

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5179218号  
(P5179218)

(45) 発行日 平成25年4月10日(2013.4.10)

(24) 登録日 平成25年1月18日(2013.1.18)

(51) Int.Cl.		F I
<b>HO4L 12/851 (2013.01)</b>		HO4L 12/56 200E
<b>HO4L 12/927 (2013.01)</b>		GO6F 13/38 340A
<b>GO6F 13/38 (2006.01)</b>		HO4L 12/56 100A
<b>HO4L 12/707 (2013.01)</b>		

請求項の数 7 (全 15 頁)

(21) 出願番号	特願2008-37692 (P2008-37692)	(73) 特許権者	000004226
(22) 出願日	平成20年2月19日 (2008.2.19)		日本電信電話株式会社
(65) 公開番号	特開2009-200611 (P2009-200611A)		東京都千代田区大手町二丁目3番1号
(43) 公開日	平成21年9月3日 (2009.9.3)	(73) 特許権者	504176911
審査請求日	平成23年1月26日 (2011.1.26)		国立大学法人大阪大学
			大阪府吹田市山田丘1番1号
		(74) 代理人	100086232
			弁理士 小林 博通
		(74) 代理人	100104938
			弁理士 鶴澤 英久
		(74) 代理人	100140361
			弁理士 山口 幸二
		(74) 代理人	100096459
			弁理士 橋本 剛

最終頁に続く

(54) 【発明の名称】 iSCSIセッションのTCPコネクション数制御方法、iSCSIホスト装置、およびiSCSIイニシエータの構成プログラム

(57) 【特許請求の範囲】

【請求項1】

iSCSIホスト装置とiSCSIデバイス装置とがiSCSIセッションによって接続されるときに

iSCSIホスト装置に構築されたイニシエータ機能をもってTCPコネクション数を制御する方法であって、

iSCSIセッションをとおして送受信する一定のデータ量と当該データ量を送受信するのに要した時間との比を参照して、次に送受信する一定のデータ量を送受信するために使用するTCPコネクション数を決定する第1ステップと、

前記決定に基づき次に送受信するために使用するTCPコネクション数を増減制御する第2ステップと、

を有することを特徴とするiSCSIセッションのTCPコネクション数の制御方法。

【請求項2】

前記第1ステップは、TCPコネクション毎に一定のデータ量に振り分けて前記iSCSIデバイス装置と送受信を行うステップと、

前記iSCSIデバイス装置側から送信されたTCPコネクション毎の送受信完了通知に基づき前記送受信に要した時間を求めるステップと、

を有することを特徴とする請求項1記載のiSCSIセッションのTCPコネクション数の制御方法。

【請求項3】

10

20

前記第1ステップは、TCPコネクション数を増加させながらデータ送受信のスループット値を計測し、該スループット値が前回値よりも低下したときにTCPコネクション数の数値巾を抽出するステップと、

前記数値巾に黄金探索法を適用してTCPコネクション数を絞り込んで、スループット値が最大となるTCPコネクション数の最適値を決定するステップと、

を有することを特徴とする請求項1または2いずれか記載のiSCSIセッションのTCPコネクション数の制御方法。

【請求項4】

iSCSIデバイス装置とiSCSIセッションによって接続されるイニシエータ機能が構築されたiSCSIホスト装置であって、

iSCSIセッションをとおして送受信する一定のデータ量と当該データ量を送受信するのに要した時間との比を参照して、次に送受信する一定のデータ量を送受信するために使用するTCPコネクション数を決定する手段と、

前記決定に基づき次に送受信するために使用するTCPコネクション数を増減制御する手段と、

して機能することを特徴とするiSCSIホスト装置。

【請求項5】

TCPコネクション毎に一定のデータ量に振り分けて前記iSCSIデバイス装置と送受信を行う手段と、

前記iSCSIデバイス装置側から送信されたTCPコネクション毎の送受信完了通知に基づき前記送受信に要した時間を求める手段と、

して機能することを特徴とする請求項4記載のiSCSIホスト装置。

【請求項6】

TCPコネクション数を増加させながらデータ送受信のスループット値を計測し、該スループット値が前回値よりも低下したときにTCPコネクション数の数値巾を抽出する手段と、

前記数値巾に黄金探索法を適用してTCPコネクション数を絞り込んで、スループット値が最大となるTCPコネクション数の最適値を決定する手段と、

して機能することを特徴とする請求項4または5いずれか記載のiSCSIホスト装置。

【請求項7】

請求項4～6記載のiSCSIホスト装置として、コンピュータを機能させることを特徴とするiSCSIイニシエータの構成プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ファイルサーバなどのiSCSIホスト装置とストレージなどのiSCSIデバイス装置とをIPネットワークを介して接続するiSCSIの技術に関し、特にスループットを向上させる技術に関する。

【背景技術】

【0002】

(1) iSCSI (Internet SCSI)

単体、もしくは複数のハードディスクドライブ装置の集合体、あるいは専用の制御部で複数のハードディスクドライブを制御するディスクアレイストレージ装置等から構成されるストレージ装置が実用化されている。

【0003】

ストレージ装置(デバイス装置)とホスト装置(サーバ, PC)間を接続するインタフェース技術としてはファイバチャネル(FC), インフィニバンド(InfiniBand), SCSI (Small Computer Systems Interface) などがある。

10

20

30

40

50

## 【 0 0 0 4 】

SCSI インタフェースは近距離間を経済的に接続する用途として優れ、既に広く普及している。当該インタフェースでは、クライアントに相当する機能をイニシエータ、サーバに相当する機能をターゲットと呼び、イニシエータとターゲット間で、SCSI コマンド (CDB: COMMAND Descriptor Block) を交換する。

## 【 0 0 0 5 】

iSCSI (internet SCSI) 技術は、非特許文献 1 に示すように、2004 年にネットワークプロトコルである TCP/IP の上位プロトコルとして、IETF (The Internet Engineering Task Force) で標準化された (RFC 3720)。この iSCSI 接続は、ファイバチャネルなどの他のネットワークストレージプロトコルと違ってコンピュータにイーサネット (登録商標) インタフェースがあれば使用することができる。

10

## 【 0 0 0 6 】

現在、iSCSI 技術は、ギガビット以上の回線品質、特にデータセンター内、データセンター間や大規模 LAN システム等の基幹回線で利用されている。非特許文献 1 によれば、iSCSI 技術は TCP コネクションを複数利用する方法を規定しているが、オプション機能のため必ずしも実装されているわけではない。

## (2) iSCSI 高速化技術

一方で、広域・広帯域の IP ネットワーク上で、ネットワーク遅延による iSCSI スループットの低下を防ぐ方法として、TCP コネクションの複数利用の有効性は知られている (非特許文献 2、3)。また、TCP コネクションの複数利用の実現手段としてマルチリンクを用いる手法の非特許文献 4 やトランスポート層のプロトコルを変更する非特許文献 5 など提案されている。

20

## 【 0 0 0 7 】

例えば、マルチリンクを用いる手法に関する非特許文献 4 では、VPN のマルチホーミング機能を用いて、並列 TCP コネクションを複数の経路で確立させ、スループット向上を実現する手法が提案されている。

## 【 0 0 0 8 】

また、トランスポート層のプロトコルを変更する手法に関する非特許文献 5 では、複数の LAN ポートを利用して物理的にリンクを多重化し、TCP コネクションを異なる経路で確立することによりスループットの向上を図っている。

30

## 【 0 0 0 9 】

しかしながら、非特許文献 4 では、マルチリンクを利用できる環境に制限があり、非特許文献 5 ではマルチベンダー機器への接続が困難であるという問題を有している。

【非特許文献 1】RFC 370 “Internet Small Computer Systems Inter-face (iSCSI)”, Apr 2004

【非特許文献 2】Q. K. Yang, “On performance of parallel iSCSI protocol for networked storage systems,” in Proceedings of the 20th International Conference on Advanced Information Networking and Applications (AINA 2006), vol. 1, pp. 629 - 636, Apr. 2006.

40

【非特許文献 3】G. Motwani and K. Gopinath, “Evaluation of advanced TCP stacks in the iSCSI environment using simulation model,” Proceedings of the 22nd IEEE/13th NASA Goddard Conference on Mass Storage Systems and Technologies (MSST 05)}, pp. 210 - 217, 2005.

【非特許文献 4】B. K. Kancheria, G. M. Narayan, and

50

K. Gopinath, "Performance evaluation of Multiple TCP connection in iSCSI," in Proceedings of the 24th IEEE Conference on Mass Storage Systems and Technologies, pp. 239 - 244, IEEE Computer Society, Sept. 2007

【非特許文献5】千島 望, 山口 実靖, 小口 正人, "iSCSIストレージにおけるVPN複数経路アクセス時の性能とTCPパラメータ解析," 電子情報通信学会第18回データ工学ワークショップ DEWS2007, Feb. 2007.

【非特許文献6】L. Qiu, Y. Zhang, and S. Keshav, "On individual and aggregate TCP Performance," in Proceedings of Internet Conference on Network Protocols, pp. 203 - 212, Oct. 1999.

【非特許文献7】T. Ito, H. Ohsaki, and M. Imase, "On parameter tuning of data transfer protocol GridFTP in wide-area Grid computing," in Proceedings of Second International Workshop on Networks for Grid Applications (GridNets 2005)}, pp. 415 - 421, Oct. 2005.

【非特許文献8】T. Ito, H. Ohsaki, and M. Imase, "GridFTP - APT: Automatic parallelism tuning mechanism for data transfer protocol GridFTP," in Proceedings of 6th IEEE International Symposium on Cluster Computing and the Grid (CCGrid2006)}, pp. 454 - 461, May. 2006.

【発明の開示】

【発明が解決しようとする課題】

【0010】

しかしながら、広域・広帯域のIPネットワークにおいてiSCSIのスループットが低下する問題は広く知られているものの、この問題を従来は効果的に解決できなかった。

【0011】

すなわち、広域・広帯域のIPネットワークを前提としたiSCSI接続では、スループット向上のため、利用する回線の帯域や遅延情報に基づき予めiSCSI/TCPなどのパラメータをチューニングし、あるいはサーバとストレージ間にプロトコル変換機を設置してスループットを改善するなど、運用設定が煩雑で、かつ経済的な運用が困難なおそれがあった。

【0012】

また、このような改善は、固定的な情報に基づく設定であるため、帯域の変動に追従できず、効率的なデータ転送が難しく、またベストエフォート型の経済性の高い回線を利用できないおそれもあった。

【0013】

そこで本発明は、広域・広帯域のIPネットワークにおいてiSCSIのスループットが低下する問題を、(1)既存のストレージ装置を利用しつつ、(2)ネットワークの遅延、利用帯域などの利用環境を限定することなく、(3)経済的に解決することを課題としている。

【課題を解決するための手段】

【0014】

10

20

30

40

50

i S C S I 技術は、一つの i S C S I セッション内に複数の T C P コネクションを確立して、複数のコマンドを並列に発行し、各 T C P コネクションで同時にデータの送受信を行う機能を有しているものの、前記従来例には T C P コネクション数をどのように決定するかは記述されていない。

【 0 0 1 5 】

本発明は、並列 T C P コネクションのスループットが多重度に関して上に凸の関数になる性質に着目して創作された技術的思想であり、T C P コネクションの多重度をスループットに対して最適化させて前記課題を解決している。

【 0 0 1 6 】

具体的には、請求項 1 記載の発明は、i S C S I ホスト装置と i S C S I デバイス装置とが i S C S I セッションによって接続されるときに i S C S I ホスト装置に構築されたイニシエータ機能をもって T C P コネクション数を制御する方法であって、i S C S I セッションをとおして送受信する一定のデータ量と当該データ量を送受信するのに要した時間との比を参照して、次に送受信する一定のデータ量を送受信するために使用する T C P コネクション数を決定する第 1 ステップと、前記決定に基づき次に送受信するために使用する T C P コネクション数を増減制御する第 2 ステップと、を有することを特徴としている。

10

【 0 0 1 7 】

請求項 2 記載の発明は、前記第 1 ステップは、T C P コネクション毎に一定のデータ量に振り分けて前記 i S C S I デバイス装置と送受信を行うステップと、前記 i S C S I デバイス装置側から送信された T C P コネクション毎の送受信完了通知に基づき前記送受信に要した時間を求めるステップと、を有することを特徴としている。

20

【 0 0 1 8 】

請求項 3 記載の発明は、前記第 1 ステップは、T C P コネクション数を増加させながらデータ送受信のスループット値を計測し、該スループット値が前回値よりも低下したときに T C P コネクション数の数値巾を抽出するステップと、前記数値巾に黄金探索法を適用して T C P コネクション数を絞り込んで、スループット値が最大となる T C P コネクション数の最適値を決定するステップと、を有することを特徴としている。

【 0 0 1 9 】

請求項 4 記載の発明は、i S C S I デバイス装置と i S C S I セッションによって接続されるイニシエータ機能が構築された i S C S I ホスト装置であって、i S C S I セッションをとおして送受信する一定のデータ量と当該データ量を送受信するのに要した時間との比を参照して、次に送受信する一定のデータ量を送受信するために使用する T C P コネクション数を決定する手段と、前記決定に基づき次に送受信するために使用する T C P コネクション数を増減制御する手段と、して機能することを特徴としている。

30

【 0 0 2 0 】

請求項 5 記載の発明は、T C P コネクション毎に一定のデータ量に振り分けて前記 i S C S I デバイス装置と送受信を行う手段と、前記 i S C S I デバイス装置側から送信された T C P コネクション毎の送受信完了通知に基づき前記送受信に要した時間を求める手段と、して機能することを特徴としている。

40

【 0 0 2 1 】

請求項 6 記載の発明は、T C P コネクション数を増加させながらデータ送受信のスループット値を計測し、該スループット値が前回値よりも低下したときに T C P コネクション数の数値巾を抽出する手段と、前記数値巾に黄金探索法を適用して T C P コネクション数を絞り込んで、スループット値が最大となる T C P コネクション数の最適値を決定する手段と、して機能することを特徴としている。

【 0 0 2 2 】

請求項 7 記載の発明は、請求項 4 ~ 6 記載の i S C S I ホスト装置として、コンピュータを機能させることを特徴としている。

【 発明の効果 】

50

## 【0023】

請求項1～7記載の発明によれば、TCPコネクション数の最適値が自動的に算出され、当該値でTCPコネクション数が運用される。したがって、運用設置が簡略化され、帯域変動にも追従でき、これにより与えられた回線で常に効率的なデータ転送が実現可能となり、経済性も向上する。

## 【0024】

特に、請求項7記載の発明によれば、既存のiSCSIホスト装置にプログラムをインストールすれば足り、システム構築・更新の経済性が一層向上する。

## 【発明を実施するための最良の形態】

## 【0025】

以下、図1～図5に基づき本発明の実施形態を説明する。図1は、iSCSIイニシエータ1とiSCSIターゲット2とがIP(internet protocol)ネットワークを介して、iSCSI接続される基本的な形態を示している。

## 【0026】

前記iSCSIイニシエータ1は、iSCSIホスト装置を構成するファイルサーバに実装される一方、前記iSCSIターゲット2は、iSCSIデバイス装置を構成するディスクストレージ装置に実装されている。ここではiSCSI接続(iSCSI session)の中にTCPコネクションを複数形成することができ、TCPコネクション数を「N」と表記する。

## 【0027】

前記iSCSIイニシエータ1の機能は、ファイルサーバの処理部(CPU「Central Processing Unit」など)が、イニシエータ構成プログラムのプログラムコードを読み込んで、通信デバイス(例えばイーサネットインタフェースなど)を通じて実現構築されている。

## 【0028】

なお、前記ファイルサーバは、コンピュータにより構成され、処理データなどを一時記憶可能なメモリ(RAM)、キーボードやマウスなどの入力デバイス、モニタなどの表示部などを有している。

## 【0029】

図2は、前記iSCSIイニシエータ1の機能にTCPコネクション数Nの制御機能(図2中のiSCSI-APT「Automatic Parallelism Tuning」3)が追加された構成を示している。この制御機能は、前記iSCSIイニシエータ1の機能と同様に、ファイルサーバの前記処理部などが前記イニシエータ構成プログラムのプログラムコードを読み込んで実現構築されている。

## 【0030】

ここではiSCSIイニシエータ1を実装したファイルサーバが、iSCSIターゲット2を実装したストレージ装置に書き込む(転送する)データを、TCPコネクション数Nが変動してもほぼ一定のサイズのデータ(チャンク)に分割しながら転送する。

## 【0031】

このときチャンクは、「iSCSI PDU」サイズの整数倍となるため、TCPコネクション数「N」に対して、各TCPコネクションにほぼ均等に「n」個を振り分けるとすると以下の式1が成立する。

## 【0032】

## 【数1】

$$\text{チャンクサイズ} = \text{iSCSI PDU 数} \times n \times N \quad \text{式(1)}$$

## 【0033】

以下、図2に基づき動作例S1～S7を説明する。

## (1) S1～S7の動作例

S1の処理：まず、iSCSIイニシエータ1の配置されるファイルサーバ内の上位層、例えばファイルシステムをとおして、iSCSI層に書き込み要求のあった転送データを、TCPコネクション数Nに依存しない均一サイズのデータ(チャンク)に分割する。チャンクの大きさは、小さすぎるとTCP/IPプロトコルの特性上、スループットを正確に反映できないため、スループットをほぼ正確に算出できる十分な大きさとする。

## 【0034】

S2の処理：本発明の処理機能であるiSCSI-APT3が、iSCSIターゲット2(ストレージ)に対してiSCSI接続するために、初期のTCPコネクションを確立し、「iSCSI login」認証を行う(通常ここで接続されるTCPコネクションは1本である。)

10

## 【0035】

S3の処理：S1の処理で生成したチャンクを、iSCSIのデータ転送単位である「iSCSI PDU」サイズに分割し、分割した「iSCSI PDU」をTCPコネクションN本に対して、ほぼ均等「n個」に振り分ける。

## 【0036】

## 【表1】

N	n:iSCSI PDU数/TCPコネクション数N
1	840
2	420
3	280
4	210
5	168
6	140
7	120
8	105

20

Table 1: TCPコネクション数Nとチャンク分割の例

30

## 【0037】

表1は、「Table 1」のTCPコネクション数Nとチャンク分割の例を示している。ここでは「チャンク(n)：iSCSI PDU数/TCPセッション」を例に、TCPコネクション「N」の増加に伴い、1本のTCPコネクションあたり「1/N」に減少することが示されている。

## 【0038】

S4の処理：振り分けられた各TCPコネクション毎に転送する「iSCSI PDU」(約n個)を各々のTCPコネクションをとおして転送する。このときiSCSI-APT3は、図3及び図4に示すように、すべてのTCPコネクションの中で一番最初に転送を始めた時刻を記録する。

40

## 【0039】

S5の処理：SCSIターゲット2が、図3に示すように、各TCPコネクション毎に「iSCSI PDU」受信完了通知、即ちR2T(Ready to transfer)を逐次送信する。

## 【0040】

S6の処理：iSCSI-APT3は、すべてのTCPコネクションの中で、最後にR2Tを受信した時刻が記録し、図4に示すように、当該チャンクの転送時間から当該チャンク転送でのスループットG(N)を算出する。直後に、次のチャンク転送に設定するTCPコネクション数Nを後述の内蔵アルゴリズム処理に基づいて決定する。

## 【0041】

50

S 7 の処理：次のチャンク転送に設定する T C P コネクション数  $N$  に変更するため，T C P コネクションの増減制御を行う。

( 2 ) 内蔵アルゴリズム処理

ここでは S 6 の処理で行う T C P コネクション数  $N$  の決定例を、図 5 に基づき説明する。ここでは並列 T C P コネクションのスループット（チャンク / 転送時間  $T$ ）が、コネクション数  $N$  の多重度に関して上に凸の関数となる性質が示されている。

【 0 0 4 2 】

【表 2】

TCPコネクション数 $N$ の初期値 $N_0$	1
制御パラメータ $k$	2

Table 2: 図5に適用するパラメータ

10

【 0 0 4 3 】

表 2 は、図 5 の処理例に使用される「Table 2」を示している。この「Table 2」は、T C P コネクション数  $N$  の初期値  $N_0$  が「1」に設定され、制御パラメータ  $k$  が「2」に設定されている。

【 0 0 4 4 】

T C P コネクション数  $N$  を最適化する処理は、以下のステージ 1、2 で構成されている。このステージ 1 には図 5 に示す「A ~ E」の推移が該当し、ステージ 2 には図 5 に示す「E ~ H」までの推移が該当する。

20

( 1 ) ステージ 1

ステージ 1 は、T C P コネクション数  $N$  を初期値  $N_0$  から、制御パラメータ  $k$  を乗じて増加させ、最適値  $N_{opt}$  の含まれるコネクション数（図 5 の横軸）のおおまかな巾（以降、ブラケットと呼ぶ）を探索する。この状態を式（2）に示す。

【 0 0 4 5 】

【数 2】

初期値  $N_0$

$$N \leftarrow N_0 \quad \text{式(2)}$$

30

【 0 0 4 6 】

そのうえで、S 3 ~ S 6 の処理で記載した方法に基づき、チャンク転送のスループット  $G(N)$  を算出する。算出したスループット  $G(N)$  が前回の算出値との間で以下の式（3）の関係となるまで、

【 0 0 4 7 】

【数 3】

$$G(N) < G(N_{-1}) \quad N_{-1} : \text{前回の } N \text{ 値} \quad \text{式(3)}$$

40

【 0 0 4 8 】

以下の式（4）に基づき、

【 0 0 4 9 】



【数4】

$$N \leftarrow k \times N \quad (k > 1) \quad \text{式(4)}$$

【0050】

TCPコネクション数 $N$ を増加させながら、チャンクのスループット $G(N)$ を計測し、式(3)の評価(前回値との比較)を行う。式(3)が満たされた(TRUE)の場合(前回値よりスループットが低下した場合)、最初のブラケットには、式(5)に示すように、過去3回分の履歴にある3つの $N$ 値を組として定め、ステージ2へ移行する。

10

【0051】

【数5】

$$\begin{aligned} & \text{(前々回の } N \text{ 値, 前回の } N \text{ 値, 今回の } N \text{ 値)} \\ & = (N_{-2}, N_{-1}, N) = (l, m, r) \end{aligned} \quad \text{式(5)}$$

【0052】

図5中、「A」は初期値であり、「Table 2」に記載のとおり、「 $N = N_0 = 1$ 」とする。次いで「Table 2」に記載のパラメータ $k$ を用いて、TCPコネクションを2倍(図5中の「B」)にする。さらに次のチャンク転送では、 $N = 4$ (図5中の「C」)、 $N = 8$ (図5中の「D」)と増やした段階で式(3)の評価が、式(6)となり、最初(1st)のブラケット(2, 4, 8)を得る(図5中の「E」)。

20

【0053】

【数6】

$$\text{TRUE} \leftarrow \{G(8) < G(4)\} \quad \text{式(6)}$$

【0054】

(2)ステージ2

30

ステージ2では、ステージ1で算出した最適値 $N_{opt}$ (チャンクのスループットの唯一の極大点を示す $N$ の値)の含まれる最初のブラケットの中(図5の横軸)に黄金探索法を適用して絞り込み、最適値 $N_{opt}$ を抽出する。これは主に非特許文献8の手法が用いられる。この非特許文献8は、iSCSIではなく、GridFTPという別の技術に黄金探索法が用いられている。以下、ステージ2の絞り込みを具体的に説明する。

【0055】

まず、最初に式(7)と式(8)の黄金比(Golden Ratio)とを用いて、ステージ1で抽出した最初のブラケットから新たなTCPコネクション数となる「 $N$ 」を算出する。

【0056】

40

【数7】

$$\text{TRUE or FALSE} \leftarrow \{G(N) > G(m)\} \quad \text{式(7)}$$

【0057】

【数 8】

$$\text{Golden Ratio: } \nu = \frac{3 - \sqrt{5}}{2} \approx 0.382 \quad \text{式(8)}$$

【0058】

ブラケット  $(l, m, r)$  } に対して、新たな TCP コネクション数  $N$  を算出するため、式(9)を適用する(小数点以下は切り上げ)。

【0059】

【数 9】

10

$$N \leftarrow \begin{cases} l + (m - l)\nu & \text{if } m - l > r - m \\ m + (r - m)\nu & \text{上記以外} \end{cases} \quad \text{式(9)}$$

【0060】

次に、S3 ~ S6 の処理で記載した方法に基づき、チャンク転送のスループット  $G(N)$  を算出する。式(7)の関係が成り立つ(TRUE)場合、式(10)により、新たなブラケットを導入する。

【0061】

【数 10】

20

$$(l, m, r) \leftarrow \begin{cases} (m, N, r) & \text{if } m < N; \text{ TRUE} \\ (l, N, m) & \text{上記以外; TRUE} \end{cases} \quad \text{式(10)}$$

【0062】

式(7)の関係が成り立たない場合(FALSE)は、式(11)により、新たなブラケットを導出し、新たに導出したブラケットを用いて、式(9)から新たな TCP コネクション数  $N$  を算出する。

【0063】

【数 11】

30

$$(l, m, r) \leftarrow \begin{cases} (l, m, N) & \text{if } m < N; \text{ FALSE} \\ (N, m, r) & \text{上記以外; FALSE} \end{cases} \quad \text{式(11)}$$

【0064】

以上の過程を繰り返し、 $(l, m, r)$  が連続する整数となった場合、式(12)とする。

【0065】

【数 12】

40

$$N_{opt} \leftarrow m \quad \text{式(12)}$$

【0066】

実際に、ステージ1で抽出した1stブラケット(2, 4, 8)を式(9)に当てはめると式(13)となる。

【0067】

【数 1 3】

$$N = 6 \leftarrow 5.528 \leftarrow (2, 4, 8) \quad \text{式(13)}$$

【0 0 6 8】

N = 6 で式 ( 7 ) を図 5 の結果 ( 縦軸 ) で評価すると、式 ( 1 4 ) が成立する。

【0 0 6 9】

【数 1 4】

10

$$\text{TRUE} \leftarrow \{G(6) > G(4)\} \quad \text{式(14)}$$

【0 0 7 0】

よって、式 ( 1 1 ) で式 ( 1 5 ) の 2nd ブラケットを構成する ( 図 5 の「E」 ) 。

【数 1 5】

$$\text{2nd ブラケット: } (4, 6, 8) \leftarrow (2, 4, 8) \quad \text{式(15)}$$

20

【0 0 7 1】

再び式 ( 9 ) で「N」を算出すると、式 ( 1 6 ) となる。

【0 0 7 2】

【数 1 6】

$$N = 7 \leftarrow 6.764 \leftarrow (4, 6, 8) \quad \text{式(16)}$$

【0 0 7 3】

N = 7 で式 ( 7 ) を評価 ( 図 5 参照 ) し、式 ( 1 7 ) の結果から式 ( 1 0 ) で式 ( 1 8 ) の 3rd ブラケットを構成する ( 図 5 の「F」 ~ 「G」 ) 。

【0 0 7 4】

【数 1 7】

$$\text{FALSE} \leftarrow \{G(7) < G(6)\} \quad \text{式(17)}$$

【0 0 7 5】

【数 1 8】

40

$$\text{3rd ブラケット: } (4, 6, 7) \leftarrow (4, 6, 8) \quad \text{式(18)}$$

【0 0 7 6】

再び式 ( 9 ) でNを算出すると式 ( 1 9 ) となる。

【0 0 7 7】

【数 19】

$$N = 5 \leftarrow 4.382 \leftarrow (4, 6, 7) \quad \text{式(19)}$$

【0078】

N = 5 で式 (7) を評価 (図 5 参照) し、式 (20) の結果から式 (11) で式 (21) の 4rd ブラケットを構成する。

【0079】

【数 20】

$$\text{FALSE} \leftarrow \{G(5) < G(6)\} \quad \text{式(20)}$$

【0080】

【数 21】

$$4\text{th ブラケット}: (5, 6, 7) \leftarrow (4, 6, 7) \quad \text{式(21)}$$

【0081】

ここでブラケット内の整数がすべて連続となったため、ステージ 2 を終了する。スループットを最大にする TCP コネクション数の最適値  $N_{opt}$  として、式 (22) に示す最後のブラケットの中央の数字である「6」が抽出できる。

【0082】

【数 22】

$$N_{opt} = 6 \leftarrow (5, 6, 7) \quad \text{式(22)}$$

【0083】

整理すると、図 5 の例では、TCP コネクション数の初期値  $N_0$  から最適値  $N_{opt}$  に移行するまで、式 (23) に示す遷移を行う。

【0084】

【数 23】

$$N_0 = 1 \rightarrow 2 \rightarrow 4 \rightarrow 8 \rightarrow 6 \rightarrow 7 \rightarrow 5 \rightarrow N_{opt} = 6 \quad \text{式(23)}$$

【0085】

このように並列 TCP コネクションのスループット (チャンク / 転送時間  $T$ ) が、コネクション  $N$  の多重度に関して上に凸の関数になるという性質 (図 5 の「A」～「H」) を利用して、TCP コネクションの多重度をスループットに対して最適化させている (スループットが最大となる TCP コネクション数を探索し、当該数の TCP コネクション数にまで増減制御している。 )。

【0086】

すなわち、iSCSI イニシエータ 1 に組み込まれた TCP コネクション数  $N$  を自動調整する iSCSI - APT 3 の制御機能によって、予め使用するネットワークの遅延特性や帯域等の情報が無い場合でも、iSCSI - APT 3 の処理機能を定期的に動作させることで、動作時のネットワークに最適化した TCP コネクション数  $N_{opt}$  に自動追従さ

10

20

30

40

50

せている。

【0087】

これにより i S C S I 接続されるファイルサーバとストレージとの間のスループットの高速化を実現させている。この結果、ファイルサーバなどの運用設置を大幅に簡略でき、帯域変動にも追従可能となることから、与えられた回線で常に効率的なデータ転送が実現可能となる。

【0088】

また、ファイルサーバの実装的には i S C S I イニシエータの修正で済むので、T C P コネクションの多重化オプションをサポートする市販のターゲット装置の利用でき、この意味で経済的なシステム構築・更新が可能である。

10

【0089】

なお、本発明は、上記実施形態に限定されるものではなく、各請求項に記載した範囲において各種の変形を行うことが可能である。例えばデータ転送（送信）だけではなく、データの読み込み（受信）に関しても同様の効果が期待できる。

【0090】

また、前記 i S C S I - A T P 3 の制御機能を内蔵した i S C S I イニシエータ 1 として、コンピュータを機能させるプログラムとして構築することもできる。

【0091】

この場合には、コンピュータの処理部（例えば C P U など）がプログラムコードを読み出し、通信デバイスを通じてステージ 1、2 の処理が実行される。

20

【0092】

したがって、既存の i S C S I システムにおいてファイルサーバに前記プログラムをインストールするだけで、本実施形態に係る i S C S I イニシエータ機能を有するファイルサーバが構築でき、この点でシステム構築・更新の経済性が一層向上する。

【0093】

このプログラムコードは、例えば C D - R O M , D V D - R O M , C D - R , C D - R W , D V D - R , D V D - R W , M O , H D D などの記録媒体に格納される。また、前記プログラムを、インターネットサイトからダウンロードしてコンピュータに提供してもよい。

【符号の説明】

30

【0094】

- 1 ... i S C S I イニシエータ
- 2 ... i S C S I ターゲット ( i S C S I デバイス )
- 3 ... i S C S I - S P T
- N ... T C P コネクション数
- G ( N ) ... スループット値

【図面の簡単な説明】

【0095】

【図1】 T C P コネクション数を複数利用する i S C S I 接続の概略図。

【図2】 本発明の実施形態に係る T C P コネクション数制御機能 ( i S C S I - A P T ) を i S C S I イニシエータの i S C S I 接続の概略図。

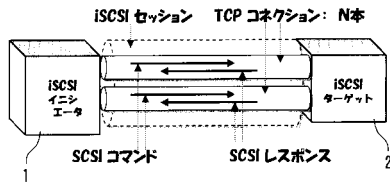
40

【図3】 1つの T C P コネクションでの「 i S C S I P D U 」転送のシーケンス例を示し、チャンク転送時間算出のための時刻定義の概略図。

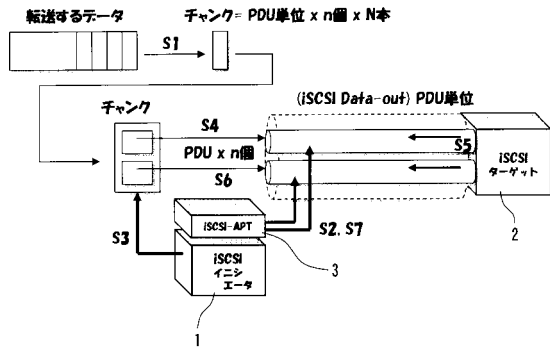
【図4】 複数の T C P コネクションを用いたチャンク転送での転送時間の定義を示す概略図。

【図5】 チャンク転送のスループットの観測から最適な T C P コネクション数を導出する処理過程例を示す概略図。

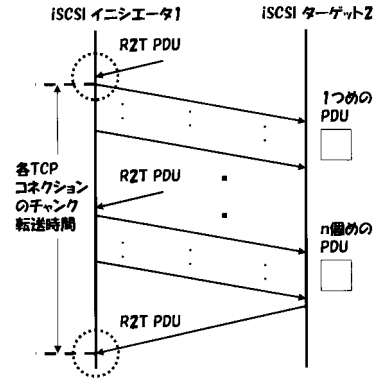
【 図 1 】



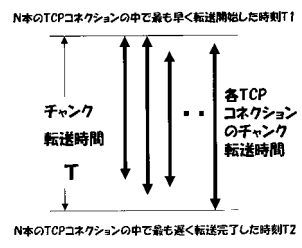
【 図 2 】



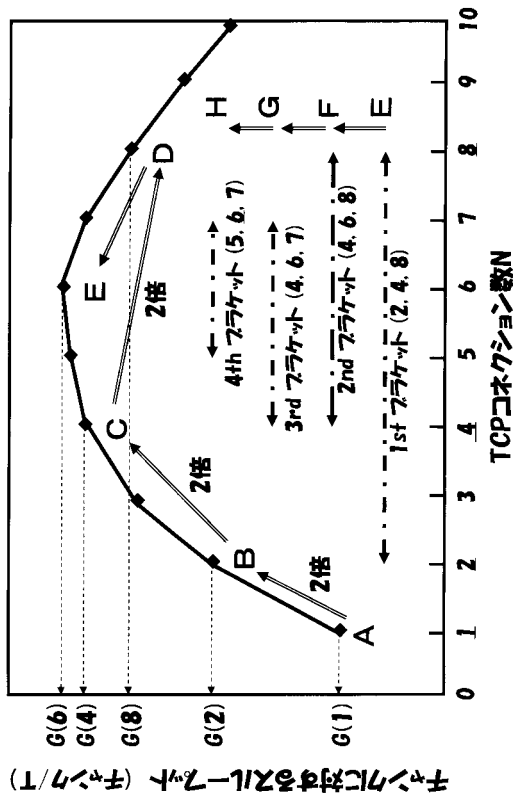
【 図 3 】



【 図 4 】



【 図 5 】



## フロントページの続き

- (72)発明者 野本 義弘  
東京都千代田区大手町二丁目3番1号 日本電信電話株式会社内
- (72)発明者 大崎 博之  
大阪府吹田市山田丘1番1号 国立大学法人大阪大学内
- (72)発明者 今瀬 真  
大阪府吹田市山田丘1番1号 国立大学法人大阪大学内
- (72)発明者 井上 史斗  
大阪府吹田市山田丘1番1号 国立大学法人大阪大学内

審査官 矢頭 尚之

- (56)参考文献 特開2006-270303(JP,A)  
武田 裕子 ほか, iSCSI over VPN環境における複数経路アクセス適応制御手法の提案と評価, 電子情報通信学会 第18回データ工学ワークショップ論文集 DEWS2007 E7-3, 日本, 電子情報通信学会データ工学研究専門委員会, 2007年 6月 1日, [online] DEWS2007 HIROSHIMA  
伊藤 建志 ほか, GridFTP - APT: データ転送プロトコルGridFTPの並列TCPコネクション数調整機構, 電子情報通信学会技術研究報告 IN2005-113, 日本, 社団法人電子情報通信学会, 2005年12月 8日, 第105巻 第472号, p.19-24  
武田 裕子 ほか, VPN上のiSCSI環境における複数経路アクセス適応制御手法の提案と評価, 日本データベース学会Letters, 日本, 日本データベース学会, 2007年 6月, 第6巻 No.1, p.129-132  
藤原 啓成 ほか, 広域IP網を介したiSCSI通信におけるプロトコルチューニングの一検討, 第68回(平成18年)全国大会 講演論文集(1) アーキテクチャ ソフトウェア科学・工学, 日本, 情報処理学会, 2006年 3月 7日, 5A-5, p.1-55

## (58)調査した分野(Int.Cl., DB名)

H04L 12/56  
G06F 13/00  
H04L 13/00