

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5803908号
(P5803908)

(45) 発行日 平成27年11月4日(2015.11.4)

(24) 登録日 平成27年9月11日(2015.9.11)

(51) Int.Cl.		F I			
G06F 12/00	(2006.01)	G06F 12/00	545A		
G06F 3/06	(2006.01)	G06F 3/06	301Z		
G06F 13/10	(2006.01)	G06F 13/10	340A		

請求項の数 6 (全 15 頁)

(21) 出願番号	特願2012-512877 (P2012-512877)	(73) 特許権者	000004237
(86) (22) 出願日	平成23年4月27日(2011.4.27)		日本電気株式会社
(86) 国際出願番号	PCT/JP2011/060238		東京都港区芝五丁目7番1号
(87) 国際公開番号	W02011/136261	(74) 代理人	100106909
(87) 国際公開日	平成23年11月3日(2011.11.3)		弁理士 棚井 澄雄
審査請求日	平成26年2月13日(2014.2.13)	(74) 代理人	100134544
(31) 優先権主張番号	特願2010-103859 (P2010-103859)		弁理士 森 隆一郎
(32) 優先日	平成22年4月28日(2010.4.28)	(74) 代理人	100149548
(33) 優先権主張国	日本国(JP)		弁理士 松沼 泰史
		(74) 代理人	100162868
			弁理士 伊藤 英輔
		(72) 発明者	鳥居 隆史
			東京都港区芝五丁目7番1号 日本電気株式会社内

最終頁に続く

(54) 【発明の名称】 ストレージシステム及びストレージシステムの制御方法

(57) 【特許請求の範囲】

【請求項1】

複数のストレージノードと管理サーバを具備するストレージシステムにおいて、
前記各ストレージノードは、
データを格納する第1の記憶部と、
ネットワーク的な距離が所定の範囲内にある前記ストレージノードからなるネットワークグループと、前記データに対応するデータ識別情報に基づくストレージグループと、
に関するノード情報を格納する第2の記憶部と、を具備し、
前記管理サーバは、

複数の前記ストレージグループの何れに属するかを複数の前記ストレージノードのそれぞれに対して特定し、前記ネットワーク的な距離が所定の範囲内にある複数の前記ネットワークグループの何れに属するかを複数の前記各ストレージノードに対して特定し、複数の前記各ストレージノードに、前記特定した前記ストレージグループの識別情報と前記特定した前記ネットワークグループの識別情報とを割り当て、

前記各ストレージノードは、該ストレージノードが属する前記ネットワークグループ内の全ての前記ストレージノードの情報である当該全ての前記ストレージノードの識別情報と、該ストレージノードが属する前記ストレージグループ内の全ての前記ストレージノードの情報である当該全ての前記ストレージノードの識別情報とを含む前記ノード情報を前記第2の記憶部に記憶し、

前記各ストレージノードは、前記データのアクセス要求を受信したときに、前記データ

10

20

のアクセス要求に基づいて、前記ノード情報に含まれる前記ネットワークグループと前記ストレージグループとに関する情報である当該データのアクセスに用いるデータ識別子を分散関数に入力して得られる値と、前記ノード情報とを参照して前記データへのアクセス処理を行う、ストレージシステム。

【請求項 2】

前記各ストレージノードは、該ストレージノードが属する前記ネットワークグループ内の全ての前記ストレージノードの情報と、該ストレージノードが属する前記ストレージグループ内の全ての前記ストレージノードの情報とを含む前記ノード情報を前記第 2 の記憶部に記憶する、請求項 1 に記載のストレージシステム。

【請求項 3】

前記各ストレージノードは、前記データへのアクセス要求に基づいて、当該データのアクセスに用いるデータ識別子を分散関数に入力して得られる値に基づいて前記ストレージグループの識別情報を特定し、当該ストレージグループの識別情報と自ノードと同一のネットワークグループの識別情報が割り当てられた他ストレージノードに対して、前記アクセス要求の対象となるデータのリクエスト要求を行う、

請求項 1 または請求項 2 に記載のストレージシステム。

【請求項 4】

データを格納する複数のストレージノードを含むストレージシステムにおけるストレージシステムの制御方法であって、

前記データのアクセス要求を受信したときに、ネットワーク的な距離が所定の範囲内にある前記ストレージノードからなるネットワークグループと、前記データに対応するデータ識別情報に基づくストレージグループと、に関するノード情報を参照するステップと、前記ノード情報に基づいて前記データへのアクセス処理を行うステップと、

を具備し、

前記アクセス要求を受信した前記ストレージノードにおいて、前記データ識別情報に基づいて前記ストレージグループを判定するステップと、

前記ノード情報を参照し、前記アクセス要求を受信した前記ストレージノードが属している前記ネットワークグループ内で前記判定した前記ストレージグループに属する前記ストレージノードを検出するステップと、

前記検出されたストレージノードに前記データへのアクセスを要求するステップと、前記検出されたストレージノードからの応答に基づいて前記アクセス要求に対して応答するステップと、

を更に具備する、ストレージシステムの制御方法。

【請求項 5】

前記検出されたストレージノードにおいて、前記アクセス要求に対応するデータを前記検出されたストレージノードが格納している場合は、前記アクセス要求に対応するデータを、前記アクセス要求を受信した前記ストレージノードに送信するステップと、

前記検出されたストレージノードにおいて、前記アクセス要求に対応するデータを前記検出されたストレージノードが格納していない場合は、前記ノード情報を参照して、前記検出されたストレージノードが属する前記ストレージグループ内の他の前記ストレージノードに対して前記アクセス要求に対応するデータを要求し、前記アクセス要求を受信した前記ストレージノードに前記アクセス要求に対応するデータを送信するステップと、

を更に具備する、請求項 4 に記載のストレージシステムの制御方法。

【請求項 6】

前記検出されたストレージノードにおいて、前記アクセス要求に対応するデータを前記検出されたストレージノードが格納している場合は、前記アクセス要求を受信した前記ストレージノードに前記アクセス要求に対応するデータを送信するステップと、

前記検出されたストレージノードにおいて、前記アクセス要求に対応するデータを前記検出されたストレージノードが格納していない場合は、前記ノード情報を参照して、前記検出されたストレージノードが属する前記ストレージグループ内の前記アクセス要求に対

10

20

30

40

50

応するデータを格納している他の前記ストレージノードを検索し、前記検索されたストレージノードを、前記アクセス要求を受信した前記ストレージノードに通知するステップと

、
を更に具備する、請求項4に記載のストレージシステムの制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、大量のデータを保持するストレージシステム及びストレージシステムの制御方法に関する。

本願は、2010年4月28日に、日本に出願された特願2010-103859号に
基づき優先権を主張し、その内容をここに援用する。 10

【背景技術】

【0002】

ネット上で提供されるデータ量の増加に伴って、大量のデータを保持するストレージが必要とされている。例えば、ウェブ検索サービスを提供する企業では、複数のサーバを並列に並べた分散ストレージ技術を採用している。この分散ストレージ技術は、数千のノード（「ピア」ともいう）にデータを分散して配置し、全体として一つの大きなストレージを構成する技術である。また、分散ストレージ技術は、高価なストレージ専用装置ではなく、比較的安価なサーバを複数並べることによって大容量のストレージを実現することができる技術として、扱うデータ量が増大しているエンタープライズやキャリアの事業分野でも注目されている技術である。一部のペタ（ 10^{15} ）バイトを超えた大容量のデータを格納する事業分野では、ストレージ専用装置においてデータを格納することができる容量がボトルネックとなってしまうため、大量のデータ格納を実現するための方法として、分散ストレージ技術を使用するしか解がないというケースも出始めている。 20

【0003】

しかし、分散ストレージ技術においては、データが複数のノードに分散している。このため、データにアクセスしようとするクライアントは、まず、データを持っているノードの位置を知る必要がある。従って、近年注目されている分散ストレージ技術においては、データを保持しているノードの位置を知るための方法が、技術的なポイントとなっている。 30

【0004】

データを保持しているノードの位置を知るための一つの方法として、データの位置情報を管理するメタサーバを設けるメタサーバ方式がある。このメタサーバ方式では、ストレージシステムの構成が大規模になることに伴って、データを格納しているノードの位置を検出するメタサーバの性能がボトルネックになることがある。

【0005】

そのため、データを保持しているノードの位置を知るための別の方法として、分散関数（例えば、ハッシュ関数）を用いてデータの位置を求める方法が注目されている。これは、以前からP2P（Peer to Peer：ピアツーピア）ネットワークの分野で用いられている分散ハッシュテーブル（Distributed Hash Table：DHT）の技術を利用しており、キーバリューストレージ（Key-Valueストレージ：KVS）と呼ばれている。このKVSは、データにアクセスする識別子をキー（Key）、データをバリュー（Value）とし、キー（Key）を分散関数にかけて、すなわち、キーを分散関数の入力値として、データを格納しているノード（以下、「データ格納ノード」という）の位置を算術的に求める。 40

【0006】

KVSのメタサーバ方式との違いは、メタサーバ方式が全てのデータの位置情報をもたなければならないのに対し、KVSでは、全てのクライアントで分散関数とノードリストとを共有すればよいだけなので、分散関数とノードリストとを共有する際のコストが小さく、性能的なボトルネックがないことである。KVSを用いれば、メタサーバ方式のよう 50

なメタサーバの性能によるボトルネックがなく、ストレージシステムが大規模になった場合でも、性能拡張性（スケーラビリティ）のある大容量ストレージを実現することができる（特許文献1～3参照）。

【先行技術文献】

【特許文献】

【0007】

【特許文献1】特開2009-289161号公報

【特許文献2】特開2009-151403号公報

【特許文献3】特開2008-269141号公報

【発明の概要】

10

【発明が解決しようとする課題】

【0008】

既存のKVS技術では、分散関数としてハッシュ関数（あるいは、それに類似の分散関数）を用いて、算術的にデータ格納ノードが決定されている。また、DHT技術では、分散ハッシュテーブルのルーティングにより、オーバーレイルーティングが行われている。このため、既存のKVS技術を用いたストレージシステムでは、大きく分けて、例えば、以下のような2つの点がみられる。

【0009】

まず、1つ目の点は、既存のKVS技術ではデータ配置の自由度が高いとはいえないことがある点である。KVS技術は、大規模のストレージシステムで利用されるため、必然的に各ストレージが大規模ネットワークに分散して配置されることとなる。すると、データにアクセスするノード（クライアント）とデータを保持しているデータ格納ノードとのネットワーク上の距離が遠くなるケースがあり得る。クライアントとデータ格納ノードとの間の遅延が大きければ、ストレージシステムにおける処理の速度も遅くなる。このため、クライアントとデータ格納ノードとをできるだけ近くに配置することによってストレージシステムの処理速度を向上することが求められている。そのためには、データ配置を自由に変更できることが求められる。

20

また、データを分散させて保持するのではなく、逆にデータを予め定めたデータ格納ノードに集めて、空いているデータ格納ノードを作ることによって、ストレージシステムの省電力化を実現することも考えられる。このようなときにも、データ配置を自由に変更することができれば、ストレージシステムの省電力モードによる制御も可能となる。しかしながら、既存のKVS技術では、データ格納ノードがハッシュ関数によって決まってしまうため、例えば、あるデータを特定のデータ格納ノードに配置しようとしても、データ配置を自由に制御することができない可能性がある。

30

【0010】

そして、2つ目の点は、分散ハッシュテーブルのオーバーレイルーティングが実際のネットワークトポロジーと連動していないため、効率的でないルーティングがなされて、結果的に性能が高いものではないことがある点である。極端な例では、東京から大阪に行くのにアメリカを経由するというようなこともあり得る。

【0011】

40

従来、スケーラブルなストレージシステムにおいて、データ配置の自由度を向上させることができるストレージシステムの制御方法が求められていた。

【課題を解決するための手段】

【0012】

本発明は、例えば、以下の側面を有しても良い。ただし、以下の記述は本発明を限定するものではない。

第1の側面はストレージシステムの制御方法であって、データを格納する複数のストレージノードによって構成されたストレージシステムの制御方法において、該ストレージシステムに含まれる前記複数のストレージノードを、該ストレージシステム内においてネットワーク的な距離が予め定められた距離の範囲内であるストレージノードからなる第1の

50

グループと、データを格納しているストレージノードの位置の情報を共有しているストレージノードからなる第2のグループと、にグループ化し、前記第2のグループ毎に該第2のグループを識別する論理的な空間上の識別子を割り当て、データ識別子を分散関数の入力値として前記論理的な空間上の位置を算出し、前記算出された位置に対応する識別子が割り当てられた前記第2のグループに属する前記ストレージノードに該データ識別子に対応するデータを格納する。

【0013】

また、上記複数のストレージノードのそれぞれは、1つ以上の前記第1のグループと1つ以上の前記第2のグループとに必ず属し、該ストレージノードが属している前記第2のグループ内の他の全てのストレージノードのリストと、該ストレージノードが属している前記第1のグループ内の他の全てのストレージノードのリストとを、記憶してもよい。

10

【0014】

また、上記ストレージシステム内でデータのアクセス要求をするストレージノードは、データ識別子を分散関数の入力値として、前記論理的な空間上の位置を算出し、該算出された位置に対応する識別子が割り当てられた前記第2のグループを選択し、該ストレージノードに記憶しているノードリストから、該ストレージノードが属している前記第1のグループ内で前記選択した第2のグループに属している他のストレージノードを検索し、前記検索した他のストレージノードにデータのアクセス要求を出力してもよい。

【0015】

また、上記ストレージシステム内でデータのアクセス要求を受けたストレージノードは、前記アクセス要求を受けたデータを該ストレージノード内に格納している場合は、該データを要求した前記ストレージノードに要求されたデータを出力し、前記アクセス要求を受けたデータを該ストレージノード内に格納していない場合は、該ストレージノードに記憶しているノードリストから、該ストレージノードが属している前記第2のグループ内で要求されたデータを格納している他のストレージノードを検索し、前記検索した他のストレージノードに、前記データのアクセス要求を転送してもよい。

20

【0016】

また、上記ストレージシステム内でデータのアクセス要求を受けたストレージノードは、前記アクセス要求を受けたデータを該ストレージノード内に格納している場合は、該データを要求した前記ストレージノードに要求されたデータを出力し、前記アクセス要求を受けたデータを該ストレージノード内に格納していない場合は、該ストレージノードに記憶しているノードリストから、該ストレージノードが属している前記第2のグループ内で要求されたデータを格納している他のストレージノードを検索し、前記検索した他のストレージノードを、前記アクセス要求をするストレージノードに通知してもよい。

30

【発明の効果】

【0017】

上記側面によれば、例えば、スケーラブルなストレージシステムにおいて、データ配置の自由度を向上させることができるという効果が得られる。

【図面の簡単な説明】

【0018】

40

【図1】一実施形態によるストレージシステムの概略構成を示したブロック図である。

【図2】従来のストレージシステムにおける既存のKVS技術を説明する図である。

【図3】一実施形態のストレージシステムにおける論理的なノード構成の一例を示した図である。

【図4】一実施形態のストレージシステムにおけるKVS技術を説明する図である。

【図5】一実施形態のストレージシステムにおける物理構成の一例を示した図である。

【図6】一実施形態のストレージシステムにおけるデータアクセス処理の流れを示したシーケンス図である。

【発明を実施するための形態】

【0019】

50

以下、一実施形態について、図面を参照して説明する。図1は、本実施形態によるストレージシステムの概略構成を示したブロック図である。図1において、ストレージシステム3は、データ格納ノードである複数のストレージノード1を備えている。各ストレージノード1には、1つもしくは複数のIDが付与されている。そして、これら複数のストレージノード1が、管理サーバ10によってストレージシステム3内でマッピングされて、1つのグローバルネームスペースとして構成されている。なお、ストレージシステム3内の各ストレージノード1が設置されている物理的な位置は1カ所ではなく、複数の場所に設置されている各ストレージノード1が、例えば、ネットワークなどによって接続されることによって、1つのグローバルネームスペースを構成している。

【0020】

10

クライアント2は、ストレージシステム3のデータにアクセスするノードである。そして、クライアント2は、ストレージシステム3を、1つの大きなストレージとみなしてアクセスする。

【0021】

ここで、既存のKVS技術について説明する。図2は、従来のストレージシステムにおける既存のKVS技術を説明する図である。図2では、それぞれ「a」、「b」、「c」、「d」というIDが付与されたデータ格納ノードが、論理的な空間の円周上にマッピングされている場合を示している。既存のKVS技術では、データ識別子であるKeyを、分散関数Fにかけて $F(\text{Key})$ が求められる。そして、この円周上で、 $F(\text{Key})$ の位置から右回りに最も近いIDを持つデータ格納ノードに、 $F(\text{Key})$ に対応するデータが保持される。図2においては、 $a < F(\text{Key}) < b$ を満たす $F(\text{Key})$ に対応するデータが、IDとして「b」が付与されたデータ格納ノードに格納されることを示している。

20

【0022】

この既存のKVS技術によるデータ格納の方法は、クライアントが、分散関数Fとデータ格納ノードのリストを共有するのみであるので、クライアントが共有する情報が少なく済むという利点がある。しかし、データ識別子であるKeyは、一度データに付与された後に変更することができないため、任意のデータ格納ノードにデータを移動することができず、データ配置の自由度がない。

【0023】

30

次に、本実施形態のストレージシステム3におけるデータ格納の方法について説明する。図3は、本実施形態のストレージシステム3における論理的なノード構成の一例を示した図である。図3に示すように、本実施形態のストレージシステム3では、ストレージシステム3内の各ストレージノード1を、ストレージグループ4およびネットワークグループ5という2種類のグループにグループ化する。このストレージシステム3内における各ストレージノード1のグループ化は、管理サーバ10によって行われる。

【0024】

ストレージグループ4は、ストレージシステム3で用いるKVS技術に基づいてデータを格納しているストレージノード1の位置の情報を共有するストレージノード1で構成されたグループである。

40

【0025】

また、ネットワークグループ5は、ストレージシステム3におけるネットワーク的な距離に基づいて、管理サーバ10によって決定されたグループであり、ネットワーク的な距離が予め定められた距離の範囲内で比較的近いストレージノード1で構成されたグループである。すなわち、ネットワークグループ5に属する任意の2つのストレージノード1間のネットワーク的な距離は、予め定められた距離の範囲内の距離となる。

【0026】

ストレージシステム3における各ストレージノード1は、管理サーバ10によって管理されるいずれかのストレージグループ4に属していると同時に、いずれかのネットワークグループ5に属している。

50

【0027】

なお、図3において白抜きの は、ストレージノード1を示しており、各ストレージノード1を表す 内には、ストレージグループ4における識別番号とネットワークグループ5における識別番号とを表している。より具体的には、各ストレージノード1を表す 内の2桁の符号のうち、左側の符号は、1つのストレージグループ4（図3においては、符号“Y”が付与されたストレージグループ4）におけるストレージノード1の識別番号（1, 2, …, X, …, m）を表し、右側の符号は、1つのネットワークグループ5（図3においては、符号“X”が付与されたネットワークグループ5）におけるストレージノード1の識別番号（1, 2, …, Y, …, n）を表している。

【0028】

次に、本実施形態のストレージシステム3におけるKVS技術について説明する。図4は、本実施形態のストレージシステム3におけるKVS技術を説明する図である。図4に示すように、本実施形態のストレージシステム3におけるKVS技術では、複数のストレージノード1をグループ化することによって、データ配置の自由度を増加させる。そして、ストレージシステム3では、管理サーバ10によって、各ストレージグループ4に1つもしくは複数のIDを付与され、IDを付与されたストレージグループ4を、ストレージシステム3内にマッピングする。図4では、それぞれ「A」、「B」、「C」、「D」というIDが付与されたストレージグループ4が、図2に示した既存のKVS技術と同様に、論理的な空間の円周上にマッピングされている場合を示している。

【0029】

そして、本実施形態のストレージシステム3におけるKVS技術では、図2に示した既存のKVS技術と同様に、データ識別子であるKeyを、分散関数FにかけてF(Key)を求める。そして、この円周上で、F(Key)の位置から右回りに最も近いIDを持つストレージグループ4を、データを保持するストレージグループ4と決定する。続いて、決定されたストレージグループ4内のどのストレージノード1がデータを保持するかを決定し、決定されたストレージノード1に、F(Key)に対応したデータを保持する。

【0030】

次に、本実施形態のストレージシステム3におけるネットワークグループ5について説明する。図5は、本実施形態のストレージシステム3における物理構成の一例を示した図である。上記に述べたとおり、ネットワークグループ5は、ネットワーク的な距離が比較的近いストレージノード1をグループ化したものである。このネットワーク的な距離とは、例えば、ネットワーク経路上のスイッチの段数と考えることができる。より具体的には、図5に示したように、複数のストレージノード1とスイッチ7とから構成される複数のラック6と、それぞれのラック6を束ねる上位スイッチ8からなるストレージシステム3を想定する。この場合、各ラック6がそれぞれのネットワークグループ5に相当する。

【0031】

ストレージシステム3における各ストレージノード1は、上記に述べたとおり、必ず1つ以上のネットワークグループ5に属していると同時に、必ず1つ以上のストレージグループ4に属している。また、各ストレージノード1が属するストレージグループ4は、ネットワークグループ5内に属しているストレージノード1から、ストレージシステム3における全てのストレージグループ4をたどることができるように、管理サーバ10によって割り当てられている。つまり、ある1つのストレージノード1が属するネットワークグループ5内の全てのストレージノード1の和集合をとれば、全てのストレージグループ4をカバーすることができるように、各ストレージノード1のストレージグループ4が割り当てられている。

【0032】

なお、各ストレージグループ4やネットワークグループ5に属するストレージノード1の数は、それぞれのストレージグループ4やネットワークグループ5で異なる数であってもよい。例えば、1つのストレージノード1が、複数のストレージグループ4やネットワークグループ5に属するように、このストレージノード1のストレージグループ4やネッ

10

20

30

40

50

トワークグループ5を割り当てることもできる。

【0033】

各ストレージノード1は、自ストレージノード1が属しているストレージグループ4内の他の全てのストレージノード1のリストと、ネットワークグループ5内の他の全てのストレージノード1のリストとを、ノードリストとして記憶している。それぞれのノードリストには、各ストレージノード1が属しているストレージグループ4およびネットワークグループ5のIDと、各ストレージノード1のアドレス(位置)情報と、各ストレージノード1が格納しているデータ(例えば、データの一覧など)の情報が含まれている。

【0034】

次に、本実施形態のストレージシステム3において各ストレージノード1が記憶しているノードリストの数について説明する。上記に述べたとおり、全てのストレージノード1は、自ストレージノード1が属するネットワークグループ5内の全てのストレージノード1のリストと、自ストレージノード1が属するストレージグループ4内の全てのストレージノード1のリストとを自ストレージノード1のメモリ上に記憶している。ストレージシステム3において各ストレージノード1が記憶しているストレージノード1のリストの総数は、従来のストレージシステムに比べて非常に少ない数となるため、ストレージシステム内におけるメモリ容量の削減と、メンテナンスコストの削減とを実現することができる。

10

【0035】

より具体的には、例えば、1000台のストレージノードで構成されたストレージシステムを考える。従来のストレージシステムにおいて全てのストレージノードのノードリストを記憶する場合、各ストレージノードは、1000個のリストをノードリストとして記憶する必要がある。

20

【0036】

これに対して、本実施形態のストレージシステム3では、自ストレージノード1が属するネットワークグループ5内の全てのストレージノード1の個数のリストと、自ストレージノード1が属するストレージグループ4内の全てのストレージノード1の個数のリストとを、ノードリストとして記憶するのみである。例えば、1000台のストレージノード1が、Nグループのストレージグループ4、およびMグループのネットワークグループ5にそれぞれグループ化され、各ストレージノード1が、それぞれ、1つのストレージグループ4および1つのネットワークグループ5に属している場合を想定する。この場合、各ストレージノード1は、自ストレージノード1が属するネットワークグループ5内のN個のストレージノード1のリストと、ストレージグループ4内のM個のストレージノード1のリストのみを記憶するのみであるため、各ストレージノード1が記憶するノードリストは、 $N + M - 1$ 個のリストとなる。ここで、 -1 個としたのは、図3からもわかるように、自ストレージノード1のリストは、ストレージグループ4とネットワークグループ5とで重複しており、この重複を回避するためである。より具体的には、ストレージグループ4が100グループ、ネットワークグループ5が10グループであった場合には、 $100 + 10 - 1 = 109$ 個のリストを各ストレージノード1が記憶するのみとなる。

30

【0037】

これは、従来のストレージシステムにおいて各ストレージノード1が1000個のリストを記憶していたのに対して、本実施形態のストレージシステム3のストレージノード1が記憶するリストの数は、約10分の1の数であり、各ストレージノード1内でノードリストの記憶に使用するメモリ容量の削減を実現していることとなる。

40

【0038】

また、一般的に、ストレージシステムでは、データにアクセスできない時間を極力減らすため、ストレージシステムの死活監視を定期的に行っている。このストレージシステムの死活監視では、ストレージシステム内の各ストレージノードのエラーをなるべく早く検出する必要があり、ノードリストに含まれているストレージノードが正常に稼働しているか否かという稼働状況をチェックすることによって行われている。もし、ストレージシ

50

テム内のいずれかのストレージノードにエラーが発生している、ストレージシステムのネットワークが不通になっている、などの原因によって、ストレージシステムが正常に稼働していない場合には、ノードリストの変更が必要となる。この稼働状況のチェックにかかるコストは、ノードリストに含まれるリストの数に比例して大きくなるため、リストの数が多くなると、ストレージシステム全体のスケーラビリティを大きく損なう要因となってしまう。そのため、ノードリストに含まれるリストの数を少なく保つことは、スケーラブルなストレージシステムにとっては重要な項目である。本実施形態のストレージシステム3では、各ストレージノード1がノードリストに記憶しているストレージノード1のリストの数が少ないため、メンテナンスコストの削減を実現することができる。

【0039】

次に、本実施形態のストレージシステム3において各ストレージノード1がデータを保持しているストレージノード1の検索方法について説明する。図6は、本実施形態のストレージシステム3におけるデータアクセス処理の流れを示したシーケンス図である。図6では、図3に示した符号“X”が付与されたネットワークグループ5（以下、「ネットワークグループNG_X」という）に属する識別番号“1”が付与されたストレージノード1（以下、「ストレージノードX1」という）がクライアント2となってデータにアクセスする場合について説明する。

【0040】

まず、ストレージノードX1は、データ識別子（Key）を分散関数FにかけてF（Key）を求める（ステップS10）。そして、F（Key）に対応するデータを保持しているストレージノード1が属しているストレージグループ4（以下、「ストレージグループSG_Y」という）を求める（ステップS20）。例えば、「A」、「B」をストレージグループ4のIDとすると、 $A < F(\text{Key}) < B$ を満たす「B」がIDとして付与されたストレージグループ4が求められる（図4参照）。

【0041】

そして、自ストレージノードX1が属するネットワークグループNG_X内で、ストレージグループSG_Yに属するストレージノード1をノードリストから求める（ステップS30）。図6では、ネットワークグループNG_X内でストレージグループSG_Yに属するストレージノードXY（図3参照）が求められたものとする。そして、ストレージノードX1は、ストレージノードXYにデータの要求（リクエスト）を送信する（ステップS40）。

【0042】

続いて、リクエストを受信したストレージノードXYは、要求されたデータが、自ストレージノードXYが保持しているデータであるか否かを検索する（ステップS50）。要求されたデータが自ストレージノードXYの保持しているデータである場合、ストレージノードXYは、ストレージノードX1からのリクエストに応答し、要求されたデータをストレージノードX1に送信する（ステップS60）。そして、ストレージノードX1が、ストレージノードXYから送信されてきたデータを受信することによって、ストレージノードX1によるデータのアクセスを完了する。

【0043】

また、ステップS50において、要求されたデータが自ストレージノードXYの保持しているデータでない場合、ストレージノードXYは、自ストレージノードXYが属しているストレージグループSG_Y内の他のストレージノード1に、要求されたデータが分散されていると判断する。そして、ストレージノードXYは、ノードリストから要求されたデータを格納しているストレージグループSG_Y内の他のストレージノード1を求め、ストレージノードX1からのリクエストを、ストレージグループSG_Y内の他のストレージノード1に転送する（ステップS61）。図6では、ストレージグループSG_Yに属するストレージノード2Y（図3参照）にリクエストが転送された場合を示している。なお、ストレージノードXYがストレージグループ4内の他のストレージノード1にリクエストを転送する方法は、ストレージグループ4内におけるデータ分散方法に依存する。

10

20

30

40

50

このデータ分散方法に関しては、後述する。

【0044】

そして、リクエストが転送されたストレージノード2 Yは、自ストレージノード2 Yが保持しているデータから、要求されたデータを検索する(ステップS70)。そして、ストレージノード2 Yは、ストレージノードXYから転送されてきたストレージノードX1からのリクエストに回答し、要求されたデータをストレージノードXYに送信する(ステップS80)。そして、ストレージノードXYは、ストレージノード2 Yからのリクエストの回答とデータとを、同じネットワークグループNG_X内のストレージノードX1に転送する(ステップS81)。そして、ストレージノードX1が、ストレージノードXYから転送されてきたデータを受信することによって、ストレージノードX1によるデータのアクセスを完了する。

10

【0045】

なお、図6の説明においては、ステップS50において、リクエストを受信したストレージノード1が要求されたデータを保持していない場合に、要求されたデータを格納しているストレージノード1にリクエストを転送する例を説明したが、ストレージノード1が他のストレージノード1にリクエストを転送せず、要求元のストレージノード1に他のストレージノード1を通知する方法とすることもできる。より具体的には、まず、ステップS50において、ストレージノードXYは、ノードリストから要求されたデータを格納しているストレージグループSG_Y内の他のストレージノード1であるストレージノード2 Yを求める。そして、ステップS60において、ストレージノードXYが要求されたデータをストレージノードX1に送信する代わりに、要求したデータがストレージグループSG_Y内のストレージノード2 Yに格納されていること通知する。そして、ストレージノードX1は、通知されたストレージノード2 Yに直接データの要求(リクエスト)を送信(再送信)し、ストレージノード2 Yから送信されてきた要求したデータを受信する。このように、ストレージノードX1が、ストレージノード2 Yにリクエストを送信し直すことによって、ストレージノード2 Yから送信されてきたデータを直接受信することができる。

20

【0046】

次に、本実施形態のストレージシステム3におけるストレージグループ4内のデータ分散方法について、2つの方法を説明する。まず、1つ目の方法は、ストレージグループ4内の全てのデータ配置を管理する1つのストレージノード1(以下、「メタサーバ」という)を決め、そのメタサーバがストレージグループ4内の全てのデータ配置を集中して管理する集中メタサーバ方式である。この集中メタサーバ方式は、データ配置を集中して管理するため、データの移動や複製の管理が容易である。

30

【0047】

この集中メタサーバ方式では、クライアント2がストレージシステム3内のデータをアクセスする際に、必ずメタサーバに対しての問い合わせが行われる。しかし、ストレージシステム3において、データを移動する頻度はそれほど多くないため、例えば、クライアント2からアクセスされるデータが特定のデータに集中している(局所性がある)可能性がある。この場合には、各ストレージノード1が、自ストレージノード1内でアクセスされたデータの位置情報を、例えば、キャッシュメモリなどに一時記憶(キャッシュ)しておくことによって、メタサーバに対しての問い合わせを行わずに、ストレージシステム3内のデータをアクセスすることもできる。

40

【0048】

なお、ストレージシステム3内のデータをアクセスする際に、メタサーバに対しての問い合わせを行う場合には、メタサーバの性能がボトルネックになる可能性がある。また、ストレージシステム3の構成が大規模である場合には、ストレージシステム3内のストレージグループ4がネットワークに分散しているため、ネットワークトラフィックが増えてしまう可能性がある。この場合においても、各ストレージノード1におけるキャッシュの機能と組み合わせることによって、上記に述べた欠点を解決することができる可能性が高

50

い。

【0049】

また、本実施形態のストレージシステム3におけるストレージグループ4内のデータ分散方法の2つ目の方法は、ストレージグループ4間のデータ配置と同様に、分散関数によってストレージノード1を決定するハッシュ方式である。このハッシュ方式では、ハッシュ値の範囲と対応するストレージノード1とが組となったハッシュテーブルに基づいて、ハッシュ値からストレージノード1を求める。この方法は、従来のハッシュテーブルを用いた方法と同様である。ただし、本実施形態のストレージシステム3においては、データ配置の自由度を高めるために、ハッシュ値の範囲を細かい単位（最も小さくなるハッシュ値）に分割していることが異なる。

10

【0050】

例えば、あるデータを別のストレージノード1に移動するときには、ハッシュテーブル上で移動するデータのハッシュ値が対応するストレージノード1を変えることによって、データを別のストレージノード1に移動することができる。このデータの移動方法は、従来のハッシュテーブルを用いたデータの移動においても同様である。しかし、従来のハッシュテーブルでは、ハッシュ値の範囲を分割していないため、データの移動を行う際に、ハッシュ値の範囲に含まれる全てのデータを移動する必要があり、データの移動におけるコストがかかりすぎる場合がある。これに対して、本実施形態のストレージシステム3では、ハッシュ値の範囲を細かい単位に分割しているため、分割したハッシュ値の範囲に含まれるデータのみを移動することができ、従来のハッシュテーブルを用いたデータの移動のように、全てのデータを移動する必要がない。

20

【0051】

なお、本実施形態のストレージシステム3では、ハッシュ値の範囲を細かい単位に分割することによって、移動するデータの量を少なくすることができるが、ハッシュ値の範囲の分割によってハッシュテーブルが大きくなる。しかし、ハッシュテーブルのサイズに上限を設け、ハッシュテーブルのサイズが上限に達したときに、アクセスの少ないハッシュ範囲を隣接するハッシュ範囲と融合することによって、ハッシュテーブルのサイズを小さく（圧縮）することができる。

【0052】

次に、本実施形態のストレージシステム3において、データを移動させる場合の一例を説明する。データの移動は、同一のストレージグループ4内のストレージノード1間で行う。例えば、図3に示したノード構成の一例では、符号“Y”が付与されたストレージグループ4に属する識別番号“1”が付与されたストレージノード1（1Y）から“m”が付与されたストレージノード1（mY）までのいずれか1つのストレージノード1にデータを移動する。このように、本実施形態のストレージシステム3では、同一のストレージグループ4内のいずれか1つのストレージノード1にデータを移動することができる。

30

【0053】

上記に述べたとおり、本実施形態によるストレージシステム3によれば、同一のストレージグループ4内のストレージノード1間でデータを移動させることによって、データ配置の自由度を向上させることができる。また、データ配置の変更を行う理由としては、ストレージシステム3のデータにアクセスするクライアント2とアクセスされるデータを保持しているストレージノード1とにおけるネットワーク上の距離が遠いことが考えられるが、本実施形態によるストレージシステム3では、ネットワークトポロジーに応じてネットワークグループ5を設定しているため、ストレージグループ4内でクライアント2とネットワークグループ5が同じストレージノード1を必ず見つけることができる。このように、ストレージグループ4内でデータを移動することができる自由度があれば、十分にストレージシステム3の効率を向上することができる。そして、データ配置の自由度が向上することによって、省電力モードによる制御も可能となり、ストレージシステム3の省電力化を実現することもできる。

40

【0054】

50

なお、ストレージシステム内の任意のストレージノードにデータを移動することができれば、データ配置の自由度が最大となるが、その場合には、データ配置が変更されたことをストレージシステム内の全てのストレージノードが把握する必要がある。しかし、これは、ストレージシステムの構成が大規模になるにしたがって、全てのストレージノードがデータ配置の変更を把握するためのコストが上昇し、結果としてストレージシステム全体のスケラビリティを大きく損なう要因となる。それに対して本実施形態のストレージシステム3における制御方法であれば、データ配置を変更したストレージグループ4内のストレージノード1のみがデータ配置の変更を把握していればよい。そして、データ配置の変更をしていない他のストレージグループ4内のストレージノード1にとっては、データ配置を変更したストレージグループ4であってもストレージグループとしての変更がされていないため、データ配置の変更を把握する必要がない。従って、データ配置の変更を把握するためのコストを低く抑えることができ、結果としてストレージシステム全体のスケラビリティが向上することとなる。そして、これは、ストレージシステム3のデータにアクセスするクライアント2とアクセスされるデータを保持しているストレージノード1とにおけるネットワーク上の距離を近くするという目的を、十分に達成しているということができる。

10

【0055】

上記実施形態におけるストレージシステムの制御方法では、スケラブルなストレージシステムにおいて、ストレージグループとネットワークグループという直交したグループでストレージノードをグルーピングする。ストレージグループの中ではネットワークポロジを考慮したデータ配置を行えるようにして、データ配置の自由度を向上させることができる。これにより、柔軟なデータ配置をすることができるため、例えば、処理ノードとデータとを近くに配置することによって、ストレージシステムの効率を向上することができる。また、上記実施形態のストレージシステムの制御方法では、データ配置の自由度が向上するため、例えば、データを予め定めたストレージノードに集めることによって、ストレージシステムの低消費電力化を実現する省電力の制御を行うことや、不要なトラフィックを防いでアクセス速度を向上させることができる。

20

【0056】

以上、一実施形態について、図面を参照して説明してきたが、具体的な構成はこの実施形態に限定されるものではなく、種々の変更も含まれる。

30

【0057】

なお、上記実施の形態におけるストレージノード1は、CPU(Central Processing Unit)、メモリ、ハードディスク、ネットワークインタフェース、などを有してもよい。上記ストレージノード1は、一般に広く用いられるサーバコンピュータ、パソコン、などであってもよい。また、上記の実施形態はソフトウェア、もしくは、ハードウェアによって実装されてもよい。

【産業上の利用可能性】

【0058】

上記の実施形態は、例えば、大量のデータを保持するストレージシステムに適用が可能である。

40

【符号の説明】

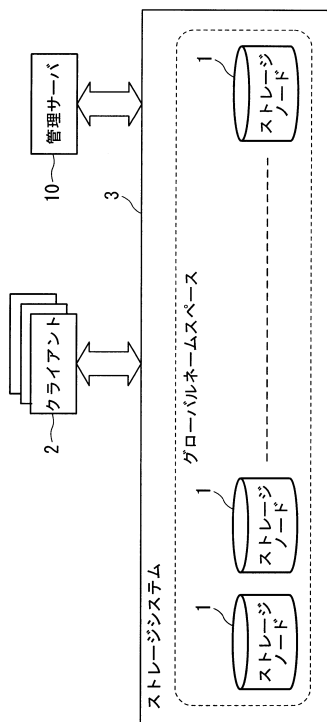
【0059】

- 1・・・ストレージノード
- 2・・・クライアント
- 3・・・ストレージシステム
- 4・・・ストレージグループ
- 5・・・ネットワークグループ
- 6・・・ラック
- 7・・・スイッチ
- 8・・・上位スイッチ

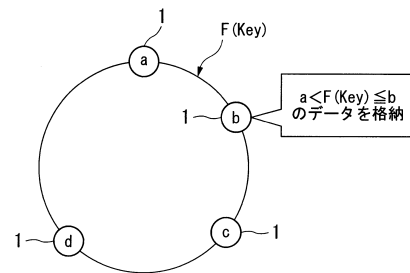
50

10・・・管理サーバ

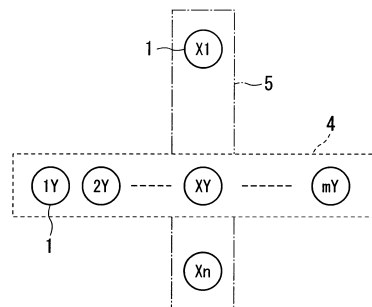
【図1】



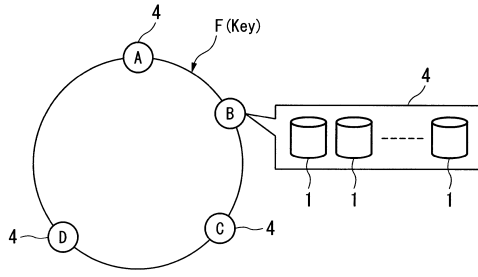
【図2】



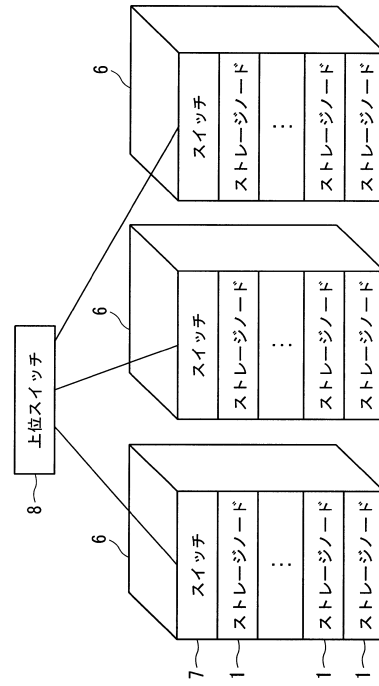
【図3】



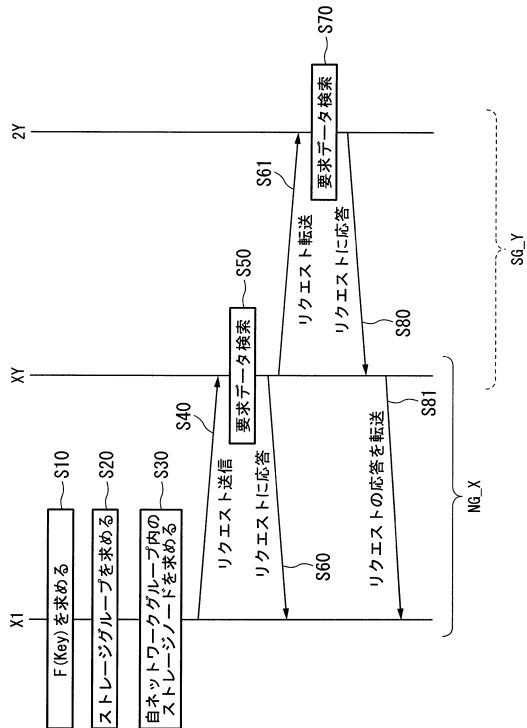
【図4】



【図5】



【図6】



フロントページの続き

(出願人による申告) 国等の委託研究の成果に係る特許出願(平成21年度 独立行政法人新エネルギー・産業技術総合開発機構「グリーンネットワーク・システム技術研究開発プロジェクト(グリーンITプロジェクト) / エネルギー利用最適化データセンタ基盤技術の研究開発 / サーバの最適構成とクラウド・コンピューティング環境における進化するアーキテクチャーの開発 / クラウド・コンピューティング技術」委託研究、産業技術力強化法第19条の適用を受ける特許出願)

前置審査

審査官 加内 慎也

(56)参考文献 特開2008-011330(JP,A)
特開2008-146119(JP,A)
特開2007-280304(JP,A)
特開2010-074604(JP,A)
特開2007-148545(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00
G06F 3/06
G06F 13/10