

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2010-134948
(P2010-134948A)

(43) 公開日 平成22年6月17日(2010.6.17)

(51) Int.Cl. F I テーマコード(参考)
G06F 12/00 (2006.01) G06F 12/00 531M 5B082

審査請求 有 請求項の数 21 O L (全 30 頁)

(21) 出願番号 特願2010-5548 (P2010-5548)
(22) 出願日 平成22年1月14日(2010.1.14)
(62) 分割の表示 特願2006-501087 (P2006-501087) の分割
原出願日 平成16年1月21日(2004.1.21)
(31) 優先権主張番号 60/441,810
(32) 優先日 平成15年1月21日(2003.1.21)
(33) 優先権主張国 米国(US)
(31) 優先権主張番号 10/761,884
(32) 優先日 平成16年1月20日(2004.1.20)
(33) 優先権主張国 米国(US)

(71) 出願人 505052076
イコールロジック, インク.
EQUALLOGIC, INC.
アメリカ合衆国 ニューハンプシャー州
03063 ナシュア タウンゼンド ウ
ェスト 9
9 Townsend West, N
ashua, NH 03063 (US
).
(74) 代理人 100136630
弁理士 水野 祐啓
(72) 発明者 コーニング, ポール, ジー.
アメリカ合衆国 ニューハンプシャー州
03070 ニュー ポストン ジョー
イングリッシュ ロード 408
最終頁に続く

(54) 【発明の名称】 データ記憶管理システム

(57) 【要約】 (修正有)

【課題】 クライアントのリソース要求に適切な応答時間を実現し、クライアントと初期サーバとの長期接続を維持しつつ、クライアント負荷をサーバシステムに迅速に分散する。

【解決手段】 各(等価)サーバ161は負荷モニタ・プロセス22Aを備えており、他の負荷モニタ・プロセスと通信して、サーバシステムへのクライアント負荷及び各サーバへのクライアント負荷を測定し、測定したシステム負荷に回答して一組のリソースを再区分することにより、クライアント負荷を再分散するリソース分散プロセスを更に含む。更に、各サーバは、この区分リソースサーバ上で維持されている各リソースへの参照を含む経路指定テーブル20Aを含むこともできる。クライアントからの要求は、対象となるリソースを維持するか或いは管理しているサーバにそうした要求を経路指定する経路指定テーブルの関数として処理される。

【選択図】 図10

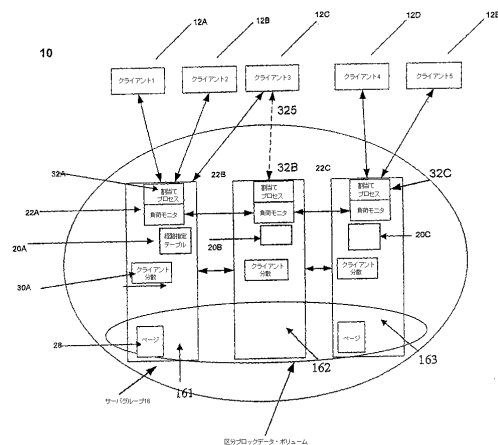


FIG. 10

【特許請求の範囲】**【請求項 1】**

区分記憶サービスを提供するシステムであって、
少なくとも2つのサーバと、
前記少なくとも2つのサーバにわたり区分された記憶ボリュームと、
前記少なくとも2つのサーバのそれぞれで動作する少なくとも2つのスナップショット・プロセスであって、前記区分記憶ボリュームの状態を表す状態情報を生成するため他のスナップショット・プロセスと動作を調整可能な、スナップショット・プロセスを含む、システム。

【請求項 2】

前記スナップショット・プロセスが調整プロセスを含み、当該調整プロセスが、少なくとも他の1つのスナップショット・プロセスと動作を調整して前記区分記憶ボリュームの状態を表す状態情報を生成するためのコマンドを生成する、請求項1に記載のシステム。

【請求項 3】

前記調整プロセスが、スナップショット処理を生成するコマンドにタイムスタンプを打刻するためのタイムスタンプ・プロセスを含む、請求項2に記載のシステム。

【請求項 4】

前記スナップショット・プロセスが、当該サーバにより受信された要求を処理するための要求制御プロセスを含む、請求項1に記載のシステム。

【請求項 5】

前記要求制御プロセスが、当該サーバによる要求の処理を一時中断するための一時中断プロセスを含む、請求項4に記載のシステム。

【請求項 6】

前記要求制御プロセスが、当該サーバにより受信された要求にタイムスタンプを打刻するためのタイムスタンプ・プロセスを含む、請求項4に記載のシステム。

【請求項 7】

前記スナップショット・プロセスが、選択した時刻後に受信された要求を特定するため、一時中断された要求を分析するためのプロセスを含む、請求項1に記載のシステム。

【請求項 8】

前記状態情報を用いて前記記憶ボリュームのコピーを作成するためのアーカイブ・プロセスを更に含む、請求項1に記載のシステム。

【請求項 9】

前記少なくとも2つのサーバにわたって区分された複数の記憶ボリュームを更に含む、請求項1に記載のシステム。

【請求項 10】

区分記憶サービスを提供する方法であって、
少なくとも2つのサーバと当該少なくとも2つのサーバにわたり区分された記憶ボリュームとを提供する段階と、

少なくとも2つのスナップショット・プロセスを前記少なくとも2つのサーバのそれぞれで動作させる段階であって、当該スナップショット・プロセスが、前記区分記憶ボリュームの状態を表す状態情報を生成するため他のスナップショット・プロセスと動作を調整可能な、動作させる段階とを含む、方法。

【請求項 11】

前記区分記憶ボリュームの状態を表す状態情報を生成するため、少なくとも他の1つのスナップショット・プロセスと動作を調整する段階を更に含む、請求項10に記載の方法。

【請求項 12】

調整する前記段階が、スナップショット処理を生成するコマンドにタイムスタンプを打刻する段階を含む、請求項11に記載の方法。

【請求項 13】

10

20

30

40

50

スナップショット・プロセスを動作させる前記段階が、当該サーバにより受信された要求を処理するための要求制御プロセスを動作させる段階を含む、請求項10に記載の方法。

【請求項14】

前記要求制御プロセスが、当該サーバによる要求の処理を一時中断するための一時中断プロセスを含む、請求項13に記載の方法。

【請求項15】

前記要求制御プロセスが、当該サーバにより受信された要求にタイムスタンプを打刻する、請求項13に記載の方法。

【請求項16】

選択した時刻後に受信された要求を特定するため、一時中断された要求を分析する段階を更に含む、請求項13に記載の方法。

【請求項17】

少なくとも2つのサーバにわたって分散されている記憶ボリュームのスナップショットを生成する方法であって、

前記少なくとも2つのサーバ上でスナップショット・プロセスを実行する段階と、

前記スナップショット・プロセスのうち第1のスナップショット・プロセスに管理コマンドを与えて、当該スナップショット・プロセスに、前記区分記憶ボリュームの状態を表す状態情報を生成するよう指示する段階と、

前記第1スナップショット・プロセスに保留要求の実行を停止させると共に、少なくとも第2のスナップショット・プロセスに指示させて、保留クライアント要求の実行を停止させる段階と、

前記第2スナップショット・プロセスに、要求の実行が停止されていることを表明させる段階と、

前記第1スナップショット・プロセスに、そのサーバに維持されている記憶区分の状態を表す状態情報を生成させ、前記第2スナップショット・プロセスにそのサーバに維持されている記憶区分の状態を表す状態情報を生成させるスナップショット・コマンドを生成させる段階とを含む、方法。

【請求項18】

前記管理コマンドが、スナップショットが作成されているデータ・ボリュームをサポートしている第2サーバへの準備コマンドを含む、請求項17に記載の方法。

【請求項19】

前記状態情報を処理して、前記記憶ボリュームの所定期間保存コピーを生成する段階を更に含む、請求項17に記載の方法。

【請求項20】

前記第1及び第2スナップショット・プロセスに、前記状態情報が生成された後に保留中の要求を解放させる段階を更に含む、請求項17に記載の方法。

【請求項21】

ストレージ・エリア・ネットワークであって、

少なくとも2つのサーバを備えたデータ・ネットワークと、

前記少なくとも2つのサーバにわたり区分された記憶ボリュームと、

前記少なくとも2つのサーバのそれぞれで動作する少なくとも2つのスナップショット・プロセスであって、前記区分記憶ボリュームの状態を表す状態情報を生成するため他のスナップショット・プロセスと動作を調整可能な、スナップショット・プロセスとを含む、ストレージ・エリア・ネットワーク。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンピュータ・ネットワークにおいてデータ記憶を管理するためのシステム

及び方法に関し、より詳細には、複数サーバにわたってデータ・リソースを記憶し、複数サーバにわたってデータブロックにバックアップを提供するシステムに関する。

【背景技術】

【0002】

クライアント・サーバ・アーキテクチャは、情報技術における非常に成功した革新の1つである。クライアント・サーバ・アーキテクチャにより、複数クライアントがサーバにより管理されるサービス及びデータ・リソースにアクセス可能となる。サーバはクライアントからの要求をリッスンし、要求に応じて要求を満足することができるかを判断し、必要に応じてクライアントに応答する。典型的な例のクライアント・サーバ・システムは、データファイルを記憶する「ファイルサーバ」設定及びサーバと通信可能な多くのクライアントを備えている。典型的には、クライアントは、サーバが、ファイルサーバにより維持される様々なデータファイルへのアクセスを許可するように要求する。データファイルが利用可能で、クライアントがそのデータファイルへのアクセスを許可されていれば、サーバは要求されたデータファイルをサーバへ引き渡すことによりクライアントの要求を満足する。

10

【0003】

クライアント・サーバ・アーキテクチャは素晴らしい働きをしてきたが、幾つかの欠点を抱えている。例えば、サーバに連絡するクライアントの数及び個別クライアントによる要求の数は、時間の経過と共により大きく変動することがある。従って、クライアントの要求に応答するサーバには、満足できないか或いはほとんど満足できないような要求量が殺到することもある。この問題に対処するため、ネットワーク管理者は、サーバにはクライアント要求の予想ピークレベルに対応できるだけのデータ処理資産を確保してきた。従って、例えば、ネットワーク管理者は、サーバが、着信しうるクライアント・トラフィックの量を処理できるメモリと記憶空間を備えた十分な数の中央処理装置(CPU)を必ず備えるようにしている。

20

【0004】

更に、大容量記憶システムの動作時には、データをどのようにこのシステム上に記憶するかに関する情報を定期的に収集し、記憶データのバックアップ・コピーを時々作成する。こうした情報を収集すると、回復不能な障害が発生した場合の回復を含め、多くの理由で有益である。

30

【0005】

大容量記憶システムのバックアップには、このシステム上に記憶されたデータを読み出して、それを磁気テープに書き込み記憶データの所定期間保存コピーを作成する。

【0006】

しかし、こうした所定期間保存コピーを作成するのは大きな負担となることがある。従来技術による多くのバックアップ作成方法では、バックアップ・コピーの保全性及び無矛盾性を保証するため、システムを進行中の(オンラインの)作業から切り離す必要がある。この理由は、通常のバックアップ技法が、大容量記憶システムからブロックを順次にリニアアクセス・テープにコピーするか、第1ディレクトリの第1ファイルの第1ブロックから開始して、最終ディレクトリの最終ファイルの最終ブロックまで順に進みつつ、この大容量記憶システムのファイルシステムを処理していくかの何れかだからである。何れの場合も、このバックアップ・プロセスは、データがテープに書き込まれる際にアップデートが実行されつつあることは気付かない。

40

【0007】

従って、バックアップ処理を実行しつつ継続的なオンライン作業を許容する場合に、バックアップ処理が進行中にデータが修正変更されるようなことがあると、矛盾が発生する。継続的な記憶作業から記憶システムを切り離すと、システム動作時に矛盾が発生する危険を排除できる。しかし、バックアップ処理には長時間を要することがあるので、システムを作業から切り離すのは望ましくない。

【0008】

50

この問題に対処する1つのアプローチとしては、1つのディスクのデータのミラーすなわち同一コピーを作成することであった。バックアップ処理が必要な時は、ミラーディスクを記憶装置の静的イメージとして用いることができる。この静的イメージが不要になれば（例えば、テープ・バックアップが完了すれば）、ミラーリングがアクティブでなかった時間に起こった変更をミラーディスクにコピーすることで2つのディスクを再同期し、その後、ミラーリングを再開する。

【0009】

リラーリングは有効だが、システムに記憶されているデータを正確に入手する必要がある。しかし、今日では、集中型記憶管理システムを使用しない新たな分散形記憶システムが開発されている。これら分散形システムは、より柔軟でスケーラブルな分散形サーバ・アーキテクチャの利点を利用する。これら記憶システムは非常に素晴らしいが、従来の記憶システムにはなかった難問を提示している。こうした難問の一つは、独立して動作する複数サーバに分散したデータ・ボリュームの信頼性が高く確かな所定期間保存コピーを生成する能力である。

10

【0010】

本開示では、「リソース」という用語は、ファイル、データブロック若しくはページ、アプリケーション、又はサーバからクライアントに提供されるサービス若しくは機能を含むがそれらに限定されないことに注目されたい。本開示では、「資産」という用語は、ハードウェア、メモリ、記憶装置、及びクライアントの要求に应答するためにサーバが使用可能な他の要素を含むがそれらに限定されない。

20

【0011】

必要なシステム・リソースを研究した上で決定しても、クライアント負荷の変動が、単一サーバ又は1つのシステムとして協調しているサーバグループに負担を掛けることがある。例えば、仮に十分なハードウェア資産がサーバシステムに設けられていても、クライアントの要求が特定のファイル、あるファイル内のデータブロック、又はサーバが維持する多のリソースに集中する場合もありうる。従って、上述の例を続けると、クライアントの要求が、ファイルサーバにより維持されているデータファイルの小さな部分に極度に集中することは珍しくない。従って、ファイルサーバが一定量のクライアント要求に対応できるだけのハードウェア資産を持っていたとしても、これらの要求が特定のデータファイルなど特定のリソースに集中すると、目標となっているデータファイルをサポートする資産に過大な負担が掛かる一方、ファイルサーバのほとんどの資産は遊休していることになる。

30

【0012】

この問題に対処するため、ネットワーク技術者達は、クライアント要求を個別の資産に分散するため、利用可能資産にわたってクライアントの要求を分散する負荷バランシング・システムを開発してきた。これを達成するため、負荷バランシング・システムは、要求を利用可能なサーバ資産に均等に分散しようクライアント要求をラウンド・ロビン式に分散できる。多の実現例では、ネットワーク管理者は、特定資産が突然に大量のクライアント要求を受けた時を識別し、その対象となったリソースを複製して、より多くのサーバ資産がそのリソースへのクライアント要求をサポートできるようにするため複製システムを設定している。

40

【0013】

更に、サーバはデータを上手く記憶するが、サーバの資産は限られている。サーバ資産を拡張するために今日用いられている一般的な一技法は、テープライブラリ、RAIDディスクアレイ、及びオプションの記憶システムに依存することである。これらの記憶装置はサーバに適切に接続すれば、データをオンラインでバックアップし、大量の情報を記憶するのに有効である。サーバにこうした装置を多数接続することで、ネットワーク管理者は、かなりの量のデータを記憶可能な「サーバファーム」（多数のサーバ装置及び付属の記憶装置からなる）を構築できる。こうした付属の記憶装置は、ネットワーク接続ストレージ（NAS）システムと集合的に呼ばれる。

50

【 0 0 1 4 】

しかし、サーバファームのサイズが増大し、マルチメディアなどのデータ集中度が高いアプリケーションへの企業の依存度が増大すると、こうした従来の記憶モデルは有用性を維持できなくなる。この理由は、これらの周辺装置へのアクセスが遅くなることがあり、全てのユーザが、常に各記憶装置に容易且つ透過的にアクセスできるとは限らないからである。

【 0 0 1 5 】

この欠点に対処するため、多くのベンダーが、ストレージ・エリア・ネットワーク (SAN) と呼ばれるアーキテクチャを開発している。SANは、NAS型の周辺装置への非常に高速なアクセスを含んだより多くのオプションをネットワーク記憶に提供する。更に、SANは、大量のデータを処理するための別個のネットワークを形成する柔軟性も提供する。

10

【 0 0 1 6 】

SANは、複数ユーザの大きなネットワークに代わって、様々な種類のデータ記憶装置を関連付けられたデータサーバに相互接続する高速の特殊目的ネットワーク又はサブネットワークである。典型的には、ストレージ・エリア・ネットワークは、企業の計算資産のネットワーク全体の一部である。SANは、ディスク・ミラー化、バックアップ及び復元、データの記録及び記録データの取り出し、1つの記憶装置から他の記憶装置へのデータ移送、並びにネットワーク内の異なるサーバ間でのデータ共有をサポートする。SANは、NASシステムを含むサブネットワークを組み込み可能である。

20

【 0 0 1 7 】

SANは、通常は、メインフレームのような他の計算リソースに近接してクラスタ化されているが、非同期転送モード (ATM) 又は同期光通信ネットワーク (SONET) などの広域通信ネットワーク技術を用いて、バックアップ及び超大容量記憶用の遠隔地まで延びることもある。SANは、光ファイバESCON又はファイバチャンネル技術などの既存の通信技術を用いて記憶周辺機器又はサーバに接続することもできる。

【 0 0 1 8 】

SANには大きな将来性があるが、大きな課題に直面している。端的に言って、消費者は自分たちのデータ記憶システムに多くを期待している。具体的には、消費者は、SANがネットワーク・レベルのスケラビリティ、サービス、及び柔軟性を提供する一方、サーバファームに太刀打ちできる速度でデータアクセスを実現することを要求している。

30

【 0 0 1 9 】

これは大きな課題となるかもしれず、とりわけ、特定の情報又は特定のファイルへのアクセスを望むクライアントを、要求した情報又はファイルを持つサーバにリダイレクトする仕組みのマルチサーバ環境では大きな課題となりうる。リダイレクト後に、クライアントは、リダイレクト先のサーバへの新たな接続を確立し、元々通信していたサーバへの接続を切断する。しかし、このアプローチでは、クライアントと最初のサーバとの間に長期間の接続を維持するという利点が生かされない。

【 0 0 2 0 】

もう一つのアプローチは「記憶装置仮想化」或いは「記憶域区分化 (原語: storage partitioning)」であり、中間デバイスをクライアントと一組の物理 (或いは論理) サーバとの間に配置して、中間デバイスが要求の経路指定を実行するというものである。この方法では、何れのサーバも区分されたサービス全体の一部のみを提供していることを意識していないし、何れのクライアントもデータ・リソースが多数のサーバにまたがって記憶されていることを意識しない。言うまでもなく、こうした中間デバイスを追加すると、システムの複雑性が増加してしまう。

40

【 0 0 2 1 】

上述の技法は一定のクライアント・サーバ・アーキテクチャでは上手く機能するが、これらは、クライアント要求とデータ移動とを調整して負荷のバランスをとるため、クライアントとサーバ資産との間に付加的な装置又はソフトウェア (或いは両方) を必要とする

50

。従って、この中央トランザクション・ポイントは、クライアント要求へのサーバ応答を遅くするボトルネックとなってしまうことがある。

【0022】

更に、リソースは、クライアント要求に応答して、待ち時間を極めて最小限にして連続的に供給されなければならない。従って、本発明の分野では、着信するクライアントのリソース要求に適切な応答時間を実現し、クライアントと初期サーバとの長期接続を維持しつつ、クライアント負荷をサーバシステムに迅速に分散するための方法に対する必要性が存在する。更に、本発明の分野では、システム内の異なるサーバにわたって維持されているデータ・ボリュームの確実なスナップショットを提供可能な分散形記憶システムに対する必要性も存在する。

10

発明の概要

【0023】

本発明の様態による、本明細書に記載したシステム及び方法は、複数クライアントからの一組のリソースへのアクセス要求への応答を管理するシステムを含む。一実施形態では、このシステムは、複数の随意選択で等価のサーバを含み、上述した一組のリソースはこれら複数サーバ間で区分されている。各等価サーバは負荷モニタ・プロセスを備えており、各負荷モニタ・プロセスは、他の負荷モニタ・プロセスと通信して、サーバシステムへのクライアント負荷及び各サーバへのクライアント負荷の大きさ測定を生成する。更に、このシステムは、測定したシステム負荷に応答し、上述の一組のリソースを再区分可能とすることにより、クライアント負荷を再分散するリソース分散プロセスを含むことができる。

20

【0024】

随意選択で、このシステムは、測定したシステム負荷に応答し、クライアント接続をサーバシステム間で再区分することにより、クライアント負荷を再分散するクライアント分散プロセスを含むことができる。

【0025】

従って、本明細書に記載したシステム及び方法は、区分サービスとともに動作可能なクライアント分散システムを含む。この区分サービスは、複数の等価サーバによりサポートされており、それぞれの等価サーバは、これらサーバにわたり区分されているサービスの一部を担当する。一実施形態では、各等価サーバは、そのサーバが通信している各クライアントが、当該システム及びそのサーバに掛けている相対負荷を監視できる。従って、各等価サーバは、クライアントがサービスに対して相対的な負担となっている時を特定できる。しかし、区分サービスに関しては、各クライアントは、当該クライアントが求めるリソースを担当する等価サーバと通信する。従って、一実施形態では、本明細書に記載したシステム及び方法は、上記複数サーバにわたってリソースを再分散することで、クライアント負荷を再分散する。

30

【0026】

本発明の別の様態によれば、本明細書に記載のシステム及び方法はサーバグループを含み、このサーバグループは、当該グループの個別サーバにわたって区分されているサービス又はリソースをサポートする。一実施形態では、このシステム及び方法は、記憶サービスを複数のクライアントに提供する区分記憶サービスを提供する。この実施形態では、データ・ボリュームは複数のサーバにわたって区分され、各サーバがデータ・ボリュームの一部を担当する。こうした区分記憶サービスでは、記憶「ボリューム」は、従来の記憶システムにおけるディスク・ドライブと類似したものと理解できる。しかし、この区分サービスでは、データ・ボリュームは幾つかのサーバに分散していて、各サーバが当該ボリューム内のデータの一部を保持している。

40

【0027】

耐故障性、データ・バックアップ、及び他の利点を得るため、本明細書に記載した区分記憶サービスは、記憶装置の管理者に記憶ボリュームの状態のコピーを作成するスナップショット処理及びシステムを提供する。典型的には、このスナップショット処理によって

50

第2の記憶ボリュームが作成される。第2記憶ボリュームは、所与の時刻における記憶システムの状態のアーカイブとして機能する。記憶装置の管理者は、このアーカイブを、元々の記憶ボリュームが後に故障した場合は回復ツールとして、オフライン・バックアップ用のバックアップ・ツールとして、或いはその他の任意適切な理由で使用できる。

【0028】

別の実施形態では、本明細書に記載したシステム及び方法は、記憶資産を企業に提供するために利用できるストレージ・エリア・ネットワーク・システム(SAN)を含む。本発明のSANは、複数のサーバ及び/又はネットワーク・デバイスを含む。これらサーバ及びネットワーク・デバイスの少なくとも一部は、それぞれのサーバ又はネットワーク・デバイスに掛けられたクライアント負荷を監視する負荷監視プロセスを含む。この負荷監視プロセスは、このストレージ・エリア・ネットワーク上で動作する他の負荷監視プロセスと通信することもできる。各負荷監視プロセスは、このストレージ・エリア・ネットワークに掛けられたクライアント負荷を示す全システム負荷分析を生成可能としてもよい。更に、負荷監視プロセスは、そのサーバ及び/又はネットワーク・デバイスに掛けられたクライアント負荷の分析を生成可能としてもよい。負荷監視プロセスが観察したクライアント負荷情報に基づいて、このストレージ・エリア・ネットワークは、クライアント負荷を再分散してクライアント要求に対する応答性を向上できる。これを達成するため、一実施形態では、このストレージ・エリア・ネットワークは、クライアント負荷を再分散するため記憶リソースを再区分できる。別の実施形態では、このストレージ・エリア・ネットワークは、クライアント負荷を当該ストレージ・エリア・ネットワークで再分散するため、システムがサポートするクライアント接続を移動できる。

10

20

【図面の簡単な説明】

【0029】

本発明の上述及び他の目的及び利点は、添付図面を参照すれば次の記載からより完全に理解されるはずである。

【図1】ストレージ・エリア・ネットワーク上に維持されたりソースにアクセスを提供する従来技術システムの構成を概略的に示す。

【図2】本発明によるシステムの機能ブロック図を示す。

【図3】図2のシステムをより詳細に示す。

【図4】サーバグループとして編成されたサーバを備えたクライアント/サーバ・アーキテクチャの概略図である。

30

【図5】クライアントから見たサーバグループの概略図である。

【図6】クライアントと、あるグループのサーバとの間での情報の流れを詳細に示す。

【図7】区分リソース環境におけるリソースの取り出しに関する処理のフローチャートである。

【図8】本発明によるシステムの第1実施形態をより詳細に機能ブロック図として示す。

【図9】図4のシステムと共に使用するのに適した経路指定テーブルの一例を示す。

【図10】本発明によるシステムの第2実施形態をより詳細に機能ブロック図として示す。

【図11】本発明によるシステムの第3実施形態をより詳細に機能ブロック図として示す。

40

【図12】図1のシステムによりサポートされる記憶ボリュームのスナップショットを生成するための処理を示す。

【図13】記憶ボリュームのスナップショットを生成する代替的な処理を示す。

【発明を実施するための最良の形態】

【0030】

本発明の全般的理解のために、幾つかの例示的な実施形態をこれから説明する。しかし、通常の実施形態を備えた当業者であれば、本明細書に記載のシステム及び方法は、分散形ファイルシステム、データベース応用例、及び/又はリソースが区分又は分散される他の用途など他の応用例においてリソースを再分散するために適合及び修正可能であることは理

50

解するはずである。そうした他の追加及び修正は、本発明の範囲に入る。

【0031】

図1には、ローカル・エリア・ネットワーク24を介して通信する複数のクライアント12からのリソース要求をサポートする従来のネットワーク・システムを示した。特に、図1は、複数のクライアント12と、ローカル・エリア・ネットワーク(LAN)24と、クライアントからの要求を処理してサーバ22にそれらを渡す中間装置16を含む記憶システム14とを示す。一実施形態では、中間装置16はスイッチである。このシステムは、マスタ・データテーブル18及び複数サーバ22a乃至22nも含む。記憶システム14は、記憶リソースをLAN24上で動作するクライアント12に提供するストレージ・エリア・ネットワーク(SAN)を提供できる。図1に示したように、各クライアント12は、SAN14に維持されているリソースへの要求20を発することができる。それぞれの要求20はスイッチ16に送信され、このスイッチがそれを処理する。処理時に、クライアント12は、LAN24を介してリソースを要求でき、更にスイッチ16は、マスタ・データテーブル18を用いて、複数サーバ22a乃至22nのどのサーバがクライアント12に要求されているリソースを備えているかを識別する。

10

【0032】

図1では、マスタ・データテーブル18はデータベース・システムとして示されているが、代替的な実施形態では、スイッチ16は、このスイッチが維持するフラットファイル・マスタ・データテーブルを用いてもよい。何れの場合も、スイッチ16は、マスタ・データテーブル18を利用して複数サーバ22a乃至22nの内どのサーバがどのリソースを維持しているかを特定する。従って、マスタ・データテーブル18は、システム14により維持される様々なリソースと、基礎となるサーバ22a乃至22nの何れがどのリソースを担当しているかと、を列記した索引として機能する。

20

【0033】

図1に更に示したように、いったんスイッチ16が要求されたリソースを得るための適切なサーバ22a乃至22nを特定すると、取り出したリソースを識別されたサーバからスイッチ16を介してLAN24に送り、適切なクライアント12に引き渡しできる(矢印21で示した)。従って図1は、システム14が、スイッチ16を、LAN24からの全要求の処理に関わる中央ゲートウェイとして使用することを示している。この中央ゲートウェイ・アーキテクチャを採用すると、クライアント12により要求されたリソースを記憶システム14から引き渡す時間が比較的長くなることがあり、システム14が維持するリソースへの需要増大による待ち時間の増加に従って、この引き渡し時間は増加することがある。

30

【0034】

図2を参照すると、本発明によるシステム10を示した。特に、図2は、複数のクライアント12と、ローカル・エリア・ネットワーク(LAN)24と、複数のサーバ32A乃至32Nを含むサーバグループ30とを示す。図2に示したように、クライアント12はLAN24を介して通信する。図2に示したように、各クライアント12は、サーバグループ30に維持されているリソースを要求できる。ある応用例では、サーバグループ30は、クライアント12にネットワーク記憶リソースを提供するストレージ・エリア・ネットワーク(SAN)である。従って、クライアント12は、図2に要求34として示したように、(LAN)24を介してサーバ(例えば、SAN30のサーバ32Bとして示した)に送信される要求を出すことができる。

40

【0035】

図示したSAN30は、複数の等価サーバ32A乃至32Nを含む。これらサーバは、それぞれ別個のIPアドレスを備えており、従って、システム10は、複数の異なるIPアドレスを含む1つのストレージ・エリア・ネットワークとして見え、それぞれのIPアドレスは、SAN30により維持される記憶リソースにアクセスするためクライアント12が使用できる。

【0036】

50

図示したSAN30は、複数サーバ32A乃至32Nを利用してこのストレージ・エリア・ネットワークにわたってリソースを区分して、区分リソース・セットを形成できる。従って、個別サーバそれぞれは、SAN30が維持するリソースの一部を担当できる。動作時には、サーバ32Bにより受信されたクライアント要求34は、サーバ32Bによって処理され、クライアント12が求めるリソースを特定し、複数サーバ32A乃至32Nのどれがこのリソースを担当しているかを特定する。図2及び3に示した例では、SAN30は、サーバ32Aがクライアント要求34で識別されたリソースを担当することを特定する。図2に更に示したように、随意選択だが、SAN30は、元々のサーバ32Bをクライアント要求34に応答させるのではなく、担当サーバを要求クライアント12に直接的に応答させるというショートカット手法を使ったシステムを採用してもよい。従って、

10

【0037】

上述したように、図2に示したSAN30は、複数の等価サーバを含む。等価サーバは、これに限定するわけではないが、クライアント12などの1つ又は複数クライアントに一樣のインターフェースを提示するサーバシステムであると理解される。これは、図2に示したシステムをより詳細に示す図3に部分的に示されており、図3では、クライアント12からの要求が複数のサーバにより処理可能で、これらサーバは、図示した実施例では適切なクライアントに応答を返す。各等価サーバは、任意のクライアント12が発した要求に同一様態で応答する。そして、クライアント12は、これらサーバの内のどれ(1つ又は複数サーバ)がその要求を処理し、応答を生成するかを知る必要はない。従って、各サーバ32A乃至32Nはクライアント12に同一の応答を与えるので、クライアント12にとっては、サーバ32A乃至32Nの内どれが要求に応答しているかは重要ではない。

20

【0038】

図示したサーバ32A乃至32Nは、それぞれカリフォルニア州サンタクララ所在のサン・マイクロシステムズ社(原語:Sun Microsystems Inc.)が市販するサーバシステムの何れかなどの、従来のコンピュータ・ハードウェア・プラットフォームを含むことができる。各サーバは、1つ又は複数のソフトウェア・プロセスを実行して、このストレージ・エリア・ネットワークを実現する。SAN30は、ファイバチャネル・ネットワーク、アービットレーテッド・ループ、又はストレージ・エリア・ネットワークを提供するのに適したそれ以外の任意種類のネットワーク・システムを使用できる。図2に更に示したように、各サーバはそれ自身の記憶リソースを維持してもよいし、それ自身に接続された1つ又は複数の付加的な記憶装置を含むこともできる。これら記憶装置は、RAIDディスクアレイ・システム、テープライブラリ・システム、ディスクアレイ、又はクライアント12に記憶リソースを提供するのに適したその他の任意装置を含むことができるが、それらに限定されない。

30

【0039】

通常の技能を備えた当業者であれば、本発明のシステム及び方法はストレージ・エリア・ネットワークの応用例に限定されるものではなく、第1サーバが要求を受信し、第2サーバがその要求に対する応答を生成且つ送信するのがより効率的な他の応用例にも適用できることは理解するはずである。他の応用例には、分散形ファイルシステム、データベース応用例、アプリケーション・サービスプロバイダ応用例、又はこの技術から利益を得られるその他の任意応用例が含まれる。

40

【0040】

図4を参照すると、1つ又は複数のクライアント12が、例えばインターネット、イントラネット、WAN、又はLANなどのネットワーク24を介して、或いは直接接続によってサーバグループ116の一部であるサーバ161、162、及び163に接続されている。

【0041】

上述のように、図示したクライアント12は、PCワークステーション、手持ち型計算装置、ワイヤレス通信装置、又はこのサーバグループ116と情報交換するためサーバ

50

ループ 1 1 6 にアクセスして、このサーバと対話可能なネットワーク・クライアント・プログラムを装備した他の装置を含む任意適切なコンピュータ・システムでよい。

【 0 0 4 2 】

システム 1 1 0 が用いるサーバ 1 6 1、1 6 2、及び 1 6 3 は、上述のような、従来の市販サーバ・ハードウェア・プラットフォームでよい。しかし、任意適切なデータ処理プラットフォームを用いてもよい。更に、サーバ 1 6 1、1 6 2、又は 1 6 3 は、テープライブラリ或いはその他の装置のような、ネットワーク 2 4 を介して他のサーバ及クライアントとネットワーク接続しているネットワーク記憶装置を含むことができるのは理解されるはずである。

【 0 0 4 3 】

各サーバ 1 6 1、1 6 2、及び 1 6 3 は、それら動作及び本明細書に記載したトランザクションを実行するソフトウェア構成要素を含むこともでき、又、サーバ 1 6 1、1 6 2、及び 1 6 3 のソフトウェア・アーキテクチャは、用途に従って変更してもよい。特定の実施形態では、サーバ 1 6 1、1 6 2、及び 1 6 3 は、当該サーバのオペレーティング・システムか、デバイスドライバか、アプリケーション・レベル・プログラムか、周辺装置（テープライブラリ、R A I D 記憶システム又は他の記憶装置、或いはそれらの任意の組合せなど）上で動作するソフトウェア・プロセスかに後述するプロセスの一部を組み込むソフトウェア・アーキテクチャを利用してもよい。何れの場合も、通常の技能を備えた当業者であれば、本明細書に記載したシステム及び方法は、多くの異なる実施形態を介して実現でき、更に、採用した実施例及び実現例は対象とする用途の関数として異なることは理解するはずである。従って、これら全ての実施形態及び実現例は本発明の範囲に入る。

【 0 0 4 4 】

動作時には、クライアント 1 2 は、サーバグループ 1 1 6 にわたって区分されたリソースを必要とするはずである。従って、各クライアント 1 2 は要求をサーバグループ 1 1 6 に送信する。典型的には、クライアント 1 2 は独立して動作し、従って、サーバグループ 1 1 6 に掛かるクライアント負荷は時間と共に変化する。こうした典型的な動作では、クライアント 1 2 は、例えばサーバ 1 6 1 などの何れかのサーバに連絡を取り、データブロック、ページ（複数ブロックを含む）、ファイル、データベース、アプリケーション、又は他のリソースなどのリソースにアクセスする。連絡を受けたサーバ 1 6 1 自体が要求されたリソースを保持しておらず、それを管理もしていないこともある。しかし好適な実施形態では、要求を最初に受信したサーバがどれであれ、サーバグループ 1 1 6 は、クライアント 1 2 による全ての区分リソースの利用を可能とするように構成されている。例示目的で、図 4 には、3 つのサーバ全て（サーバ 1 6 1、1 6 2、1 6 3）にわたって区分されている 1 つのリソース 1 8 0 と、これら 3 つのサーバの内の 2 つにわたって区分されている他のリソース 1 7 0 との 2 つのリソースが示されている。システム 1 1 0 がブロックデータ記憶システムであるこの代表的な応用例では、各リソース 1 7 0 及び 1 8 0 は区分ブロックデータ・ボリュームでよい。

【 0 0 4 5 】

従って、図示したサーバグループ 1 1 6 は、複数の等価サーバであるサーバ 1 6 1、1 6 2、及び 1 6 3 からなるストレージ・エリア・ネットワーク（S A N）として動作できるブロックデータ記憶サービスを提供する。各サーバ 1 6 1、1 6 2、及び 1 6 3 は、区分ブロックデータ・ボリューム 1 7 0 及 1 8 0 の 1 つ又は複数部分をサポートできる。図示したサーバグループ 1 1 6 では、2 つのデータ・リソース（例えばボリューム）と 3 つのサーバが存在するが、サーバの数は特に限定されるものではない。同様に、リソース又はデータ・ボリュームの数にも特に制限はない。更に、各リソースは単一サーバ上に全てが収容されていてもよいし、各データ・ボリュームは、サーバグループの全てのサーバ又はサーバグループの部分集合など、幾つかのサーバにわたって区分されていてもよい。

【 0 0 4 6 】

実際には、もちろん、サーバ 1 6 1、1 6 2、及び 1 6 3 に利用できるメモリ資産の量やサーバ 1 6 1、1 6 2、及び 1 6 3 の計算処理上の制限など、実現に関わる事情による

10

20

30

40

50

制限がありうる。更に、一実現例では、グループ分け自体（すなわち、どのサーバがグループを構成するかという決定）が運営上の決定に関わることもある。典型的なシナリオでは、1つのグループが、始めは2、3のサーバのみか或いはたった1つのサーバしか含まないこともありうる。システム管理者は、必要な性能のレベルを確保する必要性に合わせ、サーバをグループに追加していくことになる。サーバを増やせば、記憶されるリソースのためのスペース（メモリ、ディスク記憶装置）が増加し、クライアント要求を処理するCPU処理能力が増加し、クライアントからの要求及びクライアントへの応答を伝送するネットワーク能力（ネットワーク・インターフェース）が増大する。当業者であれば、本明細書に記載したシステムは、追加サーバをグループ116に加えることにより容易にスケール変更して、増大したクライアント需要に対処できることは理解するはずである。しかし、クライアント負荷が変動するにつれ、サーバグループ116はクライアント負荷を再分散して、サーバグループ116内で利用可能な資産をよりよく活用できる。

10

20

30

40

50

【0047】

このため、一実施形態では、サーバグループ116は複数の等価サーバを含む。各等価サーバは、サーバグループ116にわたって区分されたリソースの一部をサポートする。クライアント要求がこれら等価サーバに引き渡されると、等価サーバは互いに動作を調整してシステム負荷の大きさ測定を生成し、各等価サーバに対するクライアント負荷の大きさ測定を生成する。好適な実現例では、この調整はクライアント12には透過的であり、又、これらサーバは、交互にリソースへアクセスさせたり、リソースへアクセスする方法を変更させたりすることなく、互いに負荷を分散できる。

【0048】

図5を参照すると、サーバ161（図4）に接続しているクライアント12は、サーバグループ116を、それが複数IPアドレスを備えた単一サーバであるかのように見ることになる。クライアント12は、サーバグループ116が場合によっては多数のサーバ161、162、163から構築されていることを認識しないし、ブロックデータ・ボリューム170及び180が幾つかのサーバにわたって区分されていることを必ずしも認識しない。あるクライアント12は、単一サーバのみにその固有のIPアドレスを介してアクセスすることもある。結果として、サーバの数及びリソースがサーバ間で区分される様態は、クライアント12が認識するネットワーク環境に影響を与えることなく変更できる。

【0049】

図6は、図5のリソース180がサーバ161、162、及び163にわたって区分されていることを示す。区分サーバグループ116において、任意のボリュームを、サーバグループ116内の任意数のサーバにわたって分散してよい。図4及び5に示したように、1つのボリューム170（リソース1）は、サーバ162、163にわたり分散されており、別のボリューム180（リソース2）は、サーバ161、162、163にわたって分散されている。有利なことに、それぞれのボリュームは、「ページ」とも呼ばれる複数ブロックからなる固定サイズのグループで構成してもよく、代表的な1ページは8192個のブロックを含む。他の適切なページサイズを用いてもよい。又、可変数の（固定数でなく）ブロックを含むページを使用してもよい。

【0050】

代表的な実施形態では、グループ116内の各サーバは、各ボリューム用の経路指定テーブル165を含んでおり、経路指定テーブル165は、特定ボリュームの特定ページが存在するサーバを識別する。例えば、サーバ161が、ボリューム3、ブロック93847への要求をクライアント12から受け取ると、サーバ161は、そのページ番号（例えば、ページサイズが8192個であればページ11）を計算し、経路指定テーブル165においてページ11を含むサーバの位置すなわちサーバ番号をルックアップする。仮にサーバ163がページ11を含んでいる場合は、この要求はサーバ163に転送され、このサーバがデータを読み出して、そのデータをサーバ161に返す。次に、サーバ161は、この要求されたデータをクライアント12に送る。この応答は、常にクライアント12から要求を受け取ったものと同サーバ161を介してクライアント12に返してもよい

。或いは、上述のショートカット・アプローチを用いてもよい。

【0051】

従って、どのサーバ161、162、163がクライアント12が求めるリソースを持っているかは、クライアント12にとっては重要でない。上述のように、サーバ161、162、及び163は経路指定テーブルを用いてクライアント要求に応じ、クライアント12は、どのサーバが要求リソースに関連付けられているかを予め知っている必要はない。これにより、リソースの複数部分が、異なるサーバに存在できるようになる。又、クライアント12を区分サーバグループ116に接続させたまま、リソース又はその部分を移動できる。後者のタイプのリソース再区分を、データブロック又はページからなるリソース部分を移動する場合は「ブロックデータ移送」と本明細書では呼ぶ。通常の技能を備えた当業者であれば、他の種類のリソース（本明細書の他の部分で述べた）からなるリソース部分も同様の手段で移動してよい。従って、本発明は、いかなる特定種類のリソースにも限定されない。

10

【0052】

データの移動は、管理者の命令により又は本明細書で述べた記憶負荷バランシング機構により自動的に実行してもよい。典型的には、データ・リソースのこうした移動又は移送は、ページと呼ぶブロックからなるグループ単位で行われる。

【0053】

ページを1つの等価サーバから別の等価サーバへ移動する時には、応答の待ち時間を発生させたり増加させたりしないように、移動中のページのデータを含む全てのデータをクライアントに継続的にアクセス可能とすることが重要である。手動による移動の場合は、今日の幾つかのサーバで実現されているように、手動移送はクライアントへのサービスを中断してしまう。これは一般に好ましくないと考えられているので、サービス中断を引き起こさない自動移動が好ましい。こうした自動的移送では、移動はクライアントに透過的でなければならない。

20

【0054】

本発明の一実施形態によれば、移送するページは、その移動中は発信サーバ（すなわち、当該データが元々記憶されているサーバ）によって元々「所有」されていると考えられている。クライアントの読み出し要求の経路指定は、引き続きこの発信サーバを介して行われる。

30

【0055】

新たなデータを目的ページに書き込む要求は特別に処理される。すなわち、データは、発信サーバにおけるページ位置と、宛先サーバにおける新たな（コピー）ページ位置との両方に書き込まれる。こうすることで、例え複数の書き込み要求がこの移動時に処理されても、ページの無矛盾イメージが宛先サーバでもたらされる。一実施形態では、図8に示したリソース移送プロセス240がこの処理を実行する。ページが大きくなれば、より綿密なアプローチを用いればよい。こうした場合は、移送は複数部分に分けて実行できる。すなわち、既に移動された部分への書き込みを宛先サーバにリダイレクトし、現在移動中の部分への書き込みは以前のように両方のサーバに向ける。もちろん、まだ移動されていない部分への書き込みは発信サーバが処理すればよい。

40

【0056】

こうした書き込み処理アプローチは、移動中に停電などの障害が発生した場合に必要な動作をサポートするのに必要である。ページが一単位として移動される場合は、打ち切られた（失敗した）書き込みは最初から再開できる。ページが複数部分に分けて移動する場合は、障害発生時に移動中であつた部分からこの移動処理を再開できる。発信サーバと宛先サーバとの両方にデータ書き込む必要があるのは、再開する可能性があるからである。

【0057】

テーブル1は、サーバAからサーバBへの単位ブロックデータ移動に関する一連のブロックデータ移送段階を示す。テーブル2は、部分毎のデータブロック移動に関して同様の情報を示す。

50

【 0 0 5 8 】

【表 1】

| 段階 | 再開段階 | 宛先 | 動作 | 読み出し | 書き込み |
|----|-------|------|------------|------|---------|
| 1 | 該当しない | サーバA | 開始されていない | サーバA | サーバA |
| 2 | 2 | サーバA | 移動開始 | サーバA | サーバA及びB |
| 3 | 2 | サーバA | 移動完了 | サーバA | サーバA及びB |
| 4 | 該当しない | サーバB | 経路指定テーブル更新 | サーバB | サーバB |

10

【 0 0 5 9 】

【表 2】

| 段階 | 再開段階 | 宛先 | 動作 | 読み出し | 書き込み |
|-----|---------------------|------|------------|----------------------------|---|
| 1 | 該当しない | サーバA | 開始されていない | サーバA | サーバA |
| 2 | 2 | | 移動開始、部分1 | サーバA | 部分1に関してはサーバA及びB、他に関してはサーバA |
| 3 | 2 | サーバA | 移動完了、部分1 | 部分1に関してはサーバB、他に関してはサーバA | 部分1に関してはサーバB、他に関してはサーバA |
| 4 | 4 | | 移動開始、部分2 | 部分1に関してはサーバB、他に関してはサーバA | 部分1に関してはサーバB、部分2に関してはサーバA及びB、他に関してはサーバA |
| 5 | 4 | | 移動完了、部分2 | 部分1及び2に関してはサーバB、他に関してはサーバA | 部分1及び2に関してはサーバB、他に関してはサーバA |
| ... | (必要に応じて全部分に処理を繰り返す) | | | | |
| n | 該当しない | サーバB | 経路指定テーブル更新 | サーバB | サーバB |

20

30

【 0 0 6 0 】

リソースが移動されると、経路指定テーブル165(図9を再度参照する)は(本発明の分野では周知の手段により)必要に応じて更新され、その後のクライアント要求は、その要求を現時点で処理する責任を負うサーバに転送されることになる。少なくとも同一リソース170又は180を含むサーバの中では、経路指定テーブル165は、伝播遅延の影響は受けるが同一となりうる。

40

【 0 0 6 1 】

実施形態によっては、経路指定テーブルが一旦更新されると、発信サーバ(又は「ソース」サーバ)におけるページが標準的な手段によって削除される。更に、発信ページ位置に関してフラグ又は他のマーカを発信サーバにセットして、そのデータが有効でないことを少なくとも一時的に示すようにする。発信サーバ宛てのこの時点で潜在的読み出し又は書き込み要求は、そのサーバ上の期限切れデータを読み出すのではなく、エラーとそれに続く再試行をトリガする。こうした再試行が返される時点では、こうした再試行は更新済み

50

の経路指定テーブルに遭遇し、宛先サーバに正しく導かれる。ブロックデータの複製、複写、又は影コピー（これらは本発明の分野では公知の用語である）がサーバグループに残されることはない。随意選択だが、他の実施形態では、発信サーバは、宛先サーバへのポインタ又は他の標識を保持してもよい。発信サーバは、選択した一定期間にわたり、読み出し及び書き込み要求を含むがそれに限定されない要求を、宛先サーバに転送してもよい。この随意選択の実施形態では、こうした要求が非常に遅いか、発信サーバに到着しなくても、クライアント12はエラーを受信しないのは、サーバグループ内の幾つかの経路指定テーブルがまだ更新されていないからである。要求は、発信サーバと宛先サーバとの両方で処理できる。この遅延更新処理は、クライアント要求の処理を経路指定テーブル更新と同期化する必要性を無くすか又は減少させる。経路指定テーブルの更新は背景で実行される。

10

【0062】

図7は、区分サーバ環境でクライアント要求に対応するための代表的な要求対応処理400を示す。要求対応処理400は、ファイル又はファイルのブロックなどのリソースへの要求を受け取ること（ステップ420）により、ステップ410で開始する。要求対応処理400は、ステップ430において経路指定テーブルを調べ、要求されたリソースがどのサーバに位置しているかを特定する。もし要求されたリソースが最初のサーバに存在すれば、ステップ480で最初のサーバが、要求されたリソースをクライアント12に返し、処理400はステップ490で終了する。反対に、要求されたリソースがこの最初のサーバに存在しなければ、ステップ450でこのサーバは、経路指定テーブルからのデータを用いてどのサーバがクライアントに要求されたリソースを実際に保持しているかを特定する。すると、ステップ460で、この要求は要求されたリソースを保持しているサーバに転送され、ステップ480で、このサーバが要求されたリソースを最初のサーバに返す。上述と同様に、処理400はここでステップ480へ進み、最初のサーバが、要求されたリソースをクライアント12へ転送し、ステップ490で処理400は終了する。

20

【0063】

従って、通常の技能を備えた当業者であれば、本明細書に記載したシステム及び方法は、1つ又は複数の区分リソースを複数サーバにわたって移送可能で、従って複数クライアントからの要求を処理可能なサーバグループを提供できることが分かるはずである。幾つかのサーバにこうして移送されるリソースは、ディレクトリ、ディレクトリ内の個別のファイル、又はファイル内のブロック、又はそれらの任意の組合せであってもよい。他の区分サービスも実現可能である。例えば、データベースを類似の様態で区分したり、分散ファイルシステム、或いはインターネットを介して配信されるアプリケーションをサポートする分散サーバ又は区分サーバを提供したりできる。一般に、このアプローチは、クライアント要求がリソース全体の部分への要求であると解釈できる任意のサービスに適用できる。

30

【0064】

図8を参照すると、より効率的なサービスを提供するため、クライアント負荷を再分散可能なシステム500の一実施形態を示す。特に、図8は、クライアント12A乃至12Eがサーバブロック116と通信するシステム500を示す。サーバブロック116は、3つの等価サーバである等価サーバ161、162、及び163を含み、それぞれサーバは、クライアントからの同一要求に実質的に同一の応答を提供できる。典型的には、各サーバは、伝播遅延又は応答タイミングによる差異の影響を受けるが同一の応答を生成する。従って、クライアント12から見れば、サーバグループ116は、クライアント12A乃至12Eと通信するための複数ネットワーク又はIPアドレスを提供する単一のサーバシステムに見える。

40

【0065】

各サーバは、経路指定テーブル200A、200B、及び200Cとして示した経路指定テーブルと、それぞれ負荷モニタ・プロセス220A、220B、及び220Cと、クライアント割当てプロセス320A、320B、及び320Cと、クライアント分散プ

50

ロセス 300A、300B、及び 300C と、それぞれリソース移送プロセス 240A、240B、及び 240C とを含む。更に、例示目的のみだが、図 8 は、リソースを、1 つのサーバから別のサーバへ移送可能な複数ページのデータ 280 として示している。

【0066】

図 8 に矢印で示したように、各経路指定テーブル 200A、200B、及び 200C は、情報を共有する目的で互いと通信できる。上述のように、経路指定テーブルは、個別の等価サーバの内の何れがサーバグループ 116 により維持されている特定リソースを担当するかを探知できる。各等価サーバ 161、162、及び 163 は、クライアント 12 からの同一要求に同一応答を提供できるので、経路指定テーブル 200A、200B、及び 200C (それぞれ) は互いと動作を調整して、異なるリソースとこれらリソースを担当する等価サーバとのグローバル・データベースを提供する。

10

【0067】

図 9 は、経路指定テーブル 200A の一例とそこに記憶されている情報とを示す。図 9 に示したように、各経路指定テーブルは、区分データブロック記憶グループ 116 をサポートする各等価サーバ 161、162、及び 163 の識別子を含む。更に、各経路指定テーブルは、各等価サーバに関連付けられたデータブロックを識別するテーブルも含む。図 9 に示した経路指定テーブルの実施形態では、等価サーバは 2 つの区分ボリュームをサポートする。最初のボリュームは、3 つの等価サーバ 161、162、及び 163 にわたり分散すなわち区分されている。第 2 の区分ボリュームは、2 つの等価サーバ (それぞれサーバ 162 及び 163) にわたって区分されている。

20

【0068】

動作時には、図示した各サーバ 161、162、及び 163 は、サーバグループ 116 に掛けられた全負荷と、各クライアントからの負荷及びそれぞれのサーバ 161、162、及び 163 により処理されている個別のクライアント負荷とを監視できる。これを実行するため、各サーバ 161、162、及び 163 は、それぞれ負荷モニタ・プロセス 220A、220B、及び 220C を含む。上述のように、負荷モニタ・プロセス 220A、220B、及び 220C は互いに通信できる。これは図 8 に、異なるサーバ 161、162、及び 163 の負荷モニタ・プロセスを繋ぐ両方向線で図示した。

【0069】

図示した各負荷モニタ・プロセスは、それぞれのサーバ上で実行し且つそれぞれのサーバが処理しているクライアント要求を監視するソフトウェア・プロセスでよい。これら負荷モニタは、それぞれのサーバが処理している個別クライアント 12 の数、それぞれ及び全てのクライアント 12 が処理している要求の数、及び / 又はデータアクセス・パターン (主として順次データアクセス、主としてランダム・データアクセス、又はその何れでもない) などの他の情報を監視すればよい。

30

【0070】

従って、負荷モニタ・プロセス 220A は、サーバ 161 に掛かるクライアント負荷を表す情報を生成でき、更に、サーバ 162 の負荷モニタ 220B と通信できる。一方、サーバ 162 の負荷モニタ・プロセス 220B は、サーバ 163 の負荷モニタ・プロセス 220C と通信でき、負荷モニタ・プロセス 220C はプロセス 220A と通信できる (図示しない)。異なる負荷モニタ・プロセス 220A、220B、及び 220C 間での通信を可能にすることで、これら負荷モニタ・プロセスは、クライアント 12 によりサーバグループ 116 に掛けられる全システム負荷を特定できる。

40

【0071】

この例では、クライアント 12C は同一リソースへのアクセスを連続的に要求しているかもしれない。例えば、こうしたリソースは、サーバ 161 が維持するページ 280 かもしれない。他の全ての要求とこの負荷が非常に大きく、サーバ 161 が全システム・トラフィックの大きな部分を負担している一方で、サーバ 162 は予期した程度未満しか負担していないこともある。従って、負荷モニタ・プロセス及びリソース割当てプロセスは、ページ 280 をサーバ 162 に移動すべきだと判断し、クライアント分散プロセス 30

50

0 A は、ページ 280 をサーバ 161 からサーバ 162 へ移送するブロックデータ移送プロセス 350 (上述した) を起動できる。従って、図 8 に示した実施形態では、クライアント分散プロセス 300 A は、リソース移送プロセス 240 A と協働して、クライアント 12 C にサーバ 161 ではなくサーバ 162 へ連続的に要求させる状態でリソースを再区分する。

【0072】

一旦、リソース 280 がサーバ 162 に移送されると、経路指定テーブル 200 B はそれ自身を (本発明の分野では周知の標準的手段を用いて) 更新でき、更に、経路指定テーブル 200 A 及び 200 C を再び本発明の分野では周知の標準的手段を用いて更新できる。こうすることで、これらリソースは、クライアント負荷が適切に再分散される可能性が高くなるようにサーバ 161、162、及び 163 にわたって再区分できる。

10

【0073】

図 4 を再び参照すると、これらシステム及び方法は、区分サービスをより効率的に運用するためにも利用できる。

【0074】

この実施形態では、サーバグループ 16 は、複数の等価サーバであるサーバ 161、162、及び 163 からなるストレージ・エリア・ネットワーク (SAN) として動作できるブロックデータ記憶サービスを提供する。各サーバ 161、162、及び 163 は、区分ブロックデータ・ボリューム 188 及 170 の 1 つ又は複数部分をサポートできる。図示したシステム 110 では、2 つのデータ・ボリュームと 3 つのサーバが存在するが、サーバの数は特に限定されるものではない。同様に、リソース又はデータ・ボリュームの数にも特に制限はない。更に、各データ・ボリュームは単一サーバ上に全てが収容されていてもよいし、各データ・ボリュームは、サーバグループの全てのサーバ又はサーバグループの部分集合など、幾つかのサーバにわたって区分されていてもよい。実際には、もちろん、サーバ 161、162、及び 163 に利用できるメモリの量やサーバ 161、162、及び 163 の計算処理上の制限など、実現に関わる事情による制限がありうる。更に、一実現例では、グループ分け自体 (すなわち、どのサーバがグループを構成するかという決定) が運営上の決定となることもある。典型的なシナリオでは、1 つのグループが、始めは 2、3 のサーバのみか或いはたった 1 つのサーバしか含まないこともありうる。システム管理者は、必要なサービスのレベルを確保する必要性に合わせ、サーバをグループに追加していくことになる。サーバを増やせば、記憶されるリソースのためのスペース (メモリ、ディスク記憶装置) が増加し、クライアント要求を処理する CPU 処理能力が増加し、クライアントからの要求及びクライアントへの応答を伝送するネットワーク能力 (ネットワーク・インターフェース) が増大する。当業者であれば、本明細書に記載したシステムは、追加サーバをグループ 116 に加えることにより容易にスケール変更して、増大したクライアント需要に対処できることは理解するはずである。しかし、クライアント負荷が変動するにつれ、後述するように、システム 110 はクライアント負荷を再分散して、サーバグループ 116 内で利用可能な資産をよりよく活用できる。この目的のため、一実施形態では、システム 110 は複数の等価サーバを含む。各等価サーバは、サーバグループ 116 にわたって区分されたリソースの一部をサポートする。クライアント要求がこれら等価サーバに引き渡されると、等価サーバは互いに動作を調整してシステム負荷の大きさ測定を生成し、各等価サーバに対するクライアント負荷の大きさ測定を生成する。好適な一実現例では、この調整はクライアント 12 にとって透過的な様態で行われるので、クライアント 12 は、クライアント 12 とサーバグループ 116 との間で伝送される要求及び応答のみを認識する。

20

30

40

【0075】

図 5 を再び参照すると、サーバ 161 (図 4) に接続しているクライアント 12 は、サーバグループ 116 を、それが複数 IP アドレスを備えた単一サーバであるかのように見ることになる。クライアント 12 は、サーバグループ 116 が場合によっては多数のサーバ 161、162、163 から構築されていることを認識しないし、ブロックデータ・ボ

50

リユーム170、180が幾つかのサーバ161、162、163にわたって区分されていることも認識しない。結果として、サーバの数及びリソースがサーバ間で区分される様態は、クライアント12が認識するネットワーク環境に影響を与えることなく変更できる。

【0076】

図6を参照すると、区分サーバグループ116において、任意のボリュームを、グループ116内の任意数のサーバにわたって分散してよい。図4及び5に示したように、1つのボリューム170（リソース1）は、サーバ162、163にわたり分散されており、別のボリューム180（リソース2）は、サーバ161、162、163にわたって分散されている。有利なことに、それぞれのボリュームは、「ページ」とも呼ばれる複数ブロックからなる固定サイズのグループで構成されており、代表的な1ページは8192個のブロックを含む。他の適切なページサイズを用いてもよい。代表的な実施形態では、グループ116内の各サーバは、各ボリューム用の経路指定テーブル165を含んでおり、経路指定テーブル165は、特定ボリュームの特定ページが存在するサーバを識別する。例えば、サーバ161が、ボリューム3、ブロック93847への要求をクライアント12から受け取ると、サーバ161は、そのページ番号（例えば、ページサイズが8192個であればページ11）を計算し、経路指定テーブル165においてページ11を含むサーバの位置すなわちサーバ番号をルックアップする。仮にサーバ163がページ11を含んでいる場合は、この要求はサーバ163に転送され、このサーバがデータを読み出して、そのデータをサーバ161に返す。次に、サーバ161は、この要求されたデータをクライアント12に送る。言い換えると、この応答は、常にクライアント12から要求を受け取ったものと同サーバ161を介してクライアント12に返される。

10

20

【0077】

クライアント12にとっては、どのサーバ161、162、163に接続しているかは透過的である。実際は、クライアントは、これらサーバをサーバグループ116としか見えず、クライアントはサーバグループ116にリソースを要求する。クライアント要求の経路指定は、それぞれの要求毎に別々に実行されることは理解すべきである。これにより、リソースの複数部分が、異なるサーバに存在できるようになる。又、これによって、クライアントがサーバグループ116に接続している間に、リソース又はその部分を移動することが可能である。もしこれが行われた場合は、経路指定テーブル165は必要に応じて更新され、その後のクライアント要求は、現時点でその要求の処理を担当するサーバに転送される。少なくともリソース170又は180内部では、経路指定テーブル165は同一である。ここで説明する本発明は「リダイレクト」機構とは異なる。リダイレクト機構では、クライアントからの要求を処理できないことはサーバが決定し、クライアントをこの処理が可能なサーバにリダイレクトする。すると、クライアントは別のサーバと新たな接続を確立する。接続確立は比較的効率が悪いので、リダイレクト機構は頻繁な要求の処理には適していない。

30

【0078】

図7は、区分サーバ環境でクライアント要求に対応するための代表的な要求対応処理400を示す。要求対応処理400は、ファイル又はファイルのブロックなどのリソースへの要求を受け取ること（ステップ420）により、ステップ410で開始する。ステップ430で、要求対応処理400は、要求されたリソースがクライアント12から要求を受信した最初のサーバに存在するかを調べ、ステップ430で、経路指定テーブルを調べてどのサーバに要求されたリソースが存在するかを特定する。もし要求されたリソースが最初のサーバに存在すれば、ステップ480で最初のサーバが、要求されたリソースをクライアント12に返し、処理400はステップ490で終了する。反対に、要求されたリソースがこの最初のサーバに存在しなければ、ステップ440でこのサーバは経路指定テーブルを調べ、経路指定テーブルからのデータを用いてどのサーバがクライアントに要求されたリソースを実際に保持しているかを特定する（ステップ450）。すると、ステップ460で、この要求は要求されたリソースを保持しているサーバに転送され、ステップ4

40

50

80で、このサーバが要求されたリソースを最初のサーバに返す。上述と同様に、処理400はここでステップ480へ進み、最初のサーバが、要求されたリソースをクライアント12へ転送し、ステップ490で処理400は終了する。

【0079】

幾つかのサーバに分散されているリソースは、ディレクトリ、ディレクトリ内の個別のファイル、又はファイル内のブロックであってもよい。他の区分サービスを考慮することも可能である。例えば、データベースを類似の様態で区分したり、分散ファイルシステム、或いはインターネットを介して配信されるアプリケーションをサポートする分散サーバ又は区分サーバを提供したりできる。一般に、このアプローチは、クライアント要求がリソース全体の部分への要求であると解釈でき、且つリソースの部分に対する処理が、全ての部分の間におけるグローバル調整（原語：coordination）を必要としないような任意のサービスに適用できる。

10

【0080】

図10を参照すると、ブロックデータ・サービスシステム10の一実施形態を示す。特に、図10は、クライアント12がサーバグループ16と通信するシステム10を示す。このサーバグループ16は、3つのサーバ161、162、及び163を含む。各サーバは、経路指定テーブル20A、20B、及び20Cとして示した経路指定テーブルを含む。各等価サーバ161、162、及び163は、これら経路指定テーブルに加え、図10に示したようにそれぞれ負荷モニタ・プロセス22A、22B、及び22Cを含む。

20

【0081】

図10に示したように、各等価サーバ161、162、及び163は、経路指定テーブル20A、20B、及び20Cを含むことができる。図10に示したように、各経路指定テーブル20A、20B、及び20Cは、情報を共有する目的で互いと通信できる。上述のように、経路指定テーブルは、個別の等価サーバの内の何れがサーバグループ16により維持されている特定リソースを担当するかを探知できる。図10に示した実施形態では、サーバグループ16はSAN又はSANの一部とすることができ、このネットワークでは、各等価サーバ161、162、及び163は、クライアント12がこのSAN上のこの等価サーバにアクセスするのに利用できる個別のIPアドレスを備えている。上述したように、各等価サーバ161、162、及び163は、クライアント12からの同一要求に同一の応答を提供できる。それを達成するため、個別の等価サーバ161、162、及び163の経路指定テーブルは互いに動作を調整して、異なるリソースと（この代表的な実施形態では、データブロック、ページ、或いはデータブロックの他の編成）、それぞれのデータブロック、ページ、ファイル、又は他の記憶編成を担当する個別の等価サーバとのグローバル・データベースを提供する。

30

【0082】

図9を参照すると、代表的な経路指定テーブルを示した。サーバグループ16におけるテーブル20Aのような各経路指定テーブルは、区分データブロック記憶サービスをサポートする各等価サーバ161、162、及び163の識別子（サーバID）を含む。更に、各経路指定テーブルは、各等価サーバに関連付けられたデータブロック、ページを識別するテーブルも含む。図9に示した実施形態では、等価サーバは2つの区分ボリュームをサポートする。最初のボリュームであるボリューム18は、3つの等価サーバ161、162、及び163にわたり分散すなわち区分されている。第2の区分ボリュームであるボリューム17は、2つの等価サーバ（それぞれサーバ162及び163）にわたって区分されている。

40

【0083】

これら経路指定テーブルはシステム10が使用して、利用可能なサーバにわたりクライアント負荷のバランスをとる。

【0084】

各負荷モニタ・プロセス22A、22B、及び22Cは、それぞれの等価サーバに到着する要求パターンを監視して、クライアント12からのパターン又は要求がSANに転送

50

されているか、又、クライアントのサーバへの接続構成を変更することで、これらパターンにより効率的又は確実に応じることができるかを判断する。一実施形態では、負荷モニタ・プロセス 22 A、22 B、及び 22 C は、それぞれの等価サーバに到着するクライアント要求を単に監視する。一実施形態では、各負荷モニタ・プロセスは、個別の要求モニタ・プロセスが認識した異なる要求を表すテーブルを構築する。各負荷モニタ・プロセス 22 A、22 B、及び 22 C は、各等価サーバが認識した要求のグローバル・データベースを構築するために互いに通信可能である。従って、この実施形態では、各負荷モニタ・プロセスは、各等価サーバ 161、162、及び 163 からの要求データを統合して、ブロックデータ記憶システム 16 全体が認識する要求トラフィックを表すグローバル・データベースを生成できる。一実施形態では、このグローバル要求データベースをクライアント分散プロセス 30 A、30 B、及び 30 C が利用可能として、より効率的又は信頼性が高いクライアント接続が可能かどうかを判断するのに使用できるようにする。

10

【0085】

図 10 は、サーバグループ 16 が、クライアント 12 C (サーバ 161 と元々は通信していた) をサーバ 162 に再分散することにより、クライアント負荷を再分散できることを図示している。このため、図 10 は、サーバ 161 がクライアント 12 A、12 B、及び 12 C と通信している初期状態を示す。これは、サーバ 161 をクライアント 12 A、12 B、及び 12 C に繋げる両方向矢印で示した。図 10 で更に示したように、初期状態では、クライアント 12 D 及び 12 E はサーバ 163 と通信しており、サーバ 162 と通信しているクライアントはない (初期状態では)。したがって、この初期状態時には、サーバ 161 は、3 つのクライアント (クライアント 12 A、12 B、及び 12 C) からの要求をサポートする。サーバ 162 は、何れのクライアントからの要求に応じても応答してもしない。

20

【0086】

したがって、この初期状態では、サーバグループ 16 は、サーバ 161 に大きな負担が掛かっているか資産が逼迫していると判断できる。この判断は、サーバ 161 が利用可能な資産からすると、このサーバが過剰に使用されているという分析から導き出される。例えば、ことによると、サーバ 161 のメモリは限られており、クライアント 12 A、12 B、及び 12 C が生成する要求が、サーバ 161 が利用できるメモリ資産に過大な負荷を掛けているのかもしれない。従って、サーバ 161 は、許容限度を下回る動作レベルでクライアント要求に応答しているのかもしれない。或いは、許容レベルで動作し且つクライアント要求に応答してはいるが、サーバ 161 には、サーバ 162 が負担するクライアント負荷 (又は帯域幅) に比べて過大な負担が掛かっているのかもしれない。従って、サーバグループ 16 のクライアント分散プロセス 30 は、全体的効率を向上するには、クライアント負荷を初期状態からサーバ 162 がクライアント 12 C の要求に応じる状態に再分散すればよいと判断するかもしれない。負荷バランス決定を行うのに考慮すべき要件は様々であり、幾つかの例としては経路指定を減少したいという要望がある。すなわち、例えば 1 つのサーバが、リソースの一部 (例えば、ボリューム) が存在する他のサーバよりもかなり多い要求の宛先となっていれば、そのサーバに接続を移動した方が有利となることもある。或いは、サーバ通信負荷のバランスをとることが要望かもしれない。すなわち、任意サーバに対する全通信負荷が他のサーバよりもかなり大きい場合は、この高負荷が掛かったサーバから接続の一部を負荷が軽いサーバに移動すると良いかもしれない。更に、リソース・アクセス負荷 (例えば、ディスク入出力負荷) のバランスをとることも以前の通りだが、通信負荷よりもディスク入出力負荷とする。これは、多数の次元に関わる最適化処理であり、任意組の測定値に関する決定は、管理方針、クライアント活動に関する履歴データ、様々なサーバ及びネットワーク構成要素の能力などに左右される。

30

40

【0087】

これを達成するため、図 10 は、クライアント負荷のこの再分散を、クライアント 12 C とサーバ 162 との連結 325 (両方向の破線矢印で示した) で示している。このクライアント負荷の再分散を実行した後は、クライアント 12 C とサーバ 161 との間の通信

50

路は終了できることは理解されるはずである。

【0088】

クライアント負荷のバランスは、新たなクライアントからの新たな接続にも適用される。クライアント12Fは、それ自身がサーバグループ16により提供されるリソースにアクセスする必要があると判断すると、そのグループとの初期接続を確立する。この接続は、サーバ161、162、又は163の何れかで終端する。このグループはこのクライアントには単一システムに見えるので、161、162、及び163のアドレスの差を意識しない。従って、接続終端点の選択は無作為、ラウンド・ロビン、又は固定でよいが、グループ16内のサーバにおける現在の負荷パターンには応答しない。

【0089】

この初期クライアント接続が受信されると、受信サーバはその時点でクライアント負荷バランス決定を行うことができる。これが行われると、より適切なサーバが選択されることもあり、その場合はこの新たな接続は終了して、このクライアント接続がそれに従って移動される。この場合の負荷バランス決定は、様々なサーバにおける負荷の一般的なレベルや、クライアント12Fが接続を確立した時にクライアント12Fが要求したリソースのカテゴリや、サーバ12Fからのそれまでのアクセス・パターンに関連した、サーバグループ16の負荷モニタが利用可能な履歴データや、サーバグループ16の管理者が設定した方針パラメータなどに基づくことができる。

【0090】

初期クライアント接続を扱う際の別の考慮すべき点は、要求されているリソースの分散である。上述のように、あるリソースは、サーバグループの真部分集合上に分散されているかもしれない。その場合は、クライアント12Fが接続のために最初に選んだサーバは、要求リソースには全く関わりがないかもしれない。こうした接続を受け入れることは可能だが、その場合はこのクライアントからの要求の一部でなく全てが転送を必要とするので、これは特に効率的な構成ではない。そのため、初期クライアント接続のためのサーバを、新たなクライアント12Fが要求するリソースの少なくとも一部に実際に応じるサーバグループ16中のサーバの部分集合から選ぶのが有用である。

【0091】

この決定は、第2の経路指定データベースを導入することにより効率的に行うことができる。上述した経路指定データベースは、対象となっているリソースの別個に移動可能な各部分の正確な位置を指定する。この経路指定データベースのコピーを、そのクライアントが当該リソースへアクセスを要求しているクライアント接続を終端とする各サーバで利用可能にする必要がある。その接続バランス経路指定データベースは、所与のリソース全体に関して、サーバグループ16のどのサーバが現時点でそのリソースの一部を提供するかを単に示す。例えば、図1示したリソース配置を記述する接続バランス経路指定データベースは、2つの項目からなる。リソース17用のものはサーバ162及び163を列記し、リソース18用のものはサーバ161、162、及び163を列記する。

【0092】

図4乃至7を再び参照すると、通常 of 技能を備えた当業者であれば、これらシステム及び方法は本明細書に記載したシステム及び方法に使用可能で、1つ又は複数のリソースを複数サーバにわたって区分可能で、従って複数クライアントからの要求を処理可能なサーバグループを提供できることが分かるはずである。更に、本明細書に記載したシステム及び方法がリソースを再分散又は再区分して、リソースの部分のサーバグループにわたる配分又は分散状況を変更できることが本明細書には記述されている。幾つかのサーバにこうして分散されるリソースは、ディレクトリ、ディレクトリ内の個別のファイル、又はファイル内のブロック、又はそれらの任意の組合せであってもよい。他の区分サービスも実現可能である。例えば、データベースを類似の様態で区分したり、分散ファイルシステム、或いはインターネットを介して配信されるアプリケーションをサポートする分散サーバ又は区分サーバを提供したりできる。一般に、このアプローチは、クライアント要求がリソース全体の部分への要求であると解釈できる任意のサービスに適用してよい。

10

20

30

40

50

【 0 0 9 3 】

図 1 1 を参照すると、サーバ 1 6 1、1 6 2、及び 1 6 3 にわたり区分されている記憶ボリューム 1 8 の分散形スナップショットを生成可能なシステム 1 0 の一実施形態を示す。特に、図 1 1 は、複数のクライアント 1 2 がサーバグループ 1 6 と通信するシステム 1 0 を示す。このサーバグループ 1 6 は、3 つのサーバ 1 6 1、1 6 2、及び 1 6 3 を含む。図 1 1 の実施形態では、サーバ 1 6 1、1 6 2、及び 1 6 3 は、それぞれがクライアントからの同一要求に概ね同一のリソースを提供するという点では等価サーバである。従って、クライアント 1 2 から見れば、サーバグループ 1 6 は、クライアント 1 2 と通信するための複数ネットワーク又は IP アドレスを提供する単一のサーバシステムに見える。各サーバは、経路指定テーブル 2 0 A、2 0 B、及び 2 0 C として示した経路指定テーブルと、スナップショット・プロセス 2 2 A、2 2 B、及び 2 2 C とをそれぞれ含む。更に、例示目的のみだが、図 1 1 は、リソースを、元々の記憶ボリューム 1 8 のイメージである第 2 の記憶ボリュームを生成するためコピー可能な複数ページのデータ 2 8 として示している。

10

【 0 0 9 4 】

図 1 1 に示したように、各経路指定テーブル 2 0 A、2 0 B、及び 2 0 C は、情報を共有する目的で互いと通信できる。上述のように、経路指定テーブルは、個別の等価サーバの内の何れがサーバグループ 1 6 により維持されている特定リソースを担当するかを探知できる。図 1 1 に示した実施形態では、サーバグループ 1 6 は SAN を形成することができ、このネットワークでは、各等価サーバ 1 6 1、1 6 2、及び 1 6 3 は、クライアント 1 2 がこの SAN 上のその等価サーバにアクセスするのに利用できる個別の IP アドレスを備えている。上述したように、各等価サーバ 1 6 1、1 6 2、及び 1 6 3 は、クライアント 1 2 からの同一要求に同一の応答を提供できる。それを達成するため、個別の等価サーバ 1 6 1、1 6 2、及び 1 6 3 の経路指定テーブル 2 0 A、2 0 B、及び 2 0 C は互いに動作を調整して、異なるリソースと、これらリソースを担当する等価サーバとのグローバル・データベースを提供する。

20

【 0 0 9 5 】

図 9 に示したように、各経路指定テーブルは、区分データブロック記憶サービスをサポートする各等価サーバ 1 6 1、1 6 2、及び 1 6 3 の識別子 (サーバ ID) を含む。更に、各経路指定テーブルは、各等価サーバに関連付けられたデータページを識別するテーブルも含む。図 9 に示したように、等価サーバは 2 つの区分ボリュームをサポートする。最初のボリュームであるボリューム 1 8 は、3 つの等価サーバ 1 6 1、1 6 2、及び 1 6 3 にわたり分散すなわち区分されている。第 2 の区分ボリュームであるボリューム 1 7 は、2 つの等価サーバ (それぞれサーバ 1 6 2 及び 1 6 3) にわたって区分されている。

30

【 0 0 9 6 】

図 1 1 を再び参照すると、各サーバ 1 6 1、1 6 2、及び 1 6 3 は、それぞれスナップショット・プロセス 2 2 a、2 2 b、及び 2 2 c を含んでいるのが分かる。各スナップショット・プロセスは、当該サーバシステム上で動作し、記憶ボリュームのそれぞれのサーバが維持する部分のスナップショットを生成するように設計されたコンピュータ・プロセスでよい。従って、図 5 に示したスナップショット・プロセス 2 2 a は、記憶ボリューム 1 8 のサーバ 1 6 1 が維持する部分のコピーを生成する役割を担うことができる。図 1 1 では、この動作をページ 2 8 及びページのコピー 2 9 として少なくとも部分的に示した。

40

【 0 0 9 7 】

動作時には、各等価サーバ 1 6 1、1 6 2、及び 1 6 3 は、概して独立して動作可能である。従って、スナップショット・プロセス 2 2 a、2 2 b、及び 2 2 c は、ある特定時点における記憶ボリューム 1 8 の正確なスナップショットを作成するには動作を調整する必要がある。この調整の必要が発生するのは、書き込み要求が、何れかのクライアント 1 2 a 乃至 1 2 e からサーバ 1 6 1、1 6 2、及び 1 6 3 に随時出されることがあるのが少なくとも部分的にはその理由である。従って、書き込み要求は、スナップショット処理が開始された時に個別のサーバ 1 6 1、1 6 2、及び 1 6 3 により受信される。スナップシ

50

ショット処理が不的確又は予期していない結果を出すのを防ぐため、スナップショット・プロセス22a、22b、及び22cは互いと動作を調整して、特定時点における区分記憶ボリューム18の状態を表す状態情報を生成する。具体的には、一実現例では、スナップショットを作成せよという命令が出された直後の時刻「T」が存在するように時間パラメータを選んで、「T」以前に完了がクライアント12に対して表示される全ての書き込み動作が当該スナップショットに含まれ、「T」以降に完了が表示される書き込み動作は当該スナップショットには含まれないようにする。

【0098】

このため、各スナップショット・プロセス22a、22b、及び22cは、管理者から記憶ボリューム18のスナップショットを作成せよとの要求を受信できる。スナップショット・プロセスは調整プロセスを含み、この調整プロセスは、管理者が対象としている記憶ボリュームをサポートしている他のサーバ上で動作しているスナップショット・プロセスの活動及び動作を調整するためのコマンドを出す。図11に示した例では、管理者は、サーバ162上で動作するスナップショット・プロセス22bにスナップショット・コマンドを出すことができる。このスナップショット・コマンドは、スナップショット・プロセス22bに記憶ボリューム18のスナップショットの作成を要求できる。スナップショット・プロセス22bは経路指定テーブル22bにアクセスして、サーバグループ16の中のサーバで、記憶ボリューム18内のデータブロックの少なくとも一部をサポートしているサーバを特定できる。スナップショット・プロセス22bは、次に記憶ボリューム18の一部をサポートしているサーバそれぞれにコマンドを出すことができる。図11の例では、各サーバ161、162、及び163は記憶ボリューム18の一部をサポートしている。従って、スナップショット・プロセス22bは、スナップショット・プロセス22a及び22bそれぞれにスナップショットを作成する準備をせよとのコマンドを出すことができる。同時に、スナップショット・プロセス22bは、記憶ボリューム18のサーバ162に維持されている部分のスナップショットを作成する準備を開始できる。

【0099】

一実現例では、図7に示したように、スナップショット作成準備のコマンドをスナップショット・プロセス22bから受信したことに応答して、各スナップショット・プロセス22a、22b、及び22cは、実行が差し迫ったクライアントからの要求を一時中断できる。これには、書き込み及び読み出し要求と、これに関わる他の全ての要求を含むことができる。これを実行するため、各スナップショット・プロセス22a、22b、及び22cは要求制御プロセスを含むことができ、この要求制御プロセスが、当該スナップショット・プロセスに、そのサーバにより実行中の要求を処理させる一方、他の要求の実行を一時中断させることで、記憶ボリューム18の状態を変更しかねない書き込み動作を一時停止させる。

【0100】

スナップショット・プロセスは、要求の処理を一時中断した時点で、サーバが記憶ボリューム18のスナップショットを撮る準備ができたことを知らせる応答を、調整役のスナップショット・プロセス22bに出すことができる。調整役のスナップショット・プロセス22bがサーバ22a及び22cから作動可能信号を受信し、自分自身もスナップショット実行の準備が完了していると判断すると、調整役のスナップショット・プロセス22bは、各サーバにスナップショット・コマンドを出すことができる。このスナップショット・コマンドに応答して、サーバは、随意選択で、そのサーバが維持するボリューム18のデータブロックのコピーを表す状態情報を生成するアーカイブ・プロセスを起動できる。一実現例及び一実施形態では、「書き込み時のコピー（原語：copy on write）」プロセスを使ってミラーイメージを作成して、スナップショット作成時から変更されていないボリュームの部分（ページ）が一度記録されるようにする。このミラーイメージは、所望なら後でテープ又は他の超大容量記憶装置に移してもよい。こうした技法は本発明の分野では公知であり、採用する技術は、用途に合わせ又ミラーイメージの量及び他の類似の判断基準に合わせて変更すればよい。

10

20

30

40

50

【0101】

状態情報が一旦作成されると、スナップショット・プロセスは終了され、サーバは一時中断又は保留中の要求を解放して処理できる。

【0102】

図12は、サーバ161、162、及び163にわたり区分されているデータ・ボリュームのスナップショット・イメージを生成するための、本発明による処理を示す。詳しく後述するように、図12に示した分散形スナップショット70により、記憶装置の管理者は、特定時点における記憶ボリューム18の状態を表す情報を生成できる。生成される状態情報には、ファイル構造、記憶データに関するメタデータ、区分記憶ボリュームが維持するデータのコピー又は記憶ボリュームの部分のコピー、或いはその他のこうした情報が含まれる。従って、本明細書に記載したスナップショット・プロセスは、様々な用途が考えられると理解されるはずである。例えば、区分データ・ボリュームの構造に関する情報が作成され、それが後の利用のため記憶されるもの、又、区分記憶ボリュームの完全な所定期間保存対象コピーが作成されるような用途である。本明細書に記載する分散形スナップショット・プロセスを他の用途で用いてもよく、こうした他の応用例も本発明の範囲に入るものと理解されるはずである。

10

【0103】

図12は、1つ又は複数の区分記憶ボリュームの状態情報を生成するためのスナップショット要求を実行する一連の動作を示す時間/空間ダイアグラムを示す。具体的には、図12は、記憶ボリュームの無矛盾の分散形スナップショットを作成する多段階処理70を示す。このため、図12は、図5に示した3つのサーバ162、162、及び163を表す3本の垂線を示す。矢印72乃至78は、1つ又は複数クライアント12からの書き込み要求を示し、矢印82乃至88は、対応するサーバ161、162、および163からの応答を表す。

20

【0104】

図12に示したように、処理70は、スナップショット・コマンドが管理者から出された時に開始される。この例では、スナップショット・コマンドは管理者から出され、サーバ162に渡される。このスナップショット・コマンドは、サーバ162に向けた矢印90として示されている。図12に示したように、サーバ162上で動作するスナップショット・プロセスは、他のサーバ161及び163の動作を調整するコマンドを発することでこのスナップショット・コマンドに対応する。これらコマンドは、サーバ161及び163上で実行するスナップショット・プロセスの動作を調整し、それぞれのサーバが維持しているデータの状態を表す状態情報を記憶ボリューム18の一部として生成する。

30

【0105】

図12に更に示したように、サーバ162上で動作するスナップショット・プロセスは、他のサーバ161及び163に対して準備コマンド92及び94を出す。これらそれぞれのサーバ161及び163上で動作するスナップショット・プロセスは、「準備」コマンドの到着前にクライアントから受信した保留状態の要求(例えば、要求78)と「準備」コマンドの後に受信した要求(例えば、要求76)の実行を停止しておくことで上述の準備コマンドに回答する。

40

【0106】

要求の実行が停止されると、サーバ161及び163は、準備コマンドを出したサーバ162に対する応答として、サーバ161及び163が全ての保留要求の実行を一時停止したことを伝える。調整役のサーバ162は、次にスナップショット・コマンドを各サーバに出す。これは図12で矢印98及び100として示した。

【0107】

このスナップショット・コマンドに回答して、サーバ162に加えサーバ161及び163も、データ・ボリュームのそれぞれのサーバが維持する部分のスナップショットを作成する。次に、このスナップショット情報は、それぞれのサーバのデータファイルに記憶される。随意選択の実現例では、サーバ161、162、及び163それぞれのスナップ

50

ショット・プロセスは、データ・ボリュームの所定期間保存コピーを生成できる。この所定期間保存コピーは、テープ記憶装置又は他の大容量記憶装置に移送できる。

【0108】

生成したスナップショットは、領域104で完了した全ての要求を含むが、領域110で完了した要求は含まない。

【0109】

図13は、記憶ボリュームのスナップショットを生成する処理の代替的实施形態を示す。具体的には、図13は、処理120が3つの期間にわたって起こることを示す空間-時間ダイアグラムである。これら3つの期間は、図13ではこの空間-時間ダイアグラムにおいて異なる陰影を付けた領域として示し、期間122、124、及び126として表示されている。期間122は、管理者がスナップショット要求を出す時刻の前の期間であり、期間124は、このスナップショット要求が出された時刻とスナップショット処理が開始される時刻との間の期間であり、期間128はスナップショットが作成された後の期間である。スナップショットの要求は矢印140で示し、異なる書き込み要求は矢印130乃至138で示した。これら書き込み要求への応答は、矢印131、133、135、137、及び139で示した。図12と同様に、図4に示したシステム10の3つのサーバは、それぞれサーバ161、162、及び163として表示した3本の垂線で示した。

10

【0110】

図13に示した処理120は、タイムスタンプ及び同期システム・クロックの使用を介した無矛盾の分散形スナップショットの作成を示す。具体的には、処理120は、サーバ161、162、及び163が複数の書き込み要求(それぞれが何れかのサーバに随時到着可能)を受信できることを示している。図13では、これを時期122に発生する書き込み要求130、132、及び136として示した。図13に更に示したように、書き込み要求134は時期124においてに到着し、書き込み要求138は時期128において到着できる。従って、図13に示した処理120は、スナップショット処理の前、最中、その後発生する書き込み要求に対処できるように設計されている。

20

【0111】

このスナップショット処理は、スナップショット要求140がサーバ161、162、及び163の少なくとも何れかに受信された時点で開始する。図13は、スナップショット要求140が管理者からサーバ162に送信されていることを示す。スナップショット要求140が受信された時点で、サーバ162上で動作するスナップショット・プロセスは、スナップショットの作成対象であるデータ・ボリュームをサポートする他のサーバに、「準備」コマンドを出すことができる。この準備コマンドは矢印142で示されており、サーバ162からサーバ161及び163に送られる。この準備コマンドが受信されると、サーバ162に加えサーバ161及び163もスナップショット作成の準備をする。この例では、サーバで保留中の要求は保留状態を継続する必要はないので、そのまま進行させ、完了した時点で確認できる。保留する代わりに、サーバ161、162、及び163は、そうした要求が処理された時間を特定し、それぞれの要求にタイムスタンプを打刻できる。図13に示した例では、このタイムスタンプを要求136、134、及び138に打刻する。これら要求は保留中のものか、スナップショット要求140をサーバ162が受信した後に受信されたものである。調整役のサーバ162が各サーバ161及び162から「作動可能」応答を受信すると、調整役サーバ162はスナップショットを撮るコマンドを生成し、このコマンドを待ちサーバ161及び162に伝送する。このコマンドは、時刻が現在のタイムスタンプを含む。これは図13で、サーバ161及び163へのコマンドを表す矢印160及び162として示した。サーバ161及び163がこのコマンドを受信すると、これらサーバは、コマンド161及び162と共に伝送された時間よりも早いタイムスタンプが打刻された書き込み要求をスナップショットに含める。スナップショットを撮れというコマンド160及び162のタイムスタンプより遅いタイムスタンプ付きの書き込み要求は、ここで生成するスナップショットには含まれない。図13に示した例では、書き込み要求136及び134はここで生成するスナップショットに含

30

40

50

れるが、書き込み要求 138 はこのスナップショットに含まれない。このスナップショット情報が生成されると、処理 120 は、図 12 に関して述べた処理 70 と同様に進行する。

【0112】

ハードウェア、ソフトウェア（本発明の分野におけるこれらの用語の現在の定義による）、或いはその任意の組み合わせで本発明の方法を実行できる。特に、任意のタイプの 1 台のコンピューターが複数のコンピューター上で実行されるソフトウェア、ファームウェア、或いは、マイクロコードによって、本方法を実行してもよい。加えて、本発明を具体化するソフトウェアは、任意の形式（例えば、ソースコード、オブジェクトコード、インタープリタコードなど）で任意のコンピューター読み取り可能メディア（例えば、ROM、RAM、磁気メディア、パンチテープ或いはカード、任意形式のコンパクト・ディスク（CD）、DVD など）に格納したコンピューター命令を含んでもよい。その上、こうしたソフトウェアは、インターネットに接続されたデバイス間で転送される周知のウェブページ内に存在するような、搬送波に組み入れられたコンピューター・データ信号の形式をとっていてもよい。従って、本開示で特記しない限り、本発明は、いかなる特定のプラットフォームにも限定されない。

10

【0113】

更に、図示したシステム及び方法は、従来のハードウェア・システムから構築してよく、特別に開発されたハードウェアは必要ない。例えば、図示したシステムでは、クライアントは、ネットワークサーバと情報交換するためこのサーバにアクセスして、このサーバと対話可能なネットワーククライアント・ハードウェア及び/又はソフトウェアを装備した PC ワークステーション、手持ち型計算装置、ワイヤレス通信装置、又は他の装置を含む任意適切なコンピュータ・システムでよい。随意選択だが、これらクライアント及びサーバは、遠隔サーバのサービスにアクセスするにあたって安全が保証されていない通信路に依存してもよい。通信路を安全にするためには、これらクライアント及びサーバは、クライアントとサーバとの間に信頼できるパスを提供するセキュア・ソケット・レイヤー（SSL）安全保護システムなどの安全保護システムを利用すればよい。或いは、これらクライアント及びサーバは、ネットワークを介してデータを伝送する安全なチャンネルを遠隔ユーザに提供するために開発されている他の従来の安全保護システムを用いてもよい。

20

【0114】

更に、本明細書に記載したシステムで使用するネットワークは、インターネットに限定するわけではないがそれを含む、現在知られている或いは将来開発される従来又は将来のコンピュータ間通信システムを含むことができる。

30

【0115】

サーバのサポートには、任意バージョンのユニックス・オペレーティングシステムを実行し、クライアントと接続してデータを交換できるサーバ・プログラムを実行する、サン・マイクロシステムズ社（原語：Sun Microsystems, Inc.）のスパーク（原語：Sparc）（商標）システムなどの市販のサーバプラットフォームを使用してもよい。

【0116】

当業者であれば、ここに記載した実施形態及び実現例の多くの等価物を理解し、或いは、通常の実験を行うだけでそれらを特定できるはずである。例えば、サーバ 161、162、及び 163 の処理或いは入出力機能は同一でよく、又、割当てプロセス 220 は、リソース移送決定を下す際にこれを考慮する。更に、システムネットワーク・トラフィック、入出力要求率、及びデータアクセス・パターン（例えば、アクセスが主として順次アクセスか、主としてランダム/アクセスかなど）における「負荷」の大きさとなるパラメータを幾つか設定してよい。割当てプロセス 220 は、これらパラメータ全てを入力として移送決定で考慮する。

40

【0117】

上述のように、ここに記載した本発明は、ユニックス・ワークステーションなどの従来のデータ処理システム上で動作するソフトウェア構成要素としても実現できる。そうした

50

実施形態では、上述のショートカット応答機構は、C言語コンピュータ・プログラム又はC++、C#、パスカル、フォートラン、Java（登録商標）、又はベーシックを含んだ任意の高レベル言語で書かれたコンピュータ・プログラムとして実装できる。更に、マイクロコントローラ又はデジタル信号プロセッサ（DSP）が使用される実施形態では、これらショートカット応答機構は、マイクロコードで記述したコンピュータ・プログラムとして実現してもよいし、高レベル言語で記述して、使用するプラットフォーム上で実行可能なマイクロコードにコンパイルするコンピュータ・プログラムとして実現してもよい。こうしたコードの開発は当業者には公知であり、そうした技法は、例えば「TMS320ファミリーを用いたデジタル信号処理の応用例、第1、2、及び3巻、テキサス・インスツルメンツ社（1990年）（Digital Signal Processing Applications with the TMS320 Family, Volumes I, II, and III, Texas Instruments (1990)）」に記載されている。更に、高レベルプログラム作成の一般的な技法は公知であり、例えば「スティーブン・G・コーチャン、C言語でのプログラミング、ハイデン・パブリッシング（1983）（Stephen G. Kochan, Programming in C, Hayden Publishing (1983)）」に記載されている。

10

【0118】

以上本発明の特定の実施形態について示し記述してきたが、本発明の種々なる態様から逸脱することなく変更及び修正を行ってもよいことは、当業者に明白となるはずである。従って、添付した特許請求の範囲は、本発明の要旨を逸脱しない範囲に入るものとしてこうした変更及び修正全てを包含することとなる。

20

【0119】

異なる図面において同じ参照符号を使用することで、同様或いは同一の品目を示す。

【図1】

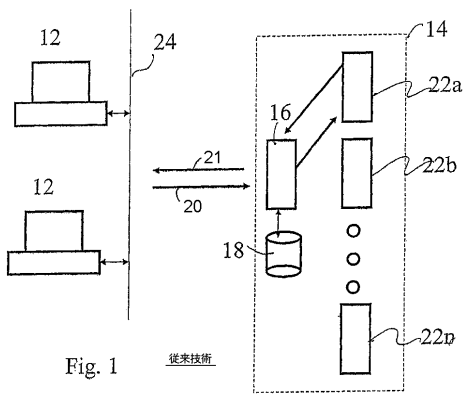


Fig. 1

【図2】

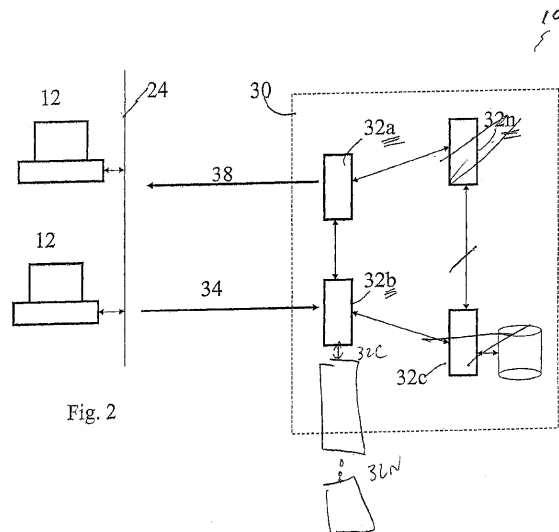


Fig. 2

【 図 3 】

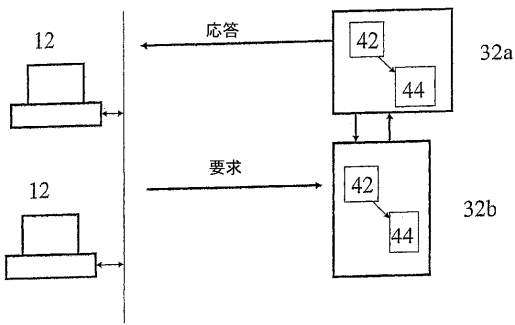


Fig. 3

【 図 4 】

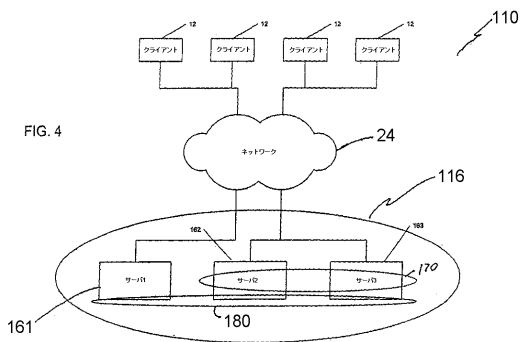
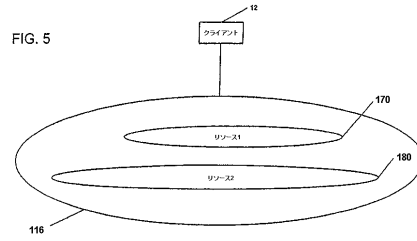


FIG. 4

【 図 5 】



【 図 6 】

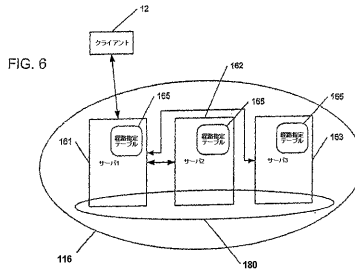


FIG. 6

【 図 7 】

400

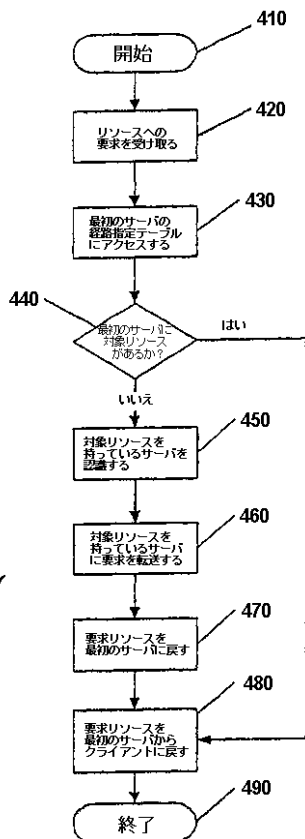


FIG. 7

【 図 8 】

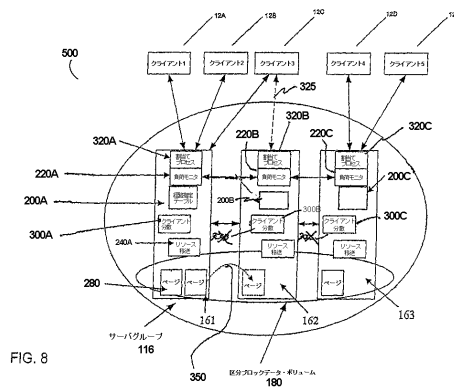


FIG. 8

【 図 9 】

| ボリューム 17 | | ボリューム 18 | |
|----------|-------|----------|-------|
| ページ | サーバID | ページ | サーバID |
| 0 | 1 | 0 | 1 |
| 1 | 2 | 1 | 3 |
| 2 | 2 | 2 | 2 |
| 3 | 1 | 3 | 1 |
| ... | ... | ... | ... |
| 9197 | 2 | 7942 | 1 |
| ... | ... | ... | ... |

| サーバID | サーバ | サーバID | サーバ |
|-------|-----|-------|-----|
| 1 | 162 | 1 | 161 |
| 2 | 163 | 2 | 162 |
| | | 3 | 163 |

FIG. 9

【図 10】

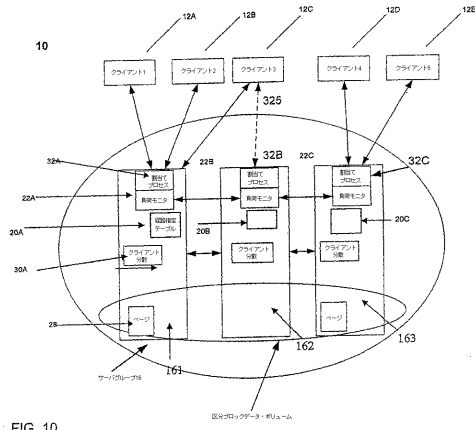


FIG. 10

【図 11】

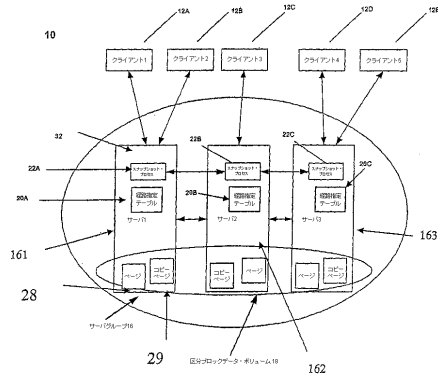


FIG. 11

【図 12】

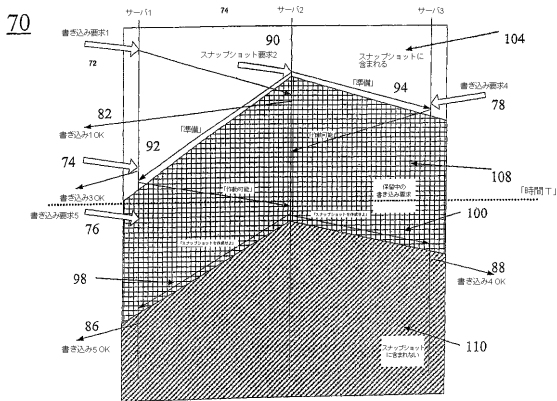


FIG. 12

【図 13】

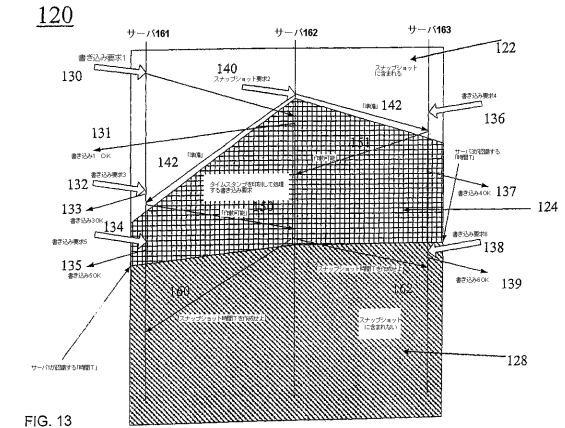


FIG. 13

フロントページの続き

- (72)発明者 ハイデン, ピーター, シー.
アメリカ合衆国 ニューハンプシャー州 03057 マウント パーノン パーゲイトリー ロ
ード 17
- (72)発明者 ロング, ポーラ
アメリカ合衆国 ニューハンプシャー州 03049 ホリス ウィンチェスター ドライブ 2
5
- (72)発明者 スマン, ダニエル, イー.
アメリカ合衆国 マサチューセッツ州 01886 ウェストフォード スイート 201 グリ
ズリー ベアー サークル 11
- (72)発明者 リー, シン, エイチ.
アメリカ合衆国 ニューハンプシャー州 03063 ナシュア タウンゼンド ウェスト 9
- Fターム(参考) 5B082 DA02 DE06