

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号  
特許第6089884号  
(P6089884)

(45) 発行日 平成29年3月8日 (2017.3.8)

(24) 登録日 平成29年2月17日 (2017.2.17)

(51) Int.Cl.

F I

G O 6 F 11/20 (2006.01)

G O 6 F 11/22 (2006.01)

G O 6 F 11/20 6 9 2

G O 6 F 11/22 6 7 3 J

請求項の数 9 (全 52 頁)

(21) 出願番号	特願2013-71904 (P2013-71904)	(73) 特許権者	000005223
(22) 出願日	平成25年3月29日 (2013.3.29)		富士通株式会社
(65) 公開番号	特開2014-197266 (P2014-197266A)		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26年10月16日 (2014.10.16)	(74) 代理人	100092978
審査請求日	平成27年12月4日 (2015.12.4)		弁理士 真田 有
		(74) 代理人	100112678
			弁理士 山本 雅久
		(72) 発明者	田村 雅寿
			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	三坂 敏夫
		最終頁に続く	

(54) 【発明の名称】 情報処理システム、情報処理装置、情報処理装置の制御プログラム、及び情報処理システムの制御方法

(57) 【特許請求の範囲】

【請求項 1】

相互に接続される複数の情報処理装置を有し、前記複数の情報処理装置間で通信を行なう情報処理システムにおいて、

前記複数の情報処理装置の各々が、

前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信する受信処理部と、

前記受信処理部が前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定する判定部と、

前記判定部が判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信する送信処理部と、を有し、

前記判定部は、

前記受信処理部が受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、前記他の情報処理装置の各々からの前記状態情報の受信状況とに基づいて、前記複数の情報処理装置の各々の状態を判定し、

前記送信処理部は、

第1所定時間ごとに、前記送信用状態情報を、前記他の情報処理装置の各々へ送信し、

前記判定部は、

前記第1所定時間以上の時間である第2所定時間内に前記状態情報を受信しなかった他

の情報処理装置の状態を、停止の可能性を示す第 1 状態と判定し、

前記受信処理部が受信した前記状態情報に基づいて、第 1 所定数以上の前記複数の情報処理装置で前記第 1 状態であると判定された情報処理装置の状態、又は、前記他の情報処理装置の少なくとも 1 つから停止を示す第 2 状態であると判定された情報処理装置の状態を、前記第 2 状態と判定することを特徴とする、情報処理システム。

【請求項 2】

前記判定部は、

前記受信処理部が受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、前記判定部が含まれる自情報処理装置の状態に関する状態情報に関する自己状態情報とに基づいて、前記複数の情報処理装置の各々の状態を判定することを特徴とする、請求項 1 記載の情報処理システム。

10

【請求項 3】

前記判定部は、

前記受信処理部が受信した前記状態情報に基づいて、前記第 1 所定数以上の数である第 2 所定数以上の前記複数の情報処理装置で前記第 2 状態であると判定された情報処理装置を、リカバリ処理中を示す第 3 状態と判定し、

前記複数の情報処理装置のうちの 1 以上の情報処理装置はさらに、

前記判定部が前記第 3 状態と判定した情報処理装置に対して、リカバリ処理を実行するリカバリ処理部を有することを特徴とする、請求項 1 又は請求項 2 記載の情報処理システム。

20

【請求項 4】

前記複数の情報処理装置が複数のグループに分割され、

前記複数のグループの各々における代表情報処理装置はさらに、

前記複数のグループのうちの自グループ以外の他のグループの各々における他の代表情報処理装置から、前記他の代表情報処理装置により判定された前記複数のグループの代表情報処理装置の各々の状態に関する代表状態情報を受信するグループ間受信処理部と、

前記グループ間受信処理部が前記他の代表情報処理装置の各々から受信した前記代表状態情報に基づいて、前記複数の代表情報処理装置の各々の状態を判定するグループ間判定部と、

前記グループ間判定部が判定した前記複数の代表情報処理装置の各々の状態に関する送信用代表状態情報を、前記他の代表情報処理装置の各々へ送信するグループ間送信処理部と、を有し、

30

前記複数の情報処理装置の各々において、

前記送信処理部は、

前記送信用状態情報を、前記自グループにおける他の情報処理装置の各々へ送信し、

前記判定部は、

前記受信処理部が前記自グループにおける他の情報処理装置の各々から受信した前記状態情報に基づいて、前記自グループにおける情報処理装置の各々の状態を判定することを特徴とする、請求項 1 ~ 3 のいずれか 1 項記載の情報処理システム。

【請求項 5】

40

前記複数のグループの各々における代表情報処理装置はさらに、

前記自グループにおける情報処理装置の数が第 4 所定値を超えた場合、前記自グループから、複数の情報処理装置を分割して新たなグループを作成する管理部を有することを特徴とする、請求項 4 記載の情報処理システム。

【請求項 6】

前記管理部は、

前記自グループにおける情報処理装置の数が第 5 所定値未満の場合、前記自グループと前記他のグループのうちのいずれかのグループとを統合することを特徴とする、請求項 5 記載の情報処理システム。

【請求項 7】

50

相互に接続される複数の情報処理装置の各々において、

前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信する受信処理部と、

前記受信処理部が前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定する判定部と、

前記判定部が判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信する送信処理部と、を有し、

前記判定部は、

前記受信処理部が受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、前記他の情報処理装置の各々からの前記状態情報の受信状況とに基づいて、前記複数の情報処理装置の各々の状態を判定し、

前記送信処理部は、

第1所定時間ごとに、前記送信用状態情報を、前記他の情報処理装置の各々へ送信し、

前記判定部は、

前記第1所定時間以上の時間である第2所定時間内に前記状態情報を受信しなかった他の情報処理装置の状態を、停止の可能性を示す第1状態と判定し、

前記受信処理部が受信した前記状態情報に基づいて、第1所定数以上の前記複数の情報処理装置で前記第1状態であると判定された情報処理装置の状態、又は、前記他の情報処理装置の少なくとも1つから停止を示す第2状態であると判定された情報処理装置の状態を、前記第2状態と判定することを特徴とする、情報処理装置。

#### 【請求項8】

相互に接続される複数の情報処理装置の各々を制御する情報処理装置に、

前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信し、

前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定し、

判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信する、処理を実行させ、

前記判定において、

受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、前記他の情報処理装置の各々からの前記状態情報の受信状況とに基づいて、前記複数の情報処理装置の各々の状態を判定し、

前記送信において、

第1所定時間ごとに、前記送信用状態情報を、前記他の情報処理装置の各々へ送信し、

前記判定において、

前記第1所定時間以上の時間である第2所定時間内に前記状態情報を受信しなかった他の情報処理装置の状態を、停止の可能性を示す第1状態と判定し、

受信した前記状態情報に基づいて、第1所定数以上の前記複数の情報処理装置で前記第1状態であると判定された情報処理装置の状態、又は、前記他の情報処理装置の少なくとも1つから停止を示す第2状態であると判定された情報処理装置の状態を、前記第2状態と判定することを特徴とする、情報処理装置の制御プログラム。

#### 【請求項9】

相互に接続される複数の情報処理装置を有し、前記複数の情報処理装置間で通信を行なう情報処理システムの制御方法において、

前記複数の情報処理装置の各々が、

前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信し、

前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定し、

判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信し、

前記判定において、

受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、前記他の情報処理装置の各々からの前記状態情報の受信状況とに基づいて、前記複数の情報処理装置の各々の状態を判定し、

前記送信において、

第1所定時間ごとに、前記送信用状態情報を、前記他の情報処理装置の各々へ送信し、

前記判定において、

前記第1所定時間以上の時間である第2所定時間内に前記状態情報を受信しなかった他の情報処理装置の状態を、停止の可能性を示す第1状態と判定し、

受信した前記状態情報に基づいて、第1所定数以上の前記複数の情報処理装置で前記第1状態であると判定された情報処理装置の状態、又は、前記他の情報処理装置の少なくとも1つから停止を示す第2状態であると判定された情報処理装置の状態を、前記第2状態と判定することを特徴とする、情報処理システムの制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本件は、情報処理システム、情報処理装置、情報処理装置の制御プログラム、及び情報処理システムの制御方法に関する。

【背景技術】

【0002】

従来、複数のノード（ストレージ装置、情報処理装置）をそなえ、データを複数のノードに分散させて保持する分散ストレージシステム（ストレージシステム、情報処理システム）が知られている。

分散ストレージシステムにおいて、例えば、複数のノードのいずれかのノードに故障が発生した場合、分散ストレージシステムを使用するクライアントは、故障したノードへアクセスをすることができなくなる。

【0003】

また、故障したノードが他のノードと冗長化されていた場合には、クライアントは、故障したノードの代わりに冗長化されたノードへアクセスをすることができる。しかし、冗長化されたノードをそなえる分散ストレージシステムは、ノードに故障が発生する前のデータの多重化状態を回復するリカバリ処理及び故障ノードの交換が行なわれるまで、冗長度が低下した信頼性の低い状態になる。

【0004】

従って、分散ストレージシステムは、複数のノードの状態を監視し、ノードの故障を速やかに検出することが好ましい。

しかし、分散ストレージシステムでは、ノード又はノード間のリンクの故障により、複数のノードが分断され、分断された一方のノードと他方のノードとが、ノードの故障について異なる判断をすることがある。この状態をスプリットブレイン（Split Brain）状態という。スプリットブレイン状態の一例としては、一方のノードと他方のノードとの間でリンクの故障が発生することにより互いにアクセスができなくなる状態が発生するが、双方のノードは、互いにアクセスができなくなった相手のノードが故障したと判断する場合が挙げられる。

【0005】

例えば、一方及び他方のノードが、同一データの冗長データを互いに保持する場合に、スプリットブレイン状態に陥ると、双方のノードは、それぞれの保持する冗長データを個別に更新したり、他のノードへリカバリをし、冗長データの一貫性を崩す可能性がある。

10

20

30

40

50

分散ストレージシステムにおいて、スプリットブレイン状態に陥ることを防止する手法としては、以下に例示する手法が知られている。

(1) 複数のノードの各々が、複数のノードのうちの所定のノード(コントロールノード)へ自ノードの構成情報及び生存報告を通知する。コントロールノードは、複数のノードの各々から得た情報を集約して複数のノードを監視し、監視結果から故障したノードを検出すると、リカバリを行ない、管理者等へノードの故障を通知する。

(2) 複数のノードの各々が、互いに生存報告のやり取りを行ない(情報交換フェーズ)、どのノードが監視及び故障ノードの検出を行なうかを、他のノードとの間で合意を取ることで選定する。合意を得たノード(決定ノード)は、複数のノードの各々の状態を監視し、監視結果から故障したノードを検出すると、リカバリを行ない、管理者等へノードの故障を通知する。

10

(3) 複数のノードの各々が、所定のノードへ生存報告を行なう。故障ノードは所定のノードにより即座に検出はされず、管理者等が、所定のノードを参照し手動で故障ノードの検出及びリカバリ等の対応を行なう。

#### 【0006】

上記(1)の手法では、コントロールノードが故障ノードの検出を行ない、上記(2)の手法では、合意を得た決定ノードが故障ノードの検出を行なう。また、上記(3)の手法では、管理者等が故障ノードの検出を行なう。従って、上記(1)~(3)の手法によれば、複数のノードで判断が行なわれるのではなく、特定のノード又は管理者が判断を行なうため、スプリットブレイン状態に陥ることを防止できる。

20

#### 【0007】

なお、関連する技術として、分散ストレージシステムにおいて、データの消失を防ぐため、コンピュータが、複数のストレージノードから収集した属性に基づき、ストレージノードを2以上のグループに分ける技術が知られている(例えば、特許文献1参照)。この技術では、コンピュータは、作成した各グループ内において、データを分散した分散データと、当該データと同一内容の冗長データを分散した冗長分散データとが存在しないように、分散データ及び冗長分散データを各グループに割り当てる。

#### 【0008】

また、関連する他の技術として、管理サーバが、データを保持する全てのストレージで同一のデータ・プールを構成し、異なるデータをできるだけプール内の複数の異なるストレージに分散して保持させる技術が知られている(例えば、特許文献2参照)。

30

さらに、関連する他の技術として、ネットワーク監視装置が、複数ノードをグループ単位に分割し、分割したグループの1つのノードから論理回線状態を取得して、論理回線の監視を行なう技術が知られている(例えば、特許文献3参照)。

#### 【0009】

また、関連する他の技術として、ネットワーク管理システムが、ノードの自装置の情報、ホップ数等の情報に基づいて形成されたグループ毎に、グループ内のノードを監視するグループ管理装置をそなえる技術が知られている(例えば、特許文献4参照)。

#### 【先行技術文献】

#### 【特許文献】

40

#### 【0010】

【特許文献1】国際公開第WO2008/114441号パンフレット

【特許文献2】特表2011-505617号公報

【特許文献3】特開2010-258614号公報

【特許文献4】特開2011-055231号公報

#### 【発明の概要】

#### 【発明が解決しようとする課題】

#### 【0011】

上記(1)の手法では、複数のノードの各々の情報が1点(コントロールノード)に集約されるため、コントロールノードがSPOF(Single Point Of Failure; 単一障害点

50

）となる。従って、コントロールノードが故障した場合、クライアントは、コントロールノードが復旧するまで分散ストレージシステムの利用が制限されるという課題がある。

上記（２）の手法では、複数のノード間で合意を形成するために複雑な手順が行なわれるため、上記（１）の手法と比較して、合意を形成するまでの時間が余計にかかる場合がある。また上記（３）の手法では、管理者等による人為的な判断が行なわれるため、ノードの故障が発生してからノードの故障が検出され、リカバリ処理が行なわれるまでに、上記（１）及び（２）の手法と比較して時間がかかる場合がある。つまり、上記（２）及び（３）の手法では、障害が発生したノードに対するリカバリ処理等の開始が遅くなり、クライアントが分散ストレージシステムの利用を制限される時間が長くなるという課題がある。

10

#### 【００１２】

なお、上述した関連する技術は、いずれも、上記（１）の手法のように管理装置が複数のノードを管理するものであり、上述した課題については考慮されていない。

このように、複数のストレージ装置をそなえるストレージシステムにおいて、複数のストレージ装置の各々の状態を判断する上述した技術では、ストレージシステムの可用性が低下するという課題がある。

#### 【００１３】

ここまで、情報処理システムがストレージシステム（分散ストレージシステム）であるものとして説明したが、これに限定されるものではない。上述した課題は、情報処理システムがそなえる複数の情報処理装置の各々が、分散データではなく他の情報処理装置とは異なるデータを保持する場合であっても、同様に生じ得る。

20

１つの側面では、本発明は、複数の情報処理装置をそなえる情報処理システムにおいて、可用性の低下を抑止することを目的とする。

#### 【００１４】

なお、前記目的に限らず、後述する発明を実施するための形態に示す各構成により導かれる作用効果であって、従来の技術によっては得られない作用効果を奏することも本発明の他の目的の１つとして位置付けることができる。

#### 【課題を解決するための手段】

#### 【００１５】

本件の情報処理システムは、相互に接続される複数の情報処理装置を有し、前記複数の情報処理装置間で通信を行なう情報処理システムにおいて、前記複数の情報処理装置の各々が、前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信する受信処理部と、前記受信処理部が前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定する判定部と、前記判定部が判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信する送信処理部と、をそなえてよい。また、前記判定部は、前記受信処理部が受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、前記他の情報処理装置の各々からの前記状態情報の受信状況とに基づいて、前記複数の情報処理装置の各々の状態を判定し、前記送信処理部は、第１所定時間ごとに、前記送信用状態情報を、前記他の情報処理装置の各々へ送信してよい。さらに、前記判定部は、前記第１所定時間以上の時間である第２所定時間内に前記状態情報を受信しなかった他の情報処理装置の状態を、停止の可能性を示す第１状態と判定し、前記受信処理部が受信した前記状態情報に基づいて、第１所定数以上の前記複数の情報処理装置で前記第１状態であると判定された情報処理装置の状態、又は、前記他の情報処理装置の少なくとも１つから停止を示す第２状態であると判定された情報処理装置の状態を、前記第２状態と判定してよい。

30

40

#### 【発明の効果】

#### 【００１６】

第１実施形態及び第２実施形態によれば、複数の情報処理装置をそなえる情報処理シス

50

テムにおいて、可用性の低下を抑止することができる。

【図面の簡単な説明】

【0017】

【図1】第1実施形態の一例としてのストレージシステムの構成例を示す図である。

【図2】第1実施形態の一例としてのノードのハードウェア構成例を示す図である。

【図3】第1実施形態の一例としてのノードの機能構成例を示す図である。

【図4】第1実施形態の一例としてのノードが送受信するノード状態情報を例示する図である。

【図5】第1実施形態の一例としてのノードが管理するノード状態管理情報を例示する図である。

10

【図6】第1実施形態の一例としての新規ノードが送信する情報を例示する図である。

【図7】第1実施形態の一例としての新規ノードが受信する情報を例示する図である。

【図8】第1実施形態の一例としてのノードが他ノードの状態を判定するときの状態遷移の一例を示す図である。

【図9】第1実施形態の一例としての複数のノードによるノード状態情報の送受信処理の一例を説明する図である。

【図10】第1実施形態の一例としてのノードが自ノードの状態を判定するときの状態遷移の一例を示す図である。

【図11】第1実施形態の一例としての新規ノードによる起動後の動作例を説明するフローチャートである。

20

【図12】第1実施形態の一例としてのノードによる他ノードの状態を判定する動作例を説明するフローチャートである。

【図13】第1実施形態の一例としてのノードによる自ノードの状態を判定する動作例を説明するフローチャートである。

【図14】第2実施形態の一例としてのノードの機能構成例を示す図である。

【図15】第2実施形態の一例としてのノードが管理するパーティ管理情報を例示する図である。

【図16】第2実施形態の一例としての複数のノードによる代表ノード状態情報及びノード状態情報の送受信処理の一例を説明する図である。

【図17】第2実施形態の一例としてのノードが送受信するノード状態情報を例示する図である。

30

【図18】第2実施形態の一例としてのノードが送受信する代表ノード状態情報を例示する図である。

【図19】第2実施形態の一例としてのノードが管理するノード状態管理情報を例示する図である。

【図20】第2実施形態の一例としてのストレージシステムにノードが追加される例を示す図である。

【図21】図20に示すストレージシステムにおけるパーティの分割処理の一例を説明する図である。

【図22】図21に示すストレージシステムにおけるノードの削除処理及びパーティの統合処理の一例を説明する図である。

40

【図23】第2実施形態の一例としてのストレージシステムにおけるパーティの分割処理の具体例を説明する図である。

【図24】第2実施形態の一例としての代表ノードによる他の代表ノードの状態を判定する動作例を説明するフローチャートである。

【図25】第2実施形態の一例としてのノードによるパーティ内の他ノードが停止した場合の動作例を説明するフローチャートである。

【図26】第2実施形態の一例としてのノードによるパーティの分割処理及び統合処理の動作例を説明するフローチャートである。

【発明を実施するための形態】

50

## 【 0 0 1 8 】

以下、図面を参照して実施の形態を説明する。

## 〔 1 〕 第 1 実施形態

## 〔 1 - 1 〕 ストレージシステムの構成

以下、図 1 及び図 2 を参照して、第 1 実施形態の一例としてのストレージシステム 1 の構成について説明する。

## 【 0 0 1 9 】

図 1 は、第 1 実施形態の一例としてのストレージシステム 1 の構成例を示す図であり、図 2 は、図 1 に示すノード 1 0 - 1 ~ 1 0 - 5 のハードウェア構成例を示す図である。

図 1 に示すように、第 1 実施形態に係るストレージシステム（情報処理システム）1 は、複数（例えば 5 つ）のノード 1 0 - 1 ~ 1 0 - 5 及び複数（例えば 3 つ）のスイッチ 2 0 - 1 ~ 2 0 - 3 をそなえる。

## 【 0 0 2 0 】

なお、以下、ノード 1 0 - 1 ~ 1 0 - 5 を区別しない場合には、単にノード 1 0 といい、スイッチ 2 0 - 1 ~ 2 0 - 3 を区別しない場合には、単にスイッチ 2 0 という。

ストレージシステム 1 は、複数のノード 1 0 及びスイッチ 2 0 により、S A N（Storage Area Network）を形成し、相互に接続される複数のノード 1 0 間で通信を行なう。また、ストレージシステム 1 は、図示しないクライアントに接続され、クライアントに対してノード 1 0 が有する記憶領域（リソース）を提供する。

## 【 0 0 2 1 】

ストレージシステム 1 としては、分散ストレージシステム又はクラスタファイルシステム等の、データを複数のノード 1 0 に分散させて保持する種々のストレージシステムが例として挙げられる。例えば、ストレージシステム 1 は、W e b サーバのデータベースやクラウドストレージ等に用いられることがある。

なお、複数のノード 1 0 の各々が、分散データではなく他のノード 1 0 とは異なるデータを保持してもよい。

## 【 0 0 2 2 】

ノード（ストレージ装置、ノード装置、情報処理装置）1 0 は、クライアント（端末装置、図示省略）からの各種要求に応じて、ノード 1 0 がそなえる記憶部 1 0 c（図 2 参照）に対する各種処理を行なう。なお、ノード 1 0 としては P C（Personal Computer）サーバ等の情報処理装置が挙げられる。

ノード 1 0 は、図 2 に示すように、C P U（Central Processing Unit）1 0 a、メモリ 1 0 b、記憶部 1 0 c、ネットワークインタフェース 1 0 d、入出力部 1 0 e、記録媒体 1 0 f、及び読取部 1 0 g をそなえる。なお、ノード 1 0 - 1 ~ 1 0 - 5 は、互いに同様のハードウェアをそなえることができるため、以下、任意のノード 1 0 がそなえるハードウェアについて説明する。

## 【 0 0 2 3 】

C P U 1 0 a は、メモリ 1 0 b、記憶部 1 0 c、ネットワークインタフェース 1 0 d、入出力部 1 0 e、記録媒体 1 0 f、及び読取部 1 0 g と接続され、種々の制御や演算を行なう演算処理装置（プロセッサ）である。C P U 1 0 a は、メモリ 1 0 b、記憶部 1 0 c、記録媒体 1 0 f、読取部 1 0 g に接続又は挿入された記録媒体 1 0 h、又は図示しない R O M（Read Only Memory）等に格納されたプログラムを実行することにより、ノード 1 0 における種々の機能を実現する。なお、C P U 1 0 a に限らず、プロセッサとしては、M P U（Micro Processing Unit）等の電子回路が用いられてもよい。

## 【 0 0 2 4 】

メモリ 1 0 b は、種々のデータやプログラムを格納する記憶装置である。C P U 1 0 a は、プログラムを実行する際に、メモリ 1 0 b にデータやプログラムを格納し展開する。なお、メモリ 1 0 b としては、例えば R A M（Random Access Memory）等の揮発性メモリが挙げられる。

記憶部 1 0 c は、例えば H D D（Hard Disk Drive）等の磁気ディスク装置、S S D（S

10

20

30

40

50



olid State Drive)等の半導体ドライブ装置、又はフラッシュメモリ等の不揮発性メモリ等の、種々のデータやプログラム等を格納する1以上のハードウェアである。記憶部10cが有する記憶領域は、クライアントにより用いられる。

#### 【0025】

ネットワークインタフェース10dは、スイッチ20を介したノード10又はクライアントとの間の接続及び通信の制御を行なうコントローラである。ネットワークインタフェース10dとしては、例えば、LAN(Local Area Network)、ファイバチャネル(Fibre Channel; FC)、又はインフィニバンド(InfiniBand)(登録商標)等に準拠したインタフェースカードが挙げられる。なお、ネットワークインタフェース10dは、LANに準拠する場合、iSCSI(Internet Small Computer System Interface)に対応する

10

#### 【0026】

入出力部10eは、例えばマウスやキーボード等の入力装置及びディスプレイやプリンタ等の出力装置の少なくとも一方を含んでよい。例えば入出力部10eは、ストレージシステム1の管理者等により、後述するノード情報の設定又は参照、ログの参照、その他種々の作業に用いられる。

記録媒体10fは、フラッシュメモリやROM等の記憶装置であり、種々のデータやプログラムを記録する。読取部10gは、光ディスクやUSB(Universal Serial Bus)メモリ等のコンピュータ読取可能な記録媒体10hに記録されたデータやプログラムを読み出す装置である。

20

#### 【0027】

記録媒体10f及び10hの少なくとも一方には、第1実施形態に係るノード10(及び後述する第2実施形態に係るノード10A)の機能を実現する制御プログラムが格納されてもよい。すなわち、CPU10aは、記録媒体10f、又は読取部10gを介して記録媒体10hから出力された制御プログラムを、メモリ10b等の記憶装置に展開して実行することにより、ノード10の機能を実現する。

#### 【0028】

なお、上述した各ハードウェアは、互いにバスを介して通信可能に接続される。例えば、CPU10a、メモリ10b、及びネットワークインタフェース10dは、システムバスに接続される。また、例えば、記憶部10c、入出力部10e、記録媒体10f、及び読取部10gは、I/O(Input/Output)インタフェース等を介してシステムバスに接続される。なお、記憶部10cは、SCSI(Serial Attached SCSI)、ファイバチャネル、SATA(Serial Advanced Technology Attachment)等に準拠したバス(ケーブル)で、DI(Disk Interface)等のI/Oインタフェースに接続される。

30

#### 【0029】

なお、ノード10の上述したハードウェア構成は例示である。従って、ノード10内のハードウェアの増減や分割等は適宜行なわれてもよい。

スイッチ(接続装置)20は、複数のノード10間又は他のスイッチ20間に接続され、スイッチ20に接続されたノード10間でやり取りされるコマンド又はデータ等の情報を中継する。スイッチ20としては、例えばL2(Layer 2)スイッチ、FCスイッチ等のハードウェアスイッチが挙げられる。

40

#### 【0030】

図1に例示するストレージシステム1では、スイッチ20-1は、スイッチ20-2及び20-3に接続される。また、スイッチ20-2はスイッチ20-1及びノード10-1及び10-2に、スイッチ20-3はスイッチ20-1及びノード10-3~10-5に、それぞれ接続される。なお、スイッチ20は、図1に示すものに限られず、ノード10の数等に応じて、多段に接続されてもよいし、1つのスイッチ20が用いられてもよい。

#### 【0031】

なお、クライアントがインターネット又はイントラネット等のネットワークを介してス

50

ストレージシステム 1 に接続される場合、スイッチ 20 とクライアントとの間にルータが介設されてもよい。ルータとしては、例えば、ソフトウェアルータの他、L3 スwitch 等のハードウェアルータ等が挙げられる。

#### 〔 1 - 2 〕 ノードの説明

第 1 実施形態の一例としてのストレージシステム 1 は、上述のように、相互に接続される複数のノード 10 を有し、複数のノード 10 間で通信を行なう。

#### 【 0032 】

具体的には、第 1 実施形態の一例としてのノード 10 の各々は、以下の ( a ) ~ ( c ) の処理を行なう。

( a ) 複数のノード 10 のうちの自ノード 10 以外の他ノード 10 の各々から、他ノード 10 により判定された複数のノード 10 の各々の状態に関するノード状態情報 T1 ( 図 4 参照 ) を受信する。

10

#### 【 0033 】

( b ) 他ノード 10 の各々から受信したノード状態情報 T1 に基づいて、複数のノード 10 の各々の状態を判定する。

( c ) 判定した複数のノード 10 の各々の状態に関するノード状態情報 T1 を、他ノード 10 の各々へ送信する。

なお、ノード 10 の状態とは、ノード 10 が正常に動作しているか否かを示す種別であり、詳細は後述する。

#### 【 0034 】

20

ノード 10 は、上記 ( a ) ~ ( c ) の処理を繰り返す。つまり、ノード 10 の各々は、自ノード 10 が判定 ( 生成 ) したノード状態情報 T1 を、自ノードが正常に稼働していることを示すハートビートとして定期的に他ノード 10 の各々へ送信する。そして、ノード 10 の各々は、他ノード 10 からハートビートとして送信されてきたノード状態情報 T1 を受信し、自ノード 10 が保持する管理情報を更新する。これにより、ノード 10 は、ストレージシステム 1 内の複数のノード 10 間で各ノード 10 の状態を共有し、他ノード 10 からのノード状態情報 T1 に基づいて自律的に複数のノード 10 の各々の状態を判定することができる。

#### 【 0035 】

なお、ストレージシステム 1 において、複数のノード 10 の接続形態は、図 1 に例示したものに限定されないが、複数のノード 10 間が離間するほど、ノード状態情報 T1 の送受信においてレイテンシ又はパケットロス等が生じ得る。従って、ストレージシステム 1 において、ノード状態情報 T1 を送受信する複数のノード 10 は、ネットワークの品質が互いに均一となることが望ましい。

30

#### 【 0036 】

#### 〔 1 - 3 〕 ノードの構成

次に、図 3 ~ 図 10 を参照して、第 1 実施形態の一例としてのノード 10 の構成について説明する。

図 3 は、第 1 実施形態の一例としてのノード 10 の機能構成例を示す図である。図 4 は、ノード 10 が送受信するノード状態情報 T1 を例示する図であり、図 5 は、ノード 10 ( 特にノード 10 - 1 ) が管理するノード状態管理情報 T2 を例示する図である。

40

#### 【 0037 】

ノード 10 は、上述した処理を行なうため、図 3 に例示するように、ノード状態保持部 11、受信処理部 12、ノード状態決定部 13、送信処理部 14、リカバリ処理部 15、及び停止処理部 16 をそなえる。なお、ノード 10 - 1 ~ 10 - 5 は、互いに同様の機能をそなえることができるため、以下、任意のノード 10 がそなえる機能について説明する。

#### 【 0038 】

#### 〔 1 - 3 - 1 〕 ノード状態保持部

ノード状態保持部 11 は、図 5 に示すノード状態管理情報 T2 を保持する記憶領域であ

50

り、例えば上述したメモリ 10b により実現される。

#### 〔 1 - 3 - 2 〕 受信処理部

受信処理部 12 は、上記 (a) の処理を行なう。具体的には、受信処理部 12 は、複数のノード 10 のうちの自ノード 10 以外の他ノード 10 の各々から、図 4 に例示するノード状態情報 T1 を受信し、ノード状態保持部 11 が保持するノード状態管理情報 T2 (図 5 参照) を更新する。

#### 【 0039 】

ノード状態情報 (状態情報) T1 は、送信元のノード 10 で判定された各ノード 10 の状態を含む情報である。例えば、自ノード 10 が送信するノード状態情報 T1 には、自ノード 10 が判定した各ノード 10 の状態が含まれ、自ノード 10 が受信するノード状態情報 T1 には、受信するノード状態情報の送信元である他ノード 10 で判定された各ノード 10 の状態が含まれる。なお、ノード 10 は、図 4 に示すようにノード状態情報 T1 をテーブルとして生成し、送受信することができる。

#### 【 0040 】

図 4 に示すように、ノード状態情報 T1 は、ノード 10 の識別情報の一例であるノード ID、ノード 10 ごとの状態、ノード 10 のアドレスの一例である IP (Internet Protocol) アドレス、及びノード 10 のポート番号を含む。図 4 に示すノード状態情報 T1 は、ノード 10 - 1 ~ 10 - 5 に対応するノード ID “ 1 ” ~ “ 5 ” の状態を含む。

一例として、ノード ID “ 1 ” には、状態 “ Alive ”、IP アドレス “ 192.168.0.1 ”、ポート番号 “ 12345 ” が対応付けられる。

#### 【 0041 】

なお、ノード 10 の識別情報として、ノード ID を例に挙げたが、これに限定されるものではない。識別情報は、各ノード 10 を特定できるユニークな情報であればよい。例えば、識別情報として、ノード 10 の IP アドレス、シリアル番号、又はネットワークインタフェース 10d の MAC (Media Access Control) アドレス等が用いられてもよい。

また、ノード 10 のアドレスとして、IP アドレスを例に挙げたが、これに限定されるものではない。アドレスは、IP 以外のプロトコルにおいてノード 10 を特定可能な種々のアドレスが用いられてもよい。

#### 【 0042 】

ノード状態管理情報 T2 は、自ノード 10 及び他ノード 10 で判定された複数のノード 10 の各々の状態を管理する情報である。例えば、ノード状態管理情報 T2 は、自ノード 10 が各ノード 10 の状態をどう判断しているか、他ノード 10 が各ノード 10 をどう判断しているか、及び最後に各ノード 10 からハートビートとしてのノード状態情報 T1 を受信したのはいつかといった情報を含む。なお、ノード 10 は、図 5 に示すようにノード状態管理情報 T2 をテーブルとして生成し、管理することができる。

#### 【 0043 】

以下、図 5 の説明においては、自ノード 10 がノード 10 - 1 であるものとする。

図 5 に示すように、ノード状態管理情報 T2 は、図 4 に示すノード状態情報 T1 と同様に、ノード 10 の識別情報の一例としてのノード ID、ノード 10 ごとの状態、ノード 10 のアドレスの一例としての IP アドレス、及びノード 10 のポート番号を含む。また、ノード状態管理情報 T2 はさらに、他のノード 10 から受信したノード状態情報 T1 に含まれるノード 10 ごとの状態 (図 5 中、“ by 2 ” ~ “ by 5 ” と表記)、及び他のノード 10 ごとの最終更新情報を含む。図 5 に示すノード状態管理情報 T2 は、ノード 10 - 1 ~ 10 - 5 に対応するノード ID “ 1 ” ~ “ 5 ” の状態を含む。

#### 【 0044 】

一例として、ノード ID “ 1 ” には、自ノード 10 - 1 が判定した状態 “ Alive ”、他ノード 10 - 2 ~ 10 - 5 がそれぞれ判定した状態 “ Alive ”、最終更新情報 “ 1 sec ago ” (1 秒前)、IP アドレス “ 192.168.0.1 ”、ポート番号 “ 12345 ” が対応付けられる。つまり、ノード状態管理情報 T2 には、受信処理部 12 が受信したノード状態情報 T1 が示す複数のノード 10 の各々の状態が含まれる。また、ノード状態管理情報 T2 には、ノード状態決定

10

20

30

40

50

部 1 3 が含まれる自ノード 1 0 の状態に関するノード状態情報 T 1 に関する自己状態情報が含まれる。

【 0 0 4 5 】

受信処理部 12 は、他ノード 10 の各々から上述したノード状態情報 T1 を受信すると、受信したノード状態情報 T1 に含まれるノード 10 ごとの状態を、ノード状態管理情報 T2 における対応する他ノード 10 の列に設定する。つまり、図 5 に例示する “by 2” ~ “by 5” (自ノード 10 がノード 10 - 1 の場合) の状態は、対応する他ノード 10 からの情報に基づき設定される。なお、ノード ID “4” の状態の説明は後述する。

【 0 0 4 6 】

また、ノード 10 - 1 の受信処理部 12 は、ノード 10 - 2 からノード状態情報 T1 を受信すると、ノード状態管理情報 T2 において、ノード状態情報 T1 に含まれるノード 10 ごとの状態を、“by 2”の列に設定する。また、受信処理部 12 は、ノード 10 - 2 に対応するノード ID “2”の最終更新情報を更新する。

なお、最終更新情報は、最後にハートビートを受信したのがいつであることを示す情報であり、図5に示す例では、最終更新情報として、現在時刻と最後に受信を行なった時刻（最終受信時刻）との差を示しているが、これに限定されるものではない。例えば、ノード10は、最終更新情報に最終受信時刻そのものを設定することで、最終更新情報を更新してもよい。また、ノード10は、ノード10ごとに、時間の経過に応じて値が変化（例えば増加）するタイマを実行し、ノード状態管理情報T2の最終更新情報では、対応するタイマ値を参照してもよい。最終更新情報にタイマ値が用いられる場合、ノード10は、最終更新情報の更新の際に、タイマのカウント値をリセットすることで、最終更新情報を更新することができる。

【 0 0 4 7 】

受信処理部１２は、他ノード１０からノード状態情報Ｔ１を受信した都度、受信したノード状態情報Ｔ１に基づきノード状態管理情報Ｔ２を更新してもよい。また、受信処理部１２は、受信したノード状態情報Ｔ１を送信元のノード１０の識別情報と対応付けてメモリ１０ｂ等に保持しておき、後述する第１所定時間ごとに、メモリ１０ｂが保持するノード状態情報Ｔ１に基づきノード状態管理情報Ｔ２を更新してもよい。

【 0 0 4 8 】

また、受信処理部 12 は、上述したノード状態情報 T1 の受信に加え、ノード 10 の IP アドレス及びポート番号を受信することができる。

図 6 は、第 1 実施形態の一例としての新規ノード 10 が送信する情報を例示する図であり、図 7 は、新規ノード 10 が受信する情報を例示する図である。

ノード１０（後述する送信処理部１４）は、起動後、つまりストレージシステム１に追加されると、自ノード１０のＩＰアドレス及びポート番号を含む情報を全てのノード１０へ通知する。具体的には、ストレージシステム１に追加されたノード（新規ノード）１０は、図６に例示する送信情報Ｔ３を、ストレージシステム１内の全てのノード１０へブロードキャスト等により通知する。

【 0 0 4 9 】

図 6 に示すように、新規ノード 10 が送信する送信情報 T3 は、新規ノード 10 の識別情報の一例としてのノード ID、新規ノード 10 の状態、新規ノード 10 のアドレスの一例としての IP アドレス、及び新規ノード 10 のポート番号を含む。例えば、図 6 に示す送信情報 T3 は、新規ノード 10 に対応するノード ID “6” の状態を含む。

一例として、ノード ID “ 6 ” には、新規ノード 10 が判定した状態 “ Alive ”、IP アドレス “ 192.168.0.6 ”、ポート番号 “ 12345 ” が対応付けられる。

【 0 0 5 0 】

他ノード１０の受信処理部１２は、追加された新規ノード１０から送信情報Ｔ３を通知されると、送信情報Ｔ３に含まれるＩＰアドレス及びポート番号、並びに送信元のノードＩＤの情報をノード状態管理情報Ｔ２に追加する。以後、ノード１０（送信処理部１４）は、追加した新規ノード１０のＩＰアドレス及びポート番号に対してもハートビートを送

信する。

【 0 0 5 1 】

また、新規ノード 1 0 (受信処理部 1 2) は、新規ノード 1 0 が通知した送信情報 T 3 を受け取った他ノード 1 0 の各々から、順次ハートビート (ノード状態情報 T 1) を受信する。なお、新規ノード 1 0 が受け取るノード状態情報 T 1 は、図 4 に示すノード状態情報 T 1 と同様のデータ構造であるが、新規ノード 1 0 の情報が追加されているため、便宜上、ノード状態情報 T 1 と表記する。

【 0 0 5 2 】

図 7 に示すように、新規ノード 1 0 が受信するノード状態情報 T 1 は、図 4 に示すノード状態情報 T 1 に加えて、新規ノード 1 0 に対応するノード ID “ 6 ” の状態を含む。一例として、ノード ID “ 6 ” には、他ノード 1 0 が判定した新規ノード 1 0 の状態 “ Alive ”、IP アドレス “ 192.168.0.6 ”、ポート番号 “ 12345 ” が対応付けられる。

10

新規ノード 1 0 (受信処理部 1 2) は、受信したノード状態情報 T 1 に含まれる他ノード 1 0 の IP アドレス及びポート番号、並びに送信元のノード ID の情報からノード状態管理情報 T 2 を作成又は更新する。これにより、新規ノード 1 0 は、ノード状態情報 T 1 をハートビートとして定期的に送信する送信処理部 1 4 のサービスを開始することができる。

【 0 0 5 3 】

〔 1 - 3 - 3 〕 ノード状態決定部

ノード状態決定部 (判定部) 1 3 は、上記 (b) の処理を行なう。具体的には、ノード状態決定部 1 3 は、ノード状態管理情報 T 2 を参照してノード 1 0 ごとの状態を判定し、ノード状態管理情報 T 2 に設定する。より具体的に、ノード状態決定部 1 3 は、受信処理部 1 2 が受信したノード状態情報 T 1 が示す複数のノード 1 0 の各々の状態と、他ノード 1 0 の各々からのノード状態情報 T 1 の受信状況とに基づいて、複数のノード 1 0 の各々の状態を判定する。

20

【 0 0 5 4 】

ここで、ノード 1 0 の状態及び状態遷移について説明する。

図 8 は、第 1 実施形態の一例としてのノード 1 0 が他ノード 1 0 の状態を判定するときの状態遷移の一例を示す図であり、図 9 は、複数のノード 1 0 によるノード状態情報 T 1 の送受信処理の一例を説明する図である。図 1 0 は、ノード 1 0 が自ノード 1 0 の状態を判定するときの状態遷移の一例を示す図である。

30

【 0 0 5 5 】

なお、図 9 に示す例においては、説明の簡略化のため、ノード 1 0 間の接続状態のみを示し、スイッチ 2 0 の図示を省略している。

〔 1 - 3 - 3 - 1 〕 ノード状態決定部が他ノードについて判定する各状態の説明

はじめに、ノード 1 0 (ノード状態決定部 1 3) が他ノード 1 0 について判定する各状態について説明する。図 8 に示すように、ノード 1 0 が他ノード 1 0 について判定する状態には、Alive、Suspect、Down、及び Zombie が含まれる。

【 0 0 5 6 】

Alive は、ノード 1 0 が正常に動作している状態 (稼動中) を示す。ノード状態決定部 1 3 は、ノード状態管理情報 T 2 を参照して、最終更新情報が第 2 所定時間内であり、且つ第 1 所定数以上のノード 1 0 から Suspect と判定されていない他ノード 1 0 の状態を、Alive と判定する。

40

なお、他ノード 1 0 がストレージシステム 1 に追加された場合、ノード状態決定部 1 3 は、追加された他ノード 1 0 に関する最初の判定において、追加された他ノード 1 0 の状態を初期状態である Alive と判定する (図 8 の矢印 (I) 参照)。

【 0 0 5 7 】

ここで、第 2 所定時間としては、ノード 1 0 がノード状態情報 T 1 を送信する時間周期である第 1 所定時間以上の時間とすることができる。例えば、各ノード 1 0 がノード状態情報 T 1 を 1 秒 (第 1 所定時間) ごとに送信する場合、第 2 所定時間は、ノード 1 0 の負

50

荷による送信処理の遅延又は通信経路の輻輳等を考慮して、数倍～数十倍程度の時間（例えば20秒）とすることができる。

【0058】

また、第1所定数としては、例えば過半数とすることができる。

以下、第1所定時間は1秒であり、第2所定時間は20秒であり、第1所定数はノード10の数の過半数であるものとして説明する。

Suspect（第1状態）は、ノード10が故障（停止）している疑いのある状態（停止の可能性）を示す。ノード状態決定部13は、ノード状態管理情報T2を参照して、最終更新情報が第2所定時間よりも前である他ノード10の状態、つまり第2所定時間内にノード状態情報T1を受信しなかった他ノード10の状態を、Suspectと判定する。すなわち、ノード状態決定部13は、ハートビートの不達時間が閾値（第2所定時間）を超えた他ノード10の状態を、Suspectと判定する。

10

【0059】

例えば、ノード状態決定部13は、Aliveの状態と判定した他ノード10から、20秒よりも長くノード状態情報T1を受信できない場合、当該他ノード10の状態をAliveからSuspectに遷移させる（図8の矢印（II）参照）。

また、ノード状態決定部13は、Suspectの状態と判定した他ノード10の状態が自ノード10又は他ノード10によりDownと判定される前に、当該他ノード10からノード状態情報T1を受信する場合がある。この場合、ノード状態決定部13は、当該他ノード10の状態をSuspectからAliveに遷移させる（図8の矢印（III）参照）。

20

【0060】

Down（第2状態）は、ノード10において故障等の障害が発生している状態（停止中）を示す。ノード状態決定部13は、第1所定数以上のノード10でSuspectと判定された他ノードの状態、又は他ノード10の少なくとも1つからDownと判断された他ノード10の状態を、Downと判定する。

例えば、ノード状態決定部13がAlive又はSuspectの状態と判定した他ノード10の状態について、過半数以上のノード10でSuspectと判定される場合、又は他ノード10のうちのいずれかがDownと判定される場合がある。この場合、ノード状態決定部13は、Alive又はSuspectの状態と判定した当該他ノード10の状態を、Downと判定する（図8の矢印（IV）又は（V）参照）。

30

【0061】

一例として、図9に示すように、ノード10-1がノード10-2、10-3、及び10-5からノード状態情報T1を1秒ごとに受け取る一方、ノード10-4からノード状態情報T1を30秒間受け取っていない場合を考える。このとき、ノード状態管理情報T2は、図5に例示する状態になる。

つまり、自ノード10-1は、ノード10-4から20秒よりも長くノード状態情報T1を受け取っていないため、ノード10-4の状態をSuspectと判定する。また、他ノード10-3及び10-5も、ノード10-4から20秒よりも長くノード状態情報T1を受け取っておらず、他ノード10-3及び10-5によるノード10-4の状態の判定結果もSuspectとなる。この場合、ノード状態決定部13は、ノード10-4の状態が過半数のノード10によりSuspectと判定されたため、ノード10-4の状態をDownに遷移させる。

40

【0062】

このように、他ノード10に障害等が発生した場合、ノード状態管理情報T2では、障害等が発生した他ノード10の状態が、当該他ノード10のノードIDの行方向（図5中、横軸方向）に順にSuspectに遷移する（図5中、ノードID“4”参照）。そして、ノード状態決定部13は、Suspectになったノード10の数が過半数に達した場合に、当該他ノード10の状態をDownと判定するのである。

【0063】

50

なお、図 9 に示す例では、ノード 10 - 1 は、ノード 10 - 2 ~ 10 - 5 へノード状態情報 T1 を送信するが、ノード 10 - 4 は故障（停止）している（又は疑いのある）状態であるため、ノード状態情報 T1 はノード 10 - 4 により受信されない。

Zombie（第 3 状態）は、ノード 10 が後述するリカバリ処理部 15 によりリカバリ処理が行なわれている状態（リカバリ処理中）を示す。Zombie は、ノード 10 に故障等の障害が発生した後、障害が発生したノード 10 のノード情報が削除されるまでの暫定状態である。クライアント及びリカバリ処理に係わるノード 10 以外のノード 10 は、Zombie の状態のノード 10 へのアクセスが制限される。

#### 【0064】

具体的には、ストレージシステム 1 は、障害が発生したノード 10 について、障害が発生したノード 10 が保持するデータに関連するデータを持つノード 10 により、リカバリ処理を実行させる。リカバリ処理は、上述のように、障害が発生したノード 10 内のデータの冗長データを保持するノード 10 から、冗長データを他ノード 10 へコピーし、データの多重化状態を回復する処理である。

#### 【0065】

例えば、リカバリ処理部 15 によるリカバリ処理中に、障害が発生したノード 10 が復旧する場合、又は同一のノード名でストレージシステム 1 に追加される場合があり得る。この場合、ストレージシステム 1 上で障害が発生したノード 10 に古いデータが存在する状態で、古いデータと独立してリカバリ処理が実行される状態が発生し、データの一貫性が崩れる可能性がある。

#### 【0066】

クライアントは、ストレージシステム 1 内でどのノード 10 にデータが格納されているかを管理するテーブルを保持するが、このテーブルでは、ノード 10 に障害が発生したことを即座に検知できない場合がある。仮に、クライアントが、障害が発生したノード 10 からデータ（古いデータ）を取得してしまうと、取得したデータと、リカバリ処理が行なわれ、他ノード 10 にコピーされた冗長データとの間で不整合が生じることになる。

#### 【0067】

以上の理由から、ノード状態決定部 13 は、障害が発生したノード 10 の状態を、リカバリ処理が完了する（古いデータが削除される）まで、Zombie と判定する。これにより、ノード状態決定部 13 は、Zombie のノード 10 に対して、クライアント及びリカバリ処理に係わるノード 10 以外のノード 10 からアクセスできないようにし、データの一貫性が崩れることを防止する。従って、Zombie の状態である期間は、リカバリ処理が完了するまで、障害が発生したノード 10 から古いデータが読み出されることを抑止するガード期間であるといえる。

#### 【0068】

ノード状態決定部 13 は、第 2 所定数以上のノード 10 で Down と判定された他ノードの状態を、Zombie と判定する。

ここで、第 2 所定数としては、第 1 所定数以上の数、好ましくは、全てのノード 10 の数とすることができる。以下、第 2 所定数は全てのノード 10 の数であるものとして説明する。

#### 【0069】

例えば、ノード状態決定部 13 は、自ノード 10 を含め全てのノード 10（Down 又は Zombie の状態のノード 10 を除く）で Down の状態と判定されたノード 10 の状態を、Down から Zombie に遷移させる（図 8 の矢印（VI）参照）。

ノード状態決定部 13 は、全てのノード 10 から Down と判定されたノード 10 を Zombie とすることで、全てのノード 10 の共通認識によって、障害が発生したノード 10 をリカバリすべきノード 10 であると確実に決定することができる。

#### 【0070】

なお、リカバリ処理が完了すると、障害が発生したノード 10 以外のノード 10 のノード状態決定部 13 は、自ノード 10 が保持するノード状態管理情報 T2 から、障害が発生

10

20

30

40

50

したノード10に関する情報を削除する(図8の矢印(VII)参照)。

以上のように、ノード状態決定部13は、受信処理部12が受信したノード状態情報T1が示す複数のノード10の各々の状態、つまり図5のノード状態管理情報T2における“他ノードからの情報”に基づいて、複数のノード10の各々の状態を判定する。

【0071】

また、ノード状態決定部13は、以下に説明するように、ノード状態決定部13が含まれる自ノード10の状態に関して判断を行なった自己状態情報(図5の“自ノードでの判断”参照)にさらに基づいて、複数のノード10の各々の状態を判定してもよい。

〔1-3-3-2〕ノード状態決定部が自ノードについて判定する各状態の説明

次に、ノード10(ノード状態決定部13)が自ノード10について判定する各状態について説明する。図10に示すように、ノード10が自ノード10について判定する状態には、Alive、Isolate、及びDownが含まれる。

【0072】

Alive(初期状態)は、ノード状態決定部13が他ノード10について判定するAliveと同様の状態であり、自ノード10が正常に動作している状態(稼動中)を示す。

自ノード10が起動したとき、ノード状態決定部13は、自ノード10に関する最初の判定において、自ノード10の状態をAliveと判定する(図10の矢印(i)参照)。

【0073】

Isolate(第4状態)は、自ノード10がストレージシステム1から切り離された状態を示す。Isolateの状態になる場合としては、例えば自ノード10からスイッチ20までの経路で障害が発生した場合や、自ノード10のネットワークインタフェース10dが故障した場合等が挙げられる。

ノード状態決定部13は、ノード状態管理情報T2を参照して、第2所定時間内に第3所定数以上の他ノード10からノード状態情報T1を受信しなかった場合、自ノード10の状態を、AliveからIsolateに遷移させる。すなわち、ノード状態決定部13は、ハートビートの不達時間が閾値(第2所定時間)を超えた他ノード10の数が第3所定数以上である場合、自ノード10の状態をIsolateと判定するのである。

【0074】

ここで、第3所定数としては、第1所定数と同様、例えばノード10の数の過半数とすることができる。

以下、第3所定数はノード10の数の過半数であるものとして説明する。

例えば、ノード状態決定部13は、ハートビートの不達時間が閾値を超えたノード10の数が過半数に達した場合に、自ノード10の状態をIsolateに遷移させる(図10の矢印(ii)参照)。

【0075】

なお、自ノード10が経路障害等によりストレージシステム1から切り離された場合、受信処理部12は、他ノード10からハートビートを受信しない。その結果、ノード状態管理情報T2では、自ノード10で判定したノード10ごとの状態が列方向(図5中、縦軸方向)に順にSuspectに遷移する。そして、ノード状態決定部13は、Suspectになったノード10の数が過半数に達した場合に、自ノード10の状態をIsolateと判定するのである。

【0076】

また、自ノード10の状態がIsolateに遷移した場合、後述する停止処理部16による停止処理により、自ノード10は停止する(図10の矢印(iii)参照)。

ところで、ノード10は、自ノード10の状態がIsolateに遷移した場合、ストレージシステム1から切り離されているため、自ノード10の状態がIsolateであることを他のノードへノード状態情報T1により伝えることができない。また、ノード10は、他ノード10の状態がIsolateになった場合にも、当該他ノード10はスト

10

20

30

40

50



レンジシステム 1 から切り離されているため、他ノード 10 の状態が *I s o l a t e* になったことをノード状態情報 T 1 により検知することができない。

【 0 0 7 7 】

自ノード 10 の状態が *I s o l a t e* に遷移した場合、他ノード 10 間でやり取りされるノード状態情報 T 1 内では、自ノード 10 の状態として *S u s p e c t*、*D o w n*、*Z o m b i e* の順で遷移する。換言すれば、ノード 10 は、他ノード 10 の状態を *S u s p e c t*、*D o w n*、又は *Z o m b i e* と判定する場合、当該他ノード 10 自身で判定した状態は *I s o l a t e* である可能性がある。

【 0 0 7 8 】

*D o w n* (第 2 状態) は、ノード状態決定部 13 が他ノード 10 について判定する *D o w n* と同様の状態であるが、*D o w n* に遷移するまでの判定内容が、他ノード 10 について判定する場合と異なる。ノード 10 (例えばノード状態決定部 13) は、自ノード 10 内で所定の障害が発生したことを検出した場合、自ノード 10 の状態を *A l i v e* から *D o w n* に遷移させる。

【 0 0 7 9 】

所定の障害としては、例えば自ノード 10 による復旧が不可能又は困難な障害であり、ハードウェア障害等が挙げられる。なお、ノード 10 による自ノード 10 での障害の発生  
の検出は、既知の種々の手法により行なうことが可能であり、その詳細な説明は省略する。

ノード状態決定部 13 は、自ノード 10 に例えば復旧不可能な障害が発生した場合、自  
ノード 10 の状態を、*D o w n* と判定する (図 10 の矢印 (iv) 参照)。

【 0 0 8 0 】

また、自ノード 10 の状態が *D o w n* に遷移した場合、後述する停止処理部 16 による  
停止処理により、自ノード 10 は停止する (図 10 の矢印 (v) 参照)。

なお、ノード 10 は、自ノード 10 の状態を *I s o l a t e* 又は *D o w n* と判定した場合、他ノード 10 で判定される自ノード 10 の状態は、*S u s p e c t*、*D o w n*、*Z o m b i e* の順で遷移する。

【 0 0 8 1 】

他ノード 10 により、自ノード 10 の状態が *Z o m b i e* と判定されると、上述のよう  
に、自ノード 10 に対するリカバリ処理が実行され、自ノード 10 以外のノード 10 が保  
持するノード状態管理情報 T 2 から、自ノード 10 に関する情報が削除される。

ノード状態決定部 13 は、上述のように、自ノード 10 及び他ノード 10 の状態を判定  
し、ノード状態管理情報 T 2 を更新する。

【 0 0 8 2 】

具体的には、ノード状態決定部 13 は、自ノード 10 及び他ノード 10 の各々について  
判定した状態を、図 5 に例示するノード状態管理情報 T 2 における“状態”の列に設定す  
る。

ノード状態決定部 13 は、以上のようにして、ノード状態管理情報 T 2 に基づき複数の  
ノード 10 の各々の状態を判定することができる。つまり、ノード状態決定部 13 は、受  
信処理部 12 が受信したノード状態情報 T 1 が示す複数のノード 10 の各々の状態と、  
ノード状態決定部 13 が含まれる自ノード 10 の状態に関するノード状態情報 T 1 に関する  
自己状態情報とに基づいて、上記判定を行なう。

【 0 0 8 3 】

なお、ノード状態決定部 13 による上述した判定は、第 1 所定時間置きに全ノード 10  
について一括で行なわれてもよいし、ノード 10 ごとに異なるタイミングで、第 1 所定時  
間置きに行なわれてもよい。

また、ノード 10 は、ノード状態決定部 13 により自ノード 10 の状態が *D o w n* 又は  
*I s o l a t e* と判定された場合、ノード状態保持部 11 が保持するノード状態管理情報  
T 2 を、記録媒体 10 f 等の不揮発性メモリに保存してもよい。これにより、リカバリ処  
理後、作業等者は、ノード 10 の停止要因が復旧不可能又は困難な障害 (*D o w n*) によ

10

20

30

40

50

るものか、ストレージシステム 1 から切り離されたこと ( I s o l a t e ) によるものかを判断でき、障害復旧を迅速に行なうことができる。

【 0 0 8 4 】

〔 1 - 3 - 4 〕 送信処理部

送信処理部 1 4 は、上記 ( c ) の処理を行なう。具体的には、送信処理部 1 4 は、第 1 所定時間ごとに、ノード状態決定部 1 3 が判定した複数のノード 1 0 の各々の状態に関するノード状態情報 T 1 を、他ノード 1 0 の各々へ送信する。

より具体的に、送信処理部 1 4 は、ノード状態管理情報 T 2 を参照して、IP アドレス及びポート番号を取得し、他ノード 1 0 へ送信するノード状態情報 T 1 の宛先ノードを判定する。また、送信処理部 1 4 は、ノード状態管理情報 T 2 を参照して、自ノード 1 0 が判定した各ノード 1 0 についてのノード ID、状態、IP アドレス、及びポート番号の情報からノード状態情報 T 1 を生成する。そして、送信処理部 1 4 は、生成したノード状態情報 T 1 を、ハートビートとして他ノード 1 0 の各々へ送信する。

【 0 0 8 5 】

以下、受信処理部 1 2 が受信するノード状態情報 T 1 及び送信処理部 1 4 が送信するノード状態情報 T 1 は、同様のデータ構造であるが、便宜上、送信処理部 1 4 が送信するノード状態情報 T 1 を送信用ノード状態情報 ( 送信用状態情報 ) T 1 という場合がある。

また、送信処理部 1 4 は、ノード状態情報 T 1 の送信に加え、上述のように、自ノード 1 0 の起動後、送信情報 T 3 ( 図 6 参照 ) をストレージシステム 1 内の全てのノード 1 0 へブロードキャスト等により通知する。

【 0 0 8 6 】

〔 1 - 3 - 5 〕 リカバリ処理部

リカバリ処理部 1 5 は、他ノード 1 0 に対してリカバリ処理を実行する。具体的には、リカバリ処理部 1 5 は、ノード状態決定部 1 3 が Z o m b i e と判定したノード 1 0 に対して、リカバリ処理を実行する。

なお、リカバリ処理は、全てのノード 1 0 により行なわれなくても、ノード状態管理情報 T 2 において Z o m b i e と判定されたノード 1 0 に関係するノード 1 0 により行なわれればよい。

【 0 0 8 7 】

例えば、Z o m b i e と判定されたノード 1 0 内のデータの冗長データ又は関連するデータを保持するノード 1 0 のリカバリ処理部 1 5 が、上記冗長データ又は関連するデータを他ノード 1 0 の記憶部 1 0 c へコピーすればよい。又は、上記冗長データ又は関連するデータのコピー先のノード 1 0 のリカバリ処理部 1 5 が、上記冗長データ又は関連するデータを保持するノード 1 0 からデータを取得し、自ノード 1 0 の記憶部 1 0 c 等へ保存してもよい。

【 0 0 8 8 】

リカバリ処理部 1 5 は、リカバリ処理においてコピーが完了すると、Z o m b i e と判定されたノード 1 0 内のデータを削除して、リカバリ処理を終わる。なお、Z o m b i e と判定されたノード 1 0 が停止した場合等には、リカバリ処理部 1 5 は、Z o m b i e と判定されたノード 1 0 内のデータを削除できない可能性がある。この場合、リカバリ処理部 1 5 は、Z o m b i e と判定されたノード 1 0 内のデータの削除を行わずに、リカバリ処理を終わってもよい。また、リカバリ処理部 1 5 は、リカバリ処理が完了すると、リカバリ処理の完了をノード状態決定部 1 3 へ通知する。

【 0 0 8 9 】

ノード状態決定部 1 3 は、リカバリ処理の完了が通知されると、Z o m b i e と判定したノード 1 0 に関するノード ID、自ノード 1 0 及び各ノード 1 0 で判定された状態、最終更新情報、IP アドレス、及びポート番号をノード状態管理情報 T 2 から削除する。これにより、各ノード 1 0 は、Z o m b i e と判定されたノード 1 0 をストレージシステム 1 から完全に切り離すことができる。

【 0 0 9 0 】

10

20

30

40

50

なお、ノード状態管理情報 T 2 から情報を削除されたノード 10 は、例えば作業者等により、修理又は交換によりストレージシステム 1 へ再組み込み可能な状態になると、起動され、自ノード 10 の状態を A l i v e と判定する（図 10 の矢印（i）参照）。このとき、上述のように、新規ノード 10 は IP アドレス及びポート番号を他ノード 10 へ通知し、各ノード 10 のノード状態管理情報 T 2 に新規ノード 10 の情報が追加されて、使用可能な状態になる。

#### 【0091】

ところで、ストレージシステム 1 の運用において、作業者等は、障害が発生したノード（障害ノード）10 が保持する情報の初期化、又は故障箇所等の交換（ノード 10 全体又は部品交換等の場合）を行ない、障害ノード 10 の復旧を行なう。そして、作業者等は、復旧した障害ノード 10 をストレージシステム 1 へ再組み込みすることで、他ノード 10 に新規ノード 10 として認識させることができる。従って、障害により低下した、障害ノード 10 に関するデータの多重度及びノード 10 の冗長度を回復させることができる。

#### 【0092】

また、復旧前の障害ノード 10 が使用していた IP アドレスは、他ノード 10 のノード状態管理情報 T 2 から削除されているため、復旧後の障害ノード 10 は、再組み込み後も同じ IP アドレスを使いまわすことができる。従って、ストレージシステム 1 の管理者は、ストレージシステム 1 における IP アドレスの管理を容易に行なうことができ、利便性が高い。

#### 【0093】

##### 〔1-3-6〕停止処理部

停止処理部 16 は、自ノード 10 に所定の障害が発生し、ノード状態決定部 13 が自ノード 10 の状態を D o w n と判定した場合、又は、ノード状態決定部 13 が自ノード 10 の状態を I s o l a t e と判定した場合、自ノード 10 を停止させる処理を行なう。

なお、停止処理部 16 による停止処理は、リカバリ処理部 15 によるリカバリ処理が完了した後、具体的にはリカバリ処理において故障ノード 10 内のデータが削除された後に行なわれることが好ましい。

#### 【0094】

また、他ノード 10 のリカバリ処理部 15 が、障害が発生したノード 10 に対するリカバリ処理の完了後又はリカバリ処理の過程で、障害が発生したノード 10 の停止処理を行なってもよい。この場合、停止処理部 16 を省略することができる。

以上のように、第 1 実施形態の一例としてのストレージシステム 1 によれば、複数のノード 10 により、メッシュ状に、ノード 10 間でハートビートが行なわれる。ハートビートには、各ノード 10 で判定された複数のノード 10 の各々の状態が含まれ、複数のノード 10 間でノード 10 の各々の状態が共有される。

#### 【0095】

従って、ストレージシステム 1 は、個々のノード 10 が自律的に判定した他ノード 10 の状態の判定結果に基づき、各ノード 10 の状態について、信頼性の高い判定結果を得ることができる。つまり、特定のノード又は監視装置等がノードの状態を集中的に監視する場合、特定のノード等により他ノードの状態について誤った判定がされる場合がある。これに対し、ストレージシステム 1 によれば、各ノード 10 は、複数のノード 10 から見た各ノード 10 の状態を考慮して、自ノード 10 及び他ノード 10 の状態を判定することができるため、特定のノード等により誤った判定がされることを防止できる。

#### 【0096】

また、各ノード 10 は、判定結果を共有し、信頼性の高い判定結果を得ることができるため、スプリットブレイン状態に陥ることを抑止できる。なお、ノード 10 は、仮にスプリットブレイン状態に陥ったとしても、自ノード 10 が I s o l a t e になると自律的に停止するため、冗長データの不整合が発生することを抑止できる。

さらに、各ノード 10 は、ノード 10 の状態をハートビート等の簡素な手法により共有するため、従来の手法と比較して、高速に、且つ容易に、自ノード 10 及び他ノード 10

10

20

30

40

50

の状態を判定することができる。

#### 【0097】

従って、各ノード10は、例えば、障害が発生したノード10を高速に検出することが可能となり、クライアントからストレージシステム1へのアクセスの停止時間の短縮や、信頼性が低下する時間の短縮を図ることが可能となる。

#### 〔1-4〕動作例

次に、図11～図13を参照して、上述の如く構成された第1実施形態の一例としてのノード10による動作例を説明する。図11は、第1実施形態の一例としての新規ノード10による起動後の動作例を説明するフローチャートである。図12は、ノード10による他ノード10の状態を判定する動作例を説明するフローチャートであり、図13は、ノード10による自ノード10の状態を判定する動作例を説明するフローチャートである。

10

#### 【0098】

##### 〔1-4-1〕新規ノードによる起動後の動作例

はじめに、図11を参照して、新規ノード10による起動後の動作例を説明する。

図11に示すように、ノード（新規ノード）10が起動し（ステップS1）、ストレージシステム1内のネットワークに接続されると、新規ノード10のノード状態決定部13により、自ノード10の状態がAliveと判定される（ステップS2）。

#### 【0099】

次いで、送信処理部14により、自ノード10のIPアドレス及びポート番号等のノード情報が収集され、送信情報T3（図6参照）が生成される。そして、送信処理部14により、生成した送信情報T3がブロードキャスト等によりストレージシステム1内の全てのノード10へ送信される（ステップS3）。

20

送信情報T3を受信した他ノード10の各々は、新規ノード10のノード情報をノード状態管理情報T2に追加し、宛先に新規ノード10を含めてハートビート（ノード状態情報T1）を送信する。

#### 【0100】

新規ノード10では、受信処理部12により、ハートビートが待ち受けられる（ステップS4、ステップS4のNoルート）。他ノード10からハートビートが受信されると（ステップS4のYesルート）、受信処理部12により、受信したノード状態情報T1（図7参照）から他ノード10の各々のノード情報が抽出され、ノード状態管理情報T2が作成される（ステップS5）。

30

#### 【0101】

そして、新規ノード10では、ノード状態管理情報T2に基づいて、送信処理部14による第1所定時間ごとにハートビートを送信するサービスが開始され（ステップS6）、新規ノード10による起動後に行なわれる処理が終了する。

##### 〔1-4-2〕ノードによる他ノードの状態を判定する動作例

次に、図12を参照して、ノード10による他ノード10の状態を判定する動作例を説明する。

#### 【0102】

なお、図12に示すステップS11～S23の処理は、ノード10の各々において、ノード状態決定部13により一のノード10の状態が判定される際に行なわれる処理である。従って、ステップS11～S23の処理は、各ノード10のノード状態決定部13により、他ノード10の各々について、定期的（第1所定時間ごと）に実行される。

40

図12に示すように、ノード状態決定部13により、ノード状態管理情報T2内の“状態”が参照され、判定対象のノード10について直前に判定した状態がどの状態であるかが判定される（ステップS11、S16、S19）。

#### 【0103】

判定対象のノード10について直前に判定した状態がAliveである場合（ステップS11のYesルート）、ノード状態決定部13により、判定対象のノード10からのハートビートの不達時間が閾値を超えたか否かが判定される（ステップS12）。このとき

50

、ノード状態決定部 13 は、ノード状態管理情報 T2 の“最終更新情報”の時間が第 2 所定時間よりも長いかなかを判定する。

【0104】

ハートビートの不達時間が閾値を超えた場合（ステップ S12 の Yes ルート）、ノード状態決定部 13 により、判定対象のノード 10 の状態が S u s p e c t と判定され（ステップ S13）、処理が終了する。このとき、ノード状態決定部 13 は、判定対象のノード 10 について、ノード状態管理情報 T2 内の“状態”に S u s p e c t を設定する。そして、ノード状態決定部 13 は、次の判定対象のノード 10 がある場合、次の判定対象のノード 10 に係る状態の判定処理に移行する。

【0105】

一方、ステップ S12 において、ノード状態決定部 13 により、ハートビートの不達時間が閾値を超えていないと判定された場合（ステップ S12 の No ルート）、処理がステップ S14 に移行する。ステップ S14 では、ノード状態決定部 13 により、判定対象のノード 10 の状態が、過半数（第 1 所定値）のノード 10 から S u s p e c t と判定されたか否か、又は複数のノード 10 のいずれかのノード 10 により D o w n と判定されたかが判定される。

【0106】

判定対象のノード 10 の状態が、過半数のノード 10 から S u s p e c t と判定されておらず、複数のノード 10 のいずれかのノード 10 により D o w n と判定されていない場合（ステップ S14 の No ルート）、判定対象のノード 10 に対する処理が終了する。一方、判定対象のノード 10 の状態が、過半数のノード 10 から S u s p e c t と判定された又は複数のノード 10 のいずれかのノード 10 により D o w n と判定された場合（ステップ S14 の Yes ルート）、処理がステップ S15 に移行する。

【0107】

ステップ S15 では、ノード状態決定部 13 により、判定対象のノード 10 の状態が D o w n と判定され、処理が終了する。このとき、ノード状態決定部 13 は、判定対象のノード 10 について、ノード状態管理情報 T2 内の“状態”に D o w n を設定する。

また、判定対象のノード 10 について直前に判定した状態が S u s p e c t である場合（ステップ S11 の No ルートからステップ S16 の Yes ルート）、処理がステップ S17 に移行する。ステップ S17 では、ノード状態決定部 13 により、判定対象のノード 10 から新たなハートビートが受信されたか否か、つまりハートビートの不達時間が閾値未満となったか否かが判定される。このとき、ノード状態決定部 13 は、ノード状態管理情報 T2 の“最終更新情報”の時間が第 2 所定時間未満であるか否かを判定する。

【0108】

新たなハートビートが受信されていない場合（ステップ S17 の No ルート）、処理がステップ S14 に移行する。一方、新たなハートビートが受信された場合（ステップ S17 の Yes ルート）、ノード状態決定部 13 により、判定対象のノード 10 の状態が A l i v e と判定され（ステップ S18）、処理が終了する。このとき、ノード状態決定部 13 は、判定対象のノード 10 について、ノード状態管理情報 T2 内の“状態”に A l i v e を設定する。

【0109】

判定対象のノード 10 について直前に判定した状態が D o w n である場合（ステップ S11 の No ルート、ステップ S16 の No ルートからステップ S19 の Yes ルート）、処理がステップ S20 に移行する。ステップ S20 では、ノード状態決定部 13 により、判定対象のノード 10 の状態が第 2 所定値の数のノード 10（例えば全ノード 10）により D o w n と判定されたか否かが判定される。

【0110】

全ノード 10 により D o w n と判定されていない場合（ステップ S20 の No ルート）、判定対象のノード 10 に対する処理が終了する。一方、全ノード 10 により D o w n と判定された場合（ステップ S20 の Yes ルート）、ノード状態決定部 13 により、判定

10

20

30

40

50

対象のノード10の状態がZombieと判定される。また、自ノード10が保持するデータが判定対象のノード10が保持するデータに関連する場合、リカバリ処理部15により、判定対象のノード10に対するリカバリ処理が実行され(ステップS21)、処理が終了する。このとき、ノード状態決定部13は、判定対象のノード10について、ノード状態管理情報T2内の“状態”にZombieを設定する。

#### 【0111】

判定対象のノード10について直前に判定した状態がZombieである場合(ステップS11のNoルート、ステップS16のNoルートからステップS19のNoルート)、処理がステップS22に移行する。ステップS22では、ノード状態決定部13により、判定対象のノード10についてリカバリ処理が完了したか否かが判定される。リカバリ処理が完了していない場合、判定対象のノード10に対する処理が終了する。一方、リカバリ処理が完了した場合(ステップS22のYesルート)、ノード状態決定部13により、ノード状態管理情報T2から、判定対象のノード10に関する情報が削除され(ステップS23)、処理が終了する。

10

#### 【0112】

以上のように、ノード10により、一のノード10の状態の判定処理が行なわれる。

〔1-4-3〕ノードによる自ノードの状態を判定する動作例

次に、図13を参照して、ノード10による自ノード10の状態を判定する動作例を説明する。

なお、図13に示すステップS31~S34の処理は、ノード10の各々において、ノード状態決定部13により自ノード10の状態が判定される際に行なわれる処理である。従って、ステップS31~S34の処理は、各ノード10のノード状態決定部13により、定期的(第1所定時間ごと)に実行される。

20

#### 【0113】

図13に示すように、ノード状態決定部13により、自ノード10内で所定の障害の発生、例えば修復不可能な障害の発生が検出されたか否かが判定される(ステップS31)。

所定の障害の発生が検出されると(ステップS31のYesルート)、ノード状態決定部13により、自ノード10の状態がDownと判定され(ステップS32)、処理が終了する。このとき、ノード状態決定部13は、自ノード10について、ノード状態管理情報T2内の“状態”にDownを設定する。

30

#### 【0114】

一方、所定の障害の発生が検出されない場合(ステップS31のNoルート)、ノード状態決定部13により、ハートビートの不達時間が閾値を超えたノード数が過半数に達したか否かが判定される(ステップS33)。このとき、ノード状態決定部13は、ノード状態管理情報T2の“最終更新情報”の時間が第2所定時間よりも長い他ノード10が第3所定値以上の数であるか否かを判定する。

#### 【0115】

ハートビートの不達時間が閾値を超えたノード数が過半数である場合(ステップS33のYesルート)、ノード状態決定部13により、自ノード10の状態がIsolateと判定され(ステップS34)、処理が終了する。このとき、ノード状態決定部13は、自ノード10について、ノード状態管理情報T2内の“状態”にIsolateを設定する。

40

#### 【0116】

一方、ステップS33において、ノード状態決定部13により、ハートビートの不達時間が閾値を超えたノード数が過半数未満であると判定された場合(ステップS33のNoルート)、自ノード10の状態に係る判定処理が終了する。そして、ノード状態決定部13は、次の判定対象のノード10がある場合、次の判定対象のノード10に係る状態の判定処理に移行する。

#### 【0117】

50

なお、ステップ S 3 2 又は S 3 4 において、ノード状態決定部 1 3 により、自ノード 1 0 の状態が D o w n 又は I s o l a t e と判定されると、自ノード 1 0 は、他ノード 1 0 のリカバリ処理部 1 5 からリカバリ処理を受ける。そして、自ノード 1 0 は、停止処理部 1 6 により、又は、他ノード 1 0 のリカバリ処理部 1 5 により、停止処理が行なわれる。

以上のように、ノード 1 0 により、自ノード 1 0 の状態の判定処理が行なわれる。

#### 【 0 1 1 8 】

〔 1 - 5 〕 第 1 実施形態のまとめ

このように、第 1 実施形態の一例としてのストレージシステム 1 によれば、複数のノード 1 0 の各々において、受信処理部 1 2 は、他ノード 1 0 の各々から、ノード状態情報 T 1 を受信する。また、ノード状態決定部 1 3 は、受信処理部 1 2 が他ノード 1 0 の各々から受信したノード状態情報 T 1 に基づいて、複数のノード 1 0 の各々の状態を判定する。さらに、送信処理部 1 4 は、ノード状態決定部 1 3 が判定した結果に基づき送信用ノード状態情報 T 1 を、他ノード 1 0 の各々へ送信する。

10

#### 【 0 1 1 9 】

従って、各々のノード 1 0 は、特定のノード又は監視装置等によりノード 1 0 の状態を集中的に監視するのではなく、他ノード 1 0 が判定した複数のノード 1 0 の状態に基づいて、自ノード 1 0 及び他ノード 1 0 を監視することができる。従って、特定のノード又は監視装置等の故障により、ストレージシステム 1 の利用が制限されるといった点を解消できる。また、各々のノード 1 0 が自ノード 1 0 及び他ノード 1 0 を自律的に監視するため、監視を行なうノードを決定せずに済み、さらに管理者等を介入させずに済むため、ノード 1 0 の故障後にストレージシステム 1 の利用が制限される時間を短縮できる。

20

#### 【 0 1 2 0 】

このように、第 1 実施形態の一例としてのストレージシステム 1 によれば、複数のノード 1 0 をそなえるストレージシステム 1 において、複数のノード 1 0 の状態の監視に伴う可用性の低下を抑止することができる。

また、ノード状態決定部 1 3 は、受信処理部 1 2 が受信したノード状態情報 T 1 が示す複数のノード 1 0 の各々の状態と、他ノード 1 0 の各々からのノード状態情報 T 1 の受信状況とに基づいて、複数のノード 1 0 の各々の状態を判定する。また、送信処理部 1 4 は、第 1 所定時間ごとに、送信用ノード状態情報 T 1 を、他ノード 1 0 の各々へ送信する。

#### 【 0 1 2 1 】

30

これにより、各々のノード 1 0 は、第 1 所定時間ごとの他ノード 1 0 の各々からのノード状態情報 T 1 の受信状況に応じて、複数のノード 1 0 の各々の状態を判定することができ、ノード状態情報 T 1 を送信できないノード 1 0 の異常を容易に検出することができる。

さらに、ノード状態決定部 1 3 は、第 2 所定時間内にノード状態情報 T 1 を受信しなかった他ノード 1 0 の状態を、S u s p e c t と判定する。また、ノード状態決定部 1 3 は、第 1 所定数以上の複数のノード 1 0 で S u s p e c t であると判定されたノード 1 0 の状態、又は、他ノード 1 0 の少なくとも 1 つから D o w n であると判定されたノード 1 0 の状態を、D o w n と判定する。

#### 【 0 1 2 2 】

40

これにより、各々のノード 1 0 は、自ノード 1 0 でノード状態情報 T 1 が不達になったノード 1 0 を直ちに障害等が発生したノード 1 0 であると判断せず、他ノード 1 0 の判断結果を考慮して、障害等が発生したノード 1 0 を判定することができる。これにより、ノード 1 0 は、各ノード 1 0 の状態について、信頼性の高い判定結果を得ることができる。

また、ノード状態決定部 1 3 は、第 2 所定数以上の複数のノード 1 0 で D o w n であると判定されたノード 1 0 を、Z o m b i e と判定する。また、リカバリ処理部 1 5 は、ノード状態決定部 1 3 が Z o m b i e と判定したノード 1 0 に対して、リカバリ処理を実行する。

#### 【 0 1 2 3 】

これにより、リカバリ処理部 1 5 は、第 2 所定数以上、例えば全てのノード 1 0 が D o

50

w nであると判定したノード10について、リカバリ処理を行なうため、誤った判断でリカバリ処理が行なわれることを抑止できる。また、障害等が発生したノード10の状態がリカバリ処理中を示すZombie状態になることで、クライアント又はリカバリ処理を行なわないノード10が古いデータを保持するZombie状態のノード10へアクセスすることを抑止できる。

#### 【0124】

さらに、ノード状態決定部13は、自ノード10に所定の障害が発生した場合、自ノード10の状態をDownと判定する。また、ノード状態決定部13は、第2所定時間内に第3所定数以上の他ノード10からノード状態情報T1を受信しなかった場合、自ノード10の状態を、Isolateと判定する。さらに、停止処理部16は、ノード状態決定部13が自ノード10の状態をDown又はIsolateと判定した場合、自ノード10を停止させる。

10

#### 【0125】

これにより、クライアント又はリカバリ処理を行なわないノード10が、自ノード10が保持する古いデータへアクセスすることを抑止できる。また、Isolateになったノード10が自律的に停止するため、スプリットブレイン状態に陥ったとしても、冗長データの不整合の発生を抑止できる。

#### 〔2〕第2実施形態

##### 〔2-1〕ノードの説明

次に、第2実施形態の一例としてのノード10Aについて説明する。

20

#### 【0126】

第1実施形態及び第2実施形態に係るストレージシステム1は、多数（例えば、数十から数千台）のノードをそなえることがある。

上述のように、第1実施形態に係るストレージシステム1は、全ノード10対全ノード10の完全なメッシュ状態でハートビートの通信を行なう。

一方、第2実施形態に係るストレージシステム1は、ノード10Aをある程度（例えば数～数十台程度）のまとまり（以下、パーティという）に分割し、パーティ内のノード10A間では完全メッシュのハートビートの通信を行なう。一方、パーティ間では、各パーティの代表のノード10A（代表ノード10A）同士による完全メッシュのハートビートの通信を行なう。

30

#### 【0127】

このように、第2実施形態の一例としてのストレージシステム1は、複数のノード10Aにより、階層的なノード10Aでの情報交換を行なう。これにより、ストレージシステム1は、全ノード10Aによる完全メッシュのハートビートの通信を行なうよりも、ストレージシステム1における通信負荷及び処理負荷を低減させることができる。特に、ストレージシステム1が、例えば数千台もの多数のノード10Aをそなえる場合に有効である。

#### 【0128】

##### 〔2-2〕ノードの構成

次に、図14～図23を参照して、第2実施形態の一例としてのノード10Aの構成について説明する。

40

図14は、第2実施形態の一例としてのノード10Aの機能構成例を示す図である。

第2実施形態に係るノード10Aは、第1実施形態に係るノード10と比べて、パーティ情報保持部101、パーティ間受信処理部102、パーティ間ノード状態決定部103、パーティ間送信処理部104、及びパーティ管理部105をさらにそなえる。

#### 【0129】

また、第2実施形態に係るノード10Aは、第1実施形態に係るノード10がそなえるノード状態保持部11及び受信処理部12とは一部の機能が異なるノード状態保持部11A及び受信処理部12Aをそなえる。

さらに、第2実施形態に係るノード10Aは、第1実施形態に係るノード10がそなえ

50



るノード状態決定部 13 及び送信処理部 14 とは一部の機能が異なるノード状態決定部 13A 及び送信処理部 14A をそなえる。

【0130】

なお、ノード 10A は、上述した以外の点については、以下の説明において特に言及しない限り、ノード 10 と同様の構成をそなえる。従って、以下、ノード 10A の説明において、ノード 10 がそなえる構成と同一の符号の構成についての重複した説明は省略する。

〔2-2-1〕パーティ情報保持部及びノード状態保持部

パーティ情報保持部 101 は、図 15 に示すパーティ管理情報 T4 を保持する記憶領域であり、例えば上述したメモリ 10b により実現される。

10

【0131】

図 15 は、第 2 実施形態の一例としてのノード 10A が管理するパーティ管理情報を例示する図である。

上述のように、第 2 実施形態の一例としてのストレージシステム 1 は、複数のノード 10A を数～数十台程度の複数のパーティに分割する。

パーティ管理情報 T4 は、複数のパーティとパーティに属するノード 10A とを対応付けて管理する情報である。なお、ノード 10A は、図 15 に示すようにパーティ管理情報 T4 をテーブルとして生成し、送受信することができる。

【0132】

図 15 に示すように、パーティ管理情報 T4 は、パーティの識別情報の一例であるパーティ ID、パーティに属するノード 10A の識別情報の一例であるノード ID、及びパーティのバージョン番号を含む。図 15 に示すパーティ管理情報 T4 は、パーティ ID “A” ～ “E” についての情報を含む。

20

一例として、パーティ ID “A” には、ノード ID “1～10”、バージョン番号 “1” が対応付けられる。

【0133】

なお、パーティの識別情報として、パーティ ID を例に挙げたが、これに限定されるものではない。識別情報は、各パーティを特定できるユニークな情報であればよい。例えば、識別情報として、アルファベットの他、数値、ノード ID の範囲の最小値又は最大値、IP アドレスのマスク等が用いられてもよい。

30

また、ノード 10 の識別情報として、ノード ID を例に挙げたが、これに限定されるものではなく、第 1 実施形態において既述のように、ノード 10A を特定できるユニークな情報であればよい。

【0134】

なお、図 15 に例示するパーティ管理情報 T4 において、ノード ID にはパーティに属するノード 10A のノード ID の範囲（最小値～最大値）が設定されているが、これに限定されるものではない。例えば、ノード ID には、パーティに属するノード 10A のノード ID が複数の範囲、又は一つずつ設定されてもよい。

バージョン番号は、ノード 10A において、自ノード 10A が持つパーティ管理情報 T4 が最新の情報であるか否かを判断するために用いられる。例えば、後述するパーティ管理部 105 により、パーティが分割又は統合される場合がある。この場合、分割又は統合が行なわれたパーティに属するノード ID も変化するため、各ノード 10A は、バージョン番号を参照して、最新のパーティ管理情報 T4 を識別するのである。

40

【0135】

ノード状態保持部 11A は、図 19 に示すノード状態管理情報 T7 を保持する記憶領域であり、例えば上述したメモリ 10b により実現される。

〔2-2-2〕パーティ間受信処理部及び受信処理部

次に、図 16～図 19 を参照して、パーティ間受信処理部 102 及び受信処理部 12 について説明する。

【0136】

50

図 16 は、第 2 実施形態の一例としての複数のノード 10A による代表ノード状態情報 T5 及びノード状態情報 T6 の送受信処理の一例を説明する図である。図 17 は、ノード 10A が送受信する代表ノード状態情報 T5 を例示する図であり、図 18 は、ノード 10A が送受信するノード状態情報 T6 を例示する図である。図 19 は、ノード 10A が管理するノード状態管理情報 T7 を例示する図である。

【0137】

なお、図 16 に示す例においては、説明の簡略化のため、ノード 10A 間の接続状態のみを示し、スイッチ 20 の図示を省略している。

図 16 に例示するように、代表ノード（代表ストレージ装置、代表情報処理装置）10A は、複数のパーティのうちの自パーティ以外の他のパーティの各々における他の代表ノード 10A との間で、代表ノード状態情報 T5 を送受信する。また、代表ノード 10A は、自パーティのパーティメンバであるメンバノード 10A へ代表ノード状態情報 T5 を送信し、メンバノード 10A は、自パーティの代表ノード 10A へノード状態情報 T6 を送信する。

【0138】

なお、図 16 に示す例において、丸で囲われた数字は、ノード ID を示す。以下、例えばノード ID “1” の代表ノード 10A を特定する場合には、代表ノード 10A - 1 又はノード 10A - 1 と表記する。また、例えばノード ID “2” のメンバノード 10A を特定する場合には、メンバノード 10A - 2 又はノード 10A - 2 と表記する。

代表ノード 10A 及びメンバノード 10A は、特に言及しない限り互いに同様の機能を提供することができるため、以下の説明において、任意のノード 10A が提供する機能について説明する。

【0139】

パーティ間受信処理部（グループ間受信処理部）102 は、自ノード 10A がパーティの代表ノード 10A である場合に、他のパーティの各々の代表ノード 10A から、図 17 に例示する代表ノード状態情報 T5 を受信する。そして、代表ノード 10A のパーティ間受信処理部 102 は、受信した代表ノード状態情報 T5 に基づいて、ノード状態保持部 11A が保持するノード状態管理情報 T7（図 19 参照）を更新する。

【0140】

受信処理部 12A は、自ノード 10A が属するパーティ内の自ノード 10A 以外の他ノード 10A（自パーティ内の代表ノード 10A を含む）の各々から、代表ノード状態情報 T5 又は図 18 に例示するノード状態情報 T6 を受信する。そして、受信処理部 12A はノード状態保持部 11A が保持するノード状態管理情報 T7（図 19 参照）を更新する。

代表ノード状態情報（代表状態情報）T5 は、送信元の代表ノード 10A により判定された複数のパーティの代表ノード 10A の各々の状態に関する情報である。例えば、代表ノード 10A が送信する代表ノード状態情報 T5 には、代表ノード 10A が判定した自パーティ内のメンバノード 10A の状態と、他のパーティの代表ノード 10A から取得した他のパーティに属する全てのノード 10A の状態が含まれる。なお、代表ノード 10A は、図 17 に示すように代表ノード状態情報 T5 をテーブルとして生成し、送受信することができる。

【0141】

例えば、図 17 に示す例では、図 16 に示す代表ノード 10A - 1 は、他の代表ノード 10A - 11 及び 10A - 21 へ送信する代表ノード状態情報 T5 に、自パーティ内で判定した自パーティ内の各ノード 10A - 1 ~ 10A - 3 の状態を含める。また、代表ノード 10A - 1 は、代表ノード状態情報 T5 に、他の代表ノード 10A - 11 及び 10A - 21 から受信した他のパーティ内のノード 10A - 11 ~ 10A - 13 及び 10A - 21 ~ 10A - 23 の状態を含める。

【0142】

また、代表ノード 10A は、自パーティ内のメンバノード 10A - 2 及び 10A - 3 に対しても他の代表ノード 10A へ送信するものと同様の代表ノード状態情報 T5 を送信し

10

20

30

40

50

、メンバノード１０Ａ - ２及び１０Ａ - ３からはノード状態情報Ｔ６を受信する。

つまり、パーティ内の代表ノード１０Ａ及びメンバノード１０Ａは、互いにパーティ内のノード１０Ａの状態の判定結果をハートビートで通知し合い、代表ノード１０Ａは、自パーティ内での判定結果を全パーティの代表ノード１０Ａへ伝達する。

#### 【０１４３】

なお、代表ノード状態情報Ｔ５のデータ構造は、図４に示すノード状態情報Ｔ１と基本的に同様であるため、詳細な説明は省略する。

ノード状態情報（状態情報）Ｔ６は、送信元のノード１０Ａで判定された自パーティにおける他ノード（メンバノード）１０Ａの各々の状態を含む情報である。例えば、図１８に示す例では、図１６に示すメンバノード１０Ａ - ２は、自パーティに属するノード１０  
10  
Ａ - １及び１０Ａ - ３へ送信するノード状態情報Ｔ６に、自パーティ内で判定した各ノード１０Ａ - １～１０Ａ - ３の状態を含める。なお、ノード１０Ａは、図１８に示すようにノード状態情報Ｔ６をテーブルとして生成し、送受信することができる。

#### 【０１４４】

なお、ノード状態情報Ｔ６のデータ構造は、図４に示すノード状態情報Ｔ１と基本的に同様であるため、詳細な説明は省略する。

以下、代表ノード状態情報Ｔ５及びノード状態情報Ｔ６を、単にノード状態情報Ｔ５及びＴ６と表記する場合がある。

ノード状態管理情報Ｔ７は、自ノード１０Ａ及び全パーティの全ノード１０Ａで判定された複数のノード１０Ａの各々の状態を管理する情報である。なお、ノード１０Ａは、図  
20  
１９に示すようにノード状態管理情報Ｔ７をテーブルとして生成し、管理することができる。

#### 【０１４５】

以下、図１９の説明においては、自ノード１０Ａが代表ノード１０Ａ - １であるものとする。

図１９に示すように、ノード状態管理情報Ｔ７は、図５に示すノード状態管理情報Ｔ２と同様に、ノード１０ＡのノードＩＤ、ノード１０Ａごとの状態、ノード１０ＡのアドレスのＩＰアドレス、及びノード１０Ａのポート番号を含む。また、ノード状態管理情報Ｔ  
30  
７はさらに、他のノード１０Ａから受信したノード状態情報Ｔ５又はＴ６に含まれるノード１０Ａごとの状態、及び他のノード１０Ａごとの最終更新情報を含む。例えば、他のノード１０Ａから受信したノード状態情報Ｔ５又はＴ６に含まれるノード１０Ａごとの状態には、“by 2”、“by 3”、“by 11”～“by 13”、及び“by 21”～“by 23”が含まれる。

#### 【０１４６】

図１９に示すノード状態管理情報Ｔ７は、ノード１０Ａ - １～１０Ａ - ３、１０Ａ - 11～１０Ａ - 13、及び１０Ａ - 21～１０Ａ - 23に対応するノードＩＤ“1”～“3”、“11”～“13”、及び“21”～“23”の状態を含む。

一例として、ノードＩＤ“1”には、自ノード１０Ａが判定した状態“Alive”、他ノード１０Ａ - 2、１０Ａ - 3、１０Ａ - 11、及び１０Ａ - 21がそれぞれ判定した状態  
40  
“Alive”、最終更新情報“1 sec ago”が対応付けられる。また、ノードＩＤ“1”にはさらに、ＩＰアドレス“192.168.0.1”、ポート番号“12345”が対応付けられる。

#### 【０１４７】

パーティ間受信処理部１０２は、他の代表ノード１０Ａの各々から上述した代表ノード状態情報Ｔ５を受信すると、ノード状態管理情報Ｔ７を更新する。また、受信処理部１２Ａは、自パーティ内の他ノード１０Ａの各々から上述したノード状態情報Ｔ５又はＴ６を受信すると、ノード状態管理情報Ｔ７を更新する。具体的には、パーティ間受信処理部  
50  
１０２及び受信処理部１２Ａは、受信したノード状態情報Ｔ５又はＴ６に含まれるノード１０Ａごとの状態を、ノード状態管理情報Ｔ７における対応する他ノード１０Ａの列に設定する。つまり、図１９に例示する他ノード１０Ａが判定した状態は、対応する他ノード１０Ａからの情報に基づき設定される。

## 【 0 1 4 8 】

なお、パーティ間受信処理部 1 0 2 及び受信処理部 1 2 A による、ノード状態管理情報 T 7 の更新は、第 1 実施形態に係る受信処理部 1 2 による処理と同様であるため、重複した説明は省略する。

パーティ間受信処理部 1 0 2 及び受信処理部 1 2 A は、受信処理部 1 2 と同様に、ノード状態情報 T 5 又は T 6 を受信した都度、又は第 1 所定時間ごとに、ノード状態管理情報 T 7 を更新する。

## 【 0 1 4 9 】

なお、受信処理部 1 2 A は、上述したノード状態情報 T 5 又は T 6 の受信に加え、図 6 及び図 7 を用いて上述したように、新規に追加されたノード 1 0 A の I P アドレス及びポート番号を受信することができる。

10

また、パーティ間受信処理部 1 0 2 は、上述した代表ノード状態情報 T 5 の受信に加え、図 1 5 に示すパーティ管理情報 T 4 を受信することができる。

## 【 0 1 5 0 】

パーティ間受信処理部 1 0 2 は、代表ノード 1 0 A からパーティ管理情報 T 4 を受信すると、ノード状態保持部 1 1 A が保持するパーティ管理情報 T 4 と比較する。そして、パーティ間受信処理部 1 0 2 は、受信したパーティ管理情報 T 4 に、新たに追加されたパーティ ID、又はバージョン番号が更新されたパーティ ID がある場合、当該パーティ ID の情報を用いて自ノード 1 0 A が保持するパーティ管理情報 T 4 を更新する。

## 【 0 1 5 1 】

20

ところで、各パーティの代表ノード 1 0 A は、所定のルールに基づいて決定される。例えば、代表ノード 1 0 A は、各ノード 1 0 A が保持するパーティ管理情報 T 4 及びノード状態管理情報 T 7 等に基づいて求められる。

一例として、代表ノード 1 0 A は、パーティに属するノード 1 0 A の中で、最も小さいノード ID を持つノード 1 0 A とすることができる。このように、各ノード 1 0 A が保持する情報から判断可能な所定のルールを予め定めておくことで、各ノード 1 0 A は、代表ノード 1 0 A を容易に選出することができる。

## 【 0 1 5 2 】

これにより、代表ノード 1 0 A に障害等が発生した場合であっても、パーティ内のノード 1 0 A は、所定のルールに基づき次の代表ノード 1 0 A を選出することができる。また、代表ノード 1 0 A は、他のパーティの代表ノード 1 0 A が停止した場合であっても、他のパーティの新たな代表ノード 1 0 A を推定できるため、新たな代表ノード 1 0 A との間で、パーティ間のハートビートの通信を継続することができる。

30

## 【 0 1 5 3 】

〔 2 - 2 - 3 〕パーティ間ノード状態決定部及びノード状態決定部

パーティ間ノード状態決定部（グループ間判定部）1 0 3 は、パーティ間受信処理部 1 0 2 が他の代表ノード 1 0 A の各々から受信した代表ノード状態情報 T 5 に基づいて、複数の代表ノード 1 0 A の各々の状態を判定する。

なお、パーティ間ノード状態決定部 1 0 3 による、代表ノード 1 0 A 間でのノード 1 0 A の各々の状態の判定手法は、第 1 実施形態に係るノード状態決定部 1 3 によるノード 1 0 間でのノード 1 0 の各々の状態の判定と同様である。

40

## 【 0 1 5 4 】

例えば、パーティ間ノード状態決定部 1 0 3 は、受信した代表ノード状態情報 T 5 が示す複数の代表ノード 1 0 A の各々の状態と、他の代表ノード 1 0 A の各々からの代表ノード状態情報 T 5 の受信状況とに基づいて、代表ノード 1 0 A の各々の状態を判定する。

なお、代表ノード 1 0 A は、他のパーティの代表ノード 1 0 A の状態が全ての代表ノード 1 0 A から D o w n であると判定された場合、パーティ管理情報 T 4 及びノード状態管理情報 T 7 から、当該他のパーティにおいて次に代表ノード 1 0 A となるべきノード 1 0 A を判断する。この判断は、上述のように、代表ノード 1 0 A を選出する所定のルールに基づいて行なわれる。

50

## 【 0 1 5 5 】

そして、代表ノード 1 0 A は、次の代表ノード 1 0 A と判断した他のパーティのノード 1 0 A へハートビートを送信する。代表ノード 1 0 A は、ハートビートが疎通すると（他のパーティのノード 1 0 A からハートビートを受信すると）、当該他のパーティのノード 1 0 A を新たな代表ノード 1 0 A と判断する。一方、他のパーティのノード 1 0 A からのハートビートの不達時間が閾値を超えると、さらに次の代表ノード 1 0 A となるべきノード 1 0 A を判断する。

## 【 0 1 5 6 】

パーティ間ノード状態決定部 1 0 3 は、他のパーティ内の全ノード 1 0 A に対してハートビートが疎通しなかった場合、当該他のパーティに属する全ノード 1 0 A が停止したと判断する。この場合、パーティ間ノード状態決定部 1 0 3 は、当該他のパーティに属する全ノード 1 0 A の状態を Z o m b i e と判定し、リカバリ処理部 1 5 にリカバリ処理を実行させる。

10

## 【 0 1 5 7 】

ノード状態決定部（判定部）1 3 A は、受信処理部 1 2 A が自パーティにおける他ノード 1 0 A の各々から受信したノード状態情報 T 5 又は T 6 に基づいて、自パーティにおけるノード 1 0 の各々の状態を判定する。

なお、ノード状態決定部 1 3 A による、自パーティ内のノード 1 0 A 間でのノード 1 0 A の各々の状態の判定手法は、第 1 実施形態に係るノード状態決定部 1 3 によるノード 1 0 間でのノード 1 0 の各々の状態の判定と同様である。

20

## 【 0 1 5 8 】

例えば、ノード状態決定部 1 3 A は、受信したノード状態情報 T 5 又は T 6 が示す複数のノード 1 0 A の各々の状態と、他ノード 1 0 A の各々からのノード状態情報 T 5 又は T 6 の受信状況とに基づいて、自パーティ内のノード 1 0 A の各々の状態を判定する。

なお、ノード状態決定部 1 3 A は、自パーティの代表ノード 1 0 A の状態を D o w n と判定した場合、自パーティ内で生存している（A l i v e 状態の）ノード 1 0 A 間で、上述した代表ノード 1 0 A を選出する所定のルールを適用する。

## 【 0 1 5 9 】

そして、各ノード 1 0 A は、自ノード 1 0 A が代表ノード 1 0 A に昇格するか否かを判断し、昇格すると判断した場合、代表ノード 1 0 A として、他のパーティの代表ノード 1 0 A との間でハートビートの通信を開始する。

30

ここで、パーティ間ノード状態決定部 1 0 3 及びノード状態決定部 1 3 A による、ノード状態管理情報 T 7 の参照箇所及び更新箇所について説明する。なお、この説明では、パーティ間ノード状態決定部 1 0 3 及びノード状態決定部 1 3 A は、ノード 1 0 A - 1 にそなえられるものとする。

## 【 0 1 6 0 】

図 1 9 に示すように、ノード状態管理情報 T 7 における“状態”の列の、二重線で囲われた領域は、自パーティ以外の他のパーティにおいて判定された状態である。従って、ノード 1 0 A - 1 ~ 1 0 A - 3 がそなえるパーティ間ノード状態決定部 1 0 3 及びノード状態決定部 1 3 A は、二重線で囲われた領域（破線で四角く囲われた領域を除く）については基本的に判定及び更新を行なわない。

40

## 【 0 1 6 1 】

また、図 1 9 に示すように、ノード状態管理情報 T 7 における“状態”の列の、破線で四角く囲われた領域は、複数のパーティの各代表ノード 1 0 A で判定された状態である。従って、ノード 1 0 A - 1 がそなえるパーティ間ノード状態決定部 1 0 3 は、破線で四角く囲われた領域を、判定により更新する。

例えば、パーティ間ノード状態決定部 1 0 3 は、他の代表ノード 1 0 A の最終更新情報を参照し、ハートビート（代表ノード状態情報 T 5 ）の到達の有無に応じて A l i v e 又は S u s p e c t の判定を行なう。また、パーティ間ノード状態決定部 1 0 3 は、他の代表ノード 1 0 A について、ノード状態管理情報 T 7 における破線で丸く囲われた領域を参

50

照し、Suspect、Down、又はZombieの判定を多数決等により行なう。

【0162】

さらに、図19に示すように、ノード状態管理情報T7における“状態”の列の、実線で四角く囲われた領域は、自パーティ内の各ノード10Aで判定された状態である。従って、ノード10A-1がそなえるノード状態決定部13Aは、実線で四角く囲われた領域を、判定により更新する。

ノード状態決定部13Aは、他ノード10Aの最終更新情報を参照し、ハートビート(ノード状態情報T5又はT6)の到達の有無に応じてAlive又はSuspectの判定を行なう。また、ノード状態決定部13Aは、他ノード10Aについて、ノード状態管理情報T7における実線で角丸の四角で囲われた領域を参照し、Suspect、Down、又はZombieの判定を多数決等により行なう。

10

【0163】

なお、パーティ間ノード状態決定部103及びノード状態決定部13Aによる判定の基準は、第1実施形態において既述のものと同様であり、詳細な説明は省略する。

また、パーティ間ノード状態決定部103及びノード状態決定部13Aは、上述のように、ノード10Aの状態を判定すると、ノード状態管理情報T7を更新する。

具体的には、パーティ間ノード状態決定部103及びノード状態決定部13Aは、自ノード10A及び他ノード10Aの各々について判定した状態を、図19に例示するノード状態管理情報T7における“状態”の列に設定する。

【0164】

20

なお、パーティ間ノード状態決定部103及びノード状態決定部13Aによる上述した判定は、第1所定時間置きに判定対象の全ノード10Aについて一括で行なわれてもよいし、ノード10Aごとに異なるタイミングで、第1所定時間置きに行なわれてもよい。

〔2-2-4〕パーティ間送信処理部及び送信処理部

パーティ間送信処理部(グループ間送信処理部)104は、第1所定時間ごとに、パーティ間ノード状態決定部103が判定した複数の代表ノード10Aの各々の状態に関する代表ノード状態情報T5を、他の代表ノード10Aの各々へ送信する。

【0165】

具体的には、パーティ間送信処理部104は、パーティ管理情報T4及びノード状態管理情報T7を参照して、上述のように所定のルールに基づき、他のパーティの代表ノード10Aを特定する。そして、パーティ間送信処理部104は、ノード状態管理情報T7から、他の代表ノード10AのIPアドレス及びポート番号を取得し、代表ノード状態情報T5の宛先ノードを判定する。

30

【0166】

また、パーティ間送信処理部104は、ノード状態管理情報T7を参照して、全ノード10AについてのノードID、状態、IPアドレス、及びポート番号の情報から代表ノード状態情報T5を生成する。そして、パーティ間送信処理部104は、生成した代表ノード状態情報T5を、ハートビートとして他の代表ノード10Aの各々へ送信する。

また、パーティ間送信処理部104は、代表ノード状態情報T5の送信に加え、後述するパーティ管理部105によりパーティ管理情報T4が更新された場合には、パーティ管理情報T4(図14参照)をストレージシステム1内の全てのノード10Aへ通知する。なお、この通知は、ブロードキャスト等により行なわれてもよい。

40

【0167】

また、パーティ間送信処理部104は、パーティ管理情報T4が更新されたタイミングに限らず、パーティ管理情報T4を代表ノード状態情報T5とともにハートビートとして、他の代表ノード10Aへ送信してもよい。

送信処理部14Aは、送信用ノード状態情報T6を、自パーティにおける他ノード10Aの各々へ送信する。

【0168】

具体的には、送信処理部14Aは、パーティ管理情報T4及びノード状態管理情報T7

50

を参照して、自パーティ内の他ノード10Aを特定する。そして、送信処理部14Aは、ノード状態管理情報T7から、自パーティ内の他ノード10AのIPアドレス及びポート番号を取得し、ノード状態情報T6の宛先ノードを判定する。

また、送信処理部14Aは、ノード状態管理情報T7を参照して、自ノード10Aが判定した各ノード10AについてのノードID、状態、IPアドレス、及びポート番号の情報からノード状態情報T6を生成する。そして、送信処理部14Aは、生成したノード状態情報T6を、ハートビートとして自パーティ内の他ノード10Aの各々へ送信する。

#### 【0169】

また、送信処理部14Aは、ノード状態情報T6の送信に加え、上述のように、自ノード10Aの起動後、送信情報T3(図6参照)をストレージシステム1内の全てのノード10Aへブロードキャスト等により通知する。

10

なお、パーティ間受信処理部102が受信する代表ノード状態情報T5及びパーティ間送信処理部104が送信する代表ノード状態情報T5は、同様のデータ構造である。また、受信処理部12Aが受信するノード状態情報T6及び送信処理部14Aが送信するノード状態情報T6は、同様のデータ構造である。以下、便宜上、パーティ間送信処理部104が送信する代表ノード状態情報T5を送信用代表ノード状態情報(送信用代表状態情報)T5といい、送信処理部14Aが送信するノード状態情報T6を送信用ノード状態情報(送信用状態情報)T6という場合がある。

#### 【0170】

〔2-2-5〕パーティ管理部

20

次に、図20～図23を参照して、パーティ管理部105について説明する。

図20は、第2実施形態の一例としてのストレージシステム1にノード10Aが追加される例を示す図であり、図21は、図20に示すストレージシステム1におけるパーティの分割処理の一例を説明する図である。図22は、図21に示すストレージシステム1におけるノード10Aの削除処理及びパーティの統合処理の一例を説明する図である。図23は、第2実施形態の一例としてのストレージシステム1におけるパーティの分割処理の具体例を説明する図である。

#### 【0171】

なお、図20～図22に示す例においては、説明の簡略化のため、ノード10A間の接続状態のみを示し、スイッチ20の図示を省略している。

30

パーティ管理部(管理部)105は、自ノード10Aが属するパーティに関する管理を行なう。

具体的には、パーティ管理部105は、自パーティにおけるノード10Aの追加又は削除により、自パーティに属するノード10Aの数が所定の上限又は所定の下限を超えた場合、自パーティの分割又は統合を行なう。

#### 【0172】

例えば、ストレージシステム1の運用が開始されたとき等の初期状態において、パーティが1つ又は複数ある場合、ストレージシステム1の運用に応じてパーティにノード10Aが追加される場合がある。ノード10Aの追加により、パーティを構成するノード10Aの数が多くなると、パーティ内でのハートビートの通信によりノード10Aの処理負荷及びネットワークの負荷が高まり、ストレージシステム1の性能が低下する可能性がある。

40

#### 【0173】

そこで、パーティ管理部105は、自パーティにおけるノード10Aの数が予め決められた上限(第4所定値)を上回った場合、自パーティから複数のノード10Aを分割し、新たなパーティを作成する。

また、逆に、パーティ管理部105は、パーティを構成するノード数が下限(第5所定値)を下回った場合、パーティを統合する。パーティを統合する理由は、少数のノード10Aを含むパーティが多数あると、代表ノード10A間のハートビートの通信による代表ノード10Aの処理負荷及びネットワーク負荷が高まるためである。また、パーティ管理

50

情報 T 4 が肥大化し、パーティの管理に係る処理負荷が増大することも理由の一つである。

【 0 1 7 4 】

なお、予め定められた上限及び下限としては、ストレージシステム 1 の規模やポリシー等によって異なるが、例えば上限を数十～数百台程度とし、下限を数～数十台程度とすることができる。以下、説明の簡略化のため、上限を 5 台、下限を 2 台として説明する。

パーティ管理部 1 0 5 によるパーティの分割又は統合に伴うパーティ管理情報 T 4 の変更は、各パーティに所属する代表ノード 1 0 A が、自パーティのエントリについて行なうことができる。代表ノード 1 0 A がそなえるパーティ管理部 1 0 5 は、パーティ管理情報 T 4 を変更すると、パーティ間送信処理部 1 0 4 を介して、ハートビートに載せて全ノード 1 0 A へ伝達する。

10

【 0 1 7 5 】

なお、パーティ管理情報 T 4 は、ブロードキャスト等により全ノード 1 0 A へ伝達されてもよいし、代表ノード状態情報 T 5 とともにハートビートとして各代表ノード 1 0 A へ伝達されてもよい。パーティ管理情報 T 4 が各代表ノード 1 0 A へ伝達される場合、パーティ管理情報 T 4 を受け取った代表ノード 1 0 A は、自パーティのメンバノード 1 0 A へ転送することが好ましい。

【 0 1 7 6 】

以下、パーティ管理部 1 0 5 によるパーティの分割処理及び統合処理について説明する。

20

図 2 0 の紙面左上に示すように、ストレージシステム 1 が、パーティ A 及び B をそなえる場合を例に挙げて説明する。なお、パーティ A は、ノード ID “ 1 ”、“ 3 ”、“ 5 ”、“ 7 ”、及び “ 9 ” の 5 つのノード 1 0 A を有し、パーティ B は、ノード ID “ 1 1 ”、“ 1 3 ”、“ 1 5 ”、“ 1 7 ”、及び “ 1 9 ” の 5 つのノード 1 0 A をそなえるものとする。また、図 2 0 の紙面右側に示すように、パーティ管理情報 T 4 には、パーティ ID “ A ” にノード ID “ 1 ~ 1 0 ” が、パーティ ID “ B ” にノード ID “ 1 1 ~ 2 0 ” が、それぞれ対応付けられているものとする。

【 0 1 7 7 】

なお、パーティ A 及び B の代表ノード 1 0 A は、それぞれノード ID “ 1 ” 及び “ 1 1 ” のノード 1 0 A (以下、代表ノード 1 0 A - 1 及び 1 0 A - 1 1 という)である。

30

以上の例において、パーティ A にノード ID “ 8 ” のノード 1 0 A が追加される場合を想定する(図 2 0 の紙面左下及び図 2 1 の紙面左上参照)。この場合、パーティ A には、6 つのノードが含まれる。なお、ノード ID “ 8 ” は、パーティ A に対応付けられたノード ID の範囲内であるため、パーティ管理情報 T 4 に変更はない。

【 0 1 7 8 】

代表ノード 1 0 A - 1 がそなえるパーティ管理部 1 0 5 は、自パーティ A に属するノード 1 0 A の数が上限である 5 つを超えるため、パーティ A を分割する。

図 2 1 の紙面左下に示すように、代表ノード 1 0 A - 1 のパーティ管理部 1 0 5 は、パーティ管理情報 T 4 及びノード状態管理情報 T 7 を参照して、パーティ A をノード数が 2 分の 1 になるように分割する。例えば、パーティ管理部 1 0 5 は、パーティ A に属するノード ID のうち、ノード ID が小さい順に 3 つのノード 1 0 A をパーティ A に残し、それ以外の 3 つのノードをパーティ C として分割する。つまり、パーティ管理部 1 0 5 は、パーティ A を、ノード ID “ 1 ”、“ 3 ”、及び “ 5 ” をそなえる新たなパーティ A と、ノード ID “ 7 ” ~ “ 9 ” をそなえるパーティ C とに分割する。

40

【 0 1 7 9 】

なお、パーティ管理部 1 0 5 は、パーティの分割において、ノード数を 2 分の 1 にする際に余りが出る場合、余りのノード 1 0 A を分割に係る 2 つのパーティのいずれかに割り振る。

代表ノード 1 0 A - 1 のパーティ管理部 1 0 5 は、パーティ A を分割すると、パーティ管理情報 T 4 のパーティ ID “ A ” のエントリにおいて、ノード ID を “ 1 ~ 5 ” に設定

50



し、バージョン番号を“ 2 ”に変更する。また、代表ノード 10A - 1 のパーティ管理部 105 は、パーティ管理情報 T4 にパーティ ID “ C ” のエントリを追加し、ノード ID “ 6 ~ 10 ”、バージョン番号 “ 1 ” を対応付ける。

【 0180 】

そして、代表ノード 10A - 1 のパーティ管理部 105 は、変更したパーティ管理情報 T4 を、パーティ間送信処理部 104 を介して全ノード 10A へ通知する。

なお、ノード 10A - 1 は、新たなパーティ A において、引き続き代表ノード 10A としてパーティ管理情報 T4 のパーティ ID “ A ” のエントリの管理を行なう。一方、パーティ C では、ノード ID “ 7 ” ~ “ 9 ” のノード 10A で、代表ノード 10A を選出する所定のルールが適用され、例えばノード ID “ 7 ” のノード 10A (以下、代表ノード 10A - 7 という) が、代表ノード 10A になる。代表ノード 10A - 7 は、代表ノード 10A - 1、10A - 11 とともに、代表ノード 10A 間でのハートビートの通信を行なうとともに、パーティ C のエントリを管理する。

10

【 0181 】

以上のように、代表ノード 10A がそなえるパーティ管理部 105 により、パーティの分割処理が行なわれる。

次いで、図 22 の紙面左上に示すように、パーティ A 内のノード ID “ 3 ” 及び “ 5 ” のノード 10A が障害等の発生により停止した場合を想定する。

代表ノード 10A - 1 がそなえるパーティ管理部 105 は、自パーティ A に属するノード 10A の数が、ノード 10A の停止に伴い下限である 2 つを超える (下回る) ため、パーティ A を他のパーティと統合する。

20

【 0182 】

図 22 の紙面左下に示すように、代表ノード 10A - 1 のパーティ管理部 105 は、パーティ管理情報 T4 及びノード状態管理情報 T7 を参照して、パーティ A と統合する他のパーティを決定する。パーティ A と統合する他のパーティとしては、例えばノード数が最も少ないパーティが挙げられる。この場合、代表ノード 10A - 1 のパーティ管理部 105 は、自パーティ A 以外でノード数が最も少ないパーティ C を統合対象のパーティに決定する。

【 0183 】

代表ノード 10A - 1 のパーティ管理部 105 は、統合対象のパーティを決定すると、パーティ管理情報 T4 のパーティ ID “ A ” のエントリにおいて、ノード ID をパーティ C とマージさせて、“ 1 ~ 10 ” に設定し、バージョン番号を “ 3 ” に変更する。また、代表ノード 10A - 7 のパーティ管理部 105 は、パーティ管理情報 T4 からパーティ ID “ C ” のエントリを削除する。

30

【 0184 】

そして、代表ノード 10A - 1 のパーティ管理部 105 は、変更したパーティ管理情報 T4 を、パーティ間送信処理部 104 を介して全ノード 10A へ通知する。

なお、ノード 10A - 1 は、新たなパーティ A において、引き続き代表ノード 10A としてパーティ管理情報 T4 のパーティ ID “ A ” のエントリの管理を行なう。一方、パーティ C の代表ノード 10A - 7 は、新たなパーティ A において代表ノード 10A を選出する所定のルールの敗者であるので、メンバノード 10A - 7 に降格する。

40

【 0185 】

また、図 22 に示す例では、パーティ A とパーティ C とが統合されたため、ノード ID も “ 1 ~ 5 ” と “ 6 ~ 10 ” とがマージされて “ 1 ~ 10 ” になった。しかし、パーティ管理情報 T4 の状態によっては、統合する 2 つのパーティのノード ID の範囲が離れ、間に存在するノード ID が他のパーティを構成する場合も考えられる。このような場合、統合後のパーティに属するノード ID は、1 つの範囲ではなく、複数の範囲で又は 1 つずつ設定されてもよい。

【 0186 】

以上のように、代表ノード 10A がそなえるパーティ管理部 105 により、パーティの

50

統合処理が行なわれる。

なお、代表ノード10Aのパーティ管理部105は、所定時間ごとに、自パーティのノード10Aの数が上限に達したか否か、及び下限を下回ったか否かを判定することができる。

【0187】

また、代表ノード10Aのパーティ管理部105は、ストレージシステム1に追加された新規ノード10Aから送信される送信情報T3を受信したことを契機に、自パーティのノード10Aの数が上限に達したか否かを判定してもよい。

さらに、代表ノード10Aのパーティ管理部105は、自パーティ内のノード10Aに障害等が発生し、当該ノード10Aのリカバリ処理が完了したことを契機に、自パーティのノード10Aの数が下限を下回ったか否かを判定してもよい。

10

【0188】

上述したパーティ管理部105の説明では、パーティ管理部105は、パーティの分割及び統合に係るノード10Aの選定を、ノードIDの値に基づいて行なうものとして説明した。しかし、ストレージシステム1において、ノード10A間のハートビートの通信は、ノード10A間の距離に応じたレイテンシやパケットロスの影響を受ける。

そこで、パーティ管理部105は、以下で説明するように、パーティの分割及び統合に係るノード10Aの選定を、例えばノード10Aが接続されるスイッチ20に基づいて行なうことが好ましい。なお、以下の説明は、管理者等による、ストレージシステム1の運用開始前のパーティの初期設定の際や、運用中においてパーティの構成が複雑化したこと

20

【0189】

一例として、パーティの初期設定の際に、1つのスイッチ20に接続されるノード10A群が同じパーティに設定されることが考えられる。図23の紙面上側に示す例では、スイッチ20に、ノード10A-1～10A-4が接続され、これらノード10A-1～10A-4が1つのパーティを構成する。なお、スイッチ20のポート数は4つであるものとする。

【0190】

ストレージシステム1へのノード10A-5及び10A-6の追加に伴い、ノード10Aの数が1つのスイッチ20に収まらなくなると、作業等により、スイッチ20の増設及びトポロジの調整が行なわれる。例えば、図23の紙面下側に示すように、スイッチ20-1にノード10A-1～10A-3が接続され、スイッチ20-2にノード10A-4～10A-6が接続される。また、スイッチ20-1及び20-2間が接続される。

30

【0191】

代表ノード10Aのパーティ管理部105は、ノード10A-5及び10A-6が追加されると(ノード10Aが図23の紙面下側に示す接続状態になると)、ノード10A及びスイッチ20の接続関係に関する情報を取得する。例えば、パーティ管理部105は、スイッチ20が保持する各ポートの接続先の情報等を取得することで、ノード10A及びスイッチ20の接続関係に関する情報を取得(推定)することができる。なお、スイッチ20からの接続先の情報等の取得は、既知の種々の手法により行なうことが可能であり、その詳細な説明は省略する。また、パーティ管理部105は、作業等から、入出力部10eを介してノード10A及びスイッチ20の接続関係に関する情報を入力されてもよい。

40

【0192】

そして、パーティ管理部105は、取得したノード10A及びスイッチ20の接続関係から、例えば、パーティを、スイッチ20-1に接続されるノード10A群と、スイッチ20-2に接続されるノード10A群とに分割する。

このように、パーティ管理部105は、自パーティにおけるノード10A及びスイッチ20の物理的な接続関係に関する情報に基づいて、パーティから分割するノード10Aを決定することができる。

50

## 【0193】

なお、パーティ管理部105は、ノード10A及びスイッチ20の接続関係に関する情報として、代表ノード10Aからパーティ内の他ノード10Aの各々までのホップ数を検出してもよい。これは、ホップ数が近いノード10A同士は、同じスイッチ20に接続されている可能性が高いと推測できるからである。

ここまで、図23を参照してパーティ管理部105によるパーティの分割処理について説明したが、パーティ管理部105によるパーティの統合処理についても同様である。

## 【0194】

すなわち、パーティ管理部105は、自パーティと統合する他のパーティとして、ノード数が少ないパーティを選択するのではなく、ノード10A及びスイッチ20の接続関係に基づき選択してもよい。

## 〔2-3〕動作例

次に、図24～図26を参照して、上述の如く構成された第2実施形態の一例としてのノード10Aによる動作例を説明する。図24は、第2実施形態の一例としての代表ノード10Aによる他の代表ノード10Aの状態を判定する動作例を説明するフローチャートである。図25は、ノード10Aによるパーティ内の他ノード10Aが停止した場合の動作例を説明するフローチャートである。図26は、ノード10Aによるパーティの分割処理及び統合処理の動作例を説明するフローチャートである。

## 【0195】

## 〔2-3-1〕代表ノードによる他の代表ノードの状態を判定する動作例

はじめに、図24を参照して、代表ノード10Aによる他の代表ノード10Aの状態を判定する動作例を説明する。

なお、図24に示すステップS41～S55の処理は、代表ノード10Aの各々において、パーティ間ノード状態決定部103により他の一の代表ノード10Aの状態が判定される際に行なわれる処理である。従って、ステップS41～S55の処理は、各代表ノード10Aのパーティ間ノード状態決定部103により、他の代表ノード10Aの各々について、定期的（第1所定時間ごと）に実行される。

## 【0196】

また、図24に示すステップS41～S49、S52、及びS53の処理は、図12に示すステップS11～S19、S22、及びS23の処理と比較して、判定対象のノード10（10A）が代表ノード10Aである点が異なる。以下、ステップS41～S49、S52、及びS53の処理の説明において、図12に示すステップS11～S19、S22、及びS23の処理と同様な部分の詳細は省略する。

## 【0197】

図24に示すように、パーティ間ノード状態決定部103により、ノード状態管理情報T7内の“状態”が参照され、判定対象の代表ノード10Aについて直前に判定した状態がどの状態であるかが判定される（ステップS41、S46、S49）。

判定対象の代表ノード10Aについて直前に判定した状態がAliveである場合（ステップS41のYesルート）、処理がステップS42に移行する。ステップS42では、パーティ間ノード状態決定部103により、判定対象の代表ノード10Aからのハートビートの不達時間が閾値を超えたか否かが判定される。

## 【0198】

ハートビートの不達時間が閾値を超えた場合（ステップS42のYesルート）、パーティ間ノード状態決定部103により、判定対象の代表ノード10Aの状態がSuspectと判定され（ステップS43）、処理が終了する。このとき、パーティ間ノード状態決定部103は、判定対象の代表ノード10Aについて、ノード状態管理情報T7内の“状態”にSuspectを設定する。そして、パーティ間ノード状態決定部103は、次の判定対象の代表ノード10Aがある場合、次の判定対象の代表ノード10Aに係る状態の判定処理に移行する。

## 【0199】

一方、ステップS 4 2において、パーティ間ノード状態決定部1 0 3により、ハートビートの不達時間が閾値を超えていないと判定された場合(ステップS 4 2のN o ルート)、処理がステップS 4 4に移行する。ステップS 4 4では、パーティ間ノード状態決定部1 0 3により、判定対象の代表ノード1 0 Aの状態が、過半数(第1所定値)の代表ノード1 0 AからS u s p e c tと判定されたか否かが判定される。又は、パーティ間ノード状態決定部1 0 3により、判定対象の代表ノード1 0 Aの状態が、複数の代表ノード1 0 Aのいずれかの代表ノード1 0 AによりD o w nと判定されたか否かが判定される。

【0 2 0 0】

判定対象の代表ノード1 0 Aの状態が、過半数の代表ノード1 0 AからS u s p e c tと判定されておらず、いずれかの代表ノード1 0 AによりD o w nとも判定されていない場合(ステップS 4 4のN o ルート)、代表ノード1 0 Aに対する処理が終了する。一方、判定対象の代表ノード1 0 Aの状態が、過半数の代表ノード1 0 AからS u s p e c tと判定された又はいずれかの代表ノード1 0 AによりD o w nと判定された場合(ステップS 4 4のY e s ルート)、処理がステップS 4 5に移行する。

【0 2 0 1】

ステップS 4 5では、パーティ間ノード状態決定部1 0 3により、判定対象の代表ノード1 0 Aの状態がD o w nと判定され、処理が終了する。このとき、パーティ間ノード状態決定部1 0 3は、判定対象の代表ノード1 0 Aについて、ノード状態管理情報T 7内の“状態”にD o w nを設定する。

また、判定対象の代表ノード1 0 Aについて直前に判定した状態がS u s p e c tである場合(ステップS 4 1のN o ルートからステップS 4 6のY e s ルート)、処理がステップS 4 7に移行する。ステップS 4 7では、パーティ間ノード状態決定部1 0 3により、判定対象の代表ノード1 0 Aから新たなハートビートが受信されたか否かが判定される。

【0 2 0 2】

新たなハートビートが受信されていない場合(ステップS 4 7のN o ルート)、処理がステップS 4 4に移行する。一方、新たなハートビートが受信された場合(ステップS 4 7のY e s ルート)、パーティ間ノード状態決定部1 0 3により、判定対象の代表ノード1 0 Aの状態がA l i v eと判定され(ステップS 4 8)、処理が終了する。このとき、パーティ間ノード状態決定部1 0 3は、判定対象の代表ノード1 0 Aについて、ノード状態管理情報T 7内の“状態”にA l i v eを設定する。

【0 2 0 3】

判定対象の代表ノード1 0 Aについて直前に判定した状態がD o w nである場合(ステップS 4 1のN o ルート、ステップS 4 6のN o ルートからステップS 4 9のY e s ルート)、処理がステップS 5 0に移行する。ステップS 5 0では、パーティ間ノード状態決定部1 0 3により、判定対象の代表ノード1 0 Aの状態が第2所定値の数の代表ノード1 0 A(例えば全代表ノード1 0 A)によりD o w nと判定されたか否かが判定される。

【0 2 0 4】

全代表ノード1 0 AによりD o w nと判定されていない場合(ステップS 5 0のN o ルート)、判定対象の代表ノード1 0 Aに対する処理が終了する。一方、全代表ノード1 0 AによりD o w nと判定された場合(ステップS 5 0のY e s ルート)、処理がステップS 5 4に移行する。ステップS 5 4では、ノード状態決定部1 3により、該当パーティに他ノード1 0 Aが生存しているか否か、つまり該当パーティにA l i v eと判定された他ノード1 0 Aが存在するか否かが判定される。

【0 2 0 5】

ステップS 5 4において、該当パーティに生存している他ノード1 0 Aが存在しないと判定された場合(ステップS 5 4のN o ルート)、処理がステップS 5 1に移行する。ステップS 5 1では、パーティ間ノード状態決定部1 0 3により、該当パーティに所属する全ノード1 0 Aの状態がZ o m b i eと判定される。また、自ノード1 0 Aが保持するデータが該当パーティに所属するいずれかのノード1 0 Aが保持するデータに関連する場合

、リカバリ処理部 15 により、該当ノード 10 A に対するリカバリ処理が実行され、処理が終了する。このとき、ノード状態決定部 13 は、該当パーティに所属する全ノード 10 A について、ノード状態管理情報 T7 内の“状態”に Z o m b i e を設定する。

【0206】

一方、ステップ S54 において、該当パーティに生存している他ノード 10 A が存在すると判定された場合（ステップ S54 の Y e s ルート）、処理がステップ S55 に移行する。ステップ S55 では、パーティ間ノード状態決定部 103 により、生存している他ノード 10 A のうちの次点のノード 10 A が新たな判定対象の代表ノード 10 A と判定され、処理が終了する。なお、この判定は、代表ノード 10 A を選定する所定のルール（例えばノード ID が最も小さいノード 10 A）に基づいて行なわれる。代表ノード 10 A は、ステップ S55 において選定した新たな判定対象の代表ノード 10 A に対して、次回以降のハートビートの通信を行なう。

【0207】

判定対象の代表ノード 10 A について直前に判定した状態が Z o m b i e である場合（ステップ S41 の N o ルート、ステップ S46 の N o ルートからステップ S49 の N o ルート）、処理がステップ S52 に移行する。ステップ S52 では、パーティ間ノード状態決定部 103 により、判定対象の代表ノード 10 A についてリカバリ処理が完了したか否かが判定される。リカバリ処理が完了していない場合、判定対象の代表ノード 10 A に対する処理が終了する。一方、リカバリ処理が完了した場合（ステップ S52 の Y e s ルート）、パーティ間ノード状態決定部 103 により、ノード状態管理情報 T7 から、判定対象の代表ノード 10 A に関する情報が削除され（ステップ S53）、処理が終了する。

【0208】

以上のように、代表ノード 10 A により、他の一の代表ノード 10 A の状態の判定処理が行なわれる。

〔2-3-2〕ノードによるパーティ内の他ノードが停止した場合の動作例

次に、図 25 を参照して、ノード 10 A によるパーティ内の他ノード 10 A が停止した場合の動作例を説明する。

【0209】

なお、図 25 に示すステップ S61 ~ S63 の処理は、メンバノード 10 A の各々において、ノード状態決定部 13 A により自パーティの代表ノード 10 A の状態が判定される際に行なわれる処理である。従って、ステップ S61 ~ S63 の処理は、メンバノード 10 A のノード状態決定部 13 A により、自パーティの代表ノード 10 A について、定期的（第 1 所定時間ごと）に実行される。

【0210】

図 25 に示すように、メンバノード 10 A による自パーティ内の他ノード 10 A の状態の判定により、自パーティの代表ノード 10 A が停止したか否かが判定される（ステップ S61）。

自パーティの代表ノード 10 A が停止していない場合（ステップ S61 の N o ルート）、処理が終了する。ノード状態決定部 13 A は、次の判定対象のメンバノード 10 A がある場合、次の判定対象のメンバノード 10 A に係る状態の判定処理に移行する。

【0211】

一方、自パーティの代表ノード 10 A が停止した場合（ステップ S61 の Y e s ルート）、自ノード 10 A が代表ノード 10 A になるか否かが判定される（ステップ S62）。なお、この判定は、代表ノード 10 A を選定する所定のルールに基づいて行なわれる。

ノード 10 A により、ステップ S62 において自ノード 10 A が代表ノード 10 A になると判定された場合（ステップ S62 の Y e s ルート）、他のパーティの代表ノード 10 A の各々との間のハートビートの通信が開始され（ステップ S63）、処理が終了する。

【0212】

一方、ノード 10 A により、ステップ S62 において自ノード 10 A が代表ノード 10 A にならないと判定された場合（ステップ S62 の N o ルート）、処理が終了する。

以上のように、ノード10Aによるパーティ内の他ノード10Aが停止した場合の処理が終了する。

〔2-3-3〕ノードによるパーティの分割処理及び統合処理の動作例

次に、図26を参照して、ノード10Aによるパーティの分割処理及び統合処理の動作例を説明する。

【0213】

図26に示すように、代表ノード10Aのパーティ管理部105により、例えば所定時間ごとに、自パーティのノード10Aの数が上限を上回ったか否かが判定される(ステップS71)。

ノード10Aの数が上限を上回った場合(ステップS71のYesルート)、代表ノード10Aのパーティ管理部105により、パーティが2つに分割され、パーティ管理情報T4が更新される(ステップS72)。なお、パーティ管理部105は、上述のように、自パーティをノード数が2分の1になるように分割し、余りが出る場合、余りのノード10Aを分割に係る2つのパーティのいずれかに割り振る。また、パーティ管理部105は、自パーティのノード10Aを分割後のいずれのパーティに割り当てるかを、ノードID、ノード10A及びスイッチ20の接続関係に基づき決定する。

【0214】

ステップS71において、ノード10Aの数が上限以下である場合(ステップS71のNoルート)、パーティ管理部105により、自パーティのノード10Aの数が下限未満であるか否かが判定される(ステップS73)。

ノード10Aの数が下限未満である場合(ステップS73のYesルート)、代表ノード10Aのパーティ管理部105により、自パーティと他のパーティとの統合が行なわれる。具体的には、代表ノード10Aのパーティ管理部105は、自パーティと統合する他のパーティを、ノードID、ノード10A及びスイッチ20の接続関係、又は他のパーティの代表ノード10Aまでのホップ数等に基づき決定する。

【0215】

そして、代表ノード10Aのパーティ管理部105により、統合後に自ノード10Aが代表ノード10Aになるか否かが判定される(ステップS74)。具体的には、代表ノード10Aのパーティ管理部105は、パーティ管理情報T4及びノード状態管理情報T7を参照して、自ノード10AのノードIDと決定した他のパーティの代表ノード10AのノードIDとを比較する。そして、パーティ管理部105は、自ノード10AのノードIDが他のパーティの代表ノード10AのノードIDよりも小さいか否かを判定する。

【0216】

統合後に自ノード10Aが代表ノード10Aにならないと判定された場合(ステップS74のNoルート)、代表ノード10Aのパーティ管理部105により、パーティ管理情報T4の自パーティのエントリが削除され(ステップS75)、処理が終了する。

一方、統合後に自ノード10Aが代表ノード10Aになると判定された場合(ステップS74のYesルート)、代表ノード10Aのパーティ管理部105により、自パーティと他のパーティとが統合される。具体的には、代表ノード10Aのパーティ管理部105により、パーティ管理情報T4の自パーティのエントリのノードIDに、統合する他のパーティのノードIDがマージされて、パーティ管理情報T4が更新される(ステップS76)。そして、パーティ管理部105による処理が終了する。

【0217】

なお、上述のように、ステップS71の処理は、ストレージシステム1に追加された新規ノード10Aから送信される送信情報T3を受信したことを契機に開始されてもよい。

また、ステップS73の処理は、自パーティ内のノード10Aに障害等が発生し、当該ノード10Aのリカバリ処理が完了したことを契機に開始されてもよい。

さらに、ステップS71及びS72の処理と、ステップS73～S76の処理とは、互いに独立して実行されてもよいし、処理順序が変更されてもよい。

【0218】

## 〔 2 - 4 〕 第 2 実施形態のまとめ

上述のように、第 2 実施形態の一例としてのノード 10A によれば、第 1 実施形態に係るノード 10 と同様の効果を奏することができる。

また、第 2 実施形態の一例としてのノード 10A によれば、複数のノード 10A が複数のパーティに分割される。そして、各パーティの代表ノード 10A の各々において、パーティ間受信処理部 102 は、他のパーティの各々における他の代表ノード 10A から、代表ノード状態情報 T5 を受信する。また、パーティ間ノード状態決定部 103 は、代表ノード状態情報 T5 に基づいて、複数の代表ノード 10A の各々の状態を判定する。さらに、パーティ間送信処理部 104 は、パーティ間ノード状態決定部 103 が判定した複数の代表ノード 10A の各々の状態に関する送信用代表ノード状態情報 T5 を、他の代表ノード 10A の各々へ送信する。

10

## 【 0219 】

さらに、複数のノード 10A の各々において、送信処理部 14A は、送信用ノード状態情報 T6 を、自パーティにおける他ノード 10A の各々へ送信する。また、ノード状態決定部 13A は、受信処理部 12A が自パーティにおける他ノード 10A の各々から受信したノード状態情報 T6 に基づいて、自ノード 10A におけるノード 10A の各々の状態を判定する。

## 【 0220 】

これにより、メンバノード 10A により、自パーティ内のノード 10A の状態が判定され、代表ノード 10A により、パーティ間（代表ノード 10A 間）の状態が判定される。

20

従って、ストレージシステム 1 においてノード 10A の数が増大した場合でも、ノード間でやり取りされる情報の直接の送信先を絞ることができるため、ハートビートの送受信のコストの増大を抑えることができる。

## 【 0221 】

つまり、ストレージシステム 1 は、全ノード 10A による完全メッシュのハートビートの通信を行なうよりも、ストレージシステム 1 における通信負荷及び処理負荷を低減させることができる。

また、各代表ノード 10A において、パーティ管理部 105 は、自パーティにおけるノード 10A の数が第 4 所定値を超えた場合、自パーティから、複数のノード 10A を分割して新たなパーティを作成する。

30

## 【 0222 】

これにより、パーティ内でのハートビートの通信によるノード 10A の処理負荷及びネットワークの負荷に起因した、ストレージシステム 1 の性能低下を抑止することができる。

さらに、パーティ管理部 105 は、自パーティにおけるノード 10A 及びスイッチ 20 の接続関係に関する情報に基づいて、自パーティから分割する複数のノード 10A を決定する。

## 【 0223 】

これにより、ノード 10A 間の距離に応じたレイテンシやパケットロスの影響を抑止することができる。

40

また、パーティ管理部 105 は、自パーティにおけるノード 10A の数が第 5 所定値未満の場合、自パーティと他のパーティのうちのいずれかのパーティとを統合する。

これにより、多数の代表ノード 10A 間のハートビートの通信による代表ノード 10A の処理負荷及びネットワーク負荷に起因した、ストレージシステム 1 の性能低下を抑止することができる。

## 【 0224 】

## 〔 3 〕 その他

以上、本発明の好ましい実施形態について詳述したが、本発明は、係る特定の実施形態に限定されるものではなく、本発明の趣旨を逸脱しない範囲内において、種々の変形、変更して実施することができる。

50

例えば、第１及び第２実施形態に係るストレージシステム１がそなえるノード１０及び１０Ａ、並びにスイッチ２０の構成及び台数は、上述したものに限定されず、任意の構成及び台数とすることができる。

【０２２５】

また、第１及び第２実施形態においては、ストレージシステム１がそなえるノード１０及び１０Ａにおける処理について説明したが、これに限定されるものではない。ノード１０及び１０Ａは、ストレージ装置のほか、情報に対する処理を行なうサーバ等の情報処理装置であってもよく、ストレージシステム１は、複数の情報処理装置をそなえる情報処理システムであってもよい。

【０２２６】

また、第１及び第２実施形態においては、ノード１０及び１０Ａは、例えばストレージシステム１によるサービスの提供に用いられるＩＰラインを介してハートビートを行なうものとして説明したが、これに限定されるものではない。例えば、ノード１０及び１０Ａは、ＬＡＮケーブル等の専用の制御線を介して相互に接続され、専用線を用いてハートビートを行なってもよい。これにより、ＩＰラインのネットワークの負荷を軽減させることができる。なお、ノード１０及び１０Ａは、ＩＰラインを用いる場合、ノード１０及び１０Ａ間のバスの障害検出を行なうことができるため、専用線を用いるよりも監視範囲を拡張することができる。

【０２２７】

さらに、第２実施形態においては、ノード１０Ａは、１段のパーティを構成するものとして説明したが、これに限定されるものではなく、多段のパーティを構成してもよい。つまり、代表ノード１０Ａが数百～数千台等の多数存在する場合、代表ノード１０Ａを複数の上位パーティに分割し、上位パーティ間でハートビートの通信を行なうとともに、各上位パーティ内で、代表ノード１０Ａ間のハートビートの通信を行なってもよい。

【０２２８】

また、第２実施形態においては、全てのノード１０Ａが代表ノード１０Ａになる可能性があったが、これに限定されるものではない。例えば、ノード１０Ａ間で、特定の処理を行なうノード１０Ａ等の処理負荷を増やしたくないノード１０Ａについて、代表ノード１０Ａの候補から除外するＮＧリストを共有してもよい。この場合、各ノード１０Ａは、ＮＧリストに含まれるノード１０Ａについては代表ノード１０Ａに選出しないようにする。

【０２２９】

さらに、第１及び第２実施形態に係るノード１０及び１０Ａがそなえる各機能は、適宜省略してもよく、分割又は統合してもよい。例えば、第２実施形態に係るパーティ間受信処理部１０２及び受信処理部１２Ａを統合し、１つの受信処理部としてもよく、パーティ間送信処理部１０４及び送信処理部１４Ａを統合し、１つの送信処理部としてもよい。また、パーティ間ノード状態決定部１０３及びノード状態決定部１３Ａを統合し、１つのノード状態決定部（判定部）としてもよい。

【０２３０】

また、第２実施形態に係るノード１０Ａは、代表ノード１０Ａとして動作する際に、パーティ間受信処理部１０２、パーティ間ノード状態決定部１０３、パーティ間送信処理部１０４、パーティ管理部１０５の機能を実行する。従って、ノード１０Ａは、代表ノード１０Ａとして動作しない場合（例えば上記ＮＧリストに自ノード１０が登録される場合等）には、これらの機能を無効化又は省略してもよい。

【０２３１】

さらに、第１及び第２実施形態の一例における各処理フローのステップの実行順序を、適宜変更してもよい。

また、第１実施形態に係るノード１０及び第２実施形態に係るノード１０Ａの各種機能の全部もしくは一部は、コンピュータ（ＣＰＵ，情報処理装置，各種端末を含む）が所定のプログラムを実行することによって実現されてもよい。

【０２３２】



そのプログラムは、例えばフレキシブルディスク、ＣＤ、ＤＶＤ、ブルーレイディスク等のコンピュータ読取可能な記録媒体（例えば図２に示す記録媒体１０ｈ）に記録された形態で提供される。なお、ＣＤとしては、ＣＤ－ＲＯＭ、ＣＤ－Ｒ、ＣＤ－ＲＷ等が挙げられる。また、ＤＶＤとしては、ＤＶＤ－ＲＯＭ、ＤＶＤ－ＲＡＭ、ＤＶＤ－Ｒ、ＤＶＤ－ＲＷ、ＤＶＤ＋Ｒ、ＤＶＤ＋ＲＷ等が挙げられる。この場合、コンピュータはその記録媒体からプログラムを読み取って内部記憶装置または外部記憶装置に転送し格納して用いる。

#### 【０２３３】

ここで、コンピュータとは、ハードウェアとＯＳ（Operating System）とを含む概念であり、ＯＳの制御の下で動作するハードウェアを意味している。また、ＯＳが不要でアプリケーションプログラム単独でハードウェアを動作させるような場合には、そのハードウェア自体がコンピュータに相当する。ハードウェアは、少なくとも、ＣＰＵ等のマイクロプロセッサと、記録媒体に記録されたコンピュータプログラムを読み取る手段とをそなえている。上記プログラムは、上述のようなコンピュータに、第１実施形態に係るノード１０又は第２実施形態に係るノード１０Ａの各種機能を実現させるプログラムコードを含んでいる。また、その機能の一部は、アプリケーションプログラムではなくＯＳによって実現されてもよい。

10

#### 【０２３４】

##### 〔４〕付記

以上の第１及び第２実施形態に関し、更に以下の付記を開示する。

20

##### （付記１）

相互に接続される複数の情報処理装置を有し、前記複数の情報処理装置間で通信を行なう情報処理システムにおいて、

前記複数の情報処理装置の各々が、

前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信する受信処理部と、

前記受信処理部が前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定する判定部と、

前記判定部が判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信する送信処理部と、を有することを特徴とする、情報処理システム。

30

#### 【０２３５】

##### （付記２）

前記判定部は、

前記受信処理部が受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、前記他の情報処理装置の各々からの前記状態情報の受信状況とに基づいて、前記複数の情報処理装置の各々の状態を判定し、

前記送信処理部は、

第１所定時間ごとに、前記送信用状態情報を、前記他の情報処理装置の各々へ送信することを特徴とする、付記１記載の情報処理システム。

40

#### 【０２３６】

##### （付記３）

前記判定部は、

前記受信処理部が受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、前記判定部が含まれる自情報処理装置の状態に関する状態情報に関する自己状態情報とに基づいて、前記複数の情報処理装置の各々の状態を判定することを特徴とする、付記１記載の情報処理システム。

#### 【０２３７】

##### （付記４）

50

前記判定部は、

前記第 1 所定時間以上の時間である第 2 所定時間内に前記状態情報を受信しなかった他の情報処理装置の状態を、停止の可能性を示す第 1 状態と判定し、

前記受信処理部が受信した前記状態情報に基づいて、第 1 所定数以上の前記複数の情報処理装置で前記第 1 状態であると判定された情報処理装置の状態、又は、前記他の情報処理装置の少なくとも 1 つから停止を示す第 2 状態であると判定された情報処理装置の状態を、前記第 2 状態と判定することを特徴とする、付記 2 記載の情報処理システム。

【 0 2 3 8 】

( 付記 5 )

前記判定部は、

前記受信処理部が受信した前記状態情報に基づいて、前記第 1 所定数以上の数である第 2 所定数以上の前記複数の情報処理装置で前記第 2 状態であると判定された情報処理装置を、リカバリ処理中を示す第 3 状態と判定し、

前記複数の情報処理装置のうちの 1 以上の情報処理装置はさらに、

前記判定部が前記第 3 状態と判定した情報処理装置に対して、リカバリ処理を実行するリカバリ処理部を有することを特徴とする、付記 4 記載の情報処理システム。

【 0 2 3 9 】

( 付記 6 )

前記判定部は、

前記自情報処理装置に所定の障害が発生した場合、前記自情報処理装置の状態を前記第 2 状態と判定し、

前記第 2 所定時間内に第 3 所定数以上の前記他の情報処理装置から前記状態情報を受信しなかった場合、前記自情報処理装置の状態を、前記他の情報処理装置から切り離されたことを示す第 4 状態と判定し、

前記複数の情報処理装置の各々はさらに、

前記自情報処理装置に所定の障害が発生し、前記判定部が前記自情報処理装置の状態を前記第 2 状態と判定した場合、又は、前記判定部が前記自情報処理装置の状態を前期第 4 状態と判定した場合、前記自情報処理装置を停止させる処理を行なう停止処理部を有することを特徴とする、付記 4 又は付記 5 記載の情報処理システム。

【 0 2 4 0 】

( 付記 7 )

前記複数の情報処理装置が複数のグループに分割され、

前記複数のグループの各々における代表情報処理装置はさらに、

前記複数のグループのうちの自グループ以外の他のグループの各々における他の代表情報処理装置から、前記他の代表情報処理装置により判定された前記複数のグループの代表情報処理装置の各々の状態に関する代表状態情報を受信するグループ間受信処理部と、

前記グループ間受信処理部が前記他の代表情報処理装置の各々から受信した前記代表状態情報に基づいて、前記複数の代表情報処理装置の各々の状態を判定するグループ間判定部と、

前記グループ間判定部が判定した前記複数の代表情報処理装置の各々の状態に関する送信用代表状態情報を、前記他の代表情報処理装置の各々へ送信するグループ間送信処理部と、を有し、

前記複数の情報処理装置の各々において、

前記送信処理部は、

前記送信用状態情報を、前記自グループにおける他の情報処理装置の各々へ送信し、

前記判定部は、

前記受信処理部が前記自グループにおける他の情報処理装置の各々から受信した前記状態情報に基づいて、前記自グループにおける情報処理装置の各々の状態を判定することを特徴とする、付記 1 ~ 6 のいずれか 1 項記載の情報処理システム。

【 0 2 4 1 】

10

20

30

40

50

(付記 8)

前記複数のグループの各々における代表情報処理装置はさらに、

前記自グループにおける情報処理装置の数が第 4 所定値を超えた場合、前記自グループから、複数の情報処理装置を分割して新たなグループを作成する管理部を有することを特徴とする、付記 7 記載の情報処理システム。

【 0 2 4 2 】

(付記 9)

前記情報処理システムはさらに、

前記複数の情報処理装置間に介設され、前記複数の情報処理装置間で送受信される情報を中継する接続装置を有し、

前記管理部は、

前記自グループにおける情報処理装置及び前記接続装置の接続関係に関する情報に基づいて、前記自グループから分割する複数の情報処理装置を決定することを特徴とする、付記 8 記載の情報処理システム。

【 0 2 4 3 】

(付記 10)

前記管理部は、

前記自グループにおける情報処理装置の数が第 5 所定値未満の場合、前記自グループと前記他のグループのうちのいずれかのグループとを統合することを特徴とする、付記 8 又は付記 9 記載の情報処理システム。

【 0 2 4 4 】

(付記 11)

相互に接続される複数の情報処理装置の各々において、

前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信する受信処理部と、

前記受信処理部が前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定する判定部と、

前記判定部が判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信する送信処理部と、を有することを特徴とする、情報処理装置。

【 0 2 4 5 】

(付記 12)

相互に接続される複数の情報処理装置の各々を制御する情報処理装置の制御プログラムにおいて、

前記情報処理装置に、

前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信させ、

前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定させ、

判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信させることを特徴とする、情報処理装置の制御プログラム。

【 0 2 4 6 】

(付記 13)

前記情報処理装置に判定させる処理は、

受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、前記他の情報処理装置の各々からの前記状態情報の受信状況とに基づいて行なわれ、

前記情報処理装置に送信させる処理は、

第 1 所定時間ごとに行なわれることを特徴とする、付記 12 記載の情報処理装置の制御

10

20

30

40

50

プログラム。

【 0 2 4 7 】

( 付記 1 4 )

前記情報処理装置に判定させる処理は、

受信した前記状態情報が示す前記複数の情報処理装置の各々の状態と、自情報処理装置の状態に関する状態情報に関する自己状態情報とに基づいて、前記複数の情報処理装置の各々の状態を判定させることを特徴とする、付記 1 2 記載の情報処理装置の制御プログラム。

【 0 2 4 8 】

( 付記 1 5 )

前記情報処理装置に、

前記第 1 所定時間以上の時間である第 2 所定時間内に前記状態情報を受信しなかった他の情報処理装置の状態を、停止の可能性を示す第 1 状態と判定させ、

受信した前記状態情報に基づいて、第 1 所定数以上の前記複数の情報処理装置で前記第 1 状態であると判定された情報処理装置の状態、又は、前記他の情報処理装置の少なくとも 1 つから停止を示す第 2 状態であると判定された情報処理装置の状態を、前記第 2 状態と判定させることを特徴とする、付記 1 3 記載の情報処理装置の制御プログラム。

【 0 2 4 9 】

( 付記 1 6 )

前記情報処理装置に、

受信した前記状態情報に基づいて、前記第 1 所定数以上の数である第 2 所定数以上の前記複数の情報処理装置で前記第 2 状態であると判定された情報処理装置を、リカバリ処理中を示す第 3 状態と判定させ、

前記第 3 状態と判定した情報処理装置に対して、リカバリ処理を実行させることを特徴とする、付記 1 5 記載の情報処理装置の制御プログラム。

【 0 2 5 0 】

( 付記 1 7 )

前記複数の情報処理装置が複数のグループに分割された前記情報処理システムにおける前記情報処理装置に、

前記複数のグループのうちの自グループ以外の他のグループの各々における他の代表情報処理装置から、前記他の代表情報処理装置により判定された前記複数のグループの代表情報処理装置の各々の状態に関する代表状態情報を受信させ、

前記他の代表情報処理装置の各々から受信した前記代表状態情報に基づいて、前記複数の代表情報処理装置の各々の状態を判定させ、

判定した前記複数の代表情報処理装置の各々の状態に関する送信用代表状態情報を、前記他の代表情報処理装置の各々へ送信させ、

前記情報処理装置に前記送信用状態情報を送信させる処理は、

前記送信用状態情報を、前記自グループにおける他の情報処理装置の各々へ送信させ、

前記情報処理装置に前記複数の情報処理装置の各々の状態を判定させる処理は、

前記自グループにおける他の情報処理装置の各々から受信した前記状態情報に基づいて、前記自グループにおける情報処理装置の各々の状態を判定させることを特徴とする、付記 1 2 ~ 1 6 のいずれか 1 項記載の情報処理装置の制御プログラム。

【 0 2 5 1 】

( 付記 1 8 )

前記情報処理装置に、

前記自グループにおける情報処理装置の数が第 4 所定値を超えた場合、前記自グループから、複数の情報処理装置を分割して新たなグループを作成させることを特徴とする、付記 1 7 記載の情報処理装置の制御プログラム。

【 0 2 5 2 】

( 付記 1 9 )

10

20

30

40

50

前記自グループにおける情報処理装置と、前記複数の情報処理装置間に介設され前記複数の情報処理装置間で送受信される情報を中継する接続装置との接続関係に関する情報に基づいて、前記自グループから分割する複数の情報処理装置を決定させることを特徴とする、付記 18 記載の情報処理装置の制御プログラム。

【0253】

(付記 20)

前記情報処理装置に、

前記自グループにおける情報処理装置の数が第 5 所定値未満の場合、前記自グループと前記他のグループのうちのいずれかのグループとを統合させることを特徴とする、付記 18 又は付記 19 記載の情報処理装置の制御プログラム。

10

【0254】

(付記 21)

相互に接続される複数の情報処理装置を有し、前記複数の情報処理装置間で通信を行なう情報処理システムの制御方法において、

前記複数の情報処理装置の各々が、

前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信し、

前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定し、

20

判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信することを特徴とする、情報処理システムの制御方法。

【0255】

(付記 22)

相互に接続される複数の情報処理装置の各々において、

プロセッサを有し、

前記プロセッサが、

前記複数の情報処理装置のうちの自情報処理装置以外の他の情報処理装置の各々から、前記他の情報処理装置により判定された前記複数の情報処理装置の各々の状態に関する状態情報を受信し、

30

前記他の情報処理装置の各々から受信した前記状態情報に基づいて、前記複数の情報処理装置の各々の状態を判定し、

判定した前記複数の情報処理装置の各々の状態に関する送信用状態情報を、前記他の情報処理装置の各々へ送信することを特徴とする、情報処理装置。

【符号の説明】

【0256】

1 ストレージシステム (情報処理システム)

10, 10-1 ~ 10-5, 10A, 10A-1 ~ 10A-6, 10A-11 ~ 10A-13, 10A-21 ~ 10A-23 ノード (ストレージ装置, 情報処理装置)

10a CPU (プロセッサ)

40

10b メモリ

10c 記憶部

10d ネットワークインタフェース

10e 入出力部

10f, 10h 記録媒体

10g 読取部

11, 11A ノード状態保持部

12, 12A 受信処理部

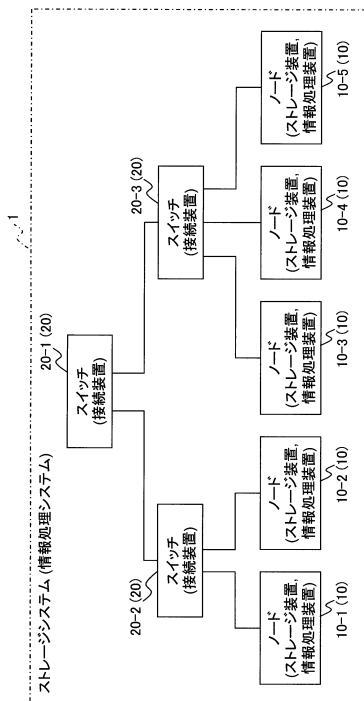
13, 13A ノード状態決定部 (判定部)

14, 14A 送信処理部

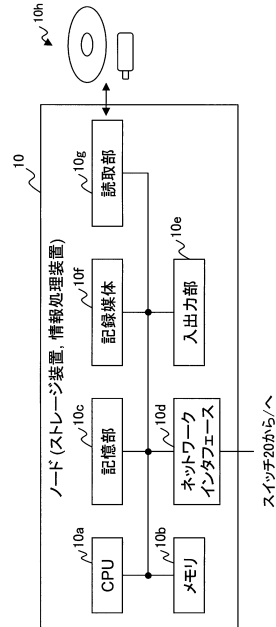
50

- 1 5      リカバリ処理部
- 1 6      停止処理部
- 1 0 1    パーティ情報保持部
- 1 0 2    パーティ間受信処理部（グループ間受信処理部）
- 1 0 3    パーティ間ノード状態決定部（グループ間判定部）
- 1 0 4    パーティ間送信処理部（グループ間送信処理部）
- 1 0 5    パーティ管理部（管理部）
- 2 0 , 2 0 - 1 ~ 2 0 - 3    スイッチ（接続装置）

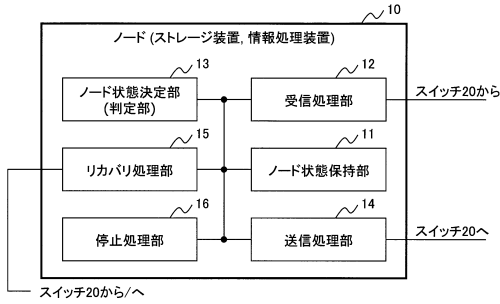
【図 1】



【図 2】



【図 3】



【図 4】

ノード状態情報 (状態情報) ↖ T1

ノードID	状態	IPアドレス	ポート番号
1	Alive	192.168.0.1	12345
2	Alive	192.168.0.2	12345
3	Alive	192.168.0.3	12345
4	Suspect	192.168.0.4	12345
5	Alive	192.168.0.5	12345

【図 5】

ノード状態管理情報 ↖ T2

ノードID	状態	by	最終更新情報	IPアドレス	ポート番号
1	Alive	by 5	1sec ago	192.168.0.1	12345
2	Alive	by 4	1sec ago	192.168.0.2	12345
3	Alive	by 3	1sec ago	192.168.0.3	12345
4	Suspect	by 2	30 sec ago	192.168.0.4	12345
5	Alive	by 1	1sec ago	192.168.0.5	12345

他ノードからの情報  
自ノードでの判断

【図 6】

送信情報 ↖ T3

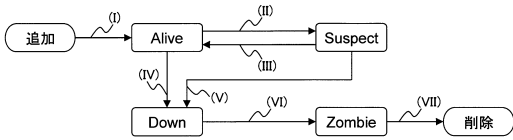
ノードID	状態	IPアドレス	ポート番号
6	Alive	192.168.0.6	12345

【図 7】

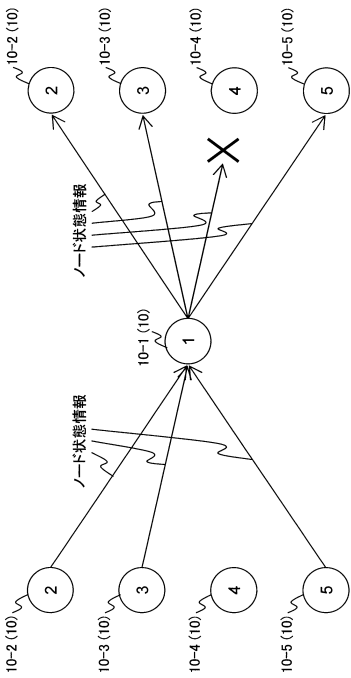
ノード状態情報 (状態情報) ↖ T1'

ノードID	状態	IPアドレス	ポート番号
1	Alive	192.168.0.1	12345
2	Alive	192.168.0.2	12345
3	Alive	192.168.0.3	12345
4	Suspect	192.168.0.4	12345
5	Alive	192.168.0.5	12345
6	Alive	192.168.0.6	12345

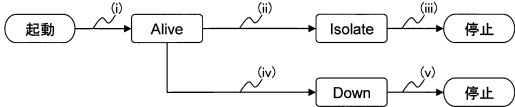
【図 8】



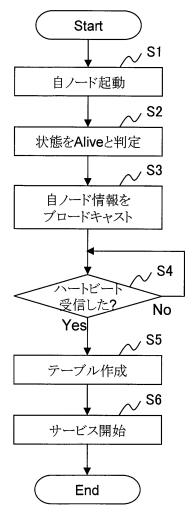
【図 9】



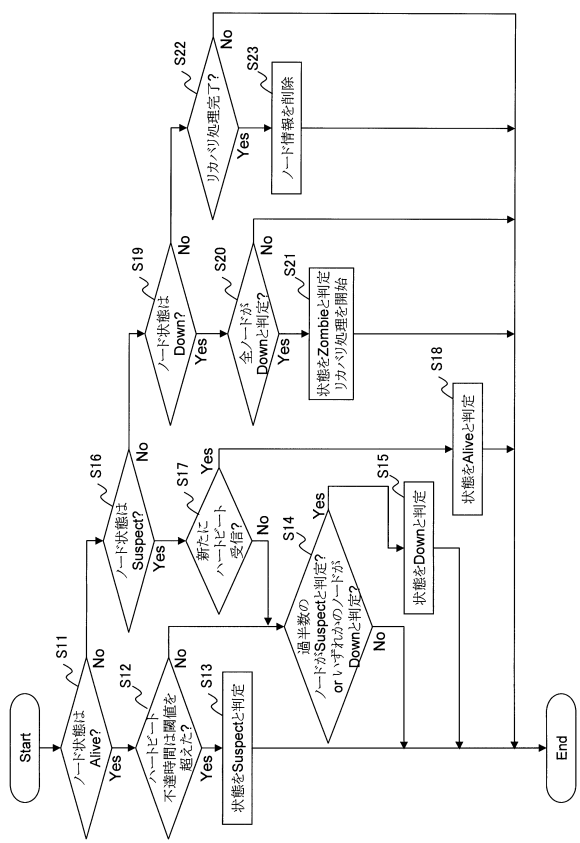
【図 10】



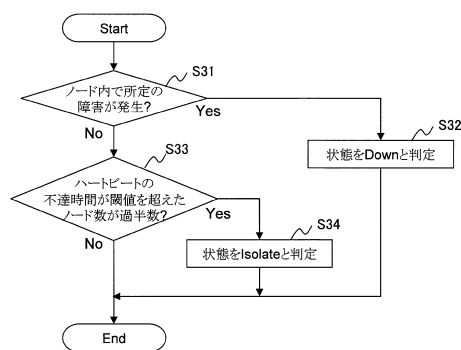
【図 1 1】



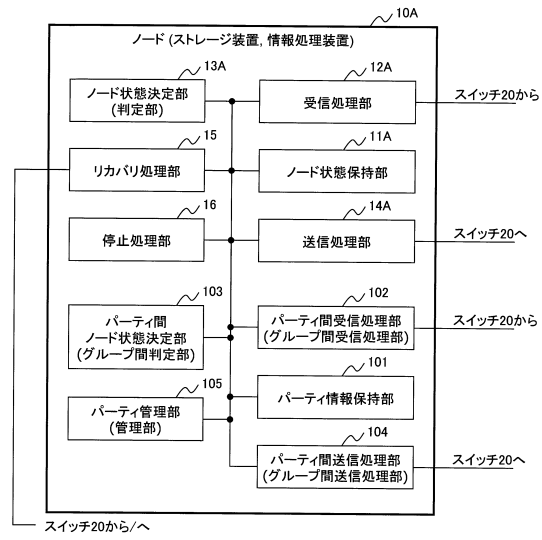
【図 1 2】



【図 1 3】



【図 1 4】



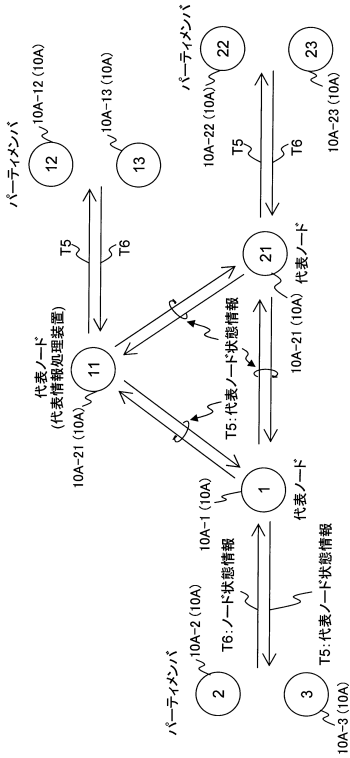
【図 1 5】

パーティ管理情報 T4

パーティID	ノードID	バージョン番号
A	1~10	1
B	11~20	1
C	21~30	3
D	31~40	4
E	41~50	2



【図 16】



【図 17】

代表ノード状態情報 (代表状態情報)  
(代表ノードからメンバノード、代表ノードから代表ノード)

ノードID	状態	IPアドレス	ポート番号
1	Alive	192.168.0.1	12345
2	Alive	192.168.0.2	12345
3	Alive	192.168.0.3	12345
11	Alive	192.168.0.11	12345
12	Alive	192.168.0.12	12345
13	Alive	192.168.0.13	12345
21	Alive	192.168.0.21	12345
22	Alive	192.168.0.22	12345
23	Alive	192.168.0.23	12345

【図 18】

ノード状態情報 (状態情報)  
(メンバノードから代表ノード、メンバノードからメンバノード)

ノードID	状態	IPアドレス	ポート番号
1	Alive	192.168.0.1	12345
2	Alive	192.168.0.2	12345
3	Alive	192.168.0.3	12345

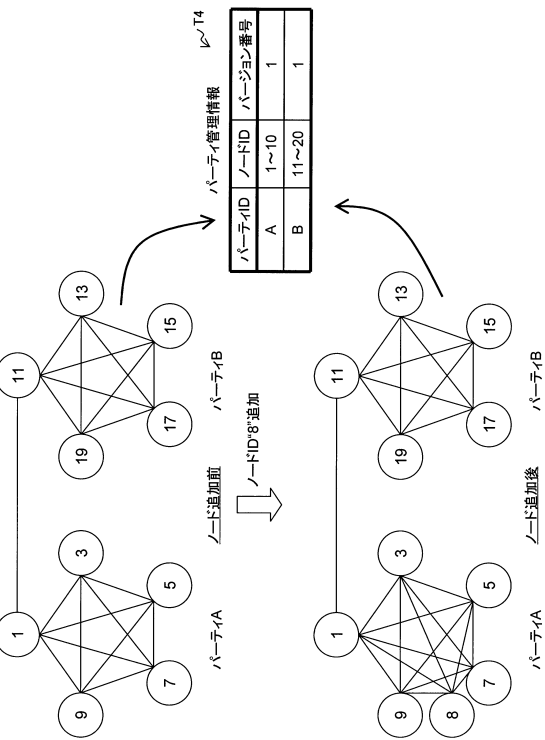
【図 19】

ノード状態管理情報

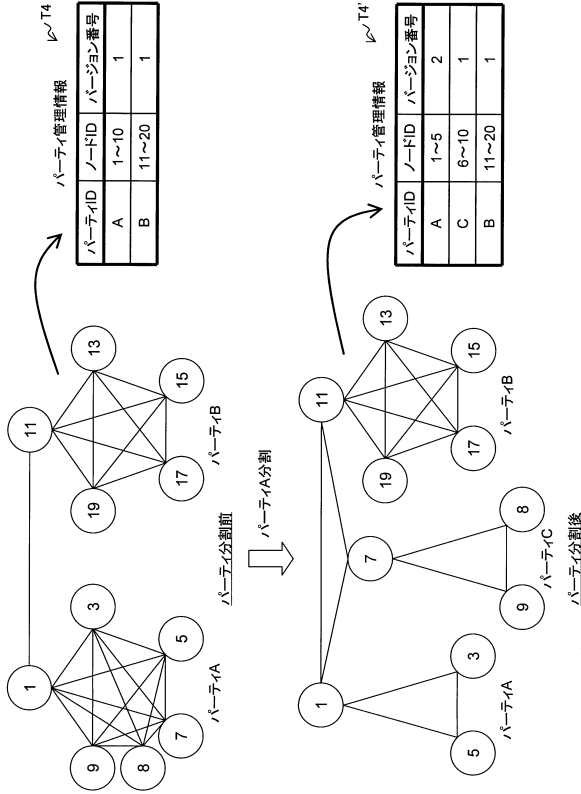
ノードID	状態	by 2	by 3	by 11	by 12	by 13	by 21	by 22	by 23	最終更新情報	IPアドレス	ポート番号
1	Alive	Alive	Alive	Alive			Alive			1sec ago	192.168.0.1	12345
2	Alive	Alive	Alive							1sec ago	192.168.0.2	12345
3	Alive	Alive	Alive							1sec ago	192.168.0.3	12345
11	Alive	Alive	Alive	Alive	Alive	Alive	Alive			1sec ago	192.168.0.11	12345
12	Alive			Alive	Alive	Alive					192.168.0.12	12345
13	Alive			Alive	Alive	Alive					192.168.0.13	12345
21	Alive			Alive			Alive	Alive	Alive	1sec ago	192.168.0.21	12345
22	Alive						Alive	Alive	Alive		192.168.0.22	12345
23	Alive						Alive	Alive	Alive		192.168.0.23	12345

自ノードでの判断  
代表ノード間でハートビート不達  
他ノードでの判断

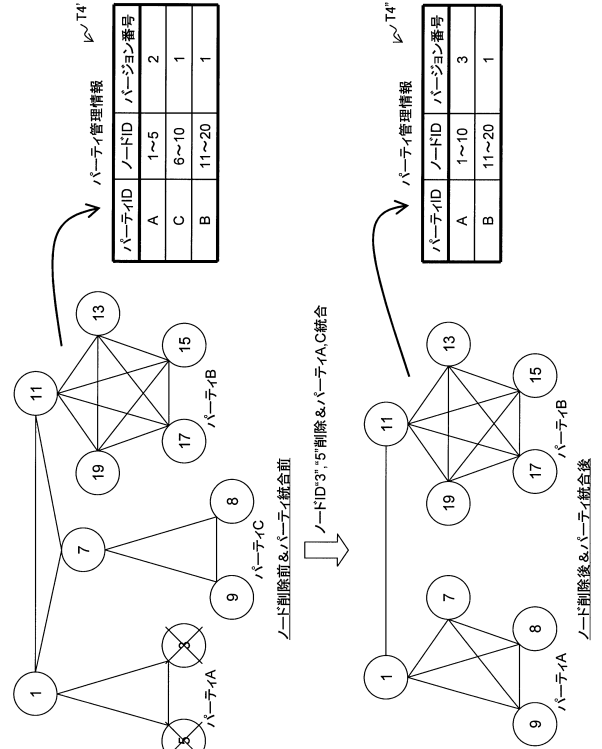
【図 20】



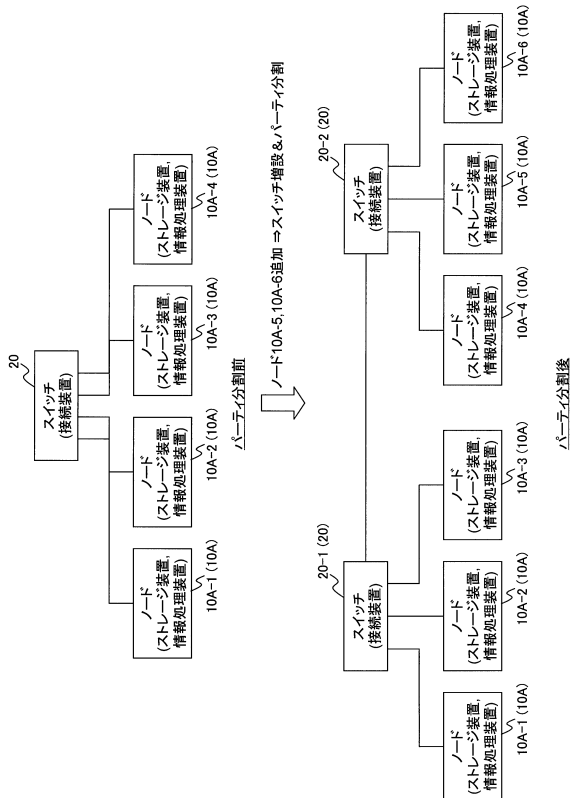
【図 2 1】



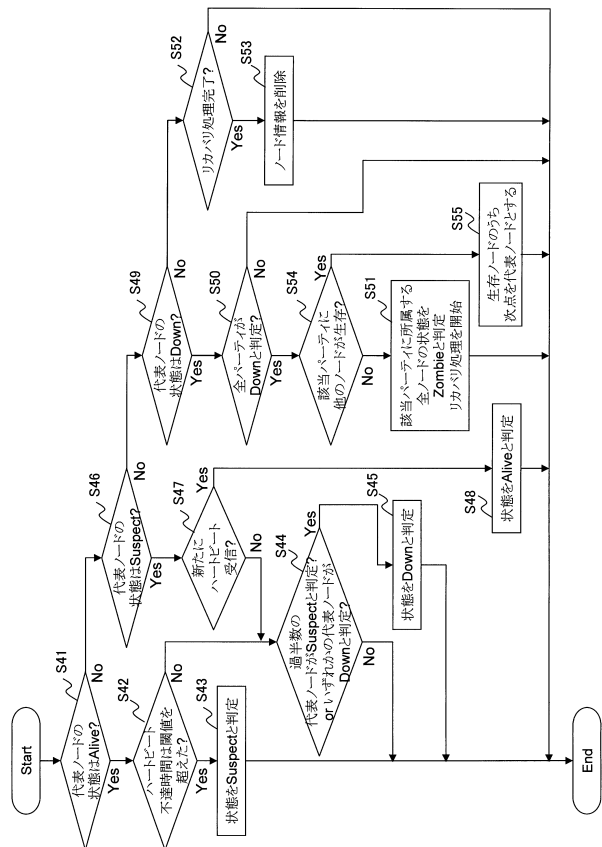
【図 2 2】



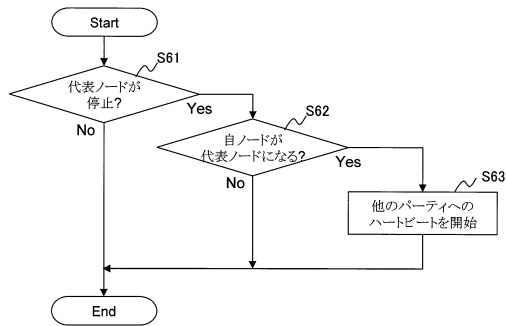
【図 2 3】



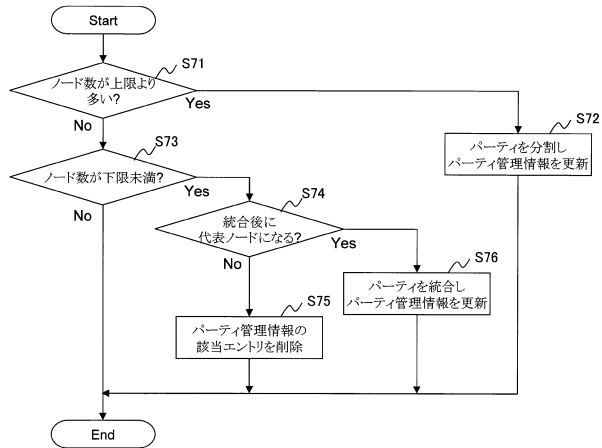
【図 2 4】



## 【図 25】



## 【図 26】



---

フロントページの続き

(56)参考文献 特開2002-132535(JP,A)  
特開昭57-201945(JP,A)  
特開2009-217504(JP,A)  
特開2002-312199(JP,A)

(58)調査した分野(Int.Cl., DB名)  
G06F 11/20  
G06F 11/22