

(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2014-59750
(P2014-59750A)

(43) 公開日 平成26年4月3日(2014.4.3)

(51) Int.Cl.	F 1	テーマコード (参考)
G06F 11/20 (2006.01)	G06F 11/20 310C	5B034
G06F 9/46 (2006.01)	G06F 9/46 350	
G06F 11/18 (2006.01)	G06F 11/18 310F	

審査請求 有 請求項の数 15 O L (全 33 頁)

(21) 出願番号	特願2012-204608 (P2012-204608)	(71) 出願人	000006507 横河電機株式会社 東京都武蔵野市中町2丁目9番32号
(22) 出願日	平成24年9月18日 (2012.9.18)	(72) 発明者	吉田 善貴 東京都武蔵野市中町2丁目9番32号 横河電機株式会社内
		(72) 発明者	植原 正太 東京都武蔵野市中町2丁目9番32号 横河電機株式会社内
		(72) 発明者	大野 毅 東京都武蔵野市中町2丁目9番32号 横河電機株式会社内
		Fターム(参考)	5B034 BB17 DD06

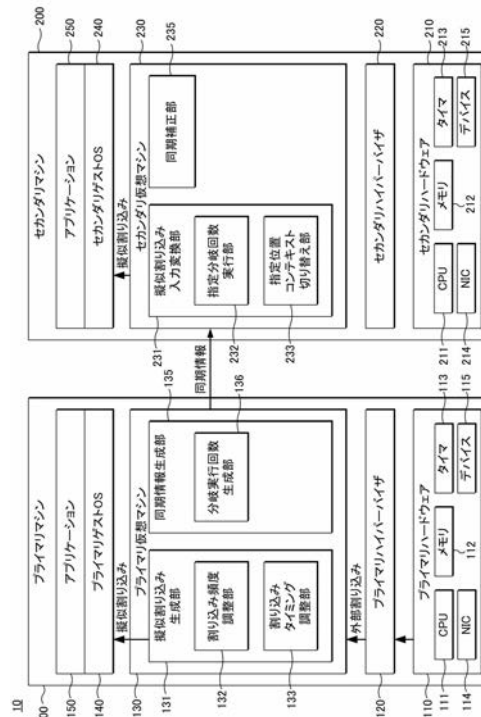
(54) 【発明の名称】 フォールトトレラントシステム

(57) 【要約】

【課題】 2台のコンピュータ上に形成された2台の仮想マシンを並列同期動作させるフォールトトレラントシステムにおいて、それぞれのコンピュータのハードウェアの差異を適切に調整する。

【解決手段】 プライマリ仮想マシンと、セカンダリ仮想マシンとを備えたフォールトトレラントシステムであって、プライマリ仮想マシンは、セカンダリ仮想マシンから、動作性能情報を収集し、プライマリ仮想マシンの動作能力と、セカンダリ仮想マシンの動作能力が等しくなるように、プライマリ仮想マシンの動作能力とセカンダリ仮想マシンの動作能力とを設定する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

ブラマリハイパーバイザが稼働するプライマリマシン上に形成され、前記ブラマリハイパーバイザからの外部割り込みに基づく擬似割り込みをブラマリゲストOSに入力するプライマリ仮想マシンと、

セカンダリハイパーバイザが稼働するセカンダリマシン上に形成され、前記プライマリ仮想マシンから伝達された前記擬似割り込みのタイミング情報に基づいて、擬似割り込みをセカンダリゲストOSに入力するセカンダリ仮想マシンと、を備えたフォールトトレラントシステムであって、

前記プライマリ仮想マシンは、

前記セカンダリ仮想マシンから、動作性能情報を収集し、前記プライマリ仮想マシンの動作能力と、前記セカンダリ仮想マシンの動作能力が等しくなるように、前記プライマリ仮想マシンの動作能力と前記セカンダリ仮想マシンの動作能力とを設定することを特徴とするフォールトトレラントシステム。

【請求項 2】

前記プライマリ仮想マシンは、

設定した前記セカンダリ仮想マシンの動作能力を前記セカンダリ仮想マシンに通知し、

前記セカンダリ仮想マシンは、

通知された動作能力に従って動作を行なうことを特徴とする請求項 1 に記載のフォールトトレラントシステム。

【請求項 3】

前記動作能力は、動作周波数を含み、

前記プライマリ仮想マシンは、

前記プライマリ仮想マシンが使用可能な動作周波数と、前記セカンダリ仮想マシンが使用可能な動作周波数とを収集して、前記プライマリ仮想マシンの動作周波数と、前記セカンダリ仮想マシンの動作周波数とを設定することを特徴とする請求項 1 または 2 に記載のフォールトトレラントシステム。

【請求項 4】

前記プライマリ仮想マシンは、

前記プライマリ仮想マシンが使用するデバイスの動作仕様と、前記セカンダリ仮想マシンが使用するデバイスの動作仕様とが異なる場合に、

前記プライマリ仮想マシンが使用するデバイスに対する処理が、前記セカンダリ仮想マシンが使用するデバイスに対する処理と同じになるようにエミュレートすることを特徴とする請求項 1 ~ 3 のいずれか 1 項に記載のフォールトトレラントシステム。

【請求項 5】

前記プライマリ仮想マシンは、

前記デバイスの動作仕様の違いが所定の機能の有無の場合は、

前記所定の機能を無効にして、前記プライマリ仮想マシンが使用するデバイスに対する処理が、前記セカンダリ仮想マシンが使用するデバイスに対する処理と同じになるようにエミュレートし、

前記デバイスの動作仕様の違いが能力あるいは容量の場合は、動作に相違が生じる処理を前記ブラマリゲストOSの外部で実行させることを特徴とする請求項 4 に記載のフォールトトレラントシステム。

【請求項 6】

前記セカンダリ仮想マシンは、

前記プライマリ仮想マシンが使用するデバイスの動作仕様と、前記セカンダリ仮想マシンが使用するデバイスの動作仕様とが異なる場合に、

前記セカンダリ仮想マシンが使用するデバイスに対する処理が、前記プライマリ仮想マシンが使用するデバイスに対する処理と同じになるようにエミュレートすることを特徴とする請求項 1 ~ 3 のいずれか 1 項に記載のフォールトトレラントシステム。

10

20

30

40

50

【請求項 7】

前記セカンダリ仮想マシンは、

前記デバイスの動作仕様の違いが所定の機能の有無の場合は、

前記所定の機能を無効にして、前記セカンダリ仮想マシンが使用するデバイスに対する処理が、前記プライマリ仮想マシンが使用するデバイスに対する処理と同じになるようにエミュレートし、

前記デバイスの動作仕様の違いが能力あるいは容量の場合は、動作に相違が生じる処理を前記セカンダリゲストOSの外部で実行させることを特徴とする請求項6に記載のフォールトトレラントシステム。

【請求項 8】

前記セカンダリ仮想マシンは、

前記動作能力の設定後であっても、前記プライマリ仮想マシンより実行速度が遅い場合には、

遅れの度合いを取得し、取得した値が予め定められた遅延許容時間を超えると、同期初期化処理を行なうことを特徴とする請求項1～7のいずれか1項に記載のフォールトトレラントシステム。

【請求項 9】

前記セカンダリ仮想マシンは、

前記取得した遅れの度合いが、想定される量より大きい場合に、故障として検出することを特徴とする請求項8に記載のフォールトトレラントシステム。

【請求項 10】

前記セカンダリ仮想マシンは、

前記動作能力の設定後であっても、前記プライマリ仮想マシンより実行速度が遅い場合には、前記プライマリゲストOSのアイドル時に、擬似割り込み入力タイミングを前倒しすることで、前記プライマリ仮想マシンに対する遅れを取り戻すことを特徴とする請求項1～7のいずれか1項に記載のフォールトトレラントシステム。

【請求項 11】

前記プライマリ仮想マシンは、

前記動作能力の設定後であっても、前記セカンダリ仮想マシンより実行速度が速い場合には、前記セカンダリ仮想マシンの実行速度に合わせて、ウェイトを入れながら動作を行なうことを特徴とする請求項1～7のいずれか1項に記載のフォールトトレラントシステム。

【請求項 12】

前記セカンダリ仮想マシンは、

前記動作能力の設定後であっても、前記プライマリ仮想マシンより実行速度が速い場合には、前記プライマリ仮想マシンの実行速度に合わせて、ウェイトを入れながら動作を行なうことを特徴とする請求項1～7のいずれか1項に記載のフォールトトレラントシステム。

【請求項 13】

前記プライマリ仮想マシンおよび前記セカンダリ仮想マシンは、

それぞれのストレージデバイスの内容を、それぞれの仮想マシンに割り当てられたメモリ上に展開し、ストレージデバイスアクセスをメモリアccessに置き換えることを特徴とする請求項1～12のいずれか1項に記載のフォールトトレラントシステム。

【請求項 14】

前記プライマリ仮想マシンおよび前記セカンダリ仮想マシンは、

前記メモリアccessを行なうと、アクセス完了の外部割り込みを待つことなく、それぞれのゲストOSに擬似割り込みを入力することを特徴とする請求項13に記載のフォールトトレラントシステム。

【請求項 15】

前記プライマリ仮想マシンおよび前記セカンダリ仮想マシンは、

10

20

30

40

50

それぞれのストレージデバイスの内容を、それぞれのゲストOSで設定されたRAMディスクに展開し、ストレージデバイスアクセスをRAMディスクに置き換えることを特徴とする請求項1～12のいずれか1項に記載のフォールトトレラントシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、仮想マシンを利用したフォールトトレラントシステムに関する。

【背景技術】

【0002】

フォールトトレラントシステムは、その構成の一部に問題が生じてシステム全体が機能停止することなく動作を継続できるシステムであり、特に、連続可動性が要求されるシステム等に適用される。例えば、フォールトトレラントを適用したサーバコンピュータは、ハードウェアに障害が発生しても、外部機器のクライアントアプリケーションからのネットワークアクセスに対して、通信エラーとならずに、正しいデータを出力することができる。

10

【0003】

特許文献1に記載されているように、通信可能な2台のコンピュータ上でそれぞれ稼働する仮想マシンでフォールトトレラントシステムを実現する手法が知られている。仮想マシンを利用したフォールトトレラントシステムは、2台の仮想マシンに実行状態を同期させて同じ動作をさせており、一方のコンピュータに障害が発生した場合でも、他方のコンピュータで稼働する仮想マシンが処理を引き継ぐことで、システムのサービスを途切れることなく連続して提供することができる。

20

【0004】

仮想マシンを利用したフォールトトレラントシステムは、一方の仮想マシンをプライマリ仮想マシンとし、他方の仮想マシンをセカンダリ仮想マシンとする。そして、プライマリ仮想マシンを先行して動作させるとともに、外部との入出力のやり取りはプライマリ仮想マシンのみが行なうようにする。

【0005】

一般に、同じプログラムを実行している2台のコンピュータは、外部からの入力に正確に同じタイミングで入力されると、同じ動作を行なって同じ出力を行なう。そこで、仮想マシンを利用したフォールトトレラントシステムは、プライマリ仮想マシンで、外部入力に基づく割り込みが発生すると、割り込みを行なったタイミングを同期情報としてセカンダリ仮想マシンに伝達する。そして、遅れて動作を行なっているセカンダリ仮想マシンが、同期情報で通知されたタイミングと同じタイミングで擬似的に割り込みを行なうことで、プライマリ仮想マシンとセカンダリ仮想マシンは、同期して同じ動作を行なうことになる。

30

【0006】

図6は、仮想マシンを利用した従来のフォールトトレラントシステムの構成を示すブロック図である。本図に示すように、フォールトトレラントシステム60は、ネットワーク接続されたプライマリマシン600とセカンダリマシン700とを備えている。

40

【0007】

プライマリマシン600は、物理的なコンピュータ環境であるプライマリハードウェア610上でプライマリハイパーバイザ620が動作しており、プライマリ仮想マシン630が形成されている。プライマリ仮想マシン630では、プライマリゲストOS640が動作しており、プライマリゲストOS640上で、アプリケーション650が実行されている。

【0008】

プライマリハードウェア610は、CPU、メモリ、ネットワークインタフェースカード(NIC)、ストレージ等の各種デバイス等を備えている。

【0009】

50

プライマリ仮想マシン 630 は、プライマリハードウェア 610 の一部が割り当てられ、外部との入出力を行なって動作する仮想的なコンピュータ環境である。プライマリ仮想マシン 630 は、プライマリハイパーバイザ 620 により管理される。

【0010】

セカンダリマシン 700 も同様に、物理的なコンピュータ環境であるセカンダリハードウェア 710 上でセカンダリハイパーバイザ 720 が動作しており、セカンダリ仮想マシン 730 が形成されている。セカンダリ仮想マシン 730 では、セカンダリゲスト OS 740 が動作しており、セカンダリゲスト OS 740 上で、アプリケーション 750 が実行されている。

【0011】

セカンダリハードウェア 710 は、CPU、メモリ、ネットワークインタフェースカード (NIC)、ストレージ等の各種デバイス等を備えている。

【0012】

セカンダリ仮想マシン 730 は、セカンダリハードウェア 710 の一部が割り当てられ、プライマリ仮想マシン 630 に同期して動作する仮想的なコンピュータ環境である。セカンダリ仮想マシン 730 は、セカンダリハイパーバイザ 720 により管理される。

【0013】

従来のフォールトトレラントシステム 60 において、プライマリ仮想マシン 630 とセカンダリ仮想マシン 730 との実行状態の同期は以下のような手順で行なわれている。

【0014】

プライマリハイパーバイザ 620 が、プライマリハードウェア 610 から外部割り込みを受けると、プライマリ仮想マシン 630 に外部割り込みの入力を行なう。

【0015】

次いで、プライマリ仮想マシン 630 が、プライマリゲスト OS 640 に擬似割り込みの入力を行なう。ここで、プライマリ仮想マシン 630 からプライマリゲスト OS 640 への擬似割り込みの入力について説明する。

【0016】

プライマリゲスト OS 640 の処理中に、外部割り込み、特権命令、例外等の仮想マシンコンテキスト切り替えイベントが発生すると、プライマリゲスト OS 640 の処理を中断して、ゲスト OS コンテキストから仮想マシンコンテキストに切り替わり、プライマリ仮想マシン 630 に処理が遷移する。

【0017】

プライマリ仮想マシン 630 は、各種イベントに応じてそのタイミングでプライマリゲスト OS 640 に擬似割り込みを行なう必要がある場合は、擬似割り込みの設定を行なう。擬似割り込みの設定を行なった場合は、プライマリ仮想マシン 630 が処理を終えて、イベント発生時に割り込まれたプライマリゲスト OS 640 に処理が戻るときにプライマリゲスト OS 640 に擬似割り込みが入力される。

【0018】

プライマリ仮想マシン 630 は、プライマリゲスト OS 640 に擬似割り込みの入力を行なうと、セカンダリ仮想マシン 730 に同期情報の伝達を行なう。同期情報には、擬似割り込みの識別情報と擬似割り込みを入力する同期タイミング情報とが含まれる。

【0019】

同期タイミング情報は、プライマリゲスト OS 640 に入力された擬似割り込みと同じタイミングでセカンダリゲスト OS 740 に擬似割り込みを入力するための情報であり、CPU 固有の実行中断位置と実行命令数とが含まれる。

【0020】

実行中断位置は、例えば、擬似割り込みが入力されたときに実行していた命令のアドレスを示すプログラムカウンタの値を用いることができる。また、実行命令数は、CPU が備えるパフォーマンスカウンタの CPU 実行命令数カウンタで計測することができる。

【0021】

10

20

30

40

50

この場合、プライマリ仮想マシン630がプライマリゲストOS640に擬似割り込みの入力を行なうと、プライマリゲストOS640の実行を再開する前に、CPU実行命令数カウンタをゼロクリアし、CPU実行命令数カウンタを有効化する。これにより、前回の擬似割り込み入力からプライマリゲストOS640が実行した命令数をカウントすることができる。

【0022】

同期タイミング情報として実行中断位置だけを用いると、実行中断位置の示す命令がループ処理に含まれる場合や条件分岐先にある場合に、ループや条件分岐の度にその命令を通過するため、擬似割り込み入力のタイミングを特定することができない。

【0023】

また、同期タイミング情報として実行命令数だけを用いると、パイプライン機能などの高速化技術により、セカンダリゲストOS740が指定された実行命令数で停止できず、その実行命令数を超えて停止してしまい、プライマリゲストOS640と同じタイミングで擬似割り込みを入力することができない場合がある。

【0024】

そこで、同期タイミング情報として実行中断位置と実行命令数とを組み合わせ、セカンダリゲストOS740において、実行中断位置が示す命令を通る度に、実行命令数を確認することで、セカンダリゲストOS740を、プライマリゲストOS640で擬似割り込みが入力されたのと同じタイミングで停止させるようにしている。

【0025】

従って、プライマリ仮想マシン630は、プライマリゲストOS640の実行を中断した後、プライマリゲストOS640に擬似割り込みの入力を行なうと、プライマリゲストOS640が実行を中断した時点のプログラムカウンタの値と、CPU実行命令数カウンタの値とを取得し、同期タイミング情報を生成する。そして、擬似割り込みの識別情報と同期タイミング情報とを同期情報としてセカンダリ仮想マシン730に伝達する。

【0026】

同期情報を受信したセカンダリ仮想マシン730は、同期タイミング情報に従ってセカンダリゲストOS740の実行を中断させる。このときのセカンダリ仮想マシン730の動作について、図7のフローチャートを参照して説明する。

【0027】

同期タイミング情報の実行中断位置で指定されたプログラム位置にブレーク命令を埋め込んで(S301)、セカンダリゲストOS740を再開させる(S302)。そして、セカンダリゲストOS740が停止すると(S303: Yes)、CPU実行命令数カウンタを確認して、指定された実行命令数と一致していたら(S304: Yes)、その停止位置でセカンダリゲストOS740を中断させる(S305)。指定された実行命令数と一致していなければ(S304: No)、再度セカンダリゲストOS740を再開させ(S302)、実行命令数の確認処理を繰り返す。

【0028】

セカンダリ仮想マシン730は、セカンダリゲストOS740を中断させると、同期情報の擬似割り込み識別情報に従って擬似割り込みの設定を行なって、(S306)、セカンダリゲストOS740を再開させる(S307)。これにより、セカンダリゲストOS740に、プライマリゲストOS640と同じタイミングで擬似割り込みが入力され、プライマリ仮想マシン630とセカンダリ仮想マシン730との実行状態が同期する。

【0029】

なお、プライマリマシン600、セカンダリマシン700のいずれかでハードウェアの故障が生じると、実行状態の同期が乱れる。実行状態の同期が乱れると、プライマリ仮想マシン630とセカンダリ仮想マシン730とで、出力データの値が異なってくる。そこで、従来のフォールトトレラントシステム60は、プライマリ仮想マシン630の出力とセカンダリ仮想マシン730の出力とを照合し、出力データの値が異なる場合は、どちらかのハードウェアで故障が発生したと判断する。

10

20

30

40

50

【0030】

この故障判断処理を行なうため、セカンダリ仮想マシン730に出力データ照合部731を設け、プライマリ仮想マシン630の出力データを収集して、セカンダリ仮想マシン730の出力データと照合するようにしている。

【先行技術文献】

【特許文献】

【0031】

【特許文献1】特開2009-80695号公報

【発明の概要】

【発明が解決しようとする課題】

10

【0032】

フォールトトレラントシステム60において、プライマリマシン600のプライマリハードウェア610とセカンダリマシン700のセカンダリハードウェア710は、CPU、メモリ、デバイス等、種々のハードウェアで構成されているが、双方のハードウェアの性能は必ずしも一致するとは限られない。

【0033】

例えば、CPUのスペックが異なっていたり、メモリの容量が異なっている場合も想定される。また、フォールトトレラントシステム60の運用過程においてデバイスが交換される等により、双方のデバイスが異なっている場合や、同じデバイスであっても、経年劣化等により性能差が生じる場合もある。

20

【0034】

このため、プライマリ仮想マシン630とセカンダリ仮想マシン730とで処理速度に差が生じる場合がある。上述のように、フォールトトレラントシステム60は、プライマリ仮想マシン630が先行して動作を行なうようにしているが、例えば、セカンダリ仮想マシン730の方が、処理速度が速いと、動作の逆転現象が生じ、プライマリ仮想マシン630からの同期情報に基づく実行状態の同期が行なえなくなる。プライマリ仮想マシン630の方が、処理速度が速いと、セカンダリ仮想マシン730の遅れが蓄積されていく。

【0035】

また、デバイスが異なっていると、プライマリゲストOS640とセカンダリゲストOS740のドライバ処理が異なって、両者の動作に相違が生じる結果、同期ずれが生じるおそれがある。

30

【0036】

このため、ハードウェアの差異を埋める仕組みが必要となる。そこで、本発明は、2台のコンピュータ上に形成された2台の仮想マシンを並列同期動作させるフォールトトレラントシステムにおいて、それぞれのコンピュータのハードウェアの差異を適切に調整することを目的とする。

【課題を解決するための手段】

【0037】

上記課題を解決するため、本発明のフォールトトレラントシステムは、プライマリハイパーバイザが稼働するプライマリマシン上に形成され、前記プライマリハイパーバイザからの外部割り込みに基づく擬似割り込みをプライマリゲストOSに入力するプライマリ仮想マシンと、セカンダリハイパーバイザが稼働するセカンダリマシン上に形成され、前記プライマリ仮想マシンから伝達された前記擬似割り込みのタイミング情報に基づいて、擬似割り込みをセカンダリゲストOSに入力するセカンダリ仮想マシンと、を備えたフォールトトレラントシステムであって、前記プライマリ仮想マシンは、前記セカンダリ仮想マシンから、動作性能情報を収集し、前記プライマリ仮想マシンの動作能力と、前記セカンダリ仮想マシンの動作能力が等しくなるように、前記プライマリ仮想マシンの動作能力と前記セカンダリ仮想マシンの動作能力とを設定することを特徴とする。

40

ここで、前記プライマリ仮想マシンは、設定した前記セカンダリ仮想マシンの動作能力

50

を前記セカンダリ仮想マシンに通知し、前記セカンダリ仮想マシンは、通知された動作能力に従って動作を行なうことができる。

また、前記動作能力は、動作周波数を含み、前記プライマリ仮想マシンは、前記プライマリ仮想マシンが使用可能な動作周波数と、前記セカンダリ仮想マシンが使用可能な動作周波数とを収集して、前記プライマリ仮想マシンの動作周波数と、前記セカンダリ仮想マシンの動作周波数とを設定することができる。

また、前記プライマリ仮想マシンは、前記プライマリ仮想マシンが使用するデバイスの動作仕様と、前記セカンダリ仮想マシンが使用するデバイスの動作仕様とが異なる場合に、前記プライマリ仮想マシンが使用するデバイスに対する処理が、前記セカンダリ仮想マシンが使用するデバイスに対する処理と同じになるようにエミュレートすることができる。

10

このとき、前記プライマリ仮想マシンは、前記デバイスの動作仕様の違いが所定の機能の有無の場合は、前記所定の機能を無効にして、前記プライマリ仮想マシンが使用するデバイスに対する処理が、前記セカンダリ仮想マシンが使用するデバイスに対する処理と同じになるようにエミュレートし、前記デバイスの動作仕様の違いが能力あるいは容量の場合は、動作に相違が生じる処理を前記プライマリゲストOSの外部で実行させることができる。

また、前記セカンダリ仮想マシンは、前記プライマリ仮想マシンが使用するデバイスの動作仕様と、前記セカンダリ仮想マシンが使用するデバイスの動作仕様とが異なる場合に、前記セカンダリ仮想マシンが使用するデバイスに対する処理が、前記プライマリ仮想マシンが使用するデバイスに対する処理と同じになるようにエミュレートすることができる。

20

あるいは、前記セカンダリ仮想マシンは、前記デバイスの動作仕様の違いが所定の機能の有無の場合は、前記所定の機能を無効にして、前記セカンダリ仮想マシンが使用するデバイスに対する処理が、前記プライマリ仮想マシンが使用するデバイスに対する処理と同じになるようにエミュレートし、前記デバイスの動作仕様の違いが能力あるいは容量の場合は、動作に相違が生じる処理を前記セカンダリゲストOSの外部で実行させることができる。

また、前記セカンダリ仮想マシンは、前記動作能力の設定後であっても、前記プライマリ仮想マシンより実行速度が遅い場合には、遅れの度合いを取得し、取得した値が予め定められた遅延許容時間を超えると、同期初期化処理を行なうことができる。

30

このとき、前記セカンダリ仮想マシンは、前記取得した遅れの度合いが、想定される量より大きい場合に、故障として検出することができる。

また、前記セカンダリ仮想マシンは、前記動作能力の設定後であっても、前記プライマリ仮想マシンより実行速度が遅い場合には、前記プライマリゲストOSのアイドル時に、擬似割り込み入力タイミングを前倒しすることで、前記プライマリ仮想マシンに対する遅れを取り戻すことができる。

また、前記プライマリ仮想マシンは、前記動作能力の設定後であっても、前記セカンダリ仮想マシンより実行速度が速い場合には、前記セカンダリ仮想マシンの実行速度に合わせて、ウェイトを入れながら動作を行なうことができる。

40

また、前記セカンダリ仮想マシンは、前記動作能力の設定後であっても、前記プライマリ仮想マシンより実行速度が速い場合には、前記プライマリ仮想マシンの実行速度に合わせて、ウェイトを入れながら動作を行なうことができる。

また、前記プライマリ仮想マシンおよび前記セカンダリ仮想マシンは、それぞれのストレージデバイスの内容を、それぞれの仮想マシンに割り当てられたメモリ上に展開し、ストレージデバイスアクセスをメモリアクセスに置き換えることができる。

このとき、前記プライマリ仮想マシンおよび前記セカンダリ仮想マシンは、前記メモリアクセスを行なうと、アクセス完了の外部割り込みを待つことなく、それぞれのゲストOSに擬似割り込みを入力することができる。

また、前記プライマリ仮想マシンおよび前記セカンダリ仮想マシンは、それぞれのスト

50

レージデバイスの内容を、それぞれのゲストOSで設定されたRAMディスクに展開し、ストレージデバイスアクセスをRAMディスクに置き換えることができる。

【発明の効果】

【0038】

本発明によれば、2台のコンピュータ上に形成された2台の仮想マシンを並列同期動作させるフォールトトレラントシステムにおいて、それぞれのコンピュータのハードウェアの差異を適切に調整することができる。

【図面の簡単な説明】

【0039】

【図1】本実施形態に係るフォールトトレラントシステムの構成を示すブロック図である。

10

【図2】プライマリ仮想マシンとセカンダリ仮想マシンとの並列同期動作の概要について説明するフローチャートである。

【図3】分岐回数と実行中断位置とによる擬似割り込み入力位置特定を説明する図である。

【図4】CPUアーキテクチャの詳細が判別している場合の擬似割り込みタイミング調整動作の具体的な手順について説明するフローチャートである。

【図5】CPUアーキテクチャの詳細が不明の場合の擬似割り込みタイミング調整動作の具体的な手順について説明するフローチャートである。

【図6】仮想マシンを利用した従来のフォールトトレラントシステムの構成を示すブロック図である。

20

【図7】同期情報を受信したセカンダリ仮想マシンが、同期タイミング情報に従ってセカンダリゲストOSの実行を中断させる動作を説明するフローチャートである。

【発明を実施するための形態】

【0040】

本発明の実施の形態について図面を参照して説明する。図1は、本実施形態に係るフォールトトレラントシステムの構成を示すブロック図である。本図に示すように、フォールトトレラントシステム10は、ネットワーク接続されたプライマリマシン100とセカンダリマシン200とを備えている。プライマリマシン100、セカンダリマシン200は、サーバコンピュータ、パーソナルコンピュータ等の汎用的なコンピュータを用いて構成することができる。

30

【0041】

プライマリマシン100は、物理的なコンピュータ環境であるプライマリハードウェア110上でプライマリハイパーバイザ120が動作しており、プライマリハイパーバイザ120が動作することにより、プライマリ仮想マシン130が形成されている。プライマリ仮想マシン130では、プライマリゲストOS140が動作しており、プライマリゲストOS140上で、アプリケーション150が実行されている。

【0042】

プライマリハードウェア110は、CPU111、メモリ112、タイマ113、ネットワークインタフェースカード(NIC)114、ストレージ等の各種デバイス115を備えている。

40

【0043】

プライマリ仮想マシン130は、プライマリハードウェア110の一部が割り当てられ、外部との入出力を行なって動作する仮想的なコンピュータ環境である。プライマリ仮想マシン130は、プライマリハイパーバイザ120により管理され、プライマリハイパーバイザ120は、複数個のプライマリ仮想マシン130を管理することができる。

【0044】

セカンダリマシン200も同様に、物理的なコンピュータ環境であるセカンダリハードウェア210上でセカンダリハイパーバイザ220が動作しており、セカンダリハイパーバイザ220が動作することにより、セカンダリ仮想マシン230が形成されている。セ

50

カンダリ仮想マシン230では、セカンダリゲストOS240が動作しており、セカンダリゲストOS240上で、アプリケーション250が実行されている。

【0045】

セカンダリハードウェア210は、CPU211、メモリ212、タイマ213、ネットワークインタフェースカード(NIC)214、ストレージ等の各種デバイス215を備えている。

【0046】

セカンダリ仮想マシン230は、セカンダリハードウェア210の一部が割り当てられ、プライマリ仮想マシン130に同期して動作する仮想的なコンピュータ環境である。セカンダリ仮想マシン230は、セカンダリハイパーバイザ220により管理され、セカンダリハイパーバイザ220は、複数個のセカンダリ仮想マシン230を管理することができる。

10

【0047】

フォールトトレラントシステム10では、従来と同様に、外部との入出力のやり取りはプライマリ仮想マシン130のみが行ない、プライマリ仮想マシン130が先行して動作を行なうものとする。

【0048】

本実施形態において、プライマリ仮想マシン130は、擬似割り込み生成部131と同期情報生成部135とを備え、セカンダリ仮想マシン230は、擬似割り込み入力変換部231と同期補正部235とを備えている。

20

【0049】

プライマリ仮想マシン130の擬似割り込み生成部131は、プライマリハイパーバイザ120から実際に入力された外部割り込みに対して、割り込みの頻度やタイミングを調整した擬似割り込みを生成して、プライマリゲストOS140に出力する。このため、擬似割り込み生成部131は、割り込みの頻度を調整する割り込み頻度調整部132と、割り込みタイミングを調整する割り込みタイミング調整部133とを備えている。

【0050】

プライマリ仮想マシン130の同期情報生成部135は、タイミング調整され、プライマリゲストOS140に入力された擬似割り込みのタイミングを、実行中断位置と分岐実行回数とで特定する。このため、同期情報生成部135は、擬似割り込みが入力されるまでのプライマリゲストOS140におけるプログラムの分岐回数をカウントする分岐実行回数生成部136を備えている。

30

【0051】

また、同期情報生成部135は、プライマリ仮想マシン130とセカンダリ仮想マシン230の並列同期動作に先立ち、セカンダリ仮想マシン230から動作性能情報を収集する。そして、プライマリハードウェア110とセカンダリハードウェア210とで、ハードウェア構成の差異がある場合等に、速度差を補正するため等の同期補正情報を生成し、セカンダリ仮想マシン230の同期補正部235に伝達する。また、必要に応じてプライマリ仮想マシン130の動作の調整を行なう。

【0052】

セカンダリ仮想マシン230の擬似割り込み入力変換部231は、プライマリ仮想マシン130から送られた同期情報の同期タイミング情報に従って、セカンダリゲストOS240を擬似割り込み入力位置まで進めて中断させ、擬似割り込みを入力する。擬似割り込み入力変換部231は、セカンダリゲストOS240を擬似割り込み入力位置まで進めて中断させるため、指定分岐回数実行部232と、指定位置コンテキスト切り替え部233とを備えている。

40

【0053】

セカンダリ仮想マシン230の同期補正部235は、並列同期動作に先立ち、プライマリ仮想マシン130の同期情報生成部135から伝達された同期補正情報に基づいて、セカンダリ仮想マシンの動作の調整を行なう。また、必要に応じて、並列同期動作時にプラ

50

イマリ仮想マシン 130 とセカンダリ仮想マシン 230 との速度差の調整を行なう。

【0054】

(並列同期動作の概要)

次に、上記構成のフォールトトレラントシステム 10 における、プライマリ仮想マシン 130 とセカンダリ仮想マシン 230 の実行状態を同期させる並列同期動作の概要について図 2 のフローチャートを参照して説明する。具体的には、プライマリハイパーバイザ 120 が、プライマリハードウェア 110 から外部割り込みを受けた場合の一連の動作について説明する。

【0055】

プライマリハイパーバイザ 120 が、プライマリハードウェア 110 から外部割り込みを受けると (S101)、プライマリ仮想マシン 130 に外部割り込みの入力を行なう (S102)。

【0056】

プライマリハイパーバイザ 120 からの外部割り込みが入力されたプライマリ仮想マシン 130 は、プライマリゲスト OS 140 に擬似割り込みの入力を行なう。このとき、入力されたすべての外部割り込みに対して擬似割り込みの入力を行なうのではなく、割り込み頻度調整部 132 が、擬似割り込みの頻度を調整する (S103)。

【0057】

仮想マシンを利用したフォールトトレラントシステム 10 では、プライマリゲスト OS 140 に擬似割り込みが入力されると、実行状態を同期させるためにセカンダリゲスト OS 240 にも同じタイミングで擬似割り込みを入力する処理を行なわなくてはならない。このため、プライマリ仮想マシン 130 で頻繁に擬似割り込みが発生すると実行状態同期のための処理負荷が増え、本来のサービスに割り当てる CPU 時間や通信帯域等の資源が減少してしまう。

【0058】

そこで、本実施形態のフォールトトレラントシステム 10 では、割り込み頻度調整部 132 が、擬似割り込みの頻度を調整することで、実行状態同期のための処理負荷等を抑制し、本来のサービスのための処理能力低下を防止するようにしている。割り込み頻度調整部 132 が行なう割り込み頻度調整動作の詳細については後述する。

【0059】

割り込み頻度調整部 132 が擬似割り込みの頻度を調整した結果 (S103)、擬似割り込みを行なうと判定した場合 (S104: Yes) は、プライマリゲスト OS 140 に擬似割り込みの入力を行なう。擬似割り込みを行なわないと判定した場合 (S104: No) は、本処理を終了し、プライマリハイパーバイザ 120 からの次の外部割り込みの入力を待つ。

【0060】

プライマリゲスト OS 140 に擬似割り込みの入力を行なう場合は、セカンダリゲスト OS 240 でも確実に同じ位置で中断できるような位置でプライマリゲスト OS 140 を中断させて、擬似割り込みの入力を行なう必要がある。

【0061】

例えば、外部割り込み等の仮想マシンコンテキスト切り替えイベントにより、プライマリゲスト OS 140 が中断してプライマリ仮想マシン 130 に制御が移ったときに、プライマリゲスト OS 140 の中断位置がクリティカルセクションである場合が起こり得る。なお、クリティカルセクションは、単一のリソースに対して、複数の処理が同時期に実行されると、破綻をきたす部分であり、プログラムにより、ブレーク命令無効化等の排他制御が行なわれる部分である。

【0062】

このような場合、その中断位置でプライマリゲスト OS 140 に擬似割り込みを入力すると、クリティカルセクション内であるため、セカンダリゲスト OS 240 では、同じ位置にブレーク命令を埋め込んだとしても、同じ位置で中断することができず、クリティカ

10

20

30

40

50

ルセクションを抜け出してから中断する。このため、擬似割り込みの位置がずれてしまい、実行状態の同期がずれてしまう。

【0063】

そこで、本実施形態のフォールトトレラントシステム10では、割り込みタイミング調整部133が、擬似割り込みタイミングの調整を行ない(S105)、クリティカルセクションを避けてプライマリゲストOS140を中断させることで、セカンダリゲストOS240でも確実に同じ位置で中断できるようにしている。割り込みタイミング調整部133が行なう擬似割り込みタイミング調整動作の詳細については後述する。

【0064】

擬似割り込みタイミングが調整されると、プライマリ仮想マシン130はプライマリゲストOS140に擬似割り込みの入力を行なう(S106)。そして、同期情報生成部135が同期情報を生成する(S107)。

10

【0065】

従来、同期情報は、擬似割り込みの識別情報と同期タイミング情報とを含み、同期タイミング情報は、プログラムカウンタが示す実行中断位置とパフォーマンスカウンタで測定したCPU実行命令数とを含んでいた。

【0066】

しかしながら、CPU固有の特殊な命令によっては、パフォーマンスカウンタのカウンタ漏れ等が発生して、プライマリゲストOS140が擬似割り込み入力までに実行したCPU実行命令数を正確にカウントできない場合がある。例えば、Intel社CPUのPentium(登録商標)では、外部割り込みの発生によりREP命令の実行が中断された場合にREP命令の実行がカウントされないため、CPU実行命令数カウンタは実際の実行命令数よりも少なくなる。

20

【0067】

このように、CPU実行命令数はカウント値が不正確の場合があるため、本実施形態のフォールトトレラントシステム10は、同期タイミング情報にCPU実行命令数を用いないようにしている。

【0068】

しかしながら、一般に、プログラムはループ文やジャンプ文、条件分岐文を多く含んでいるため、同じ命令を何度も通過する。このため、実行中断位置だけでは、擬似割り込み入力のタイミングを特定することができない。

30

【0069】

そこで、本実施形態のフォールトトレラントシステム10は、同期タイミング情報に実行中断位置と分岐回数とを含めるようにしている。分岐回数は、前回の擬似割り込みから今回の擬似割り込みまでに実行した分岐命令の数であり、命令コードの実行順序が順次実行から変更された回数である。

【0070】

すなわち、プライマリゲストOS140で、分岐回数をカウントしておけば、セカンダリゲストOS240において、その回数分の分岐の後に最初に通過する実行中断位置を、擬似割り込み入力位置と特定することができる。

40

【0071】

例えば、図3(a)に示すように、命令アドレスa1、a2、a3、...に命令コードCodeA、CodeB、CodeC、...が割り当てられており、CodeCからCodeGまでが3回ループされる場合に、プライマリゲストOS140において、プログラム実行順16番目の「a6」で、擬似割り込みが入力されたとする。この同期タイミングは、分岐回数のカウント値「2」、実行中断位置「a6」と表わすことができる。

【0072】

この場合、セカンダリ仮想マシン230において、分岐回数を2回カウントし、最初の実行中断位置「a6」でセカンダリゲストOS240を中断させることで、プライマリゲストOS140と同じタイミングで擬似割り込みを入力することができる。

50

【 0 0 7 3 】

このような分岐回数のカウント処理を行なうために、同期情報生成部 1 3 5 は、分岐実行回数生成部 1 3 6 を備えている。命令コードレベルでは、ループ文やジャンプ文、条件分岐文は、すべて分岐命令に置き換えられるため、分岐実行回数生成部 1 3 6 は、実行中のアプリケーション 1 5 0 において、プライマリゲスト OS 1 4 0 が実行した分岐命令の数をカウントすることにより分岐回数のカウントを行なうことができる。分岐回数のカウントは、例えば、パフォーマンスカウンタを利用することができる。

【 0 0 7 4 】

図 2 の説明に戻って、プライマリ仮想マシン 1 3 0 は、擬似割り込みの識別情報と同期タイミング情報とを含んだ同期情報を生成するとセカンダリ仮想マシン 2 3 0 に伝達する (S 1 0 8) 。

10

【 0 0 7 5 】

同期情報を受信したセカンダリ仮想マシン 2 3 0 は、同期タイミング情報に従ってセカンダリゲスト OS 2 4 0 の実行を中断させるために、まず、指定分岐回数実行部 2 3 2 が、セカンダリゲスト OS 2 4 0 の実行時に、分岐回数をカウントし、同期タイミング情報の分岐回数が示す値になるとセカンダリゲスト OS 2 4 0 を中断させる。すなわち、指定された分岐回数までセカンダリゲスト OS 2 4 0 を実行させる (S 1 0 9) 。

このため、指定された分岐回数に達するまでは、実行中断位置や CPU 実行命令数を監視する必要がなく、処理負荷が軽減する。指定分岐回数実行部 2 3 2 が行なう指定分岐回数実行動作の詳細については後述する。

20

【 0 0 7 6 】

そして、指定位置コンテキスト切り替え部 2 3 3 が、同期タイミング情報の実行中断位置が示す位置で、セカンダリゲスト OS 2 4 0 の実行を中断させる (S 1 1 0) 。

【 0 0 7 7 】

セカンダリゲスト OS 2 4 0 の実行を中断させる際に、従来のようにブレーク命令を埋め込んでブレークポイントを設定する手法を用いると、仮想マシン (1 3 0 、 2 3 0) 上で、ゲスト OS (1 4 0 、 2 4 0) やアプリケーション (1 5 0 、 2 5 0) のデバッグを行なっている場合に、デバッグのブレークポイントと干渉し、相互の処理に悪影響を与えてしまう。同期状態の仮想マシン (1 3 0 、 2 3 0) 上でアプリケーション (1 5 0 、 2 5 0) 動作の検証等を行なうニーズもあるため、ブレークポイントの干渉は避けることが望ましい。

30

【 0 0 7 8 】

また、仮想化支援機能を持っていない CPU の場合は、ブレークポイントの設定でセカンダリゲスト OS 2 4 0 からセカンダリ仮想マシン 2 3 0 に処理を遷移することができない。

【 0 0 7 9 】

そこで、本実施形態のフォールトトレラントシステム 1 0 では、指定位置コンテキスト切り替え部 2 3 3 が、ブレーク命令を用いずに、コンテキスト切り替えを行なうことにより、セカンダリゲスト OS 2 4 0 の実行を中断させるようにしている。指定位置コンテキスト切り替え部 2 3 3 が行なう指定位置コンテキスト切り替え動作の詳細については後述する。

40

【 0 0 8 0 】

セカンダリ仮想マシン 2 3 0 は、セカンダリゲスト OS 2 4 0 を中断させると、同期情報の擬似割り込み識別情報に従って擬似割り込みの設定を行なって、セカンダリゲスト OS 2 4 0 を再開させる。これにより、セカンダリゲスト OS 2 4 0 に、プライマリゲスト OS 1 4 0 と同じタイミングで擬似割り込みが入力され (S 1 1 1) 、プライマリ仮想マシン 1 3 0 とセカンダリ仮想マシン 2 3 0 との実行状態が同期する。

【 0 0 8 1 】

なお、プライマリマシン 1 0 0 、セカンダリマシン 2 0 0 のいずれかでハードウェアの故障が生じると、実行状態の同期が乱れる。本実施形態では、パフォーマンスカウンタの

50

カウント漏れ等の影響を受けずに、正確に実行状態の同期を行なうことができる。このため、ハードウェアが正常であれば、実行状態の同期が乱れることはないので、実行状態の同期の乱れを検出することでハードウェアの故障を発見することができる。

【0082】

従来は、実行状態の同期の乱れが故障によるものか、実行命令数のカウント漏れの影響によるか判別できないため、出力データの照合が必要であったが、本実施形態のフォールトトレラントシステム10では、実行状態の同期の乱れにより故障を検出することができるため、故障検出のための出力データの照合が不要となる。これにより、出力データの照合のためのCPU処理負荷や通信負荷を削減することができる。また、故障の発生を迅速に検出できる。

10

【0083】

なお、実行状態の同期の乱れの検出は、例えば、セカンダリ仮想マシン230の擬似割り込み入力変換部231において、同期タイミング情報で指定された実行中断位置を通過する前に、分岐回数が、同期タイミング情報で指定された回数を超過した場合に同期が乱れたと判定することができる。

【0084】

(各ブロックの動作の詳細)

次に、フォールトトレラントシステム10の各ブロックの動作について第1実施例、第2実施例、第3実施例により詳細に説明する。

【0085】

20

従来のフォールトトレラントシステムは、プライマリマシンのCPUとセカンダリマシンのCPUとが同じアーキテクチャであることを前提としていた。これは、従来の同期合わせの手法が、CPUアーキテクチャやCPU動作仕様に依存しており、CPUの種類が異なれば、双方の実行命令を合わせることができず、同期ずれが発生する場合があったからである。

【0086】

しかしながら、本発明は、第1実施例、第2実施例で説明するように、プライマリマシン100のCPU111とセカンダリマシン200のCPU211とが同じアーキテクチャのみならず、異なるアーキテクチャの場合でも適用することができる。

【0087】

30

一般に、長年に渡ってフォールトトレラントシステムを運用する場合には、故障や経年劣化による機器の交換が余儀なくされるが、この際に、従前の機器と同一仕様の機器を用意することは困難である。このため、本実施形態のフォールトトレラントシステム10は、CPUアーキテクチャが異なる場合であっても、並列同期動作を継続することができる仕組みを備えるようにしている。

【0088】

<第1実施例>

第1実施例として、プライマリマシン100のCPU111とセカンダリマシン200のCPU211とが同じアーキテクチャの場合の動作について説明する。ここで、CPUアーキテクチャが同じとは、CPU命令セットに互換性があることを意味する。

40

【0089】

この場合、プライマリ仮想マシン130とセカンダリ仮想マシン230とは、CPU命令レベルの同期合わせによる並列同期動作を行なう。このため、両仮想マシン(130、230)において、ゲストOS(140、240)とアプリケーション(150、250)のプログラムは、通常どおり、CPU(111、211)に対してネイティブのCPU命令コードにコンパイルされたコード形式を使用する。この並列同期動作を、「CPU命令同期モード」と称するものとする。

【0090】

(割り込み頻度調整動作)

まず、割り込み頻度調整部132が行なう割り込み頻度調整動作(S103)について

50

説明する。プライマリ仮想マシン 130 の割り込み頻度調整部 132 は、すべての外部割り込みに対して擬似割り込みの入力を行なうのではなく、アプリケーション 150 に要求される性能や精度、許容通信ディレイ等を満たす範囲で、擬似割り込みの頻度を調整する。

【0091】

一般に、外部割り込みは、タイマ 113 等による定周期割り込みと、ネットワークインタフェースカード (NIC) 114、デバイス 115 等による非定周期割り込みとに区別することができるが、割り込み頻度調整部 132 は、頻度の調整を、定周期割り込みと非定周期割り込みとで分けて行なう。

【0092】

定周期の外部割り込みの場合は、間引きを行なって、定周期の外部割り込み数回に 1 回の割合で擬似割り込みをプライマリゲスト OS 140 に入力するようにする。

【0093】

例えば、定周期のタイマ割り込みをプライマリゲスト OS 140 に入力することによって、プライマリゲスト OS 140 は正確な時刻を計時することができるが、アプリケーション 150 で要求される時刻の精度が高くない場合は、タイマ割り込みの入力を間引いても影響はない。

【0094】

例えば、プライマリハイパーバイザ 120 からのタイマ割り込みの周期が 1ms、アプリケーション 150 で要求される時刻の精度が 100ms の場合、擬似割り込みの周期を 10ms に間引いても 100ms の精度を維持することができる。このとき、擬似割り込みの頻度は 1/10 に減少するため、実行状態の同期処理も 1/10 に減少する。

【0095】

なお、アプリケーション 150 の要求精度に応じて、予め間引く割合を設定しておいてもよいし、パラメータで調整できるようにしてもよい。また、定周期の外部割り込みに対して一律に間引く割合を設定してもよいし、定周期の外部割り込みの要因ごとに間引く割合を設定してもよい。

【0096】

非定周期の外部割り込みの場合は、処理途中に入力される外部割り込みに対する擬似割り込みを省き、処理が完了した最後の外部割り込みのみを擬似割り込みとしてプライマリゲスト OS 140 に入力する。

【0097】

例えば、NIC 114 からのネットワークの受信割り込みは、NIC 114 にデータが到着する度に発生する。このため、データとして意味のある単位を受信する間に、複数回の外部割り込みが発生する。従来は、外部割り込みの発生ごとにプライマリゲスト OS 140 に擬似割り込みを行なって、NIC 114 からプライマリゲスト OS 140 に受信データを渡していたが、割り込み頻度調整部 132 は、データ受信途中の外部割り込みに対する擬似割り込みを省き、データ受信完了後の 1 回の外部割り込みに対して擬似割り込みをプライマリゲスト OS 140 に入力する。これにより、実行状態の同期処理の回数が減少するとともに、まとまったサイズのデータを 1 度に処理することが可能となる。

【0098】

なお、本手法を適用できる非定周期の外部割り込みとしては、ネットワークの送受信割り込み、ストレージの読み書き割り込み、シリアル通信の送受信割り込み、アナログ/デジタル入出力ボード等の各種 I/O デバイスの読み書き割り込み等が挙げられる。

【0099】

(擬似割り込みタイミング調整動作)

次に、割り込みタイミング調整部 133 が行なう擬似割り込みタイミング調整動作 (S105) について説明する。擬似割り込みタイミング調整は、セカンダリゲスト OS 240 が同じ位置で中断できるように、プライマリゲスト OS 140 のクリティカルセクションを避けて擬似割り込みを入力するための調整である。

10

20

30

40

50

【 0 1 0 0 】

擬似割り込みタイミング調整動作は、CPUアーキテクチャの詳細が判別している場合と、CPUアーキテクチャの詳細が不明の場合とで処理が異なる。なお、CPUアーキテクチャの詳細が判別している場合とは、中断位置のCPU命令や状態フラグ等から、中断位置がクリティカルセクションかどうかを判別できる場合をいう。

【 0 1 0 1 】

まず、CPUアーキテクチャの詳細が判別している場合の具体的な動作手順について、図4のフローチャートを参照して説明する。

【 0 1 0 2 】

プライマリ仮想マシン130は、プライマリハイパーバイザ120から外部割り込みの入力を受けると(S201)、プライマリゲストOS140の実行の中断を待つ(S202)。

【 0 1 0 3 】

プライマリゲストOS140の実行が中断し、プライマリ仮想マシン130に処理が遷移すると、プライマリ仮想マシン130は、プライマリ仮想マシン130の状態を取得する(S203)。取得する状態は、クリティカルセクションかどうかを判別するための情報であり、例えば、CPU111の割り込みマスクの設定情報、フラグレジスタの内容等である。

【 0 1 0 4 】

そして、取得した情報を元に、中断位置がクリティカルセクションかどうかを判断する(S204)。例えば、割り込みマスクで割り込みの入力が禁止されていないこと、ブレーク命令の設定を無効化するフラグがフラグレジスタに設定されていないこと等で、中断位置がクリティカルセクションでないとは判断することができる。

【 0 1 0 5 】

中断位置がクリティカルセクションでない場合(S204:No)は、中断位置を擬似割り込みの入力タイミングに設定する(S205)。

【 0 1 0 6 】

中断位置がクリティカルセクションの場合(S204:Yes)は、プライマリ仮想マシン130の状態を確認しながらプライマリゲストOS140の実行を進める(S206)。このとき、例えば、ステップ実行を行なうことで、1命令ごとに状態を確認する(S203)。そして、クリティカルセクションかどうかを判断し(S204)、クリティカルセクションを抜け出すと(S204:No)、擬似割り込みの入力タイミングに設定する(S205)。

【 0 1 0 7 】

次に、CPUアーキテクチャの詳細が不明の場合の具体的な動作手順について、図5のフローチャートを参照して説明する。CPUアーキテクチャの詳細が不明の場合、プライマリ仮想マシン130の状態を取得しても、クリティカルセクションかどうかを判別できないため、中断のためのイベントを設定し、実際に中断できた位置を、擬似割り込みを入力する中断位置とする。

【 0 1 0 8 】

プライマリ仮想マシン130は、プライマリハイパーバイザ120から外部割り込みの入力を受けたら(S301)、プライマリゲストOS140の実行の中断を待つ(S302)。

【 0 1 0 9 】

プライマリゲストOS140の実行が中断し、プライマリ仮想マシン130に処理が遷移すると、プライマリ仮想マシン130は、プライマリゲストOS140に、所定命令実行後、例えば、1命令実行後に中断する旨のイベントの設定を行なう(S303)。

【 0 1 1 0 】

そして、ステップ実行でプライマリゲストOS140の実行を再開し(S304)、1命令実行ごとに中断の要因を取得する(S305)。これは、中断の要因が、必ずしも設

10

20

30

40

50

定した中断イベントによるものとは限られないからである。

【0111】

中断の要因が、設定した中断イベントによるものの場合（S306：Yes）は、中断位置がクリティカルセクションではなく、セカンダリゲストOS240でも同じ位置で中断できるものと判断し、擬似割り込みの入力タイミングに設定する（S307）。

【0112】

中断の要因が、設定した中断イベントによるものでない場合（S306：No）は、中断位置がクリティカルセクションであり、セカンダリゲストOS240では同じ位置で中断できないものと判断し、再度、1命令実行後に中断する旨のイベントの設定を行なって（S303）、プライマリゲストOS140の実行再開（S304）以降の処理を繰り返して、クリティカルセクションを抜け出すと（S306：Yes）、擬似割り込みの入力タイミングに設定する（S307）。

10

【0113】

（同期情報生成動作）

次に、同期情報生成部135が行なう同期情報の生成動作（S107）について説明する。プライマリマシン100のCPU111とセカンダリマシン200のCPU211とが同じアーキテクチャの場合、同期情報の同期タイミング情報における実行中断位置は、CPU命令コードの実行停止位置で示され、プログラムカウンタ値を用いることができる。

【0114】

また、同期タイミング情報における分岐実行回数は、前回の同期タイミング情報生成時からプライマリゲストOS140で実行された分岐回数で示される。分岐回数は、CPU命令コードレベルで実行した分岐命令をカウントすることで得ることができる。

20

【0115】

（指定分岐回数実行動作）

次に、指定分岐回数実行部232が行なう指定分岐回数実行動作（S109）について説明する。指定分岐回数実行動作は、セカンダリ仮想マシン230が、同期タイミング情報の分岐回数で示された回数でセカンダリゲストOS240を中断させる動作である。

【0116】

指定分岐回数実行部232は、CPU命令コードレベルで実行した分岐命令をカウントし、カウント値が指定された分岐回数に達すると、その旨をセカンダリ仮想マシン230に通知する。これにより、セカンダリ仮想マシン230にコンテキストが切り替えられ、セカンダリゲストOS240の実行が中断する。

30

【0117】

なお、CPU211のパイプライン機能等の高速化により、指定された分岐回数の実行直後に停止できないCPUアーキテクチャの場合には、指定された分岐回数より少ない値に設定し、指定された分岐回数で確実に中断できるようにする。

【0118】

（指定位置コンテキスト切り替え動作）

次に、指定位置コンテキスト切り替え部233が行なう指定位置コンテキスト切り替え動作（S110）について説明する。指定位置コンテキスト切り替え動作は、セカンダリゲストOS240を指定分岐回数で中断した後、指定された実行中断位置で、ブレーク命令を用いずに、セカンダリ仮想マシン230にコンテキストを切り替える動作である。

40

【0119】

セカンダリ仮想マシン230に専用のコンテキスト切り替え命令が用意されている場合は、セカンダリゲストOS240の実行中断位置に、コンテキスト切り替え命令を埋め込むようにする。コンテキスト切り替え命令は、例えば、Intel社CPUのVMCALL命令である。

【0120】

セカンダリゲストOS240が、コンテキスト切り替え命令を実行すると、コンテキス

50

トがセカンダリゲストOS 240からセカンダリ仮想マシン230に切り替わり、セカンダリゲストOS 240の実行が中断される。

【0121】

専用のコンテキスト切り替え命令が用意されていない場合は、実行中断位置で擬似割り込みやCPU例外のイベントを発生させてコンテキスト切り替えが行なわれるように設定する。コンテキスト切り替えが行なわれると、コンテキストがセカンダリゲストOS 240からセカンダリ仮想マシン230に切り替わり、セカンダリゲストOS 240の実行が中断される。

【0122】

いずれの場合も、セカンダリゲストOS 240の実行が中断されると、セカンダリ仮想マシン230は、擬似割り込みの設定をして、セカンダリゲストOS 240の実行を再開させ、そのタイミングで擬似割り込みがセカンダリゲストOS 240に入力される。

10

【0123】

<第2実施例>

第2実施例として、プライマリマシン100のCPU111とセカンダリマシン200のCPU211とが異なるアーキテクチャの場合の動作について説明する。ここで、CPUアーキテクチャが同じとは、CPU命令セットに互換性がないことを意味する。

【0124】

この場合、プライマリ仮想マシン130とセカンダリ仮想マシン230とは、中間言語レベルの同期合わせによる並列同期動作を行なう。ここで、「中間言語」は、「中間コード」、「バイトコード」等とも称され、高水準言語のソースコードと機械語のコードの中間にあたる中間表現の言語である。この並列同期動作を「中間言語同期モード」と称するものとする。

20

【0125】

両仮想マシン(130、230)間で、CPU命令セットの互換性がないため、ゲストOS(140、240)とアプリケーション(150、250)のプログラムは、ネイティブのCPU命令コードの代わりに、Java(登録商標)やC#等を変換した中間言語のプログラムを使用する。

【0126】

両仮想マシン(130、230)は、インタープリタやJITコンパイラ等の既存のCPU命令変換手法を用いて、中間言語を仮想マシンCPUに対してネイティブのCPU命令に変換して実行する。

30

【0127】

なお、プライマリマシン100のCPU111とセカンダリマシン200のCPU211とが同じアーキテクチャの場合も、中間言語同期モードによる並列同期動作を行なうことができるが、CPU命令同期モードによる並列同期動作の方が、実行速度が速くなる。

【0128】

次に、第2実施例における各ブロックの動作について説明するが、大部分は第1実施例のCPU命令同期モードと同様であるため、第1実施例と異なる処理についてのみ説明する。

40

【0129】

(擬似割り込みタイミング調整動作)

第1実施例のCPU命令同期モードにおいて、擬似割り込みタイミング調整動作(S105)は、CPUアーキテクチャの詳細が判別している場合と、CPUアーキテクチャの詳細が不明の場合とで処理が異なっていたが、第2実施例の中間言語同期モードでは、CPUアーキテクチャの詳細が判別している場合と同じ処理を行なう。

【0130】

すなわち、中断位置のCPU命令や状態フラグ等から、中断位置がクリティカルセクションかどうかを判別できるため、図4に示したフローチャートに従って、クリティカルセクションを避けてプライマリゲストOS 140の中断位置を設定し、擬似割り込みを入力

50

する。

【0131】

また、中間言語同期モードでは、中断位置のCPU命令や状態フラグに加え、中間言語レベルでクリティカルセクションかどうかを判断することができる。例えば、中間言語がJava（登録商標）のバイトコードの場合は、クリティカルセクションはsynchrononizedの排他制御を使って実装することが想定される。このため、実行済みの中間言語の命令を確認することで、プライマリ仮想マシン130は、中断したJava（登録商標）の中間言語がクリティカルセクションにいるかを把握できる。

【0132】

（同期情報の生成動作）

10

第1実施例のCPU命令同期モードでは、同期情報の生成動作（S107）において、同期タイミング情報に含まれる実行中断位置は、CPU命令コードの実行停止位置で示し、プログラムカウンタ値を用いた。また、同期タイミング情報に含まれる分岐実行回数は、前回の同期タイミング情報生成時から実行されたCPU命令コードの分岐回数で示した。

【0133】

これに対し、第2実施例の中間言語同期モードでは、実行中断位置は、擬似割り込みを入力した時点で実行した中間言語レベルの実行命令位置を用いる。また、分岐実行回数は、前回の同期タイミング情報生成時から実行された中間言語レベルの分岐回数を用いる。

【0134】

20

（指定位置コンテキスト切り替え動作）

第1実施例のCPU命令同期モードでは、指定位置コンテキスト切り替え動作（S110）において、セカンダリ仮想マシン230に専用のコンテキスト切り替え命令が用意されている場合は、セカンダリゲストOS240の実行中断位置に、コンテキスト切り替え命令を埋め込み、専用のコンテキスト切り替え命令が用意されていない場合は、実行中断位置で擬似割り込みやCPU例外のイベントを発生させてコンテキスト切り替えが行なわれるように設定した。第2実施例の中間言語同期モードでも、中間言語レベルの実行中断位置で、CPU命令変換により同様の設定を行なう。

【0135】

<第3実施例>

30

次に、第3実施例として、プライマリマシン100とセカンダリマシン200とのハードウェア構成に相違がある場合の並列同期動作について説明する。プライマリマシン100のプライマリハードウェア110とセカンダリマシン200のセカンダリハードウェア210は、CPU（111、211）、メモリ（112、212）、ストレージ等のデバイス（115、215）等、種々のハードウェアで構成されているが、双方のハードウェアの機能や性能は必ずしも一致するわけではない。

【0136】

例えば、CPU（111、211）のスペックが異なっていたり、メモリ（112、212）の容量が異なっている場合も想定される。また、フォールトトレラントシステム10の運用過程においてデバイス215が交換される等により、双方のデバイス（115、215）が異なっている場合もある。

40

【0137】

このため、プライマリ仮想マシン130とセカンダリ仮想マシン230とで、動作は同じだが、処理速度に差が生じる場合がある。上述のように、フォールトトレラントシステム10は、プライマリ仮想マシン130が先行して動作を行なうようにしているが、例えば、セカンダリ仮想マシン230の方が、処理速度が速いと、動作の逆転現象が生じ、プライマリ仮想マシン130からの同期情報に基づく実行状態の同期が行なえなくなる。プライマリ仮想マシン130の方が、処理速度が速いと、セカンダリ仮想マシン230の遅れが蓄積されていく。

【0138】

50

また、異なるデバイス(115、215)により、プライマリゲストOS140とセカンダリゲストOS240のドライバ処理が異なって、動作に相違が生じる結果、同期ずれが生じるおそれがある。

【0139】

そこで、フォールトトレラントシステム10では、プライマリハードウェア110とセカンダリハードウェア210との差異に基づく処理速度差や同期ずれを補正する仕組みを設けている。具体的には、プライマリ仮想マシン130の同期情報生成部135とセカンダリ仮想マシン230の同期補正部235とが同期補正動作を行なう。同期補正動作は、並列同期動作実行前の同期補正動作と、並列同期動作実行時の同期補正動作とに分けることができる。

10

【0140】

(並列同期動作実行前の同期補正動作)

プライマリ仮想マシン130の同期情報生成部135は、プライマリ仮想マシン130とセカンダリ仮想マシン230の並列同期動作に先立ち、セカンダリ仮想マシン230からセカンダリ仮想マシン230に関する動作性能情報を収集する。

【0141】

動作性能情報は、セカンダリ仮想マシン230におけるセカンダリゲストOS240に対するインタフェースの動作仕様や、セカンダリ仮想マシン230のプログラム実行速度の指標となるベンチマーク情報等を含んだ情報である。同期情報生成部135は、使用可能な範囲で複数項目の動作性能情報を取得するようにする。

20

【0142】

同期情報生成部135は、プライマリ仮想マシン130自身の動作性能情報も取得して、セカンダリ仮想マシン230の動作性能情報と対照する。そして、プライマリ仮想マシン130の並列同期動作における動作能力と、セカンダリ仮想マシン230の並列同期動作における動作能力とがほぼ等しくなるように、プライマリ仮想マシン130とセカンダリ仮想マシン230が使用する動作能力を設定する。

【0143】

動作能力は、例えば、動作周波数、メモリサイズ、キャッシュサイズ、TLB(Translation Look-aside Buffer)/EPT(Enhanced Page Table)仕様等の動作仕様、MIPS値やSPEC値等のベンチマーク情報である。中間言語同期モードの場合は、中間言語のベンチマーク値を利用することができる。例えば、中間言語としてJava(登録商標)のバイトコードを用いた場合は、ベンチマーク値としてSPECjbb、CaffeineMark等が挙げられる。

30

【0144】

そして、セカンダリ仮想マシン230が使用する動作能力を同期補正情報としてセカンダリ仮想マシン230の同期補正部235に通知する。セカンダリ仮想マシン230の同期補正部235は、通知された同期補正情報に従って、並列動作実行時におけるセカンダリ仮想マシン230の動作能力を設定する。

【0145】

動作周波数の調整による同期補正を例に説明する。なお、上述のように、プライマリマシン100とセカンダリマシン200とのハードウェア構成に相違があるときは、動作は同じだが、処理速度に差が生じる場合と、動作自体が異なる場合とがある。本例は、動作は同じだが、処理速度に差が生じる場合の同期補正の説明である。

40

【0146】

仮想マシン(130、230)の動作周波数の調整には、省電力のために用意された電力制御機能を利用することができる。例えば、標準の電力制御機能にはACPI(Advanced Configuration and Power Interface)を使用し、動作周波数を変更する機能には、Linux(登録商標)のcpufreqサブシステム等を使用することができる。この機能を利用して、双方の仮想マシン(130、230)は、使用可能な動作周波数を複数個用意しておく。

50

【0147】

そして、プライマリ仮想マシン130の同期情報生成部135は、セカンダリ仮想マシン230で使用可能な動作周波数の一覧をセカンダリ仮想マシン230から取得するとともに、プライマリ仮想マシン130自身で使用可能な動作周波数の一覧を取得して、両者がほぼ一致する動作周波数を選択する。選択した動作周波数はセカンダリ仮想マシン230の同期補正部235に同期補正情報として通知する。

【0148】

プライマリ仮想マシン130とセカンダリ仮想マシン230の並列同期動作時には、それぞれの仮想マシン(130、230)は、選択された動作周波数で動作するようにすることで、プライマリハードウェア110とセカンダリハードウェア210との動作周波数の差異に基づく処理速度差が解消される。

10

【0149】

なお、プライマリ仮想マシン130の使用可能な動作周波数と、セカンダリ仮想マシン230の使用可能な動作周波数が一致しない場合は、例えば、セカンダリ仮想マシン230の方を速い動作周波数とし、後述するように、セカンダリ仮想マシン230の動作にウェイトを入れるようにすることができる。

【0150】

中間言語同期モードの場合は、ベンチマーク値の違いによりプライマリ仮想マシン130とセカンダリ仮想マシン230の処理時間に差が生じる。従って、取得した双方のベンチマーク値に基づいて、動作周波数と同様の調整を行なうことによって処理速度差を補正する。

20

【0151】

次に、ハードウェア構成の相違により、プライマリ仮想マシン130とセカンダリ仮想マシン230とで動作自体が異なる場合の同期補正について説明する。

【0152】

例えば、CPUアーキテクチャが同じでもデバイス115等の動作仕様が異なる場合には、双方のゲストOS(140、240)が同じドライバ処理を実行するように変更することで、双方のゲストOS(140、240)が同じ動作を行ない、同期ずれを防ぐことができる。具体的には、並列同期動作の開始前に、プライマリ仮想マシン130またはセカンダリ仮想マシン230のデバイスエミュレータに動作仕様の違いを吸収する機能を組み込むようにする。

30

【0153】

この機能は、補正の対象によって下記の2種類に分類できる。

【0154】

(1) ハードウェア機能の有無の違いに対する補正

あるハードウェア機能を一方の仮想マシン(130、230)だけが持っている場合は、機能を持っている方の仮想マシン(130、230)がその機能を利用しないようにし、双方の仮想マシン(130、230)がそのハードウェア機能を持たない状態で動作する。

40

【0155】

例えば、FPU(浮動小数点演算装置)機能の有無が異なる場合、双方の仮想マシン(130、230)ともFPU機能を持たない状態で動作させる。このとき浮動小数点演算はゲストOS(140、240)内のエミュレーションコードで実行される。具体的には、浮動小数点演算実行時に例外が発生し、ゲストOS(140、240)の例外ハンドラがこの例外を捕捉してソフトウェアで演算をエミュレートする。

【0156】

あるいは、機能を持っていない方の仮想マシン(130、230)が、その機能をゲストOS(140、240)の外部でソフトウェアの処理によってエミュレートすることにより、違いを吸収するようにしてもよい。ゲストOS(140、240)の外部とは、仮想マシン(130、230)やハイパーバイザ(120、220)などである。

50

【 0 1 5 7 】

例えば、F P U機能を持っていない方の仮想マシン（ 1 3 0、 2 3 0 ）が、仮想マシン（ 1 3 0、 2 3 0 ）内のエミュレーションコードで浮動小数点演算を実行する。具体的には、浮動小数点演算実行時に発生する例外によって仮想マシンコンテキストへの切り替えを行ない、仮想マシン（ 1 3 0、 2 3 0 ）内のソフトウェアで演算をエミュレートした後、ゲストOSコンテキストに戻る。ゲストOS（ 1 4 0、 2 4 0 ）の中では例外を捕捉しないため、ゲストOS（ 1 4 0、 2 4 0 ）内のエミュレーションコードは実行されない。

【 0 1 5 8 】

（ 2 ）ハードウェア機能の能力や容量の違いに対する補正

ハードウェア機能の能力や容量が異なる場合は、動作に相違が生じる処理をゲストOS（ 1 4 0、 2 4 0 ）の外部で実行する。ゲストOS（ 1 4 0、 2 4 0 ）の外部とは、仮想マシン（ 1 3 0、 2 3 0 ）やハイパーバイザ（ 1 2 0、 2 2 0 ）、ハードウェア（ 1 1 0、 2 1 0 ）などである。

10

【 0 1 5 9 】

例えば、T L Bのキャッシュサイズが異なる場合、一方の仮想マシン（ 1 3 0、 2 3 0 ）のみがキャッシュにヒットせず、T L Bミスが発生する場合がある。そこで、T L Bミスに対するリカバリー処理を、ゲストOS（ 1 4 0、 2 4 0 ）の外部で実行する。I n t e l社C P Uの場合は、ハードウェアの機能でリカバリー処理を行なうことができるが、C P U（ 1 1 1、 2 1 1 ）にハードウェアのリカバリー機能がない場合は、仮想マシン（ 1 3 0、 2 3 0 ）やハイパーバイザ（ 1 2 0、 2 2 0 ）のソフトウェアでリカバリー処理を行なう。

20

【 0 1 6 0 】

中間言語同期モードの場合は、中間言語で表現されたゲストOSプログラムを仮想マシン（ 1 3 0、 2 3 0 ）が仮想メモリ上でインタープリタなどのソフトウェア機能を使って実行する。このため、動作仕様の違いによってゲストOS（ 1 4 0、 2 4 0 ）の動作に相違が生じることはない。

【 0 1 6 1 】

（ 並列同期動作実行時の同期補正動作 ）

並列同期動作の開始前の調整によってもプライマリ仮想マシン 1 3 0 とセカンダリ仮想マシン 2 3 0 の処理速度差が解消しきれない場合には、同期情報生成部 1 3 5、同期補正部 2 3 5 は、並列同期動作時に、以下のような同期補正を行なう。

30

【 0 1 6 2 】

まず、動作能力調整後のプライマリ仮想マシン 1 3 0 が動作能力調整後のセカンダリ仮想マシン 2 3 0 よりも高速である場合の同期補正について説明する。

【 0 1 6 3 】

同期情報生成部 1 3 5 は、セカンダリ仮想マシン 2 3 0 に通知する同期補正情報の動作速度情報として、プライマリ仮想マシン 1 3 0 の動作能力とするか、セカンダリ仮想マシン 2 3 0 の動作能力とするかを選択する。この選択は、例えば、パラメータ設定に基づいて行なうことができる。

【 0 1 6 4 】

プライマリ仮想マシン 1 3 0 の動作能力とした場合は、プライマリ仮想マシン 1 3 0 の先行動作が拡大していくため、セカンダリ仮想マシン 2 3 0 の遅れを解消するための処理を行なう。このとき、セカンダリ仮想マシン 2 3 0 の遅れが想定される量よりも大きい場合に、マシンの故障が発生したと検出するようにしてもよい。

40

【 0 1 6 5 】

ここでは、遅れを解消するための処理として、2つの補正方法を説明する。まず、第1の補正方法である遅れをリセットして再同期する方法について説明する。

【 0 1 6 6 】

セカンダリ仮想マシン 2 3 0 の同期補正部 2 3 5 は、セカンダリ仮想マシン 2 3 0 のプログラムの実行速度の実測遅延時間が、同期補正情報から得られる想定遅延時間と等しい

50

を確認する。

【0167】

ここで、実測遅延時間について説明する。セカンダリ仮想マシン230は、プライマリ仮想マシン130から、タイマによる定周期疑似割り込みの同期タイミング情報が伝達されるが、この同期タイミング情報は、プライマリ仮想マシン130の実行速度で実行命令数に換算した周期時間のことであり、セカンダリ仮想マシン230にとっては理論値である。一方、セカンダリ仮想マシン230にタイマによる定周期割り込みを発生させた時の周期時間は、実測値である。実測値から理論値を引いた値が実測遅延時間である。

【0168】

例えば、プライマリ仮想マシン130で実際に発生させたタイマ割り込みの周期値が10msのとき、セカンダリ仮想マシン230では、プライマリ仮想マシン130と同じ動作を行なった後に疑似割り込みを入力すると、セカンダリ仮想マシン230でのタイマ割り込みの周期の実測値が15msであったとする。この場合、理論値が10ms、実測値が15msとなるため、タイマ割り込みは5msの遅れと算出される。

【0169】

実測遅延時間が、同期補正情報から得られる想定遅延時間と同程度であれば、正常な遅れと判断し、処理を継続する。処理を継続することにより、実測遅延時間は積算されるため、許容範囲となる遅延許容時間を予め決めておき、実測遅延時間が遅延許容時間を越えた時点で、並列同期動作の開始時と同じ同期初期化処理を再度行なって、遅れをリセットする。

【0170】

例えば、プライマリ仮想マシン130のプログラム実行速度が1500MIPS (Million Instructions Per Second)、セカンダリ仮想マシン230のプログラム実行速度が1000MIPSの場合、セカンダリ仮想マシン230では、1秒ごとに500MI (5億命令数) を処理する時間分の遅延が発生していれば、正常の遅れと判断する。また、10秒間継続して蓄積された5000MI (50億命令数) を処理する時間は、セカンダリ仮想マシン230では、5秒間に相当し、遅延許容時間を越えた場合には、再同期して遅れをリセットする。

【0171】

次に、第2の補正方法である時間を進めて遅れを取り戻す方法について説明する。これは、プライマリゲストOS140がアイドル状態になり、次の疑似割り込みが入力されるまで動作が停止している場合、つまり、CPU処理の負荷が100%より小さい場合は、セカンダリゲストOS240が再開する時間を進めることで遅れを取り戻す。

【0172】

プライマリゲストOS140のアイドル状態とは、実行する処理がなく、プライマリゲストOS140の動作が停止している状態である。プライマリゲストOS140は、疑似割り込みが入力されて割り込みハンドラ処理を実行するまでは動作を再開しない。それまでは、プライマリ仮想マシン130が外部割り込みの入力を待っている状態である。

【0173】

プライマリ仮想マシン130では、実際に外部割り込みが入力されるまで待ち続け、セカンダリ仮想マシン230では、疑似割り込みを入力するために同期情報を受信するまで待ち続ける。

【0174】

セカンダリ仮想マシン230が遅れている場合は、同期情報が溜まっているため、次に入力する疑似割り込みを知ることができる。そのため、セカンダリ仮想マシン230は、プライマリ仮想マシン130と同じ時間待つ必要はなく、セカンダリゲストOS240が再開する時間を進めて疑似割り込みを入力することができる。

【0175】

例えば、第1の補正方法と同様に、プライマリ仮想マシン130のプログラム実行速度を1500MIPS、セカンダリ仮想マシン230のプログラム実行速度を1000MI

10

20

30

40

50

PSとする。プライマリゲストOS 140のCPU負荷が20%の場合、1秒間に300MIの処理を行ない、1秒間の残りの時間だけプライマリ仮想マシン130を停止させる。プライマリ仮想マシン130の停止時間は、プライマリ仮想マシン130が、1200MIの処理を行なう時間分である。

【0176】

一方、セカンダリ仮想マシン230は、1秒間にプライマリ仮想マシン130と同じ300MIの処理を行なった後に、プライマリ仮想マシン130と同じ1200MIの処理を行なう時間分停止するのではなく、1秒間の残りの時間、つまり、セカンダリ仮想マシン230が700MIの処理を行なう時間分停止する。このように、セカンダリ仮想マシンの停止時間を1200MIから700MI処理時間分に短縮することにより、セカンダリ仮想マシン230の遅れを取り戻すことが可能となる。

10

【0177】

次に、プライマリ仮想マシン130がセカンダリ仮想マシン230よりも高速である場合に、セカンダリ仮想マシン230に通知する同期補正情報の動作速度情報として、セカンダリ仮想マシン230の動作能力とした場合について説明する。

【0178】

この場合、プライマリ仮想マシン130がセカンダリ仮想マシン230の動作能力で動作するため、並列同期動作時に、プライマリ仮想マシン130の同期情報生成部135は、同期補正処理を行ない、プログラム実行速度を調整する。具体的には、セカンダリ仮想マシン230のプログラム実行速度に合わせて、ウェイトを入れながら動作を行なう。

20

【0179】

次に、動作能力調整後のセカンダリ仮想マシン230が動作能力調整後のプライマリ仮想マシン130よりも高速である場合の同期補正について説明する。同期情報生成部135は、セカンダリ仮想マシン230に通知する同期補正情報の動作速度情報を、プライマリ仮想マシン130の動作能力とする。

【0180】

そして、プライマリ仮想マシン130は、並列同期動作時には、同期補正処理を行わず、セカンダリ仮想マシン230側で、プライマリ仮想マシン130のプログラム実行速度に合わせて、ウェイトを入れながら動作させる。

【0181】

例えば、プライマリ仮想マシン130のプログラム実行速度が1000MIPS、セカンダリ仮想マシン230の実行速度が1500MIPSの場合、セカンダリ仮想マシン230は、定周期の擬似割り込みごとに、1000MIの処理を行なった後で、500MIの処理をする時間分のウェイトを入れる。これにより、双方の仮想マシン(130、230)は、等しく1000MIPSの処理を行なうことができる。

30

【0182】

(アクセス時間が不確定なストレージデバイスに対する対応)

仮想マシン(130、230)から磁気ディスク、SSD、CD、DVD、テープ等のストレージデバイス(115、215)にアクセスを行なう場合、ストレージデバイス(115、215)の個体差や経年劣化等によりアクセス要求に対する完了待ち時間にばらつきが生じる。このばらつきによる同期ずれを防ぐためには、遅い方の完了待ち時間に合わせる必要があり、同期待ちが発生して実行速度の低下を招くことになる。

40

【0183】

本実施形態のフォールトトレラントシステム10では、デバイス(115、215)がストレージデバイスの場合、仮想マシン(130、230)が、デバイス(115、215)をメモリ上にエミュレーションして、デバイス(115、215)へのアクセスをメモリアクセスに置き換えることにより、デバイス(115、215)への読み書き要求に対する完了の待ち時間をゼロにすることができる。

【0184】

この場合、デバイス(115、215)への読み書き要求に対する完了割り込みは、実

50

際の完了割り込みを待つ必要はなく、要求時のメモリアクセスと同時に擬似割り込みをゲストOS (2 4 0、2 4 0) に入力することによって、読み書き要求と完了擬似割り込みを共通の同期タイミングで処理することができる。これにより、完了割り込みによる擬似割り込みの同期処理が不要となり、処理負荷が軽減されるとともに、同期待ちによる実行速度の低下を防ぐことができる。

【 0 1 8 5 】

メモリ上にエミュレーションしたデバイス (1 1 5、2 1 5) の内容の変更は、ゲストOS (1 4 0、2 4 0) からのアクセスとは非同期にデバイス (1 1 5、2 1 5) に反映する。対応方法は、メモリを用意する場所によって処理が異なる。

【 0 1 8 6 】

(1) 仮想マシン (1 3 0、2 3 0) 内のメモリを使用する場合

仮想マシン (1 3 0、2 3 0) は、ストレージデバイス (1 1 5、2 1 5) の内容を仮想マシン (1 3 0、2 3 0) に割り当てられたメモリ上に展開し、ストレージへのアクセスをメモリアクセスに置き換える。このメモリを「仮想ストレージ」と称する。仮想ストレージの内容が更新された場合は、ゲストOS (1 4 0、2 4 0) からのアクセスとは非同期にストレージデバイス (1 1 5、2 1 5) に反映させる。

【 0 1 8 7 】

具体的には、以下の手順によりデータの読み込み、書き込みを行なう。まず、仮想ストレージの初期化として、ストレージデバイス (1 1 5、2 1 5) の内容を仮想マシン (1 3 0、2 3 0) に割り当てられたメモリ上に展開しておく。

【 0 1 8 8 】

データの読み込みの場合は、ゲストOS (1 4 0、2 4 0) が仮想ストレージにデータ読み込みを要求すると、仮想マシン (1 3 0、2 3 0) が、仮想ストレージ上の内容をゲストOS (1 4 0、2 4 0) に通知するとともに、読み込み完了の擬似割り込みをゲストOS (1 4 0、2 4 0) に入力する。

【 0 1 8 9 】

データの書き込みの場合は、ゲストOS (1 4 0、2 4 0) が仮想ストレージにデータ書き込みを要求すると、仮想マシン (1 3 0、2 3 0) が、仮想ストレージにデータを書き込むとともに、書き込み完了の擬似割り込みをゲストOS (1 4 0、2 4 0) に入力する。そして、仮想マシン (1 3 0、2 3 0) は、仮想ストレージの変更箇所を、非同期でストレージデバイス (1 1 5、2 1 5) に書き込む。ストレージデバイス (1 1 5、2 1 5) への書き込み完了後、完了の外部割り込みが仮想マシン (1 3 0、2 3 0) に入力されるが、この完了の外部割り込みに基づく擬似割り込みは行なう必要がない。

【 0 1 9 0 】

(2) ゲストOS (1 4 0、2 3 0) 内のメモリを使用する場合

ゲストOS (1 4 0、2 3 0) がRAMディスクを利用できる場合は、ストレージデバイス (1 1 5、2 1 5) のエミュレーションにゲストOS (1 4 0、2 3 0) で設定されたRAMディスクを使用する。RAMディスクとは、読み書きに対する待ち時間がゼロであり、完了割り込みが発生しないデバイスである。例えば、Linux (登録商標) などのOSでRAMディスクを利用できる。

【 0 1 9 1 】

上述の仮想ストレージへの読み書きはゲストOS (1 4 0、2 3 0) 内のRAMディスクに対して行なわれ、仮想マシン (1 3 0、2 3 0) に対して読み書きの要求は行なわれない。ゲストOS (1 4 0、2 3 0) の内部で処理が完結しているため、読み書き要求や完了割り込みに対する同期処理が不要となる。

【 0 1 9 2 】

具体的には、以下の手順によりRAMディスクに対するデータの読み込み、書き込みを行なう。まず、RAMディスクの初期化として、ストレージデバイス (1 1 5、2 1 5) の内容をゲストOS (1 4 0、2 3 0) のメモリに展開しておく。

【 0 1 9 3 】

10

20

30

40

50

データの読み込みの場合は、ゲストOS (1 4 0、2 3 0) が R A M ディスクにデータ読み込みを要求すると、R A M ディスク上のストレージデバイス (1 1 5、2 1 5) の内容がゲストOS (1 4 0、2 3 0) に通知される。読み込み完了の擬似割り込みは不要である。

【 0 1 9 4 】

データの書き込みの場合は、ゲストOS (1 4 0、2 3 0) が R A M ディスクにデータ書き込みを要求すると、R A M ディスクにデータが書き込まれる。書き込み完了の擬似割り込みは不要である。そして、R A M ディスクの内容を定期的にストレージデバイス (1 1 5、2 1 5) に保存するようにする。

【符号の説明】

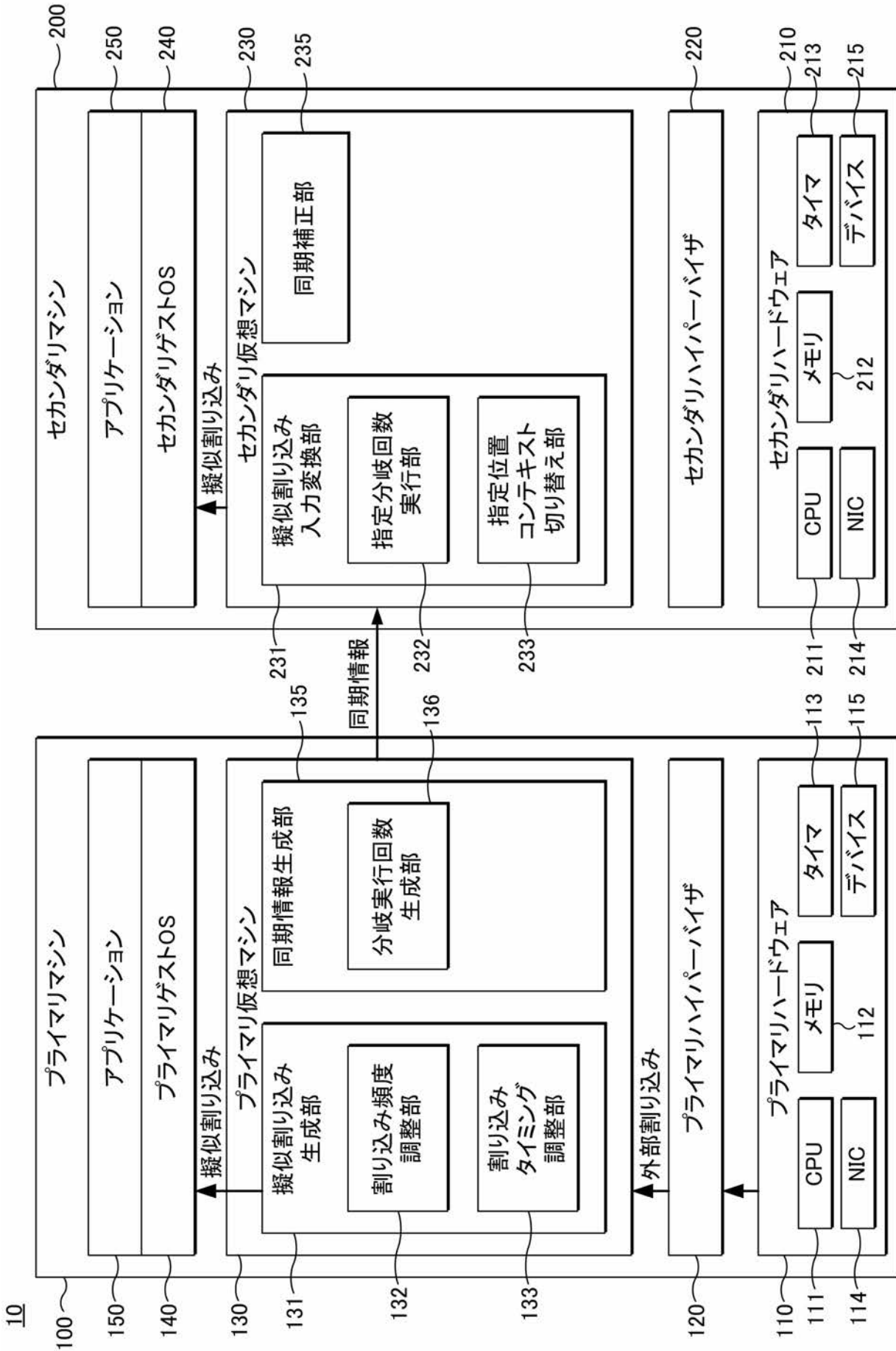
10

【 0 1 9 5 】

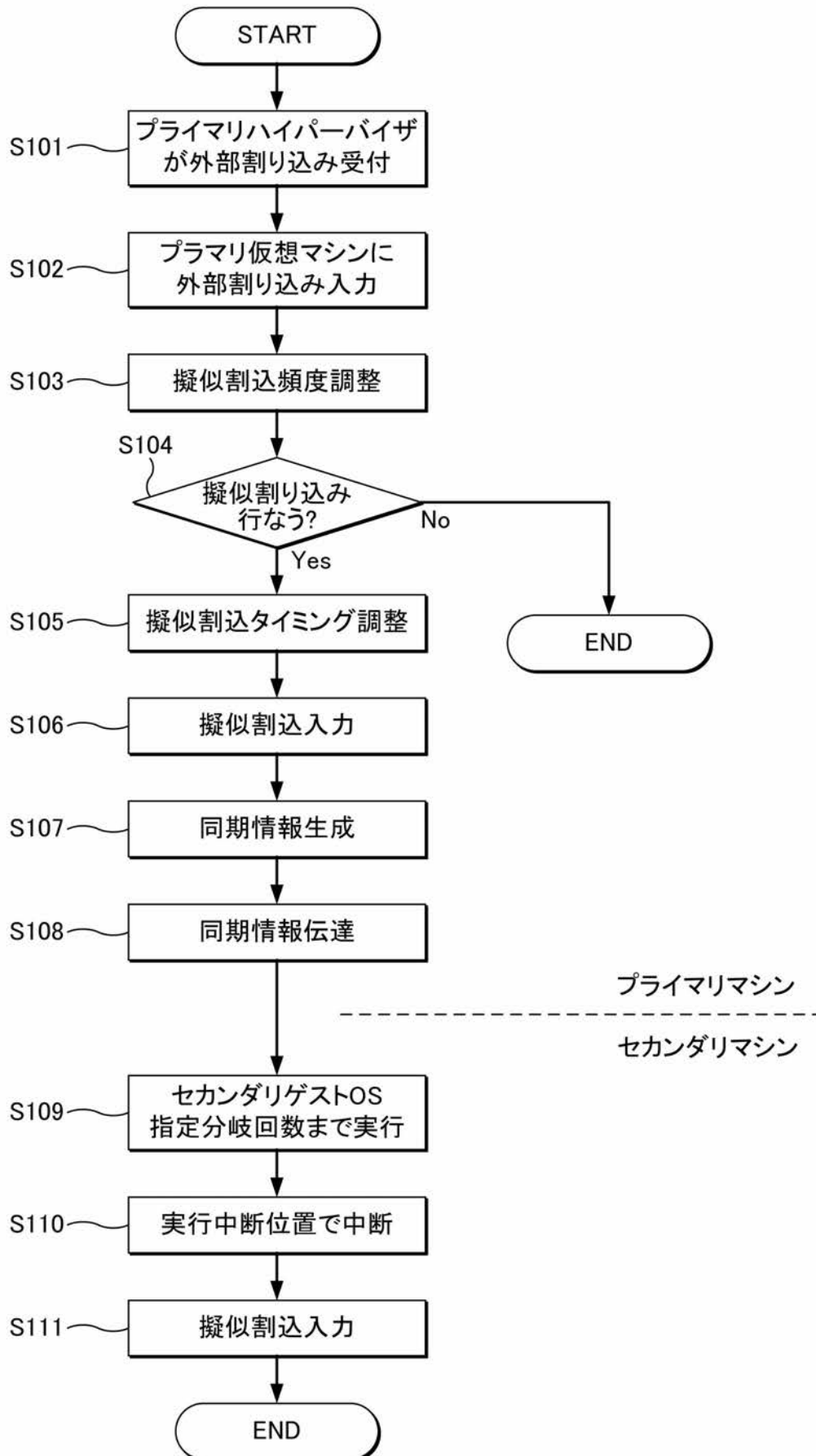
1 0 ... フォールトトレラントシステム、6 0 ... フォールトトレラントシステム、1 0 0 ... プライマリマシン、1 1 0 ... プライマリハードウェア、1 1 1 ... C P U、1 1 2 ... メモリ、1 1 3 ... タイマ、1 1 4 ... N I C、1 1 5 ... デバイス、1 2 0 ... プライマリハイパーバイザ、1 3 0 ... プライマリ仮想マシン、1 3 1 ... 擬似割り込み生成部、1 3 2 ... 割り込み頻度調整部、1 3 3 ... 割り込みタイミング調整部、1 3 5 ... 同期情報生成部、1 3 6 ... 分岐実行回数生成部、1 4 0 ... プライマリゲストOS、1 5 0 ... アプリケーション、2 0 0 ... セカンダリマシン、2 1 0 ... セカンダリハードウェア、2 1 1 ... C P U、2 1 2 ... メモリ、2 1 3 ... タイマ、2 1 4 ... N I C、2 1 5 ... デバイス、2 2 0 ... セカンダリハイパーバイザ、2 3 0 ... セカンダリ仮想マシン、2 3 1 ... 擬似割り込み入力変換部、2 3 2 ... 指定分岐回数実行部、2 3 3 ... 指定位置コンテキスト切り替え部、2 3 5 ... 同期補正部、2 4 0 ... セカンダリゲストOS、2 5 0 ... アプリケーション、6 0 0 ... プライマリマシン、6 1 0 ... プライマリハードウェア、6 2 0 ... プライマリハイパーバイザ、6 3 0 ... プライマリ仮想マシン、6 4 0 ... プライマリゲストOS、6 5 0 ... アプリケーション、7 0 0 ... セカンダリマシン、7 1 0 ... セカンダリハードウェア、7 2 0 ... セカンダリハイパーバイザ、7 3 0 ... セカンダリ仮想マシン、7 3 1 ... 出力データ照合部、7 4 0 ... セカンダリゲストOS、7 5 0 ... アプリケーション

20

【図1】



【図 2】



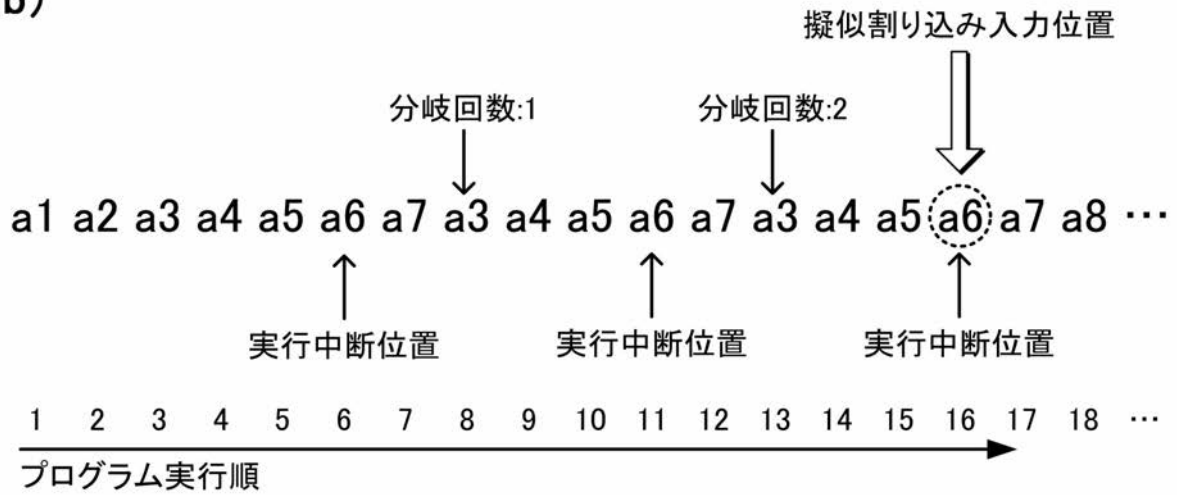
【図3】

(a)

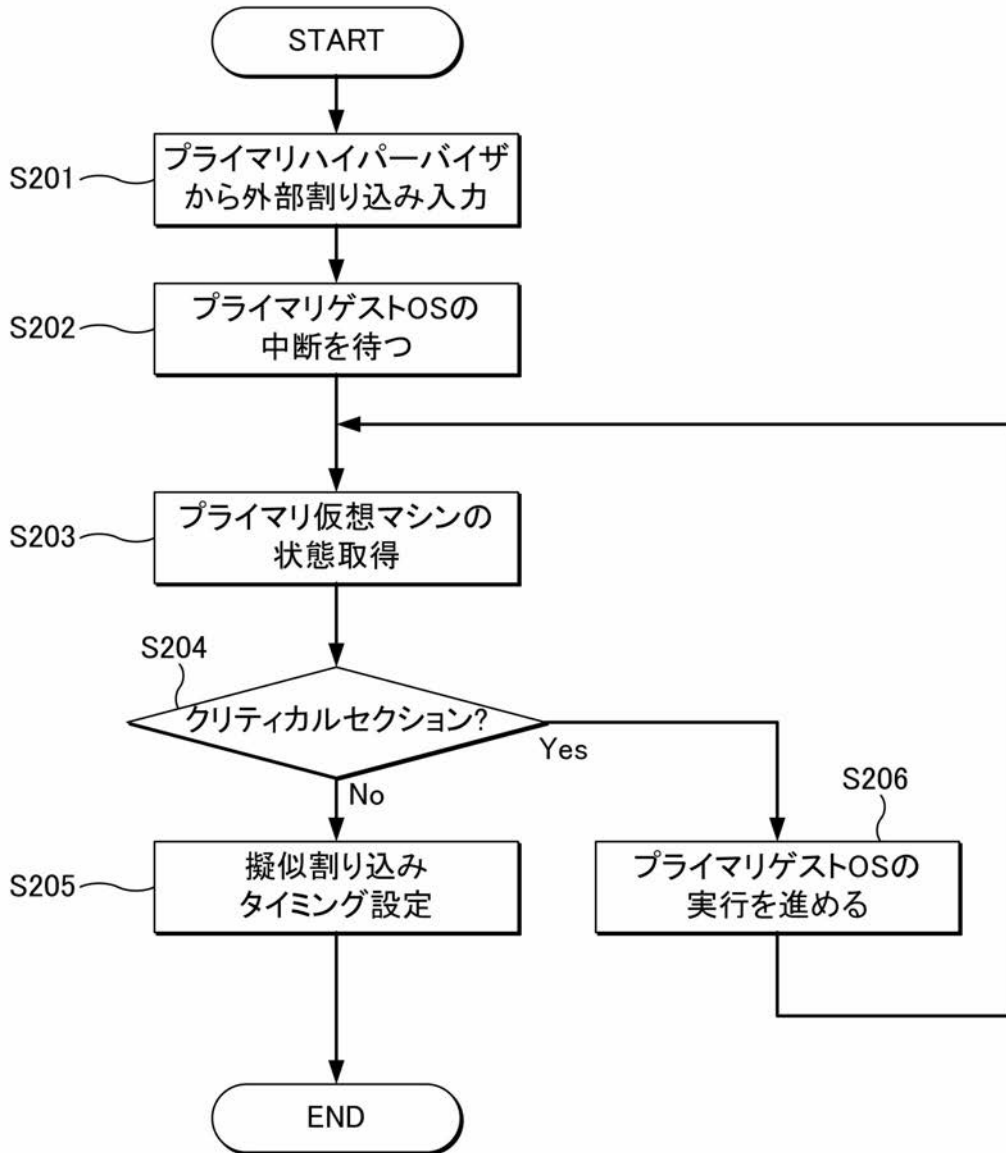
アドレス	命令コード
a1	CodeA
a2	CodeB
a3	CodeC
a4	CodeD
a5	CodeE
a6	CodeF
a7	CodeG
a8	CodeH
・	・
・	・
・	・

3回ループ

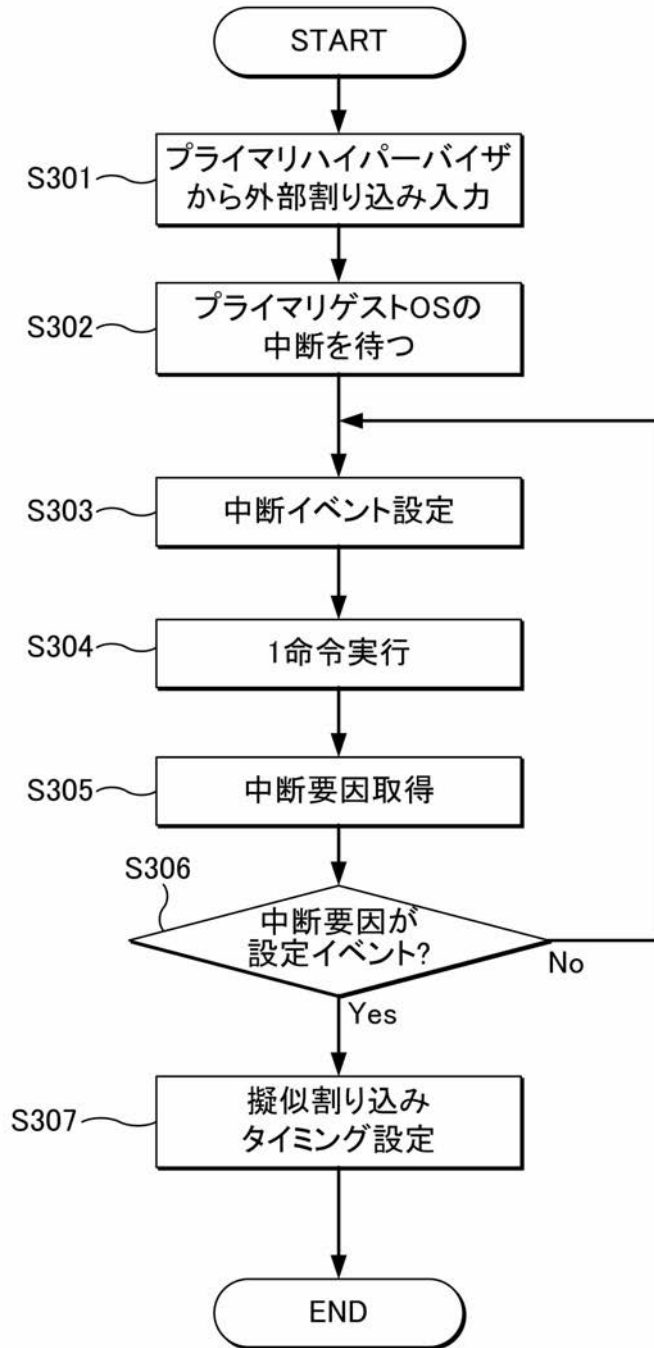
(b)



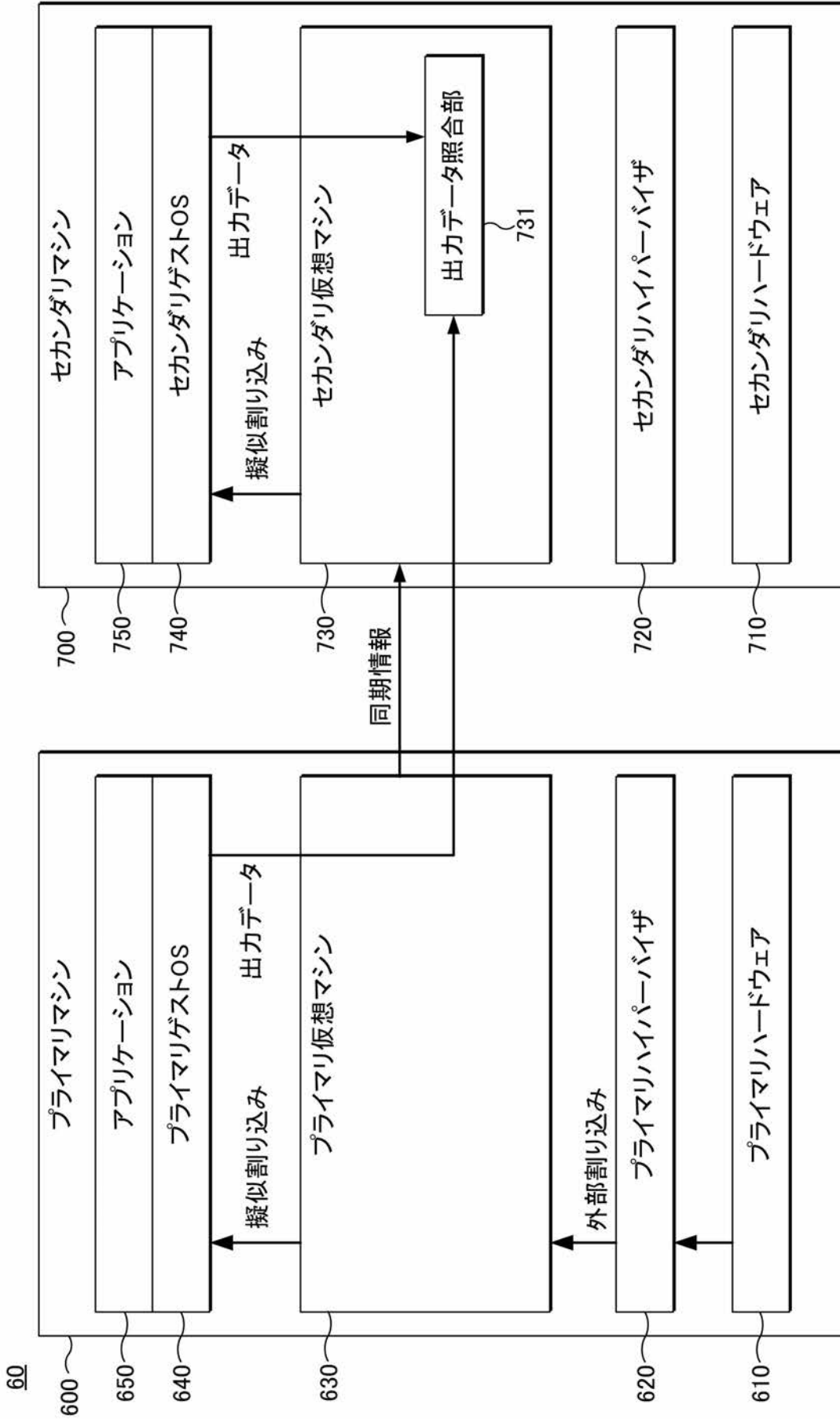
【 図 4 】



【 図 5 】



【図6】



60

【 図 7 】

