アダプタ) 130 も有する。各 ECA 130 は、処理装置 102 がその一部を構成する MP サブシステムの外部の (または任意選択でサーバ・コンピュータ・システム 100 の外部の) 装置またはシステムとの I/O 通信をサポートする。様々な I/O 通信オプションを実現するために、処理装置 102 には、それぞれあるいは全体として (たとえばイーサネット (R)、S O N E T、PCIExpress、InfiniBandなどの) 様々な通信プロトコルを実施する ECA 130 を設けることができる。

10

#### 【0032】

好ましい実施形態では、IMC 112、IFI 114、および ECA 130 のそれぞれは、1 つまたは複数のオペレーティング・システムによって割り当てられた実効 (または実) アドレスを有するメモリ・マッピングされたリソースである。かかる実施形態では、処理装置 102 は、IMC 112、IFI 114、および ECA 130 へのアドレスの割当てを記録する MM (memory map、メモリ・マップ) 122 を備えている。したがって、各処理装置 102 は、コマンドのタイプまたはメモリ・マップ 122 内で提供されるアドレス・マッピングあるいはその両方に基づいて、IMC 112、IFI 114、および ECA 130 のどれかにコマンド (たとえば、I/O 書込みコマンドまたはメモリ読出し要求) を経路指定することが可能である。好ましい実施形態では、IMC 112 および ECA 130 は、同じダイ中に一体化された個々のプロセッサ・コア 108 とはどのような密接な関係ももたず、任意の処理装置 102 の任意のプロセッサ・コア 108 からアクセス可能になっていることに留意されたい。さらに、ECA 130 は、I/O 読出し動作および I/O 書込み動作を実施するために、サーバ・コンピュータ・システム 100 内の任意のメモリ 104 にアクセスすることができる。

20

#### 【0033】

より具体的に ECA 130 について検討すると、各 ECA 130 は、少なくとも DTL (data-transfer logic、データ転送ロジック) 133 およびプロトコル・ロジック 134 を含み、任意選択の I/O MC (I/O memory controller、I/Oメモリ制御装置) 131 をさらに含むことができる。DTL 133 は、プロセッサ・コア 108 間で通信リンク 150 へのアクセスの調停を行い、プロセッサ・コア 108 による I/O 読出しコマンドおよび I/O 書込みコマンドに回答して、通信リンク 150 とメモリ 104 の間でのデータ転送を制御する制御回路を含む。メモリ 104 にアクセスするために、DTL 133 は、任意の IMC 112 にメモリ読出し要求およびメモリ書込み要求を発行し、あるいは、専用 I/O MC 131 にかかるメモリ・アクセス要求を発行することにより、メモリ 104 にアクセスすることができる。I/O MC 131 は、複数のメモリ・アクセス要求または着信 I/O データもしくは発信 I/O データあるいはその両方をバッファするための任意選択のバッファ記憶装置 132 を含むことができる。

30

40

#### 【0034】

各 ECA 130 の DTL 133 はさらに、TLB (Translation Lookaside Buffer、変換索引バッファ) 124 に結合され、これらのバッファは、プロセッサ・コア 108 によって使用される実効アドレス (EA) を実アドレス (RA) に変換するのに利用される PTE (Page Table Entries、ページ・テーブル・エントリ) のサブセットのコピーをバッファする。本明細書中では、実効アドレス (EA) は、仮想アドレス空間にマッピングされたメモリ記憶位置またはその他のリソースを識別するアドレスとして定義される。他方、実アドレス (RA) は、本明細書において、実メモリ記憶位置またはその他の実リソースを識別する、実アドレス空間内のアドレスとして定義される。TLB 124 は、1 つまたは複数のプロセッサ・コア 108 で共用することができ、あるいは 1 つまたは複数の D

50

TL133によって専用で使用される別個のTLBを備えることができる。

【0035】

本発明の重要な一態様によれば、図7を参照して以下に詳細に説明するように、DTL133は、TLB124にアクセスして、I/O動作で転送すべきI/Oデータのソース・アドレスまたは宛先アドレスとしてプロセッサ・コア108によって指定される目標EAを、RAに変換する。したがって、本発明によって、I/Oアドレス変換を実施するために従来技術のようにTCE24(図1参照)を使用する必要が全くなくなり、システム・メモリ中でTCEを作成し管理するための付随的なOSオーバーヘッドが完全に取り除かれる。

【0036】

再びECA130を参照すると、プロトコル・ロジック134は、着信I/Oデータおよび発信I/Oデータをバッファするための複数のエントリ136を含むデータ・キュー135を含む。以下で説明するように、これらのハードウェア・キューは、バッファ132またはメモリ104あるいはその両方の中の仮想キューで補完することができる。さらに、プロトコル・ロジック134は、通信リンク150のレイヤ2プロトコルに則って発信I/Oデータを処理し、着信I/Oデータを処理して、たとえばレイヤ2ヘッダを取り除き、その他のデータ・フォーマット化を実施するLLC(link layer controller、リンク・レイヤ制御装置)138を含む。典型的な応用例では、プロトコル・ロジック134は、通信リンク150上で伝送すべき発信データを直列化し、通信リンク150から受信する着信データを非直列化するSER/DES(serializer/deserializer、並直列変換器/直並列変換器)140をさらに含む。

【0037】

理解を容易にするために、各ECA130は完全に別個の回路を有するものとして図4に示してあるが、一部の実施形態では、ダイ面積の効率的な使用を促進するために複数のECA130が共通の回路を共用できることを理解されたい。たとえば、複数のECA130が、1つのI/O MC131を共用することができる。その代わりにまたはそれに加えて、プロトコル・ロジック134の複数のインスタンスを、DTL133の1つのインスタンスによって制御し、1つのインスタンスに接続することもできる。かかる代替実施形態も本発明の範囲に含まれると理解すべきである。

【0038】

図4にさらに示すように、処理装置102内に集積された各ECA130の部分は実装特有(implementation-specific)であり、本発明の異なる実施形態によって変わるものである。たとえば、実施形態の一例では、ECA130aのI/O MC131およびDTL133は、処理装置102内に一体化されるが、ECA130aのプロトコル・ロジック134は、処理装置102のピン数およびダイ・サイズを減少させるためにオフチップのASIC(Application Specific Integrated Circuit、特定用途向けIC)として実施される。それとは対照的に、ECA130nはすべて、処理装置102の基板内に一体化される。

【0039】

各ECA130が、従来技術のI/Oアダプタ(たとえば、図1のPCI I/Oアダプタ50)と比較してかなり簡略化されていることに留意されたい。具体的には、従来技術のI/Oアダプタは、一般にSMPバス・インターフェース・ロジック、ならびに様々なアクティブ・セッションおよび「動作中(in flight)の」バス・トランザクションの状態を維持するための1つまたは複数のハードウェアまたはファームウェアの状態機械を含む。I/O通信は、従来のSMPバス上では経路指定されないため、ECA130は、従来のSMPバス・インターフェース回路を必要としない。さらに、図5および図7に関して以下に詳細に論じるように、かかる状態機械は、メモリ内にセッション状態の情報をI/Oデータと共に記憶することによって、ECA130内では規模縮小または除去される。

【0040】

10

20

30

40

50

さらに、処理装置 102 内に I/O ハードウェアを組み込むことにより、I/O データ通信を、スイッチング機構 106 上のデータ通信と同様に完全にキャッシュ・コヒーレントとすることができることに留意されたい。すなわち、各処理装置 102 内のキャッシュ階層構造 110 が、キャッシュ可能なデータを転送する I/O 読出しおよび書込み動作の検出に应答して、キャッシュされたデータ・グラニューールのコヒーレンシ状態を適切に更新することが好ましい。たとえば、キャッシュ階層構造 110 は、I/O 読出し動作中に指定されるアドレスと一致するアドレスを有するキャッシュ済みのデータ・グラニューールを無効にする。同様にキャッシュ階層構造 110 は、キャッシュされたデータ・グラニューールのアドレスと一致するアドレスを指定する I/O 書込み動作に应答して、キャッシュ階層構造 110 内にキャッシュされるデータ・グラニューールのコヒーレンシ状態を、排他的なキャッシュ・コヒーレンシ状態（たとえば、MESI の排他的状態または修正状態）から共用状態（たとえば MESI の共用状態）に更新する。さらに、I/O 書込み動作中に伝送されるデータ・グラニューールは、共用状態および無効状態だけに限定されるのではなく、修正状態（たとえば MESI の修正状態）または排他的状態（たとえば、MESI の排他的状態または修正状態）として伝送することもできる。かかるデータ転送のスヌープに应答して、キャッシュ階層構造 110 は、対応するキャッシュ・ラインを無効化（あるいは対応するキャッシュ・ラインのコヒーレンシ状態を更新）することになる。

#### 【0041】

多くの場合、キャッシュ済みのデータのコヒーレンシ状態に影響を及ぼす I/O 通信は、スイッチング機構 106 を介してメモリ 104 と ECA 130 の間で I/O データが通信されるために、複数の処理装置 102 のキャッシュ階層構造 110 によってスヌープされることになる。しかしながら、場合によっては、特定の I/O 通信セッションに關与する ECA 130 およびメモリ 104 が、共に同じ処理装置 102 に関連づけられることもある。したがって、I/O セッション中の I/O 読出し動作および I/O 書込み動作は、処理装置 102 内で内部伝送され、他の処理装置 102 からは見えないものとなる。そうした場合、I/O データ転送のマスタ（たとえば、ECA 130）またはスヌープ機構（たとえば、IFI 114 または IMC 112）のいずれかが、スイッチング機構 106 上に 1 つまたは複数のアドレスのみのデータ・キル（data kill）動作またはデータ共用のコヒーレンシ動作を伝送して、他の処理装置 102 中のキャッシュ階層構造 110 に、I/O データに関連するディレクトリ・エントリを適切なキャッシュ・コヒーレンシ状態に更新させることが好ましい。

#### 【0042】

次に図 5 を参照すると、サーバ・コンピュータ・システム 100 内の処理装置 102 に結合されたメモリ 104 の記憶内容のより詳細なブロック図が示されている。メモリ 104 は、たとえば 1 つまたは複数の DRAM（dynamic random access memory、ダイナミック・ランダム・アクセス・メモリ）デバイスを備えることができる。

#### 【0043】

図に示すように、ハードウェアまたはソフトウェアあるいはその両方により、メモリ 104 中の利用可能な記憶域を、関連する処理装置 102 のプロセッサ・コア 108 に割り付けられた少なくとも 1 つのプロセッサ領域 249 と、関連する処理装置 102 の 1 つまたは複数の ECA 130 に割り付けられた少なくとも 1 つの I/O 領域 250 と、サーバ・コンピュータ・システム 100 内の処理装置 102 のすべてに割り付けられ、処理装置 102 のすべてからアクセス可能な共用領域 252 とに分割することが好ましい。プロセッサ領域 249 は、関連する処理装置 102 の各プロセッサ・コア 108 が実行する命令を列挙する任意選択の命令トレース・ログ 260 を記憶する。所望の実装形態によっては、すべてのプロセッサ・コア 108 の命令トレース・ログを、同じプロセッサ領域 249 に記憶することができる。あるいは、各プロセッサ・コア 108 が、それ自体の専用プロセッサ領域 249 中にそれぞれその命令トレース・ログ 260 を記憶することもできる。

#### 【0044】

I/O 領域 250 は、それぞれが I/O データ転送用のパラメータを指定する 1 つまた

は複数のDTCB(Data Transfer Control Blocks、データ転送制御ブロック)253を記憶することができる。I/O領域250はさらに、各ECA130または各I/Oセッションごとに、プロトコル・ロジック134内の物理ハードウェア・キュー135を補完する仮想キュー254と、着信I/Oデータまたは発信I/Oデータの一時記憶を実現するI/Oデータ・バッファ255と、I/OセッションまたはECA130の制御状態情報をバッファする制御状態バッファ256とを含むことが好ましい。たとえば、制御状態バッファ256は、かかるコマンドがDTL133によって処理される準備ができるまで、1つまたは複数のI/Oコマンドをバッファすることができる。さらに、セッション状態の概念を使用するI/O接続では、制御状態バッファ256は、セッション状態情報を、多分、ポインタ、またはI/Oデータ・バッファ255中に記憶されるI/Oデータと他の構造化関連付けとともに記憶することができる。

10

#### 【0045】

図5にさらに示すように、共用領域252は、様々な処理装置102が実行できるソフトウェア158の少なくとも一部分と、処理装置102の1つが受信または送信するI/Oデータ262を含むことができる。さらに、共用領域252は、以前に論じたように、実効アドレス(EA)と実アドレス(RA)の間で変換するために利用されるPTE(Page Table Entries、ページ・テーブル・エントリ)の少なくとも一部分を含むOS作成ページ・テーブル264をさらに含む。

#### 【0046】

次に図6を参照すると、図2～図3のサーバ・コンピュータ・システム100のソフトウェア構成158の一例のソフトウェア・レイヤ図が示されている。図に示すように、ソフトウェア構成は、その最下位レベルに、システム・スーパーバイザ(またはハイパーバイザ)160を有し、このスーパーバイザ160は、データ処理システム8中で同時実行される1つまたは複数のオペレーティング・システム162にリソースを割り付ける。オペレーティング・システム162の各インスタンスに割り付けられるリソースをパーティションと呼ぶ。したがって、たとえば、ハイパーバイザ160は、2つの処理装置102をオペレーティング・システム162aのパーティションに、4つの処理装置102をオペレーティング・システム162bのパーティションに、別の処理装置102を(タイム・スライシングまたはマルチスレディングによって)複数のパーティションに、などと割り付けることができ、ある範囲の実アドレス空間および実効アドレス空間を各パーティションにそれぞれ割り付けることができる。

20

30

#### 【0047】

オペレーティング・システム162、ミドルウェア163、およびアプリケーション・プログラム164が、ハイパーバイザ160上で実行される。当業者には良く理解されるように、各オペレーティング・システム162は、ハイパーバイザ160が各オペレーティング・システム162に割り付けたリソースのプールからのアドレスおよび他のリソースを様々なハードウェア構成要素およびソフトウェア・プロセスに割り付け、そのパーティションに割り付けられたハードウェアの動作を独立に制御し、ページ・テーブル264を作成し管理し、そのアプリケーション・プログラム164がオペレーティング・システム・サービスにアクセスするための様々なAPI(application programming interfaces、アプリケーション・プログラム・インターフェース)を提供する。これらのOS APIには、ソケット・インターフェース、およびI/Oデータ転送をサポートするその他のAPIが含まれる。

40

#### 【0048】

アプリケーション・プログラム164は、広範な様々な計算処理、制御、通信、データ管理、およびプレゼンテーションの諸機能のどれかを実施するようにプログラムすることができ、いくつかのユーザ・レベル・プロセス166を含む。先に指摘したように、I/Oデータ転送を実施するために、プロセス166は、OS APIを介してその下にあるOS162にコールを行って、I/Oデータ転送をサポートする様々なOSサービスを要求する。

50

## 【0049】

図7を参照すると、本発明による、I/Oデータ通信方法の一例の高レベルの論理フローチャートが示されている。図7に示すプロセスを、図4に示すハードウェアおよび図5に示すメモリ・ダイアグラムをさらに参照して説明することにする。

## 【0050】

図に示すように、図7のプロセスは、ブロック180で開始し、次いでブロック181へと進む。このブロック181は、I/O読出し動作およびI/O書込み動作のためのI/O要求を発行する要求プロセス（たとえば、アプリケーション、ミドルウェア、またはOSプロセス）を示す。重要なことであるが、要求プロセスが要求されたI/O動作のためにアダプタまたはメモリのロックを取得する必要はない。というのは、ECA130が処理装置102内に一体化されていること、およびそれが提供する通信により、ECA130は、I/Oコマンドのサービスができるまでプロセッサ・コア108によるI/Oコマンドを「ホールド・オフ」することが可能であり、その代わりにまたはそれに加えて、それに続く処理のために多数のI/Oコマンドを、バッファ132または制御状態バッファ256あるいはその両方にバッファすることが可能だからである。以下に論じるように、この「ホールド・オフ」時間がある場合、バッファ132または255の1つにI/Oデータをローカルにバッファすることにより、それを最小にすることができる。

## 【0051】

所望のプログラミング・モデルに応じて、要求プロセスからのI/O要求をOSの関与の下で、またはOSの関与なしに処理することができる（このOSの関与についてはI/O要求内のフィールドに応じて選択的なものとして行うことができる）。I/O要求をOSが処理すべき場合、I/O要求は、OS162にI/O通信サービスを要求するAPIコールであることが好ましい。APIコールに 응답して、OS162は、ブロック182に示すように、要求されたI/O転送用のパラメータを指定するDTCB（Data Transfer Control Block、データ転送制御ブロック）を構築する。次いで、OS162は、DTCBの記憶位置（たとえば、ベースEA）の指示を要求側プロセスに戻すことができる。

## 【0052】

あるいは、I/O要求をOSの関与なしに処理すべき場合には、ブロック182に示すように、プロセスがDTCBを構築することが好ましく、この構築はブロック181でI/O要求を発行するより前またはそれと同時に行うことができる。この場合、I/O要求は、ECA130にDTCBのベースEAを供給するために、選択されたECA130のDTL133にプロセッサ・コア108から伝送されるI/Oコマンドであることが好ましい。

## 【0053】

図5に示すように、DTCBを、処理装置102のローカル・メモリ104内の参照番号253の位置に構築することができる。あるいは、DTCBをプロセッサ・コア108内の専用の記憶位置または汎用のレジスタ・セット中に構築することもできる。実施形態の一例では、DTCBは、少なくとも、（1）I/Oデータ転送が、着信I/OデータのI/O読出しかそれとも発信I/OデータのI/O書込みかを示すフィールドと、（2）I/O動作によりI/Oデータが転送される転送元または転送先の（たとえば、システム・メモリ104中の）1つまたは複数の記憶位置を識別する1つまたは複数の実効アドレス（EA）を示すフィールドと、（3）I/Oデータを受信または供給するリモート装置、システム、またはメモリ位置を識別するフォーリン・アドレス（たとえば、IP（Internet Protocol、インターネット・プロトコル）アドレス）の少なくとも一部分を示すフィールドとを含む。

## 【0054】

その後、図7に示すプロセスは、ブロック183に進み、このブロックでは、選択されたECA130のDTL133にDTCBを渡す。理解されるように、プロセッサ・コア108がDTCBをDTL133に「プッシュ」することができ、あるいは、たとえばI/O MC131またはIMC112に対して1つまたは複数のメモリ読出し動作を発行

10

20

30

40

50

することにより、DTL133がDTCBを「プル」することもできる。(かかるメモリ読出し動作には、TLB124を利用したEA-RA変換が必要となることもある。)DTCBの受信に応答して、DTL133はDTCBを検査して、要求されたI/O動作がI/O読出しか、それともI/O書込みかを判定する。DTCBがI/O読出し動作を指定する場合、図7に示すプロセスは、ブロック184から以下に説明するブロック210へと進む。しかし、DTCBがI/O書込み動作を指定する場合には、図7のプロセスは、ブロック184からブロック190へと進む。

#### 【0055】

ブロック190では、DTCB内で指定されるI/Oデータの1つまたは複数のEAを、1つまたは複数のメモリ104中のI/Oデータにアクセスするために利用できるRAに変換するために、DTL133のTLB124(図4参照)にアクセスする。実効アドレスから実アドレスへの変換を実施するのに必要なPTEがTLB124中に存在する場合、ブロック192でTLBのヒットが起こり、TLB124はDTL133に対応するRAを供給する。次いで、プロセスは、ブロック192から以下に説明するブロック200へと進む。しかし、必要なPTEが、TLB124に現在バッファされていない場合には、ブロック192でTLBミスが起こり、プロセスはブロック194に進む。ブロック194で、OSは、実効アドレス-実アドレス変換を実施するのに必要なPTEをページ・テーブル264からTLB124中にロードするために、従来型のTLB再ロード動作を実施する。次いで、プロセスは、ブロック200に進む。

#### 【0056】

ブロック200で、ローカル・メモリ104からI/Oデータを得るために、実アドレスを含む読出し要求をI/O MC131(または、I/O MCが実装されていない場合はIMC112)に発行することにより、また他のメモリ104からI/Oデータを得るために、実アドレスを含む読出し要求をIFI114に発行することにより、DTL133は、DTCB中で識別される、システム・メモリ104からI/Oデータにアクセスする。I/Oデータが伝送を待つ間、DTL133は、発信I/Oデータを1つまたは複数のバッファ132および255に一時的にバッファすることができる。重要なことであるが、このようにデータをバッファすることにより、DTL133(または要求中のプロセス)がI/Oデータ用のロックを取得することを必要とせずに、バッファ済みのI/Oデータが、伝送以前に変更されないように保護され、それによって、システム・メモリ104中のデータのコピーに1つまたは複数のプロセスがアクセスしそれを変更することが可能になる。その後、ブロック202に示すように、DTL133は、プロトコル特有のデータグラムおよびメッセージを利用して、発信I/Oデータをキュー135およびLLC138(および必要ならSER/DES140)を介して通信リンク150に伝送する。かかる伝送は、DTCBが指定するデータがすべて送信されるまで継続する。その後、プロセスは、以下に説明するブロック242に進む。

#### 【0057】

再び図7のブロック184を参照すると、DTL133が、DTCB中で指定されたI/O動作がI/O読出し動作であると判定したのに応答して、プロセスはブロック210に進む。このブロック210では、I/Oデータの受信の準備が整ったことを示すために、DTL133が、プロトコル・ロジック134および通信リンク150を介してネットワーク74上にI/O読出し要求を送信開始(launch)する。次いで、ブロック212では、データグラムがネットワーク74から受信されるまでプロセスを繰り返す。

#### 【0058】

プロトコル・ロジック134がネットワーク74からデータグラムを受信するのに応答して、データグラムがDTL133に渡されるが、このDTL133が、バッファ132、255の1つにデータグラムをバッファすることが好ましい。さらに、ブロック214に示すように、DTL133は、TLB124にアクセスして、データグラムが指定するEAの変換結果を得ようとする。EAを変換するための関連PTEがTLB124中にバッファされている場合、TLBヒットがブロック216で起こり、DTL133は目標メ

10

20

30

40

50

メモリ位置のRAを受け取り、プロセスは、以下に説明するブロック240に進む。しかし、ブロック216でTLBミスであった場合はそれに応答して、プロセスはブロック220に進む。ブロック220で、指定されたEAを変換するために必要なPTEを得るために、OSがシステム・メモリ104中のページ・テーブル264にアクセスする。ブロック230~232に示すように、TLBの再ロード動作の完了を待つ間、I/O読出しを停止することができ、あるいは、着信データを1つまたは複数のバッファ132および255にバッファしながらI/O読出しを続行することもできる。TLBの再ロード動作が完了し、I/O読出し動作のためのRAが得られた後、プロセスは、ブロック240に進む。このブロック240では、RAを指定する1つまたは複数のメモリ書込み動作を発行することにより、DTL133はI/O読出しデータを(たとえば、1つまたは複数のバッファ132、255から)メモリ104の1つに記憶する。

10

**【0059】**

たとえば、I/O読出し動作で大量のデータを読み出す場合、スイッチング機構106が多用される場合、あるいはスイッチング機構106を介したメモリ記憶動作に関連する待ち時間が望ましくないほど長い場合など、一部のケースでは、スイッチング機構106を介して伝送されるI/Oデータの量を最小限に抑えるのが望ましいこともある。したがって、ブロック214~240に示すアドレス変換プロセスの改善策として、OSは、I/Oデータを強制的にECA130にとってローカルなメモリ104に記憶させるように選択的に決定することができる。その場合には、OSは、着信I/Oデータグラムに関連するEAを変換するためのページ・テーブル264をローカル・メモリ104中の記憶位置に関連するRAで更新する。その結果、ブロック240に示す記憶ステップでは、ブロック214と232の一方で行われるEA-RA変換に基づいて、着信I/Oデータのすべてをローカル・メモリ104の共用メモリ領域252内のメモリ位置に記憶する。

20

**【0060】**

プロセスは、ブロック202またはブロック240のいずれかからブロック242に進む。ブロック242で、ECA130がI/Oデータの完了の指示を要求プロセスに提供する。完了指示は、たとえば、DTCB内の完了フィールド、ECA130中のメモリ・マッピングされた記憶位置、あるいはプロセッサ・コア108内の条件レジスタ・ビットなどその他の完了指示を含むことができる。要求プロセスは、I/Oデータ転送が完了したことを検出するために(たとえば、読出し要求を発行することにより)完了指示をポーリングすることができ、あるいは、完了指示の状態変化によって、ローカルな(すなわち、オンチップの)割込みをトリガすることもできる。重要なことであるが、本発明では、I/Oデータ転送が完了したことを、要求プロセスに信号で知らせるための従来からのI/O割込みは必要でない。その後、図7に示すプロセスは、ブロック250で終了する。

30

**【0061】**

次に図8を参照すると、本発明によるプロセッサ・コア108の実施形態の一例のより詳細なブロック図が示されている。図に示すように、プロセッサ・コア108は、ISU(instruction sequence unit、命令シーケンス・ユニット)270およびいくつかの実行ユニット282~290を含む、命令パイプラインを含む。ISU270は、IMMU(instruction memory management unit、命令メモリ管理ユニット)272によって実施されるERAT(effective-to-real address translation、実効-実アドレス変換)によって得られる実アドレスを利用して、処理用の命令をL1 I-キャッシュ274からフェッチする。もちろん、要求された命令のキャッシュ・ラインがL1 I-キャッシュ274に存在しない場合には、ISU270は、I-キャッシュ再ロード・バス276を介してキャッシュ階層構造110(またはより下のレベルの記憶装置)中のL2キャッシュに関連する命令キャッシュ・ラインを要求する。

40

**【0062】**

命令がフェッチされ、前処理が行われる場合はそれが実施された後で、ISU270は、命令タイプに基づいて、命令バス280を介して実行ユニット282~290に、命令を、おそらく順不同にディスパッチする。すなわち、条件-レジスタ-修正(condition-

50

register-modifying) 命令および分岐命令は、それぞれ C R U (condition register unit、条件レジスタ・ユニット) 2 8 2 および B E U (branch execution unit、分岐実行ユニット) 2 8 4 にディスパッチされ、固定小数点命令およびロード/ストア命令は、それぞれ F X U (fixed-point unit、固定小数点ユニット) 2 8 6 および L S U (load-store unitロード・ストア・ユニット) 2 8 8 にディスパッチされ、浮動小数点命令は、F P U (floating-point unit、浮動小数点ユニット) 2 9 0 にディスパッチされる。

#### 【0063】

好ましい一実施形態では、各ディスパッチされた命令はさらに、関連メモリ 1 0 4 (図 5 参照) 中の命令トレース・ログ 2 6 0 に記録するために、トレース・バス 2 8 1 を介して I M C 1 1 2 に伝送される。代替実施形態では、I S U 2 7 0 は、プロセッサ・コア 1 0 8 の設計 (architected) 状態にコミットした完了命令だけをトレース・バス 2 8 1 を経由して伝送することができ、あるいは、命令トレース・ログ 2 6 0 に記録するためにどの命令 (たとえば、命令なし、ディスパッチされた命令または完了命令または特定の命令タイプのみあるいはこれらの組合せ) をメモリ 1 0 4 に伝送するかの選択を可能にする、関連するソフトウェアまたはハードウェア選択可能なモード・セレクタ 2 7 3 を有することができる。さらなる改良では、トレース・バス 2 8 1 がディスパッチされた命令のすべてをメモリ 1 0 4 に伝えること、および I S U 2 7 0 がディスパッチされた命令のどれが実際に完了したかを示す完了指示をメモリ 1 0 4 に伝送する。これらの実施形態のすべてで、アプリケーションその他のソフトウェア・プログラムの完全な命令トレースを、あまり回路に手を入れずに (non-intrusively)、またプロセッサ・コア 1 0 8 の性能を実質的に低下させずに達成することができる。

#### 【0064】

可能なキューイングおよびバッファリングの後、I S U 2 7 0 によってディスパッチされた命令が、機会があれば実行ユニット 2 8 0 ~ 2 9 0 によって実行される。命令「実行」は、本明細書では、プロセッサの論理回路が、命令オペレーション・コード (opcode)、および関連オペランドがある場合はそれも調査し、それに応答して、データまたは命令をデータ処理システム中 (たとえば、システム・メモリ位置間、レジスタもしくはバッファとメモリの間など) で移動し、またはデータに対して論理演算または数値演算を実施するプロセスと定義する。メモリ・アクセス (すなわちロード・タイプまたはストア・タイプの) 命令では、実行は、一般的に命令オペランドから目標 E A を計算することを含む。

#### 【0065】

実行ユニット 2 8 2 ~ 2 9 0 の 1 つの中での実行中に、命令は、入力オペランドがある場合はそれを、実行ユニットに結合されたレジスタ・ファイル 3 0 0 ~ 3 0 4 内の 1 つまたは複数の設計レジスタまたはリネーム・レジスタあるいはその両方から受け取ることができる。同様に、命令実行のデータ結果 (すなわち宛先オペランド) がある場合はそれが、レジスタ・ファイル 3 0 0 ~ 3 0 4 内の命令で指定される位置に実行ユニット 2 8 2 ~ 2 9 0 によって書き込まれる。たとえば、F X U 2 8 6 は、入力オペランドを G P R F (general-purpose register file、汎用レジスタ・ファイル) 3 0 2 から受け取り、宛先オペランド (すなわちデータ結果) を G P R F 3 0 2 に記憶し、F P U 2 9 0 は入力オペランドを F P R F (floating-point register file、浮動小数点レジスタ・ファイル) 3 0 4 から受け取り、宛先オペランドを F P R F 3 0 4 に記憶し、L S U 2 8 8 は、入力オペランドを G P R F 3 0 2 から受け取り、データを L 1 D - キャッシュ 3 0 8 と G P R F 3 0 2 および F P R F 3 0 4 との間で転送させる。同様に、条件 - レジスタ - 修正命令または条件 - レジスタ - 依存 (condition-register-dependent) 命令を実行するときは、C R U 2 8 2 および B E U 2 8 4 は、C R F (control register file、制御レジスタ・ファイル) 3 0 0 にアクセスし、この C R F 3 0 0 は、好ましい一実施形態では、それぞれの条件レジスタ、リンク・レジスタ、カウント・レジスタ、およびリネーム・レジスタを含む。B E U 2 8 4 は、バス・アドレスを得るための条件分岐を解決するために条件レジスタ、リンク・レジスタおよびカウント・レジスタの値にアクセスする。B E U 2 8 4 は、このバス・アドレスを、それが指示するバスに沿って命令フェッチを開始するために命

令シーケンス・ユニット 270 に供給する。実行ユニットは、命令の実行を終了した後、命令シーケンス・ユニット 270 に知らせ、この命令シーケンス・ユニット 270 では、プログラム順の命令の完了、および、データ結果がある場合はそのデータ結果のプロセッサ・コア 108 の設計状態へのコミットメントをスケジュールする。

【0066】

図 8 にさらに示すように、プロセッサ・コア 108 はさらに、取込みロジック 322 と、バイパス CAM (content addressable memory、連想記憶装置) 324 を備える命令バイパス回路 320 を含む。図 10 を参照して以下に説明するように、バイパス回路 320 は、プロセッサ・コア 108 が、I/O データ転送を実施するために利用されるものを含めて、反復的コード・シーケンスをバイパスするのを可能にし、それによってシステム性能がかなり改善される。

10

【0067】

次に図 9 を参照すると、命令バイパス CAM 324 のより詳細なブロック図が示されている。図に示すように、命令バイパス CAM 324 は、命令ストリーム・バッファ 340、ユーザ・レベル設計状態 CAM 343、およびメモリ・マップ・アクセス CAM 346 を含む。

【0068】

命令ストリーム・バッファ 340 は、それぞれがスヌープ・キル・フィールド 341 および命令アドレス・フィールド 342 を含むいくつかのバッファ・エントリを含む。命令アドレス・フィールド 342 は、コード・シーケンス中の命令のアドレス(または少なくとも上位アドレス・ビット)を含み、スヌープ・キル・フィールド 341 は、命令アドレスを対象とするストア・オペレーションまたはその他の無効オペレーションが、I/O チャンネル 150、ローカル・プロセッサ・コア 108 またはスイッチング機構 106 からスヌープされたかどうかを示す。したがって、命令ストリーム・バッファ 340 の内容は、命令シーケンス中の命令のどれかが、その最後の実行以来変更されたかどうかを示す。

20

【0069】

ユーザ・レベル設計状態 CAM 343 は、それぞれがプロセッサ・コア 108 のユーザ・レベルの設計状態の一部を形成するそれぞれのレジスタに対応するいくつかの CAM エントリを含む。各 CAM エントリはレジスタ値フィールド 345 を含み、このレジスタ値フィールド 345 は、命令ストリーム・バッファ 340 中に記録されるコード・シーケンスの始めと終りの時点での(たとえば、レジスタ・ファイル CRF 300、GPRF 302、および FPRF 304 内の)対応するレジスタの値を記憶する。したがって、CAM エントリのレジスタ値フィールドは、1 つがコード・シーケンスの始めで取り込まれ、第 2 のものがコード・シーケンスの終りで取り込まれる、プロセッサ・コア 108 のユーザ・レベル設計状態の 2 つの「スナップ・ショット」を含む。各 CAM エントリには、使用フラグ 344 が関連づけられており、これは、レジスタ値フィールド 345 内の関連レジスタ値が書き込み前のコード・シーケンスの途中で読み出されたかどうか(すなわち、初期レジスタ値が、コード・シーケンスの正しい実行にとって重要であるかどうか)を示す。この情報は、後で、CAM 343 中のどの設計値を比較する必要があるかどうかを判定するために使用される。

30

40

【0070】

メモリ・マップ・アクセス CAM 346 は、メモリ・アクセス命令および I/O 命令の目標アドレスおよびデータを記憶するためのいくつかの CAM エントリを含む。各 CAM エントリは、アクセス(たとえば、ロード・タイプまたはストア・タイプの)命令の目標アドレス、および目標アドレスによって識別される記憶位置またはリソースに書き込まれまたはそこから読み出されるデータを記憶するための目標アドレス・フィールド 348 およびデータ・フィールド 352 を有する。この CAM エントリはさらに、L/S (load/store、ロード/ストア)フィールド 349、および I/O フィールド 350 を含み、これらはそれぞれ、関連するメモリ・アクセス命令がロード・タイプ命令か、ストア・タイプ命令か、また関連アクセス命令が I/O 装置に割り付けられたアドレスを対象としたも

50

のかどうかを示す。メモリ・マップ・アクセスCAM346内の各CAMエントリはさらに、スヌープ・キル・フィールド347を含み、このフィールドは、その目標アドレスを対象とするストア・オペレーションまたはその他の無効オペレーションが、I/Oチャネル150、ローカル・プロセッサ・コア108またはスイッチング機構106からスヌープされたかどうかを示す。したがって、命令ストリーム・バッファ340の内容は、命令ストリーム・バッファ340に記録される命令シーケンスによって実施される作業が命令シーケンスを最後に実行してから、修正されたかどうかを示す。

#### 【0071】

図9は、1つの命令シーケンスに関連したバイパスCAM324内のリソースを示すが、任意数の多分繰返し可能な命令シーケンスの記憶を実現するためにかかるリソースを複製できることを理解されたい。

10

#### 【0072】

次に図10を参照すると、本発明による、プログラムの実行中に反復的コード・シーケンスをバイパスする方法の一例の高レベル論理フローチャートが示されている。図に示すように、プロセスはブロック360から開始し、このブロックでは、プロセッサ・コア108が、あるプロセス(たとえば、アプリケーション・プロセス、ミドルウェア・プロセスまたはオペレーティング・システム・プロセス)内の任意の点で命令を実行している。図8に示すプロセッサ・コアの実施形態では、命令バイパス回路320中の取込みロジック322は、ISU270によって生成される命令アドレスを受け取り、任意選択でまたは追加として、ISU270によってフェッチまたはディスパッチあるいはその両方が行われる命令を受け取るように結合される。たとえば、一実施形態では、取込みロジック322は、ISU270のIAR(instruction address register、命令アドレス・レジスタ)271に含まれる次の命令フェッチ・アドレスを受け取るように結合することができる。図9のブロック352に示すように、取込みロジック322は、繰返し実行されるコード・シーケンスの始めに一般に見出される、OS APIコールなどの命令があるかどうか、ISU270内の命令アドレスまたはオペレーション・コードあるいはその両方を監視する。取込みロジック322が反復的コード・シーケンスを開始するものと認識する、1つまたは複数の命令アドレスまたは命令オペレーション・コード(opcode)あるいはその両方に基づいて、取込みロジック322は、多分反復的コード・シーケンスが検出されたことを命令バイパスCAM324に知らせる「コード・シーケンス開始」指示を命令バイパスCAM324に伝送する。他の実施形態では、各命令アドレスをバイパスCAM324に単に提供するだけでよい。

20

30

#### 【0073】

ブロック364に示すように、「コード・シーケンス開始」信号または命令アドレスに回答して、命令バイパスCAM324は、多分反復的なコード・シーケンスをバイパスすべきか否かを判定する。この判定を行う際に、好ましい一実施形態では、バイパスCAM324は4つの要因を考慮に入れる。第1に、バイパスCAM324は、命令ストリーム・バッファ340を参照して、検出された命令アドレスが命令ストリーム・バッファ340内に記録された開始命令アドレスと一致するかどうかを判定する。第2に、バイパスCAM324は、ユーザ・レベルの設計状態CAM343を参照して、使用フィールド344がセットされた始めの各ユーザ・レベル設計状態レジスタの値が、検出命令の実行後のプロセッサ・コア108中の対応するレジスタの値と一致するかどうかを判定する。この比較をする際に、使用フィールド344がリセットされたレジスタ(すなわち、命令シーケンス中で使用されないレジスタ、または読み出される前に書き込まれたレジスタ)は考慮に入れない。第3に、バイパスCAM324は、命令ストリーム・バッファ340のスヌープ・キル・フィールド341を参照して、命令シーケンス中の命令のどれかがスヌープ・キル・オペレーションによって修正または無効化されたかどうかを判定する。第4に、バイパスCAM324は、メモリ・マップ・アクセスCAM346のスヌープ・キル・フィールド347を参照して、命令シーケンス中のアクセス命令の目標アドレスのどれかがスヌープ・キル・オペレーションの対象になったかどうかを判定する。

40

50

## 【 0 0 7 4 】

一実施形態では、バイパスCAM324が、4つの条件すべてが満たされたと判断する場合、すなわち検出された命令アドレスが記憶されたコード・シーケンスの最初の命令アドレスと一致し、ユーザ・レベル設計状態が一致し、命令シーケンスの命令アドレスまたは目標アドレスがスヌープ・キルを受けていない場合には、検出済みのコード・シーケンスをバイパスすることができる。より好ましい一実施形態では、第4の条件は、たとえスヌープ・キル・フィールド347により、ストア・タイプの命令の目標アドレスに対して（ただしロード・タイプではそうではない）、1つまたは複数のスヌープ・キルが指示される場合でも、バイパスCAM324はコードをバイパスさせるように修正される。これが可能なのは、以下にさらに論ずるように、コードのバイパスをサポートするために、スヌープ・キルの影響を受けるメモリ記憶動作を実施できるからである。 10

## 【 0 0 7 5 】

バイパスCAM324が、検出された命令で始まるコード・シーケンスをバイパスすることができないと判定する場合、プロセスは以下に説明するブロック380に進む。しかしながら、バイパスCAM324が、検出されたコード・シーケンスをバイパスすることができると判定する場合には、プロセスは、処理中のコア108が反復的コード・シーケンスをバイパスするブロック368に進む。

## 【 0 0 7 6 】

反復的コード・シーケンスのバイパスは、ISU270が、処理中のコア108の命令パイプライン中にある反復的コード・シーケンスに属する任意の命令をキャンセルし、その反復的コード・シーケンス中の命令を追加してフェッチしないようにするものであることが好ましい。さらに、バイパスCAM324は、最後のユーザ・レベル設計状態をユーザ・レベル設計状態CAM343からプロセッサ・コア108のユーザ・レベル設計レジスタにロードし、I/Oフィールド350によりI/Oリソースを対象とすると指示される命令シーケンス中の各アクセス命令を実施する。I/Oストア・タイプのオペレーションでは、データ・フィールド352からのデータが使用される。最後に、ストア・タイプ・オペレーションの目標アドレスに対するスヌープ・キルが存在するときに、コード・バイパスがサポートされる場合には、バイパスCAM324は、データ・フィールド352に含まれるデータを利用して、少なくともスヌープ・キルの影響を受けるメモリ・ストア・オペレーションがあれば、そのそれぞれのオペレーション（および任意選択で命令シーケンス中のあらゆるメモリ・ストア・オペレーション）を実施する。したがって、バイパスCAM324が反復的コード・シーケンスをバイパスすることを選択した場合には、バイパスCAM324は、プロセッサ・コア108のユーザ・レベル設計状態、メモリ・イメージ、およびプロセッサ・コア108のI/Oリソースが、まるでプロセッサ・コア108の実行ユニット282～290内で反復的コード・シーケンスを実際に行うように見えるようにするために必要なオペレーションをすべて実施する。その後、ブロック368からブロック390に進むプロセスで示すように、プロセッサ・コア108は、反復的コード・シーケンスに続く命令で始まるプロセス中で、命令の通常のフェッチと実行を再開し、それによって、反復的コード・シーケンスを含む1つまたは複数の（および任意数までの）非ノーオペレーション命令を実行する必要が完全になくなる。 20 30 40

## 【 0 0 7 7 】

次に図10のブロック380を参照すると、命令バイパスCAM324が多分反復的なコード・シーケンスをバイパスできないと判定する場合には、命令バイパスCAM324は、検出されたコード・シーケンスの始めのユーザ・レベル設計状態をユーザ・レベル設計状態CAM343内に記録し、検出されたコード・シーケンス中の命令の命令アドレスを、命令ストリーム・バッファ340の命令アドレス・フィールド342内に記録することを開始し、メモリ・アクセス命令に属する目標アドレス、データ結果、およびその他の情報をメモリ・マップ・アクセスCAM346中に記録することを開始する。判断ブロック384で示すように、命令バイパスCAM324は、取込みロジック322が反復的コード・シーケンスの終りを検出するまで、検出されたコード・シーケンスに属する情報を 50

記録し続ける。命令バイパスCAM324が一杯になるか、または取込みロジック322が反復的コード・シーケンスの終りを検出するのに応答して、たとえば、1つまたは複数の命令アドレスおよびオペレーション・コードまたは割込みイベントの発生に基づいて、取込みロジック322は、「コード・シーケンス終了」信号をバイパスCAM324に送る。ブロック386に示すように、「コード・シーケンス終了」信号の受取りに応答して、バイパスCAM324は、プロセッサ・コア108の終了ユーザ・レベル設計状態をユーザ・レベル設計状態CAM343中に記録し、その後記録を継続する。その後、ブロック390で命令の実行が続行され、コード・シーケンスを次にそれが検出されたときにバイパスするのに必要な情報がバイパスCAM324にロードされる。

**【0078】**

本明細書中に記載の命令バイパスは、推論的プロセッサ、非推論的プロセッサ、および順不同実行プロセッサ中で実施することができることに留意されたい。あらゆる場合に、コード・シーケンスをバイパスすべきか否かの判断はバイパスCAM324内に記憶される非推論的な情報に基づいており、プロセッサ・コア108の設計状態にまだコミットしていない推論的な情報に基づいてはいない。

**【0079】**

本発明の命令バイパス回路320は、任意長の反復的コードをバイパス可能にすることも理解されたい。この場合、最大可能コード・バイパス長は、少なくとも部分的にはバイパスCAM324の容量によって決まる。したがって、長いコード・シーケンスのバイパスをサポートすることが望ましい実施形態では、メモリ104など、一部または全部がオフチップ・メモリの形でバイパスCAM324を実施するのが望ましいこともある。一部の実施形態では、バイパスCAM324を命令トレース・ログ260に書き込むべき命令のオンチップ「キャッシュ」として使用すること、および、たとえば命令シーケンスがバイパスCAM324からの情報で置き換えられたときには、バイパスCAM324からの情報をメモリ104に繰り返し書き込むことが好ましいこともある。かかる実施形態では、命令トレース・ログ260に書き込まれる情報は、たとえばリンクしたリスト・データ構造を利用して、ストア・オペレーションの順序が維持されるように構造化されていることが好ましい。

**【0080】**

図9～図10は、理解を容易にするために、ユーザ・レベル設計状態だけに基づいたコード・バイパスを示しているが、コード・シーケンスをバイパスするか否かを判定する際に、追加のレイヤの状態情報を含めて追加の状態情報を考慮に入れることができることを理解されたい。たとえば、命令シーケンスをバイパスすべきかを判定するために、プロセッサ・コア108の現在のスーパーバイザ・レベルの設計状態と比較するためスーパーバイザ・レベルの設計状態を状態CAM343中に記録しておくこともできるはずである。かかる実施形態では、状態CAM343中に記録されるスーパーバイザ・レベル設計状態は、必ずしも命令シーケンスの始めではなく、命令シーケンスの内部でOSコールがなされた時点での「スナップ・ショット」とすることが好ましい。記憶されたユーザ・レベル設計状態と現在のユーザ・レベル設計状態が一致し、記憶されたスーパーバイザ・レベル状態と現在のスーパーバイザ・レベル状態が一致しない場合には、命令シーケンスの部分的なバイパスを依然として実施することができ、このバイパスは、命令シーケンスがスーパーバイザ・レベル設計状態に入る前に（たとえば、OSコールの前に）終わる。

**【0081】**

以上に説明してきたように、本発明は、データ処理の改善された方法、機器、およびシステムを提供する。一態様では、集積回路が、プロセッサ・コアと、入出力通信リンクを介して入出力通信をサポートする外部通信アダプタの少なくとも一部分との両方を含んでいる。プロセッサ・コアと同じ集積回路内にI/O通信アダプタを一体化することにより、一般的にデータ処理に対するいくつかの機能改善がサポートされ、具体的にはI/O通信がサポートされる。たとえば、同じ集積回路にI/O通信アダプタおよびプロセッサ・コアを一体化することにより、ロック取得の待ち時間、プロセッサ・コアとI/O通信ア

10

20

30

40

50

アダプタとの間の通信の待ち時間、およびI/Oアドレス変換の待ち時間を含めて、I/O通信の待ち時間をもたらす複数の原因を減らしたまたは取り除くことが容易になる。さらに、プロセッサ・コアおよびその関連するキャッシュと同じ集積回路内にI/O通信アダプタを集積化することにより、修正されたキャッシュ・コヒーレンシ状態および排他的なキャッシュ・コヒーレンシ状態をI/Oデータに割り当てることを含めて、完全にキャッシュ・コヒーレントなI/O通信が容易になる。

【0082】

他の一態様では、データ処理性能が、I/O通信プロセス中で一般に見られるような反復的コード・シーケンスの実行をバイパスすることによって改善される。

【0083】

さらに他の一態様では、関連する低レベルのメモリのプロセッサ・メモリ域内で各プロセッサ・コアの命令トレースを作成することにより、データ処理の振る舞いのテスト、検証、性能評価および監視が容易になる。

【0084】

本発明を、好ましい一実施形態に関して具体的に示し説明してきたが、本発明の趣旨および範囲を逸脱することなく、形態および細部の様々な変更を行うことができることが、当業者には理解されよう。

【0085】

まとめとして、本発明の構成に関して以下の事項を開示する。

【0086】

(1) 基板と、前記基板中に形成された集積回路要素とを備える集積回路であって、前記集積回路要素が、

命令を実行するプロセッサ・コアと、

前記プロセッサ・コアに結合され、前記プロセッサ・コアと前記集積回路の外部のシステム相互接続との間の通信をサポートする相互接続インターフェースと、

前記プロセッサ・コアに結合され、入出力通信リンクを介して入出力通信をサポートする外部通信アダプタの少なくとも一部分と

を含む集積回路。

(2) 前記集積回路要素がさらに、

前記プロセッサ・コアに結合されたキャッシュ階層構造と、

前記キャッシュ階層構造に結合された集積化メモリ制御装置と

を備える、上記(1)に記載の集積回路。

(3) 前記外部通信アダプタを含めて複数の外部通信アダプタの少なくとも一部分を含み、前記複数の外部通信アダプタのうちの少なくとも2つがそれぞれ異なる入出力通信プロトコルを実施する、上記(1)に記載の集積回路。

(4) 前記集積回路要素が、通信要求で指定されるアドレスに基づいて、前記集積回路要素内の通信コマンドを前記相互接続インターフェースおよび前記外部通信アダプタに経路指定するためのメモリ・マップを含む、上記(1)に記載の集積回路。

(5) 前記外部通信アダプタに結合された変換索引バッファをさらに備え、前記変換索引バッファが、入出力(I/O)コマンドで指定される実効アドレスを実アドレスに変換する、上記(1)に記載の集積回路。

(6) 前記外部通信アダプタの前記少なくとも一部分が、I/Oデータにアクセスするためにメモリにアクセスする入出力(I/O)メモリ制御装置を含む、上記(1)に記載の集積回路。

(7) 前記外部通信アダプタの前記少なくとも一部分が、

前記プロセッサ・コアおよびリンク・レイヤ制御装置に結合され、前記プロセッサ・コアによるI/Oコマンドに応答して、入出力(I/O)データ転送を制御するデータ転送ロジックを備える、上記(1)に記載の集積回路。

(8) 前記外部通信アダプタの前記少なくとも一部分が、リンク・レイヤ制御装置をさらに備える、上記(1)に記載の集積回路。

10

20

30

40

50

( 9 ) 前記外部通信アダプタが、I/Oデータ転送の完了を示すソフトウェア・アクセス可能な標識をセットする手段を含む、上記( 1 )に記載の集積回路。

( 10 ) 前記集積回路内の前記外部通信アダプタの前記部分が、着信I/Oデータと発信I/Oデータの少なくとも一方をバッファするバッファを含む、上記( 1 )に記載の集積回路。

( 11 ) 上記( 1 )に記載の少なくとも1つの集積回路と、  
前記相互接続インターフェースに結合されたシステム相互接続と、  
前記少なくとも1つの集積回路に結合されたメモリ・システムと  
を備えるデータ処理システム。

( 12 ) 前記集積回路が、第1の集積回路を備え、  
前記外部通信アダプタの前記少なくとも一部分が、第1の部分  
を備え、  
前記データ処理システムがさらに、前記第1の集積回路に接続された第2の集積回路を含み、前記外部通信アダプタが、前記第2の集積回路内で実施される第2の部分を含む、  
上記( 11 )に記載のデータ処理システム。

10

( 13 ) 前記第2の部分が、リンク・レイヤ制御装置を備える、上記( 12 )に記載のデータ処理システム。

( 14 ) 基板と、前記基板中に形成された集積回路要素とを含む第1の集積回路であって、  
前記集積回路要素が

命令を実行するプロセッサ・コアと、  
前記プロセッサ・コアに結合された相互接続インターフェースであって、前記プロセッサ・コアが前記少なくとも1つのプロセッサ・コアと前記第1の集積回路の外部のシステム相互接続との間での通信をサポートする相互接続インターフェースと、

20

前記プロセッサ・コアに結合され、入出力通信リンクを介して入出力通信をサポートする外部通信アダプタの第1の部分とを含む、第1の集積回路と、

前記第1の集積回路のピンに接続された第2の集積回路とを備え、前記外部通信アダプタが前記第2の集積回路内に実施される第2の部分を含むシステム。

( 15 ) 前記外部通信アダプタの前記第1の部分が、前記プロセッサ・コアに結合され、プロセッサ・コアによるI/Oコマンドにตอบสนองして入出力(I/O)データ転送を制御するデータ転送ロジックを備え、

前記外部通信アダプタの前記第2の部分が、リンク・レイヤ制御装置を備える、上記( 14 )に記載のシステム。

30

( 16 ) プロセッサ・コアを含む集積回路を動作させる方法であって、

前記集積回路内の相互接続インターフェースを利用して、プロセッサ・コアと前記集積回路の外部のシステム相互接続との間での通信をサポートするステップと、

前記集積回路内の外部通信アダプタを利用して、入出力(I/O)通信リンクを介して入出力通信をサポートするステップであって、集積回路内の通信だけを利用して、前記プロセッサ・コアから前記外部通信アダプタにI/O通信コマンドを伝送することを含むステップと

を含む方法。

( 17 ) 前記プロセッサ・コアがアクセスする可能性が高いデータおよび命令を前記集積回路内のキャッシュ階層構造にキャッシュするステップと、

40

前記集積回路内の集積化メモリ制御装置を利用して外部メモリにアクセスするステップと

をさらに含む、上記( 16 )に記載の方法。

( 18 ) 前記集積回路が、第1および第2の外部通信アダプタの少なくとも一部分を含み、前記I/O通信をサポートするステップが、

前記集積回路内の前記第1の外部通信アダプタを利用して、第1のリンク・レイヤ・プロトコルを使用するI/O通信をサポートするステップと、

前記集積回路内の前記第2の外部通信アダプタを利用して、異なる第2のリンク・レイヤ・プロトコルを使用するI/O通信をサポートするステップと

50

を含む、上記(16)に記載の方法。

(19) 前記通信要求で指定されるアドレスに基づいて、前記集積回路内の通信コマンドを前記相互接続インターフェースおよび前記外部通信アダプタに経路指定するステップをさらに含む、上記(16)に記載の方法。

(20) I/O通信をサポートする前記ステップが、前記I/Oコマンド内で指定される実効アドレスを、メモリ位置を識別するために利用される実アドレスに変換するステップをさらに含む、上記(16)に記載の方法。

(21) I/O通信をサポートする前記ステップが、前記外部通信アダプタ内の入出力(I/O)メモリ制御装置を利用して、I/Oデータにアクセスするためにメモリにアクセスするステップを含む、上記(16)に記載の方法。

(22) I/O通信をサポートする前記ステップが、I/Oデータ転送の完了を示すソフトウェア・アクセス可能な標識をセットするステップを含む、上記(16)に記載の方法。

(23) 前記I/O通信をサポートするステップが、着信I/Oデータと発信I/Oデータの少なくとも一方を前記集積回路内にバッファするステップを含む、上記(16)に記載の方法。

【図面の簡単な説明】

【0087】

【図1】従来技術による対称型マルチプロセッサ(SMP)データ処理システムを示す図である。

【図2】本発明を有利に利用することができるネットワーク・システムの一例を示す図である。

【図3】本発明によるマルチプロセッサ(MP)データ処理システムの実施形態の一例を示すブロック図である。

【図4】図3のデータ処理システム内の処理装置のより詳細なブロック図である。

【図5】本発明の好ましい一実施形態による、図3に示されたMPデータ処理システム内のシステム・メモリのI/Oデータ構造およびその他の内容を示すブロック図である。

【図6】図3のMPデータ処理システム内で実行されるソフトウェアの一例を示すレイヤ図である。

【図7】本発明によるI/O通信の方法の一例を示す高レベル論理フローチャートである。

【図8】本発明の好ましい一実施形態によるプロセッサ・コアのブロック図である。

【図9】本発明の好ましい一実施形態によるバイパスCAMのより詳細な図である。

【図10】本発明による反復的コード・シーケンスの実行をバイパスする方法の一例を示す高レベルの論理フローチャートである。

【符号の説明】

【0088】

8 対称型マルチプロセッサ(SMP)データ処理システム

10 処理装置

11 SMPシステム・バス

12 外部システム・メモリ

14 プロセッサ・コア

16 キャッシュ階層構造

18 メモリ制御装置

20 専用メモリ・バス

22 コマンド制御ブロック(CCB)

24 変換制御エントリ(TCE)テーブル

30 メザニンI/Oバス制御装置

32 メザニン・バス制御装置

40 メザニン・バス

10

20

30

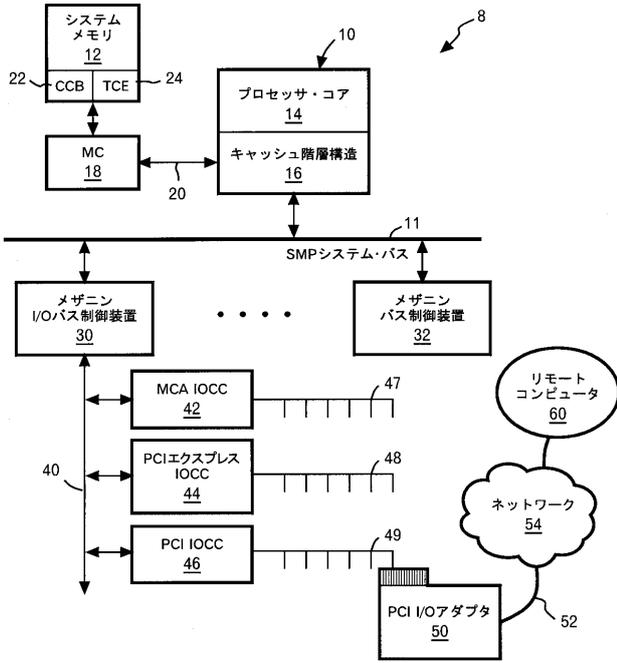
40

50

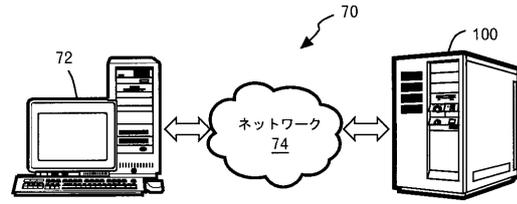
4 2	マイクロチャネル・アーキテクチャ I O C C	
4 4	P C I エクスプレス ( 3 G I O ) I O C C	
4 6	P C I I O C C	
4 7 ~ 4 9	バス	
5 0	周辺装置相互接続 ( P C I ) I / O アダプタ	
5 2	イーサネット ( R ) 通信リンク、I / O 通信リンク	
5 4	ネットワーク	
6 0	リモート・コンピュータ	
7 0	ネットワーク・システム	
7 2	ワークステーション・コンピュータ・システム	10
7 4	ネットワーク	
1 0 0	サーバ・コンピュータ・システム	
1 0 2	処理装置	
1 0 4	メモリ	
1 0 6	スイッチング機構	
1 0 8	プロセッサ・コア	
1 1 0	キャッシュ階層構造	
1 1 2	集積化メモリ制御装置 ( I M C )	
1 1 4	集積化機構インターフェース ( I F I )	
1 1 8	メモリ・バス	20
1 2 2	メモリ・マップ ( M M )	
1 2 4	変換索引バッファ ( T L B )	
1 3 0	外部通信アダプタ ( E C A )	
1 3 1	I / O メモリ制御装置 ( I / O M C )	
1 3 2	バッファ記憶装置	
1 3 3	データ転送ロジック ( D T L )	
1 3 4	プロトコル・ロジック	
1 3 5	データ・キュー	
1 3 6	エントリ	
1 3 8	リンク・レイヤ制御装置 ( L L C )	30
1 4 0	並直列変換器 / 直並列変換器 ( S E R / D E S )	
1 5 0	I / O 通信リンク	
1 5 8	ソフトウェア、ソフトウェア構成	
1 6 0	システム・スーパーバイザ ( ハイパーバイザ )	
1 6 2	オペレーティング・システム	
1 6 3	ミドルウェア	
1 6 4	ユーザ・レベル・プロセス	
1 6 4	アプリケーション・プログラム	
2 4 9	プロセッサ領域	
2 5 0	I / O 領域	40
2 5 2	共用領域	
2 5 3	データ転送制御ブロック ( D T C B )	
2 5 4	仮想キュー	
2 5 5	I / O データ・バッファ	
2 5 6	制御状態バッファ	
2 6 0	命令トレース・ログ	
2 6 2	I / O データ	
2 6 4	ページ・テーブル	
2 7 0	命令シーケンス・ユニット ( I S U )	
2 7 1	命令アドレス・レジスタ ( I A R )	50

272	命令メモリ管理ユニット ( I M M U )	
273	モード・セクタ	
274	L1 I - キャッシュ	
276	I - キャッシュ再ロード・バス	
280	命令バス	
281	トレース・バス	
282	条件レジスタ・ユニット ( C R U )	
284	分岐実行ユニット ( B E U )	
286	固定小数点ユニット ( F X U )	
288	ロード・ストア・ユニット ( L S U )	10
290	浮動小数点ユニット ( F P U )	
300	制御レジスタ・ファイル ( C R F )	
302	汎用レジスタ・ファイル ( G P R F )	
304	浮動小数点レジスタ・ファイル ( F P R F )	
308	L1 D - キャッシュ	
320	命令バイパス回路	
322	取込みロジック	
324	連想記憶装置 ( C A M )	
340	命令ストリーム・バッファ	
341	スヌープ・キル・フィールド	20
342	命令アドレス・フィールド	
343	ユーザ・レベル設計状態 C A M	
344	使用フラグ	
345	レジスタ値フィールド	
346	メモリ・マップ・アクセス C A M	
347	スヌープ・キル・フィールド	
348	目標アドレス・フィールド	
349	ロード/ストア ( L / S ) フィールド	
350	I / O フィールド	
352	データ・フィールド	30

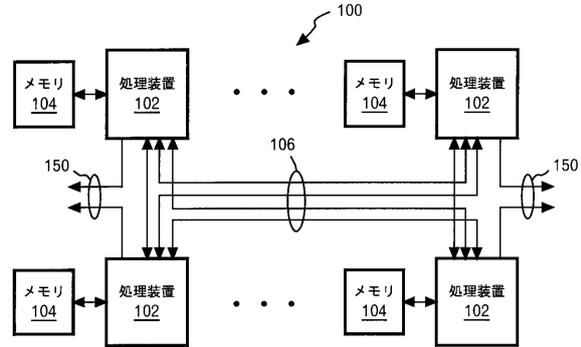
【 図 1 】



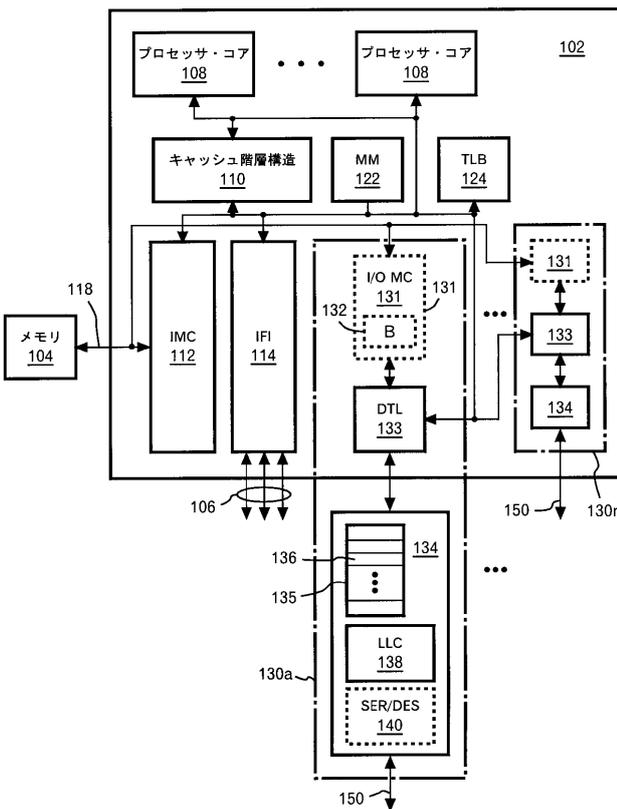
【 図 2 】



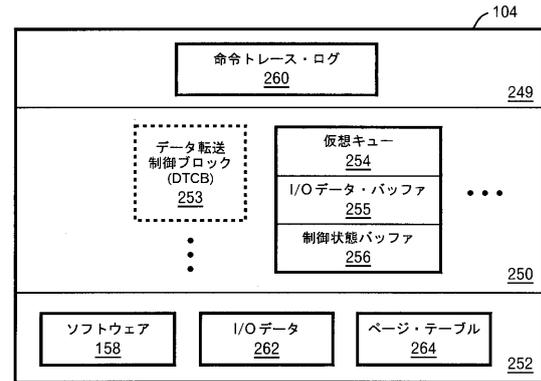
【 図 3 】



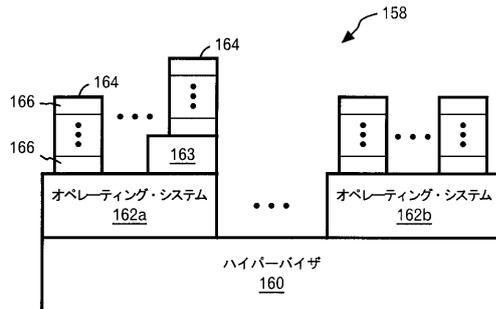
【 図 4 】



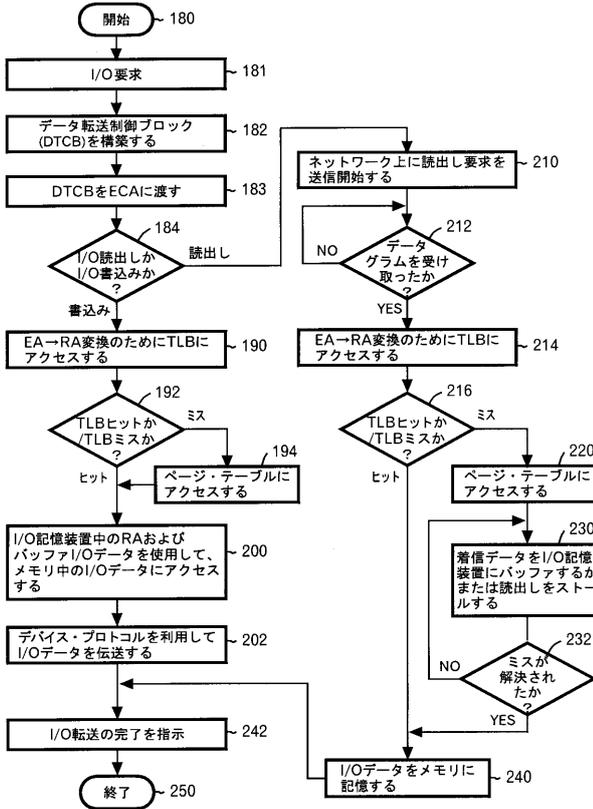
【 図 5 】



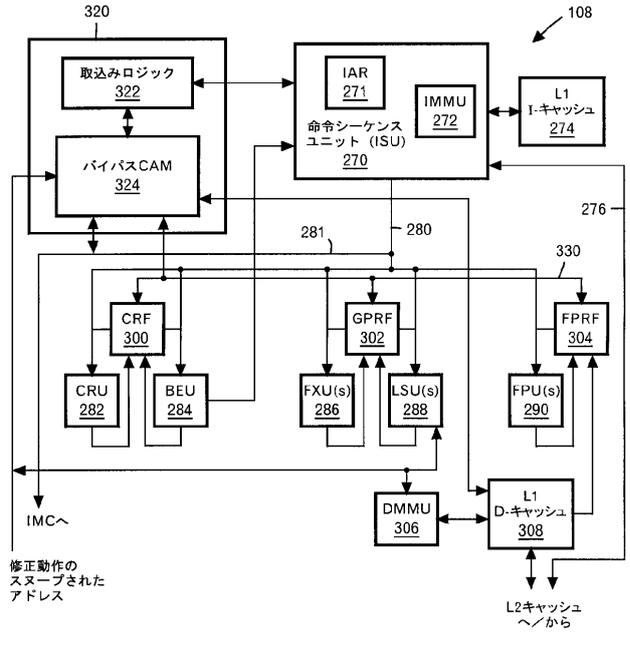
【 図 6 】



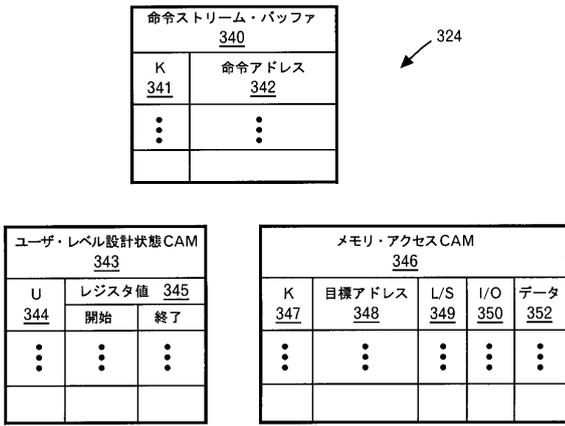
【 図 7 】



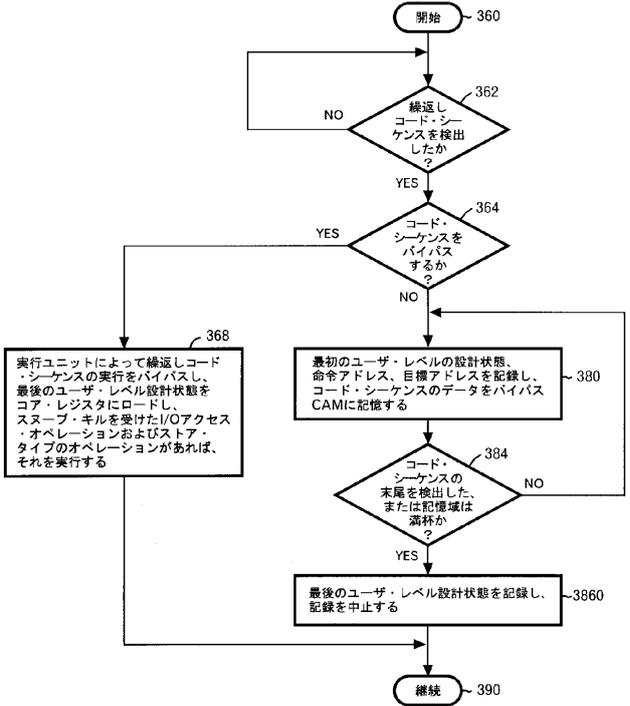
【 図 8 】



【 図 9 】



【 図 10 】



---

 フロントページの続き

(51) Int.Cl. <sup>7</sup>	F I	テーマコード(参考)
	G 0 6 F 12/10	5 0 1 Z
	G 0 6 F 12/10	5 5 5
	G 0 6 F 12/10	5 5 7
(72)発明者	ラビ・クマール・アーミミリ	
	アメリカ合衆国 7 8 7 5 9 テキサス州オースティン	スパイスブラッシュ・ドライブ 9 2 2 1
(72)発明者	ロバート・アラン・カーグノニ	
	アメリカ合衆国 7 8 7 4 6 テキサス州オースティン	キー・ウェスト・コーブ 1 9 0 4
(72)発明者	ガイ・リン・ガスリー	
	アメリカ合衆国 7 8 7 2 6 テキサス州オースティン	カラバール・ドライブ 1 1 1 4 5
(72)発明者	ウィリアム・ジョン・スターク	
	アメリカ合衆国 7 8 6 8 1 テキサス州ラウンド・ロック	シー・アッシュ・サークル 8 6 1 2
Fターム(参考)	5B005 JJ12 KK16 MM01 MM51 NN12 PP03 PP05 PP21 RR01 UU15	
	5B014 GC31 HB26	
	5B060 AB26 AC20 KA02 KA06	