

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3657665号
(P3657665)

(45) 発行日 平成17年6月8日(2005.6.8)

(24) 登録日 平成17年3月18日(2005.3.18)

(51) Int. Cl.⁷

F I

G06F 9/46

G06F 9/46 350

G06F 9/52

G06F 9/46 472A

請求項の数 4 (全 53 頁)

<p>(21) 出願番号 特願平7-260543 (22) 出願日 平成7年10月6日(1995.10.6) (65) 公開番号 特開平8-287021 (43) 公開日 平成8年11月1日(1996.11.1) 審査請求日 平成11年7月30日(1999.7.30) (31) 優先権主張番号 特願平7-25458 (32) 優先日 平成7年2月14日(1995.2.14) (33) 優先権主張国 日本国(JP)</p> <p>前置審査</p>	<p>(73) 特許権者 000005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番1号 (74) 代理人 100070150 弁理士 伊東 忠彦 (72) 発明者 村瀬 仁志 静岡県静岡市伝馬町16番地の3 株式会社 社富士通静岡エンジニアリング内 (72) 発明者 高比良 順 静岡県静岡市伝馬町16番地の3 株式会社 社富士通静岡エンジニアリング内 (72) 発明者 平石 壽▲徳▼ 静岡県静岡市伝馬町16番地の3 株式会社 社富士通静岡エンジニアリング内 最終頁に続く</p>
---	--

(54) 【発明の名称】 共用メモリに結合される複数の計算機システム及び共用メモリに結合される複数の計算機システムの制御方法

(57) 【特許請求の範囲】

【請求項1】

外部記憶装置である共用メモリに複数の仮想計算機(ゲストクラスタ)運用可能な計算機(クラスタ)が結合されたシステムであって、

共用メモリ上の領域が他のクラスタのオペレーティングシステムまたはゲストクラスタのオペレーティングシステムにロックされているか否かを判断し、ロックされていないと判断すると、前記共用メモリ上の領域のロックを獲得する、ロック獲得手段と、

前記ロック獲得手段において、他のクラスタのオペレーティングシステムまたはゲストクラスタのオペレーティングシステムが前記共用メモリ上の領域をロックしていると判断すると、ロックを行っているクラスタまたはゲストクラスタと通信を行い、応答がなければロックを獲得しているクラスタまたはゲストクラスタが停止していると判断する第1の停止監視手段と、

前記第1の停止監視手段がロックを獲得しているクラスタまたはゲストクラスタが停止していると判断すると、該停止しているクラスタまたはゲストクラスタが仮想計算機運用されているか否かを判定する仮想・実計算機確認手段と、

前記仮想・実計算機確認手段が前記停止しているクラスタが仮想計算機運用でないと判断すると前記停止しているクラスタと前記共用メモリとのアクセスパスを切断し、前記仮想・実計算機確認手段が前記停止しているゲストクラスタが仮想計算機運用されていると判定するとAVMと該停止している仮想計算機との間の論理パスを切断するパス切断手段と、

10

20

を有することを特徴とするシステム。

【請求項 2】

前記仮想計算機運用可能な計算機が、自計算機内において運用しているゲストクラスタの状態を示す情報を持つパラメータ域手段を持ち、

前記バス切断手段が、前記仮想・実計算機確認手段が前記停止しているゲストクラスタが仮想計算機運用されていると判断すると、該停止しているゲストクラスタを運用している仮想計算機運用可能な計算機が持つ前記パラメータ域手段を参照して、該停止しているゲストクラスタを特定してダウン状態にしてから、A V Mと該停止している仮想計算機との間の論理バスを切断する手段を有すること、を特徴とする請求項 1 記載のシステム。

【請求項 3】

外部記憶装置である共用メモリに複数の仮想計算機（ゲストクラスタ）運用可能な計算機（クラスタ）が結合されたシステムの制御方法であって、

共用メモリ上の領域が他のクラスタのオペレーティングシステムまたはゲストクラスタのオペレーティングシステムにロックされているか否かを判断し、ロックされていないと判断すると、前記共用メモリ上の領域のロックを獲得し、

他のクラスタのオペレーティングシステムまたはゲストクラスタのオペレーティングシステムが前記共用メモリ上の領域をロックしていると判断すると、ロックを行っているクラスタまたはゲストクラスタと通信を行い、応答がなければロックを獲得しているクラスタまたはゲストクラスタが停止していると判断し、

ロックを獲得しているクラスタまたはゲストクラスタが停止していると判断すると、該停止しているクラスタまたはゲストクラスタが仮想計算機運用されているか否かを判定し、

前記停止しているクラスタが仮想計算機運用でないとして判断すると前記停止しているクラスタと前記共用メモリとのアクセスパスを切断し、前記停止しているゲストクラスタが仮想計算機運用されていると判定すると A V M と該停止している仮想計算機との間の論理バスを切断する

ことを特徴とする方法。

【請求項 4】

前記仮想計算機運用可能な計算機が、自計算機内において運用しているゲストクラスタの状態を示す情報を持つパラメータ域手段を持ち、

前記停止しているゲストクラスタが仮想計算機運用されていると判断すると、該停止しているゲストクラスタを運用している仮想計算機運用可能な計算機が持つ前記パラメータ域手段を参照して、該停止しているゲストクラスタを特定してダウン状態にしてから、A V Mと該停止している仮想計算機との間の論理バスを切断すること、

を特徴とする請求項 3 記載の方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、共用メモリに結合される複数の計算機システム及び共用メモリに結合される複数の計算機システムの制御方法に係り、特に、共用メモリに結合された複数の計算機間で種々の制御を行うための、共用メモリに結合される複数の計算機システム、及び共用メモリに結合される複数の計算機システムの制御方法に関する。

【0002】

近年のコンピュータシステムでは、単一プロセッサの能力の伸びの鈍化、信頼性向上の強いニーズ等の理由から共用メモリを介した複数の計算機システムにより構築されたシステムが一般的になりつつある。また、共用メモリを介して1つの計算機システムを複数の仮想計算機システムとして利用することが要求されている。

【0003】

さらには、仮想計算機システムを制御するオペレーティングシステム（以下、オペレーティングシステムのことを OS と記す）である仮想計算機制御プログラム（以下、仮想計算

10

20

30

40

50

機制御プログラムのことをA V Mと記す)で運用された計算機の異常によるダウンを検出し、ホットスタンバイが可能なシステムが要求されている。

【0004】

【従来の技術】

(1) 従来の計算機システム - 1

最初に従来の第1の計算機システムを説明する。

図47は、従来の第1の計算機システムの構成例を示す。

【0005】

同図の例は、1つの実計算機(以下、クラスタと記す)10を複数の仮想計算機11-1~11-n(以下、ゲストクラスタと記す)で運用する例である。クラスタ10は、仮想計算機であるゲストクラスタ11-1~11-n制御用の制御オペレーティングシステム(以下、A V Mと記す)12を有し、当該A V Mが複数のゲストクラスタ11-1~11-nを制御する。

10

【0006】

(2) 従来の計算機システム - 2

第2に、計算機システムが外部記憶装置と接続されている場合について説明する。

図48は、従来の第2の計算機システムの構成例を示す。

【0007】

同図の例は、上記に示した1つの第1の計算機システム(クラスタ)10が外部記憶装置(以下、S S Uと記す)50に接続されている例を示している。クラスタ10とS S U50は、1台のS S U50が有する実アクセスパス60により接続され、クラスタ10はS S U50に対して情報の読出し/書き出しの処理を実行する。

20

【0008】

また、クラスタ10は、A V M12と複数のゲストクラスタ11-1~11-nを有する。A V M12とゲストクラスタ11の間には、各々論理(仮想)アクセスパス71が介在している。ゲストクラスタ11-1から11-nは、このアクセスパス71、及びA V M12を介してS S U50より情報の読出しや書込みを行う。

【0009】

図49は、従来の第2の計算機システムを説明するための図である。

同図に示すシステムは、従来の第2の計算機システムにおいて、S S U51にアクセスパス61を介してクラスタ10が接続され、アクセスパス62を介してクラスタ20が接続されている。S S U52には、アクセスパス63を介してクラスタ30が接続され、アクセスパス64を介してクラスタ40が接続されている。

30

【0010】

このうち、S S U51に接続されるクラスタ10がS S U51に対して処理を実行中であり、クラスタ20は、1つのゲストクラスタがA V Mの制御により処理待ち状態であり、クラスタ20内の他のゲストクラスタは開発に使用されている。また、S S U52に接続されるクラスタ30は、S S U52に対して処理を実行中であり、クラスタ40は、処理待ち中となっている。このように、図49に示すシステムは、ホットスタンバイにおいて、1つのS S U51(52)に対して1つのクラスタ10(30)が実行している時は、他のクラスタ20(40)は、待機中とすることにより排他制御を行っている。

40

【0011】

(3) 従来の計算機システム - 3

第3に1つのS S Uに複数のクラスタが接続されている場合について説明する。

図50は、従来の第3の計算機システムの構成例を示す。

【0012】

同図に示す計算機システムは、1つのS S U50に複数のクラスタ10、20、30、40が接続されている例である。A V M運用のクラスタ30、40内のゲストクラスタは、各々クラスタ内で相対的な計算機番号(以下、相対計算機番号と言う)を有している。例えば、クラスタ30の各ゲストクラスタの相対計算機番号は、ゲストクラスタ31-1= “

50

1”、ゲストクラスタ3 1 -2= “ 2 ”、ゲストクラスタ3 1 -3= “ 3 ”、ゲストクラスタ3 1 -4= “ 4 ”のように付与されている。また、クラスタ4 0についても同様に、ゲストクラスタ4 1 -1= “ 1 ”、ゲストクラスタ4 1 -2= “ 2 ”、ゲストクラスタ4 1 -3= “ 3 ”、ゲストクラスタ4 1 -4= “ 4 ”のように相対計算機番号が付与されている。また、クラスタ1 0、2 0についてもクラスタ1 0 = 0、クラスタ2 0 = 1、クラスタ3 0 = 2、クラスタ4 0 = 3のように予め実計算機番号が設定されている。

【0013】

ここで、仮想計算機により運用されているクラスタ3 0のゲストクラスタ3 1 -1をオペレータ8 0が指定する場合について述べる。オペレータ8 0は、クラスタ3 0の実計算機番号“ 2 ”を指定すると、実計算機番号“ 2 ”のクラスタ上のAVM 3 2が設定される。この場合、AVM 3 2は、予め定められた順序または、配列順にクラスタ3 0内のゲストクラスタ3 1 -1を示す相対計算機番号（例えば、“ 1 ”）を指定したことになる。即ち、SSUを介した複数の計算機システムにおいて、AVM運用された計算機内では、1つのゲストクラスタのみ、他の計算機と結合することができる。従って、実計算機番号と仮想計算機番号が1対1の関係であるため、実計算機番号を指定しても、その指定番号から仮想計算機を特定することが可能である。つまり、オペレータが実計算機番号“ 2 ”を指定するということは、ゲストクラスタ“ 3 1 - 1 ”を指定するのと同義である。

【0014】

(4) 従来の計算機システムの通信方式

第4に従来の計算機システムにおいて、クラスタ間で通信を行う場合について説明する。図5 1は、従来の第3の計算機システムにおける通信システムを説明するための図（その1）である。同図に示すように、1つのSSU 5 0を複数のクラスタ2 0、3 0、4 0で共用しているとき、クラスタ間でSSU 5 0を介して他計算機と通信する場合には、クラスタ同士が相手先のクラスタの実計算機番号を指定して通信する。例えば、同図において、クラスタ1 0がクラスタ2 0を指定する場合には、クラスタ1 0の実計算機番号が“ 0 ”であり、クラスタ2 0の実計算機番号が“ 1 ”であるとき、クラスタ1 0は、クラスタ2 0の実計算機番号“ 1 ”を指定することにより、通信を行う。また、クラスタ4 0は、複数のゲストクラスタ4 1 - 1 ~ 4 1 - nを有している。このとき、クラスタ4 0のゲストクラスタ4 1 - 3は、クラスタ2 0の実計算機番号“ 1 ”を指定して、AVM 4 2とSSU 5 0を介してクラスタ2 0と通信することが可能である。また、クラスタ1 0、クラスタ2 0、または、クラスタ4 0に対して通信を行う場合には、実計算機番号“ 2 ”を指定することで、SSU 5 0を介して通信することが可能である。

【0015】

(5) 従来の通信時における割り込み処理

次に、従来の通信における割り込み時の処理について説明する。

前述の1つのSSUを複数のクラスタで共用している複合システムにおいて、GSIGP命令によりシステム間の通信を行うことが可能であり、割り込みの保留状態や反映状態を正確に把握するには、割り込みの保留状態を利用して、図5 2に示すような通信処理を行う。なお、GSIGP命令の機能には、クラスタ間の通信機能と他クラスタの制御機能がある。他クラスタの制御機能とは、ダウンクラスタを制御する際に用いるものであり、CPU停止、I/Oリセット、ダンプ採取等を行うものである。

【0016】

図5 2は、従来の通信時における割り込み処理を説明するためのシーケンスチャートである。

クラスタAとクラスタBが通信を行うものとして説明する。

ステップ1) クラスタAは、通信要求としてGSIGP命令をSSUを介してクラスタBに発行する。

【0017】

ステップ2) クラスタBは、割り込みが反映状態でないため、ハードウェアで割り込みを保留しておく。

10

20

30

40

50

ステップ3) クラスタAは、次の通信要求が発生した場合、G S I G P命令をS S Uを介して発行する。

【0018】

ステップ4) S S Uは、クラスタBが割り込み保留状態であると判定された場合に、その割り込みがいずれ確実にクラスタBに反映されることを前提として、クラスタAからのG S I G P命令を受け取り、S S U上にキューイングする。ステップ5) クラスタBは、割り込みが反映可能となった時点で、ハードウェアが保留状態を解除し、割り込みを反映する。

【0019】

ステップ6) クラスタBは、割り込まれた通信要求と共にS S U上にキューイングされていた通信要求も処理する。 10

なお、S S U上に通信要求をキューイングし、1回の割り込みで複数の通信要求を処理することができる。

【0020】

(6) 従来のシステム制御

従来の実計算機で運用されるクラスタ同士でS S Uを共用した複合システムの場合、G S I G P命令(リセット)により、ダウンしたシステムを他のクラスタシステムからリセットするようなシステム制御を行うことができ、このリセット完了に基づいて、ホットスタンバイによるシステムの切り替えを行う。

【0021】

図53は、従来のシステム制御のリセット処理を説明するためのシーケンスチャートである。 20

以下の説明では、クラスタAがクラスタBをリセット制御する場合について述べる。

【0022】

ステップ10) クラスタAは、クラスタBをリセットするため、G S I G P命令(リセット)を発行する。

ステップ11) クラスタBは、S S Uを介してハードウェアでシステムのリセットを開始する。

【0023】

ステップ12) クラスタAがリセット完了/リセット中をG S I G P命令(センス)を発行することで認識できる。 30

ステップ13) クラスタAは、G S I G P命令(センス)の結果によりリセット中を認識する通知する。

【0024】

ステップ14) クラスタBは、ハードウェアでのリセットを完了する。

ステップ15) クラスタAが、G S I G P命令(センス)を発行する。

ステップ16) クラスタAは、クラスタBがリセットを完了したことをG S I G P命令(センス)の結果により認識する。

【0025】

(7) 従来のダウン時における処理 40

従来、OSがダウンを検出すると、G S I G P命令(CPU停止)やG S I G P命令(リセット)等により、ダウンクラスタの制御を行う。なお、G S I G P命令は、各クラスタ毎に配置されているサービスプロセッサが受け付け、実行する。また、サービスプロセッサは、G S I G P命令(ペンディング)が永久に続く場合を想定し、タイマ監視を行い、タイムアウトになった場合に強制リセットを行う。

【0026】

また、ゲストクラスタがセッション閉塞(D E A C T I V A T E)した場合には、A V Mがゲストクラスタのセッション閉塞を認識し、ゲストクラスタとA V M間の論理パスを切断することにより、ゲストクラスタを切り離す処理を行う。

【0027】

【発明が解決しようとする課題】

しかしながら、上記の従来の(1)～(7)の各システムには、以下のような問題がある。

1 上記図47から図49に示すシステムは、クラスタ単独または、SSUとクラスタが直接接続されている構成であり、他のクラスタ間との情報の授受を行うことは可能であるが、AVM運用した場合に1クラスタのみしか行うことができないという問題がある。

【0028】

2 複数のクラスタが同時にIPL操作により初期化を行うような場合に、同時に共用メモリに対する初期化が行われる可能性があり、データの整合性がとれなくなるという問題がある。さらに、IPL操作を介したOSが停止して、再度起動した場合には、システムの誤動作に繋がる可能性もある。

10

【0029】

3 また、図50に示すシステムは、オペレータから指定されたゲストクラスタをAVMで所定の順序で指定することは可能であるが、図54に示すようなシステムのように、クラスタ内の複数のゲストクラスタに、予め指定順序が決定されていない場合や、同時に動作するような場合には、実計算機番号だけでは仮想計算機を特定することができない。図54において、オペレータ80がクラスタ30の実計算機番号“2”を指定すると、実計算機番号“2”の計算機上のAVM32に制御が渡る。しかし、AVM32は、このクラスタ30に組み込まれているどのゲストクラスタを指定すればよいのか判断できないため、ある特定のゲストクラスタを対象した通信等の処理ができないという問題がある。

20

【0030】

4 また、従来は、SSUを複数のクラスタが使用する場合、あるクラスタを仮想計算機運用しても1つのゲストクラスタしかSSUを使用することができないという制限がある。つまり、1台のクラスタを複数のゲストクラスタで運用し、さらにそれらの構成のクラスタを何台もSSUを介して接続された複合システムを構築しようとしても、実計算機番号のみ管理しているため、配下のゲストクラスタを特定することができない。

【0031】

5 また、あるクラスタから他のAVM運用の第1のクラスタに対して通信(GSIGP命令)を発行したとき、第2のクラスタにおいて受け付けられない状態である場合に保留となる。このとき、第1のクラスタが他のAVMの第2のクラスタ宛に通信を発行した場合、保留状態であるため、この通信はキューイングされる。しかし、他のAVMの第1のクラスタが受け付け可能となったときに、SSU上にキューイングされている要求は第2のクラスタ宛であるため、一緒に処理されない。このため、AVMで要求を保留しない方法が考えられる。しかし、計算機間のSSUを共用する通信時において割り込みがある場合には、クラスタAから発行されたGSIGP命令がAVMに反映する前に、送信側のプログラムに割り込みが反映されてしまうため、ハードウェア内には割り込みが保留されない。このような場合に、他系のクラスタからは、割り込みが反映されたように見えてしまうため、仮にAVMに割り込みが反映されない状態であっても、クラスタAからクラスタBのAVMの状態を正確に把握できない。そのため、AVMが割り込み保留状態の場合でも送信側のプログラムにAVM宛の通信要求が次々に発生するが、制御プログラムは割り込みを棄却してしまい、AVMの通信が正しく行われられないという問題がある。また、保留状態によりSSU上に通信要求をキューイングしても、その保留となった割り込みが他のAVM宛である場合、シーケンシャルにキューイングされている要求が、該当するAVMに対していつまでも反映されないため、SSU上の通信要求が該当するAVMに通知できないという問題がある。

30

40

【0032】

6 また、SSUを共用するシステムが複数の仮想計算機で運用されるクラスタであった場合、実計算機で運用されるクラスタ同士の時と同じようにGSIGP命令(CPU宛)を発行すると、AVMが動作するクラスタシステムのCPUが停止してしまう。また

50

、特開平5 - 3 2 4 3 6 2号「計算機システム間の通信割り込み制御方式」を適用しても複数の仮想計算機により運用されているクラスタでSSUを共用すると、GSIGP命令（リセット）発行元でリセット完了契機が正しく認識できなくなるという問題がある。これは、ある仮想計算機により運用されているクラスタの1つのゲストクラスタ（例えば、ゲストクラスタa）が他のクラスタ（例えば、クラスタA）からのGSIGP命令（リセット）によりリセット処理中である場合には、別のクラスタ（例えば、クラスタB）からゲストクラスタaと同じ制御プログラム下にあるゲストクラスタbのリセット要求や通信要求ができなくなるからである。このような状況を解決するために、プログラムがAVMのリセットを起動する段階でプログラムからハードウェア（実計算機制御機能）に対してリセット処理中解除を指示した場合であっても、GSIGP命令（リセット）発行元でリセットの完了（リセット処理中が解除）が正しく認識できなくなるという問題がある。図55は、従来の問題点を説明するための図（その2）である。同図に付されている内の番号と以下の番号は一致するものとする。

10

【0033】

1 クラスタAから仮想計算機により運用されているクラスタVのゲストクラスタaをリセットするためにGSIGP命令（ゲストクラスタa宛リセット）を発行する。

2 クラスタVでは、クラスタAからのリセット要求が一旦ハードウェアで保留される。

【0034】

3 クラスタVのAVMがリセット要求を認識すると、AVMによりゲストクラスタaのリセット処理を行う。

20

4 クラスタBからクラスタVのゲストクラスタbをリセットするために、GSIGP命令（ゲストクラスタb宛リセット）を発行する。

【0035】

5 クラスタVのハードウェアがリセット処理中状態のため、クラスタBからゲストクラスタb宛のリセット要求が受け付けられない。

6 ゲストクラスタaのリセットが完了すると、AVMは、ハードウェアに対してリセット保留の解除を指示する。

【0036】

7 クラスタAは、クラスタVのゲストクラスタaのリセット完了を認識する。つまり、あるクラスタで発行したリセット要求により、ある仮想計算機がリセット処理中である場合には、他のクラスタから現在リセット処理中のAVM下にある別の仮想計算機からのリセット要求や通信要求ができなくなる。

30

【0037】

図56は、従来の問題点を説明するための図（その3）である。同図に示す内の番号と、以下の番号は一致するものとする。

1 クラスタAから仮想計算機により運用されているクラスタVのゲストクラスタaをリセットするためにGSIGP命令（リセット）を発行する。

【0038】

2 特開平5 - 3 2 4 3 6 2『計算機システム間の割込制御方式』の方法により、リセット要求が、クラスタVのハードウェアで保留される。

40

3 仮想計算機により運用されているクラスタVのAVMは、リセット要求を認識すると、ハードウェアに対してリセット保留の解除を指示する。

【0039】

4 AVMがゲストクラスタaのリセット処理を行う。

5 クラスタAでクラスタVのゲストクラスタaのリセット完了を誤認する。

このように、仮想計算機により運用されているクラスタのAVMがリセット処理中解除を指示すると、GSIGP命令（リセット）発行元でリセットの完了（リセット処理中解除）を誤認してしまうという問題がある。

【0040】

50

さらに、従来の共用メモリを介して複数のクラスタが接続されているシステム構成において、A V M運用されているクラスタのゲストクラスタがダウンした場合に、クラスタ毎に付設されているサービスプロセッサは、リセット処理が所定時間内に完了しない場合に行う強制リセットを行うことができない。その理由は、サービスプロセッサの強制リセットが動作する条件は、ペンディング状態を認識した時である。従って、A V Mが処理中の状態を解除するために、サービスプロセッサのタイマ監視が終了してしまい、ペンディング状態を認識できないために強制リセットができない。

【 0 0 4 1 】

また、A V M運用されているゲストクラスタがセッション閉塞 (D E A C T I V A T E) 時には、A V Mとゲストクラスタ間の論理パスを切断してしまうため、リセット等の制御がOS側よりできない。従って、オペレータがリセット等の処理を行う。このため、セッション閉塞が発生する毎に、オペレータがリセットを行わなければならないという問題がある。

10

【 0 0 4 2 】

本発明は、上記の点に鑑みなされたもので、複数のクラスタにより負荷分散して処理しなければならない大規模システムを、仮想計算機により運用されている複数のゲストクラスタを有する複数のクラスタ間で通信を行うことが可能な共用メモリに結合される複数の計算機システムを提供することを第1の目的とする。

【 0 0 4 3 】

また、本発明の第2の目的は、共用メモリの初期化を開始した場合に、他のクラスタからの更新要求を排他制御すると共に、停止したクラスタがあった場合に、再動作等によるシステムの誤動作を防止することが可能な共用メモリに結合される複数の計算機システムを提供することである。

20

【 0 0 4 4 】

また、本発明の第3の目的は、仮想計算機により運用されているクラスタのA V Mが各ゲストクラスタの計算機番号を設定し、それを各A V MのOSとゲストクラスタ制御システムが相互に認識することが可能な共用メモリに結合される複数の計算機システムを提供することである。

【 0 0 4 5 】

また、本発明の第4の目的は、S S Uを介して複数のクラスタ、仮想計算機により運用されているクラスタが接続されるシステム間において、柔軟な通信が可能な共用メモリに結合される複数の計算機システムを提供することである。

30

また、本発明の第5の目的は、A V Mへの割り込みが正しく反映されるような共用メモリに結合される複数の計算機システムを提供することである。

【 0 0 4 6 】

また、本発明の第6の目的は、複数の仮想計算機システムでS S Uを共用した場合でも、G S I G P命令 (リセット) によるリセットの完了が発行元に正しく認識可能となる共用メモリに結合される複数の計算機システムを提供することである。

【 0 0 4 7 】

また、本発明の第7の目的は、A V M自身の異常を検出し、他のクラスタに通知することが可能であり、ダウンしたクラスタ及び、A V M運用しているクラスタの配下のゲストクラスタに対する制御を他のクラスタから行うことが可能な共用メモリに結合される複数の計算機システムを提供することである。

40

【 0 0 4 8 】

また、本発明の第8の目的は、自クラスタでダウンした状況を他の共用メモリに接続されるクラスタに通知することが可能な共用メモリに結合される複数の計算機システムを提供することである。

また、本発明の第9の目的は、ダウンしているクラスタがある場合に、他のクラスタのOSから当該ダウンを認識することが可能な共用メモリに結合される複数の計算機システムを提供することである。

50

【 0 0 4 9 】

また、本発明の第 1 0 の目的は、他のクラスタからダウクラスタを認識した場合に、認識したクラスタからダウクラスタの制御を行うことが可能な共用メモリに結合される複数の計算機システムを提供することである。

また、本発明の第 1 1 の目的は、ダウンしたクラスタの CPU 停止及び I / O リセットを行い、ホットスタンバイ状態にすることが可能な共用メモリに結合される複数の計算機システムを提供することである。

【 0 0 5 0 】

また、本発明の第 1 2 の目的は、クラスタ間でダウクラスタの制御を行う場合に、A V M 運用のクラスタのダウン時に、実クラスタと同様にハードウェアより強制的にリセット

10

【 0 0 5 1 】

また、本発明の第 1 3 の目的は、オペレータ介入メッセージが表示され、オペレータが介入操作を終了した時点で直ちに当該メッセージを消去できる共用メモリに結合される複数の計算機システムを提供することである。

【 0 0 5 2 】

【課題を解決するための手段】

図 1 は、本発明の原理構成図である。

本発明は、外部記憶装置である共用メモリに複数の仮想計算機（ゲストクラスタ）運用可能な計算機（クラスタ）が結合されたシステムであって、

20

共用メモリ 1 0 0 上の領域が他のクラスタのオペレーティングシステムまたはゲストクラスタのオペレーティングシステムにロックされているか否かを判断し、ロックされていないと判断すると、共用メモリ上の領域のロックを獲得する、ロック獲得手段 1 と、

ロック獲得手段 1 において、他のクラスタのオペレーティングシステムまたはゲストクラスタのオペレーティングシステムが共用メモリ上の領域をロックしていると判断すると、ロックを行っているクラスタまたはゲストクラスタと通信を行い、応答がなければロックを獲得しているクラスタまたはゲストクラスタが停止していると判断する第 1 の停止監視手段 2 と、

第 1 の停止監視手段 2 がロックを獲得しているクラスタまたはゲストクラスタが停止していると判断すると、該停止しているクラスタまたはゲストクラスタが仮想計算機運用されているか否かを判定する仮想・実計算機確認手段 3 と、

30

仮想・実計算機確認手段 3 が停止しているクラスタが仮想計算機運用でないと判断すると停止しているクラスタと共用メモリとのアクセスパスを切断し、仮想・実計算機確認手段 3 が停止しているゲストクラスタが仮想計算機運用されていると判定すると A V M と該停止している仮想計算機との間の論理パスを切断するパス切断手段 4 と、を有する。

【 0 0 5 3 】

本発明は、仮想計算機運用可能な計算機が、自計算機内において運用しているゲストクラスタの状態を示す情報を持つパラメータ域手段を持ち、

パス切断手段 4 が、仮想・実計算機確認手段が停止しているゲストクラスタが仮想計算機運用されていると判断すると、該停止しているゲストクラスタを運用している仮想計算機運用可能な計算機が持つパラメータ域手段を参照して、該停止しているゲストクラスタを特定してダウン状態にしてから、A V M と該停止している仮想計算機との間の論理パスを切断する。

40

【 0 0 5 4 】

本発明は、外部記憶装置である共用メモリに複数の仮想計算機（ゲストクラスタ）運用可能な計算機（クラスタ）が結合されたシステムの制御方法であって、

共用メモリ上の領域が他のクラスタのオペレーティングシステムまたはゲストクラスタのオペレーティングシステムにロックされているか否かを判断し、ロックされていないと判断すると、共用メモリ上の領域のロックを獲得し、

50

他のクラスタのオペレーティングシステムまたはゲストクラスタのオペレーティングシステムが共用メモリ上の領域をロックしていると判断すると、ロックを行っているクラスタまたはゲストクラスタと通信を行い、応答がなければロックを獲得しているクラスタまたはゲストクラスタが停止していると判断し、

ロックを獲得しているクラスタまたはゲストクラスタが停止していると判断すると、該停止しているクラスタまたはゲストクラスタが仮想計算機運用されているか否かを判定し

停止しているクラスタが仮想計算機運用でないとは判断すると停止しているクラスタと共用メモリとのアクセスパスを切断し、停止しているゲストクラスタが仮想計算機運用されていると判定するとA V Mと該停止している仮想計算機との間の論理パスを切断する。

10

【 0 0 5 5 】

本発明は、仮想計算機運用可能な計算機が、自計算機内において運用しているゲストクラスタの状態を示す情報を持つパラメータ域手段を持ち、

停止しているゲストクラスタが仮想計算機運用されていると判断すると、該停止しているゲストクラスタを運用している仮想計算機運用可能な計算機が持つパラメータ域手段を参照して、該停止しているゲストクラスタを特定してダウン状態にしてから、A V Mと該停止している仮想計算機との間の論理パスを切断する。

【 0 0 8 5 】

上記の発明は、以下に示す作用を有する。

本発明は、外部記憶装置（共用メモリ）に仮想計算機により運用される実計算機を複数接続することが可能となり、ある実計算機に包含される仮想計算機と他の実計算機に包含される仮想計算機との通信が可能となる。

20

【 0 0 8 6 】

また、本発明は、共用メモリに接続される仮想計算機により運用されている実計算機（以下、クラスタと記す）がホットスタンバイ時に共用メモリをロックし、他のクラスタからのアクセスを排除制御することが可能であると共に、ロックしているクラスタが異常停止した状態を検出した場合に、自クラスタから初期化処理を行い、異常停止したクラスタと共用メモリとのアクセスパスを切断するため、共用メモリのデータ破壊を防止できる。

【 0 0 8 7 】

また、本発明は、仮想計算機（以下、ゲストクラスタと記す）により運用されているクラスタ内のゲストクラスタの初期化処理が異常停止している場合には、当該ゲストクラスタを制御するOS（A V M）と当該ゲストクラスタ間の論理パスを切断する。これにより、クラスタのアクセスパスは切断しないため、他のゲストクラスタから共用メモリにアクセスすることが可能となる。

30

【 0 0 9 2 】

また、本発明は、他のクラスタからリセット要求を仮想計算機により運用されているクラスタが受信した場合に、ゲストクラスタを制御するOSが制御対象のゲストクラスタを特定し、そのゲストクラスタの制御が完了した時点で、リセット要求の発行元に対して完了通知を送信することにより、リセット要求元では、ホットスタンバイによる切り替えが可能となる。

40

【 0 0 9 3 】

また、本発明は、クラスタ内の仮想計算機のリセット制御が完了するまでの間に他のクラスタから通信をゲストクラスタを制御するOSが受け付けることが可能である。

【 0 1 0 1 】

【 発明の実施の形態 】

図2は、本発明の計算機システム（SCMPシステム）の構成を示す。

同図に示す構成は、本発明のSCMPシステムの基本的な構成であり、複数の仮想計算機で運用される複数の実計算機（以下、クラスタと記す）200、300が共用メモリ（以下、SSUと記す）100を共用する複合システム（以下、SCMPシステム）である。

クラスタ200は、n個の仮想計算機（以下、ゲストクラスタと記す）220-1, ..., 2

50

20-n及び、クラスタ200のCPUに含まれる仮想計算機制御機構(以下、AVMと記す)210により構成され、各ゲストクラスタとAVM210は、論理バス71により接続されている。クラスタ300は3つのゲストクラスタ320-1, 320-2, 320-3及び、AVM310)より構成され、ゲストクラスタとAVM310は論理バス72により接続されている。

【0102】

上記の示すシステムにおいてクラスタ間の通信、制御、ダウン時の処理について以下に説明する。

【0103】

【実施例】

以下、本発明の実施例を図面を用いて詳細に説明する。

最初に第1の実施例として、共用メモリを用いた複合システムにおける計算機システムの制御を説明し、次に、第2の実施例として共用メモリを用いた複合システムにおけるダウンクラスタの制御について説明する。

【0104】

[第1の実施例]

また、図3に示す構成は、図2の構成に実計算機により運用されているクラスタ400を付加した構成である。

以下、図2、図3に示すような構成を用いて、以下の順に説明する。

【0105】

- i. 共用メモリの初期化処理
- ii. 通信時の識別子付与処理
- iii. 運用状態確認処理
- iv. クラスタ間の通信処理
- v. 通信割り込み処理
- vi. 完了確認

[i. 共用メモリの初期化処理]

まず、共用メモリの初期化処理の第1の例を説明する。

【0106】

SSU100の初期化において、SCMPシステム内で最初に立ち上がるクラスタの(IPLを実行する)オペレーティングシステム(OS)aがSSU100とのアクセスパスを接続することにより、SSU100上のロックを獲得する。これにより、他のクラスタ200のOSがIPLを実行しようとしても、ロックにより他のクラスタからのIPL操作が排他され、SSU100に格納されているデータを保証する。例えば、図3の例では、クラスタ300のOSからSSU100にIPL操作を行うと、クラスタ300のOSがSSU100にアクセスする。これにより、クラスタ300またはクラスタ400内のOSがその後にIPL操作を行ったとしても、そのIPL要求は排他される。

【0107】

また、上記のOSaがSSU100のロックを保持したまま、停止した場合には、他のOSのIPLを契機として、SSU100の初期化処理を行う。例えば、クラスタ300のOSがSSU100を獲得した状態で異常発生等により停止した場合には、他のクラスタがクラスタ300のOSの停止を発見して自クラスタからSSU100にIPL操作を行うことが可能である。なお、他のクラスタの停止を監視する方法及び他のクラスタからのIPL操作の方法については、後述する。

【0108】

図4は、本発明の第1の実施例の初期化の概要を説明するための図(その1)である。クラスタ400、500は、実計算機により運用されているクラスタであり、SSU100のアクセスパス61、62によりそれぞれ接続されている。最初にクラスタ400のOS401がSSU100にアクセスパス62を介してIPLを実行すると、クラスタ400のOS401がSSU100をロックしてSSU100の初期化の権利を獲得する。クラ

10

20

30

40

50

スタ400のOS401によるSSU100のロック後、IPLを実行したクラスタ400のOS401がSSU100のロックを保持したまま停止した場合に、クラスタ500のOS501が、クラスタ400の停止を検出し、IPLを実行する。これにより、SSU100に対するアクセスパスがロックに関係なくIPLを行うことで接続され、クラスタ500のOS501は、クラスタ400とSSU100間に接続されているアクセスパス61を物理的に切断し、自クラスタ500とSSU100間のアクセスパス62を接続する。

【0109】

これにより、停止と認識されたOSの誤動作によりSSU100データの破壊を防ぎ、データを保証するものである。上記のアクセスパスの切断については、特開平4-60750号『クラスタ停止装置』、及び特開平4-23149『二重化データ保全装置』に詳述されている。

10

【0110】

次に、AVM運用されているクラスタのOSにより他のクラスタの停止を検出した場合の動作を説明する。図5は、本発明の第1の実施例の初期化の概要を説明するため図(その2)である。

同図において、クラスタ300のゲストクラスタ320-2が、後述する停止監視機能によりクラスタ200のゲストクラスタ220-2が停止したと見做した場合、上記の第1の例のように、物理的にSSU100とクラスタ200間のアクセスパス61を切断するので、AVM運用されているクラスタ200の配下の他のゲストクラスタの接続も同時に切断されてしまうため、ゲストクラスタ320-2内のOSが通信を発行し、AVM310経由でAVM210との間で通信を行い、AVM210がゲストクラスタ220-2と接続するための論理パス71を切断する。

20

【0111】

図6は、本発明の第1の実施例の初期化処理におけるシステム構成を示す。同図に示すシステムは、図3に示す実計算機運用されるクラスタ400と、仮想計算機運用されるクラスタ200がSSU100に接続されている構成であり、他のクラスタの接続は説明の簡略化のため省略する。

【0112】

図6において、実計算機運用されるクラスタ400の実計算機制御部440は、OSとしてSSU100上の特定のメモリ領域を更新するメモリ更新部441、他のクラスタがSSU100を初期化中であるか否かを監視する初期化監視部442、初期化中の他のクラスタが停止(ダウン)していないかを監視する停止監視部443、停止監視部443において、停止しているクラスタ(ダウンクラスタ)を検出した際に、当該ダウンクラスタが実計算機か仮想計算機かを判断する仮想・実計算機判定部444、クラスタとSSU100とを接続するアクセスパスを切断するパス切断部445、自クラスタがSSU100に対して初期化する場合にIPL操作を行う初期化部446及び上記の各部を制御する制御部447より構成される。なお、メモリ更新部441、初期化監視部442、停止監視部443、仮想・実計算機判定部444、パス切断部445、初期化部446及び制御部447は各々OSである。

30

40

【0113】

仮想計算機運用されるクラスタ200は、ハードウェアである実計算機制御部240、ゲストクラスタ220と論理パス71で接続され、各ゲストクラスタを制御するAVM210及び複数のゲストクラスタ220-1、220-2、220-3より構成される。なお、各ゲストクラスタ220-1、220-2、220-3のOS詳細は、クラスタ400の実計算機制御部440のOS450の構成と同様であるので、図面上の記載を省略する。

【0114】

なお、仮想計算機により運用されている複数のクラスタがSSU100を分割して利用することがある。

上記の各部の動作を以下に示す。図7は、本発明の第1の実施例の初期化処理のフローチ

50

ャートである。

【0115】

初期化処理は、OS単位に行われ、同図の例では、ゲストクラスタが初期化の単位である。

ステップ100) ゲストクラスタのOSが初期化のためのSSU100のロックを獲得する。

【0116】

ステップ101) 初期化監視部は、既に別のクラスタがロックを獲得している場合には、ステップ103に移行し、自クラスタでロックを獲得できたならばステップ102に移行する。

10

ステップ102) 初期化部は、自クラスタの初期化処理を行う。

【0117】

ステップ103) 停止監視部は、ロックを獲得しているクラスタが生存しているかを確認する。確認は、OS間(クラスタ間)の通信により行うものとする(図8-種別1)。通信対象のクラスタアドレスはロックワード上に格納されているので、当該クラスタアドレスを取得して通信を行う。通信を行った結果、応答が有る場合には、クラスタが生存中であるので、ステップ104に移行し、応答がない場合にはステップ105に移行する。

【0118】

ステップ104) 制御部は、生存中の他のクラスタの初期化処理が完了するまで待機し、完了したらステップ108に移行する。

20

ステップ105) 仮想・実計算機判定部は、ロックワード上に格納されたクラスタアドレスに基づいて相手クラスタがAVM運用中か否かを判定し、AVM運用中であればステップ106に移行し、実(Native)クラスタとして運用されている場合にはステップ107に移行する。

【0119】

ステップ106) パス切断部は、相手クラスタ(ダウンしたクラスタ)がAVM運用中であれば、論理パスを切断して、ステップ108に移行する。

ステップ107) パス切断部は、相手クラスタが実クラスタ運用中であれば、アクセスパスを切断する。

【0120】

30

ステップ108) 初期化部は、初期化処理を実行する。

上記の処理において各判定処理は図8の内容に基づいて行うものとする。

なお、上記において、ダウンクラスタの制御方法については、第2の実施例で詳細に説明する。

【0121】

[ii. 通信時の識別子付与処理]

次に本発明の通信時の識別子付与処理について説明する。

本実施例は、例えば、図2におけるクラスタ200内のAVM210が、複数のゲストクラスタ220-1~220-nに対して、クラスタ200内の資源の割り当て等を行う際に、ゲストクラスタ固有の仮想計算機番号を付与するものである。これにより、AVM210は全てシステム内ではユニークな番号を有する。

40

【0122】

図9は、本発明の第1の実施例のゲストクラスタに対する識別子付与処理を説明するための図である。

同図(A)は、各クラスタ毎に付与されている実計算機番号と、クラスタが仮想計算機により運用されている各クラスタに含まれるゲストクラスタに付与される相対計算機番号を示す。同図(B)は、実計算機番号と相対計算機番号より生成される仮想計算機番号を示す。

【0123】

各クラスタのAVM310、410は、相対番号をOSから依頼があった時に通知する。

50

仮想番号は、OSにより実計算機番号と相対計算機番号により生成される。同図において、“x”で表されているゲストクラスタは、停止しているものとし、停止しているゲストクラスタには、仮想計算機番号は付与しない。

【0124】

例えば、クラスタ300において、AVM310のOSは、自クラスタ300内に保持する自クラスタの実計算機番号“01”と仮想計算機番号を付与しようとするゲストクラスタの相対計算機番号を讀出す。例えば、ゲストクラスタ320-3に付与する場合には、“01”+“3”により、“013”という仮想計算機番号が生成される。

【0125】

図10は、本発明の第1の実施例のゲストクラスタに対して仮想計算機番号を付与する処理を説明するためのシーケンスチャートである。以下の説明では、図9のクラスタ300の相対計算機番号“3”を有するゲストクラスタ300-3に仮想計算機番号を付与する例を用いて説明する。

10

【0126】

ステップ201) AVM310は、クラスタ300内にAVM310のゲストクラスタを設定する。図9の例では、ゲストクラスタ320-1, 320-3, 320-4, 320-5の4つのゲストクラスタを設定する。

ステップ202) クラスタ300のOSは、メモリ(図示せず)より実計算機番号と相対計算機番号を取得し、その2つの番号を合成して仮想計算機番号を生成する。

【0127】

20

ステップ203) ステップ202で生成された仮想計算機番号は、AVM310内のメモリ(図示せず)に格納する。

ステップ204) ゲストクラスタ320-3のOSがAVM310に仮想計算機番号の問い合わせを行う。

【0128】

ステップ205) AVM310は、相対計算機番号をメモリより読み出して、ゲストクラスタ320-3に転送する。

ステップ206) また、OSは、メモリに格納されている自クラスタ300内の全てのゲストクラスタの仮想計算機番号をSSU100の計算機番号用領域110に転送する。

【0129】

30

ステップ207) SSU100に接続されている他のクラスタがある場合には、自クラスタ300内の各ゲストクラスタの仮想計算機番号、実計算機番号を他のクラスタに転送する。

図11は、本発明の第1の実施例の仮想計算機番号の参照動作を説明するための図である。同図中の番号(Sxxx)と以下のステップ番号は対応するものとする。

【0130】

ステップ301) クラスタ300の相対計算機番号“4”を有するゲストクラスタ300-4より自クラスタ300内のAVM310に対して自ゲストクラスタ320-4の相対計算機番号を問い合わせる。

ステップ302) クラスタ300のAVM310はゲストクラスタ320-4に対して相対計算機番号“014”を通知する。

40

【0131】

ステップ303) ゲストクラスタ320-4内のOSは、SSU100の計算機番号用領域110にゲストクラスタ320-4の仮想計算機番号“014”を書き込む。

ステップ304) 他のクラスタ200やクラスタ400は、クラスタ300のゲストクラスタ320-4の仮想計算機番号を知りたい場合には、SSU100の計算機番号用領域110に問い合わせを行ってSSU100から情報の読み込みを行う。これにより、クラスタは、SSU100に登録されている、他のクラスタに属するゲストクラスタの仮想計算機番号を参照することが可能となるため、複数のクラスタに属するゲストクラスタ間の通信時に送信元・送信先のゲストクラスタを識別することが可能となる。

50

【0132】

また、あるクラスタから他のクラスタのゲストクラスタの仮想計算機番号を参照する他の例を示す。

図12は、本発明の第1の実施例の他のクラスタに仮想計算機番号を通知する他の例を示す。

【0133】

ステップ401) クラスタ300のAVM310は、他のクラスタ200が実計算機として運用されている場合には、直接当該クラスタ200に対してゲストクラスタ300-4の仮想計算機番号“014”を通知する。

ステップ402) クラスタ300のAVM310は、他のクラスタ400が仮想計算機として運用されている場合には、クラスタ400のAVM410に対してパラメータで存在する配下のゲストクラスタを指定して転送する。

【0134】

ステップ403) 受信したクラスタ400のAVM410は、配下の全ゲストクラスタに受信した仮想計算機番号“014”を通知する。

[iii. 運用状態確認処理]

次に、本発明の第1の実施例の運用状態確認処理について説明する。

【0135】

以下に、クラスタ内の各ゲストクラスタとSSUを接続する論理パスの状態を取得する例を説明する。

SSU100を介して複数の計算機システム間で通信を行う場合、通信先のクラスタが実計算機として運用中であるか否か、仮想計算機として運用されているか否かを知る必要がある。

【0136】

相手計算機が実計算機で運用されているのか、仮想計算機で運用されているかは他クラスタからの応答により把握する。相手計算機が実計算機で運用されている場合には、「アクセスパス/クラスタ番号/接続・ハード情報収集命令(DIAGNOSE命令(STGCNやSTGCF))」により把握できるが、相手計算機が仮想計算機で運用されている場合には、対象クラスタのAVMに対して、問い合わせ要求を発信する。この問い合わせ要求を受信したクラスタのAVMは、自クラスタ内のゲストクラスタ毎に、運用中か否かを示す情報と各ゲストクラスタとSSUを接続する論理パスが有効か否かをパラメータ域に設定し、問い合わせ元のクラスタに回答を返信する。

【0137】

本実施例では、SSU100に状態情報は格納されていないことを前提に、直接情報を必要とするAVMから他のクラスタに運用状態情報を問い合わせる。

図13は、本発明の第1の実施例の運用状態情報取得の概念を示す。クラスタ200とSSU100、クラスタ300とSSU100、クラスタ400とSSU100を接続するパスはそれぞれアクセスパス(物理パス)61、62、60であり、各ゲストクラスタ220とAVM210間を接続するパスは論理パス(仮想パス)241である。同図に示す例において、クラスタ200のゲストクラスタ220-1が運用状態情報の問い合わせ元であり、クラスタ100、300が問い合わせ先である。また、相手がAVMの場合、全ゲストクラスタの情報を問い合わせる。

【0138】

図14は、本発明の第1の実施例の運用状態情報取得時のシステム構成を示す。同図に示すシステムは、クラスタ200、300は仮想計算機により運用されているものとし、以下では、クラスタ200からクラスタ300に運用状態情報の問い合わせを行うものとして説明する。

【0139】

クラスタ200は、ハードウェアで構成される実計算機制御部240、AVM210、及び複数のゲストクラスタ220-1~220-3により構成される。

クラスタ200の配下の複数のゲストクラスタ220は、各々図15に示すように、パラメータ解析部2211と仮想計算機間通信依頼部2212を有する計算機間通信制御部221、送信元・送信先の実計算機番号と仮想計算機番号を有するパラメータ域222、情報収集依頼を行うタスク223より構成される。パラメータ域222には、送信元の実計算機番号(0)、相対計算機番号(2)、送信先の実計算機番号(1)、相対計算機番号(2)が設定されている。

【0140】

また、図14に示すクラスタ300の実計算機制御部340内に、図16に示すようなパラメータ域341が設けられ、パラメータとして、自クラスタ300内のゲストクラスタ(仮想計算機)の情報が設定されている。パラメータは、ゲストクラスタの数(4)、各ゲストクラスタの状態(OK, NO等)が設定される。

10

【0141】

図14に基づいて運用状態情報の問い合わせの動作を説明する。図17は、本発明の第1の実施例の運用状態情報の問い合わせ動作を説明するためのシーケンスチャートである。ステップ501) まず、クラスタ200のゲストクラスタ220-2において、タスク223から情報収集依頼が発行される。

【0142】

ステップ502) ゲストクラスタ220-2の計算機間通信制御部221は、パラメータ解析部2211にパラメータ域222の内容を解析するよう制御する。計算機制御部221のパラメータ解析部2211は、パラメータ域222を参照して、問い合わせを行う送信先の実計算機番号“0”と仮想計算機番号“2”及び送信元である自クラスタ200の実計算機番号“1”、仮想計算機番号“2”を取得して、計算機制御部221の仮想計算機間通信依頼部2212に転送する。

20

【0143】

ステップ503) 仮想計算機間通信依頼部2212は、取得した問い合わせ先の実計算機番号と相対計算機番号に対応するクラスタにAVM210を介して発信する。図14の例では、クラスタ300のゲストクラスタ320-2に問い合わせを行うものとする。

【0144】

ステップ504) クラスタ300の実計算機制御部340は、クラスタ200から送信された問い合わせ情報を解析する。

30

ステップ505) クラスタ300の実計算機制御部340は、図16に示す内容のパラメータ域341を有し、パラメータをAVMに渡し、AVMの返答と共に当該パラメータが依頼元に通知される。クラスタ200からの問い合わせに回答する。このとき、実計算機制御部340は、自クラスタ300の構成情報を収集し、パラメータ域341に編集する。実計算機制御部340は、AVM310に対してゲストクラスタ間に論理パスが接続されているゲストクラスタ識別子(仮想計算機番号)、及び運用中であるか否かの動作状態を問い合わせる。AVM310は、実計算機制御部340からの問い合わせにより、ゲストクラスタ間の論理パスを調査し、論理パスに接続されているゲストクラスタの動作状態を認識し、実計算機制御部340のパラメータ域341に渡されたパラメータ域341に情報を設定する。

40

【0145】

ステップ506) AVMがパラメータ域341の内容を返却する。

クラスタ200は、問い合わせの対象であったクラスタ300のゲストクラスタ320-2の内容を取得して、当該ゲストクラスタ320-2が動作中であれば(上記のステップ506において“OK”が返却された場合)、当該ゲストクラスタとの間の通信を行う。これらの処理は、図8-種別5で依頼し、種別6でAVMに返答するものである。

【0146】

上記の実施例により、複数の計算機システム間で通信を行う場合、通信先のクラスタが運用中であるか否か、クラスタ自体が運用中か、ゲストクラスタが運用中であるか等の情報を知ることが可能である。

50

[iv. クラスタ間の通信処理]

次に、本発明のクラスタ間の通信処理について説明する。計算機間の通信処理には、クラスタ上のOSとゲストクラスタを制御するAVM間の通信(OS対AVM)、あるクラスタと他のクラスタ間の通信(OS対OS)がある。OS間で通信を行う必要が生じた時、通信を実現させる上で、送信元のクラスタ内のゲストクラスタ及び送信先のクラスタ内のゲストクラスタを特定することが必要となる。また、AVM配下のゲストクラスタ内のOS間で通信を行う場合には、必然的にAVMが間に介在する。

【 0 1 4 7 】

図18は、本発明の第1の実施例のクラスタ間の通信処理を説明するための図である。図14と同一構成部分には同一符号を付している。クラスタ200はゲストクラスタ220-1~220-3を有し、クラスタ300はゲストクラスタ320-1~320-3を有し、クラスタ400はゲストクラスタを持たない実計算機として運用される。

10

【 0 1 4 8 】

以下の説明では、クラスタ200のゲストクラスタ220-2を送信元とし、クラスタ300のゲストクラスタ320-3を送信先として、SSU100を介して通信するものとする。

送信元のクラスタ200のゲストクラスタ220-2は、計算機間通信制御部221とタスク223を有し、タスク223は、計算機間通信依頼を計算機間通信制御部221に依頼する。このとき、タスク223は、パラメータ域222を参照して送信先の実計算機番号と相対計算機番号を取得する。

20

【 0 1 4 9 】

クラスタ200のAVM210は、送信先計算機(クラスタ300)が他の実クラスタであるかを見極め、他のクラスタであれば通信要求を発行し、自クラスタ内であれば、配下のゲストクラスタ内のOSに対して割り込みを発生させる。

クラスタ300の実計算機制御部340は、SSU100を介して通信要求を受け付け、送信元(クラスタ200)から送信された情報より相対計算機番号を取り出し、AVM310に割り込みを発生させる。

【 0 1 5 0 】

クラスタ300のAVM310は、実計算機制御部340から転送された相対計算機番号に基づいて、自クラスタ内のいずれのゲストクラスタに送信されているのかを判断し、対象のゲストクラスタ(ゲストクラスタ320-3)をディスパッチする。そして、対象のゲストクラスタ320-3内の計算機間通信制御部321に割り込みを発生させる。

30

【 0 1 5 1 】

ゲストクラスタ320-3内の計算機間通信制御部321は、送信元のクラスタ200より送信された情報の全てを受け取り、タスク323の通信の内容毎に用意された受信出口を起動する。

以下、一連の動作を図19に基づいて説明する。図19は、本発明の第1の実施例の計算機間の第1の通信処理動作のシーケンスチャートである。

【 0 1 5 2 】

ステップ601) 送信元のクラスタ200のゲストクラスタ220-2が有するタスク223は、送信先の実計算機番号と相対計算機番号を計算機間通信制御部221に転送し、当該番号を送信先として、通信依頼を行う。タスク223は、自ゲストクラスタ220内のパラメータ域222より送信先及び送信元(自クラスタ・ゲストクラスタ)の実計算機番号及び相対計算機番号を取得して転送する。また、通信依頼の内容は、通信命令、要求コード、送信先及び送信元(自クラスタ・ゲストクラスタ)の実計算機番号及び相対計算機番号よりなる。

40

【 0 1 5 3 】

なお、この例では、タスク223が計算機間通信制御部221に実計算機番号と相対計算機番号を渡しているが、前述の実施例のように、SSU100に問い合わせを行い、SSU100より取得する方法もある。

50

ステップ602) ゲストクラスタ220-2の計算機間通信制御部221は、クラスタ200のAVM210に対して、実計算機番号と相対計算機番号を渡す。AVM210は、実計算機番号と相対計算機番号に基づいて、送信先の実計算機(クラスタ300)を特定する。

【0154】

ステップ603) AVM210は、実計算機番号により、通信相手が自クラスタであるか、他クラスタであるかを判断する。

ステップ604) 自クラスタ200である場合には、AVM210は、自クラスタ200の相対計算機番号に対応するゲストクラスタに対して割り込み処理を行う。

【0155】

ステップ605) 他クラスタである場合には、AVM210は、送信先をクラスタ300とし、実計算機制御部240に制御を渡す。

ステップ606) 実計算機制御部240は、SSU100を介してクラスタ300に通信要求を発行する。

【0156】

ステップ607) クラスタ300の実計算機制御部340は、通信要求を受け付け、実計算機番号を参照して自クラスタに対する通信であるかを確認し、自クラスタ300宛の通信要求であれば、クラスタ200から送信された情報より仮想計算機番号を取り出し、AVM310に渡す。AVM310は、実計算機制御部340から相対計算機番号が割り込みを契機として渡されると、割り込みを発生させる。

【0157】

ステップ608) AVM310は、相対計算機番号に対応するゲストクラスタ320-3をディスパッチする。

ステップ609) AVM310は、対象のゲストクラスタ320-3内の計算機間通信制御部321に割り込みを発生させる。

【0158】

ステップ610) ゲストクラスタ320-3内の計算機間通信制御部321は、送信元のクラスタ200より送信された全情報を受け取り、タスク323の通信の内容毎に用意された受信出口を起動する。

上記の受信側のAVM310は、ゲストクラスタ320-3に通信するまでの間は、送信元に対して受信状態とならないように制御している。

【0159】

これは、SSU100を介した計算機間(OS間)の通信を行う場合には、ハードウェアで装備される通信命令(GSIGP命令)を使用する。この通信は、ハードウェアがSSU100の特定領域に通信テキストを書き込むことで、送信は完了する。送信元は相手を受信したか否かをテキストの全てを受け取った(SSU100からロードした)時点で完了と見做す。従って、受信状態とは、SSU100上からテキストの全てを取り出してしまおうと、実際にゲストクラスタOSが受信していなくても、送信元は受信状態と認識してしまうが、AVM310は、テキストの全てを取り出していないため、AVM310が当該処理を行っても送信元では受信したと認識されない。従って、クラスタ間の通信の送信元では、渡した情報が全てSSU100より送信先により読み込まれたか否かの情報により対象計算機(クラスタ、ゲストクラスタ)が通信要求内容を受信したかを判断する。詳細は、後述するv.の項で詳述する。

【0160】

なお、上記の例では、あるクラスタの1つのゲストクラスタに対する例を示したが、あるクラスタの複数のゲストクラスタまたは全ゲストクラスタに対する通信要求を発行することも可能である。この場合には、送信元のクラスタのゲストクラスタのタスクにおいて、送信先の仮想計算機番号を必要数だけ通信要求に定義すればよい。

【0161】

次に、あるクラスタのゲストクラスタから他のクラスタに通信要求を発行する場合につい

10

20

30

40

50

て説明する。図20は、本発明の第1の実施例の計算機間の第2の通信動作のシーケンスチャートである。以下の説明では、送信元をクラスタ200とし、送信先をクラスタ400として説明する。

【0162】

ステップ701) 送信元のクラスタ200のゲストクラスタ220-2が有するタスク223は、送信先の実計算機番号及び相対計算機番号を計算機間通信制御部221に転送し、当該番号を送信先として通信依頼を行う。タスク223は、自ゲストクラスタ220内のパラメータ域222より送信先及び送信元(自クラスタ・ゲストクラスタ)の実計算機番号及び相対計算機番号を取得して転送するものとする。

【0163】

ステップ702) ゲストクラスタ220-2の計算機間通信制御部221は、クラスタ200のAVM210に対して、実計算機番号及び相対計算機番号を渡す。AVM210は、実計算機番号に基づいて、送信先の計算機(クラスタ400)を特定する。

【0164】

ステップ703) AVM210は、実計算機番号により、通信相手が自クラスタであるか、他実クラスタであるかを判断し、自クラスタ200である場合には、処理を終了する。

ステップ704) 通信相手が他実クラスタである場合には、実計算機制御部240は、通信要求をSSU100を介して送信先のクラスタ400に送信する。

【0165】

ステップ705) クラスタ400の実計算機制御部440は、通信要求より実計算機番号を取り出し、自クラスタ400に対する通信であるかを確認する。

ステップ706) 実計算機制御部440は、計算機間通信制御部431に通信要求を転送する。

【0166】

ステップ707) 計算機間通信制御部431は、タスク423の通信用受信出口を起動させる。

これらの一連の処理により、クラスタのゲストクラスタと他クラスタのゲストクラスタ間の通信及びクラスタのゲストクラスタと他のクラスタ間の通信が可能となる。

【0167】

なお、上記の例に限定されることなく、送信元・送信先共にクラスタであったり、送信元がクラスタ(実計算機)であり、送信先が他のクラスタのゲストクラスタであっても通信が可能である。

[v. 通信割り込み処理]

次に、本発明の第1の実施例の割り込み処理について説明する。

【0168】

本実施例は、クラスタまたは、他のクラスタ上の仮想計算機制御部より発行されたシステム間の通信割り込み(GSIGP命令割り込み)がAVMに対して発行された場合、または、AVM内に他のクラスタから通信要求が発行された場合に、当該要求を受け付けキューイングするものである。

【0169】

第1の例は、受信側のクラスタのAVM内に通信要求をキューイングする機能を付与し、当該キューイングのキューの数が溢れた場合には、全てのクラスタに通信要求の送信先となっているゲストクラスタを通知し、送信元からの通信要求を停止させる機能を付与する。

【0170】

図21は、本発明の第1の実施例の割り込み処理の第1の例を説明するための図であり、図22は、本発明の第1の実施例の割り込み処理の第1の例の動作のシーケンスチャートである。

ステップ801) 送信元のクラスタ400から送信先をクラスタ200のゲストクラスタ

10

20

30

40

50

タ 2 2 0 -2とした通信要求として、G S I G P 命令 (a) が S S U 1 0 0 を介して発行される。

【 0 1 7 1 】

ステップ 8 0 2) 送信先のクラスタ 2 0 0 は、実計算機制御部 2 4 0 において、割り込みが反映状態ではないので、割り込みを保留する。

ステップ 8 0 3) A V M 2 1 0 が割り込み可能な状態になると、実計算機制御部 2 4 0 は、A V M 2 1 0 に割り込みを反映する。この時点で、割り込み保留が解除されるため、他のクラスタからは、割り込みが反映されたように見える。

【 0 1 7 2 】

ステップ 8 0 4) A V M 2 1 0 は G S I G P 命令 (a) を取得し、相対計算機番号により送信先のゲストクラスタを特定し、送信先のゲストクラスタ毎にキューイングする (この状態では、ゲストクラスタが割り込み反映不可状態であるとする) 。 10

【 0 1 7 3 】

ステップ 8 0 5) ここで、クラスタ 4 0 0 から、2 回目の G S I G P 命令 (b) が S S U 1 0 0 を介してクラスタ 2 0 0 のゲストクラスタ 2 2 0 -2 を送信先として発行される。

ステップ 8 0 6) 送信先のクラスタ 2 0 0 は、上記ステップ 8 0 2 と同様に、実計算機制御部 2 4 0 において、割り込みを保留する。

【 0 1 7 4 】

ステップ 8 0 7) A V M 2 1 0 に G S I G P 命令 (b) の割り込み要求を要求キューにキューイングする。 20

ステップ 8 0 8) A V M 2 1 0 は、キューイングされている先頭の通信要求をゲストクラスタ 2 2 0 -2 に反映する。以下順次、キュー待ち行列からキューを取り出し、ゲストクラスタに反映する。

【 0 1 7 5 】

ステップ 8 0 9) 3 回目の通信要求がクラスタ 4 0 0 から送信先をクラスタ 2 0 0 のゲストクラスタ 2 2 0 -2 とした G S I G P 命令 (c) が発行される。

ステップ 8 1 0) クラスタ 2 0 0 の実計算機制御部 2 4 0 は、当該 G S I G P 命令 (c) を保留する。

【 0 1 7 6 】

ステップ 8 1 1) A V M 2 1 0 は上記通信要求 G S I G P 命令 (c) を要求キューにキューイングする。 30

ステップ 8 1 2) A V M 2 1 0 は、ステップ 8 1 1 においてキューイングを行うが、ここで、キュー溢れが発生したものとする。

【 0 1 7 7 】

ステップ 8 1 3) このとき、A V M 2 1 0 は、送信元のクラスタ 4 0 0 に対して通信要求を停止するよう指示する。

なお、本ステップでは送信元のクラスタ 4 0 0 に通信要求停止指示を行っているが、他のクラスタから送信されている場合も同様に、送信元のクラスタの実計算機番号を取得し、通信要求停止指示を発行する。また、他のクラスタのゲストクラスタが送信元である場合には、実計算機番号以外に、相対計算機番号も取得して、送信元のクラスタのゲストクラスタ宛に通信要求停止指示を発行するものとする。 40

【 0 1 7 8 】

次に、通信割り込み処理の第 2 の例を説明する。

図 2 3 は、本発明の第 1 の実施例の通信割り込み処理の第 2 の例を説明するための図である。

上記第 1 の例では、送信先のゲストクラスタのキューイング中にキューが要求キューより溢れた場合には、異常発生として通信要求送信元に通信要求発行停止指示を送出する例を示したが、第 2 の例では、送信された通信要求の送信先のゲストクラスタに異常が発生した場合に、同一クラスタの他のゲストクラスタに割り込みを行う例である。

【 0 1 7 9 】

図24は、本発明の第1の実施例の通信割り込み処理の第2の例の動作シーケンスチャートである。図22のステップ810までは同一であるため、説明を省略する。

ステップ901) 送信先のクラスタ200のAVM210において、送信元のクラスタ400から指定されている送信先以外のゲストクラスタやクラスタに割り込みを行う。同図の例では、クラスタ400で指定されたゲストクラスタ220-1に対応する要求キューからキューが溢れている、または、他の異常が発生したため、他のゲストクラスタ220-2、220-3に割り込みを行う。

【0180】

ステップ902) 送信元のクラスタ400で指定された送信先のゲストクラスタに異常が発生した旨を、AVMとゲストクラスタ内の計算機間通信制御部(OS)とハンドシェイクして、GSIGP命令による通信割り込みによりクラスタ400に送信する。これにより、ゲストクラスタ220-1に異常が発生した場合に、当該ゲストクラスタ220-1に対する通信を停止するようにクラスタ400に対して通知することができる。

10

【0181】

ステップ903) 送信元のクラスタ400及び他のクラスタ(OS)では、GSIGP命令により送信先のゲストクラスタに異常が発生したことを認識する。

次に、通信割り込み処理の第3の例を説明する。

【0182】

第3の例は、新規割り込みが発生した場合に、SSU100上にキューイングしている通信要求が存在するかを確認し、存在する場合にはSSU100上にキューイングされている通信要求も処理する方法である。

20

図25は、本発明の第1の実施例の通信割り込み処理の第3の例を説明するための図である。同図に示す構成は、クラスタ400、ゲストクラスタ220-1(a)、220-2(b)により運用されているクラスタ200、クラスタ500より構成される。クラスタ400は、クラスタ200内のゲストクラスタ220-1宛にGSIGP命令を発行し、クラスタ500は、クラスタ200内のゲストクラスタ220-2(b)に対してGSIGP命令を発行するものとして以下で説明する。

【0183】

図26、27は、本発明の第1の実施例の通信割り込み処理の第3の通信動作のシーケンスチャートである。

30

ステップ1501) 送信元のクラスタ400は、クラスタ200のゲストクラスタ220-1(a)に対して通信を行うため、GSIGP命令を発行する。

【0184】

ステップ1502) クラスタ200は、割り込みが反映状態でないため、実計算機制御部240で割り込みを保留する。

ステップ1503) 送信元のクラスタ500は、クラスタ200のゲストクラスタ220-2(b)に対して通信を行うため、GSIGP命令を発行する。

【0185】

ステップ1504) 既にクラスタ200では、割り込みが保留中であるので、クラスタ500は保留を認識する。

40

ステップ1505) クラスタ500は、クラスタ220-2に対する通信要求をSSU100にキューイングする。

【0186】

ステップ1506) クラスタ400は、クラスタ200のゲストクラスタ200-1にGSIGP命令を発行する。

ステップ1507) この時点で、クラスタ200の実計算機制御部240では、まだ割り込みが保留中であるので、クラスタ400は保留を認識する。

【0187】

ステップ1508) クラスタ400は、クラスタ200のゲストクラスタ220-1(a)宛の通信要求をSSU100にキューイングする。

50

ステップ1509) ここで、ゲストクラスタ220-1(a)の割り込み反映が可能となったとする。

【0188】

ステップ1510) 実計算機制御部240がAVM210に割り込みを反映する。この時点で割り込み保留が解除されるため、他のクラスタからは、割り込みが反映されたように見える。

ステップ1511) クラスタ200のAVM210は、ゲストクラスタ220-1(a)の割り込み要求をAVM210内の割り込み要求キューにキューイングする。但し、この段階ではゲストクラスタ220-1(a)が割り込み不可の状態であるものとする。

【0189】

ステップ1512) 上記のステップ1511と同時に、割り込み反映を行うゲストクラスタ以外のSSU100を共有する全てのゲストクラスタに対する新規割り込みを各ゲストクラスタ毎にキューイングする。

ステップ1513) AVM210は、ゲストクラスタ220-1(a)の割り込み反映が可能となると、AVM210内のゲストクラスタ220-1(a)用の割り込み要求キューの先頭から順次反映する。

【0190】

ステップ1514) AVM210は、ゲストクラスタ220-2(b)の割り込み反映が可能となると、AVM210内のゲストクラスタ220-2(b)用の割り込み要求キューの先頭から新規割り込みによる通信要求を順次反映する。

ステップ1515) クラスタ200のゲストクラスタ220-1(a)のOSは、SSU100にキューイングされているゲストクラスタ220-1への通信要求が存在しているかを確認し、存在していれば、当該通信要求を取得し、ゲストクラスタ220-1に反映する。

【0191】

ステップ1516) クラスタ200のゲストクラスタ220-2のOSは、SSU100に、新規割り込みのためにキューイングされているゲストクラスタ220-2(b)への通信要求が存在しているかを確認し、存在していれば、当該通信要求を取得し、ゲストクラスタ220-2(b)に反映する。

【0192】

上記のように、種々割り込みの形態が異なっても通信要求を確実に送ることが可能である。

上記のシーケンスチャートを各割り込み発生事象でみた場合の例を図28に示す。

【0193】

G S I G P命令を発行するクラスタはクラスタ400とクラスタ500の2つのクラスタである。G S I G P命令を受信するクラスタは、クラスタ200とする。

クラスタ400で発行されたクラスタ200のゲストクラスタ220-1宛の通信要求が発生すると、G S I G P命令が発行され、割り込み反映処理側であるクラスタ200に送信される。

【0194】

ここで、クラスタ200が割り込みを保留しているものとする。

その間にクラスタ400で再度、別の通信要求が発生した場合に、G S I G P命令を発行して、クラスタ200に送信するが、まだ割り込みを保留しているため、この通信要求は、SSU100上のゲストクラスタ220-1宛の要求キューにキューイングされる。

【0195】

一方、クラスタ500において、クラスタ200のゲストクラスタ220-2(b)に対する通信要求が発生し、G S I G P命令が発行され、クラスタ200に送信されるが、まだ、割り込みを保留しているため、この通信要求もSSU100上のゲストクラスタ220-2(b)宛の要求キューにキューイングされる。

【0196】

10

20

30

40

50

ここで、クラスタ200のAVM210に割り込みが可能となり、割り込みを受け付けると、ATM210がゲストクラスタ220-1宛の割り込み要求を自AVM内の待ちキュー行列にキューイングする。そこで、ゲストクラスタ220-1のOSは、キューイングされている通信要求を順次処理していき、AVM210により反映される。

【0197】

AVM210が同時に他のゲストクラスタ220-2(b)宛の割り込み要求をキューイングし、ゲストクラスタ220-2(b)に反映させる。そこで、ゲストクラスタ220-2(b)のOSは、SSU100上にキューイングされている通信要求がなくなるまで、順次当該通信要求を処理していく。

【0198】

このように、ゲストクラスタで運用しているクラスタ(仮想計算機運用実計算機)のAVMと、各ゲストクラスタとの間の通信割り込みの状態を各々異なっても通信要求を確実に送信することが可能となる。

[vi. 完了処理(リセット処理)]

以下に説明するリセット処理は、あるクラスタ内のゲストクラスタが他のクラスタとSSU100を共有し、SSU100に具備されたシステム制御機能(以下、GSIGP命令(リセット))を使用して制御する場合において、システム制御の完了を認識する際に、クラスタまたは他のクラスタのゲストクラスタからGSIGP命令(リセット)で特定のゲストクラスタを指定し、リセットを依頼するものである。このとき、クラスタ内のAVMは、指定されたりセットを行うゲストクラスタを認識してリセットの制御を行う。なお、異常が発生したダウンクラスタの制御時におけるI/Oのリセットについては第2の実施例で後述する。

【0199】

図29は、本発明の第1の実施例のリセット処理を説明するための図である。同図において、クラスタ400がゲストクラスタ220-1、220-2により運用されているクラスタ200のゲストクラスタ220-1にリセット依頼(GSIGP命令)を発行し、クラスタ400において、リセットの完了を認識する処理である。また、クラスタ200がゲストクラスタ220-1をリセット中に、他のクラスタ500よりクラスタ200のゲストクラスタ220-2(b)にリセット依頼を発行すると、同図の例では、クラスタ400からのリセット依頼がクラスタ200上で起動されているため、クラスタ500からのリセット依頼はクラスタ200の実計算機制御部240上で保留され、ゲストクラスタ220-1のリセット処理が完了した段階でクラスタ500からのリセット依頼が起動される。

【0200】

さらに、リセットが完了した場合には、AVM210と依頼元クラスタ(OS)とがハンドシェイクして、リセットの依頼元のクラスタ400、500に対して、図8の種別8に示される内容が送信される。これにより、依頼元のクラスタ400、500は、リセットの完了を正しく認識できる。

【0201】

図30は、本発明の第1の実施例のリセット処理動作のシーケンスチャートである。動作の順序は、上記の図29に対応するものとする。

ステップ1601) クラスタ400の計算機間通信制御部(OS)は、クラスタ200のゲストクラスタ220-1に対して、GSIGP命令(リセット依頼)を発行する。

【0202】

ステップ1602) クラスタ200は、クラスタ400から送信されたGSIGP命令(リセット依頼)を受信し、実計算機制御部(ハードウェア)240で要求を保留する。

ステップ1603) クラスタ200のAVM210は、リセット要求を認識すると、実計算機制御部240に対してリセット保留の解除を指示する。

【0203】

ステップ1604) ここで、クラスタ500よりクラスタ200のゲストクラスタ220-2に対するGSIGP命令(リセット)依頼が発行される。

10

20

30

40

50

ステップ1605) クラスタ200では、クラスタ500から発行されたGSIGP命令を実計算機制御部240上にリセット保留としておく。

【0204】

ステップ1606) クラスタ200のAVM210は、ゲストクラスタ220-1のリセットが完了すると、AVM210が完了通知をクラスタ400に送信する。

ステップ1607) AVM210は、実計算機制御部240にリセット保留となっているクラスタ500からのリセット依頼を認識すると、実計算機制御部240に対してリセット保留を解除して、クラスタ200のゲストクラスタ220-2(b)のリセット処理を起動する。

【0205】

ステップ1608) クラスタ200のAVM210は、ゲストクラスタ220-2(b)のリセット処理が終了すると、AVM210が完了通知をクラスタ500に送信する。

図31は、本発明の第1の実施例のリセット処理の発生事象でみた場合の例を示す。同図において、最初にクラスタ400からクラスタ200のゲストクラスタ220-1に対してGSIGP命令(リセット依頼)を発行すると、クラスタ200のAVM210は、リセットを保留する。クラスタ400は、クラスタ200に受信された時点でGSIGP命令が成功したものと判定し、リセット完了報告の通知を待機する。クラスタ200のAVM210は、GSIGP命令を受信したことを認識すると、実計算機制御部のリセット保留状態を解除して、ゲストクラスタ220-1のリセット処理を起動して、リセット処理を行う。

【0206】

ここで、他のクラスタ500がクラスタ200の他のゲストクラスタ220-2(b)に対してGSIGP命令(リセット依頼)を発行する。クラスタ500は、このGSIGP命令がクラスタ200に受信されたことを認識すると、GSIGP命令が成功したものと判断して、完了通知を待機する。

【0207】

先に依頼を発行したクラスタ400のリセットが完了すると、クラスタ200のAVM210は、GSIGP命令(通信:ハンドシェイク)により所定の要求コード(図8-種別8)によるリセット完了通知を発行する。その後、クラスタ500のリセット依頼に対するリセット処理が完了すると、クラスタ400に発行した方法と同様にリセット完了通知を発行する。

【0208】

[第2の実施例]

本実施例では、複数のクラスタがSSUにより結合するSCMPシステムにおいて、クラスタに異常が発生した場合の処理について説明する。

本実施例では、以下の順序で説明するものとする。

【0209】

- i. AVMによるダウン通知処理
- ii. AVM運用クラスタのダウン検出処理
- iii. AVM運用実クラスタ制御処理
- iv. AVM運用クラスタのI/Oリセット処理
- v. 自実クラスタ制御時の待機処理
- vi. ゲストクラスタのセッションの閉塞時の処理
- vii. オペレータ介入の軽減処理

[i. AVMによるダウン通知処理]

図32は、本発明の第2の実施例の処理概要を示す。

【0210】

同図では、SSU100に3つのクラスタ200、300、400が接続されている構成である。クラスタ200、300は、AVM運用され、各々OSを有する3つずつのゲストクラスタを有するものとする。クラスタ400は、実クラスタでありOSにより運用さ

10

20

30

40

50

れる。A V M運用されるクラスタ300においてA V M310が回復不能な異常によりダウンしている状態を示す。このとき、A V M310は、自己のダウンをS S U100を介して、全クラスタに対して通知する。

【0211】

上記の状態を詳細に説明する。

図33は、本発明の第2の実施例のダウン通知時における各クラスタの処理を示す。同図は、図32に示すクラスタ配置と同様である。本来A V M運用されているA V M構成は同一であるが、説明の明瞭化のため、個々のクラスタ毎に説明する。

【0212】

S S U100は、各クラスタ毎に、クラスタ制御を獲得状況を記憶するクラスタ制御獲得フィールド120と、ダウンしたクラスタを記憶するダウンクラスタ管理フィールド130より構成される。

10

クラスタ制御獲得フィールド120は、図34に示すように、各クラスタ毎にS S U100内に具備され、制御クラスタアドレス格納域111、I P L世代域112、制御状態情報113より構成される。

【0213】

制御クラスタアドレス格納域111には、あるクラスタのダウンを他の動作中のクラスタが検出した場合に、検出したクラスタのアドレスを書き込む。つまり、この制御クラスタアドレス格納域111に書き込んだクラスタが、本クラスタ制御獲得フィールドのクラスタを制御する権利を取得したことになる。図34の例において、例えば、クラスタ300のアドレスを“ A A A B B B ”とした場合、クラスタ300がクラスタ200の制御権を取得したことになる。

20

【0214】

I P L世代域112は、ダウンを誤認しないようにするため、I P Lを行った世代番号が設定される。

制御状態情報113は、ダウンクラスタの制御クラスタを取得できなかったクラスタで、制御が完了した時点で行う回収処理のためのタイミングを通知するための情報が設定される。

【0215】

ダウンクラスタ管理フィールド130は、ダウンしたクラスタに関する情報を格納するフィールドであり、格納される主な情報として、

30

- ・ハードウェア監視機能が有効であるか否か；
- ・自クラスタがA V M運用中のゲストクラスタであるか実クラスタであるか；
- ・自クラスタを制御する際の時間；
- ・C P Uモデル

がある。これらの情報は、I P L時に予め登録しておき、自クラスタがダウンした時に、この情報に基づいて、他のクラスタより制御を行う。

【0216】

図33において、クラスタ200とクラスタ400は、クラスタ300のA V M310から障害発生の通知を受信するクラスタであるとする。クラスタ200は、A V M210と複数のゲストクラスタ220-1、220-2、220-3より構成される。A V M210は、通信受信部211を有し、自クラスタ200内の各ゲストクラスタ内のゲストクラスタO S（オペレーティングシステム）に対して、他のクラスタから受信したダウン通知を通知する。クラスタ200の配下の各ゲストクラスタ220-1、220-2、220-3は、各々、ゲストクラスタO S250を有し、各O S250は、ダウン通知受信部251、ダウン認識部252、ダウンクラスタ制御部253、資源回収処理部254及びホットスタンバイ処理部255より構成される。

40

【0217】

ダウン通知受信部251は、他のクラスタからのダウン通知をA V M210の通信送受信部211を介して受信し、通信のパラメータ解析を行う。

50

ダウン認識部 252 は、ダウクラスタの認識を行うと共に、ダウン理由の認識を行う。

【0218】

ダウクラスタ制御部 253 は、以下に示す処理を行う。詳細な動作は後述する。

- ・ダウクラスタの運用形態の認識
- ・ダウンが発生した時点において AVM が実行中か否かの検査
- ・配下のゲストクラスタがダウンしているか否かの検査
- ・ダウクラスタの制御権を獲得
- ・ダウクラスタの CPU を停止
- ・リセットが有効な状態であるかの検査（ハードウェア監視機能が有効であるかの検査）
- ・ダウクラスタの入出力のリセット
- ・ダウクラスタのダンプ採取

10

ホットスタンバイ処理部 255 は、ダウクラスタ制御部 253 において、ダウクラスタの制御の完了、または、オペレータ介入によるダウクラスタの制御の完了を以てホットスタンバイを行う。

【0219】

ホットスタンバイの具体的な処理として、ホットスタンバイを行うクラスタを特定し、回収処理を行う。回収処理として、DASD 上の資源、SSU 上の資源、システムが有する情報（ゲストクラスタがダウンしているか否か等）等をリセットする等の処理がある。

【0220】

次に、クラスタ 300 は、AVM 310 においてダウンが発生すると、当該ダウンを検出するダウン認識処理部 3110 と、ダウンの発生を SSU 100 を介して他のクラスタに通知するダウン通知処理部 3120 を AVM 310 内に有する。

20

【0221】

クラスタ 400 は、ゲストクラスタを持たない実クラスタであり、当該クラスタの OS 450 は、上記のクラスタ 200 と同様の構成のダウン通知受信部 451、ダウン認識部 452、ダウクラスタ制御部 453、資源回収処理部 454、ホットスタンバイ処理部 455 を有する。これらの各部は、クラスタ 200 の各部と同様の処理を行うため、説明を省略する。

【0222】

図 35 は、本発明の第 2 の実施例のダウンの発生を通知・認識動作を示すシーケンスチャートである。同図において、SSU 100 の記載は省略する。

30

ステップ 1710) クラスタ 300 の AVM 310 において何等かの回復不能な異常が発生したとする。

【0223】

ステップ 1720) クラスタ 300（以下ダウクラスタという）のダウン認識処理部 3110 は、ダウンが発生したことを認識して、ダウン通知処理部 3120 に通知する。

ステップ 1730) ダウクラスタ 300 の AVM 310 のダウン通知処理部 3120 は、SSU 100 を介してクラスタ 200 及びクラスタ 400 に通知する。通知の方法は、第 1 の実施例で説明したようにクラスタ間の通信機能を用いて通知する。なお、ダウンを他のクラスタに通知する際に、他のクラスタが AVM 運用中である場合には、当該他のクラスタの AVM の配下のゲストクラスタに通知するものとする。

40

【0224】

ステップ 1740) クラスタ 200 の AVM 210 の通信受信部 211 は、ダウクラスタ 300 からダウンの通知を受信すると、AVM 運用している配下のゲストクラスタ 220-1、220-2、220-3 のダウン通知受信部 251 に通知する。

【0225】

ステップ 1750) ダウ通知受信部 251 は、ダウン通知時に渡されるパラメータを解析する。

ステップ 1760) ダウ認識部 252 は、パラメータ解析の結果によりダウンしたクラスタ及びダウンした理由を特定する。

50

【0226】

ステップ1770) クラスタ制御部253は、ダウクラスタ300に対する制御を行う。詳細は、図36で詳細に説明する。

ステップ1780) クラスタ200のホットスタンバイ処理部255は、クラスタ制御部253の処理によりダウクラスタの制御が完了した場合、または、オペレータの介入によりダウクラスタ300の制御が完了した場合に、ホットスタンバイを実現する。

【0227】

ステップ1790) 本ステップからステップ1820に関しては、実クラスタ400がステップ1730において、ダウクラスタ300からダウ通知を受信した以降は、上記のステップ1750～ステップ1780の処理と同様であるので説明を省略する。

10

【0228】

図36は、本発明の第2の実施例のダウクラスタの制御処理のフローチャートである。以下に示す処理は、上記のステップ1770及びステップ1810に対応する処理であり、クラスタ200内のゲストクラスタ220がダウクラスタ300に対して行う例を示す。ステップ1810の処理ではクラスタ400が行うことになるが、クラスタ400では、GSIGP命令を発行せず、クラスタ200(クラスタ200内のゲストクラスタ内の1つ)が行う制御の完了を待つ。

【0229】

ステップ1771) まず、ダウクラスタ制御部253は、ダウクラスタ300のAVM310が運用中であるか否かを判定し、運用中である場合にはステップ1772に移行する。運用中でない場合には、当該ダウクラスタの制御処理を終了する。

20

【0230】

ステップ1772) AVM310が運用中である場合には、ダウ通知がAVM310自身の自己申告ダウであるか否かを判定し、AVM310の自己申告ダウである場合にはステップ1773に移行し、そうでないならば、運用中の実クラスタに対する制御を行う。

【0231】

ステップ1773) ダウクラスタ300のAVM310が運用中であり、かつ、自己申告ダウである場合には、ダウクラスタ300の配下のゲストクラスタを特定し、当該ゲストクラスタ内で生存中のものをダウ状態とする。

30

ステップ1774) ダウクラスタ制御部253は、ダウクラスタ300のCPUの停止や、I/Oリセット、ダンプ採取のためにダウクラスタ300の制御権を取得する。ダウクラスタ300の制御権の取得は、SSU100上のクラスタ制御獲得フィールド110に前述のパラメータ解析により取得したダウクラスタのアドレスを書き込むことにより獲得できる。なお、ダウクラスタ制御権を取得するクラスタは、GSIGP命令を発行するクラスタであり、SCMPシステム内で1クラスタのみである。他のクラスタは、制御権を有するクラスタの制御が完了するのを待機することになる。

【0232】

ステップ1775) ダウクラスタ制御権を獲得したクラスタ200は、対象となるダウクラスタ300のCPUを停止した後、SSU100を参照してダウクラスタ300のリセット機能が有効か否かの判定を行う。ハードウェア(サービスプロセッサ)は、各クラスタが初期化時において、ハードウェア(SVP)監視機能が有効か無効であるかをSSU100上に記録しておく。

40

【0233】

ステップ1776) ダウクラスタ制御部253は、ダウクラスタ300用のダウクラスタ管理フィールド120を参照して、ダウクラスタ300のリセット機能が有効である場合には、ステップ1777に移行し、無効の場合には、ステップ1778に移行する。

【0234】

ステップ1777) リセット機能が有効である場合には、ダウクラスタ300が、A

50

VM運用中の実クラスタであっても有効な強制リセットをハードウェアに対してSVPに指示し、ステップ1779に移行する。なお、一般的には、SVP監視機能は“有効”に設定しておくものとする。

【0235】

ステップ1778) また、リセット機能が無効である場合には、ハードウェアによる強制リセットは行わず、オペレータ介入によりI/Oリセットを行い、処理を終了する。この場合、オペレータ用にディスプレイ装置(図示せず)に介入要求を表示する等してオペレータに介入操作を依頼する。

【0236】

ステップ1779) ダウンクラスタ300のダンプを採取し、障害箇所の検出作業を行う。 10

[ii. AVM運用クラスタのダウン検出処理]

次に、AVM運用中のクラスタにおけるダウンを他のクラスタで認識する処理を説明する。

【0237】

図37は、本発明の第2の実施例の他のクラスタからダウンクラスタを認識する処理概要を説明するための図である。

同図に示す例は、SSU100に接続されるクラスタ200、300、400のうち、クラスタ300のゲストクラスタ320-2がダウンした場合の例である。ゲストクラスタ320-2のダウンは、クラスタ200の各ゲストクラスタ220-1、220-2、220-3及びクラスタ300のゲストクラスタ320-1、320-3、クラスタ400で検出される。 20

【0238】

このとき、ダウンしたクラスタ300のゲストクラスタ320-2に対する制御権は、ダウンを認識した全クラスタが、SSU100のクラス制御獲得フィールド110にシリアライズすることで、獲得可能である。図37の例では、クラスタ200のゲストクラスタ220-2がクラスタ300のゲストクラスタ320-2の制御権を獲得したものである。

【0239】

図38は、本発明の第2の実施例の他のクラスタのゲストクラスタのダウン認識動作のシーケンスチャートである。AVM運用のクラスタ300のゲストクラスタ320-2でダウンが発生し、当該クラスタの制御権をクラスタ200内のゲストクラスタ220-2が獲得した場合について説明する。 30

【0240】

ステップ2010) クラスタ300のゲストクラスタ320-2にダウンが発生する。
ステップ2020) ゲストクラスタ320-2は、ダウン通知を各クラスタに通知する。なお、ダウン通知は、互いに監視している場合には不要である。

【0241】

ステップ2030) ゲストクラスタ220-2のダウン通知受信部251は、パラメータを生成する。 40

ステップ2040) ゲストクラスタ220-2のダウン通知受信部251は、ダウンクラスタとそのダウン理由を特定する。

【0242】

ステップ2050) ゲストクラスタ220-2に対応するAVM210は、ダウン通知をSSU100を介して、クラスタ300のAVM310とクラスタ400のOS450に通知すると共に、自クラスタ200のAVM210を介して自クラスタ200の他のゲストクラスタ220-1、220-3にも通知する。

【0243】

ステップ2060) ゲストクラスタ220-2のダウンクラスタ制御部253は、ダウンクラスタ320-2が実計算機(Native)運用であるかAVM運用であるかを判断する 50

。ステップ2070) ダウンクラスタ3202がAVM運用であるとき、ゲストクラスタ220-2がダウンクラスタ320-2の制御権を獲得する。

【0244】

なお、上記の図38に示す動作は、クラスタ200のゲストクラスタ220-2について述べているが、他のクラスタでも同様の処理を行うものとする。

次に、ダウン通知を受信したクラスタ200内のゲストクラスタ220-2の動作を説明する。

【0245】

図39は、本発明の第2の実施例のダウン通知を受信した際のダウンクラスタの制御動作を示すシーケンスチャートである。以下の説明では、クラスタ300のゲストクラスタ320-2がダウンし、その通知をクラスタ200内のゲストクラスタ220-2が受信した場合を例として説明する。

【0246】

ステップ2080) クラスタ200のゲストクラスタ220-2のダウンクラスタ制御部253は、AVM210の通信受信部211に対して、ダウンクラスタのCPU停止並びにI/Oリセットを依頼する。

ステップ2090) これにより、AVM210の通信送受信部211は、クラスタ300のAVM310に対してゲストクラスタ320-2に対するCPU停止並びにI/Oリセットを依頼する。これは、CPU停止のG S I G P命令を発行すると、実クラスタ300のCPUが停止してしまうため、相手のクラスタのAVM310に依頼するものである。

【0247】

ステップ2100) クラスタ300のAVM310は、ダウンゲストクラスタ320-2に対してCPU停止を行う。

ステップ2110) この間、クラスタ200では、タイマ監視を行う。タイマ監視が必要な理由は、相手のAVMが行うCPU停止の処理時間を考慮しても、AVMシステムに異常が発生した場合には、応答が通知されない場合があるためである。タイマ監視により、タイムオーバーとなった場合には、相手のクラスタのAVM310がダウンしたものと認識し、AVM運用クラスタ配下のゲストクラスタを全てSCMPシステムから切り離す。なお、配下のゲストクラスタ個々に処理を行うと、制御完了までに時間がかかるため、実クラスタ300への制御を行うようにする。

【0248】

ステップ2120) CPU停止後、クラスタ200のゲストクラスタ220-2に対してクラスタ間通信機能を用いて完了の通知を行う。なお、クラスタ間通信機能については、前述の第1の実施例で説明している。

ステップ2130) リセット完了通知がない場合には、クラスタ200のダウンクラスタ制御部253は、SSU100のダウンクラスタ管理フィールド120に対して、ハードウェアによる強制的な制御が可能であることを示すリセットが有効であるか無効であるかを検査する。

【0249】

ステップ2140) 有効である場合には、AVM310の通信受信部311にI/Oリセットを依頼する。これにより、通信受信部311は、クラスタ300のAVM310に対してI/Oリセット要求を依頼する。

ステップ2150) クラスタ300のAVM310は、ダウンクラスタ320-2のI/Oをリセットする。

【0250】

ステップ2160) この間クラスタ200では、上記のステップ2110と同様にタイマ監視を行う。

ステップ2170) クラスタ300のAVM310がダウンクラスタ320-2のI/O

10

20

30

40

50

Ｏリセットが完了すると、完了通知をクラスタ２００に通知する。

【０２５１】

ステップ２１８０）クラスタ２００のダウクラスタ制御部２５３は、ダウクラスタ３２０－２のダンプを採取する。

ＯＳが検出するＡＶＭのダウは、以下に示す動作により検出される。

１ ゲストクラスタがダウする。

【０２５２】

２ ゲストクラスタに対して停止及びリセットをダウゲストクラスタを管理するＡＶＭのＯＳに依頼する。

３ ＡＶＭのＯＳから停止及びリセットの完了／失敗の通知が時間内に無かった場合にＯＳがＡＶＭのダウを検出する。 10

【０２５３】

なお、２においてＡＶＭがゲストクラスタに対して行う制御は、ＣＰＵ停止とＩ／Ｏリセットと一緒にＡＶＭに依頼し、ＡＶＭがＣＰＵ停止、Ｉ／Ｏリセットと共に完了／失敗した時点で通知する。

[iii. ＡＶＭ運用実クラスタ制御処理]

ＡＶＭ運用実クラスタの制御として、ＡＶＭがダウ（正常停止を含む）した場合、当該クラスタの配下のゲストクラスタをダウとして処理（前述のステップ１０７３の処理）する場合について説明する。

【０２５４】

ＡＶＭ運用中の実クラスタをダウと認識すると、配下の生存中のゲストクラスタの数、クラスタアドレスを認識し、対象ゲストクラスタを管理する情報（共用メモリ上及び主記憶上に保持）をダウ状態に書き換える。 20

その後、ＡＶＭ運用中の実クラスタに対して、ＣＰＵ停止及びＩ／Ｏリセットを行い、ダウした実クラスタの制御を行う。

【０２５５】

さらに、ホットスタンバイのため、ダウした実クラスタへの制御が完了した後、配下のゲストクラスタが保持する資源（共用ＤＡＳＤや共用メモリ上に保持する資源）の解放を行う。

図４０は、本発明の第２の実施例のＡＶＭ運用クラスタにおいて、ＡＶＭがダウした場合の処理を説明するための図である。同図において、クラスタ３００のＡＶＭ３１０自体にダウが発生すると、クラスタ２００内のゲストクラスタ２２０－１～２２０－３の何れかが上記のii.の項目の処理により、クラスタ３００のＡＶＭダウを認識し、ＡＶＭ３１０配下のゲストクラスタ３２０－１，３２０－２，３２０－３が生存中であっても、強制的にダウ状態にする。即ち、ＡＶＭ３１０がダウすると、配下のゲストクラスタ３２０－１，３２０－２，３２０－３がダウしていなくとも、ＡＶＭのＯＳがダウするため、配下のゲストクラスタは動作することができないため、ダウ状態とする。これにより、他のクラスタ２００では、個々のゲストクラスタに対して制御する必要がないため、ホットスタンバイの高速化が図れる。また、個々に、ＣＰＵ停止やＩ／Ｏリセット等の制御を行わなくとも実ＣＰＵの停止、実クラスタのＩ／Ｏリセットを行うため、個々の制御が不要となる。 30 40

【０２５６】

[iv. ＡＶＭ運用クラスタのＩ／Ｏリセット処理]

ＡＶＭ運用中のクラスタがダウし、そのクラスタのＣＰＵ停止やリセットを行う場合について説明する。

従来は、ＯＳからＧＳＩＧＰ命令を発行すると、ＡＶＭ運用のクラスタ受け付けのゲストクラスタに対して制御しているため、ＧＳＩＧＰ命令のＩ／Ｏリセット要求は、ゲストクラスタに対してＣＰＵ停止とＩ／Ｏリセットを行うための命令として用いられている。従って、ハードウェアは、ＡＶＭ運用されたクラスタに対し、ＧＳＩＧＰ命令のＩ／Ｏリセット要求が発行されると、ＡＶＭのＯＳが制御するものと認識し、ハードウェアは実質的 50

にI/Oリセットを行わない。

【0257】

そこで、本発明では、AVM運用のクラスタ自身のダウンをOSが認識し、実クラスタに対してのCPU停止、I/Oリセットを行うように構成する。この結果、従来のGSI GP命令のI/Oリセット要求にAVMシステム自身(実クラスタ)のリセットを行う要求が追加され、ハードウェア(SVP)も、これを認識した場合には、AVMが行うのではなく、ハードウェアが実クラスタのリセットを行う。

【0258】

図41は、本発明の第2の実施例のAVM運用の実クラスタのI/Oリセット処理を説明するための図である。各クラスタ200、300毎に、サービスプロセッサ290、390を有する。サービスプロセッサ290、390は、ハードウェアからのリセットを行う機能や回線異常等をサポートする機能を有し、一般にメンテナンスサポートを行うプロセッサである。共用メモリ100上には、ハードウェア管理情報領域130を有し、各クラスタ毎に、ハードウェアの状態情報をIPL時に登録しておく。ハードウェアの状態情報の内容としては、各サービスプロセッサ290、390による監視機能が有効であるか無効であるかが登録される。ここで、有効が指定されている場合には、他のクラスタからの制御により強制的にクラスタをリセットすることができるが、無効が指定されている場合には、他のクラスタからの制御要求が入力されても、リセットを行わない。

【0259】

図41において、クラスタ300のAVM310においてダウンが発生している場合に、クラスタ200内のゲストクラスタ220-1~220-3のいずれかが、そのダウンを認識したとする。このとき、クラスタ200内のゲストクラスタ220-1~220-3は、SSU100上のハードウェア管理情報領域130にアクセスし、クラスタ300の監視機能が有効であるか否かを参照する。ここで、監視機能が有効である場合には、OSがまず、CPU停止の依頼をサービスプロセッサ390に行い、サービスプロセッサ390によるCPU停止が完了した後、強制リセットを指示する。サービスプロセッサ390は、I/O要求が発行されないことを保証してからI/Oリセットを行う。

【0260】

なお、監視機能の有効/無効は、OSがSSU100上に記録する。図41の例において、クラスタ300の中で最初にOSのIPLを行うクラスタ(例えばクラスタ320-1)が自実クラスタのサービスプロセッサ監視機能が有効か無効かをSSU100上に記録する。つまり、AVM運用された実クラスタ内で最初にIPLされるゲストクラスタのみがSSU100上に記録する。

【0261】

[v.自実クラスタ制御時の待機処理]

本制御処理は、OSがAVMのダウンを検出した場合の制御方法に関する。

OSが検出するAVMのダウンは、ゲストクラスタのダウンがあり、そのゲストクラスタへの制御を行い、AVMから停止ならびにリセットの完了/失敗通知が行われない場合に検出される。

【0262】

図42は、本発明の第2の実施例の自実クラスタ制御時の待ち制御処理を説明するための図である。

1 まず、ゲストクラスタ220-2がダウンしたとする。

2 AVMからゲストクラスタ220-2に対する停止及びリセットの完了/失敗通知がないため、実クラスタ200のAVMがダウンしたとOSが認識する。ここで言うOSとは、クラスタ200のゲストクラスタ220-1、220-3とクラスタ300内のゲストクラスタ320-1~320-3の5つのクラスタのOSを指す。

【0263】

3 一般に、ダウンクラスタへの制御は、SCMPシステム内のどれか1クラスタでのみ行うため、ダウンを認識したクラスタ(OS)からSSU上のダウンクラスタ制御権

10

20

30

40

50

獲得フィールドをシリアルライズ（更新）する。

4 最初にクラスタ200のAVMのダウンをゲストクラスタ220-1や220-3が認識した場合、ダウンクラスタ制御権獲得を待機する。

【0264】

5 ゲストクラスタ220-1や220-3が待機することにより、クラスタ300ないのいずれか1つのゲストクラスタでクラスタ200の制御権を獲得することが可能となる。従って、クラスタ200の制御が完了する。

6 また、クラスタ300が存在していない場合には、ゲストクラスタ220-1やゲストクラスタ220-3が制御権を獲得することにより、自実クラスタ200に対し、CPU停止要求のGSIGP命令を発行するため、実CPUが停止する。

10

【0265】

このように、自実クラスタのAVMのダウンを自実クラスタ配下のゲストクラスタが認識し、制御権の獲得を待機する理由は、この処理を実施しないと、実クラスタ200に対する制御（CPU停止は完了、I/Oリセットは未完了）が完了することはなく、クラスタ300内の各OSではオペレータ介入メッセージを出力し、オペレータの処置なしでは、ホットスタンバイが実現できなくなってしまうからである。

【0266】

[v. ゲストクラスタのセッションの閉塞時の処理]

AVMのゲストクラスタがセッション閉塞（DEACTIVATE）した場合、当該クラスタのAVMから他のクラスタに対して閉塞状態となった旨を通知し、通知を受けたクラスタでは、当該通知よりどのゲストクラスタが閉塞状態となっているのかを認識する。但し、ゲストクラスタが閉塞となったクラスタのAVMは、ダウンしていないため、当該閉塞状態のゲストクラスタのCPU停止及びI/Oのリセットを行う。

20

【0267】

図43は、本発明の第2の実施例のゲストクラスタのセッション閉塞時の処理概要を説明するための図である。同図において、クラスタ300のゲストクラスタ320-2が閉塞状態となっている。ここで、クラスタ300のAVM310は、配下のゲストクラスタ320-2が閉塞状態となった旨を認識し、配下のゲストクラスタ320-2のCPU停止及びI/Oリセットを処理を行い、他のクラスタ200内OS及び400内OS及びゲストクラスタ320-1、320-3にゲストクラスタ320-2の閉塞を通知する。なお、オペレータによりログオフされた時点で閉塞状態になるが、ゲストクラスタの閉塞を制御しているのは、AVMであるので、ゲストクラスタのOSは、閉塞に関与していない。

30

【0268】

図44は、本発明の第2の実施例のゲストクラスタのセッション閉塞時の処理における各クラスタの構成図である。

同図において、図43と同一構成部分には、同一符号を付す。

図44において、クラスタ300のゲストクラスタ320-2が閉塞状態となっている。このとき、AVM310のDEACT（閉塞）認識部315が“DEACT（閉塞）コマンド”を認識する。

【0269】

AVM310のDEACT処理部316は、ゲストクラスタ320-2のCPUを停止し、I/Oをリセットする。さらに、通信送受信部311を介して、他のクラスタ200内のOS、400内のOS並びに300内の他ゲストクラスタに配下のゲストクラスタ320-2の閉塞状態を通知する。

40

【0270】

クラスタ200及びクラスタ400の通信受信部211、411は、クラスタ300からの通知により、ゲストクラスタ320-2の閉塞状態を認識する。なお、閉塞状態の場合には、閉塞となったゲストクラスタを配下とするクラスタ300のAVMによりCPUの停止及びI/Oリセットを行うため、既にゲストクラスタ320-2の制御は完了している。従って、GSIGP命令は発行しない。

50

【0271】

従って、図43において、AVM310からゲストクラスタ320-2の閉塞の通知を、各クラスタ200内OS、400内OS及び他ゲストクラスタが受け取った時にゲストクラスタ320-2が閉塞状態であると認識する。また、リセットは、クラスタ300のAVM310自体で、配下のゲストクラスタ320-2のCPU停止やI/Oリセットを行う。つまり、ゲストクラスタ220-1~220-3、320-1、320-3、420-1~420-3のいずれかのゲストクラスタが制御権を獲得しても、各々のクラスタのダウンクラスタ制御部ではSIGP命令を発行しない。

【0272】

[vi. オペレータ介入の軽減処理]

オペレータ介入の軽減処理は、クラスタへの制御処理(リセット、CPU停止等)がハードウェアの故障等により失敗した場合、オペレータの介入により、対象クラスタの制御を代替実行するためのオペレータ介入メッセージを消去する処理である。オペレータ介入の軽減の契機としては、

- ・クラスタが再度IPL処理を行った場合；
- ・オペレータの処置完了による応答があった場合；
- ・ゲストクラスタに対するオペレータ介入メッセージの場合、対象ゲストクラスタが閉塞である場合；
- ・ゲストクラスタに対するオペレータ介入メッセージの場合、対象ゲストクラスタを含むAVM運用中、実クラスタがダウンした場合；

があり、このような場合には、配下のゲストクラスタが複数であっても複数個のメッセージを削除する。

【0273】

図45は、本発明の第2の実施例のオペレータ介入メッセージが出力されている状態を示す。同図中、内の数字は、制御順序を示す。

同図において、クラスタ200とクラスタ400は、実クラスタとして運用されているクラスタであり、クラスタ300は、AVM運用されているクラスタである。

【0274】

ダウンクラスタへの制御は、SCMPシステム内で1つのクラスタのみ行う。従って、他のクラスタ(OS)では、その制御が完了するのを待つことになる。制御が完了した場合、または、失敗した場合は、SSU上に制御クラスタが表示し、制御権のないクラスタは、その表示情報を参照する。制御に失敗した場合は、全クラスタでオペレータ介入メッセージをコンソールに出力する。この処理は、各クラスタが隣接していない場合に、最初に気づいたオペレータが対応することが可能となる。

【0275】

ここで、クラスタ200のOSがクラスタ300のゲストクラスタ320-1の制御を失敗した場合に、クラスタ200に接続されている表示装置291上に、ゲストクラスタ320-1に関するオペレータ介入メッセージを表示する。さらに、ゲストクラスタ320-2の制御も失敗すると、表示装置291上にゲストクラスタ320-2のオペレータ介入メッセージを表示する。

【0276】

図46は、本発明の第2の実施例のオペレータ介入抑止時の状態を示す。

1 同図において、ゲストクラスタ320-1、320-2に対する制御が失敗している。従って、生存中のクラスタ(200, 320-3, 400)では、ゲストクラスタ320-1、320-2に対するオペレータ介入メッセージが出力されている。

【0277】

2 1の状態の時、クラスタ300のAVMがダウンしたものとす。

3 クラスタ200やクラスタ400、ゲストクラスタ320-3でAVMのダウン(2)を認識する。

4 各クラスタ(200, 320-3, 400)では、SSU上の情報に基づいて実

10

20

30

40

50

クラスタ300配下のゲストクラスタに対するオペレータ介入メッセージを出力しているかを判断する。

【0278】

5 4 において出力されていると判定された場合には、オペレータ介入メッセージを消去する。

また、上記の例は、ゲストクラスタに対するオペレータ介入メッセージを実クラスタを制御する際に消去しているが、同様の処理を閉塞時でも実現できる。図46において、クラスタ300のゲストクラスタ320-1がダウンし、何等かの理由により失敗し、OSが失敗を認識した場合、他の生存中のクラスタでは、ゲストクラスタ320-1に対するオペレータ介入メッセージが表示されている。この状態で、ゲストクラスタ320-1が閉塞し、AVM310が全クラスタに閉塞状態の旨を通知し、通知を受けたクラスタ200、400、320-2、320-3では、閉塞による制御契機であり、対象ゲストクラスタ320-1に対して、オペレータメッセージが表示されている場合には、当該メッセージを消去する。

【0279】

なお、クラスタ300のゲストクラスタ320-1に対する制御が失敗し、かつ、クラスタ400がAVM運用され、ゲストクラスタ420-1でもダウンし、やはり制御が失敗した状態となった場合には、他の生存中のクラスタ200の表示装置には、2つのオペレータ介入メッセージが表示される。このとき、クラスタ200のOSは、クラスタ300のAVMにダウンが発生した場合には、SSU100のダウンクラスタ管理フィールド130から実クラスタ300のゲストクラスタの状態を管理する情報を取得する。さらに、クラスタ300の配下のゲストクラスタ内にオペレータ介入メッセージを出力しているゲストクラスタが存在するかを判定し、存在する場合には、どのゲストクラスタに対するオペレータ介入メッセージを表示しているかを判断する。ここで、オペレータ介入メッセージを表示していれば、クラスタ300の制御依頼を契機として、クラスタ300の配下のゲストクラスタに対して出力されているオペレータ介入メッセージを消去する。

【0280】

上記により、ホットスタンバイ処理部では、オペレータの介入を軽減してホットスタンバイ処理が可能となる。

なお、本発明は、上記の実施例に限定されることなく、特許請求の範囲内で種々変更・応用が可能である。

【0281】

【発明の効果】

本発明は、複数のクラスタを共用メモリにより結合するSCMPシステムにおいて、1つ以上のクラスタが仮想計算機システムとして運用される時に、従来は、1つの実クラスタ上の1つのゲストクラスタしか稼働できなかったが、本発明によれば、複数のクラスタにより負荷分散をして処理しなければならないような大規模なシステムをより柔軟に構築することが可能となり、SCMPシステムを構築ことができ、複数のゲストクラスタが種々処理を実行することが可能となる。

【0282】

また、本発明は、共用メモリを介して情報交換を行うシステムにおいて、初期化状態の共用メモリを最初に起動させたクラスタが初期化する処理を行う場合に、他のクラスタが同時に初期化処理を行わないように排他制御することが可能であると共に、アクセスパスの切断や再接続を与えずに、共用メモリの初期化処理の競合による誤動作を防止する。

【0283】

また、本発明は、1つ以上のクラスタに1つ以上の仮想計算機により運用されているクラスタが共用メモリに接続され、通信を行う場合に、複数のゲストクラスタを一意に特定できる。

また、本発明は、各ゲストクラスタと共用メモリを接続する論理的なパスの状態や、各クラスタ、ゲストクラスタの運用状態を知ることが可能である。

【0284】

また、本発明は、実計算機により運用されるクラスタ、仮想計算機により運用されているクラスタ間において、通信要求の送信時にどのゲストクラスタ、どのクラスタに対して通信要求を発行するのかを特定することが可能である。

また、本発明は、仮想計算機により運用されているクラスタのゲストクラスタの通信の割り込み状態が各々異なっても通信要求を確実にゲストクラスタに反映させることが可能である。

【0285】

また、本発明は、複数のゲストクラスタで共用メモリを共用した場合でもリセット要求の完了が発行元に正しく認識できる。そのため、複数のゲストクラスタと実計算機のクラスタとで共用メモリを共用するシステムにおいて、各OS間のホットスタンバイシステムによる切り替えが可能となる。

【0286】

また、本発明は、AVM運用されているクラスタ内で発生したダウンをAVMから、共用メモリに接続されている各クラスタに通知することにより、他のクラスタよりダウンクラスタの制御を行うことが可能となる。

また、本発明は、AVM運用されているクラスタ内で発生したダウンを他のクラスタから認識することができるため、自動的に認識した他のクラスタからダウンクラスタを制御することが可能である。

【0287】

また、本発明は、AVM運用クラスタのAVM自体がダウンした場合に、当該クラスタをリセットしなければならぬが、このとき、AVM運用のクラスタの配下の全てのゲストクラスタをダウン状態とすることにより、他のクラスタからのタイマ監視による制御の時間が短縮される。

【0288】

また、各クラスタのハードウェアの運用情報を登録しておき、他のAVM運用のクラスタからダウンクラスタの制御を行う場合には、当該運用情報を参照して、各々のクラスタに付設されているサービスプロセッサに指示することにより、実クラスタのみならず、AVM運用されているクラスタであっても強制的なハードウェアによるリセットが可能となる。

【0289】

また、配下のゲストクラスタのセッション閉塞時にAVM自体でゲストクラスタのリセット等の制御を行い、他のクラスタにセッション閉塞を通知することにより、他のクラスタに対する通知及びリセットのための時間が短縮される。

また、オペレータ介入メッセージを当該介入処理が終了した時点で自動的に消去することにより、オペレータの介入を軽減し、ホットスタンバイを実現できる。

【0290】

このように、本発明によれば、実計算機内の2台以上の仮想計算機が共用メモリを介して複数計算機システム間での通信が可能となることにより、種々の通知やリセット等の制御が容易になる。

【図面の簡単な説明】

【図1】本発明の原理構成図である。

【図2】本発明の計算機システム(SCMPシステム)の構成図である。

【図3】本発明の第1の実施例の計算機システム(SCMPシステム)の構成図である。

【図4】本発明の第1の実施例の初期化の概要を説明するための図(その1)である。

【図5】本発明の第1の実施例の初期化の概要を説明するための図(その2)である。

【図6】本発明の第1の実施例の初期化処理におけるシステム構成図である。

【図7】本発明の第1の実施例の初期化処理を説明するためのフローチャートである。

【図8】本発明の各要求コードを示す図である。

【図9】本発明の第1の実施例のAVMに対する識別子付与処理を説明するための図であ

10

20

30

40

50

る。

【図10】本発明の第1の実施例のゲストクラスタに対して仮想計算機番号を付与する処理を説明するためのシーケンスチャートである。

【図11】本発明の第1の実施例の仮想計算機番号の参照動作を説明するための図である。

【図12】本発明の第1の実施例のあるクラスタから他のクラスタへゲストクラスタの仮想計算機番号を通知する他の例を示す図である。

【図13】本発明の第1の実施例の運用情報取得の概念図である。

【図14】本発明の第1の実施例の運用状態情報取得時のシステム構成図である。

【図15】本発明の第1の実施例のゲストクラスタ内の構成を示す図である。

10

【図16】本発明の第1の実施例の問い合わせ先のクラスタのパラメータ域の例を示す図である。

【図17】本発明の第1の実施例の運用状態の問い合わせ動作を説明するためのシーケンスチャートである。

【図18】本発明の第1の実施例のクラスタ間の通信処理を説明するための図である。

【図19】本発明の第1の実施例の計算機間の第1の通信動作のシーケンスチャートである。

【図20】本発明の第1の実施例の計算機間の第2の通信動作のシーケンスチャートである。

【図21】本発明の第1の実施例の通信割り込み処理の第1の例を説明するための図である。

20

【図22】本発明の第1の実施例の通信割り込み処理の第1の例の動作シーケンスチャートである。

【図23】本発明の第1の実施例の通信割り込み処理の第2の例を説明するための図である。

【図24】本発明の第1の実施例の通信割り込み処理の第2の例の動作のシーケンスチャートである。

【図25】本発明の第1の実施例の通信割り込み処理の第3の例を説明するための図である。

【図26】本発明の第1の実施例の通信割り込み処理の第3の通信動作シーケンスチャート(その1)である。

30

【図27】本発明の第1の実施例の通信割り込み処理の第3の通信動作のシーケンスチャート(その2)である。

【図28】本発明の第1の実施例の各割り込み発生事象でみた場合の例を示す図である。

【図29】本発明の第1の実施例のリセット処理を説明するための図である。

【図30】本発明の第1の実施例のリセット処理動作のシーケンスチャートである。

【図31】本発明の第1の実施例のリセット処理の発生事象でみた場合の例を示す図である。

【図32】本発明の第2の実施例のAVMダウン時の処理概要を示す図である。

【図33】本発明の第2の実施例のダウン通知時における各クラスタの処理を示す図である。

40

【図34】本発明の第2の実施例のSSUのクラスタ制御獲得フィールドの構成図である。

【図35】本発明の第2の実施例のダウンの発生の通知・認識動作を示すシーケンスチャートである。

【図36】本発明の第2の実施例のダウンクラスタ制御処理のフローチャートである。

【図37】本発明の第2の実施例の他のクラスタからダウンクラスタを認識する処理の概要を説明するための図である。

【図38】本発明の第2の実施例の他のクラスタのゲストクラスタのダウン認識動作のシーケンスチャートである。

50

【図 39】本発明の第 2 の実施例のダウン通知を受信した際のダウンクラスタの制御動作を示すシーケンスチャートである。

【図 40】本発明の第 2 の実施例の AVM 運用クラスタにおいて、AVM がダウンした場合の処理を説明するための図である。

【図 41】本発明の第 2 の実施例の AVM 運用クラスタの I/O リセット処理を説明するための図である。

【図 42】本発明の第 2 の実施例の自実クラスタ制御時の待ち制御処理を説明するための図である。

【図 43】本発明の第 2 の実施例のゲストクラスタのセッション閉塞時の処理概要を説明するための図である。

10

【図 44】本発明の第 2 の実施例のゲストクラスタのセッション閉塞時の処理における各クラスタの構成図である。

【図 45】本発明の第 2 の実施例のオペレータ介入メッセージが出力されている状態を示す図である。

【図 46】本発明の第 2 の実施例のオペレータ介入抑止時の状態を示す図である。

【図 47】従来の第 1 の計算機システムの構成例である。

【図 48】従来の第 2 の計算機システムの構成例である。

【図 49】従来のシステムを説明するための図である。

【図 50】従来の第 3 の計算機システムの構成例である。

【図 51】従来の第 3 の計算機システムにおける通信システムを説明するための図である

20

【図 52】従来の通信時における割り込み処理を説明するためのシーケンスチャートである。

【図 53】従来のシステムの制御のリセット処理を説明するためのシーケンスチャートである。

【図 54】従来の問題点を説明するための図（その 1）である。

【図 55】従来の問題点を説明するための図（その 2）である。

【図 56】従来の問題点を説明するための図（その 3）である。

【符号の説明】

1 ロック獲得手段

30

2 第 1 の停止監視手段

3 仮想・実計算機確認手段

4 バス切断手段

60、61、62 アクセスバス

71、72 論理バス

100 共有メモリ

110 計算機番号領域

111 制御クラスタアドレス格納域

112 IPL 世代

113 制御状態情報

40

120 クラスタ制御獲得フィールド

130 ダウンクラスタ管理フィールド

200、300、400、500 クラスタ

210、310、410 AVM、仮想計算機用制御手段

211、311、411 通信送受信部

220、320、420 ゲストクラスタ

221、431、531 計算機間通信制御部

222 パラメータ域

223、323、423 タスク

250、450 OS

50

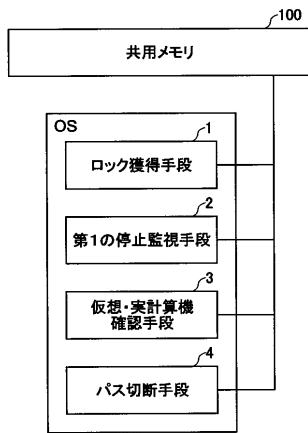
- 2 5 1 , 3 5 1 , 4 5 1 ダウン通知受信部
- 2 5 2 , 3 5 2 , 4 5 2 ダウン認証部
- 2 5 3 , 3 5 3 , 4 5 3 ダウンクラスタ制御部
- 2 5 4 , 3 5 4 , 4 5 4 資源回収処理部
- 2 5 5 , 3 5 5 , 4 5 5 ホットスタンバイ処理部
- 2 4 0 , 3 4 0 , 4 4 0 , 5 4 0 実計算機制御部（ハードウェア）、実計算機制御手段
- 2 9 0 , 3 9 0 サービスプロセッサ
- 2 9 1、3 9 1 , 4 9 1 表示装置
- 3 1 5 閉塞認識部（DEACT認識部）
- 3 1 6 閉塞処理部（DEACT処理部）
- 3 4 1 パラメータ域
- 4 0 1 , 5 0 1 OS
- 4 2 0 OS
- 4 3 0 ダウン通知部
- 4 4 1 メモリ更新部
- 4 4 2 初期化監視部
- 4 4 3 停止監視部
- 4 4 4 仮想・実計算機判定部
- 4 4 5 パス切断部
- 4 4 6 初期化部
- 4 4 7 制御部
- 2 2 1 1 パラメータ解析部
- 2 2 1 2 仮想計算機間通信依頼部
- 3 1 1 0 ダウン認識処理部
- 3 1 2 0 ダウン通知処理部

10

20

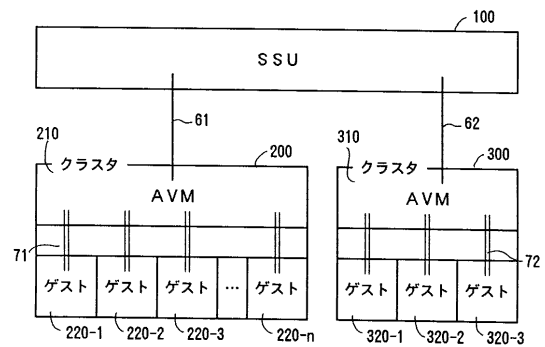
【 図 1 】

本発明の原理構成図



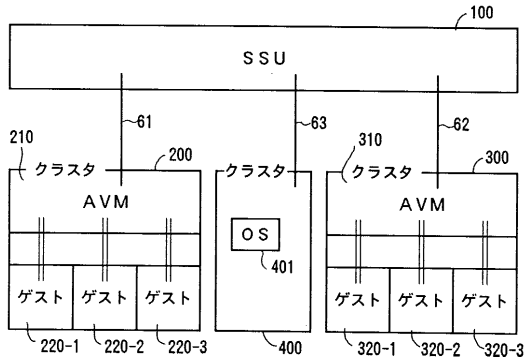
【 図 2 】

本発明の計算機システム（SCMPシステム）の構成図



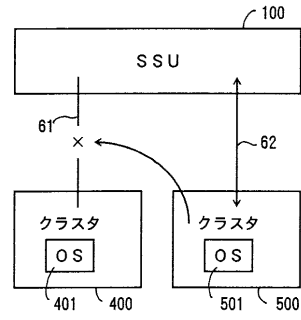
【図3】

本発明の第1の実施例の計算機システム（SCMPシステム）の構成図



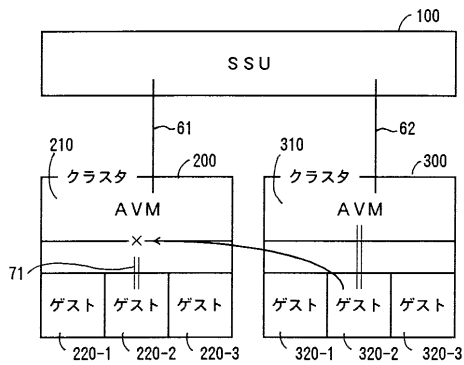
【図4】

本発明の第1の実施例の初期化の概要を説明するための図（その1）



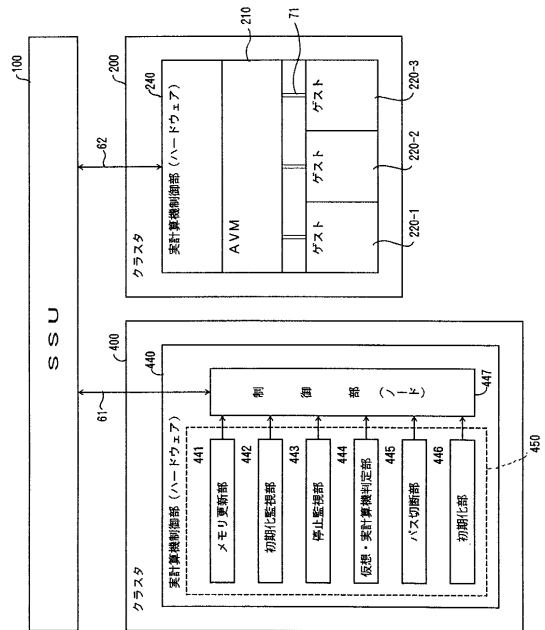
【図5】

本発明の第1の実施例の初期化の概要を説明するための図（その2）



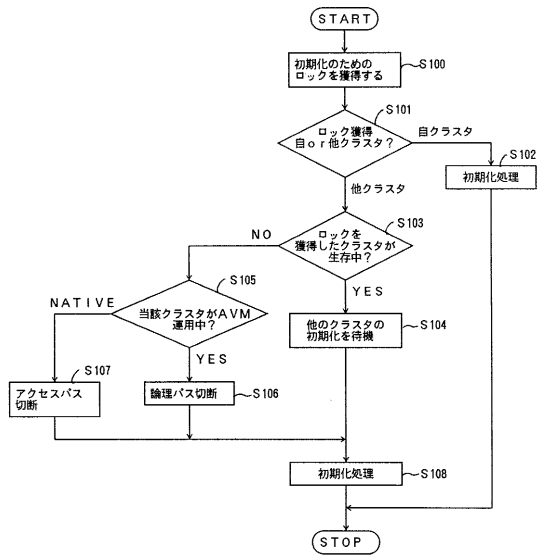
【図6】

本発明の第1の実施例の初期化処理におけるシステム構成図



【図7】

本発明の第1の実施例の初期化処理のフローチャート



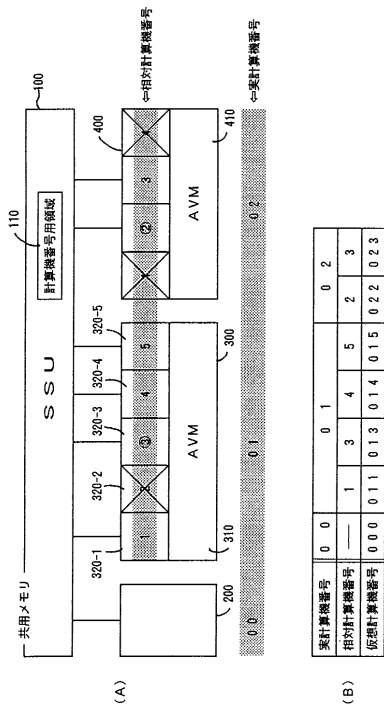
【図8】

本発明の各要求コードの例を示す図

種別	AVM-OS間通信の要求の種類	送信者およびインタフェース (カッコ内は内部コード)
1	AVM/EXクラスタが動作中かどうかを確認する	送信方向 : OS→AVM 詳細コード: AVMCHK (4011) パラメータ: なし
2	AVM/EXクラスタが動作中かどうかを確認する	送信方向 : AVM→OS 詳細コード: AVMACT (4111) パラメータ: AVM動作実クラスタ識別子
3	AVM/EX配下の仮想クラスタのSSUバスを切断する	送信方向 : OS→AVM 詳細コード: AVMCLEUT (4012) パラメータ: 切断する仮想クラスタ識別子
4	AVM/EX配下仮想クラスタが消失したことを通知する	送信方向 : AVM→OS 詳細コード: AVMGDEA (4102) パラメータ: 消失した仮想クラスタ識別子
5	AVM/EX配下の仮想クラスタのSSUバス接続・切断状態を確認する	送信方向 : OS→AVM 詳細コード: AVMPCHK (4003) パラメータ: なし
6	AVM/EX配下の仮想クラスタのSSUバス接続・切断状態を確認する	送信方向 : AVM→OS 詳細コード: AVMPINF (4103) パラメータ: 仮想クラスタバス情報
7	AVM/EXクラスタがダウンしたことを通知する	送信方向 : AVM→OS 詳細コード: AVMDOWN (4105) パラメータ: ダウンコード
8	AVM/EX配下の仮想クラスタのリセット・ダンプロードが完了したことを通知する	送信方向 : AVM→OS 詳細コード: AVMGRST (4106) パラメータ: 要求時のトークン他

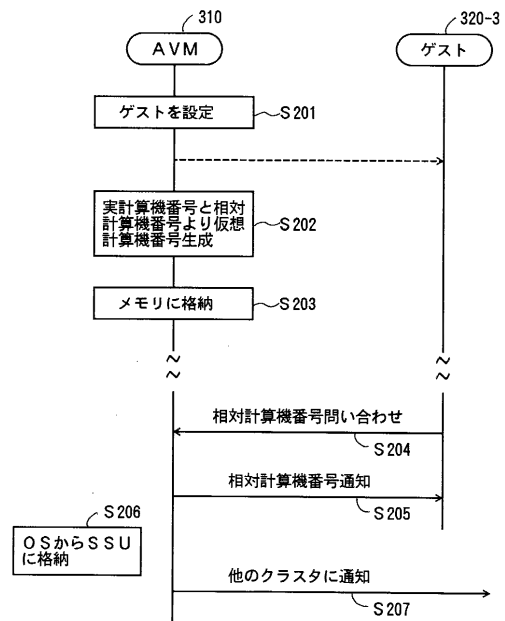
【図9】

本発明の第1の実施例のAVMに対する識別子付与処理を説明するための図



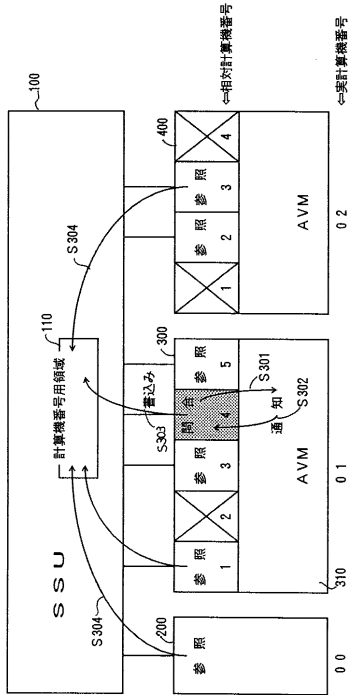
【図10】

本発明の第1の実施例のゲストに対して仮想計算機番号を付与する処理を説明するためのシーケンスチャート



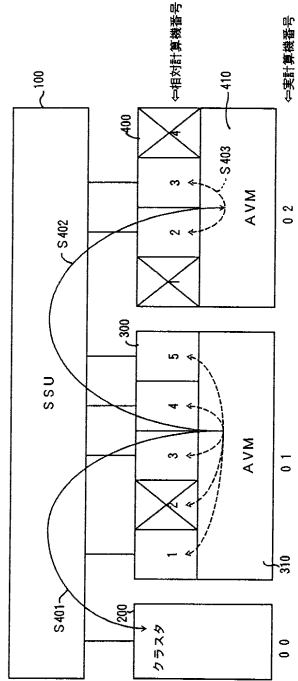
【 図 1 1 】

本発明の第1の実施例の仮想計算機番号の参照動作を説明するための図



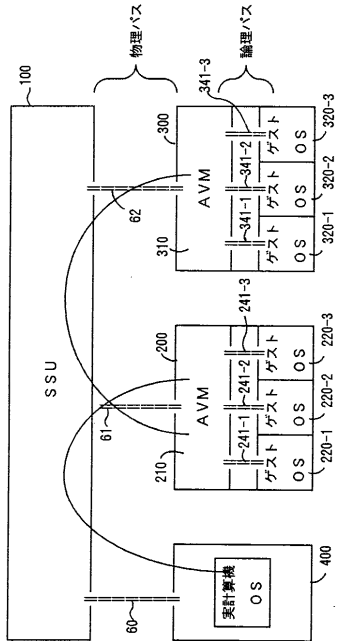
【 図 1 2 】

本発明の第1の実施例のあるクラスタから他のクラスタへゲストの仮想計算機番号を通知する他の例を示す図



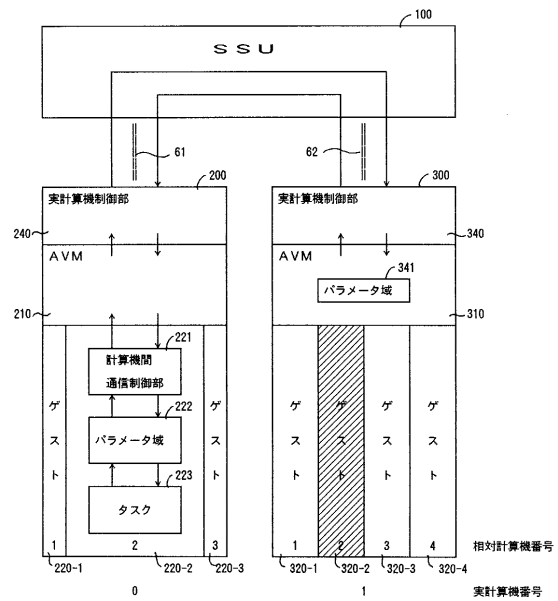
【 図 1 3 】

本発明の第1の実施例の運用情報取得の概念図



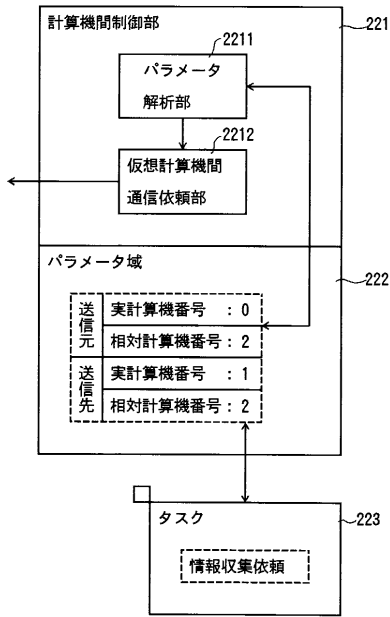
【 図 1 4 】

本発明の第1の実施例の運用状態情報取得時のシステム構成図



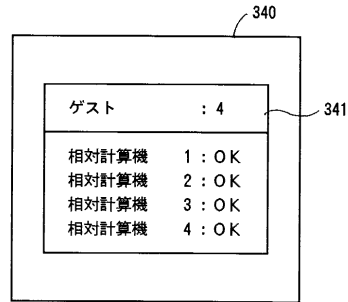
【 図 1 5 】

本発明の第1の実施例のゲスト内の構成を示す図



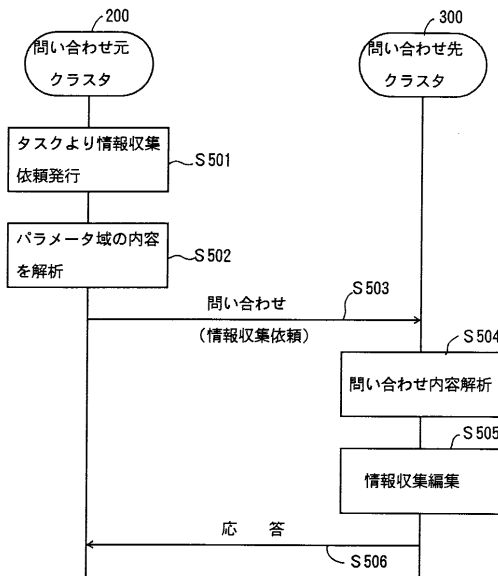
【 図 1 6 】

本発明の第1の実施例の問い合わせ先のクラスタのパラメータ域の例を示す図



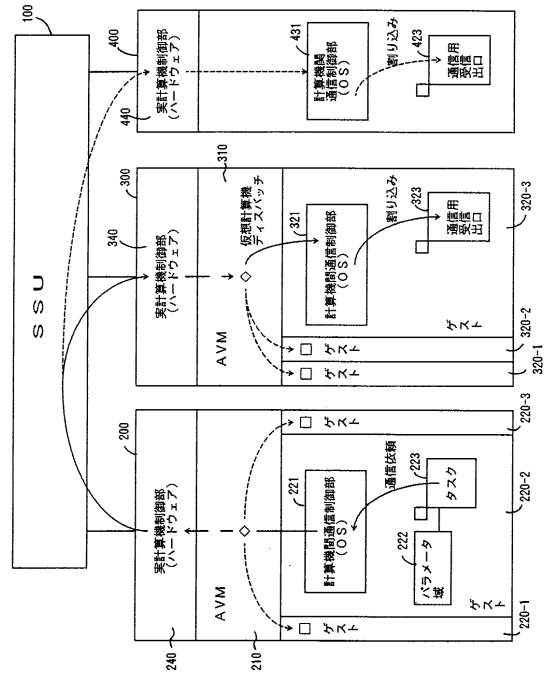
【 図 1 7 】

本発明の第1の実施例の運用状態情報の問い合わせ動作を説明するためのシーケンスチャート



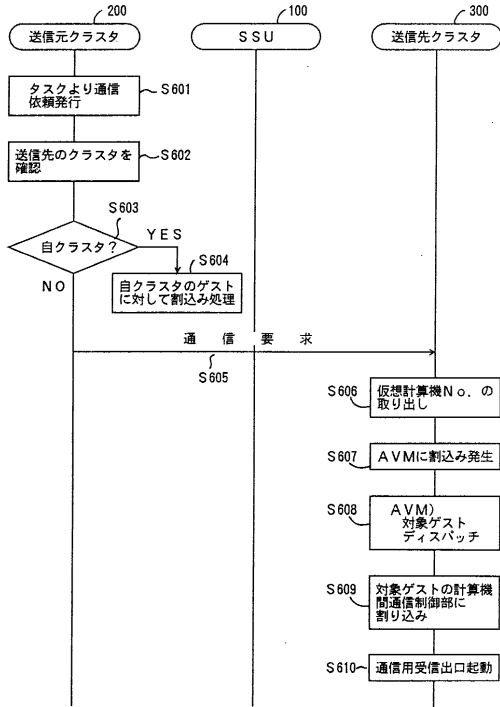
【 図 1 8 】

本発明の第1の実施例のクラスタ間の通信処理を説明するための図



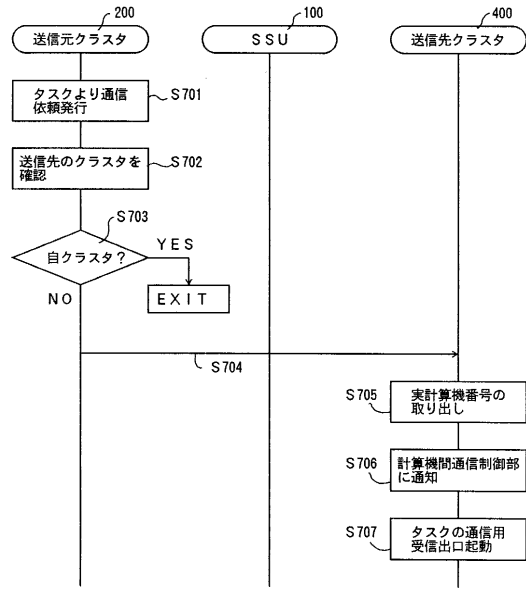
【図19】

本発明の第1の実施例の計算機間の第1の通信動作のシーケンスチャート



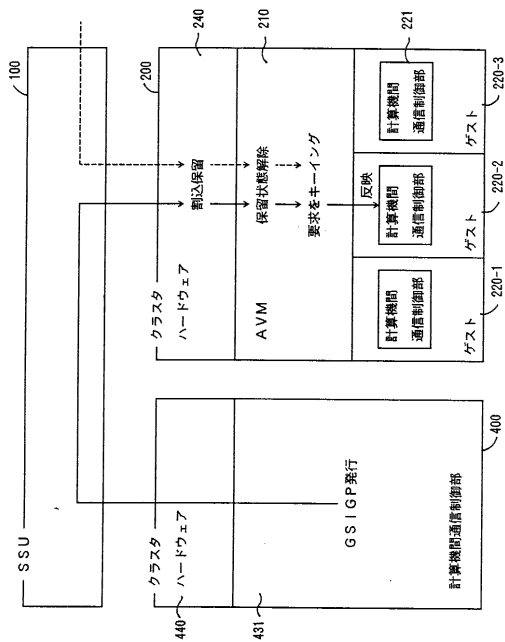
【図20】

本発明の第1の実施例の計算機間の第2の通信動作シーケンスチャート



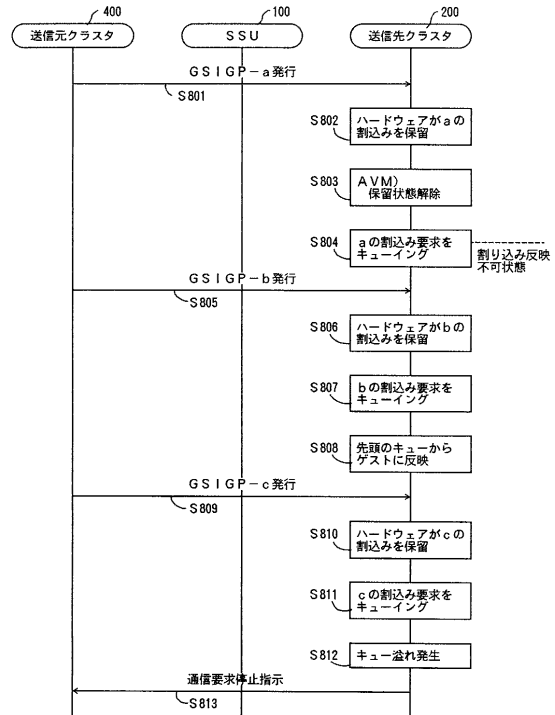
【図21】

本発明の第1の実施例の通信割り込み処理の第1の例を説明するための図



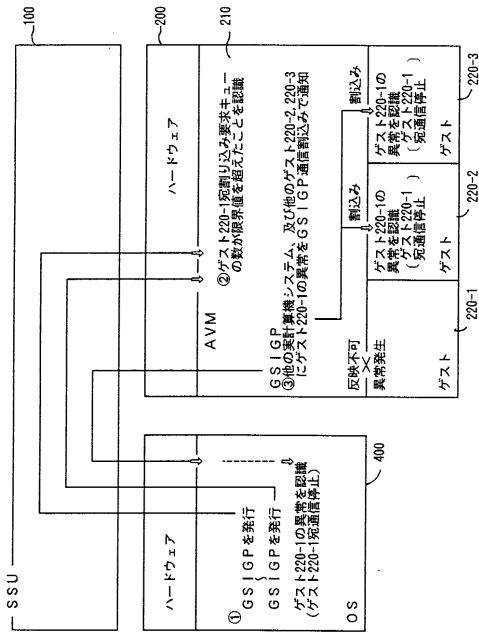
【図22】

本発明の第1の実施例の通信割り込み処理の第1の例の動作シーケンスチャート



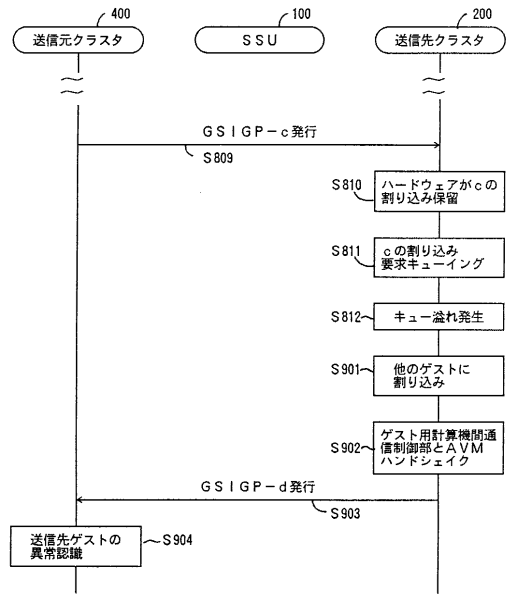
【図 23】

本発明の第1の実施例の割り込み処理の第2の例を説明するための図



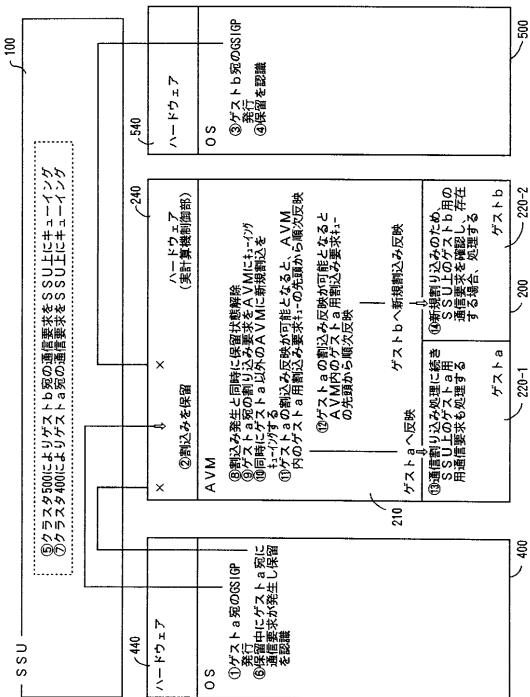
【図 24】

本発明の第1の実施例の通信割り込み処理の第2の例の動作シーケンスチャート



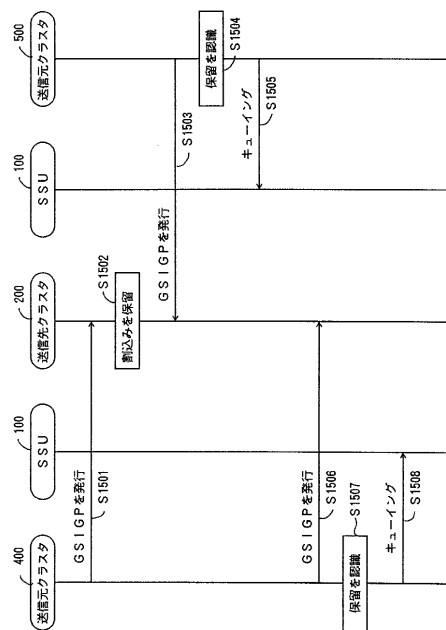
【図 25】

本発明の第1の実施例の通信割り込み処理の第3の例を説明するための図



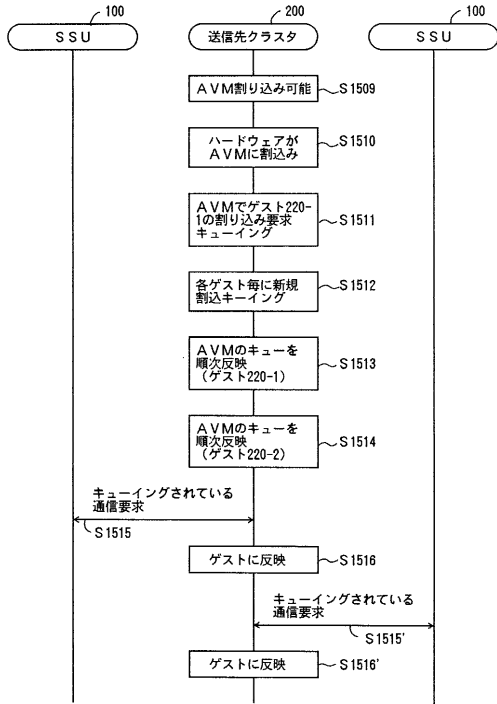
【図 26】

本発明の第1の実施例の通信割り込み処理の第3の通信動作シーケンスチャート (その1)



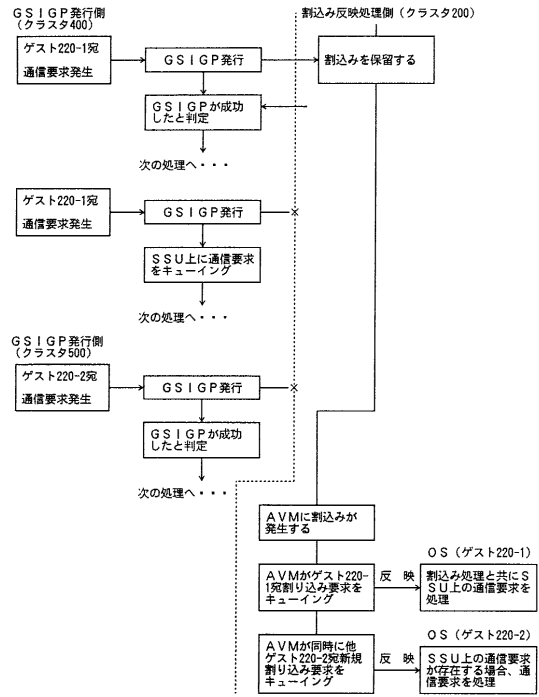
【 図 27 】

本発明の第1の実施例の通信割り込み処理の第3の通信動作シーケンスチャート(その2)



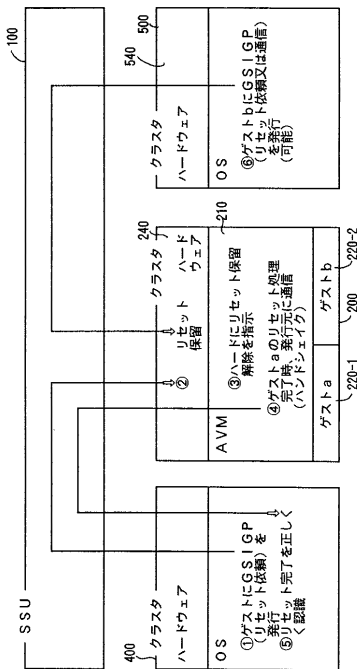
【 図 28 】

本発明の第1の実施例の各割り込み発生事象でみた場合の例を示す図



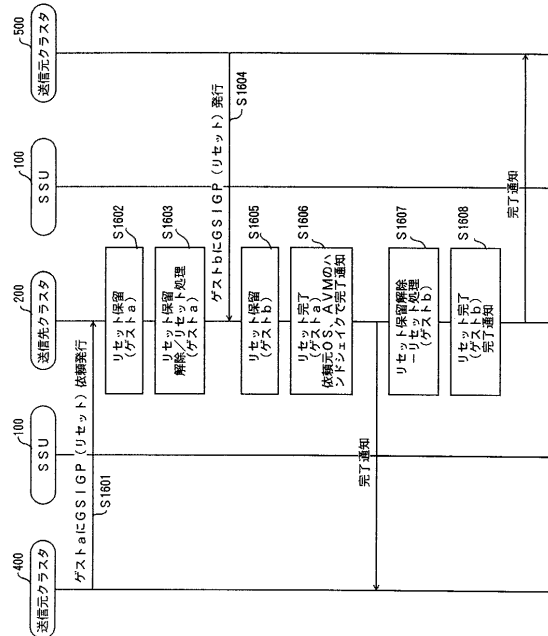
【 図 29 】

本発明の第1の実施例のリセット処理を説明するための図



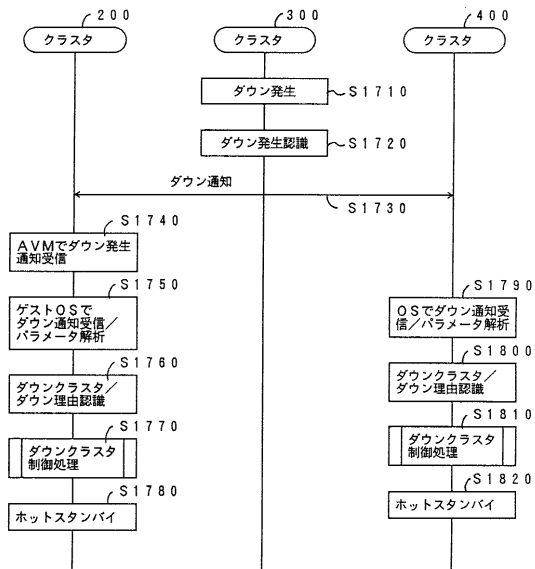
【 図 30 】

本発明の第1の実施例のリセット処理動作のシーケンスチャート



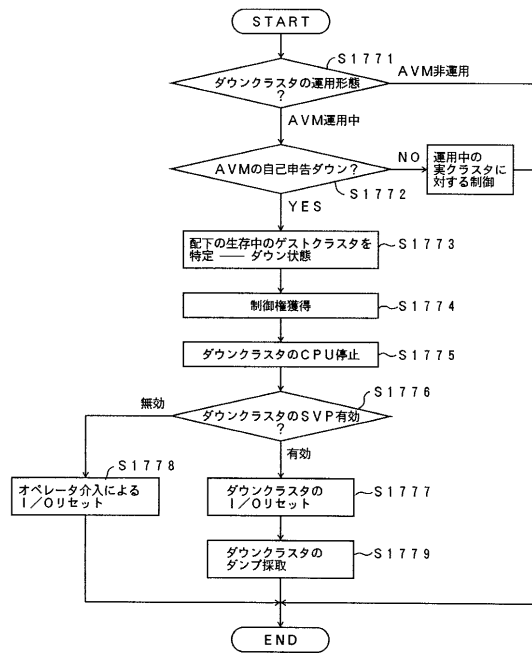
【 図 3 5 】

本発明の第2の実施例のダウンの発生・認識動作を示すシーケンスチャート



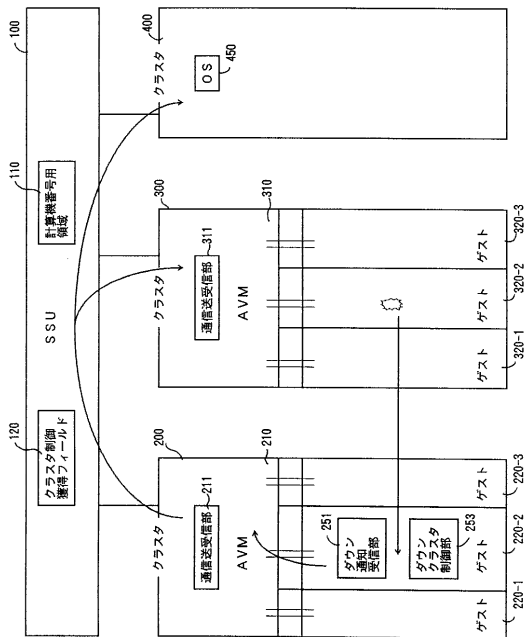
【 図 3 6 】

本発明の第2の実施例のダウンクラスタの制御処理のフローチャート



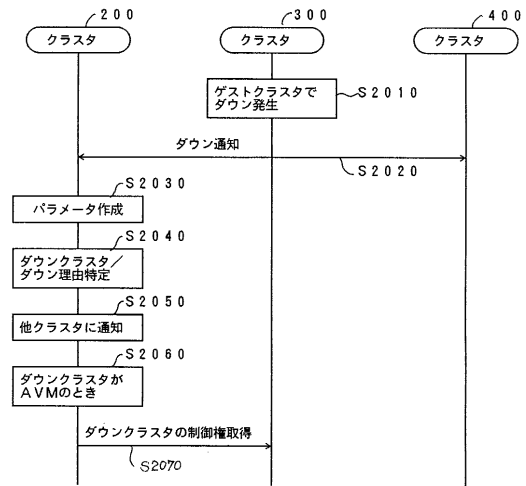
【 図 3 7 】

本発明の第2の実施例の他のクラスタからダウンクラスタを認識する処理概要を説明するための図



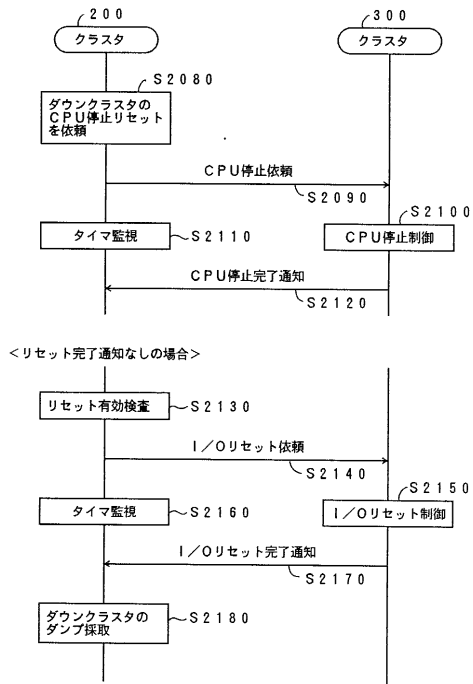
【 図 3 8 】

本発明の第2の実施例のダウンの発生・認識動作を示すシーケンスチャート



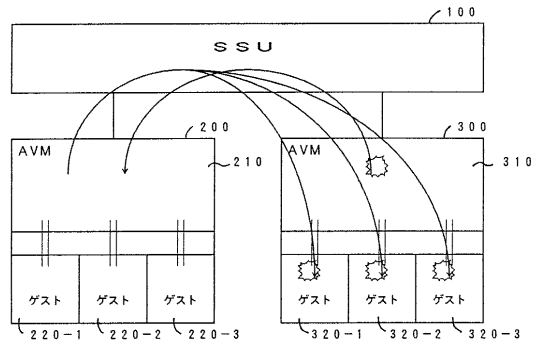
【図39】

本発明の第2の実施例のダウン通知を受信した際の
ダウンクラスタの制御動作を示すシーケンスチャート



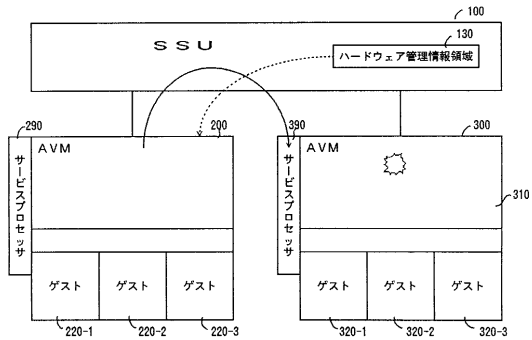
【図40】

本発明の第2の実施例のAVM通用クラスタにおいてAVMが
ダウンした場合の処理を説明するための図



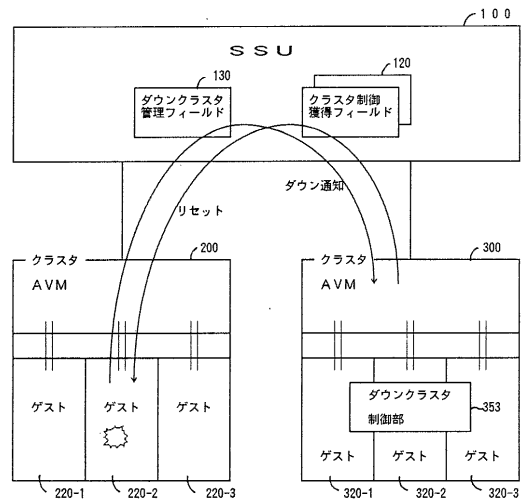
【図41】

本発明の第2の実施例のAVM通用クラスタのI/Oリセット処理を説明するための図



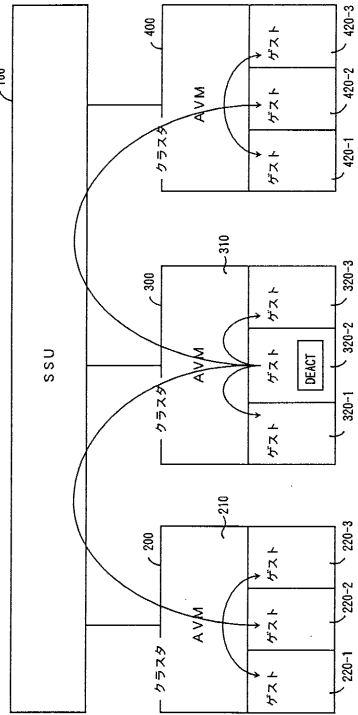
【図42】

本発明の第2の実施例の自実クラスタ制御時の待ち制御処理を説明するための図



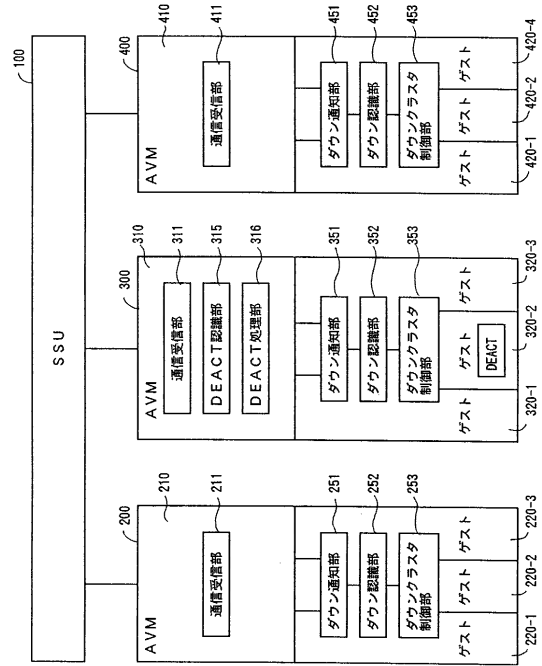
【図 4 3】

本発明の第2の実施例のゲストクラスタのセッション閉塞時の処理概要を説明するための図



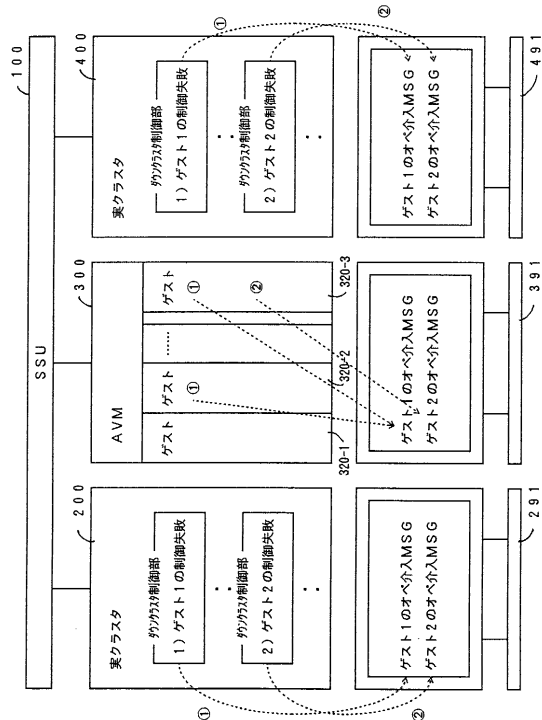
【図 4 4】

本発明の第2の実施例のゲストクラスタのセッション閉塞時の処理における各クラスタの構成図



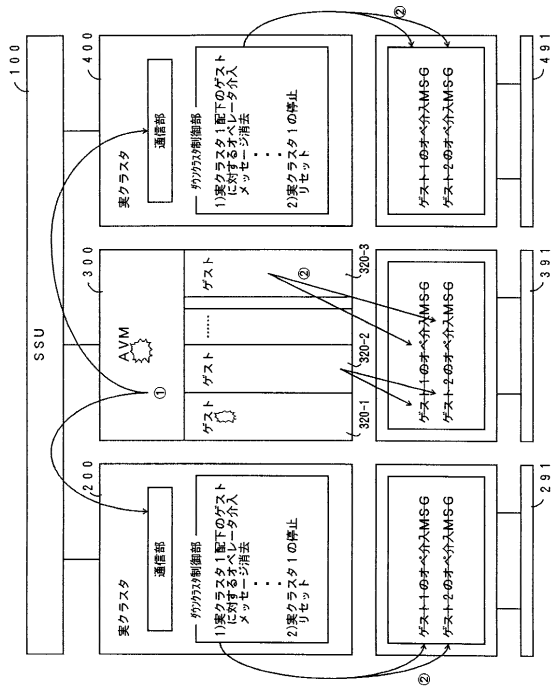
【図 4 5】

本発明の第2の実施例のオペレータ介入メッセージが出力されている状態を示す図



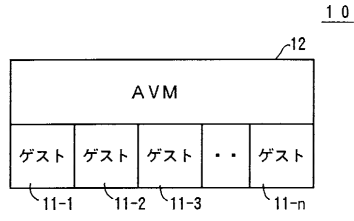
【図 4 6】

本発明の第2のオペレータ介入抑止時の状態を示す図



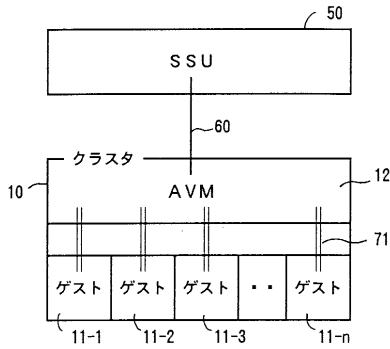
【図47】

従来の第1の計算機システムの構成例



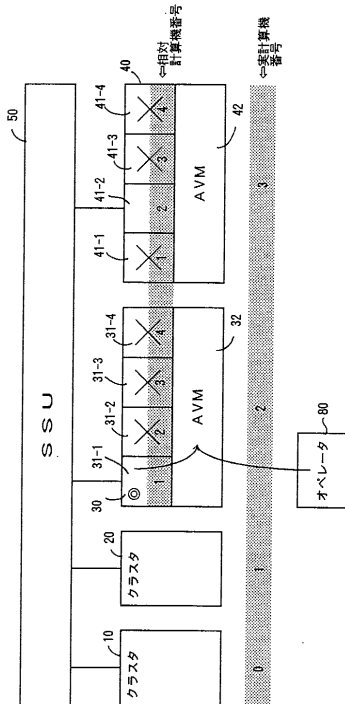
【図48】

従来の第2の計算機システムの構成例



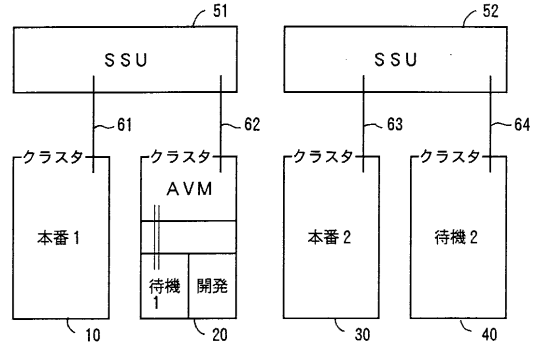
【図50】

従来の第3の計算機システムの構成例



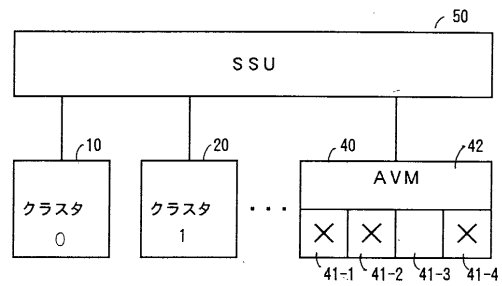
【図49】

従来の第2の計算機システムを説明するための図



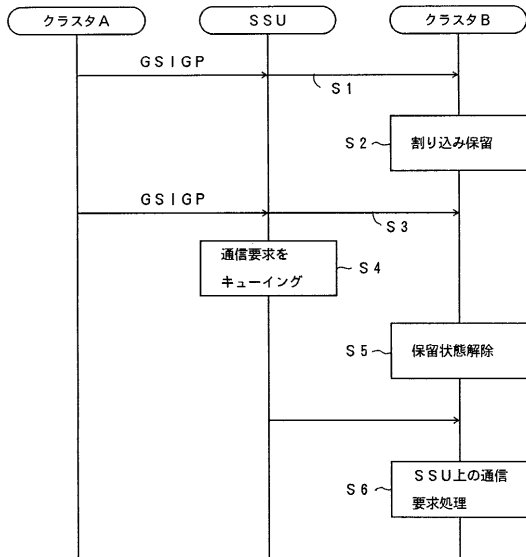
【図51】

従来の第3の計算機システムにおける通信システムを説明するための図



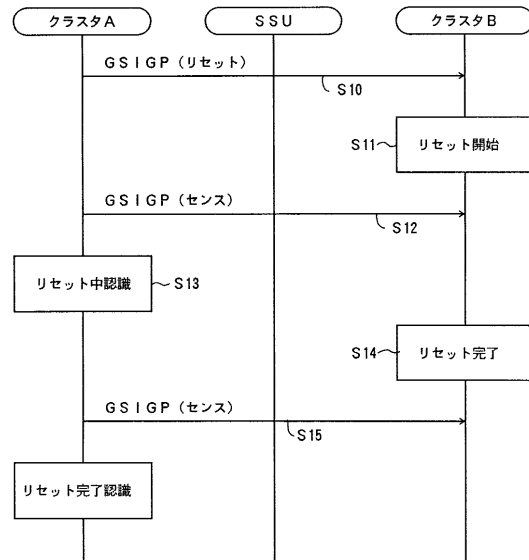
【 図 5 2 】

従来の通信時における割り込み処理を説明するためのシーケンスチャート



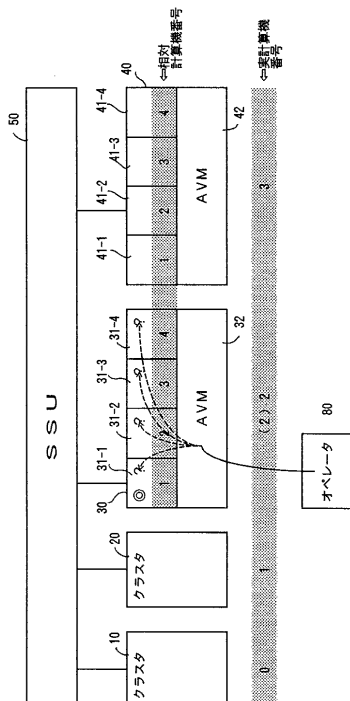
【 図 5 3 】

従来のシステム制御のリセット処理を説明するためのシーケンスチャート



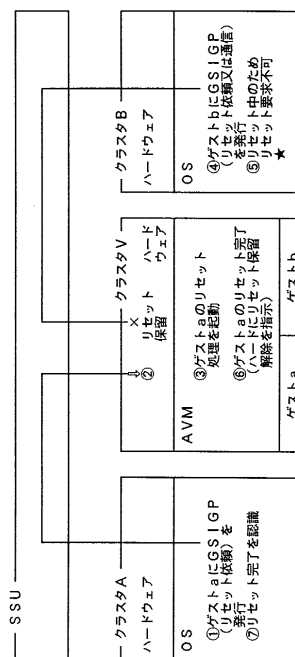
【 図 5 4 】

従来の問題点を説明するための図(その1)



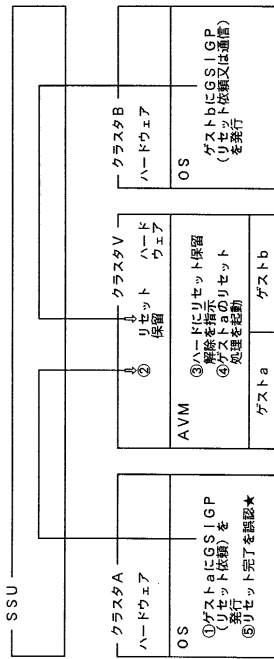
【 図 5 5 】

従来の問題点を説明するための図(その2)



【 図 5 6 】

従来の問題点を説明するための図（その3）



フロントページの続き

- (72)発明者 斎藤 優
神奈川県川崎市中原区上小田中1015番地 富士通株式会社内
- (72)発明者 下川 健一郎
神奈川県川崎市中原区上小田中1015番地 富士通株式会社内
- (72)発明者 平岡 勝則
静岡県静岡市伝馬町16番地の3 株式会社富士通静岡エンジニアリング内
- (72)発明者 堀崎 公史
静岡県静岡市伝馬町16番地の3 株式会社富士通静岡エンジニアリング内
- (72)発明者 塚本 建一
静岡県静岡市伝馬町16番地の3 株式会社富士通静岡エンジニアリング内
- (72)発明者 落合 由美
静岡県静岡市伝馬町16番地の3 株式会社富士通静岡エンジニアリング内

審査官 久保 正典

- (56)参考文献 特開平03-077143(JP,A)
特開平02-017563(JP,A)
特開平05-216855(JP,A)
特開平04-023160(JP,A)
特開平05-073519(JP,A)
特開平04-141744(JP,A)
特開平05-324362(JP,A)

- (58)調査した分野(Int.Cl.⁷, DB名)
G06F9/46-9/54
G06F15/16-15/177