

B-027

## 高遅延環境下における複数台の iSCSI Initiator 使用時の TCP 輻輳ウィンドウの考察

A Study of TCP Congestion Window using Multiple iSCSI Initiator in Long-Latency Network

豊田 真智子<sup>†</sup>  
Machiko Toyoda

山口 実靖<sup>‡</sup>  
Saneyasu Yamaguchi

小口 正人<sup>†</sup>  
Masato Oguchi

### 1. まえがき

大容量のデータを蓄積するアプリケーションの登場とインターネット技術の発展により、サーバなどで管理されるデータ容量の増大が無視できない問題となっている。そのため、サーバとストレージ間をストレージ専用的高速ネットワークによって接続する SAN (Storage Area Network) が登場し、高い実績をあげている。ファイバチャネルを用いて構築する FC-SAN に変わり、Ethernet と TCP/IP を用いて構築する IP-SAN が次世代 SAN として期待されている。この IP-SAN のプロトコルにおいて、最も注目されているのが iSCSI プロトコルである [1][2]。

iSCSI では、サーバ (Initiator) とストレージ (Target) 間のデータ通信を SCSI コマンドによって実現し、SCSI over iSCSI over TCP/IP over Ethernet という複雑なプロトコルスタックを構成する。そのため、各レイヤ間のデータ受け渡しにおけるメモリコピーなどの影響で、大幅に性能が劣化する [3]。また、ストレージアクセス時のスループットは、TCP パラメータの輻輳ウィンドウと密接に関連し、輻輳ウィンドウがのこぎり型の増加減少の変化を繰り返す場合には輻輳ウィンドウの変化に従い、スループットも不安定になることが確認されている [4]。iSCSI は遠隔ストレージへのバックアップやストレージのアウトソーシングなどの利用が想定されるため、iSCSI 使用の際の性能低下を最小限に抑えることは非常に重要となり、輻輳ウィンドウの振舞を考察することは、性能向上の糸口になると考えられる。

そこで本稿では、複数の iSCSI Initiator からストレージ (Target) にアクセスする環境を想定し、その際の輻輳ウィンドウとスループットの振舞について考察を行う。

### 2. 複数台の Initiator を用いた性能測定実験

複数台のサーバからストレージにアクセスした場合の iSCSI ネットワークの性能を測定するため、最大 4 台の Initiator マシンを用いた実験を行う。Initiator から Target の raw デバイスに対してシーケンシャルリードアクセスを行い、その時の輻輳ウィンドウ、スループットを測定する。各マシンごとの振舞を確認するため、スループットは Target 側ではなく、各 Initiator において測定を行った。

また、本実験で用いた NIC (Network Interface Card) は、通信時に TCP 層から受け取ったデータを保持するためのバッファサイズを、ディスクリプタ値を設定することにより変更することが可能となっている。そこで、

表 1: 使用計算機: Initiator

CPU	Intel Pentium III 800MHz
Main Memory	640MB
OS	Linux2.4.18-3
NIC	Intel PRO/1000MT Server Adapter

表 2: 使用計算機: Target

CPU	Intel Xeon 2.4GHz
Main Memory	512MB
OS	Linux2.4.18-3
NIC	Intel PRO/1000XT Server Adapter

表 3: 使用計算機: Dummynet

CPU	Intel Xeon 2.4GHz
Main Memory	512MB
OS	FreeBSD4.9 - RELEASE
NIC	Intel PRO/1000MT Server Adapter

複数台の Initiator から Target にシーケンシャルリードアクセスを行った場合の性能をより詳細に調べるために、データ送信側である Target の NIC ディスクリプタを変更することで送信バッファサイズを変更して実験を行った。NIC ディスクリプタは 80 から 4096 までの間で変更することができ、デフォルトは 256 に設定されている。そこで、デフォルト値である 256 と、デフォルトよりバッファサイズを大きくした場合として 4096 の 2 パターンを設定し、各バッファサイズにおける測定を行った。

#### 2.1 実験環境

本実験は以下の環境で行った。Initiator と Target 間は 1000Base-T スイッチングハブを介して Gigabit Ethernet で接続し、遠隔ストレージアクセスを想定した実験を行うため、Target とハブの接続途中に人工的な遅延装置である FreeBSD Dummynet[5] を挿入し TCP/IP 接続を確立した。実験環境の概観を図 1 に示す。Initiator, Target, Dummynet はすべて PC 上に構築した。iSCSI を利用したストレージアクセスにおけるネットワークの性能を調べるため、Target はメモリモードで動作させ、ディスクアクセスを伴わないように設定した。実験で使用した計算機的环境を表 1, 2, 3 に示す。

また、本実験で用いた iSCSI 実装において、Target にはニューハンプシャー大学 InterOperability Lab が提供する UNH IOL reference implementation ver.3 on iSCSI Draft 18[6] を用いた。この UNH 実装では、大きなブロックサイズで read コマンドを発行しても、SCSI 層におい

<sup>†</sup>お茶の水女子大学  
Ochanomizu University

<sup>‡</sup>東京大学 生産技術研究所  
Institute of Industrial Science, The University of Tokyo

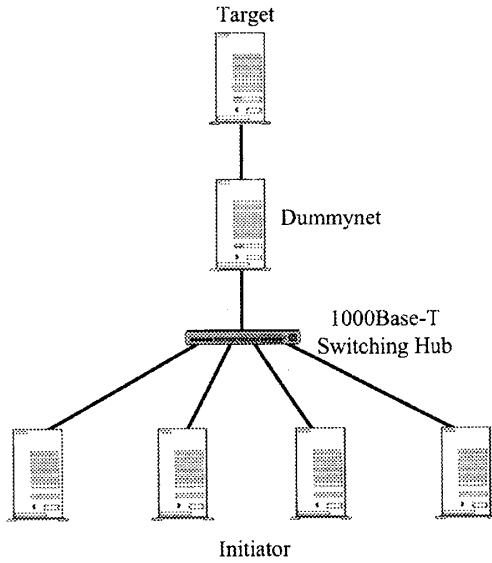


図 1: 実験環境概観図

て要求したブロックサイズより小さなブロックサイズに分割されてしまい、これにより iSCSI ストレージアクセスの性能が大きく低下してしまうことが確認されている [3]。本実験ではこの実装による性能測定への影響を避けるため、Initiator に UNH 実装を用いず、UNH 実装の Initiator と同等の機能を持ち、かつ大きなブロックサイズのデータ転送も行える自作 Initiator を用いて実験を行った。この自作 Initiator は通常のユーザ空間のアプリケーションとして動作し、iSCSI Target と TCP/IP コネクションを確立して iSCSI プロトコルで通信を行うものである。

## 2.2 実験結果

片道遅延時間 8ms (ラウンドトリップタイム: 16ms) に設定し、Initiator の台数を 1 台から 4 台の間で変更してストレージアクセスを行った実験結果として、各ディスクリプタ値における Target で観測した輻輳ウィンドウの時間変化のグラフを図 2, 3 に示す。

ディスクリプタ値が 256 の場合は、下位層のデバイスドライババッファが溢れてしまうため (CWR エラー)、TCP 実装により輻輳が起こったとみなされ、輻輳ウィンドウが急激に減少し、再度増加するというのこぎり型の変化を繰り返している。また、Initiator 数の増加に伴い、輻輳ウィンドウの成長は小さくなり、エラーが増加することが確認される。

一方、ディスクリプタ値が 4096 の場合は一度も輻輳が発生せず、輻輳ウィンドウの上限値は大きい。Linux TCP においては、通信中に一度設定された輻輳ウィンドウは、その値が使い切られない限りは変化しないという特徴を持ち、この時一定値をとる実装となっている。Initiator 台数が 4 台の場合、輻輳ウィンドウが大きく低下しているが、これは輻輳発生によるエラーではなく 4 台目の Initiator に割り当てられている輻輳ウィンドウの値である。これについては次節で詳しく述べる。

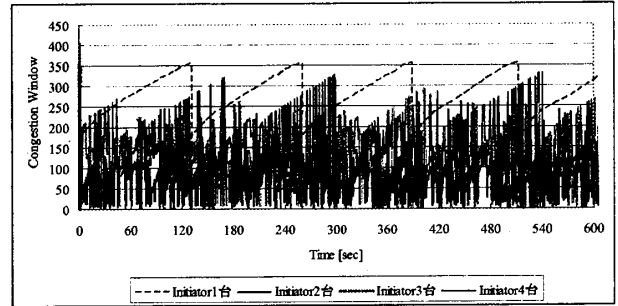


図 2: NIC ディスクリプタ 256 (デフォルト) の場合の輻輳ウィンドウの時間変化

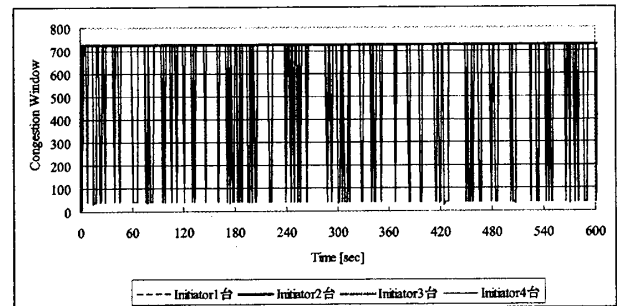


図 3: NIC ディスクリプタ 4096 の場合の輻輳ウィンドウの時間変化

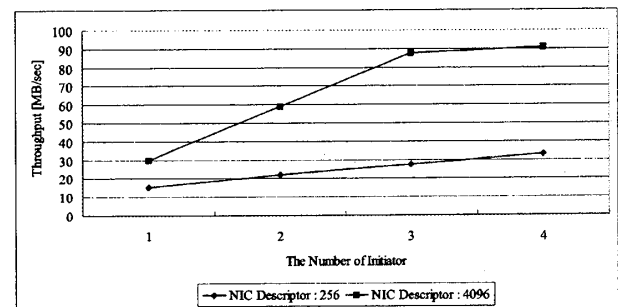


図 4: スループット測定結果

次に各 Initiator のスループットを合計したスループット測定結果を図 4 に、Initiator 台数を変化させた場合の各 Initiator の平均スループットを表 4, 5 に示す。合計スループットは、各環境における Target のスループットであるとみなすことができる。ディスクリプタを大きく設定することにより送信バッファ容量が増加し、より多くのデータを一度に送り出すことができるため、各台数におけるスループットはディスクリプタ値が大きい 4096 の場合に大きく向上する。また、ディスクリプタ値 256 の場合は Initiator の台数増加に伴いスループットも増加しているが、ディスクリプタ値 4096 の場合は 3 台以降にあまり変化が無く、ほぼ飽和状態となっていることがわかる。

表 4: NIC ディスクリプタ 256 の場合の各 Initiator におけるスループット

Initiator 台数	スループット [MB/sec]			
	Initiator1	Initiator2	Initiator3	Initiator4
1	15.24			
2	10.55	11.47		
3	9.6	9.01	9.11	
4	7.88	7.69	7.23	10.38

表 5: NIC ディスクリプタ 4096 の場合の各 Initiator におけるスループット

Initiator 台数	スループット [MB/sec]			
	Initiator1	Initiator2	Initiator3	Initiator4
1	29.88			
2	29.54	29.69		
3	29.25	29.23	29.21	
4	29.13	29.15	29.13	3.2

### 3. 考察

NIC ディスクリプタを 256 に設定した場合は Initiator 台数を増加させるに伴い、輻輳ウィンドウの上限が低下している様子が図 2 から確認される。ここで、Initiator 台数が 1 台から 2 台に増加した場合を考える。輻輳ウィンドウは一度に送信できるパケット数を意味しているため、輻輳ウィンドウが減少した場合にスループットも減少するはずである。しかし、本実験の結果から、輻輳ウィンドウが低下しているのに対し、スループットは大きく向上していることが確認された (図 4)。このことから、Target では NIC ディスクリプタにより決められる送信バッファを、同時に存在するコネクション数によって分割していることがわかる。輻輳ウィンドウは各コネクションごとに存在するため、輻輳ウィンドウの上限も Initiator 数が増加するごとに小さくなるが、その輻輳ウィンドウの値が各コネクションに適用されて通信が行われるため、Initiator 台数が増加すると合計のスループットは向上すると考えられる。

NIC ディスクリプタを 4096 に設定した場合、最も大きな変化が見られたのは、表 5 における 4 台目の Initiator のスループットの値である。他のディスクリプタ値における実験の場合も含め、どの実験においても各 Initiator のスループットは同等であったが、Initiator 台数を 4 台にした実験においては、そのうち 1 台のスループットが大きく低下した。また、輻輳ウィンドウの値も Initiator 台数が 1 台から 3 台の場合にはほぼ一定値をとり、大きな変化はみられないが、Initiator 台数を 4 台にした場合、その値が大きく減少している様子が図 3 から確認される。この輻輳ウィンドウは、輻輳が検出されたことによる低下ではない。このことから、低下した時の小さな輻輳ウィンドウの値は、4 台目の Initiator に割り当てられている輻輳ウィンドウであると考えられる。輻輳ウィンドウが小さいために Target から送信されるデータ量

が少なくなり、スループットの向上は見られなかったと言える。

高遅延環境においては応答時間が長くなるため、遅延が少ない場合に比べて性能が大幅に低下する。しかし、送信バッファを大きくすることにより輻輳ウィンドウが大きな値まで成長し、その性能は大幅に向上する。また、同時に Target にアクセスする Initiator 台数を増やすことも、性能向上の大きな要因の一つであることが確認された。

### 4. まとめと今後の課題

iSCSI ストレージアクセスにおいて、高遅延環境下で複数の Initiator からシーケンシャルリードアクセスを行い、輻輳ウィンドウとスループットの値からその振舞いを考察した。輻輳ウィンドウはコネクションごとに設定され、送信バッファもコネクション数に応じて分割されることが確認された。そのため、コネクション数の増加に伴い輻輳ウィンドウの上限は減少する。また高遅延環境においては、送信バッファ容量の増大と、アクセスする Initiator 台数の増加によりスループットが大きく向上することがわかった。

今後は、今回の実験から得られた結果を踏まえ、複数台の Initiator を用いたストレージアクセスにおける性能向上手法を検討していきたい。

### 謝辞

本研究は、一部、文部科学省科学研究費特定領域研究課題番号 13224014 によるものである。

### 参考文献

- [1] iSCSI Specification,  
<http://www.ietf.org/rfc/rfc3720.txt?number=3270/>
- [2] SCSI Specification,  
<http://www.danbbs.dk/~dino/SCSI/>
- [3] 山口実靖, 小口正人, 喜連川優: “高遅延広帯域ネットワーク環境下における iSCSI プロトコルを用いたシーケンシャルストレージアクセスの性能評価ならびにその性能向上手法に関する考察”, 電子情報通信学会論文誌 Vol.J87-D-I, No.2, pp.216-231, February 2004.
- [4] 豊田真智子, 山口実靖, 小口正人: “iSCSI ストレージアクセス時における TCP 輻輳ウィンドウとシステム性能の関連性評価”, FIT2004 第 3 回情報科学技術フォーラム, B-004, pp.107-109, September 2004.
- [5] L.Rizzo: “dummysnet,”  
[http://info.iet.unipi.it/~luigi/ip\\_dummysnet/](http://info.iet.unipi.it/~luigi/ip_dummysnet/)
- [6] InterOperability Lab  
Univ. of New Hampshire,  
<http://www.iol.unh.edu/consortiums/iscsi/>