

A-9 複数 NIC を用いた通信負荷分散によるバンド幅性能向上と並列処理への効果

福永隆文\*      上瀧剛\*\*      村上悟\*\*  
 菊池泰紀\*\*      芦原評\*\*      梅野英典\*\*

1. はじめに

近年、PC/AT 互換機等 (以後 PC) の高性能化及び Fast Ethernet や Gigabit Ethernet 等による通信性能の向上にとともに PC をクラスタ結合した並列システムが広まりつつある。しかしながら、通信の多いアプリケーションでは通信オーバーヘッドが大きくなり性能が向上しないケースも見受けられる。

我々は TCP/IP アーキテクチャを用いた Fast Ethernet 上に複数 NIC を用いて通信バンド幅を向上させる仕組み (以後 NIC 分散) を組み込み、バンド幅の向上が並列処理に与える効果を評価する。当通信負荷分散の仕組みは、Linux の Channel Bonding ドライバ (以後 Bond ドライバ) へ変更を加えることにより実装した。

複数 NIC を用いて通信バンド幅を向上させる仕組みとしては PM/Ethernet<sup>1)</sup> の Network Trunking があるが、我々が提案する方式は TCP 上で実現している点が異なる。

2. NIC 分散処理

2.1 NIC 分散処理の概要

複数 NIC を用いた NIC 分散方式のアーキテクチャを図 1 に示す。eth0~eth3 は NIC を示す。送信処理では、今回修正を加えた Bond ドライバが各 NIC ドライバをラウンドロビン方式で呼び出すことにより、送信データを複数の NIC へ交互に割り当てて送信する。1枚の NIC 使用時には、NIC 送信中は次の送信メッセージが存在しても送信できないため待ちが発生するが、複数 NIC を使うことにより送信中以外の NIC を利用して送信できる。各 NIC から送信される Ethernet フレームの宛先アドレスは、送信側 NIC 毎に対応する受信側 NIC のアドレスに書き換えるので、Switch により対応する NIC へ送信される。

2.2 NIC 分散処理の実装

我々は Linux2.4.2 上に NIC 分散を実装した。カーネル及び NIC ドライバは変更せず、Bond ドライバの修正を行った。このドライバは Linux に実装されているダイナミックロードが可能なドライバであり、送信側の負荷分散のみを行う。宛先 Ethernet アドレスは同一であり、受信側では特定の NIC が受信し、負荷分散は行われぬ。この Bond ドライバの送信関数内に宛先 Ethernet アドレスをラウンドロビン方式で交互に変更する仕組みを組み込んだ。送信元の NIC と送信先の NIC は送信先アドレステーブルで対応づけられている。送信先の NIC

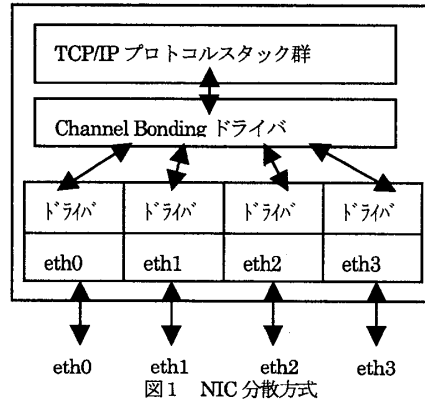


図1 NIC分散方式

表1 送信先アドレステーブルの項目

項目	説明
送信先識別子	送信先識別用
NIC 番号	送信元のNICに付けた番号で、eth0が1, eth1が2, 以下連番となる
変更アドレス	送信元NICと通信を行う送信先NICのEthernetアドレス

アドレスは事前に送信先アドレステーブルに登録する必要があります。テーブルの項目を表1に示す。Bond ドライバの送信関数内で各NICのドライバをラウンドロビンで呼び出しているが、その呼び出しの前に当該テーブルを検索し、宛先アドレスを送信元NICに対応する送信先NICのアドレスに書き換える。

TCP/IP 処理は受信側が特にオーバーヘッドが高いが、今回のコード追加は送信側のみコードを追加した。これにより、オーバーヘッドの増加による通信への影響はある程度押さえられる。

3. 評価

バンド幅の測定結果、Channel Bonding 方式との比較結果及び NAS 並列ベンチマークの性能評価結果を示す。測定環境は Celeron 533MHz, 256MB SDRAM メモリ (環境1) と Pentium IV 1.6GHz, 256MB SDRAM メモリ (環境2) を用いた。1MB は 1,024 × 1,024 バイトで計算した。

3.1 TCP/IP バンド幅性能

TCP/IP 処理はコネクション型でありオーバーヘッドが高い。今回 NIC 分散のためコードを追加するのでオーバーヘッドはさらに高くなる。そこで、NIC 分散によるバンド幅性能向上と TCP/IP 処理オーバーヘッドが NIC 分散に与える影響を示す。

図2は環境1を用いた場合のNIC分散とバンド幅の関係であり、図3は環境2を用いた場合の結果である。図2よりCPUの処理能力が低い場合、NIC2枚までは分散による効果が見られるが、それ以上NIC枚数を増やしても性能向上が見られな

\* 熊本県立技術短期大学校  
 \*\* 熊本大学 工学部

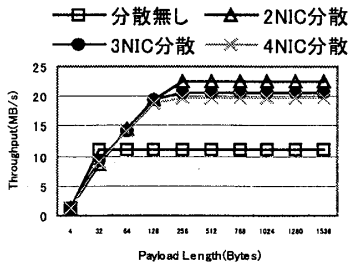


図2 バンド幅性能(環境1)

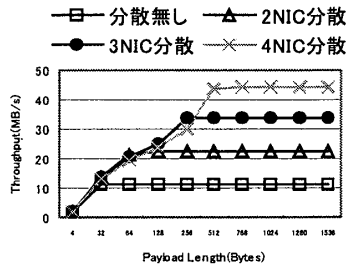


図3 バンド幅性能(環境2)

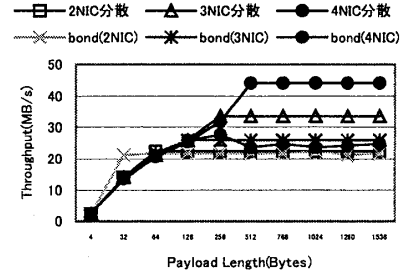


図4 バンド幅性能比較

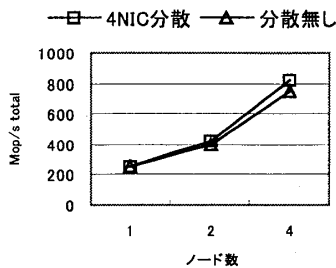


図5 LUクラスW

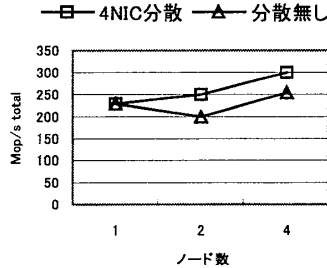


図6 MGクラスW

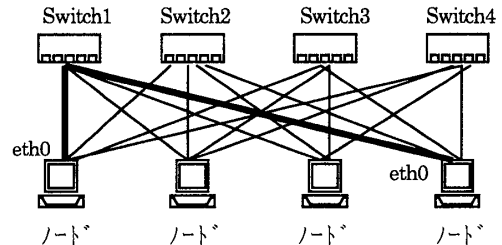


図7 Switchとの接続形態

いことが分かる。原因は CPU アイドル率が2枚に分散した時点でほぼ 0%に達しているためである。これに対し図3から分かるように CPU の処理能力が高い場合、枚数にほぼ比例して性能が向上している。CPU の処理能力向上に伴い TCP/IP 処理のオーバーヘッドは相対的に低くなっていくと考えられる。

### 3.2 Channel Bonding との比較

図4に NIC 分散方式と Channel Bonding 方式のバンド幅の比較結果を示す。測定に用いた環境は前述の環境2に加え Switch を用いた。

Channel Bonding 方式で2枚の NIC を使った場合は、1,024 バイトメッセージで 22.1MB/s の性能を示した。NIC 分散方式と同程度である。Fast Ethernet の物理的性能である 11.9MB/s の 1.8 倍以上の性能を実現した。しかし、3枚以上の NIC を使った場合の性能向上は見られない。受信側 CPU の負荷率が 60%程度であることから CPU の受信処理がボトルネックになっているわけではなく、受信側 NIC による処理がボトルネックになっていると考えられる。

図4の結果より NIC 分散方式は、3枚 NIC の場合 1,024 バイトメッセージで Channel Bonding の 1.29 倍の性能を示し、4枚 NIC の場合 1.88 倍の性能を示している。

### 3.3 NAS 並列ベンチマーク

NIC 分散方式を組み込んだ場合と組み込まない場合について NAS 並列ベンチマーク 2.3 の性能比較を行った。今回は通信が比較的多い LU と MG について測定した。クラスは W を使用した。NIC 分散方式は4枚の NIC を用いた。結果を図5、図6に示す。

結果より LU, MG とともに台数増加による速度向上が見られる。また、NIC 分散方式を利用した場合と利用しなかった場合を比較すると、NIC 分散方式が良い性能を示すことが分かる。MG

クラス W では4ノードの時、分散しなかったときの 1.16 倍の性能を実現している。LU クラス W では4ノードの時 1.09 倍の性能を実現している。

## 4. Switch 接続方式

当 NIC 分散方式は送信側の各 NIC に受信側の各 NIC を対応させ、その間で並列に送受信を行うため、図7の接続形態が可能である。この接続形態を用いれば、複数 Switch を必要とする場合も Switch 間を接続する必要がなく、Switch 間リンクのバンド幅による制限を受けず、さらに複数の Switch に通信負荷を分散するため、各 Switch 内のバンド幅の制限も受けにくい。

## 5. まとめ

本稿では NIC 分散による通信バンド幅の性能向上とそれが並列処理に与える効果について述べた。コードの変更は Channel Bonding 方式で用いられている Bond ドライバに対してのみ行った。カーネルや各 NIC ドライバに対しては変更を行っておらず、既存の TCP/IP 等の標準のプロトコルもサポートし、また、多くの NIC で利用可能である。

通信バンド幅では、4枚の NIC を使って分散した場合は分散しない場合の4倍近い性能を実現できることが分かった。NAS 並列ベンチマークにおいては、わずかではあるが、NIC 分散方式を用いた場合が分散しなかった場合より良い性能を示した。

### 参考文献

- 1) 住元真司, 堀 敦史, 手塚宏史, 原田 浩, 高橋俊行, 石川 裕: 既存 OS の枠組みを用いたクラスタシステム向け高速通信機構の提案, 情報処理学会論文誌, Vol.41, No.6, pp.1688-1696 (2000).