

Beowulf クラスタにおける高精度実行時間測定の検討と評価

早川 潔† 関口 智嗣§ 岩根 雅彦†

Beowulf クラスタシステムは、比較的安価でかつ容易に構築できるため、非常に注目を集めており、Beowulf クラスタでの通信性能をはじめとする種々の処理性能向上に関する研究が盛んに行われている。また、処理性能の高精度な測定は、より正確な評価を可能にし、通信時間の隠蔽などといった並列処理の効率化を可能にする。一般的に処理性能を測定するためには、各ノードの実行開始時刻を正確に揃えなければならない。一般的な Beowulf クラスタでは、各ノードの実行開始時刻を揃えるために、MPI などの通信ライブラリの Barrier 関数が用いられる。しかし、Beowulf クラスタに実装される Barrier では、ある程度の誤差が生じてしまう。そこで、本稿では、PC をベースとした Beowulf クラスタシステム (SCCB-Cluster system) における高精度実行時間測定システムの検討を行った。高精度な測定を可能にするために、Beowulf クラスタに高速なバリア同期を可能にする SCC ボードを搭載した。また、その SCC ボードの中にクロックカウンタを搭載し、疑似的なグローバルクロックを実装する。性能評価として、Beowulf クラスタの collective 通信性能を測定した。SCC ボードでの高速なバリア同期を用いた実行時間測定値は、Ethernet を使用した MPI Barrier を用いた測定値より、安定し、かつ、短い値を示したものであった。

Accurate Performance Measurement System for Beowulf Cluster Systems

Kiyoshi Hayakawa †, Satoshi Sekiguchi §, Masahiko Iwane †

Beowulf cluster consisted of commodity parts, such as PCs and 100base/TX LAN card, is the most remarkable parallel computer system. Collective communications using MPI are the most integral packet forwarding methods on the cluster computing. Accurate performance analysis of collective communication is useful on performance evaluation and prediction of Beowulf cluster system. In order to measure execution time accurately, each node have to take the first step with execution by barrier. But it is difficult for each node to take the first step with execution each other, since it receives the packet indicating barrier completion through Ethernet (i.e. MPI_Barrier) in different time.

This paper describes the Beowulf cluster system (SCCB cluster) that allows us to measure execution time accurately. SCC (Synchronization Communication Controller) was implemented in this cluster system. SCC is able to finish executing barrier less than 10us with 32 nodes. As the performance analysis of MPI collective communication using barrier that SCC performs (SCC_Barrier), performance of MPI collective communication of SCC_Barrier is less fluctuation than that of MPI_Barrier.

1 はじめに

PC や 100base/TX Ethernet などの汎用部品で構成された Beowulf クラスタシステムは、比較的安価でかつ容易に並列処理システムを構築できるため、注目を集めている。また、市販マイクロプロセッサの性能が急激に向上しているため、そのプロセッサを

使用する Beowulf クラスタシステムはより高速な並列処理を可能にしている。よって、Beowulf クラスタでの通信性能をはじめとする種々の処理性能に関する研究 [2][3] は重要な研究の一つであり、処理性能の高精度な測定は、より正確な評価を可能にすることに役立つ [6][5]。また、処理時間を正確に測定し、その測定値を基にしてアプリケーション内の各処理の実行時間を予測することにより、通信時間の隠蔽などといった並列処理の効率化を可能にする。

一般的に処理性能を測定するためには、各ノードの実行開始時刻を正確に揃えなければなら

†:九州工業大学 工学部 電気工学科
Department of Electronic and Computer engineering
Kyushu Institute of Technology
§: 産業技術総合研究所
National Institute of Advanced Industrial Science and Technology

い。一般的に、各ノードの実行開始時刻を揃えるために、MPIなどの通信ライブラリのBarrier関数(MPIBarrier)が用いられる。MPIなどの通信ライブラリは、Ethernetを使用してBarrierを行っている。しかし、Ethernetを使用したBarrier同期では、パケット転送時間やTCP/IPなどの各種通信レイヤでの処理時間のオーバーヘッドが大きいため、ある程度の精度(数100 μ s程度)では揃えられるが、それ以上の精度は望めない。

そこで、本稿では、PCをベースした32台のBeowulfクラスタシステム(SCCB-Cluster system)における高精度実行時間測定システムの検討を行った。高精度な測定を可能にするために、Beowulfクラスタに高速なバリア同期を可能にするSCCボードを搭載した。また、そのSCCボードの中にクロックカウンタを搭載し、疑似的なグローバルクロックを実装する。

性能評価として、SCCB-Clusterのcollective通信性能を測定した。SCCボードでの高速なバリア同期を用いた実行時間測定値とEthernetでのMPIBarrier同期を用いた実行時間測定値とを比較した。

2 SCCB-Cluster システム

図1にSCCB¹-Clusterのシステム構成を示す。SCCBクラスタはPCボード(PICMG規格のボード CPU: PentiumIII600MHz)32台を100Base/TXのSwitching HUBで結合したシステムである。各ノードには、SCC²[1]と呼ばれる高速バリアを可能にするボードが実装されている。各SCCはSync-Comm Networkと呼ばれる同期・通信用ネットワークで結合されている。OSはLinux(kernel version:2.2.15)、通信ライブラリはMPICH(version 1.2.0)である。

3 SCC ボード

SCCボードは、メッセージパッシング型の通信処理とバリアを拡張した同期処理を高速に行うためのPCIボードである[1]。SCCは、高スループットを要求する処理は従来のネットワークボードにまかせ、低レイテンシが要求される同期や小量データパケッ

トのみを処理することにより良い並列処理が実現する。本稿では、SCCの同期処理機能のみを使用して、高精度実行時間測定システムを検討する。

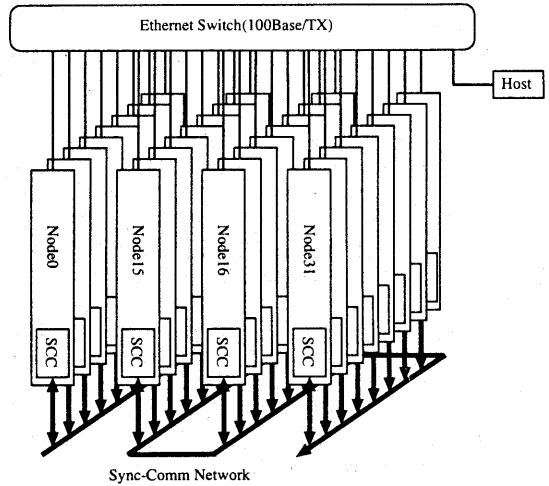


図 1: SCCB-Cluster

3.1 SCC ボードの構成

SCCボードは、Control LSIおよびsync-comm network用コネクタ(Link_AおよびLink_B)のみで構成されている。Control LSIは、同期・通信処理をハードウェアで実装することにより高速化を実現する。sync-comm network用コネクタは、sync-comm network用コネクタのLink_AとLink_Bをデージーチェーンで接続する(図2参照)ことにより、ハブ等の外部接続装置を必要とせず、sync-comm networkが構成可能である。

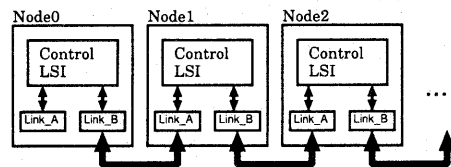


図 2: Sync-comm Network

3.2 SCCのControl LSI

図3にControl LSI内のブロック図を示す。Control LSIは、PCIバス-内部バスブリッジ(PCI

¹Sync-Comm Controller Beowulf

²Synchronization Communication Controller

BUS-Internal BUS Bridge), 同期制御部 (Sync Control), 通信制御部 (Comm Control), 疑似グローバルクロック部 (Pseudo-global clock), およびネットワーク制御部 (Network Control) で構成される。

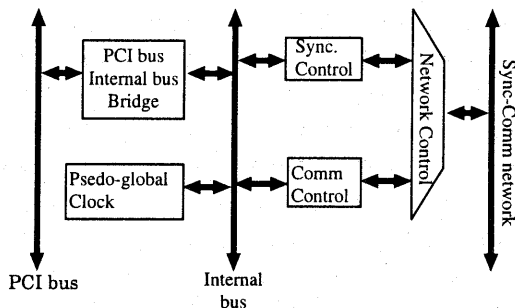


図 3: Control LSI 内のブロック図

PCIバス・内部バスブリッジは,PCIバスプロトコルと内部バスプロトコルをインターフェースし,PCIのマスター/ターゲット機能もサポートする。同期制御部は,sync-comm networkを使用した同期処理アクセスを行い,任意参加バリア,FuzzyバリアおよびRBC同期機構[1]の同期を処理する。通信制御部は,パケットデータ処理を行う。疑似グローバルクロックは,バリア制御部と連携し,各ノードの実行開始時刻を揃えると同時にクロックカウンタをスタートさせるなどの高精度実行時間測定に必要な処理をハードウェアで行う。ネットワーク制御部は,同期信号と通信信号とを sync-comm network 上で融合させるための制御を行う。

3.3 sync-comm network での同期信号

同期時における sync-comm network の信号線は,図4のように接続される。同期時の信号線は,最大16台で構成されるマイクロクラスタ内での同期信号線とそのマイクロクラスタ間で同期を行うための同期信号線に分かれる。

同期成立の検出は,まず,マイクロクラスタ内で同期成立を検出し,その後マイクロクラスタ間で同期成立を検出する。マイクロクラスタ内では, sync_n (nはマイクロクラスタ内でのノード番号)を使用し,

マイクロクラスタを構成する各ノード間の同期成立を検出する。

sync_n 信号は,マイクロクラスタ内ノード番号 n のノードが同期ポイントに到達していることを示す信号である。この信号線は,マイクロクラスタ内のみのバス接続である。

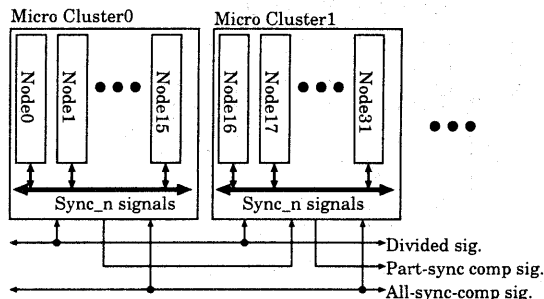


図 4: Sync-Comm Network の同期信号

マイクロクラスタ間では,Divided 信号,Part-sync-comp 信号,All-sync-comp 信号を使用して,クラスタ間の同期成立を検出する。Divided 信号はターゲットクラスタシステムがマイクロクラスタに分割されていること(ノード数が17台以上であること)を示す信号,Part-sync-comp 信号はクラスタ内の同期が成立したことを示す信号,All-sync-comp 信号は,クラスタ全体の同期が成立したことを示す信号である。Divided 信号線,All-sync-comp 信号線は,クラスタ全体のバス接続である。Part-sync-comp 信号線は,デージチェーン接続である。

3.4 同期制御部の機能

同期制御部は,同期成立検出回路(synchronization completion detector)およびR-BC回路で構成される(図5参照)。

同期成立検出回路には,マクロクラスタ内での任意参加バリアをサポートするために,同期グループレジスタ(sync-group reg)(16ビット)が用意されている。また,同期の成立を知らせるために,同期フラグ(S-flag)が用意されている。

バリア同期の動作を以下に示す。

1. バリア同期ポイントに到達したノードは,S-flag をセットする。

2. sync. completion detector がバリア同期が成立したかどうかチェックする。
3. 同期が成立したら,S-flag をリセットする。
4. バリアポイントで実行を中断しているノードは,S-flag のリセット後,実行を再開する。

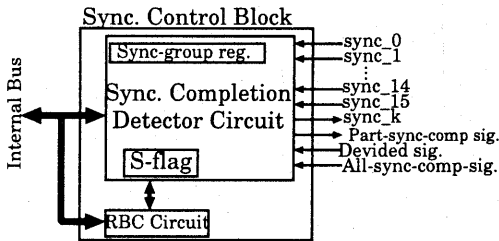


図 5: 同期制御部の構成

4 疑似グローバルクロックカウンタシステム

一般的に,実行時間測定には,各ノードのマザーボード内に搭載されているクロックカウンタを使用する。それらのカウンタを使用した実行時間測定には,ある程度誤差が生じてしまう。あるノード単独で処理する処理時間を測定するには,誤差は生じないが,通信処理などの他のノードと協調して処理する処理時間を正確に測定するには,グローバルクロックカウンタシステムを使用して実行時間を測定する必要がある。

グローバルクロックシステムを実装していないシステムの場合,各ノードの処理開始時刻を揃えた後,各ノードのクロックカウンタをスタートさせ実行時間を測定する。一般的に,各ノードの実行開始時刻を揃えるために,MPIなどの通信ライブラリのBarrier関数(MPI_Barrier)が用いられる。MPIなどの通信ライブラリは,Ethernetを使用してBarrierを行っている。しかし,Ethernetを使用したBarrier同期では,パケット転送時間やTCP/IPなどの各種通信レイヤでの処理時間のオーバーヘッドが大きく,ある程度の誤差で揃えた後に,クロックカウンタをスタートさせざるおえない。よって,各ノードのクロックカウンタによる実行時間測定には,ある程度の誤差が生じてしまう。実行時間測定の誤差を低減させるた

め,SCCに疑似グローバルシステムを実装することを検討する。

4.1 疑似グローバルクロックカウンタシステムの構成

図 6に疑似グローバルクロックカウンタシステムの構成を示す。疑似グローバルシステムは,25MHzのクロック(OSC),クロックカウンタ(32bit),タイムスタンプメモリおよびカウンタコントローラで構成されている。カウンタコントローラがバリア同期とカウンタスタートを連動して行う。つまり,カウンタコントローラがSynchronization Control内のS-flagを常に監視し,S-flagがリセットさせたら,カウンタをスタートさせる。また,ユーザープログラムの要求にしたがって,カウンタの値をタイムスタンプメモリに格納したり,読み出したりする。シミュレーション結果から,このシステムを使用することにより,誤差が240nsまで低減する。

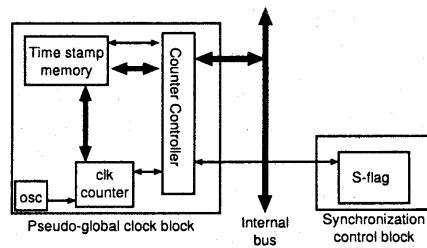


図 6: 疑似グローバルクロックシステム

5 SCC 同期ライブラリ

SCCを使用したバリア同期を行うために,SCC用の同期ライブラリを作成した(図7参照)。SCCライブラリには,SCC_Init(), SCC_finalize(), SCC_Barrier(), SCC_FuzzyBarrier()などが実装されている。SCCライブラリは,SCCのデバイスドライバに対してシステムコールを行うことにより,SCCにアクセスする。デバイスドライバは,PCIバスを介して,SCCからS-flagなどの情報を読み出し,バリアグループなどの情報を書き込む。

SCCライブラリはMPICHプログラムの中で使用することができる。つまり,MPIのプログラムの中に,SCCライブラリの各関数を挿入することができる。Ethernetを使用したバリア同期を行いたい場

合,MPI.Barrierを使用し,SCCを使用した高速バリア同期を行いたい場合,SCC.Barrierを使用する。

Applications	
SCC Interface Library	MPI(MPICH)
Device Driver for SCC Board	Device Driver for Ethernet card
Hardware	

図 7: SCC ライブラリのレイヤ構成

6 collective 通信の性能測定

高精度測定の評価として,SCCでの高速バリア同期を使用した collective 通信の測定を行った。今回の測定は,疑似グローバルクロックシステムを使用するかわりに,SCCでの高速バリアおよびマザーボードのクロックカウンタを使用した性能測定を行った。つまり,高速バリア同期を行った後にマザーボードのカウンタをスタートさせるという動作をソフトウェアで行った。

6.1 高速バリア同期の性能評価

図 8 にノード数と SCC での高速バリア同期レイテンシの関係を示す。SCC を使用した高速バリア同期 (SCC.Barrier) の性能と Ethernet を使用したバリア同期 (MPI.Barrier) の性能を比較した。SCC を使用したバリア同期のレイテンシは,32 ノードまではノード数に関係なく,1.6 μ s または,3.3 μ s である。このシステムは,ほぼ 1.5 μ s に一回の割合で同期フラグにアクセスできる。しかし,ハードウェア見積りに比べて非常に遅いので,図 8 に示す結果を得たと推測される。バリア同期のハードウェアレイテンシの見積りが 64 ノードでも 3 μ s を越えないため,32 ノード以上のバリア同期でも,3.3 μ s のレイテンシで処理できると推測される。

一方,MPI.Barrier は,ノード数が増えるにつれて増加し,32 ノードでは,本同期コントローラのバリア同期レイテンシに比べほぼ 200 倍のレイテンシを得た。

6.2 MPI collective 通信における高精度測定の評価

SCC を使用した高速バリア同期で,開始時刻を揃えて測定した MPI collective 通信の測定結果を図 9, 図 10, 図 11, および図 12 に示す。各図では,Ethernet を使用した MPI.Barrier で,開始時刻を揃えて測定した MPI collective 通信の測定結果と比較した。

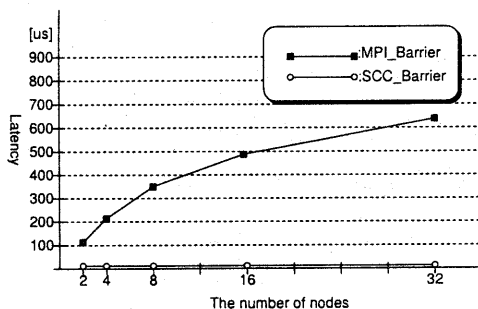


図 8: SCC を使用したバリア同期のレイテンシ

高速バリアで開始時刻を揃えて測定した測定結果は,全てほぼ $\log(n)$ カーブを示している。一方,Ethernet で開始時刻を揃えた測定結果は, $\log(n)$ カーブを示していないものもあり,プロセッサ台数が増えているにもかかわらず,実行時間が減少してしまうという不可解な結果を示している。

MPI.Bcast と MPI.Reduce の一部の結果は以外,高速バリアで開始時刻を揃えて測定した測定結果は,Ethernet で開始時刻を揃えた結果よりも実行時間が短い。

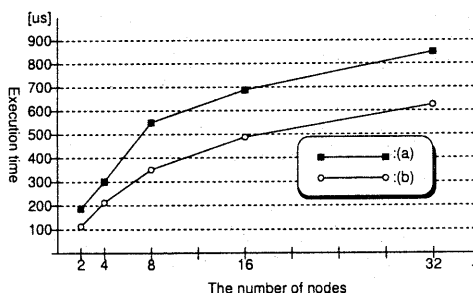


図 9: MPI.Barrier の測定結果。(a):MPI.Barrier で開始時刻を揃えて測定,(b):SCCでの高速バリア同期で開始時刻を揃えて測定。

7 おわりに

PCをベースした32台のBeowulfクラスタシステム(SCCB-Cluster system)における高精度実行時間測定システムの検討を行った。高精度な測定を可能にするために、Beowulfクラスタに高速なバリア同期を可能にするSCCボードを搭載した。また、そのSCCボードの中にクロックカウンタを搭載し、疑似的なグローバルクロックを実装することを検討した。

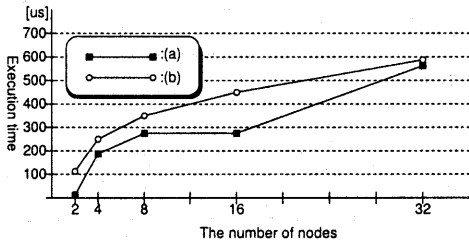


図 10: MPI_Bcast の測定結果

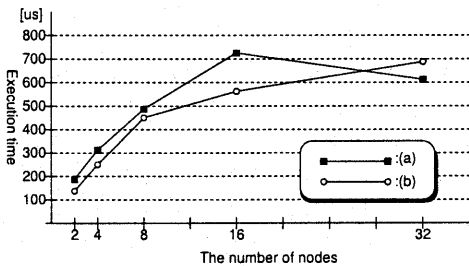


図 11: MPI_Reduce の測定結果

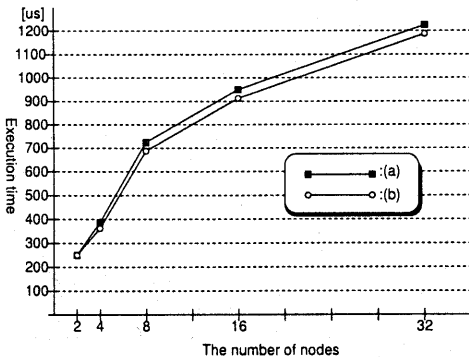


図 12: MPI_Allreduce の測定結果

性能評価として、Beowulfクラスタのcollective通信性能を測定した。高精度な性能測定に必要なSCC

での高速バリア同期のレイテンシを測定した結果、32台まで、台数に関わらず、 $1.6\mu\text{s}$ または $3.3\mu\text{s}$ の高速なバリア同期を実現できた。シミュレーションでは、32台以上でも、最大 $3.3\mu\text{s}$ の値を保持できる結果を得ている。この高速バリア同期を用いて、collective通信の性能測定を行った結果、高速バリアを用いた性能測定では、比較的誤差が少なく測定できていると推測できる。

今後の課題は、疑似グローバルクロックを実装することと、その疑似グローバルクロックの値を使用した処理時間表示ツールを作成することである。

参考文献

- [1] Kiyosh Hayakawa, Satoshi Sekiguchi "Design and Implementation of a Synchronization and Communication Controller for Cluster Computing Systems," Proc. 4th Intel. Conf. High Performance Computing in Asia-Pacific Region, vol.I, pp76-81, May.2000.
- [2] Nanette Boden, Danny Cohen, Robert Feldrman, Alan Kulawik, Chrle Seitz, Jakov Seizovic, and Wen-King Su, "Myrinet - A Gigabit per Second Local Area Network", *IEEE Micro*, vol.15, No.1, pp.29-36, Feb.1995.
- [3] H.Tezuka, A.Hori, Y.Ishikawa, and M.Sato, "PM: An Operating System Coordinated High Performance Communication Library.", *High-Performance Computing and Networking*, vol.1225, pp.708-717, Apr.1997.
- [4] S.Pakin, M.Lauria, and A.Chein, "High Performance Messaging on Workstations: Illinois Fast Messages (FM) for Myrinet", In Proc. of Supercomputing '95, 1995.
- [5] 田中良夫, 久保田和人, 佐藤三久, 関口智嗣 "並列アルゴリズムにおける Collective 通信の性能比較", 情報処理学会研究報告,96-HPC-62, pp.19-26, Aug.1996.
- [6] 益口摩紀, 建部修見, 佐藤三久, 関口智嗣, 長嶋雲兵 "並列システム性能の視覚的解析とその性能評価", 情報処理学会研究報告,99-HPC-75, pp.73-78, Mar.1999.