

(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2006-178659

(P2006-178659A)

(43) 公開日 平成18年7月6日(2006.7.6)

(51) Int. Cl.

G06F 11/20 (2006.01)

F I

G06F 11/20 310E

テーマコード(参考)

5B034

審査請求 有 請求項の数 19 O L (全 21 頁)

(21) 出願番号 特願2004-369875 (P2004-369875)

(22) 出願日 平成16年12月21日(2004.12.21)

(特許庁注: 以下のものは登録商標)

1. リナックス

(出願人による申告) 平成15年度 新エネルギー・産業技術総合開発機構 15年度新エネ電情第0619006号 半導体アプリケーションチッププロジェクト(高機能・高信頼性サーバー用半導体チップ) 次世代高可用性サーバに係わる半導体チップ及び関連ソフトウェア技術の研究開発 産業活力再生特別措置法第30条の適用を受ける特許出願

(71) 出願人 000004237

日本電気株式会社
東京都港区芝五丁目7番1号

(74) 代理人 100102864

弁理士 工藤 実

(72) 発明者 阿部 晋樹

東京都港区芝五丁目7番1号 日本電気株式会社内

Fターム(参考) 5B034 AA01 BB01 CC01 CC05 CC06
DD05 DD06

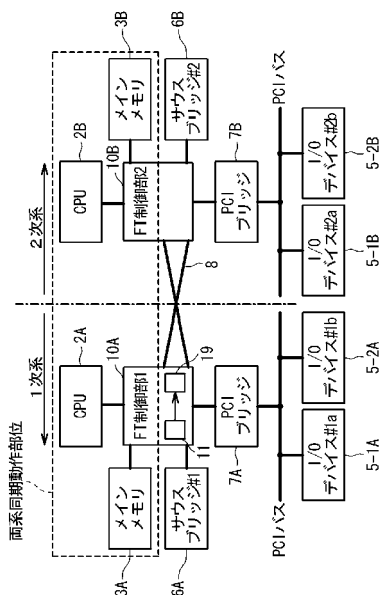
(54) 【発明の名称】 フォールト・トレラント・コンピュータシステムと、そのための割り込み制御方法

(57) 【要約】

【課題】 2つのシステムが同期して動作可能なフォールト・トレラント・コンピュータシステムを提供する。

【解決手段】 フォールト・トレラント・コンピュータシステムは、1次系システムと2次系システムとを具備する。前記1次系システムは、第1CPU(2A)と、前記第1CPUに接続された第1FT制御部(10A)と、及び前記第1FT制御部に電気的かつ動作的に接続された第1サウス・ブリッジ(6A)とを具備する。前記2次系システムは、第2CPU(2B)と、前記第2CPUに接続された第2FT制御部(10B)と、及び前記第2FT制御部に電気的に接続され、かつ動作的に接続されていない第2サウス・ブリッジ(6B)とを具備する。前記第1FT制御部と前記第2FT制御部はリンク(8)により接続され、前記リンクを用いて前記1次系システムと前記2次系システムとは、前記第2サウス・ブリッジを除き同期して動作する。

【選択図】 図2



【特許請求の範囲】

【請求項 1】

1次系システムと2次系システムとを具備し、

前記1次系システムは、第1CPUと、前記第1CPUに接続された第1FT制御部と、及び前記第1FT制御部に電気的かつ動作的に接続された第1サウス・ブリッジとを具備し、

前記2次系システムは、第2CPUと、前記第2CPUに接続された第2FT制御部と、及び前記第2FT制御部に電気的に接続され、かつ動作的に接続されていない第2サウス・ブリッジとを具備し、

前記第1FT制御部と前記第2FT制御部はリンクにより接続され、

前記リンクを用いて前記1次系システムと前記2次系システムとは、前記第2サウス・ブリッジを除き同期して動作する

フォールト・トレラント・コンピュータシステム。

10

【請求項 2】

請求項1に記載のフォールト・トレラント・コンピュータシステムにおいて、

前記第1CPUと第2CPUは同じオペレーティング・システム上で動作し、

前記第2サウス・ブリッジは前記オペレーティング・システムから不可視である
フォールト・トレラント・コンピュータシステム。

【請求項 3】

請求項1又は2に記載のフォールト・トレラント・コンピュータシステムにおいて、

前記第1FT制御部は、第1マスターIOAPIC制御部を有し、

前記第2FT制御部は、第2マスターIOAPIC制御部を有し、

前記第1サウス・ブリッジは、PIC制御部とIOAPIC制御部を有し、

前記第1サウス・ブリッジの前記IOAPIC制御部のアドレス空間は、前記第1マスターIOAPIC制御部のアドレス空間の一部と同じであり、前記第2サウス・ブリッジの前記IOAPIC制御部のアドレス空間は、前記第2マスターIOAPIC制御部のアドレス空間の一部と同じである
フォールト・トレラント・コンピュータシステム。

20

【請求項 4】

請求項1乃至3のいずれかに記載のフォールト・トレラント・コンピュータシステムにおいて、

前記第1と第2のFT制御部は、前記1次系と2次系のシステムの設定データとステータスを示すステータスデータを格納する第1と第2のステータス格納部をそれぞれ更に有し、

前記1次系システムがアクティブ・システムとして動作し、前記2次系システムがスタンバイ・システムとして前記1次系システムと同期的に動作している間に、前記第1サウス・ブリッジに障害が発生したとき、前記第1FT制御部の前記第1ステータス格納部に格納された前記設定データ・ステータスデータは前記第2FT制御部の前記第2ステータス格納部に転送され、

その後、前記2次系システムがアクティブ・システムとして動作する

フォールト・トレラント・コンピュータシステム。

40

【請求項 5】

請求項3に記載のフォールト・トレラント・コンピュータシステムにおいて、

前記1次系システムは、前記第1FT制御部に接続された第1PCIブリッジと、前記第1PCIブリッジに接続された第1I/Oデバイス群とを更に備え、

前記2次系システムは、前記第2FT制御部に接続された第2PCIブリッジと、前記第2PCIブリッジに接続された第2I/Oデバイス群とを更に備え、

前記1次系システムがアクティブ・システムとして設定され、前記2次系システムがスタンバイ・システムとして前記1次系システムと同期的に動作するように設定されているとき、オペレーティング・システムが動作していない起動時のレガシーモードにおいて、

50

前記第 1 I / O デバイス群のうちの 1 つからの第 1 割り込み要求は、前記第 1 マスター I O A P I C 制御部と前記第 1 サウス・ブリッジの P I C 制御部を介して前記第 1 C P U に転送される

フォールト・トレラント・コンピュータシステム。

【請求項 6】

請求項 5 に記載のフォールト・トレラント・コンピュータシステムにおいて、

前記第 1 割り込み要求は、前記リンクを介して前記第 2 F T 制御部の前記第 2 マスター I O A P I C 制御部を介して前記第 2 C P U に転送され、また、予め定められた遅延時間をもって前記第 1 F T 制御部の前記第 1 マスター I O A P I C 制御部に転送される

フォールト・トレラント・コンピュータシステム。

10

【請求項 7】

請求項 5 または 6 に記載のフォールト・トレラント・コンピュータシステムにおいて、

前記レガシーモード後の、前記オペレーティング・システムが動作している拡張モードにおいて、前記第 1 サウス・ブリッジからの第 2 割り込み要求は、前記第 1 I O A P I C 制御部から前記第 1 C P U に転送され、また前記リンク、前記第 2 F T 制御部の前記第 2 I O A P I C 制御部から前記第 2 C P U に転送される

フォールト・トレラント・コンピュータシステム。

【請求項 8】

請求項 5 乃至 7 のいずれかに記載のフォールト・トレラント・コンピュータシステムにおいて、

20

前記拡張モードにおいて、前記第 1 I / O デバイス群のうちの 1 つからの第 2 割り込み要求は、前記第 1 マスター I O A P I C 制御部を介して前記第 1 C P U に転送され、また前記リンクと前記第 2 マスター I O A P I C 制御部を介して前記第 2 C P U に転送される

フォールト・トレラント・コンピュータシステム。

【請求項 9】

請求項 5 乃至 8 のいずれかに記載のフォールト・トレラント・コンピュータシステムにおいて、

前記拡張モードにおいて、前記第 2 I / O デバイス群のうちの 1 つからの第 3 割り込み要求は、前記第 2 マスター I O A P I C 制御部を介して前記第 2 C P U に転送され、また前記リンクと前記第 1 マスター I O A P I C 制御部を介して前記第 1 C P U に転送される

30

フォールト・トレラント・コンピュータシステム。

【請求項 10】

請求項 5 乃至 9 に記載のフォールト・トレラント・コンピュータシステムにおいて、

前記第 1 と第 2 の F T 制御部は、前記 1 次系と 2 次系の設定データとシステムのステータスを示すステータスデータを格納する第 1 と第 2 のステータス格納部をそれぞれ更に有し、

前記第 1 I O A P I C 制御部は、前記第 1 割り込み要求以外の受信した前記割り込み要求とそのときの前記システムのステータスデータを前記第 1 ステータス格納部に格納し、

前記第 2 I O A P I C 制御部は、前記第 1 割り込み要求以外の受信した前記割り込み要求とそのときの前記システムのステータスデータを前記第 2 ステータス格納部に格納する

40

フォールト・トレラント・コンピュータシステム。

【請求項 11】

請求項 1 乃至 3 と 5 乃至 9 のいずれかに記載のフォールト・トレラント・コンピュータシステムにおいて、

前記第 1 と第 2 の F T 制御部は、前記 1 次系と 2 次系のシステムの設定データとステータスを示すステータスデータを格納する第 1 と第 2 のステータス格納部をそれぞれ更に有し、

前記第 1 C P U と前記第 2 C P U が同期して動作している間に前記第 1 と第 2 のステータス格納部の格納データの不一致が検出されたとき、S M I (システム・マネージメント・インターラプト) ハンドラーは、前記第 1 C P U と前記第 2 C P U の動作を停止させ、

50

前記第 1 と第 2 のステータス格納部の前記設定データ / ステータスデータが一致するように、前記第 1 と第 2 の F T 制御部を制御する
フォールト・トレラント・コンピュータシステム。

【請求項 1 2】

1 次系システムと 2 次系システムとを具備し、

前記 1 次系システムは、第 1 C P U と、前記第 1 C P U に接続された第 1 F T 制御部と、及び前記第 1 F T 制御部に電気的かつ動作的に接続された第 1 サウス・ブリッジと、前記第 1 F T 制御部に接続された第 1 P C I ブリッジと、前記第 1 P C I ブリッジに接続された第 1 I / O デバイス群とを具備し、

前記 2 次系システムは、第 2 C P U と、前記第 2 C P U に接続された第 2 F T 制御部と、及び前記第 2 F T 制御部に電気的に接続され、かつ動作的に接続されていない第 2 サウス・ブリッジと、前記第 2 F T 制御部に接続された第 2 P C I ブリッジと、前記第 2 P C I ブリッジに接続された第 2 I / O デバイス群とを具備し、

前記第 1 F T 制御部と前記第 2 F T 制御部はリンクにより接続され、

前記第 1 F T 制御部は、第 1 マスター I O A P I C 制御部を有し、前記第 2 F T 制御部は、第 2 マスター I O A P I C 制御部を有し、

前記第 1 C P U と前記第 2 C P U の各々は、第 1 割り込みパスと第 2 割り込みパスを有し、

起動時に前記 1 次系システム内で生成された第 1 割り込みは、前記第 1 と第 2 のマスター I O A P I C と前記第 1 割り込みパスを経由して前記第 1 と第 2 の C P U へ転送され、動作時に前記 1 次系システムまたは前記第 2 次系システム内で生成された第 2 割り込みは、前記第 1 と第 2 のマスター I O A P I C と前記第 2 割り込みパスを経由して前記第 1 と第 2 の C P U へ転送される

フォールト・トレラント・コンピュータシステム。

【請求項 1 3】

1 次系システムと 2 次系システムとを具備し、前記 1 次系システムは、第 1 C P U と、前記第 1 C P U に接続された第 1 F T 制御部と、及び前記第 1 F T 制御部に電気的かつ動作的に接続された第 1 サウス・ブリッジと、前記第 1 F T 制御部に接続された第 1 P C I ブリッジと、前記第 1 P C I ブリッジに接続された第 1 I / O デバイス群とを具備し、前記 2 次系システムは、第 2 C P U と、前記第 2 C P U に接続された第 2 F T 制御部と、及び前記第 2 F T 制御部に電気的に接続され、かつ動作的に接続されていない第 2 サウス・ブリッジと、前記第 2 F T 制御部に接続された第 2 P C I ブリッジと、前記第 2 P C I ブリッジに接続された第 2 I / O デバイス群とを具備し、前記第 1 F T 制御部は、第 1 マスター I O A P I C 制御部を有し、前記第 2 F T 制御部は、第 2 マスター I O A P I C 制御部を有するフォールト・トレラント・コンピュータシステムにおいて、

前記 1 次系システムがアクティブ・システムとして設定され、前記 2 次系システムがスタンバイ・システムとして前記 1 次系システムと同期的に動作するように設定されているとき、オペレーティング・システムが動作していない起動時のレガシーモードにおいて、第 1 I / O デバイス群のうちの 1 つからの第 1 割り込み要求を、前記第 1 マスター I O A P I C 制御部に転送するステップと、

前記第 1 割り込み要求を前記リンクを介して前記第 2 F T 制御部の前記第 2 マスター I O A P I C 制御部に転送するステップと、

前記第 1 割り込み要求を、前記第 1 マスター I O A P I C 制御部から前記第 1 サウス・ブリッジの P I C 制御部を介して、再び前記第 1 マスター I O A P I C 制御部から前記第 1 C P U に転送するステップと、

前記第 1 割り込み要求を、前記第 2 マスター I O A P I C 制御部から前記第 2 サウス・ブリッジの P I C 制御部を介して、再び前記第 2 マスター I O A P I C 制御部から前記第 2 C P U に転送するステップとを具備する割り込み制御方法。

【請求項 1 4】

10

20

30

40

50

請求項 1 3 に記載の割り込み制御方法において、

前記第 1 割り込み要求が前記第 1 マスター I O A P I C 制御部に届くタイミングと前記第 2 マスター I O A P I C 制御部に届くタイミングは同一である
割り込み制御方法。

【請求項 1 5】

請求項 1 3 又は 1 4 に記載の割り込み制御方法において、

前記レガシーモード後の、前記オペレーティング・システムが動作している拡張モードにおいて、前記第 1 サウス・ブリッジからの第 2 割り込み要求を、前記第 1 I O A P I C 制御部から前記第 1 C P U に転送するステップと、

前記第 2 割り込み要求を前記リンク、前記第 2 F T 制御部の前記第 2 I O A P I C 制御部から前記第 2 C P U に転送するステップと
を更に具備する割り込み制御方法。 10

【請求項 1 6】

請求項 1 5 に記載の割り込み制御方法において、

前記拡張モードにおいて、前記第 2 割り込み要求を、前記第 1 マスター I O A P I C 制御部を介して前記第 1 C P U に転送するステップと、

前記第 2 割り込み要求を、前記リンクと前記第 2 マスター I O A P I C 制御部を介して前記第 2 C P U に転送するステップと
を更に具備する割り込み制御方法。

【請求項 1 7】

請求項 1 3 乃至 1 6 のいずれかに記載の割り込み制御方法において、

前記拡張モードにおいて、前記第 2 I / O デバイス群のうちの 1 つからの第 3 割り込み要求を、前記第 2 マスター I O A P I C 制御部を介して前記第 2 C P U に転送するステップと、

前記第 3 割り込み要求を前記リンクと前記第 1 マスター I O A P I C 制御部を介して前記第 1 C P U に転送するステップと
を更に具備する割り込み制御方法。 20

【請求項 1 8】

請求項 1 3 乃至 1 7 に記載の割り込み制御方法において、

前記第 1 割り込み要求以外の受信された前記割り込み要求とそのときの前記システムのステータスデータを前記第 1 F T 制御部内の第 1 ステータス格納部に格納するステップと 30

、
前記第 1 割り込み要求以外の受信した前記割り込み要求とそのときの前記システムのステータスデータを前記第 2 F T 制御部内の第 2 ステータス格納部に格納するステップと
を更に具備する割り込み制御方法。

【請求項 1 9】

請求項 1 8 に記載の割り込み制御方法において、

前記第 1 C P U と前記第 2 C P U が同期して動作している間に前記第 1 と第 2 のステータス格納部の格納データの不一致が検出されたとき、S M I (システム・マネージメント・インターラプト) ハンドラーにより、前記第 1 C P U と前記第 2 C P U の動作を停止させるステップと、 40

前記第 1 と第 2 のステータス格納部の設定データと前記ステータスデータが一致するように、前記第 1 と第 2 の F T 制御部を制御するステップと、

前記第 1 と第 2 のステータス格納部の前記設定データと前期ステータスデータが一致したとき、前記第 1 C P U と前記第 2 C P U の動作を再開するステップと
を更に具備する割り込み制御方法。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本発明は、コントローラの二重化方式に関し、特に割り込み制御も二重化されているフ 50

ォールト・トレラント・コンピュータシステムに関する。

【背景技術】

【0002】

高度な信頼性を提供するコンピュータとして、フォールト・トレラント・コンピュータシステムが知られている。フォールト・トレラント・コンピュータでは、システムを構成する全てのハードウェア・モジュールが二重化され、または多重化されている。全てのハードウェア・モジュールは同期して動作し、たとえある部位で故障が発生したとしても、該ハードウェア・モジュールが切り離され、正常なハードウェア・モジュールで処理が継続される。これにより、耐故障性が向上している。

【0003】

図1は、フォールト・トレラント・コンピュータシステムの構成の一例を示している。この例のフォールト・トレラント・コンピュータシステムはフォールト・トレラント制御部（以下、FT制御部と呼ぶ）を有し、CPU、メモリ、I/Oデバイスといったハードウェア・モジュールが二重化されている。FT制御部は、ハードウェア・モジュールと接続され、同期動作処理、故障時の切り替え制御などを行う。

【0004】

図1に示されるフォールト・トレラント・コンピュータシステムでは、CPU（群）とメインメモリが一つのCPUサブシステムを構成し、これと全く同じ構成を有する他方のCPUサブシステムが存在する。これら2つのCPUサブシステムは二重化されている。同様に、同一構成のI/Oデバイス（群）も二重化され、I/Oサブシステムを構成する。FT制御部はそれらの中心に位置し、各々のモジュール（CPUサブシステム、I/Oデバイス群）を制御し、CPUサブシステムの両系の同期動作の維持、故障の検出と、故障モジュールの切り離し制御を行う。

【0005】

図1では2個のCPUサブシステムが存在するが、故障サブシステムはFT制御部により論理的に切り離され、1個のCPUサブシステムとI/Oサブシステムで処理は継続される。一般的に、フォールト・トレラント・コンピュータは、ハードウェア的に二重化制御される部分と、ソフトウェア的に二重化制御される部分とに分かれる。例えば、CPUサブシステムは、ソフトウェアが実行される基盤であり、これらはハードウェア的に二重化制御される必要がある。このためCPUサブシステム内でエラーが発生した場合、FT制御部が該当CPU又はメモリをシステムから切り離し、正常動作しているCPU及びメモリに影響を及ぼさないように制御を行う。

【0006】

一方、I/Oデバイスの故障の場合、それを検出したFT制御部は、I/Oデバイスを制御しているソフトウェア（以下I/Oデバイス・ドライバと呼ぶ）に対し、エラー通知を行うことで、I/Oデバイスの切り替えをソフトウェア的に行うことが可能である。この場合、I/Oデバイス・ドライバは、故障したI/Oデバイスの使用を中止し、代わって二重化されている別のI/Oデバイスを使用することになる。これらはI/Oサブシステム内での、使用I/Oデバイスの切り替えとなる。

【0007】

しかし、数あるI/Oデバイスのうちソフトウェア的に二重化されることができないものもある。例えば、割り込みコントローラもその一つである。割り込みコントローラは、各I/Oデバイスからの割り込み要求を受け付け、それをCPUへ知らせるのが主な役割である。I/Oデバイスからの割り込みは、オペレーティング・システム（OS）によって、IRQと呼ばれる割り込み番号に割り宛てられる。場合によっては、一つの割り込み番号に複数のI/Oデバイスが割り宛てられる場合もある。割り込みコントローラは、各デバイスからの割り込み要求を、設定された割り込み番号に変換してCPUへ通知する。このとき、CPUがある割り込み番号の割り込み処理を処理中の場合には、割り込みコントローラは、同じ番号の割り込み要求を知らせないか、または、複数のデバイスからの割り込みが失われないように管理する。従って、割り込みコントローラは、処理中の割り込みの状態の保持

10

20

30

40

50

などを内部で行っており、割り込みコントローラで障害が発生した場合、それら情報が全て失われることになる。このため、ソフトウェアにより割り込みコントローラを元の状態に復帰させることは不可能である。

【0008】

さらに、現在のウィンドウズ（Windows（登録商標））やリナックスといったOSは、複数の割り込みコントローラの存在を許しているが、動作中に割り込みコントローラが増減することには対応していない。起動時に存在した割り込みコントローラは、OSがシャットダウンされるまで存在し、正常動作しつづける必要がある。

【0009】

一般的に、現状のPCサーバなどではオープン化が進んでおり、安価なサーバを製造する場合、おのずとインテル（Intel（登録商標））系CPUや、一般市場に安価に出回っている部材が選択されることになる。また、現在PCサーバで主流であるウィンドウズやリナックスなどのOSも、これらインテル系アーキテクチャーに立脚している。しかしながら、オープン系PCサーバにおいて、安価にフォールト・トレラント・コンピュータを構成しようとする場合、以下のような多くの問題が存在する。

【0010】

例えば、PCサーバで採用されるほとんどのI/Oデバイスや、ウィンドウズといったOSは、フォールト・トレラント・コンピュータシステムを意識して設計されておらず、デバイスが二重化されても、故障時のフェイル・オーバー処理には全く対応していない。インテル系PCサーバは、サウス・ブリッジと呼ばれるレガシー（Legacy）機能が集約された特殊なI/Oデバイスに割り込み制御が負わせている。特に、割り込み制御はシステム動作の中核であるので、OSが直接サウス・ブリッジにアクセスを行い、その動作の制御を行っている。このため、一度サウス・ブリッジに障害が発生すると、OSは機能不全を起し、結果的にシステム・ダウンを引き起こすことになる。また、オープン系で主に使用されるウィンドウズといったOSに、フォールト・トレラント・コンピュータシステム用の改造を加えることは、現実的に不可能である。

上記と関連して、特開平9-251443号公報には、情報処理システムのプロセッサ障害回復処理方法が開示されている。この従来例では、情報処理システムは、複数のプロセッサを備え、少なくとも1個のプロセッサをシステム支援プロセッサとして動作させ、その他のプロセッサを命令プロセッサとして動作させるマルチプロセッサ構成の計算機システムである。プロセッサに固定障害が発生したとき、前記システム支援プロセッサの障害発生時、少なくとも1個の命令プロセッサ上で動作しているオペレーティングシステムに割り込みを発生し、前記オペレーティングシステムが、前記命令プロセッサにおいて障害が発生したことを認識し、前記命令プロセッサ上で前記割り込み発生時に動作していたアプリケーションプログラムを異常終了させ、前記命令プロセッサをシステム支援プロセッサと交代させている。

【特許文献1】特開平9-251443号公報

【発明の開示】

【発明が解決しようとする課題】

【0011】

本発明の課題は、多重化、例えば二重化された2つのシステムが同期して動作可能なフォールト・トレラント・コンピュータシステムを提供することである。

本発明の他の課題は、システムの切替え時に割り込み要求が保持されることが出来るフォールト・トレラント・コンピュータシステムを提供することにある。

本発明の他の課題は、CPUからはサウス・ブリッジの故障を隠蔽することが出来るフォールト・トレラント・コンピュータシステムを提供することにある。

本発明の他の課題は、故障したFT制御部が交換されたときでも、完全に同期状態に復帰させることが出来るフォールト・トレラント・コンピュータシステムを提供することにある。

本発明の他の課題は、フォールト・トレラント・コンピュータシステムを意識せずに作

10

20

30

40

50

成された既存のOS、既存のサウス・ブリッジを搭載したコンピュータシステム（サーバ）においても、割込みコントローラの二重化を実現することが可能となるフォールト・トレラント・コンピュータシステムを提供することにある。

【課題を解決するための手段】

【0012】

以下に、[発明の実施の形態]で使用する番号・符号を用いて、課題を解決するための手段を説明する。これらの番号・符号は、[特許請求の範囲]の記載と発明の実施の形態の記載との対応関係を明らかにするために付加されたものであるが、[特許請求の範囲]に記載されている発明の技術的範囲の解釈に用いてはならない。

【0013】

本発明の観点では、フォールト・トレラント・コンピュータシステムは、1次系システムと2次系システムとを具備する。前記1次系システムは、第1CPU(2A)と、前記第1CPUに接続された第1FT制御部(10A)と、及び前記第1FT制御部に電気的かつ動作的に接続された第1サウス・ブリッジ(6A)とを具備する。前記2次系システムは、第2CPU(2B)と、前記第2CPUに接続された第2FT制御部(10B)と、及び前記第2FT制御部に電気的に接続され、かつ動作的に接続されていない第2サウス・ブリッジ(6B)とを具備する。前記第1FT制御部と前記第2FT制御部はリンク(8)により接続され、前記リンクを用いて前記1次系システムと前記2次系システムとは、前記第2サウス・ブリッジを除き同期して動作する。これにより、第2サウス・ブリッジはコンピュータシステムに影響を与えない。

10

20

【0014】

このとき、前記第1CPUと第2CPUは同じオペレーティング・システム上で動作し、前記第2サウス・ブリッジは前記オペレーティング・システムから不可視である。これにより、第2サウス・ブリッジは、前記オペレーティング・システムの影響を受けない。

【0015】

また、前記第1FT制御部は、第1マスターIOAPIC制御部(12A)を有し、前記第2FT制御部は、第2マスターIOAPIC制御部(12B)を有し、前記第1サウス・ブリッジは、PIC制御部(34)とIOAPIC制御部(36)を有してもよい。前記第1サウス・ブリッジの前記IOAPIC制御部のアドレス空間は、前記第1マスターIOAPIC制御部のアドレス空間の一部と同じであり、前記第2サウス・ブリッジの前記IOAPIC制御部のアドレス空間は、前記第2マスターIOAPIC制御部のアドレス空間の一部と同じである。この結果、各マスターIOAPICにデータが設定されれば、そのデータは前記各サウス・ブリッジの前記IOAPIC制御部に反映されることができる。

30

【0016】

また、前記第1と第2のFT制御部は、前記1次系と2次系のシステムのステータスを示すステータスデータを格納する第1と第2のステータス格納部(22)をそれぞれ更に有してもよい。この場合、前記1次系システムがアクティブ・システムとして動作し、前記2次系システムがスタンバイ・システムとして前記1次系システムと同期的に動作している間に、前記第1サウス・ブリッジに障害が発生したとき、前記第1FT制御部の前記第1ステータス格納部に格納されたステータスデータは前記第2FT制御部の前記第2ステータス格納部に転送され、その後、前記2次系システムがアクティブ・システムとして動作する。これにより、CPU間で同期がはずれたとき、同期を再確立することができ、また両システムの一方に障害が発生したとき、障害発生部を交換したのち、他方からデータを転送することにより、再び同期動作を行うことが可能となる。

40

【0017】

また、前記1次系システムは、前記第1FT制御部に接続された第1PCIブリッジ(7A)と、前記第1PCIブリッジに接続された第1I/Oデバイス群(5A)とを更に備え、前記2次系システムは、前記第2FT制御部に接続された第2PCIブリッジ(7B)と、前記第2PCIブリッジに接続された第2I/Oデバイス群(5B)とを更に備

50

えていてもよい。この場合、前記1次系システムがアクティブ・システムとして設定され、前記2次系システムがスタンバイ・システムとして前記1次系システムと同期的に動作するように設定されているとき、オペレーティング・システムが動作していない起動時のレガシーモードにおいて、前記第1 I/Oデバイス群のうちの1つからの第1割り込み要求は、前記第1マスター I/O A P I C制御部と前記第1サウス・ブリッジの P I C制御部を介して前記第1 C P Uに転送される。こうして、両システムで割り込み処理を実行することができる。

【0018】

また、前記第1割り込み要求は、前記リンクを介して前記第2 F T制御部の前記第2マスター I/O A P I C制御部を介して前記第2 C P Uに転送され、また、予め定められた遅延時間をもって前記第1 F T制御部の前記第1マスター I/O A P I C制御部に転送される。これにより、両システムで割り込み処理を同期して実行することができる。

10

また、前記レガシーモード後の、前記オペレーティング・システムが動作している拡張モードにおいて、前記第1サウス・ブリッジからの第2割り込み要求は、前記第1 I/O A P I C制御部から前記第1 C P Uに転送され、また前記リンク、前記第2 F T制御部の前記第2 I/O A P I C制御部から前記第2 C P Uに転送される。

【0019】

また、前記拡張モードにおいて、前記第1 I/Oデバイス群のうちの1つからの第2割り込み要求は、前記第1マスター I/O A P I C制御部を介して前記第1 C P Uに転送され、また前記リンクと前記第2マスター I/O A P I C制御部を介して前記第2 C P Uに転送される。

20

また、前記拡張モードにおいて、前記第2 I/Oデバイス群のうちの1つからの第3割り込み要求は、前記第2マスター I/O A P I C制御部を介して前記第2 C P Uに転送され、また前記リンクと前記第1マスター I/O A P I C制御部を介して前記第1 C P Uに転送される。

以上により、レガシーモードにおいても、拡張モードにおいて、割り込みが同期して処理されることができる。

【0020】

また、前記第1と第2の F T制御部は、前記1次系と2次系のシステムのステータスを示すステータスデータを格納する第1と第2のステータス格納部(22)をそれぞれ更に有してもよい。前記第1 I/O A P I C制御部は、前記第1割り込み要求以外の受信した前記割り込み要求とそのときの前記システムのステータスを前記第1ステータス格納部に格納し、前記第2 I/O A P I C制御部は、前記第1割り込み要求以外の受信した前記割り込み要求とそのときの前記システムのステータスを前記第2ステータス格納部に格納する。これにより、両システムは、同一のステータスを保持することができる。

30

【0021】

また、前記第1と第2の F T制御部は、前記1次系と2次系のシステムのステータスを示すステータスデータを格納する第1と第2のステータス格納部(22)をそれぞれ更に有してもよい。前記第1 C P Uと前記第2 C P Uが同期して動作している間に前記第1と第2のステータス格納部の格納データの不一致が検出されたとき、S M I(システム・マネージメント・インターラプト)ハンドラーは、前記第1 C P Uと前記第2 C P Uの動作を停止させ、前記第1と第2のステータス格納部の格納データが一致するように、前記第1と第2の F T制御部を制御する。これにより、同期はずれが発生したときにも、あるいは故障が発生したときにも、割り込み処理が正しく継承されることができる。

40

【0022】

本発明の他の観点では、フォールト・トレラント・コンピュータシステムは、1次系システムと2次系システムとを具備する。前記1次系システムは、第1 C P U(2A)と、前記第1 C P Uに接続された第1 F T制御部(10A)と、及び前記第1 F T制御部に電気的かつ動作的に接続された第1サウス・ブリッジ(6A)と、前記第1 F T制御部に接続された第1 P C Iブリッジ(7A)と、前記第1 P C Iブリッジに接続された第1 I /

50

デバイス群(5A)とを具備する。前記2次系システムは、第2CPU(2B)と、前記第2CPUに接続された第2FT制御部(10B)と、及び前記第2FT制御部に電氣的に接続され、かつ動作的に接続されていない第2サウス・ブリッジ(6B)と、前記第2FT制御部に接続された第2PCIブリッジ(7B)と、前記第2PCIブリッジに接続された第2I/Oデバイス群(5B)とを具備する。前記第1FT制御部と前記第2FT制御部はリンク(8)により接続され、前記第1FT制御部は、第1マスターIOAPIC制御部(12A)を有し、前記第2FT制御部は、第2マスターIOAPIC制御部(12B)を有する。前記第1CPUと前記第2CPUの各々は、第1割り込みパスと第2割り込みパスを有する。起動時に前記1次系システム内で生成された第1割り込みは、前記第1と第2のマスターIOAPICと前記第1割り込みパスを経由して前記第1と第2のCPUへ転送され、動作時に前記1次系システムまたは前記第2次系システム内で生成された第2割り込みは、前記第1と第2のマスターIOAPICと前記第2割り込みパスを経由して前記第1と第2のCPUへ転送される。これにより、レガシーモードと拡張モードにおいて、割り込み要求の転送パスを変更することができる。

10

【0023】

また、本発明の他の観点では、割り込み制御方法が提供される。フォールト・トレラント・コンピュータシステムは、1次系システムと2次系システムとを具備し、前記1次系システムは、第1CPU(2A)と、前記第1CPUに接続された第1FT制御部(10A)と、及び前記第1FT制御部に電氣的かつ動作的に接続された第1サウス・ブリッジ(6A)と、前記第1FT制御部に接続された第1PCIブリッジ(7A)と、前記第1PCIブリッジに接続された第1I/Oデバイス群(5A)とを具備し、前記2次系システムは、第2CPU(2B)と、前記第2CPUに接続された第2FT制御部(10B)と、及び前記第2FT制御部に電氣的に接続され、かつ動作的に接続されていない第2サウス・ブリッジ(6B)と、前記第2FT制御部に接続された第2PCIブリッジ(7B)と、前記第2PCIブリッジに接続された第2I/Oデバイス群(5B)とを具備し、前記第1FT制御部は、第1マスターIOAPIC制御部(12A)を有し、前記第2FT制御部は、第2マスターIOAPIC制御部(12B)を有する。このフォールト・トレラント・コンピュータシステムにおいて、割り込み制御方法は、前記1次系システムがアクティブ・システムとして設定され、前記2次系システムがスタンバイ・システムとして前記1次系システムと同期的に動作するように設定されているとき、オペレーティング・システムが動作していない起動時のレガシーモードにおいて、第1I/Oデバイス群のうちの一つからの第1割り込み要求を、前記第1マスターIOAPIC制御部に転送することと、前記第1割り込み要求を前記リンクを介して前記第2FT制御部の前記第2マスターIOAPIC制御部に転送することと、前記第1割り込み要求を、前記第1マスターIOAPIC制御部から前記第1サウス・ブリッジのPIC制御部を介して、再び前記第1マスターIOAPIC制御部から前記第1CPUに転送することと、前記第1割り込み要求を、前記第2マスターIOAPIC制御部から前記第2サウス・ブリッジのPIC制御部を介して、再び前記第2マスターIOAPIC制御部から前記第2CPUに転送することにより達成される。

20

30

また、前記第1割り込み要求が前記第1マスターIOAPIC制御部に届くタイミングと前記第2マスターIOAPIC制御部に届くタイミングは同一であることが望ましい。

40

【0024】

また、割り込み制御方法は、前記レガシーモード後の、前記オペレーティング・システムが動作している拡張モードにおいて、前記第1サウス・ブリッジからの第2割り込み要求を、前記第1IOAPIC制御部から前記第1CPUに転送するステップと、前記第2割り込み要求を前記リンク、前記第2FT制御部の前記第2IOAPIC制御部から前記第2CPUに転送するステップとを更に具備してもよい。

【0025】

また、割り込み制御方法は、前記拡張モードにおいて、前記第2割り込み要求を、前記第1マスターIOAPIC制御部を介して前記第1CPUに転送するステップと、前記第

50

2 割り込み要求を、前記リンクと前記第 2 マスター I O A P I C 制御部を介して前記第 2 C P U に転送するステップとを更に具備してもよい。

また、割り込み制御方法は、前記拡張モードにおいて、前記第 2 I / O デバイス群のうちの 1 つからの第 3 割り込み要求を、前記第 2 マスター I O A P I C 制御部を介して前記第 2 C P U に転送するステップと、前記第 3 割り込み要求を前記リンクと前記第 1 マスター I O A P I C 制御部を介して前記第 1 C P U に転送するステップとを更に具備してもよい。

【 0 0 2 6 】

また、割り込み制御方法は、前記第 1 割り込み要求以外の受信された前記割り込み要求とそのときの前記システムのステータスを前記第 1 F T 制御部内の第 1 ステータス格納部に格納するステップと、前記第 1 割り込み要求以外の受信した前記割り込み要求とそのときの前記システムのステータスを前記第 2 F T 制御部内の第 2 ステータス格納部に格納するステップとを更に具備してもよい。

10

【 0 0 2 7 】

また、割り込み制御方法は、前記第 1 C P U と前記第 2 C P U が同期して動作している間に前記第 1 と第 2 のステータス格納部の格納データの不一致が検出されたとき、S M I (システム・マネージメント・インターラプト) ハンドラーにより、前記第 1 C P U と前記第 2 C P U の動作を停止させるステップと、前記第 1 と第 2 のステータス格納部の格納データが一致するように、前記第 1 と第 2 の F T 制御部を制御するステップと、前記第 1 と第 2 のステータス格納部の格納データが一致したとき、前記第 1 C P U と前記第 2 C P U の動作を再開するステップとを更に具備することが好ましい。

20

【 発明の効果 】

【 0 0 2 8 】

以上に示す通り、I O A P I C を使用する拡張モードでは、レガシーモードと同様、サウス・ブリッジ 6 が故障した場合、S M I ハンドラーがマスター I O A P I C 1 2 A と 1 2 B のコンフィグレーション/ステータス格納部 2 2 を参照して、スタンバイ側のサウス・ブリッジ 6 の I O A P I C 3 6 に対し、全く同様の設定を行うことが可能となる。結果として C P U 2 からみてサウス・ブリッジ 6 の故障を隠蔽することが可能である。

さらに、両モードにおいても、マスター I O A P I C 1 2 A と 1 2 B は常に同期して動作しているので、一方の F T 制御部 1 0 自身が故障し、C P U サブシステムが論理的に切り離されたとしても、他方の正常動作している F T 制御部 1 0 のマスター I O A P I C 1 2 により正常動作を続けることが可能である。こうして、割り込みをロストすることもない。

30

さらに、故障した F T 制御部 1 0 が交換された場合、交換後のモジュールのマスター I O A P I C 1 2 やサウス・ブリッジ 6 内の I O A P I C 3 6 の設定、状態は全て消えてしまっているが、システム・ソフトウェア (S M I ハンドラー) により動作を続行しているシステム側のマスター I O A P I C 1 2 のコンフィグレーション/ステータス格納部 2 2 を参照し、コピーすることにより、完全に同期状態に復帰させることが可能である。

以上に示す通り、F T 制御部内にコンフィグレーション/ステータスを保持する格納部を備えたマスター I O A P I C 1 2 を実装し、割り込みのルーティング制御を行うことで、割り込みコントローラを二重化することが可能である。これにより、フォールト・トレラント・コンピュータシステムを意識せずに作成された既存の O S 、既存のサウス・ブリッジ 6 を搭載したサーバにおいても、割り込みコントローラの二重化を実現することが可能となる。

40

【 発明を実施するための最良の形態 】

【 0 0 2 9 】

以下に、添付図面を参照して、本発明のフォールト・トレラント・コンピュータシステムについて詳細に説明する。本発明のフォールト・トレラント・コンピュータシステムは、例えば、サーバシステムに適用可能である。

【 0 0 3 0 】

50

図 2 は、本発明の実施形態によるフォールト・トレラント・コンピュータシステムの基本構成を示すブロック図である。図 2 に示されるように、本発明のフォールト・トレラント・コンピュータシステムは、同じ構成を有する 2 つ系、即ち 1 次系システム # 1 と 2 次系システム # 2 とを有している。1 次系システムと 2 次系システムの各々は、F T 制御部 1 0 (1 0 A、1 0 B)、C P U 2 (2 A、2 B)、メインメモリ 3 (3 A、3 B)、サウス・ブリッジ 6 (6 A、6 B)、P C I ブリッジ 7 (7 A、7 B)、I / O デバイス 5 (5 - 1 A、5 - 2 A、または 5 - 1 B、5 - 2 B) とを有している。尚、上記で添え字 A は 1 次系を示し、B は 2 次系を示す。F T 制御部 1 0 A と 1 0 B は F T リンク 8 により接続されている。この実施形態では、コントローラが二重化されている。本発明の割込みコントローラは、フォールト・トレラント (F) 制御部の中に内蔵されている。

10

【 0 0 3 1 】

故障個所の交換を可能とするため、1 次系システムと 2 次系システムは別々のボードで構成されていることが好ましい。また、C P U 2 とメインメモリ 3 を有する C P U サブシステムと I / O サブシステムも分離できるよう、4 枚以上のボードで構成するのが理想的である。2 つの C P U サブシステムの各々は、C P U 群 (この実施形態では 1 つの C P U 2)、メインメモリ 3、及び割込みコントローラを含む F T 制御部 1 0 の上半分を有する。2 つの C P U サブシステムはクロックも含め完全に同期して動作する。

【 0 0 3 2 】

I / O サブシステムは、二重化された I / O デバイス群 5 と、P C I ブリッジ 7 と、サウス・ブリッジ 6 とを備えている。I / O サブシステムは、1 次系 / 2 次系とも全く同じハードウェア構成を有している。P C I ブリッジ 7 は、各 I / O デバイス 5 と F T 制御部 1 0 を接続している。

20

【 0 0 3 3 】

図 4 は、二重化されていない、一般的な P C サーバの割込みルーティングを図式的に示している。各 I / O デバイス群 (ここでは P C I デバイス) は最大 4 本 (# A ~ # D) の割込み線を持つことが可能であり、これらは一旦 P C I ブリッジ 7 に接続される。P C I ブリッジ 7 は、複数の割込み線をワイアードオアに接続し、やはり 4 本の割込み線としてサウス・ブリッジ 6 の P I C または I O A P I C に接続される。サウス・ブリッジ 6 には通常、レガシー用の P I C と拡張用の I O A P I C が存在する。現状の P C サーバは、起動時にはレガシー状態で起動されることになっており、この際、割込みコントローラとして P I C が使用される。また、ウィンドウズやリナックスなどの O S が動作する際は、P I C の動作は停止され、より高機能の I O A P I C が使用される。

30

【 0 0 3 4 】

図 3 は、図 2 に示されるシステムでの P C I 階層構造を図式的に示す。全てのアクセス可能なデバイスは、P C I バス仕様に倣い、P C I バス番号と、P C I デバイス番号、P C I 関数番号を持ち、C P U を頂点とした階層構造を有する。サウス・ブリッジ 6 だけは、全く同じデバイス番号を持つが、通常時は片方のみが使用される。以後、使用されるサウス・ブリッジ 6 はアクティブ・サウス・ブリッジ 6 と呼ばれる。他方のサウス・ブリッジ 6 は、スタンバイ・サウス・ブリッジ 6 と呼ばれる。スタンバイ・サウス・ブリッジ 6 は、F T 制御部 1 0 から論理的に切り離されており、フェイル・オーバーが発生するまで、一切のアクセスは不可能である。

40

【 0 0 3 5 】

サウス・ブリッジ 6 は、一般にレガシーデバイスと呼ばれる、シリアルポート、パラレルポート、マウス、キーボード、タイマ、時計等 (いずれも図示せず) のシステムで唯一存在するデバイスを備え、あるいは接続されている。これらレガシーデバイスは、システム上で固定のアドレスを有し、システム上に 2 個存在することは出来ない。また、O S から直接アクセスされることが多い。サウス・ブリッジ 6 は、他の I / O デバイス 5 とは異なる割り込み制御方法を採用し、他の I / O デバイス 5 とは異なり、ソフトウェアによる二重化は不可能である。従って、本発明のフォールト・トレラント・コンピュータシステムでも、動作しているのは、1 次系と 2 次系の一方のみである。他方は、動作中のサウス

50

・ブリッジ 6 が障害を起こすまで、スタンバイ状態として、OS 側からは不可視状態にされる。

【0036】

サウス・ブリッジ 6 は、レガシー割り込み制御部としての P I C (プログラマブル・インターラプト・コントローラ) 34、割り込みコントローラ (I O アドバンスト・プログラマブル・インターラプト・コントローラ) 36、およびルーティングロジック 32 とを備えている。P I C 34 は、起動時のレガシーモードにおいて、これらレガシーデバイスの割り込み制御を行う。I O A P I C 36 は、一般的にインテル系 P C サーバで使用されるサウス・ブリッジ 6 に組み込まれ、サウス/ブリッジに関連する割り込み要求を統括する。ルーティングロジック 32 は、内部で発生した割り込み要求あるいは外部からの割り込み要求を P I C 34 あるいは I O A P I C 36 に出力する。

10

【0037】

F T リンク 8 は、1 次系システムの F T 制御部 # 1 10 A と 2 次系システムの F T 制御部 # 2 10 B とを接続する。F T リンク 8 は、1 次系システムから、2 次系システムへの、及び 2 次系システムから 1 次系システムへの I / O デバイスへのアクセスに使用される。これにより、1 次系システムの F T 制御部 # 1 10 A は、P C I ブリッジ # 1 7 A 及びその配下の I / O デバイス 5 A に対するアクセス・リクエストのみを 2 次系システムの F T 制御部 # 2 10 B に転送する。同様に 2 次系システムの F T 制御部 # 2 10 B は、P C I ブリッジ # 2 7 B 及びその配下の I / O デバイス 5 B に対するアクセスを受け持ち、それらへのアクセス・リクエストのみを 1 次系システムの F T 制御部 # 2 10 A に転送する。したがって、両系の同期チェック範囲も、上記範囲に限られる。即ち、本発明のコンピュータシステムでは、F T 制御部 10 による同期チェックは分散的に行われることになる。

20

【0038】

F T 制御部 10 は、エラー検出部 11、マスター I O A P I C 12、メッセージコンバーター 14、F T コンパレータ 15、ゲートコントローラ 16、ルーター 18、タイマー 19 を備えている。加えて、F T 制御部 10 は、図示しないが、1 次系と 2 次系の C P U サブシステムの同期動作を保証するための同期動作保証制御部を有している。

【0039】

エラー検出器 11 は、C P U あるいは I / O デバイスからのリクエストを比較し、内部あるいは I / O サブシステム等でのエラーを検出する。エラーが検出されたとき、S M I (システム・マネージメント・インターラプト) を生成する。マスター I O A P I C 12 は、C P U 2 がオペレーティング・システム上で動作している拡張モードにおいて割り込み要求を統括し、レガシーモードのとき、I / O サブシステムからの割り込み要求をサウス・ブリッジ 6 へ転送する。また、レガシーモードにおいて、サウス・ブリッジ 6 からの割り込み要求を C P U 2 にスルーで転送する。メッセージコンバータ 14 は、I / O サブシステムからの割り込み要求を割り込みメッセージに変換する。

30

【0040】

ゲート・コントローラ 16 は、マスター I O A P I C 12 からの割り込み要求をサウス・ブリッジ 6 に接続し、サウス・ブリッジ 6 からの割り込み要求をマスター I O A P I C 12 に接続する。ルーター 18 は、C P U からのデータ/コマンドをメインメモリ 3 あるいは I / O サブシステムに転送し、また、I / O サブシステムからのデータ/コマンド及び割り込み要求をメインメモリあるいは C P U に転送する。また、1 次系システムの F T 制御部 10 A のルーター 18 A は、F T リンク 8 を介して 2 次系システムの F T 制御部 10 B の I O A P I C 12 B に割り込み要求を転送する。反対も同様である。なお、他系のマスター I O A P I C 12 への通知は F T リンク 8 を介することになり、タイムラグが発生することになる。しかしながら、そのタイムラグを考慮して予め決められた時間の遅延後、上記割り込みメッセージは自系マスター I O A P I C 12 へ通知されるので、実質的に同一のタイミングで割り込みメッセージが通知されることができる。

40

【0041】

50

F T制御部 10 は、更に、モジュールの物理的な位置、即ち当該 F T制御部 10 が 1 次系システム内にあるのか 2 次系システム内にあるのかを示す外部ピン（図示せず）と、アクティブ・サウス・ブリッジ 6 のアドレス位置を示すアクティブ・サウス・ブリッジ・レジスタ（図示せず）を備えている。F Tコンパレータ 15 は、両者の値を比較し、C P U 2 からの設定コマンドをルーター 18 を介してアクティブ・サウス・ブリッジ 6 に転送する。

【0042】

図 7 は、本発明のコンピュータシステムのシステム・アドレス・マップの例を示している。マスター I O A P I C 1 2 は、例として F E C 0 _ 0 0 0 0 h ~ F E C 7 _ F F F F h にマップされ、設定などもこの空間を介して行われる。サウス・ブリッジ 6 内の I O A P I C 3 6 のアドレス空間はマスター I O A P I C 1 2 のアドレス空間の一部と重なっている。こうして、本発明のコンピュータシステムでは、アクティブ・サウス・ブリッジ 6 内の I O A P I C 3 6 は C P U 2 あるいは O S から隠蔽されて不可視状態にある。しかしながら、アクティブ・サウス・ブリッジ 6 内の I O A P I C 3 6 の設定を行うことは必要であるので、サウス・ブリッジ 6 内の I O A P I C のアドレス空間は、マスター I O A P I C 1 2 のアドレス空間により覆い被されている。

10

【0043】

図 8 に示されるように、C P U から発行された設定コマンドは、ルーター 18 によりマスター I O A P I C 1 2 へ転送され、マスター I O A P I C 1 2 に設定される。また、設定コマンドのうち、サウス・ブリッジ 6 の I O A P I C 3 6 と重複する部分に関しては、ルーター 18 により F Tコンパレータ 15 に転送される。F Tコンパレータ 15 は、設定コマンドとモジュールの物理的な位置データとアクティブ・サウス・ブリッジ・レジスタのデータの組とを比較することにより、アクティブ・サウス・ブリッジ 6 へ設定コマンドをフォワードする。これにより、サウス・ブリッジ 6 の I O A P I C 3 6 の設定、その状態が等価的にマスター I O A P I C 1 2 上に現れることになる。こうして、マスター I O A P I C 1 2 とサウス・ブリッジ 6 の I O A P I C 3 6 の重複する部分の設定は、マスター I O A P I C 1 2 だけでなくサウス・ブリッジ I O A P I C 3 6 にも同様の設定が行われることになる。つまり、全く同じ設定の I O A P I C のコピーを作り出していることになる。

20

【0044】

マスター I O A P I C 1 2 は、マスター割り込みコントローラであり、システム全体の割り込みを統括する。マスター I O A P I C 1 2 は、拡張性をもった割り込みコントローラで、割り込み要因が発生した場合、C P U 2 に対し、メッセージの形で割り込みの番号も一緒に通知する。両系の F T制御部 10 A、10 B 内の 2 個のマスター割り込みコントローラ 12 A と 12 B は、同期動作保証制御部により完全に同期して動作する。P C Iブリッジ 7 A と 7 B 側からの割り込み線 # A ~ # D 上の割り込みは、メッセージコンバーター 14 により I N T # x アサート・メッセージ及び I N T # x デアサート・メッセージに変換され、ルーター 18 により両系の I O A P I C 1 2 A と 1 2 B に同時に通知される。

30

【0045】

P I C 3 4 と I O A P I C 1 2 と 3 6 は以下の点で相違する。即ち、P I C 3 4 は過去の資産を継承するためのレガシー割り込みコントローラであり、割り込み要因が発生した場合、C P U 2 へ一本の割り込み線（I N T R）を使用して、C P U 2 へ割り込み要求を出力する。I N T R 信号を受け取った C P U 2 は、インターラプト・アクノリッジ・コマンドを P C I に対して発行し、割り込みの番号を知る。一方、I O A P I C 1 2 と 3 6 は、さらに拡張性をもった割り込みコントローラであり、割り込み要因が発生した場合、C P U 2 に対し、メッセージの形で割り込みの番号も一緒に通知する。以上の違いが存在するので、C P U 2 への割り込み通知は各系で 2 系統存在することになる。

40

現在のシステムでは、上記のように、O S が起動されるまでのレガシーモードでは P I C 3 4 が使用され、O S の起動後の拡張モードでは I O A P I C 1 2 と 3 6 が使用される。こうして、割り込み要求パスは、切り替えが行われている。

50

【 0 0 4 6 】

図 5 は、本発明の割り込みコントローラ二重化方式を採用したフォールト・トレラント・コンピュータシステムの割り込みルーティングを図式的に示している。あくまでも既存のオープン系デバイス、OS を使用しているので、上述した動作と矛盾なく二重化してある。FT 制御部 10 のマスター I O A P I C 1 2 は、OS から可視のシステムであり、唯一の割り込みコントローラである。両系の FT 制御部 10 内の 2 個のマスター割り込みコントローラ 1 2 は完全に同期して動作する。

【 0 0 4 7 】

上記のように、FT 制御部 10 にはアクティブ/スタンバイ・ゲート・コントローラ 16 が存在している。スタンバイ側では、サウス・ブリッジ 6 は電氣的に FT 制御部 10 に接続されているが、FT 制御部 10 との接続が論理的に切り離されている。この結果、スタンバイ側のサウス・ブリッジ 6 への一切の割り込み通知が遮断される。

【 0 0 4 8 】

図 6 は、FT 制御部とサウス・ブリッジの割り込み要求転送構成を示す図である。サウス・ブリッジ 6 内は一般的な構成となっており、外部及び内部デバイスからの割り込み要求を受け付け、モードによって P I C また I O A P I C に通知先を変更する割り込みルーティングロジック 3 2 と、レガシーモード時の割り込みコントローラである P I C 3 4、拡張時の割り込みコントローラである I O A P I C 3 6 とを備えている。

FT 制御部 10 の I O A P I C 1 2 は、I O A P I C 2 4 と、マスター I O A P I C 1 2 の全ての設定、状態を知ることができるコンフィグレーション/ステータス格納部 2 2 (レジスタ群)と、P C I ブリッジ 7 側から I N T # x メッセージを受け付け、モードにより I N T # x メッセージを I O A P I C 2 4 へ、あるいはゲート・コントローラ 1 6 を介してサウス・ブリッジ 6 へ転送するルーティングロジック 2 0 とを備えている。

【 0 0 4 9 】

ステータス格納部 2 2 は、故障した FT モジュールを交換し、新たなモジュールが接続されて二重化される際、その時点の割り込みコントローラの状態を、他系に完全に再現可能な情報を持っており、フェイルオーバー時にシステム・ソフトウェアから参照される。格納部 2 2 は、

- ・ I O A P I C に対する設定情報
- ・ FT 制御部 10 内 I O A P I C 制御ロジックの内部ステータス (バイナリ状態であり、システム・ソフトウェアがこの値を見て何かを判断するわけではなく、純粋に内部状態をコピーするために使用される)
- ・ P I C に対する設定情報 (FT 制御部 10 は P I C 機能は持っていないが、フェイルオーバー時のサウス・ブリッジ 6 への設定に使用)
- ・ FT 制御部内 P I C 制御ロジックの内部ステータス (バイナリー状態)
- ・ その他、割り込みコントローラに対する設定情報 (FT 制御部独自のレジスタ設定情報など)
- ・ その他、割り込みコントローラロジックの内部ステータス (バイナリー状態)

を保持している。格納部 2 2 の内容の全てを、交換された新モジュールの格納部 2 2 へコピーすることで、マスター I O A P I C 1 2 はコピー元と全く同じ設定、動作状況となり、完全に同期して動作することが可能となる。

【 0 0 5 0 】

尚、タイマー 1 9 により周期的にエラー検出部 1 1 がアクティブとされ、両系のステータス格納部 2 2 が比較されてもよい。この比較の結果、不一致が検出されたときには、S M I ハンドラーに知らされる。S M I ハンドラーは、C P U 2 の動作を停止し、両系のステータス格納部 2 2 の格納データが同じになるように、データの転送処理をする。その後、S M I ハンドラーは、C P U 2 の動作を再開する。こうして、累積による誤差を所定時間ごとに除くことができる。また、所定時間ごとに、あるいは C P U 2 または I / O デバイス 5 からのリクエストに回答して、チェックしてサウス・ブリッジ 6 あるいはその他の個所に障害の発生を検出した場合、S M I ハンドラーに知らされる。S M I ハンドラーは

、CPU 2の動作を停止する。故障個所のボードが交換された後、S M Iハンドラーは、両系のステータス格納部 2 2の格納データが同じになるように、データの転送処理をする。その後、S M Iハンドラーは、CPU 2の動作を再開する。

【0051】

レガシーモード時には、サウス・ブリッジ 6内のP I C 3 4がシステムで唯一の割り込みコントローラとして利用されることができるよう、P C Iブリッジ 7からのI N T # xメッセージを割り込み信号線 # A ~ # Dに戻してサウス・ブリッジ 6に接続するための出力が存在する。CPU 2への割り込みは、マスターI O A P I Cが行う。サウス・ブリッジ 6のP I C 3 4からの割り込み要求I N T Rと、I O A P I C 3 6からの割り込みメッセージはアクティブ/スタンバイ・ゲート・コントローラ 1 6を介してマスターI O A P I C 1 2に接続される。こうして、マスターI O A P I C 1 2は、サウス・ブリッジ 6の割り込み要求I N T RをスルーによりCPU 2に転送するI N T R割り込み線と、マスターI O A P I C 1 2で処理される割り込みメッセージの出力が存在し、共にCPU 2に接続されている。

10

【0052】

レガシーモードの後の拡張モードでは

CPU 2との割り込みに関する授受は全てマスターI O A P I C 1 2により行われる。マスターI O A P I C 1 2は、拡張モード時にはアクティブ・サウス・ブリッジ 6の割り込みと、P C Iブリッジ 7からの割り込みを統括して管理する。このため、マスターI O A P I C 1 2の一部はアクティブ・サウス・ブリッジのI O A P I C 3 6が、そのまま透過的に見える形となる。このため各サウス・ブリッジ 6内のI O A P I C 3 6は、システムからは不可視状態にされる。これは、サウス・ブリッジ 6が故障するケースを考慮したものである。アクティブ・サウス・ブリッジ 6が故障した場合、マスターI O A P I C 1 2の割り込み制御は、直ちにスタンバイ・サウス・ブリッジ 6のI O A P I C 3 6に置換される。このため、OS側からは特にI O A P I Cの増減は発生しない。

20

【0053】

次に、図 2に示される本発明のコンピュータシステムにおいて、1次系側のサウス・ブリッジ 6をアクティブ・サウス・ブリッジと呼び、通常処理で使用されるサウス・ブリッジと仮定する。

【0054】

まず、レガシーモードでの動作を説明する。レガシーモードでは、P I C 3 4が割り込み制御の中心となる。P I C 3 4を使用する場合、システムで唯一のP I C 3 4が全てのデバイスの割り込みを制御することになる。F T制御部 1 0は、P C Iブリッジ 7下のI / Oデバイス 5 - 1, 5 - 2からの割り込み状態を監視することができるが、サウス・ブリッジ 6内のデバイスの状態を把握することは不可能である。このため、結果として、P I C 3 4はアクティブ・サウス・ブリッジ 6内のP I Cを使うことになる。

30

【0055】

図 1 0において、P C Iデバイス # 1 bから割り込み信号がアサートされる。このとき、割り込み要求は、P C Iブリッジ # 1 7 Aを介してF T制御部 1 0 Aのメッセージコンバーター 1 4 Aに通知される(ステップ S 1)。メッセージコンバーター 1 4 Aは、信号線の状態、即ち割り込み要求をI N T # xアサート・メッセージに変換し、両方のマスターI O A P I C 1 2 Aと 1 2 Bに通知を行う(ステップ S 2)。F T制御部 # 2 1 0 BへはF Tリンク 8を介して通知され、F T制御部 # 1 1 0 Aは、前もって設定された遅延を経てからマスターI O A P I C 1 2 Aに通知される。これにより両系のマスターI O A P I C 1 2 Aと 1 2 Bは同時に割り込み通知を受け取り、完全に同期して動作することが可能である。

40

【0056】

両マスターI O A P I C 1 2 Aと 1 2 Bは、さらにゲート・コントローラ 1 6 Aと 1 6 Bに対し、I N T # xアサート・メッセージを送る(ステップ S 3)。ゲート・コントローラ 1 6 Aは、ボード位置ピンとアクティブ・サウス・ブリッジ・レジスタの値に基づい

50

て、自身がアクティブだと判断すると、INT# xアサート・メッセージを割り込み信号線INT# xに戻し、サウス・ブリッジ6に通知する(ステップS4)。一般的に、サウス・ブリッジ6は、図6に示される構成を有しており、外部から入力された割り込み要因はルーティングロジック32に供給される。なお、シリアルポート、パラレルポート、マウス、キーボード、タイマ、時計などの、もともとサウス・ブリッジ6内にある内部デバイスの割り込みも同様にルーティングロジック32に供給される。この場合、割り込み通知はここからのスタートとなる。

【0057】

サウス・ブリッジ6内のルーティングロジック32は、レガシーモードにあるので、割り込みをPIC34に通知する。PIC34は、INTR信号として割り込み線をアサートする(ステップS5)。ゲート・コントローラ16Sは、INTR信号をINTRアサート・メッセージに変換して、双方のマスターIOAPIC12Aと12Bに通知を行う。この時、スタンバイ側のマスターIOAPIC12Bへの通知はFTリンク8を介して行われるが、実際には先のINT# xメッセージと同じパスを通過するため、INTRアサート・メッセージは両方のマスターIOAPIC12Aと12Bに同時に通知される(ステップS6)。両方のマスターIOAPIC12Aと12Bは、INTRアサート・メッセージを受け取ると、CPU2に対して同時にINTRをアサートする(ステップS7)。

10

【0058】

レガシーモードでアクティブ・サウス・ブリッジ6が故障した場合、サウス・ブリッジ6の故障を示す割り込みが両CPU2に通知され、FT制御用のシステム・ソフトウェアがコールされる。このシステム・ソフトウェアの呼び出しのためには、最高レベルの割り込みが使用され、例えば、インテル系CPUではシステム・マネージメント・インターラプト(SMI)が使用される。これにより、CPU2上で実行されている処理は全て一旦止められる。この停止の間に、SMIハンドラーは、アクティブ・サウス・ブリッジ6の設定を全てスタンバイ側にコピーし、アクティブ・サウス・ブリッジ・レジスタの値を入れ替える。SMIハンドラーの処理の終了後、一旦停止させられたCPU2の処理が再開される。このとき、サウス・ブリッジ6が入れ替わったことは完全に隠蔽されている。

20

【0059】

次に拡張モード、つまりIOAPICが使用される場合を説明する。図9は、各種デバイスと割り込み番号(IRQ)との対応を示すテーブルである。IOAPICは、割り込み要因を受け付けるとそのIRQを直接CPUへ通知するため、このようなテーブルがIOAPIC内に設けられている。サウス・ブリッジ6のIRQテーブルは一般的な設定で、特にインテル系CPUシステムではIRQ0~IRQ15は固定的に決まっている。サウス・ブリッジ6のIRQテーブルの設定は実際にはマスターIOAPIC12に対してなされるが、図8に示されるように、ルーター18により同じ設定コマンドがアクティブ・サウス・ブリッジ6に送られるので、結果的に両者は同じ設定となる。アクティブ・サウス・ブリッジ6から発行された割り込みメッセージは、そのままマスターIOAPIC12の割り込み受け付けに置き換えられる。また、PCIブリッジ7からの割り込みは、例としてIRQ20~27に割り付けられる。

30

40

【0060】

図11は、IOAPICが使用される場合の動作を示している。I/Oデバイス#2b5-2Aが、INTR信号をアサートしたと仮定する(ステップS1)。割り込みは、PCIブリッジ#27Bを経由してFT制御部#210Bに割り込み信号、仮にINT#Cで通知されたとする。これを受け取ったメッセージコンバータ114Bは、両系のマスターIOAPIC12Aと12Bに対し、INT# cアサート・メッセージを通知する(ステップS2)。マスターIOAPIC12Aと12Bは、PCIブリッジ#27BからのINT# cをIRQ26と判断し、CPU2へ割り込みメッセージを通知する(ステップS3)。詳細は略すが、サウス・ブリッジ6からの割り込みも同様の経路を通過する。

50

【図面の簡単な説明】

【0061】

【図1】図1は、フォールト・トレラント・コンピュータシステムの構成の一例を示すブロック図である。

【図2】図2は、本発明の実施形態によるフォールト・トレラント・コンピュータシステムの基本構成を示すブロック図である。

【図3】図3は、図2に示されるシステムのP C I階層構造を図式的に示すブロック図である。

【図4】図4は、二重化されていないP Cサーバの割込みルーティングを図式的に示すブロック図である。

【図5】図5は、本発明の割込みコントローラ二重化方式を採用したフォールト・トレラント・コンピュータシステムの割込みルーティングを図式的に示すブロック図である。

【図6】図6は、F T制御部のマスターI O A P I Cの構成を示すブロック図である。

【図7】図7は、本発明のシステムのシステム・アドレス・マップの例を示す図である。

【図8】図8は、ルーターにより同じ設定コマンドがアクティブ・サウス・ブリッジに送られる様子を示す図である。

【図9】図9は、各種デバイスと割込み番号（I R Q）との対応を示すテーブルである。

【図10】図10は、レガシーモードにおける割り込み制御を示す図である。

【図11】図11は、拡張モードにおける割り込み制御を示す図である。

【符号の説明】

【0062】

10（10A、10B）：F T制御部

2（2A、2B）：C P U

3（3A、3B）：メインメモリ

6（6A、6B）：サウス・ブリッジ

7（7A、7B）：P C Iブリッジ

5（5-1（5-1A、5-1B）、5-2（5-2A、5-2B））：I / Oデバイス

8：F Tリンク

12：マスターI O A P I C

14：メッセージコンバーター

15、F Tコンパレーター

16：ゲートコントローラ

18：ルーター

20：ルーティングロジック

22：コンフィグレーション/ステータス格納部（レジスタ群）

32：割り込みルーティングロジック

34：P I C

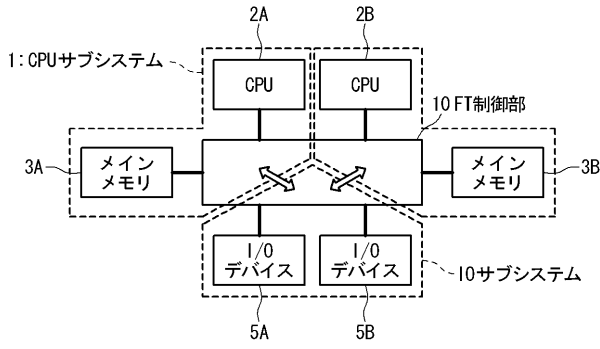
36：I O A P I C

10

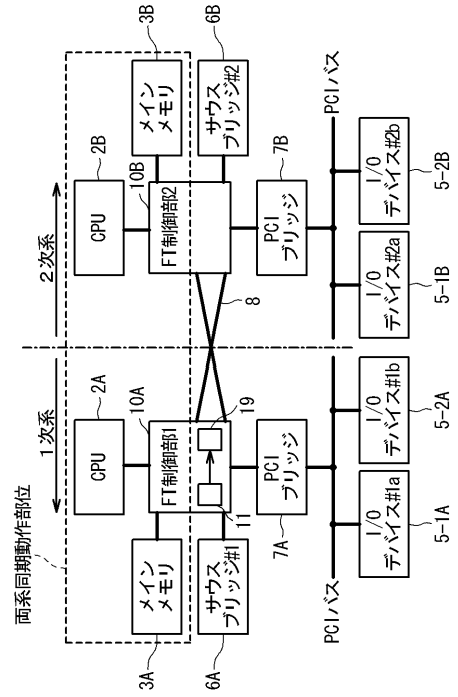
20

30

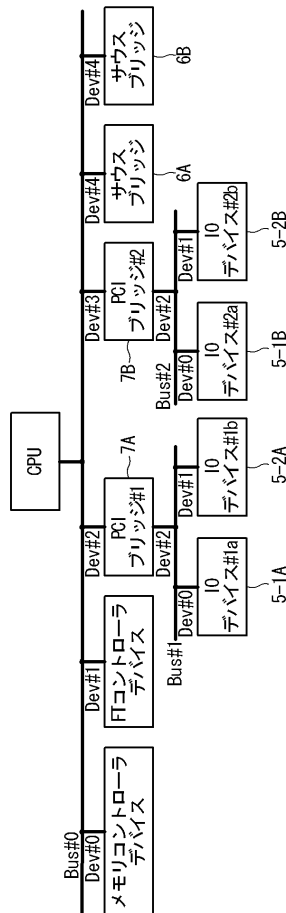
【 図 1 】



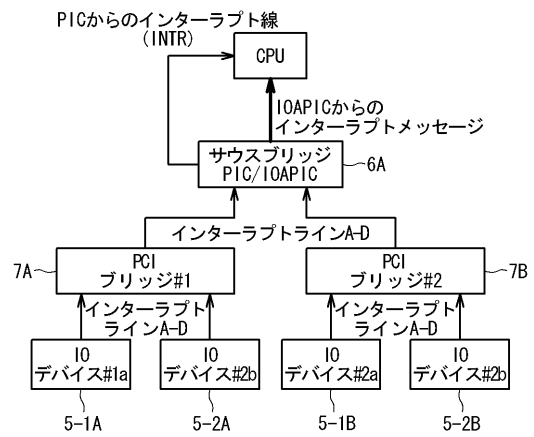
【 図 2 】



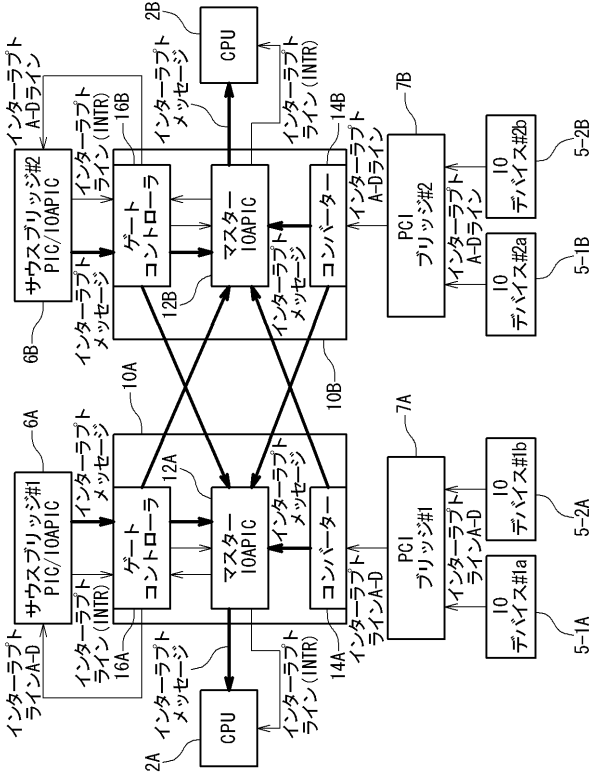
【 図 3 】



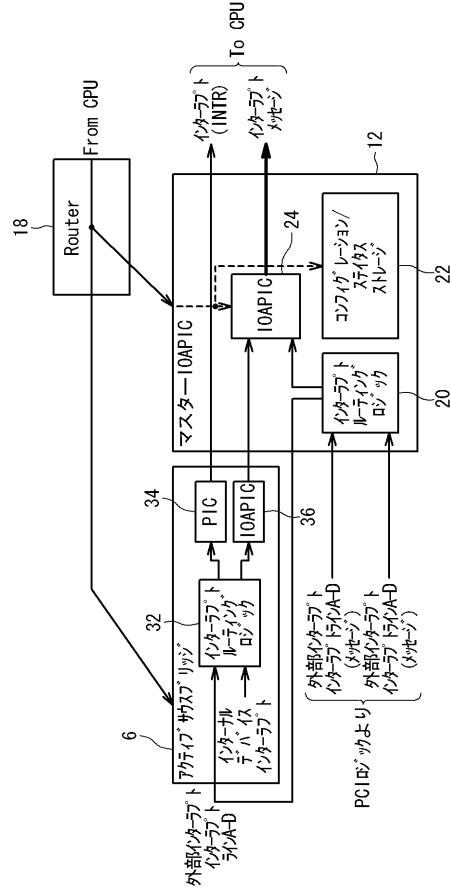
【 図 4 】



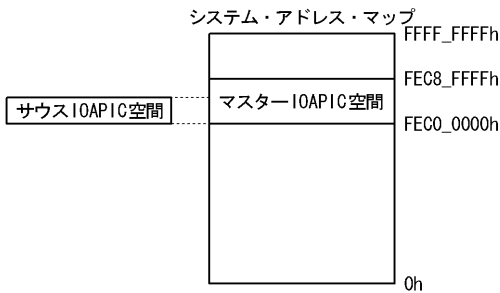
【 図 5 】



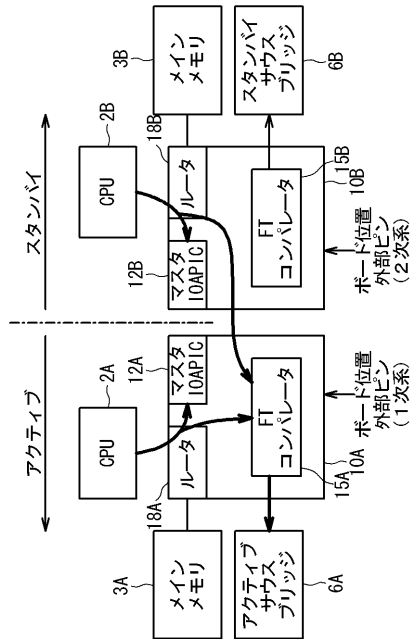
【 図 6 】



【 図 7 】



【 図 8 】

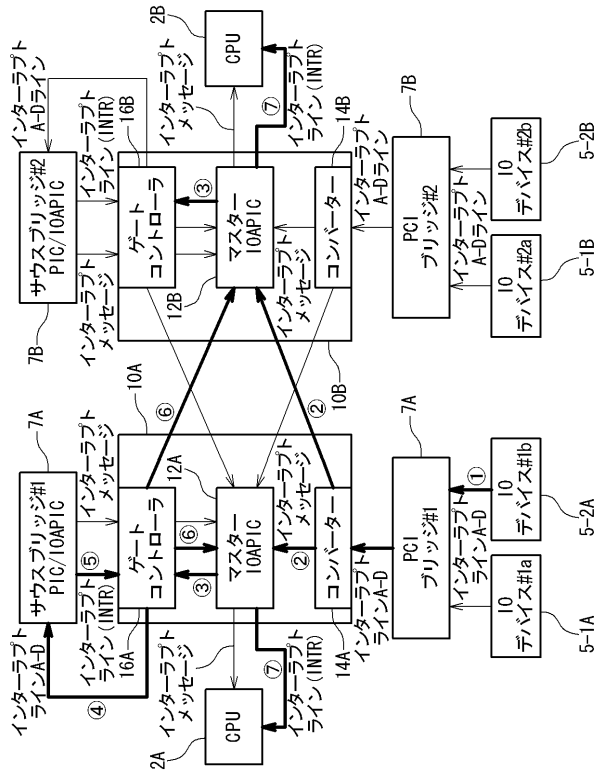


【 図 9 】

サウスブリッジ IOAPIC IRQテーブル		FTコントローラ マスタ IOAPIC IRQテーブル		
0	INTR	→	0	INTR
1	キーボード	→	1	キーボード
2	タイマ	→	2	タイマ
3	COMポート#1	→	3	COMポート#1
4	COMポート#2	→	4	COMポート#2
5	パラレルポート#1	→	5	パラレルポート#1
6	フロッピーディスク	→	6	フロッピーディスク
7	パラレルポート#2	→	7	パラレルポート#2
8	リアルタイムクロック	→	8	リアルタイムクロック
9		→	9	
10		→	10	
11		→	11	
12	マウス	→	12	マウス
13	フローティングポイントユニット	→	13	フローティングポイントユニット
14	一次系IDEディスク	→	14	一次系IDEディスク
15	二次系IDEディスク	→	15	二次系IDEディスク
16	外部INT#A	→	16	外部INT#A
17	外部INT#B	→	17	外部INT#B
18	外部INT#C	→	18	外部INT#C
19	外部INT#D	→	19	外部INT#D
20		→	20	PCIブリッジ#1INT#A
21		→	21	PCIブリッジ#1INT#B
22		→	22	PCIブリッジ#1INT#C
23		→	23	PCIブリッジ#1INT#D
24		→	24	PCIブリッジ#1INT#E
25		→	25	PCIブリッジ#1INT#F
26		→	26	PCIブリッジ#1INT#G
27		→	27	PCIブリッジ#1INT#H

PCIブリッジから外部インターラプト線入力 A-D (メッセージ)

【 図 10 】



【 図 11 】

