

1 はじめに

Ethernet や TCP/IP に代表される従来型のネットワーク技術を用いたクラスタシステムは、狭帯域、転送遅延、プロトコルのオーバーヘッドなどの問題がある。一方、ATM[3] は通信帯域が 155Mbps~622Mbps と大きく、ネットワークの転送遅延やエラー率も非常に小さい。さらに、ポイントツーマルチポイントコネクション [1] を提供し、上位のグループ通信の機能に対応することができる。

また、MPI[2] は並列プログラムのためのメッセージ通信ライブラリインタフェースの標準仕様である。MPI のグループ通信には、一つのプロセスからグループ内の全てのプロセスへ、同一データを伝送するという broadcast 機能と、一つのプロセスからグループ内の各プロセスへ、異なるデータを伝送するという scatter 機能がある。

本稿では、ATM ネットワークで接続したワークステーションクラスタ上での MPI のグループ通信ライブラリの実装及び通信性能に関する評価実験について述べる。

2 ATM 通信方式モデル

コネクションの確立、解放に要する時間を含む通信時間を定量的に与える ATM 通信方式モデルを提案した。これにより、コネクションオーバーヘッドを含む通信時間を得ることができ、適切なコネクション管理と通信方法を明らかにできる。

実行手順はコネクション要求の送信、結果の待ち受け及び結果の処理という三つ部分から構成される。モデルと実測に基づくパラメータ値に基づいて、コネクション数とデータサイズによる各通信方式を用いた場合の通信時間を計算する。その結果に基づいて適切な伝送方式を選択する。

3 MPI ライブラリの実装

ATM の AAL 層上にグループ通信プロトコルを設計して、MPI ライブラリのプロトタイプ (MPI over ATM) を実装した。設計した broadcast 機能と scatter 機能に

対応する 1 対多通信プロトコルの概要を以下に示す。

1. 受信バッファのオーバーフローを回避するために、ノード受信速度、通信バッファサイズ及び伝送データサイズによる最適な送信レートを求める [4]。broadcast 通信プロトコルは、グループ中の一番低い受信速度のノードに対応する送信レートを選ぶ。
2. 通信プロトコルは、ATM 通信方式モデルを用いた解析に基づいて、より通信時間が短くなるコネクションを利用する方針とする。それにより、broadcast 通信プロトコルはポイントツーマルチポイントコネクションでデータを伝送し、双方向通信可能なポイントツーポイントコネクションで制御メッセージのやり取りを行う。また、ポイントツーマルチポイントコネクション確立の時間のために通信時間が長くなる場合には、ポイントツーポイントコネクションでデータ伝送及び制御メッセージのやり取りを行う。scatter 通信プロトコルは双方向通信可能なポイントツーポイントコネクションでデータ伝送及び制御メッセージのやり取りを行う。
3. broadcast 通信プロトコルは、再送信が必要な場合に、再送信すべきデータに基づいて、再送信時間が短くなるようにポイントツーポイントコネクションもしくはポイントツーマルチポイントコネクションを選択する。

4 実装環境

1. ノード
 - SUN-SS5 × 4 台
MicroSparcII 70MHz, 32MB, Solaris2.5
 - SUN-SS20 × 2 台
SuperSparc 75MHz, 48MB, Solaris2.5
2. ATM スイッチ
 - 住友電気工業 (株)SUMINET-3700SH
SONET/SDH, 155.52Mbps, UTP-5, 8 ポート
3. ATM ネットワークインタフェースカード
 - Efficient Networks, Inc. ENI-155s-U5-c
155.52Mbps, Sbus adapter, 512KB Memory
Efficient Software Development Kit (SDK) 3.3.0

* Implementation and Evaluation of a Collective Communication MPI Library on a ATM Workstation Cluster. Xindan WANG, Noritaka OSAWA and Toshitsugu YUBA, Graduate School of Information Systems, the University of Electro-Communications

5 MPI ライブラリの評価

IP over ATM、IP over Ethernet(10Mbps) 及び MPI over ATM 上で、MPI プログラムを実行し、データ伝送の評価実験を行なった。その結果について考察する。

1. ノード数を6に固定した場合の MPI_Bcast 関数と MPI_Scatter 関数のデータサイズに対する実行時間を図1に示す。

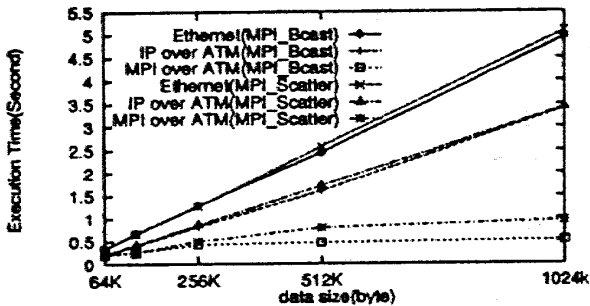


図1. MPI 関数のデータサイズに対する実行時間

IP over ATM と Ethernet 上の MPI_Bcast 関数は、ポイントツーポイント通信で実現するので、MPI_Scatter と同じデータ伝送方法である。その結果、MPI_Bcast 関数の実行時間は、MPI_Scatter とほぼ等しくなる。MPI over ATM 上の MPI_Bcast 関数は、ポイントツーマルチポイントコネクションを利用しているので、MPI_Scatter より実行時間が短くなることわかる。

2. 伝送データサイズが小さい時 (64KB) の MPI_Bcast 関数のノード数に対する実行時間を図2に示す。IP over ATM と Ethernet よりも、MPI over ATM 上での実行時間が長いことわかる。

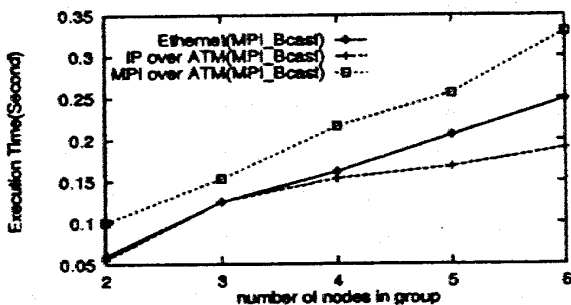


図2. MPI 関数のノード数に対する実行時間 (64KB)

MPI over ATM 上の MPI 関数の実行時間にはコネクションの確立、解放時間を含む。一方、IP over ATM は、各ノードでの IP over ATM 環境を初期化する際に、グループ中の他のノードにコネクションを確立するので、MPI 関数の実行時間がデータ伝送時間だけである。また、Ethernet でデータを伝送する際にも、他のノードに対するコネクションが必要ないので、MPI 関数の実行時間は、データ伝送時間だけである。従って、伝送データサイズが小さい

場合に、MPI over ATM 上の MPI 関数は Ethernet 及び IP over ATM よりも、実行時間が長くなる。

3. 伝送データサイズが大きい時 (1024KB) の MPI_Bcast 関数のノードに対する実行時間を図3に示す。IP over ATM と Ethernet より、MPI over ATM 上での実行時間が短いことわかる。

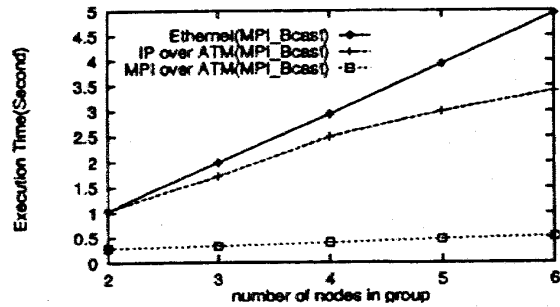


図3. MPI 関数のノード数に対する実行時間 (1024KB)

設計した通信プロトコルは、ATM ドライバの通信機能を直接利用するので、IP パケットヘッダや TCP セグメントヘッダを持たない。さらに、スループットを向上するため、データをユーザバッファに直接に送受信するので、システムバッファ間のメモリコピーがない。その結果、伝送データサイズが大きい場合に、MPI over ATM 上の MPI 関数は Ethernet 及び IP over ATM よりも、実行時間が短くなる。

6 おわりに

本稿では、ATM ワークステーションクラスタにおけるグループ通信のプロトコルを提案して、MPI ライブラリとして実装し、IP over ATM 上と Ethernet 上で動作する MPI ライブラリとの比較を行った。データサイズ、ノード数と MPI 関数の実行時間の関係を調べ、実装した MPI ライブラリの利点、欠点を考察した。今後は、コネクション管理方法の改善及びスキャン、リダクション、バリア同期などの MPI 機能の追加を行う予定である。

参考文献

- [1] The ATM Forum. *ATM User-Network Interface (UNI) Specification Version 3.1*. Prentice Hall PTR, 1995.
- [2] M.Snir, S.W.Otto, S.H.Lederman, D.W.Walker, and J.Dongarra. *MPI The Complete Reference*. The MIT Press, 1996.
- [3] マーチン・ドゥブライカー. ATM 詳解、新世代通信網構築技術. 株式会社プレントイスホール, Sep. 1996.
- [4] 大井拓哉, 大澤範高, 弓場敏嗣. ATM ワークステーションクラスタにおける通信スループットの改善. 分散システム運用技術シンポジウム論文集, pp. 13-18. 情報処理学会, Feb. 1997.