

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6135403号
(P6135403)

(45) 発行日 平成29年5月31日 (2017.5.31)

(24) 登録日 平成29年5月12日 (2017.5.12)

(51) Int.Cl.		F I			
G06F 11/07	(2006.01)	G06F 11/07	1 7 2		
G06F 11/34	(2006.01)	G06F 11/07	1 4 0 A		
		G06F 11/34	1 7 6		

請求項の数 5 (全 27 頁)

(21) 出願番号	特願2013-175250 (P2013-175250)	(73) 特許権者	000005223
(22) 出願日	平成25年8月27日 (2013.8.27)		富士通株式会社
(65) 公開番号	特開2015-45905 (P2015-45905A)		神奈川県川崎市中原区上小田中4丁目1番1号
(43) 公開日	平成27年3月12日 (2015.3.12)	(74) 代理人	100094525
審査請求日	平成28年5月10日 (2016.5.10)		弁理士 土井 健二
		(74) 代理人	100094514
			弁理士 林 恒徳
		(72) 発明者	結城 和博
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		審査官	石川 亮

最終頁に続く

(54) 【発明の名称】 情報処理システム、情報処理システムの障害処理方法

(57) 【特許請求の範囲】

【請求項 1】

複数のノード間でメモリを共有する情報処理システムにおいて、
 前記ノードの各々は、
 複数の機能回路と前記機能回路を制御する制御装置と、
 前記複数の機能回路から発生する割り込み要因を格納するレジスタとを有し、
 前記複数のノードのうちの1のノードの前記制御装置は、
 他の前記ノードの割り込み要因の発生に応じて前記レジスタの前記割り込み要因を受信し、前記割り込み要因のうち、障害として検出すべき割り込み要因を抽出して、抽出結果に応じて障害ノードを特定し、前記障害ノードの前記メモリへのアクセスを抑止後、前記他のノードから受信したログ情報に基づいて前記障害ノードの切り離し制御を行う情報処理システム。

10

【請求項 2】

請求項 1 において、
 前記 1 のノードは、網結合装置を備え、
 前記他のノードは、データ処理を実行し、前記網結合装置を介して前記メモリにアクセスする処理装置を備える情報処理システム。

【請求項 3】

請求項 1 または 2 において、
 前記 1 のノードの前記制御装置は、前記障害として検出すべき割り込み要因の波及元と

20

なる割り込み要因が発生しているか否かを判定し、発生していない場合に、前記割り込み要因に対応するノードを前記障害ノードとして特定し、発生している場合に、前記波及元となる割り込み要因に対応するノードを前記障害ノードとして特定する情報処理システム。

【請求項 4】

請求項 3 において、

前記 1 のノードは、

前記割り込み要因と、前記割り込み要因の波及元となる割り込み要因との対応関係を有する定義テーブルを有し、

前記 1 のノードの前記制御装置は、前記定義テーブルに基づいて、前記障害として検出すべき割り込み要因の波及元となる割り込み要因が発生しているか否かを判定する情報処理システム。

10

【請求項 5】

複数のノード間でメモリを共有する情報処理システムの障害処理方法において、

前記ノードの各々は、

複数の機能回路と前記機能回路を制御する制御装置と、

前記複数の機能回路から発生する割り込み要因を格納するレジスタとを有し、

前記複数のノードのうちの 1 のノードの前記制御装置は、

他の前記ノードの割り込み要因の発生に応じて前記レジスタの前記割り込み要因を受信し、前記割り込み要因のうち、障害として検出すべき割り込み要因を抽出して、抽出結果に応じて障害ノードを特定し、前記障害ノードの前記メモリへのアクセスを抑止後、前記他のノードから受信したログ情報に基づいて前記障害ノードの切り離し制御を行う情報処理システムの障害処理方法。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理システム、情報処理システムの障害処理方法に関する。

【背景技術】

【0002】

複数のノードを有する情報処理システムは、例えば、ビルディングブロック (BB: Building Block) 構造を有する。例えば、複数のノードでメモリを共有する情報処理システムは、ノード間で、クロスバーを介してメモリを共有する。情報処理システムで動作するアプリケーションは、共有されたメモリを使用することによって、システムの処理性能の向上を図る。一方、それぞれのノードで動作する OS (Operation system、以下、OS と称する) やハイパーバイザ (hypervisor) は、各ノードのローカルメモリ上で動作する。OS やハイパーバイザがローカルメモリ上で動作することにより、各ノードの独立性が高まり、システムの可用性が向上する。

30

【0003】

このような情報処理システムにおいて、一部のノードのハードウェアに障害が発生した場合、障害が発生した障害ノードを検出すると共に、障害ノードをシステムから切り離れた状態で、運用を再開することが求められる。ハードウェアの障害の検出は、例えば、特許文献 1 に記載される。

40

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特開 2011-248653 号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

情報処理システムは、障害ノードの特定やシステムからの切り離しの可否を、障害の事

50

象を有するログ情報に基づいて、順次、解析する。したがって、情報処理システムのノード数や、障害の事象の種類が増加に伴って、障害ノードの特定やシステムからの切り離しの要否に係る解析時間も増加する。また、ログ情報のデータ量が膨大であることにより、ログ情報の収集にも時間を要する。

【 0 0 0 6 】

1つの側面は、本発明は、障害発生時に障害ノードによる他ノードの影響を早急に低減する情報処理システム、情報処理システムの障害処理方法を提供することを目的とする。

【課題を解決するための手段】

【 0 0 0 7 】

第1の側面は、複数のノード間でメモリを共有する情報処理システムにおいて、前記ノードの各々は、複数の機能回路と前記機能回路を制御する制御装置と、前記複数の機能回路から発生する割り込み要因を格納するレジスタとを有し、前記複数のノードのうちの1のノードの前記制御装置は、他の前記ノードの割り込み要因の発生に応じて前記レジスタの前記割り込み要因を受信し、前記割り込み要因のうち、障害として検出すべき割り込み要因を抽出して、抽出結果に応じて障害ノードを特定し、前記障害ノードの前記メモリへのアクセスを抑止後、前記他のノードから受信したログ情報に基づいて前記障害ノードの切り離し制御を行う。

【発明の効果】

【 0 0 0 8 】

第1の側面によれば、情報処理システムは、割り込み要因のうち、障害として検出すべき割り込み要因を抽出して、抽出結果に応じて障害ノードを特定し、障害ノードのメモリへのアクセスを抑止することで、障害発生時に障害ノードによる他ノードの影響を早急に低減する。

【図面の簡単な説明】

【 0 0 0 9 】

【図1】本実施の形態例における情報処理システムの概要を説明する図である。

【図2】図1の情報処理システムの構成の一例を示す図である。

【図3】図2に示したシステムボードの構成の一例を説明する図である。

【図4】図3のレジスタを説明する図である。

【図5】図1～図3で述べた、本実施の形態例における情報処理システムの一部のノードにおいて障害が発生した場合の処理の流れを説明する図である。

【図6】図5において述べた、マスターノードのシステム制御装置におけるログ情報の解析処理の概要を説明する図である。

【図7】図6のログ情報の解析、及び、FNL解析処理に要する時間を例示する図である。

【図8】本実施の形態における情報処理システムの各ノードのソフトウェアモジュール図である。

【図9】図8において説明したFNL(Fail Node List)の一例を示す図である。

【図10】割り込み要因が発生しFNLが更新される間のマスターノードのシステム制御装置、及び、スレーブノードのシステム制御装置における、処理の流れを時系列に説明する図である。

【図11】本実施の形態例におけるFNL解析部の処理、及び、FNL更新部の処理を説明するフローチャート図である。

【図12】波及先の割り込み要因の抑止処理を説明するフローチャート図である。

【図13】FNLDB(Fail Node DB)の一例を示す図である。

【図14】アクション番号(act)を有する定義テーブルの具体例を示す図である。

【図15】具体例におけるメモリのアクセスの抑止範囲を説明する図である。

【発明を実施するための形態】

【 0 0 1 0 】

以下、図面にしたがって本発明の実施の形態を説明する。ただし、本発明の技術的範囲

10

20

30

40

50

はこれらの実施の形態に限定されず、特許請求の範囲に記載された事項とその均等物まで及ぶものである。

【 0 0 1 1 】

〔 情報処理システムの概要 〕

図 1 は、本実施の形態例における情報処理システム 1 の概要を説明する図である。図 1 に示す情報処理システム 1 は、H P C (High Performance Computing) モデル等の計算機システムである。このようなシステムは、ビルディングブロック (BB: Building Block) 構造によって構成される。各ビルディングブロック 1 0 a ~ 1 0 e は、図 1 に示すシステムボード 1 A ~ 1 E を収容し、ラックに抜き差し可能である。また、図 1 の情報処理システム 1 は、複数のシステムボード 1 A ~ 1 E と、網結合装置 (以下、クロスバスイッチと称する) 2 を備えるシステムボードとを有する。各システムボード 1 A ~ 1 E は、クロスバスイッチ 2 を介して、相互に接続する。なお、図 1 には、5 つのシステムボード 1 A ~ 1 E が示されるが、情報処理システム 1 は、例えば、1 6 台のシステムボードを有する。

【 0 0 1 2 】

また、システムボード 1 A は、複数の C P U (Central Processing Unit) 1 2 a とメモリ 3、1 1 a と、I / O (Input Output) 装置 1 3 a とを有する。また、メモリ 3、1 1 a の一部の領域は、情報処理システム 1 が有する全ての C P U が共用する共有メモリ 3 として使用され、他の領域は、C P U 1 2 a がカーネルデータ等を格納するローカル領域 1 1 a として使用される。他のシステムボード 1 B ~ 1 E も、システムボード 1 A と同様の構成を有する。以下、各システムボードをノードと称する。

【 0 0 1 3 】

また、ノード 1 A のファームウェア層 1 4 a では、例えば、ハイパーバイザ (hypervisor) と呼ばれる制御ソフトウェアが動作する。ハイパーバイザは、ノード 1 A のリソースを論理的に分割して、1 つまたは複数の論理パーティション D a、D b を生成する。複数の論理パーティション D a、D b が生成されることにより、1 つのノード上で複数の O S (Operation system、以下、O S と称する) が動作可能になる。なお、図 1 の例において、各論理パーティション D a、D b 上で動作する O S (例えば、S o l a r i s (登録商標)) は、異なる種類の O S であってもよい。

【 0 0 1 4 】

また、各論理パーティション D a ~ D h 上で動作するアプリケーション p a ~ p h は、例えば、共有メモリ 3 を使用する。即ち、本実施の形態例では、各ノードが共有メモリ 3 の一部を有し、各ノードが他ノードの共有メモリ 3 を利用する分散型共有メモリを構成する。そして、アプリケーション p a ~ p h は、共有メモリ 3 に記憶された共有の情報に基づいて、所定の処理を行う。また、ハイパーバイザや O S は、各々のローカルメモリ 1 1 a ~ 1 1 e 上で動作することにより独立性が高まり、システムの可用性が向上する。

【 0 0 1 5 】

分散型共有メモリ 3 を有する情報処理システム 1 において、例えば、ノード 1 A の C P U 1 2 a が、アプリケーション p a の実行にあたり、共有メモリ 3 上のノード 1 A とは別のノード (例えば、ノード 1 B) の共有メモリ 3 の領域にアクセスする場合、C P U 1 2 a は、クロスバスイッチ 2 を介して、ノード 1 B の共有メモリ 3 の領域にアクセスのリクエストを送信する。また、C P U 1 2 a が、自ノード 1 A の共有メモリ 3 の領域にアクセスする場合、直接接続を介して、メモリアクセスのリクエストを送信する。

【 0 0 1 6 】

〔 情報処理システムの構成 〕

図 2 は、図 1 の情報処理システム 1 の構成の一例を示す図である。図 2 において、図 1 で示したものと同一のものは、同一の記号で示す。図 2 に示すように、情報処理システム 1 は、例えば、処理装置としての 1 6 台のシステムボード (S B: System Board) 1 A ~ 1 P と、4 台のクロスバスイッチボックス 2 A B ~ 2 D B とを有する。クロスバスイッチボックス 2 A B ~ 2 D B がそれぞれ有するクロスバスイッチ 2 A ~ 2 D は、図 1 に示すクロスバスイッチ 2 に対応する。本実施の形態も、クロスバスイッチボックス 2 A B ~ 2 D

Bは、ビルディングブロック構造である。

【0017】

図2の例において、クロスバスイッチボックス2ABは、クロスバスイッチ2Aと、システム制御装置(SVP: Service Processor)V1とを有する。クロスバスイッチボックス2ABのシステム制御装置V1は、クロスバスイッチ2Aの状態監視、状態設定、及び、起動、停止制御等を行う。また、クロスバスイッチ2Aは、スイッチ2aとポートav、aw~dv、dw、qv、qw、rv、rw、sv、swとスイッチ2aとを有する。スイッチ2aは、通信経路を切り替える。他のクロスバスイッチボックス2BB~2DBの構成も同様である。

【0018】

10

また、図2の例において、それぞれのシステムボード1Aは、2つのクロスバスイッチ2Aとの接続用ポートax、ayを有する。また、クロスバスイッチ2Aも、各システムボード1Aとの2つの接続用ポートav、awを有する。即ち、各システムボード1Aは、2つの回線n1、n2によって、対応するクロスバスイッチ2Aに接続する。このように、図2に示すクロスバスイッチ2A~2Dは、接続対称との間に二重の回線を有する対称型のクロスバスイッチである。二重の回線を有するため、クロスバスイッチ2A~2Dは、片側の回線に障害が発生した場合であっても、残りの一つの回線を使用して動作することができる。

【0019】

この例において、第1、第2、第3、第4のシステムボード1A、1B、1C、1Dは、第1のクロスバスイッチ2Aに接続する。また、第5、第6、第7、第8のシステムボード1E、1F、1G、1Hは、第2のクロスバスイッチ2Bに接続する。また、第9、第10、第11、第12のシステムボード1I、1J、1K、1Lは、第3のクロスバスイッチ2Cに接続する。第13、第14、第15、第16のシステムボード1M、1N、1O、1Pは、第4のクロスバスイッチ2Dに接続する。

20

【0020】

また、図2の例において、第1のクロスバスイッチ2Aは、バスL1、L2によって、第2のクロスバスイッチ2Bと接続する。また、第1のクロスバスイッチ2Aは、バスL7、L8によって、第3のクロスバスイッチ2Cと接続する。また、第1のクロスバスイッチ2Aは、バスL9、L10によって、第4のクロスバスイッチ2Dと接続する。さらに、第2のクロスバスイッチ2Bは、バスL11、L12によって、第3のクロスバスイッチ2Cと接続し、第2のクロスバスイッチ2Bは、バスL3、L4によって、第4のクロスバスイッチ2Dと接続する。そして、第3のクロスバスイッチ2Cは、バスL5、L6によって、第4のクロスバスイッチ2Dと接続する。

30

【0021】

また、各システムボード1A~1Pも、システム制御装置(図3にて図示)を有する。情報処理システム1における各クロスバスイッチボックス2AB~2DBのシステム制御装置V1~V4、及び、各システムボード1A~1Pのシステム制御装置22は、内部バスL40によって互いに接続する。なお、図2において、情報処理システム1は、16台のシステムボード1A~1Pと4台のクロスバスイッチ2A~2Dとを有するが、システムボードの台数及びクロスバスイッチの台数は16台、4台に限定されない。続いて、各システムボード1A~1Pの構成を説明する。

40

【0022】

[システムボードの構成]

図3は、図2に示したシステムボード1A~1Pの構成の一例を説明する図である。図3の例では、システムボード1Aの構成を説明する。他のシステムボード1B~1Pの構成も、システムボード1Aと同様である。図3に示すように、システムボード1Aは、システムボードユニットB1と、サービスプロセッサボードB2とを有する。

【0023】

システムボードユニットB1は、例えば、複数のCPU(CPUチップ)12aと、シ

50

ステムコントローラ (System Controller) 15 と、I/Oコントローラ16と、PCI (Peripheral Component Interconnect) Express 17と、メモリアクセスコントローラ18と、メモリ3、11aと、MBC (Maintenance Bus Controller 以下、MBCと称する) 19とを有する。メモリ3、11aは、例えば、DRAM (Dynamic Random Access Memory) である。MBC 19は、サービスプロセッサボードB2との通信経路を制御する。

【0024】

CPU12aは、図1で説明したアプリケーションpa、pbを実行する演算処理装置である。CPU12aの各々は、システムコントローラ15に接続する。システムコントローラ15は、メモリ3、11aに接続されたメモリアクセスコントローラ18に接続する。また、システムコントローラ15は、I/Oコントローラ16に接続する。I/Oコントローラ16は、例えば、外部メモリ (大容量メモリ及び/又はストレージ装置) やネットワークインタフェースカード (NIC) が接続されたPCI Express 17と接続する。

10

【0025】

そして、システムコントローラ15は、CPU12aとメモリアクセスコントローラ18との間の転送制御を行う。また、システムコントローラ15は、接続ポートax、ayを介して、クロスバスイッチ2Aに接続し、クロスバスイッチ2AとCPU12aとの間の転送制御、及び、クロスバスイッチ2Aとメモリアクセスコントローラ18との間の転送制御を行う。例えば、システムコントローラ15は、ブリッジ回路の役割を果たす。

20

【0026】

また、図1において、前述したとおり、メモリ3、11aの一部の領域はクロスバスイッチ2Aを介して共有され、共有メモリ3 (図1) として使用され、他の一部の領域は、ローカルメモリ11aとして使用される。例えば、システムコントローラ15は、CPU12aが、別のシステムボードに搭載される共有メモリ3の領域にアクセスする場合、接続ポートax、ayを介して、クロスバスイッチ2Aに接続する。一方、CPU12aが、システムボード1Aに搭載されるメモリ3、11aの領域にアクセスする場合、システムコントローラ15は、メモリアクセスコントローラ18にアクセスする。

【0027】

また、サービスプロセッサボードB2は、システム制御装置22とMBC (Maintenance Bus Controller 以下、MBCと称する) 21を有する。システム制御装置22は、ノード内のハードウェアのアクセス制御、監視、電源投入、ログの採取、ユーザインタフェース制御 (ユーザI/F) 等の制御を行う。MBC 21は、システムボードユニットB1との通信経路を制御する。また、MBC 21は、CPU12aやメモリ3、11a、I/Oコントローラ16、システム制御装置22等のハードウェアから発生する割り込み要因を格納するレジスタrgを有する。また、図2で前述したとおり、システム制御装置22は、LAN (Local Area Network) などのネットワーク回線L40を介して、別のノードのシステム制御装置22、V1~V4と相互に接続する。

30

【0028】

なお、図3の例では、システムボード1A (1B~1P) が4台のCPU (CPUチップ) 12aを搭載する例を示したが、システムボード1Aが少なくとも1台のCPU12aを搭載する構成であっても良い。

40

【0029】

続いて、図3で説明したレジスタrgの具体例を説明する。

【0030】

[レジスタ]

図4は、図3のレジスタrgを説明する図である。図4の(A)は、プロセッサのレジスタマップrmの一例を示す図である。また、図4(B)は、それぞれの割り込み要因の説明図である。図3で示したとおり、各ノードのサービスプロセッサボードB2のMBC 21は、レジスタrgを有する。また、レジスタrgは、ノードが有する複数の機能回路

50

(CPU、メモリアクセスコントローラ、電源等を示す。以下、ハードウェアと称する)から発生する割り込み要因を格納する。図4の(A)のレジスタマップrmによると、レジスタrgは、例えば、割り込み要因CK、FE、IL、EC、SC、PM、LD、II O、IMを格納する。ただし、割り込み要因は、図4の例に限定されるものではない。レジスタrgは、それぞれの割り込み要因を、レジスタマップrmに対応する所定のビット位置に格納する。

【0031】

また、図4の(B)において、割り込み要因CKは、例えば、システム制御装置22のクロック制御エラーを示す。割り込み要因FEは、プロセッサにおいて発生した致命的な(FATAL)エラーを示す。また、割り込み要因ILは処理対象が不正である旨のエラー、割り込み要因ECはデバッグ時に使用する信号、割り込み要因SCはシステム制御装置22、V1~V4から発生したリクエスト、割り込み要因PMは電源装置から発生したリクエストを示す。また、割り込み要因LDはクロスバスイッチ2の二重レーンの縮退に係るエラー、割り込み要因II OはI/Oコントローラ16(図3)において発生するエラー、割り込み要因IMはメモリアクセスコントローラ18(図3)において発生するエラーを示す。

【0032】

続いて、障害発生時の処理を説明する。本実施の形態例では、以下に説明する障害発生処理において、レジスタrgを使用する。

【0033】

[障害発生処理]

図5は、図1~図3で述べた、本実施の形態例における情報処理システム1の一部のノードにおいて障害が発生した場合の処理の流れを説明する図である。図5において、図2、図3で示したものと同一のものは、同一の記号で示す。

【0034】

情報処理システム1の全体の障害解析を行う場合に、複数のノードのうち、1つのノードが主体となって障害解析を行う方が効率的である。効率化のために、情報処理システム1は、1つノードのシステム制御装置をマスターのシステム制御装置に、他のノードのシステム制御装置をスレーブのシステム制御装置に設定する。または、情報処理システム1は、マスターのシステム制御装置の切り替え用として、さらに、1つノードのシステム制御装置を、マスターの代替用のシステム制御装置に設定してもよい。図5の例において、例えば、マスターのシステム制御装置は、クロスバスイッチ2(図1)を有する1つのノード(図2の2AB)のシステム制御装置V1である。以下、マスターのシステム制御装置V1を、マスターノード2ABのシステム制御装置V1、スレーブのシステム制御装置22を、スレーブノード1A~1P、2BB~2DBのシステム制御装置22、V2~V4と称する。

【0035】

前述したとおり、システム制御装置22、V1~V4は各々、ノード内のハードウェアの状態の監視、及び、ハードウェアの制御を行う。また、システム制御装置22、V1~V4は、ノード内の各ハードウェアから発生する割り込み要因を格納するレジスタrg(図4)を有する。ハードウェアの障害の一例としては、メモリ3、11aのデータ破損や、プロセッサ12aの内部障害等が挙げられる。

【0036】

ハードウェアの障害が発生すると割り込み信号が発生し、割り込み要因がレジスタrgに格納される(図示の矢印x1)。システム制御装置22、V1~V4は、レジスタrgを監視することによって、ハードウェアの障害の発生を検知すると、割り込み要因の発生をマスターノード2ABのシステム制御装置V1に通知する(図示の矢印x2、x3)。続いて、マスターノード2ABのシステム制御装置V1は、割り込み要因の発生の通知を受けると、各スレーブノード1A~1P、2BB~2DBのシステム制御装置22、V2~V4に対して、ハードウェアのエラー情報を有するログ情報の送信を指示する(図示の

10

20

30

40

50

矢印×4)。各スレーブノード1A～1P、2BB～2DBのシステム制御装置22、V2～V4は、マスターノード2ABのシステム制御装置V1からの指示に応じて、ノード内のログ情報を収集しマスターノード2ABのシステム制御装置V1に送信する(図示の矢印×5)。そして、各ノードにおいて取得されたログ情報が、マスターノード2ABのシステム制御装置V1に収集される。

【0037】

続いて、マスターノード2ABのシステム制御装置V1は、ログ情報の解析処理を行う。例えば、システム制御装置V1は、各ノードのログ情報に基づいて、障害ノード1B、及び、障害ノード1Bにおける障害部品を特定する。そして、システム制御装置V1は、ログ情報の解析処理によって特定された情報に基づいて、障害に対するリアクションを行う(図示の矢印×6、×7)。リアクションとは、例えば、各ノードで動作するアプリケーションに対する障害ノード1Bが有する共有メモリ3の領域へのアクセス抑止や、障害ノード1Bのハードウェアの停止制御である。

【0038】

図5で説明してきたように、一部のノードで障害が発生した場合、マスターノード2ABのシステム制御装置V1は、各ノードにおいて収集されたログ情報を受信する。そして、マスターノード2ABのシステム制御装置V1は、取得した各ノードのログ情報を解析することによって、障害ノード、及び、障害が発生した回路の特定処理の後、障害に対するリアクションを行う。

【0039】

[ログ情報の解析]

図6は、図5において述べた、マスターノード2ABのシステム制御装置V1におけるログ情報の解析処理(S1)の概要を説明する図である。図6において、点線で囲む工程S3、S4は、本実施の形態例において付加される処理である。

【0040】

まず、ログ情報の解析処理(S1)を説明する。情報処理システム1は、障害が発生した場合、ノードの継続動作が可能な場合であっても、ノードの予防保守として障害の内容を特定する必要がある。また、情報処理システム1は、障害が発生しているASIC(Application Specific Integrated Circuit、以下、ASICと称する)部分を特定する必要がある。例えば、障害が発生しているASIC部分の特定、及び、ノードの継続動作の可否判定のために、マスターノード2ABのシステム制御装置V1は、ログ情報の解析処理を行う。この実施の形態では、ASICは、例えば、CPU、メモリアクセスコントローラ、I/Oコントローラに対応する。

【0041】

マスターノード2ABのシステム制御装置V1は、収集したログ情報に基づいてログ解析を行う(S1)。ログ情報とは、例えば、割り込み要因発生時のエラー情報を含むエラー要因情報と、エラーログ詳細情報である。エラーログ詳細情報とは、例えば、ASICの履歴情報やダンプ情報等である。エラーログ詳細情報はデータ量が膨大であるため、マスターノード2ABのシステム制御装置V1は、ログ情報(S1)内の各解析工程S61～S65と平行して、エラーログ詳細情報を受信する。

【0042】

続いて、工程S1における各解析工程を説明する。システム制御装置V1は、まず、エラー要因情報に基づいて、エラーコードの解析処理を行う(S61)。次に、システム制御装置V1は、ノードのハードウェアそれぞれを対象として、エラー要因情報に基づいて、障害の有無の判定、及び、障害部分の特定処理を行う(S62～S65)。システム制御装置V1は、例えば、CPU12a、クロスバススイッチ2、メモリ3、11a等を対象として、エラー要因情報に基づいて、各ハードウェアにおける障害部分の判定、及び、障害部分の詳細の判定処理を行う。工程S62～S65の処理により、障害ノード、及び、障害が発生している回路が特定され、他の回路が正常に動作していることが確認される。なお、図6の例において、システム制御装置V1は、CPU12a、クロスバススイッチ2

10

20

30

40

50

、メモリ 3、11a を対象として解析処理を行っているが、対象となるハードウェアは、この例に限定されるものではない。

【0043】

障害ノードの特定、及び、障害部分の特定が行われると、システム制御装置 V1 は、エラーログ詳細情報の収集の完了を待機して、エラーログ詳細情報の登録処理を行う (S66)。続いて、システム制御装置 V1 は、エラー要因情報に基づいて、障害部分に対応するログ情報を示す代表ログの登録処理を行う (S67)。エラーログ詳細情報、及び、代表ログは、障害の原因の分析や、障害の復旧に必要な情報である。ログ情報に基づく解析処理が完了すると、システム制御装置 V1 は、障害の重要度に応じて、情報処理システム 1 からの障害ノードの切り離し制御を行う (S2)。障害ノードの切り離し制御とは、例えば、障害ノードのハードウェアの電源停止を示す。

10

【0044】

マスターノード 2AB のシステム制御装置 V1 は、障害の重要度に関わらず、ログ情報の解析処理 (S1) を実行する。また、ログ情報の解析処理では、ノード内のハードウェアそれぞれを対象として、詳細に障害部分の判定処理を行うため、時間を要する。また、エラーログ詳細情報の転送処理は、エラーログ詳細情報のデータ量が膨大であるため、時間を要する。このため、ログ情報の解析処理には、数十秒～数分 (30 秒～5 分) 程度の時間がかかる。即ち、障害の発生から障害ノードの切り離し制御まで、5 分程度の時間を要する。

【0045】

20

しかしながら、情報処理システム 1 は、障害が発生してから短時間で運用を再開することが望ましい。運用の再開処理では、正常ノードが障害ノードの処理を引き継ぐため、情報処理システム 1 は、障害ノードを早急に特定する必要がある。また、複数のノード間でメモリを共有する情報処理システム 1 では、障害の発生に起因して、共有メモリ 3 の破損や不整合等の二次障害が発生する恐れがある。共有メモリ 3 に対する二次障害を抑止するために、早急に、障害ノードのメモリへのアクセス抑止を行うことが求められる。障害発生から障害ノードのメモリへのアクセス抑止まで、例えば、1 秒程度で完了することが望ましい。

【0046】

そこで、本実施の形態例において、マスターノード 2AB のシステム制御装置 V1 は、ログ情報の解析処理 (S1) の前に、FNL (Fail Node List、以下、FNL と称する) 解析処理 (S3) を行って障害ノードを特定し、障害ノードのメモリへのアクセス抑止を行う (S4)。

30

【0047】

本実施の形態例のマスターノード 2AB のシステム制御装置 V1 は、他のノードの割り込み要因の発生に応じてレジスタ rg の割り込み要因を受信し、割り込み要因のうち、障害として検出すべき割り込み要因を抽出する。そして、マスターノード 2AB のシステム制御装置 V1 は、抽出結果に応じて障害ノードを特定し、障害ノードのメモリへのアクセスを抑止後、他のノードから受信したログ情報に基づいて障害ノードの切り離し制御を行う。

40

【0048】

具体的に、マスターノード 2AB のシステム制御装置 V1 は、FNL 解析処理 (S3) として、まず、発生中の割り込み要因を各ノードから取得する (S51)。続いて、システム制御装置 V1 は、取得した割り込み要因のうち、障害として検出すべき割り込み要因を抽出する (S52)。次に、システム制御装置 V1 は、抽出した割り込み要因のうち、波及先の割り込み要因を FNL 解析の対象から除外する (S53)。即ち、システム制御装置 V1 は、抽出した割り込み要因のうち、別の割り込み要因に起因して発生した割り込み要因を、FNL 解析処理の対象外とする。

【0049】

続いて、システム制御装置 V1 は、複数の割り込み要因が抽出された場合、各割り込み

50

要因の優先度を判定する（Ｓ５４）次に、システム制御装置Ｖ１は、優先度の高い順に割り込み要因を選択し、割り込み要因に対応する障害ノードを特定する（Ｓ５５）。次に、システム制御装置Ｖ１は、障害ノードの他ノードからのメモリへのアクセス抑止処理を行う（Ｓ５６）。即ち、システム制御装置Ｖ１は、障害ノードが有する共有メモリ３の領域に対するアクセスを抑止する。各工程の詳細については、後述する。続いて、システム制御装置Ｖ１は、ログ情報の解析処理を実行し（Ｓ１）、障害ノードの情報処理システム１からの切り離し制御を行う（Ｓ２）。

【００５０】

図６で説明してきたように、本実施の形態例において、システム制御装置Ｖ１は、ＦＮＬ解析処理（Ｓ３）として、ログ情報の代わりに割り込み要因に基づいて、障害ノードを特定し、障害ノードのメモリに対する他ノードからのアクセス抑止処理を行う（図６のＳ５６）。障害ノードのメモリに対するアクセスを抑止することによって、システム制御装置Ｖ１は、共有メモリ３の二次障害を早急に抑止し、障害発生時の障害ノードによる他ノードへの影響を低減する。

【００５１】

そして、システム制御装置Ｖ１は、障害ノードのメモリへのアクセス抑止後、ログ情報の解析処理（Ｓ１）を行って、障害が発生しているＡＳＩＣ部分を特定し、ノードの継続動作の可否を判定する。そして、システム制御装置Ｖ１は、ログ情報の解析処理の結果に基づいて、障害ノードの情報処理システム１からの切り離し制御（Ｓ２）を行う。

【００５２】

図７は、図６のログ情報の解析処理（図６のＳ１）、及び、ＦＮＬ解析処理（Ｓ３）に要する時間を例示する図である。図７の（Ａ）は、ログ情報の解析処理（Ｓ１）から障害ノードの切り離し制御（Ｓ２）までの時間を示す図であって、図７の（Ｂ）は、ＦＮＬ解析処理（Ｓ３）から障害ノードのメモリへのアクセス抑止処理（Ｓ４）までの時間を示す図である。

【００５３】

図７の（Ａ）では、ログ情報の解析処理（Ｓ１）の後、障害ノードの切り離し制御（Ｓ２）が行われる。前述したとおり、ログ情報の解析処理（Ｓ１）は、ハードウェアそれぞれに対するログ情報の解析処理やエラーログ詳細情報の転送処理に伴って、時間を要する。図７の（Ａ）によると、障害発生から障害ノードの切り離し制御まで、期間 t_1 に示す時間を要する。

【００５４】

一方、図７の（Ｂ）において、ＦＮＬ解析処理（Ｓ３）では、システム制御装置Ｖ１は、障害として検出すべき割り込み要因に基づいて障害ノードを特定する。また、ＦＮＬ解析処理（Ｓ３）では、エラーログ詳細情報の転送が不要であり、割り込み要因（３２ビット程度）のデータ量は小さい。したがって、システム制御装置Ｖ１は、障害ノードを早急に特定することが可能になるため、図７の（Ｂ）によると、障害発生から障害ノードのメモリへのアクセス抑止処理までの時間 t_2 は、時間 t_1 に対して大幅に短縮される。

【００５５】

ここで、本実施の形態例におけるマスターノード２ＡＢのシステム制御装置Ｖ１、及び、スレーブノード１Ａ～１Ｐ、２ＢＢ～２ＤＢのシステム制御装置２２、Ｖ２～Ｖ４のソフトウェアモジュール図を説明する。

【００５６】

〔ソフトウェアモジュール図〕

図８は、本実施の形態における情報処理システムの各ノードのソフトウェアモジュール図である。図８は、マスターノード２ＡＢのシステム制御装置Ｖ１、及び、スレーブノード１Ａ～１Ｐ、２ＢＢ～２ＤＢのシステム制御装置２２、Ｖ２～Ｖ４のブロック図を有する。初めに、スレーブノード１Ａ～１Ｐ、２ＢＢ～２ＤＢのシステム制御装置２２、Ｖ２～Ｖ４のブロックを説明する。ここでは、スレーブノード１Ａのシステム制御装置２２について説明する。

【 0 0 5 7 】

図 8 において、スレーブノード 1 A のシステム制御装置 2 2 は、例えば、F N L (Fail Node List) ドライバ 5 4、F N L (Fail Node List) 部 5 0、ハード内制御部 6 1、R A S (Reliability Availability Serviceability、以下、R A S と称する) 6 2、X S C F (eXtended System Control Facility、以下、X S C F と称する) コマンド部 6 3、ハイパーバイザ 6 4 を有する。また、F N L 部 5 0 は、例えば、F N L (Fail Node List) 制御部 5 1、F N L (Fail Node List) 更新依頼受信制御部 5 2、F N L (Fail Node List) 更新部 5 3 を有する。

【 0 0 5 8 】

ハード内制御部 6 1 は、例えば、電源やプロセッサ (図 8 では、C P U と記す) やクロ
スバスイッチ (図 8 では、X B と記す) 等のハードウェアに対するアクセス処理を行う H
A P (Hardware Access Program、以下、H A P と称する) 6 5 を有する。そして、ハー
ド内制御部 6 1 は、当該ハードウェアに対するアクセス処理における割り込み要因の発生
を検知し、F N L 部 5 0 の F N L 更新依頼受信制御部 5 2 に通知する。また、R A S 6 2
は、A S I C における割り込み要因の発生を検知し、F N L 更新依頼受信制御部 5 2 に通
知する。また、X S C F コマンド部 4 3 は、例えば、ハイパーバイザにおける割り込み要
因の発生を検知し、F N L 更新依頼受信制御部 5 2 に通知する。

【 0 0 5 9 】

また、F N L 更新依頼受信制御部 5 2 は、各部からの割り込み要因の発生の通知を取得
して、F N L 制御部 5 1 に出力する。そして、F N L 制御部 5 1 は、F N L ドライバ 5 4
を介して、割り込み要因の発生をマスターノード 2 A B のシステム制御装置 V 1 の F N L
制御部 3 2 に通知する。また、F N L 制御部 5 1 は、マスターノード 2 A B のシステム制
御装置 V 1 からの割り込み要因の収集依頼に応答して、発生している割り込み要因を収集
し F N L ドライバ 5 4 を介してマスターノード 2 A B のシステム制御装置 V 1 に送信する
。F N L 更新部 5 3 は、マスターノード 2 A B のシステム制御装置 V 1 からの F N L 更新
指示に基づいて、F N L (Fail Node List、図 8 には図示せず) を更新する。F N L とは
、メモリを共有するノードそれぞれのアクセス処理の可否を管理するリストである。情報
処理システム 1 における各ノードは、F N L に基づいて、アクセス抑止対象のノードを検
知する。

【 0 0 6 0 】

また、図 8 において、マスターノード 2 A B のシステム制御装置 V 1 は、例えば、F N
L ドライバ 3 5、F N L 部 3 0、X S C F コマンド部 4 3、ハード内制御部 4 1、R A S
4 2、F N D B 3 6 (Fail Node DB、以下、F N D B と称する) を有する。また、F N L
部 3 0 は、例えば、F N L 解析部 3 1、F N L 制御部 3 2、F N L 更新部 3 3、F N L 更
新依頼受信制御部 3 4 を有する。ハード内制御部 4 1、R A S 4 2、X S C F コマンド部
4 3、F N L 更新依頼受信制御部 3 4 の処理は、スレーブノード 1 A のシステム制御装置
2 2 と同様である。

【 0 0 6 1 】

F N L 部 3 0 の F N L 制御部 3 2 は、スレーブノード 1 A のシステム制御装置 2 2 から
割り込み要因発生の通知を受信すると、F N L ドライバ 3 5 を介して、各スレーブノード
1 A のシステム制御装置 2 2 に対して、発生中の割り込み要因の収集を指示する。F N L
解析部 3 1 は、各ノードのシステム制御装置 2 2 から収集した割り込み要因に基づいて、
F N D B 3 6 を参照し障害ノードを特定する。F N D B 3 6 は、F N L 解析における解析
論理の定義を有するファイルである。そして、F N L 解析部 3 1 は、特定した障害ノード
の情報に基づいて、各ノードの F N L 更新部 3 3 に F N L の更新を指示する。

【 0 0 6 2 】

図 9 は、図 8 において説明した F N L (Fail Node List) 4 0 の一例を示す図である。
図 9 例において、情報処理システム 1 は、例えば、図 2 で示したように、1 6 個のノード
S B 0 0 ~ S B 1 5 を有し、各ノードはメモリを共有する。そこで、図 9 の F N L 4 0 は
、1 6 個のノードそれぞれに対するアクセス処理の可否を管理する値を有する。例えば、

10

20

30

40

50

値「0」の場合、対象ノードの共有メモリ3に対するアクセスが許可されることを示す。一方、値「1」の場合、対象ノードの共有メモリ3に対するアクセスが抑止されることを示す。

【0063】

続いて、割り込み要因が発生した後、図9において説明したFNL40が更新されるまでの処理の流れを、図8において説明したソフトウェアモジュールに対応して、時系列に説明する。

【0064】

[ソフトウェアモジュールの処理の流れ]

図10は、割り込み要因が発生しFNLが更新される間のマスターノード2ABのシステム制御装置V1、及び、スレーブノード1Aのシステム制御装置22における、処理の流れを時系列に説明する図である。図10において、図8で示したものと同一のものは同一の記号で示してある。

【0065】

図10の例において、例えば、スレーブノードの一部のASICにおいて障害が発生する。障害の発生により割り込み信号が発生し、障害が発生したハードウェアに対応する割り込み要因がレジスタrgに登録される。スレーブノード1Aのシステム制御装置22は、障害の発生を検知すると（図示の矢印g1）、障害が発生したことをマスターノード2ABのシステム制御装置V1におけるハード内制御部41に通知する（図示の矢印g2）。

【0066】

システム制御装置V1におけるハード内制御部41は、障害の発生の通知を受けて、各スレーブノード1Aのシステム制御装置22に対して割り込み要因の収集を指示する（図示の矢印g3）。マスターノード2ABのシステム制御装置V1の通知に応答して、各スレーブノード1Aのシステム制御装置22は、発生している割り込み要因を、FNL部50を介して（図示の矢印g4）、システム制御装置V1のFNL部30に送信する（図示の矢印g5）。

【0067】

この結果、各スレーブノード1Aのシステム制御装置22の割り込み要因が収集される。割り込み要因のデータ量は小さい。このため、システム制御装置V1は、短時間で、各スレーブノード1Aのシステム制御装置22の割り込み要因を取得することができる。また、各システム制御装置22、V1～V4間が高速通信を介して接続される場合、システム制御装置V1は、さらに、高速に、各システム制御装置22、V2～V4の割り込み要因を取得することができる。

【0068】

マスターノード2ABのシステム制御装置V1におけるFNL部30は、全てのスレーブノード1Aにおける割り込み要因を収集すると、FNL解析部31に解析処理を指示する（図示の矢印g6）。そして、FNL解析部31は、収集した割り込み要因に基づいて障害ノードを特定し、FNL部30に出力する（図示の矢印g7）。続いて、FNL部30は、障害ノードの情報に基づいて、各スレーブノード1Aのシステム制御装置のFNL部50にFNL40（図9）の更新を指示する（図示の矢印g8）。FNL40の更新指示を受信すると、スレーブノード1AのFNL部50は、FNL更新部53にFNL40の更新を実行させる（g9、g10）。

【0069】

図10のフローチャート図に示した処理の流れに基づいて、FNL解析処理、及び、FNLの更新処理が行われる。続いて、各処理の詳細をフローチャート図に基づいて説明する。

【0070】

[FNL解析処理、FNL更新処理]

図11は、図8の本実施の形態例におけるFNL解析部31の処理、及び、FNL更新

10

20

30

40

50

部 3 3 の処理を説明するフローチャート図である。初めに、例えば、マスターノード 2 A B のシステム制御装置 V 1 における F N L 解析部 3 1 は、電源障害が発生しているか否かを判定する (S 2 1)。電源障害が発生している場合、電源障害の対応が優先されるため、F N L 解析部 3 1 は処理を終了する。

【 0 0 7 1 】

一方、電源障害が発生していない場合 (S 2 1 の N O)、F N L 解析部 3 1 は、割り込み要因を取得する (S 2 2)。前述したとおり、F N L 解析部 3 1 は、割り込み要因が発生したノードからの通知に应答して、各ノードから割り込み要因を取得する。続いて、F N L 解析部 3 1 は、収集された割り込み要因から、障害として検出すべき割り込み要因を抽出する (S 2 3)。例えば、F N L 解析部 3 1 は、例えば、割り込み要因のうち、ノードの停止が必要となる障害に対応する割り込み要因を抽出する。即ち、F N L 解析部 3 1 は、ノードが継続して動作可能な割り込み要因を抽出の対処としない。

10

【 0 0 7 2 】

本実施の形態例において、マスターノード 2 A B のシステム制御装置 V 1 は、障害として検出すべき割り込み要因として、例えば、図 4 に例示した割り込み要因のうち、割り込み要因 C K、F E を抽出する (S 2 3)。割り込み要因 C K、F E は、C P U が停止する障害要因であって、割り込み要因 C K、F E 以外の割り込み要因については機能の一部が縮退する故障要因であるためである。ただし、この例に限定されるものではなく、システム制御装置 V 1 は、別の割り込み要因を、障害として検出すべき割り込み要因としてもよい。

20

【 0 0 7 3 】

なお、図 4 のレジスタマップ r m は、プロセッサのレジスタマップである。情報処理システム 1 では、プロセッサのレジスタマップの他に、クロスバススイッチ用のレジスタマップや M B C 用のレジスタマップが存在する。また、クロスバススイッチ用のレジスタに格納される割り込み要因については、例えば、複数の割り込み要因のうち、内部障害、及び、ポート障害に対応する割り込み要因が抽出対象となる。また、M B C 用のレジスタに格納される割り込み要因については、例えば、複数の割り込み要因のうち、割り込み要因 F E が抽出対象となる。

【 0 0 7 4 】

次に、F N L 解析部 3 1 は、波及先の割り込み要因を抑止する (S 2 4)。ここで、割り込み要因は、波及元の割り込み要因と、波及元の割り込み要因に基づいて誘発された波及先の割り込み要因とに区分される。F N L 解析部 3 1 は、波及先の割り込み要因を除外して、波及元の割り込み要因のみに絞り込む。

30

【 0 0 7 5 】

具体的に、例えば、あるノードのプロセッサにおいて障害が発生した場合、同一ノード内のクロスバススイッチの接続部や、プロセッサの他の部分に障害が波及することがある。この場合、プロセッサにおいて発生した障害に対応する割り込み要因が上位の割り込み要因、クロスバススイッチの接続部やプロセッサの他の部分において発生した障害に対応する割り込み要因が下位の割り込み要因に相当する。即ち、上位の割り込み要因が波及元の割り込み要因に該当し、下位の割り込み要因が波及先の割り込み要因に該当する。

40

【 0 0 7 6 】

F N L 解析部 3 1 は、波及先の割り込み要因の抑止処理 (S 2 4) によって、波及先の割り込み要因を F N L 解析の対象から除外するため、波及元の割り込み要因に対応するノードのみを障害ノードとして特定する。即ち、F N L 解析部 3 1 は、波及先の割り込み要因に対応するノードを障害ノードとして特定することを回避し、真に障害が発生するノードのみを、障害ノードとして特定する。

【 0 0 7 7 】

続いて、F N L 解析部 3 1 は、抽出した割り込み要因それぞれの優先度を取得する (S 2 5)。そして、F N L 解析部 3 1 は、優先度の高い割り込み要因から順に、当該割り込み要因に対応して、障害ノードを特定すると共に、障害ノードに対する制御内容を取得す

50

る（Ｓ２６）。そして、ＦＮＬ更新部３３、５３は、取得した制御内容に基づいてＦＮＬ４０を更新し、共有メモリ３上の障害ノードの領域への、他ノードからのアクセスを抑止する（Ｓ２７）。続いて、ＦＮＬ解析部３１、及び、ＦＮＬ更新部３３、５３は、次に優先度の高い割り込み要因を対象として、工程Ｓ２６、Ｓ２７の処理を行う。そして、ＦＮＬ解析部３１、及び、ＦＮＬ更新部３３、５３は、抽出した全ての割り込み要因を対象として工程Ｓ２６、Ｓ２７の処理を行うと、ＦＮＬの解析処理、及び、ＦＮＬの更新処理を終了する。

【００７８】

続いて、図１１のフローチャート図における波及先の割り込み要因の抑止処理（工程Ｓ２４）の詳細を説明する。

10

【００７９】

〔波及先割り込み要因の抑止（図１１の工程Ｓ２４）〕

図１２は、波及先の割り込み要因の抑止処理を説明するフローチャート図である。まず、ＦＮＬ解析部３１は、ＦＮＤＢ３６を参照し、抽出した割り込み要因（この例では、ＣＫ、ＦＥ）が、波及元の割り込み要因であるか否かを判定する（Ｓ１１）。ＦＮＤＢ３６については、次の図１３に基づいて説明する。そして、抽出した割り込み要因が、下位の割り込み要因（波及先の割り込み要因）ではない場合（Ｓ１２のＮＯ）、即ち、波及元の割り込み要因である場合、ＦＮＬ解析部３１は波及先の割り込み要因の消し込み処理を終了する。

【００８０】

20

一方、抽出した割り込み要因が下位の割り込み要因である場合（Ｓ１２のＹＥＳ）、ＦＮＬ解析部３１は、各ノードのレジスタｒｇを参照し、当該下位の割り込み要因に対応する上位の割り込み要因を示す、波及元の割り込み要因が発生しているか否かを判定する（Ｓ１３）。波及元の割り込み要因が発生している場合（Ｓ１４のＹＥＳ）、ＦＮＬ解析部３１は、下位の割り込み要因を抑止し、ＦＮＬ解析処理の対象外とする（Ｓ１５）。

【００８１】

一方、波及元の割り込み要因が発生していない場合（Ｓ１４のＮＯ）、ＦＮＬ解析部３１は波及先の割り込み要因の抑止処理を終了する。つまり、例えば、発生した割り込み要因が波及先の割り込み要因である場合であっても、波及元の割り込み要因が発生していない場合、ＦＮＬ解析部３１は、波及先の割り込み要因を抑止しない。

30

【００８２】

図４、図１２で述べてきたとおり、本実施の形態例では、ＦＮＬ解析部３１は、割り込み要因のうち、ノードの停止が必要となる障害に対応する割り込み要因に限定して、ＦＮＬ解析を行う（図１１のＳ２３）。また、ＦＮＬ解析部３１は、さらに、波及元となる割り込み要因に限定して、ＦＮＬ解析を行う（図１１のＳ２４）。したがって、ＦＮＬ解析部３１は、ノードの停止が必要となる障害に対応する割り込み要因であって、真に障害が発生しているノードにおける割り込み要因を抽出することができる。即ち、ＦＮＬ解析部３１は、ノードの停止が必要となる最小限の障害ノードを、効率的に特定することができる。

【００８３】

40

ここで、ＦＮＬ解析部３１が、ＦＮＬ解析処理において参照するＦＮＤＢ３６の具体例を説明する。ＦＮＤＢ３６は、ＦＮＬ解析における解析論理の定義を有する。

【００８４】

〔ＦＮＤＢの具体例〕

図１３は、ＦＮＤＢ３６の具体例を示す図である。図１３の（Ａ）は、ＦＮＬ解析における解析論理を定義する定義テーブルｔｂ１を示す図である。図１３の（Ｂ）は、定義テーブルｔｂ１に記述される各エントリの一部を説明する図である。定義テーブルｔｂ１は、例えば、共通定義フレームとデータ定義ブロックとを有する。共通定義フレームは、定義テーブルｔｂ１の版数や定義開始の宣言を有する。

【００８５】

50

図13の(A)によると、定義テーブルtb1には、例えば、割り込み要因に対応して、優先度(priority)、アクション番号(act)、エントリ抑止条件(ent_dis)等の定義を有する。また、図13の(B)によると、優先度(priority)は、割り込み要因の優先度を示す。例えば、優先度は数値で示される。また、アクション番号(act)は、割り込み要因に対応する、障害ノード及び共有メモリへのアクセス抑止処理の種別を示す。アクセス抑止処理の種別の詳細については、別の図14に基づいて説明する。

【0086】

そして、エントリ抑止条件(ent_dis)は、割り込み要因に対応して、当該割り込み要因に対して論理的に上位の割り込み要因を示す。即ち、エントリ抑止条件(ent_dis)は、割り込み要因に対応する、波及元の割り込み要因を示す。エントリ抑止条件(ent_dis)がブランクの場合、割り込み要因に対応する、波及元の割り込み要因が存在しないことを示す。

10

【0087】

FNL解析部31は、例えば、図13の定義テーブルtb1の記述cd1を参照して、波及先の割り込み要因の抑止処理(図11のS24)、優先度の取得処理(図11のS25)、障害ノードの特定、及び、アクセス抑止処理の取得処理(図11のS26)を行う。例えば、FNL解析部31は、定義テーブルtb1内のadrs(割り込み要因番号)列cd2を参照し、adrsの値が、割り込み要因に対応する割り込み要因番号と一致する行を探索する。

【0088】

20

例えば、レジスタrgに基づいて収集した割り込み要因が、割り込み要因CKである場合を例に挙げる。割り込み要因CKは、図4のレジスタマップrmによると、1ビット目の位置に位置することから、割り込み要因CKに対応する割り込み要因番号は値「0x00000001」である。そこで、FNL解析部31は、割り込み要因CKに対応して、定義テーブルtb1における2行目の定義情報を検出する。なお、割り込み要因FEに対応する割り込み要因は値「0x00000003」である。そこで、FNL解析部31は、割り込み要因FEに対応して、定義テーブルtb1における1行目の定義情報を検出する。

【0089】

続いて、FNL解析部31は、検出した2行目の定義情報のうち、エントリ抑止条件(ent_dis)に対応する項目cd5を参照する。図13の例において、2行目の定義情報におけるエントリ抑止条件(ent_dis)は、ブランクである。したがって、FNL解析部31は、割り込み要因CKを、波及元の割り込み要因として判定する(図11のS24)。また、FNL解析部31は、2行目の定義情報における優先度(priority)の項目cd3に基づいて、優先度「0x01」(図11のS25)を取得すると共に、アクション番号(act)の項目cd4に基づいてアクション番号「0x01」を取得する。なお、図13の例における優先度は、値が小さいほど高い。

30

【0090】

一方、FNL解析部31は、割り込み要因FEに対応して、1行目の定義情報を検出する。図13の例において、定義テーブルtb1は、1行目の定義情報におけるエントリ抑止条件(ent_dis)として、定義“/XBBOX/XBUX/GXB/FN_XB_SND”の記述を有する(cd5)。それぞれの定義XBBOX、XBUX、GXB、FN_XB_SNDは、割り込み要因を示し、割り込み要因FEの波及元の割り込み要因に該当する。定義XBBOXは、例えば、クロスバボックスにおける割り込み要因を、定義FN_XB_SNDは、例えば、クロスバスイッチの送信部における割り込み要因を示す。したがって、FNL解析部31は、割り込み要因FEを波及先の割り込み要因として判定し、FNL解析の対象から除外する。

40

【0091】

具体例に基づいて説明する。例えば、複数のノードSB00~SB03が、クロスバスイッチ2を備えるノードXB00と接続する情報処理システムを例示する。具体例では、

50

例えば、ノードS B 0 2においてクロック制御エラーが発生すると共に、クロスバスイッチ2を備えるノードX B 0 0において、ポート障害が発生する。

【0092】

いずれかの割り込み要因を検知すると、マスターノード2 A Bのシステム制御装置V 1は、各ノードにおいて発生する割り込み要因を収集する(S 2 2)。そして、システム制御装置V 1は、ノードS B 0 2において発生したクロック制御エラーに対応する割り込み要因C Kと、ノードX B 0 0において発生したポート障害に対応する割り込み要因を取得する。具体例において、割り込み要因C K、及び、ポート障害に対応する割り込み要因は、抽出対象の割り込み要因である(S 2 3)。

【0093】

続いて、FNL解析部3 1は、図1 3のFNDB 3 6を参照して(c d 5)、各割り込み要因が波及元の割り込み要因であるか否かを判定する。図1 3において前述したとおり、割り込み要因C Kは、波及元の割り込み要因である。また、図示していないが、具体例において、ポート障害に対応する割り込み要因は、波及元の割り込み要因である。このため、FNL解析部3 1は、ポート障害に対応する割り込み要因を抑止しない(S 2 4)。次に、FNL解析部3 1は、FNDB 3 6を参照して(c d 3)、各割り込み要因に対応する優先度を取得する(S 2 5)。図1 3において前述したとおり、割り込み要因C Kの優先度(p r i o)は、優先度「0 x 0 1」である。また、具体例において、図示していないが、ポート障害に対応する割り込み要因の優先度は、優先度「0 x 0 5」である。したがって、FNL解析部3 1は、割り込み要因C Kを、ポート障害に対応する割り込み要因よりも優先する。

【0094】

続いて、図1 3の定義テーブルt b 1のアクション番号(a c t)を説明する。アクション番号(a c t)は、割り込み要因に対応する、障害ノード及び共有メモリへのアクセス抑止処理の種別を示す。図1 3の定義テーブルt b 1の記述c d 4によると、割り込み要因C Kのアクション番号(a c t)は、「0 x 0 1」である。また、図示していないが、ポート障害に対応する割り込み要因のアクション番号(a c t)は、例えば、「0 x 1 2」である。次の図1 4に基づいて、アクション番号(a c t)に対応する制御情報について説明する。

【0095】

図1 4は、アクション番号(a c t)を有する定義テーブルt b 2の具体例を示す図である。図1 4の(A)は、アクション番号(a c t)に対応して制御情報(r u l e)の記述を有する定義テーブルt b 2を示す図であって、図1 4の(B)は、定義テーブルt b 2に記述される制御情報(r u l e)の各エントリの一部を説明する図である。FNDB 3 6は、例えば、図1 3の定義テーブルt b 1に加えて、図1 4に示す定義テーブルt b 2を有する。定義テーブルt b 2は、例えば、定義テーブルt b 2の版数や定義開始の宣言を有する共通定義フレームと、データ定義ブロックとを有する。

【0096】

図1 4の定義テーブルt b 2のデータ定義ブロックは、アクション番号(a c t)に対応して、障害ノードに対する制御情報(r u l e)の記述c d 6を有する。例えば、定義テーブルt b 2は、アクション番号(a c t)「0 x 0 1」に対応して、制御情報(r u l e)として、エントリFNL_UPDATEを有する。また、定義テーブルt b 2は、アクション番号(a c t)「0 x 0 2」に対応して、制御情報(r u l e)として、エントリFNL_UPDATE_DESTを有する。同様にして、定義テーブルt b 2は、アクション番号(a c t)「0 x 1 1」に対応して、制御情報(r u l e)として、エントリGCSM_DEGRADEを、アクション番号(a c t)「0 x 1 2」に対応して、制御情報(r u l e)として、エントリGCSM_DEGRADE_DESTを有する。

【0097】

図1 4の(B)によると、エントリFNL_UPDATEは、割り込み要因が検出されたノードを障害ノードとして特定し、当該障害ノードを停止対象のノードとしてメモリア

10

20

30

40

50

クセスの制御を行うことを示す。この場合、マスターノード2 A Bのシステム制御装置V 1は、例えば、割り込み要因が検出されたノード（障害ノード）の共有メモリ3内の領域に対する、他のノードからのアクセスを抑止する。アクセスが抑止されることにより、障害ノードのメモリが共有メモリから切り離され、情報処理システムの継続稼動が可能になる。また、図14の（B）によると、エントリFNL_UPDATE_DESTは、割り込み要因が検出されたノードに接続されたノードを障害ノードとして特定し、当該障害ノードを停止対象のノードとしてメモリアクセスの制御を行うことを示す。

【0098】

さらに、エントリGCSM_DEGRADEは、割り込み要因が検出されたノードを障害ノードとして特定し、当該障害ノードを機能縮退対象のノードとしてメモリアクセスの制御を行うことを示す。さらに、エントリGCSM_DEGRADE_DESTは、割り込み要因が検出されたノードに接続されたノードを障害ノードとして特定し、当該障害ノードを機能縮退対象のノードとしてメモリアクセスの制御を行うことを示す。障害ノードの機能縮退とは、例えば、障害ノードがクロスバスイッチ2を備えるノードである場合に、図2で説明したクロスバスイッチの二重の回線を一回線に縮退させる制御を示す。

【0099】

具体例において、前述したとおり、割り込み要因CKのアクション番号（act）は、値「0x01」であって、ポート障害に対応する割り込み要因のアクション番号（act）は、例えば、値「0x12」である。したがって、FNL解析部31は、割り込み要因CKに対応して割り込み要因が発生したノードSB02を障害ノードとして特定する。そして、FNL更新部33、53は、ノードSB02に係るメモリアクセスの制御（FNL_UPDATE）を行う。また、FNL解析部31は、ポート障害に対応する割り込み要因が発生したノードXB00に接続されるノードSB00～SB03を障害ノードとして特定する。そして、FNL更新部33、53は、ノードSB00～SB00に係る機能縮退制御（GCSM_DEGRADE_DEST）を行う。

【0100】

ただし、具体例において、割り込み要因CKは、ポート障害に対応する割り込み要因よりも優先される。そこで、FNL更新部33、53は、まず、割り込み要因CKに対応する障害ノードのメモリへのアクセス抑止処理を行う（S26、S27）。例えば、それぞれのノードのFNL更新部33、53はFNL40を更新し、ノードSB02のメモリに対する他のノードからのアクセスを抑止する。

【0101】

続いて、FNL更新部33、53は、ポート障害に対応する割り込み要因に対応して、ノードSB00～SB03のメモリへのアクセス抑止処理を行う（S26、S27）。例えば、それぞれのノードのFNL更新部33、53はFNL40を更新し、ノードXB00からノードSB00～SB03に対するアクセス処理における二重の回線を片側の回線に縮退させる。回線が縮退されたことにより、ノードSB00～SB03の共有メモリに対するアクセス経路が減少する。

【0102】

図15は、具体例におけるメモリの抑止範囲を説明する図である。具体例によると、ノードSB02においてクロック制御エラーが発生した場合、他ノードによる、ノードSB02の共有メモリ3に対するアクセスが抑止される（ac1）。一方、ノードXB00においてポート障害が発生した場合、ノードXB00とノードSB00～SB03との間の回線が片側に縮退される（ac2）。即ち、図15の例において、回線n1、n3、n5、n7が使用できない状態となる。なお、例えば、回線n1、n3、n5、n7が既に停止されている状態で、さらに、ポート障害が発生した場合、全ての回線n1～n8が使用不可状態となり、ノードSB00～SB03の共有メモリ3に対するアクセスが行えなくなる。

【0103】

図15に示すように、割り込み要因CKに対応するアクセスの抑止範囲は、クロスバス

10

20

30

40

50

イチ 2 のポート障害によるアクセスの抑止範囲より狭い。図 15 の例では、アクセスの抑止範囲のより広いポート障害に対応する割り込み要因の優先度が低く設定されることにより、ポート障害によるアクセス抑止処理は、割り込み要因 C K によるアクセス抑止処理より後から行われる。アクセスの抑止範囲の広い割り込みがより後から行われることによって、情報処理システム 1 の性能がより長時間維持される。図 15 の例のように、例えば、割り込み要因の優先度は、システム制御装置 1 の性能をより高性能に維持するために、抑止範囲がより小さい割り込み要因ほど、より高い優先度が設定される。

【 0 1 0 4 】

以上のように、本実施の形態例における情報処理システムは、ノードの各々は、複数の機能回路と機能回路を制御する制御装置と、複数の機能回路から発生する割り込み要因を格納するレジスタとを有する。また、情報処理システムにおける複数のノードのうちの 1 のノードの制御装置は、他のノードの割り込み要因の発生に応じてレジスタの割り込み要因を受信し、割り込み要因のうち、障害として検出すべき割り込み要因を抽出して、抽出結果に応じて障害ノードを特定する。そして、制御装置は、障害ノードのメモリへのアクセスを抑止後、他のノードから受信したログ情報に基づいて障害ノードの切り離し制御を行う。

10

【 0 1 0 5 】

本実施の形態例における情報処理システムは、割り込み要因に基づくことにより障害ノードを高速に特定することができる。また、本実施の形態例における情報処理システムは、複数の割り込み要因のうち、ノードの停止が必要となる、障害として検出すべき割り込み要因を対象として、障害ノードを特定するため、より効率的に、障害ノードを特定することができる。

20

【 0 1 0 6 】

また、本実施の形態例における情報処理システムは、高速に、障害ノードを特定することができるため、障害ノードのメモリへのアクセスを早急に抑止することができ、共有メモリへの二次障害を回避することができる。即ち、情報処理システムは、障害発生時に障害ノードによる他ノードの影響を早急に低減することができる。また、情報処理システムは、障害ノードを高速に特定できるため、障害の発生時における、障害ノードから正常ノードへの運用の切り替えにかかるオーバーヘッドを低減できる。

【 0 1 0 7 】

30

また、本実施の形態例における情報処理システムにおいて、1 のノードの制御装置は、障害として検出すべき割り込み要因の波及元となる割り込み要因が発生しているか否かを判定し、発生していない場合に、割り込み要因に対応するノードを障害ノードとして特定し、発生している場合に、波及元となる割り込み要因に対応するノードを障害ノードとして特定する。

【 0 1 0 8 】

本実施の形態例における情報処理システムは、波及元となる割り込み要因に対応するノードを障害ノードとして特定することにより、複数の割り込み要因が連動して発生している場合に、複数の割り込み要因のうち、波及元の割り込み要因のみを対象として、当該波及元の割り込み要因に対応する障害ノードを特定することができる。

40

【 0 1 0 9 】

また、本実施の形態例における情報処理システムにおいて、第 1 のノードの制御部は、障害として検出すべき割り込み要因を複数抽出した場合に、割り込み要因の優先度に基づいて、特定した障害ノードのメモリへのアクセスを抑止する。

【 0 1 1 0 】

本実施の形態例における情報処理システムは、割り込み要因の優先度に基づいて、障害ノードのメモリへのアクセスの抑止処理の順を制御するため、割り込み要因に応じて、障害ノードのメモリへのアクセスの抑止処理の順を調整することができる。また、情報処理システムは、メモリへのアクセス抑止範囲の広い割り込み要因の優先度を低く設定することによって、情報処理システムの性能をより長く維持することができる。

50

【 0 1 1 1 】

また、本実施の形態例における情報処理システムにおいて、1のノードの制御装置は、障害として検出すべき割り込み要因がデータ処理を実行するノードにおいて発生した割り込み要因である場合に、発生元のノードを障害ノードとして特定する。また、1のノードの制御装置は、障害として検出すべき割り込み要因が網結合装置を備えるノードにおいて発生した割り込み要因である場合に、網結合装置に接続されたノードを障害ノードとして特定する。このため、本実施の形態例における情報処理システムは、割り込み要因に基づいて、割り込み要因に対応する障害ノードを特定することができる。

【 0 1 1 2 】

また、本実施の形態例における情報処理システムにおいて、1のノードは、割り込み要因と、割り込み要因の波及元となる割り込み要因との対応関係を有する定義テーブルを有し、1のノードの制御装置は、定義テーブルに基づいて、障害として検出すべき割り込み要因の波及元となる割り込み要因が発生しているか否かを判定する。

10

【 0 1 1 3 】

本実施の形態例における情報処理システムは、割り込み要因と、割り込み要因の波及元となる割り込み要因との対応関係を有する定義テーブルを有することによって、波及元の割り込み要因であるか否かを高速に判定することができる。また、情報処理システムは、割り込み要因が増加した場合や、変更が発生した場合に、定義テーブルの更新処理を行うことで、割り込み要因の増加や変更を簡易に適用することができる。これにより、情報処理システムは、エンハンスや設計変更時におけるメンテナンス工数を小さく抑えることができる。

20

【 0 1 1 4 】

また、本実施の形態例における情報処理システムは、割り込み要因に対応して優先度を有する定義テーブルを有し、1のノードの前記制御装置は、定義テーブルに基づいて、割り込み要因の優先度を判定する。

【 0 1 1 5 】

本実施の形態例における情報処理システムは、割り込み要因に対応して優先度を有する定義テーブルを有することによって、割り込み要因の優先度を高速に取得することができる。また、情報処理システムは、割り込み要因が増加した場合や、変更が発生した場合に、定義テーブルの更新処理を行うことで、割り込み要因の増加や変更を簡易に適用することができる。これにより、情報処理システムは、エンハンスや設計変更時におけるメンテナンス工数を小さく抑えることができる。

30

【 0 1 1 6 】

以上、各ノードが共有メモリを有する分散型共有メモリの構成を例に説明したが、本実施の形態は、各ノードが共有メモリを設けておらず、各ノードとは別に共有メモリを備えるクラスター型構成にも適用可能である。

【 0 1 1 7 】

以上の実施の形態をまとめると、次の付記のとおりである。

【 0 1 1 8 】

(付記 1)

40

複数のノード間でメモリを共有する情報処理システムにおいて、

前記ノードの各々は、

複数の機能回路と前記機能回路を制御する制御装置と、

前記複数の機能回路から発生する割り込み要因を格納するレジスタとを有し、

前記複数のノードのうちの1のノードの前記制御装置は、

他の前記ノードの割り込み要因の発生に応じて前記レジスタの前記割り込み要因を受信し、前記割り込み要因のうち、障害として検出すべき割り込み要因を抽出して、抽出結果に応じて障害ノードを特定し、前記障害ノードの前記メモリへのアクセスを抑止後、前記他のノードから受信したログ情報に基づいて前記障害ノードの切り離し制御を行う情報処理システム。

50

【 0 1 1 9 】

(付記 2)

付記 1 において、

前記他のノードの前記制御装置は、前記レジスタの前記割り込み要因の発生を前記 1 のノードの制御装置に通知し、

前記 1 のノードの制御装置は、前記他のノードからの前記通知に応じて、前記他ノードの前記レジスタの割り込み要因と前記ログ情報とを収集する情報処理システム。

【 0 1 2 0 】

(付記 3)

付記 1 または 2 において、

前記 1 のノードは、網結合装置を備え、

前記他のノードは、データ処理を実行し、前記網結合装置を介して前記メモリにアクセスする処理装置を備える情報処理システム。

【 0 1 2 1 】

(付記 4)

付記 1 乃至 3 のいずれかにおいて、

前記 1 のノードの前記制御装置は、前記障害として検出すべき割り込み要因の波及元となる割り込み要因が発生しているか否かを判定し、発生していない場合に、前記割り込み要因に対応するノードを前記障害ノードとして特定し、発生している場合に、前記波及元となる割り込み要因に対応するノードを前記障害ノードとして特定する情報処理システム。

【 0 1 2 2 】

(付記 5)

付記 1 乃至 4 のいずれかにおいて、

前記第 1 のノードの前記制御部は、前記障害として検出すべき割り込み要因を複数抽出した場合に、前記割り込み要因の優先度に基づいて、前記特定した障害ノードの前記メモリへのアクセスを抑止する情報処理システム。

【 0 1 2 3 】

(付記 6)

付記 3 において、

前記 1 のノードの前記制御装置は、前記障害として検出すべき割り込み要因が前記データ処理を実行するノードにおいて発生した割り込み要因である場合に、発生元のノードを前記障害ノードとして特定し、前記障害として検出すべき割り込み要因が前記網結合装置を備えるノードにおいて発生した割り込み要因である場合に、前記網結合装置に接続されたノードを前記障害ノードとして特定する情報処理システム。

【 0 1 2 4 】

(付記 7)

付記 4 において、

前記 1 のノードは、

前記割り込み要因と、前記割り込み要因の波及元となる割り込み要因との対応関係を有する定義テーブルを有し、

前記 1 のノードの前記制御装置は、前記定義テーブルに基づいて、前記障害として検出すべき割り込み要因の波及元となる割り込み要因が発生しているか否かを判定する情報処理システム。

【 0 1 2 5 】

(付記 8)

付記 5 において、

前記 1 のノードは、

前記割り込み要因に対応して前記優先度を有する定義テーブルを有し、

前記 1 のノードの前記制御装置は、前記定義テーブルに基づいて、割り込み要因の前記

10

20

30

40

50

優先度を判定する情報処理システム。

【 0 1 2 6 】

(付記 9)

付記 1 乃至 8 のいずれかにおいて、
前記メモリは各前記ノード内に設けられた情報処理システム。

【 0 1 2 7 】

(付記 1 0)

複数のノード間でメモリを共有する情報処理システムの障害処理方法において、
前記ノードの各々は、

複数の機能回路と前記機能回路を制御する制御装置と、

前記複数の機能回路から発生する割り込み要因を格納するレジスタとを有し、

前記複数のノードのうちの 1 のノードの前記制御装置は、

他の前記ノードの割り込み要因の発生に応じて前記レジスタの前記割り込み要因を受信し、前記割り込み要因のうち、障害として検出すべき割り込み要因を抽出して、抽出結果に応じて障害ノードを特定し、前記障害ノードの前記メモリへのアクセスを抑止後、前記他のノードから受信したログ情報に基づいて前記障害ノードの切り離し制御を行う情報処理システムの障害処理方法。

10

【 0 1 2 8 】

(付記 1 1)

付記 1 0 において、

前記他のノードの前記制御装置は、前記レジスタの前記割り込み要因の発生を前記 1 のノードの制御装置に通知し、

前記 1 のノードの制御装置は、前記他のノードからの前記通知に応じて、前記他ノードの前記レジスタの割り込み要因と前記ログ情報とを収集する情報処理システムの障害処理方法。

20

【 0 1 2 9 】

(付記 1 2)

付記 1 0 または 1 1 において、

前記 1 のノードは、網結合装置を備え、

前記他のノードは、データ処理を実行し、前記網結合装置を介して前記メモリにアクセスする処理装置を備える情報処理システムの障害処理方法。

30

【 0 1 3 0 】

(付記 1 3)

付記 1 0 乃至 1 2 のいずれかにおいて、

前記 1 のノードの前記制御装置は、前記障害として検出すべき割り込み要因の波及元となる割り込み要因が発生しているか否かを判定し、発生していない場合に、前記割り込み要因に対応するノードを前記障害ノードとして特定し、発生している場合に、前記波及元となる割り込み要因に対応するノードを前記障害ノードとして特定する情報処理システムの障害処理方法。

40

【 0 1 3 1 】

(付記 1 4)

付記 1 0 乃至 1 3 のいずれかにおいて、

前記第 1 のノードの前記制御部は、前記障害として検出すべき割り込み要因を複数抽出した場合に、前記割り込み要因の優先度に基づいて、前記特定した障害ノードの前記メモリへのアクセスを抑止する情報処理システムの障害処理方法。

【 0 1 3 2 】

(付記 1 5)

付記 1 2 において、

前記 1 のノードの前記制御装置は、前記障害として検出すべき割り込み要因が前記データ処理を実行するノードにおいて発生した割り込み要因である場合に、発生元のノードを

50

前記障害ノードとして特定し、前記障害として検出すべき割り込み要因が前記網結合装置を備えるノードにおいて発生した割り込み要因である場合に、前記網結合装置に接続されたノードを前記障害ノードとして特定する情報処理システムの障害処理方法。

【 0 1 3 3 】

(付記 1 6)

付記 1 3 において、

前記 1 のノードは、

前記割り込み要因と、前記割り込み要因の波及元となる割り込み要因との対応関係を有する定義テーブルを有し、

前記 1 のノードの前記制御装置は、前記定義テーブルに基づいて、前記障害として検出すべき割り込み要因の波及元となる割り込み要因が発生しているか否かを判定する情報処理システムの障害処理方法。

10

【 0 1 3 4 】

(付記 1 7)

付記 1 4 において、

前記 1 のノードは、

前記割り込み要因に対応して前記優先度を有する定義テーブルを有し、

前記 1 のノードの前記制御装置は、前記定義テーブルに基づいて、割り込み要因の前記優先度を判定する情報処理システムの障害処理方法。

20

【 0 1 3 5 】

(付記 1 8)

付記 1 0 乃至 1 7 のいずれかにおいて、

前記メモリは各前記ノード内に設けられた情報処理システムの障害処理方法。

【 符号の説明 】

【 0 1 3 6 】

1 A ~ 1 P : システムボード、 2 A B ~ 2 D B : クロスバススイッチボックス、
B 1 : システムボードユニット、 1 2 : C P U チップ、 1 5 : システムコントローラ、 1 6 : I / O コントローラ、 1 8 : メモリコントローラ、 1 1 : メモリ、 1 9 : M B C (システムボードユニット)、

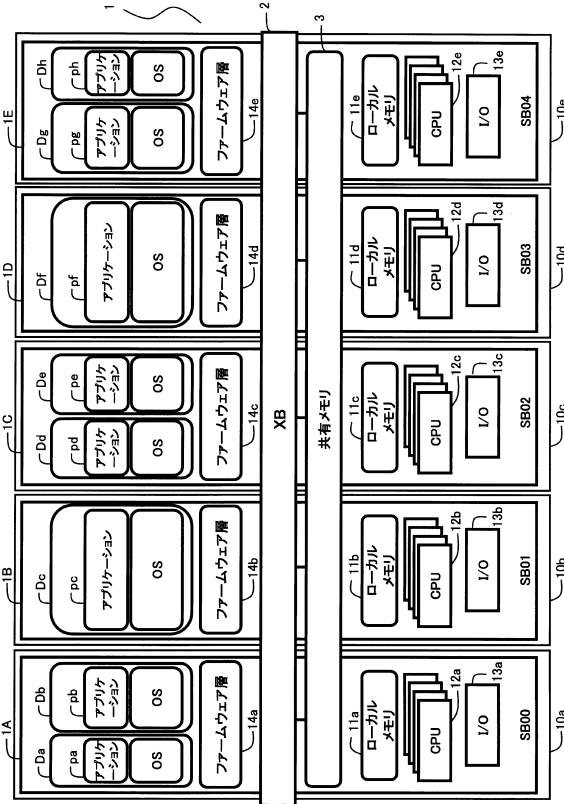
B 2 : サービスプロセッサボード、 M B C 2 1 (サービスプロセッサボードユニット)、
2 2 : システム制御装置、 r g : レジスタ、

2 A B : マスターノード、 1 V : システム制御装置、

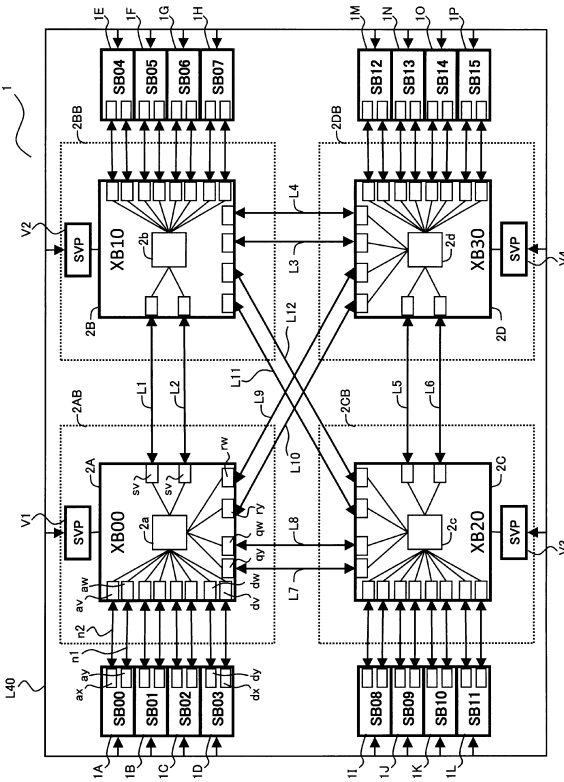
1 A : スレーブノード、 2 2 : システム制御装置

30

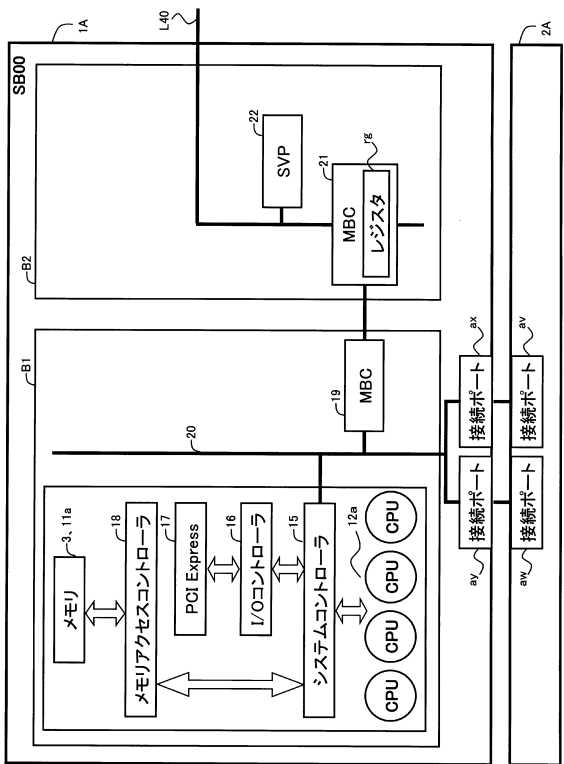
【図 1】



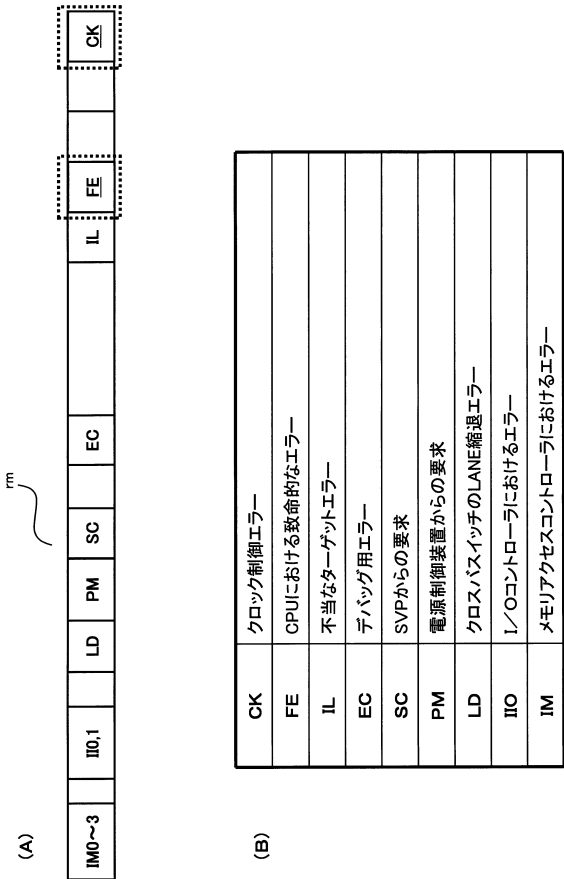
【図 2】



【図 3】

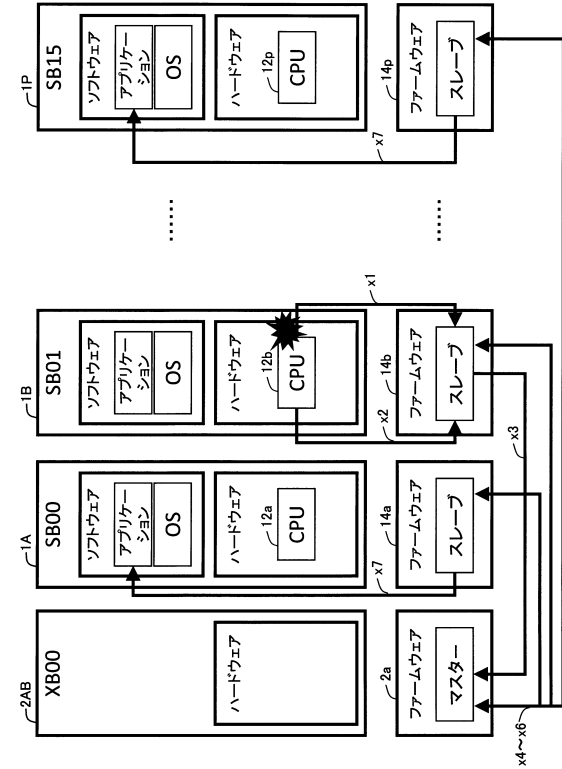


【図 4】

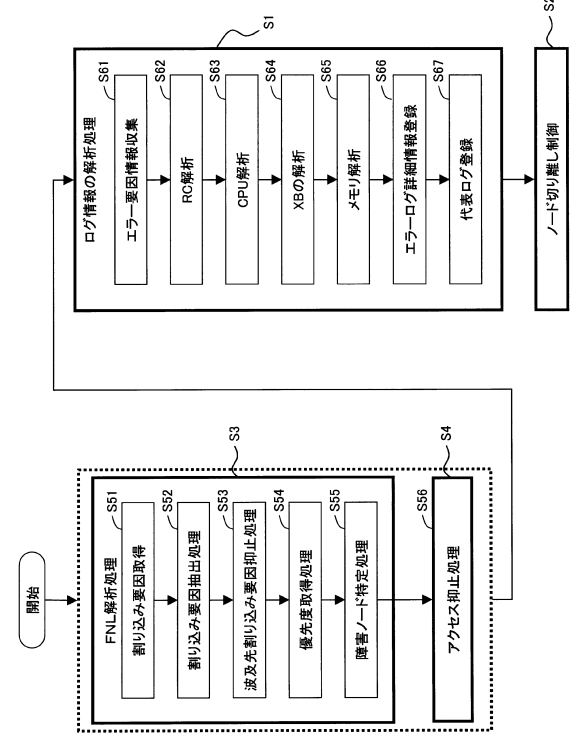


CK	クロック制御エラー
FE	CPUにおける致命的なエラー
IL	不当なターゲットエラー
EC	デバッグ用エラー
SC	SVPからの要求
PM	電源制御装置からの要求
LD	クロスバスイッチのLANE縮退エラー
IO	I/Oコントローラにおけるエラー
IM	メモリアクセスコントローラにおけるエラー

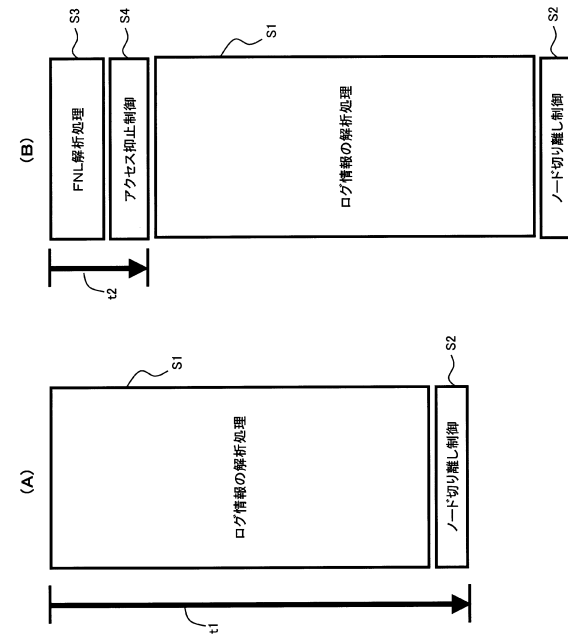
【図 5】



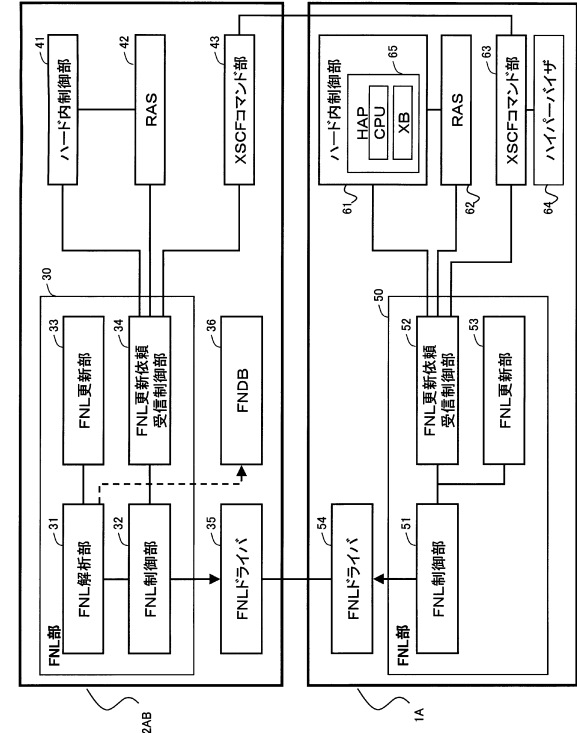
【図 6】



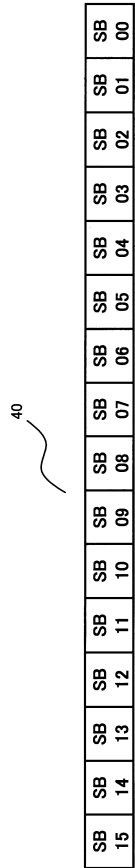
【図 7】



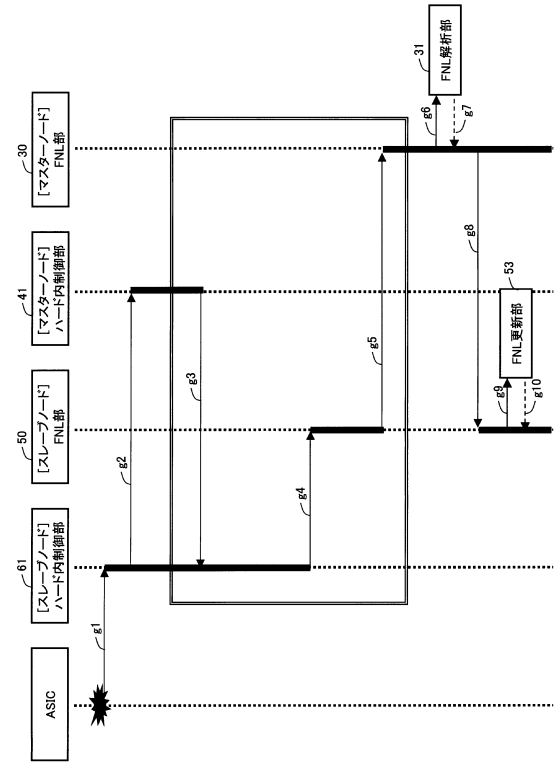
【図 8】



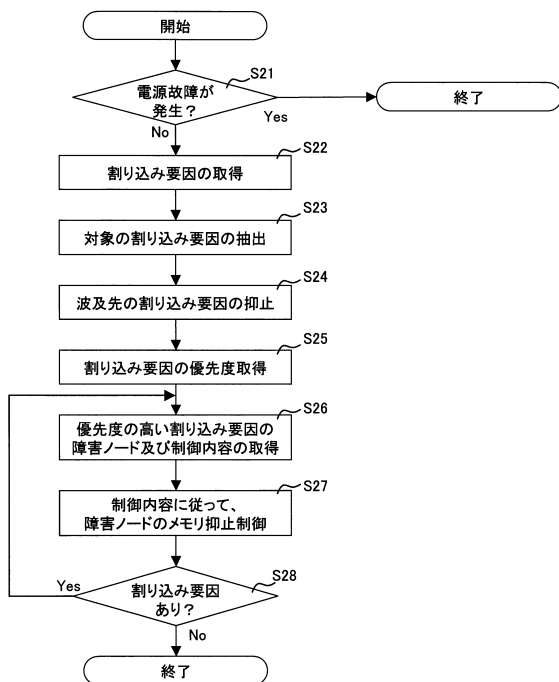
【図 9】



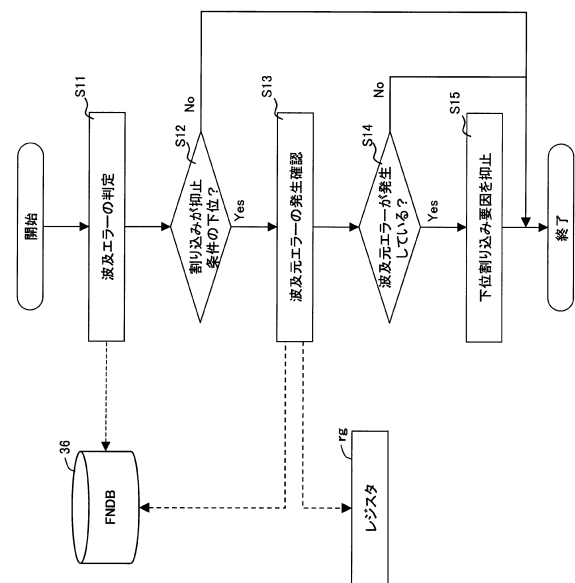
【図 10】



【図 11】



【図 12】



フロントページの続き

(56)参考文献 特開2001-134546(JP,A)
特開2013-140445(JP,A)
特開2006-178786(JP,A)
特開2004-062535(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 11/07

G06F 11/30 - 11/34