

# iSCSI を用いた IP-SAN 統合 PC クラスタの性能に関する一考察

神坂 紀久子<sup>†</sup>      山口 実靖<sup>‡</sup>      小口 正人<sup>†</sup>      喜連川 優<sup>§</sup>

<sup>†</sup>お茶の水女子大学      <sup>‡</sup>工学院大学      <sup>§</sup>東京大学生産技術研究所

## 1 はじめに

近年、コモディティなハードウェアの性能向上と低価格化により、大規模科学技術計算やデータベース処理等を PC クラスタにおいて実行することが一般的になった。従来より、PC クラスタでは、クラスタノード-ストレージ間のネットワークに Fibre Channel(FC) や Infiniband などの高速な専用回線が使用されてきた。しかし、TCP/IP と Ethernet を使用する IP-SAN(IP-Storage Area Network)[1] のプロトコルである iSCSI が登場したことにより、コモディティなネットワークだけを使用した PC クラスタの構築が可能になった。

そこで本稿では、個々に構築していたノード間のフロントエンドネットワークとノード-ストレージ間のバックエンドネットワークを、コモディティな一つのネットワークに統合した IP-SAN 統合 PC クラスタを iSCSI を用いて実現する。また、双方のネットワークを統合した環境が、それらを個々に構築した場合と比較して性能にどの程度影響を与えるかを評価した。

## 2 IP-SAN 統合 PC クラスタ

### 2.1 SAN を用いた PC クラスタ

PC クラスタにおいて、ノード-ストレージ間のバックエンドのネットワークに SAN(Storage Area Network) を用いて構築することが多くなった。SAN はサーバ機とストレージデバイスを接続する高速な専用のネットワークであり、分散されたストレージの統合と集中管理を可能にする。IP-SAN は、IP ネットワークで構築する次世代の SAN である。図 1 に示すように、FC を用いて構築する従来の SAN に代わり、バックエンドのネットワークを IP-SAN で構築することにより、安価なコストで PC クラスタの導入、運用ができる。IP-SAN の代表的なプロトコルとしては iSCSI があり、SCSI コマンドを TCP/IP パケットの中にカプセル化することにより、ブロックレベルのデータ転送を行う。

通常、SAN を使用した PC クラスタでは、ノード間のフロントエンドとノード-ストレージ間バックエンドのネットワークを個々に構築する。しかし、その場合は異なるネットワークの構築が必要になるため、運用管理の面においても容易ではない。

### 2.2 iSCSI を使用した IP-SAN 統合 PC クラスタ

PC クラスタに iSCSI を使用することで、個々に構築していたフロントエンドとバックエンドを一つのコモディティなネットワークに統合することが可能になる。そこで我々は、図 2 に示すように、ネットワーク構築コストの削減と運用管理の効率化を目的として、バックエンドのネットワークをフロントエンドに統合した

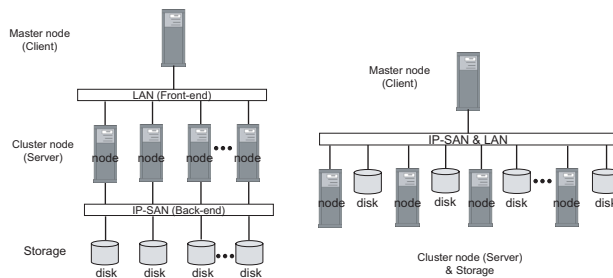


図 1: IP-SAN を用いた PC クラスタ

表 1: 実験環境：使用計算機

OS	Linux 2.6.10
CPU	Intel Pentium 4 CPU 1500MHz
Main Memory	384MB
HDD	36GB SCSI HDD
NIC	Intel(R) PRO/1000 MT

IP-SAN 統合 PC クラスタを実現した [2]。しかし、その場合、フロントエンドで行われるノード間通信とバックエンドで行われるストレージアクセスのデータが、同一の IP ネットワーク経由で混在して転送され、ネットワークへの負荷が懸念される。例えば、ストレージアクセスのバルクデータにより並列計算のためのノード間通信が多大な影響を受け、並列分散処理を実行する際の全体の性能が劣化する可能性がある。

## 3 IP-SAN 統合 PC クラスタと非統合 PC クラスタの比較実験

本稿では、IP-SAN 統合 PC クラスタが、ネットワークを個々に構築した場合と比較して並列分散処理性能にどの程度影響を及ぼすかを評価した。

### 3.1 実験環境

本実験では、ローカルストレージを用いた PC クラスタ、バックエンドに IP-SAN を用いた非統合 PC クラスタ(図 1)、IP-SAN 統合 PC クラスタ(図 2)の 3 つの環境において性能を比較する。非統合 PC クラスタにおけるフロントエンドとバックエンドのネットワーク、そして IP-SAN 統合 PC クラスタにおけるネットワークは Gigabit Ethernet を用いて構築した。実験に用いた PC のシステム環境を表 1 に示す。また IP-SAN には iSCSI を使用しているが、本稿の実験環境ではノードとストレージは 1 対 1 接続になっており、各ノードは特定のストレージデバイスに接続する構成となっている。

### 3.2 マクロベンチマーク

まずマクロベンチマークとして、一般的に使用されている並列ベンチマークである NAS Parallel Benchmark (NPB) を使用して性能を測定した。NPB は並列演算性能を測定するベンチマークであるが、本実験で使用し

Study of Performance in PC cluster system integrated with IP-SAN using iSCSI  
 Kikuko Kamisaka<sup>†</sup>, Saneyasu Yamaguchi<sup>‡</sup>,  
 Masato Oguchi<sup>†</sup> and Masaru Kitsuregawa<sup>§</sup>  
<sup>†</sup>Ochanomizu University, <sup>‡</sup>Kogakuin University, <sup>§</sup>Institute of Industrial Science, The University of Tokyo

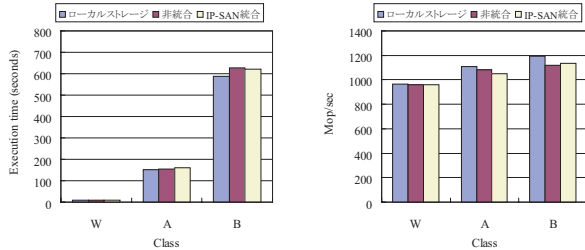


図 3: NPB I/Oの実行時間 図 4: NPB I/OのMops値

たNPB 2.4 I/Oは、対象問題BT (Block Tri-diagonal) に対し、大量のI/O処理を行うアプリケーション実行時の性能を測定できる。また本実験において実行したNPBのClassは、W、A、Bである。NPB 2.4 I/Oの実行オプションとして並列I/Oを行うepioを用いており、Class Wにおいては全ノードで22MB、Aにおいては420MB、Bにおいては1698MBのデータをストレージに書き出す。並列計算を行うノード数は4とした。

NPBの実行時間とMops (Million Operations Per Second) 値を図3、4に示す。Mops値は1秒間あたりの100万演算数である。問題サイズが小さいClass Wの場合には、実行時間とMops値ともに、3つのPCクラスタ環境でほぼ同じ値が得られた。これらのClassにおいては入出力が行われるデータ量が少なく、並列処理演算が支配的なためである。一方、Class Bの問題サイズでは、IP-SANを用いた場合はローカルストレージの場合と比較してやや性能差が大きくなる。しかし、IP-SAN統合クラスタと非統合クラスタでは、並列演算処理性能にほぼ差がないといえる。

### 3.3 マイクロベンチマーク

次に、マイクロベンチマークとして、ノード間通信とストレージアクセスがほぼ同時に行われるような並列プログラムを作成し、統合による影響をさらに詳細に評価した。

図5は、作成した並列ベンチマークの擬似コードである。本実験ではノード数を2とし、この並列プログラムを実行した際の処理時間を測定した。このプログラムでは、各ノードからそのノードに接続されているストレージデバイスに対してwriteを実行し、その後、送信ノード(ノード1)ではMPI\_Send()を、受信ノード(ノード2)ではMPI\_Recv()を実行する。fwrite関数の命令を先に発行しても、システムバッファにデータが書き出された時点ですぐに次のsend-recv関数に移行するため、ほぼ同時にwriteとsend-recvが実行されているといえる。図5における“SEND\_RECV\_ITERATION”は、send-recvの繰り返し回数、つまりメッセージ送信回数である。ノード間通信で実行されるsend-recvとストレージアクセスで実行されるwriteは性能差があると考えられるため、統合の性能への影響が大きくなるような状況を想定し、このメッセージ送信回数を増加させて性能を測定した。

図6は、I/Oサイズを1MBとし、1回に送信されるメッセージサイズを1MBにしたときの実行時間である。最終的に各ノードに接続されたストレージに作成されるファイルサイズは、各ノードそれぞれ1GBである。同図より、ローカルストレージを使用した場合には、メッセージ送信回数が2以上になると比例して実行時間が長くなる。これは、メッセージ送信回数が2

```

node1
#define SEND_RECV_ITERATION number
for(i=0; i<ITERATION; i++){
  fwrite(buf);
  for(j=0; j<SEND_RECV_ITERATION; j++){
    MPI_Send(inbuf, dst);
  }
}

node2
#define SEND_RECV_ITERATION number
for(i=0; i<ITERATION; i++){
  fwrite(buf);
  for(j=0; j<SEND_RECV_ITERATION; j++){
    MPI_Recv(outbuf, src);
  }
}

```

図 5: 並列プログラムの擬似コード

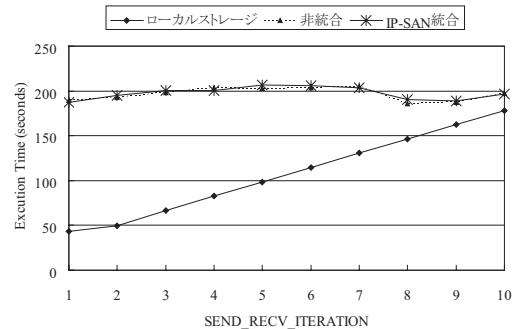


図 6: メッセージ送信回数変動による実行時間

より小さい場合にはI/Oが支配的になり、メッセージ送信回数が2以上の場合には送信するメッセージ量に依存するためである。一方、IP-SAN統合クラスタの場合には、非統合クラスタの場合と比較してほぼ同じ実行時間となった。これは、iSCSIを介したストレージアクセスの通信性能がかなり低いことが原因の一つと考えられる。つまり、iSCSIのI/Oにより送信されるパケットの頻度が低く、send-recvを行っている通信に殆ど影響を与えなかったためである。

IP-SAN統合PCクラスタは、ストレージアクセスのバルクデータとノード間通信のデータが同一ネットワーク上に混在することによる性能劣化が懸念された。しかし本稿の実験の結果から、バックエンドのIP-SANをフロントエンドのネットワークに統合した性能への影響はほとんどなく、有効であるといえる。

## 4 まとめと今後の課題

PCクラスタの構築および運用管理コストを削減するため、バックエンドのネットワークをフロントエンドに統合したIP-SAN統合PCクラスタを実現した。本稿では、フロントエンドとバックエンドのネットワークを個々に構築した場合と比較して、性能にどの程度影響を与えるかということの評価をした。その結果、バックエンドにIP-SANを持つ非統合PCクラスタと比較して、双方のネットワークを統合しても、ほぼ同等の並列分散処理性能を達成でき、有効であることがわかった。

今後の課題として、IP-SAN統合PCクラスタの性能を様々な状況においてさらに詳細に評価する。また、IP-SAN統合PCクラスタに高速なネットワークを適用することを考える。

## 参考文献

- [1] “Storage Networking Industry Association”. <http://www.snia.org/>.
- [2] 神坂紀久子, 山口実靖, 小口正人, 喜連川優: “IP-SAN統合PCクラスタを用いたトラフィック特性とI/O性能に関する考察”, 夏のデータベースワークショップ2006 (DBWS2006), 電子情報通信学会技術研究報告, DE2006-50~91.