密結合並列演算加速機構 TCA を用いた GPU 間直接通信による Collective 通信の実装と予備評価

松本 和也^{1,a)} 塙 敏博² 児玉 祐悦^{1,3} 藤井 久史³ 朴 泰祐^{1,3}

概要: 筑波大学計算科学研究センターでは,GPU クラスタにおけるノード間に跨る GPU 間通信のレイ テンシ改善を目的とした密結合並列演算加速機構 TCA (Tightly Coupled Accelerators) を独自開発してい る.本稿では,Broadcast,Scatter,Gather,Reduce,Allgather,Allreduceの6つのCollective通信の TCA による実装と,その性能を TCA 実証環境の GPU クラスタである HA-PACS/TCA において評価し た結果を述べる.TCA による実装は通信レイテンシが問題となる小さめなサイズの Collective 通信におい て,MPI による Collective 通信と比べて高速にその通信処理を行うことが可能であることを示す.実装し た Collective 通信の利用した CG 法の実装およびその性能についても述べる.本研究で用いる CG 法の並 列アルゴリズムは,AllgatherとAllreduceをその通信部分に用いるものである.TCA による Collective 通信を 用いた CG 法実装は,疎行列のサイズ (行数)が数千から数万の場合では MPI の Collective 通信を 用いた実装よりも高い性能を達成できることを示す.

1. はじめに

現在の GPU は高い演算性能とメモリバンド幅を誇り, その利点を活かした GPGPU (General Purpose GPU)処 理がさかんに行われている. GPU は消費電力あたりの演 算性能という点においても CPU を上回り, GPU を計算加 速機構として搭載した GPU クラスタも増加する一方であ る.しかし, GPU クラスタを高性能並列処理に利用する ためには,複数ノードにまたがる GPU 間データ通信が必 要であり,その通信速度は GPU の演算性能と比べて充分 に速いとは云えない.この問題は,多くの GPU クラスタ における高性能アプリケーション開発の障壁となることが 少なくない.

そこで筑波大学計算科学研究センターでは、ノードを跨 ぐ通信に関わるレイテンシとバンド幅の改善を目指して密 結合型並列演算加速機構 TCA (Tightly Coupled Accelerators)を独自開発している [1], [2]. TCA はインターコネク ト・ネットワークに関する技術であり、ノードを跨ぐ GPU などのアクセラレータ間の直接通信を可能にする. 2013 年 10 月からは、GPU クラスタ HA-PACS[3] の拡張部として、 TCA 実証環境である HA-PACS/TCA が稼働している.

TCA については Ping-pong 性能の評価 [2] を始めとし

て、いくつかの性能評価が行われている [4], [5], [6]. 実際 のアプリケーションに TCA を利用した場合の有効性につ いては, 姫野ベンチマーク [6], 格子 QCD のライブラリ QUDA[4], Conjugate Gradient 法 (CG 法) に対してこれ までに行なってきた. これらの性能評価により, TCA は ノードを跨ぐ GPU 間通信を低レイテンシで実現できるこ とが確かめられ, 通信時間が計算時間よりボトルネックと なる強スケーリングが求められるような場合の一部におい ては TCA は有効であることが示されている.

Collective 通信とは,複数ノード/プロセスが同時に関わ る通信操作のことである. MPI 標準仕様においても各種 の Collective 通信が定義されており, MPI を用いて並列 計算を行うほとんどのプログラムに Collective 通信は現れ る. Collective 通信についての研究は多様で,ネットワー クトポロジーを考慮した通信アルゴリズムの研究 [7] や通 信ライブラリの自動チューニングの研究 [8] なども行われ てきた.

本研究では、TCA を用いて各種の Collective 通信(集 団通信)を実装しその性能を評価する.TCA による通信 は片方向通信を基本としており、1対1通信が基本の MPI により記述された並列プログラムを、単純にTCA を用い る形に移植できるわけではない.本稿においては、6つ の Collective 通信(Broadcast, Scatter, Gather, Reduce, Allgather, Allreduce)のTCA を用いた実装について述べ る.そして、TCA を用いた実装と MPI の Collective 通信

¹ 筑波大学 計算科学研究センター

² 東京大学 情報基盤センター

³ 筑波大学大学院 システム情報工学研究科

^{a)} kzmtmt@ccs.tsukuba.ac.jp

実装との性能比較を行い,TCA がどのような条件で有効 であるかを示す.

また,本稿では Collective 通信実装を CG 法の実装へ 利用した結果についても記す.本研究で用いる CG 法の並 列アルゴリズムは,Allgather と Allreduce をその通信部分 に用いるものである.そのアルゴリズムについては以前の 研究会報告 [5] のものと同一であるが,以前の結果よりも Collective 通信実装の最適化を行い性能は改善されている. CG 法の GPU クラスタへの実装に関する研究はいくつか ある [9],[10] が,それらは行列サイズ(行数)が数十万以 上と大規模な疎行列に対する研究である.本研究で特に注 目する行列の行数は数千から数万であり,通常のGPUク ラスタでは並列処理性能が十分に引き出せないような小規 模な疎行列に関する評価を行うという点も本研究の寄与の 一つである.

2. 密結合並列演算加速機構 TCA

2.1 TCA & PEACH2

密結合並列演算加速機構 TCA (Tightly Coupled Accelerators) は、アクセラレータ間 (演算加速機構間)の直接結 合を実現する通信機構技術のことである [1], [2], [11]. TCA の基本的なハードウェア技術は、PCI-Express (PCIe)[12] を応用したものである. PCIe は、シリアルバス・インタ フェースであり GPU ボード, Eathernet ボード, InfiniBand ボードなどの外部機器をホストコンピュータに結合するた めに広く使われている.

PEACH2 (PCI Express Adaptive Communication Hub version 2) は, TCA 技術を実装した PCIe インタフェース・ ボードである [2]. PEACH2 同士を PCIe 外部ケーブルに より結合することにより, 計算機クラスタシステムを構築 することができる.

PEACH2 はノードを跨ぐ GPU 間の直接データ通信を GPUDirect Support for RDMA (GDR) 技術 [13] を利用す ることで実現する.現在,GDR 技術は NVIDIA の Kepler アーキテクチャの GPU ファミリにおいて利用できる. GDR を用いることで,PEACH2 や InfiniBand HCA のよ うな通信アダプタは,GPU メモリへの直接読み書きが可 能となる.GDR は不必要なシステムメモリへのデータコ ピーを削減し,CPU のオーバーヘッドを小さくし,通信レ イテンシを短くする.InfiniBand HCA も GDR による通 信をサポートするが,PEACH2 は InfiniBand に比べて更に レイテンシが小さいという利点がある.それは,PEACH2 を介した通信は PCIe のプロトコルのままで行われるので, InfiniBand を介した通信では必要な PCIe と InfiniBand 間 とのプロトコル変換に伴うオーバヘッドを削減できるため である.

PEACH2 ボードは,最大4GB/sの帯域幅でデータ通信 を行う PCIe Gen2 x8 ポートを4つ持つ.1ポートはホスト

PEACH2は, PIOとDMAの2つの通信方式を備えてい る [2]. PIO 通信は、CPUの store 操作によりリモートノー ドヘデータ書き込みを行う.通信レイテンシが小さいた め、小さなデータの通信に向いている. なお、PEACH2 は CPU 間の PIO 通信のみを提供する. それに対して DMA 通信機能は、PEACH2 チップに 4 チャネル搭載されてい る DMA コントローラにより実現される. DMA 通信は, データの読み込み元と書き込み先の PCIe アドレスおよび 書き込むデータのサイズを記述したディスクリプタに沿っ て行われる. PEACH2 は Chaining DMA 機能を備えてお り、複数のディスクリプタをポインタ連結しておけば、先 頭のディスクリプタに対する通信開始の命令を送ることで 連続した DMA 処理が可能である. DMA 通信のレイテン シは PIO 通信のレイテンシと比べて大きいが、DMA の実 測バンド幅は PIO のバンド幅より大きく,大きなデータの 通信では DMA を用いる方が高速な通信が可能である.

2.2 HA-PACS/TCA

HA-PACS (Highly Accelerated Parallel Advanced System for Computational Sciences)は、筑波大学計算科学研究 センターで開発・運用されている、アクセラレータ技術に基 づく大規模 GPU クラスタシステムである。HA-PACS は、 2012 年 2 月に運用が開始されたベースクラスタ部と、2013 年 10 月に運用が開始された TCA 部 (HA-PACS/TCA)か ら成る。HA-PACS ベースクラスタはコモディティ製品に より構成されているのに対し、HA-PACS/TCA にはコモ ディティ製品に TCA を通信機構として加えた構成である。 HA-PACS/TCA は TCA アーキテクチャの実証環境計算 機であり、PEACH2 ボードの実験環境でもある。本研究 では、HA-PACS/TCA のみを用いる。

表 1 に HA-PACS/TCA の構成仕様を記し, HA-PACS/TCA の計算ノードの構成を図 1 に示す. それ ぞれの計算ノードは 2 ソケットの Intel Xeon E5-260v2 (IvyBridge-EP) CPU と 4 つの NVIDIA K20X (Kepler GK110) GPU を演算装置として持つ. HA-PACS/TCA は計算ノードを 64 持ち, その 64 ノードは 16 ノードずつ 4 つのサブクラスタに分けられる. 16 ノードのサブクラスタ は図 9 のように 2×8 トーラスのトポロジーで PEACH2 に より接続されている. また, HA-PACS/TCA の全 64 ノー ドは, 2 ポートの InfiniBand QDR によるフルバイセクショ ンバンド幅の Fat Tree ネットワークによってもつながれ ている.

3. Collective 通信

本節では, TCA による Collective 通信の実装についてと その性能評価結果を述べる. 性能の測定は HA-PACS/TCA

表1 HA-PACS/TCA システムの構成仕様					
ノード構成					
マザーボード	SuperMicro X9DRG-QF				
CPU	Intel Xeon E5-2680 v2 2.8 GHz \times 2 (IvyBridge 10 cores / CPU)				
メモリ	DDR3 1866 MHz \times 4 ch, 128 GB (=8 \times 16 GB)				
ピーク性能	224 GFlops / CPU				
GPU	NVIDIA Tesla K20X 732 MHz \times 4 (Kepler GK110 2688 cores / GPU)				
メモリ	GDDR5 6 GB / GPU				
ピーク性能	1.31 TFlops / GPU				
インターコネクト	InfiniBand: Mellanox Connect-X3 Dual-port QDR				
	TCA: PEACH2 board (Altera Stratix-IV GX 530 FPGA)				
システム構成					
ノード数	64				
インターコネクト	InfiniBand QDR 108 ports switch \times 2 ch				
ピーク性能	364 TFlops				



図1 HA-PACS/TCA のノード構成

表 2	性能の測定条件
OS	CentOS Linux 6.4
	Linux 2.6.32-358.el6.x86_64
GPU プログラミング環境	CUDA 6.5
C コンパイラ	Intel Compiler 14.0.3
	(Composer XE 2013 sp1.3.174)
MPI 環境	MVAPICH2 GDR 2.0

の1サブクラスタ(最大16ノードまで)を用いて行って いる. HA-PACS/TCA の構成仕様は表1に,性能の測定 条件は表2に記す.本研究では,1ノードあたり1GPUの みを用いている.そのため,これ以降の記述における利用 プロセス数は利用ノード数と一致する.

TCA を用いた実装との比較のために,本節では MPI 実 装の一つである MVAPICH2-GDR 2.0 (以下 MV2GDR) [14] を用いた実装の性能も適宜示す.MV2GDR は,TCA と同様に GPU-Direct for RDMA (GDR) 技術 [13] が実装 に使われている.GDR により GPU メモリと InfiniBand ボードとの間で直接アクセスが可能となり,InfiniBand を 経由した小サイズデータ通信の際のレイテンシが通常の MVAPICH2 と比べて改善されている.MV2GDR は8 KB 以下の通信までは GDR を用い,それ以上のサイズの場合 は CPU メモリを介してパイプライン的にデータを送受信 する.

本稿で述べる Collective 通信は,表3のように通信前 と通信後のデータの状態を記述できる [15].表において, xはサイズ m のデータであり, Collective 通信によって はデータはプロセス数 p によりサイズ m/p の部分デー タ x_i (i = 0, 1, ..., p - 1) に分割される.表中の上付き 文字は,他ノードと結果を Reduce されるデータであり, $\sum_j x^{(j)}$ は Reduction された結果を表す(総和は最も頻出 する Reduction 操作であり、本稿では総和操作の結果を示 す).なお本研究では、使用するプロセス数 p は 2 のべき 乗数 (p = 2, 4, 8, or 16)のみに限定する.以降においては プロセスランク 0 を Root プロセスとみなして実装に関す る記述を行う.

3.1 Broadcast

Broadcast (One-to-all broadcast とも呼ばれる) は,特定のプロセス (Root プロセス) のデータを他の全プロセス に送る Collective 通信である. Broadcast 通信により Root プロセスのサイズ m のデータが p プロセス全てにコピー される.

Broadcast に関しては、表4のようにデータを送信する プロセス数を毎ステップごとに倍にしていき $\log_2 p$ で通信 を完了するアルゴリズムを実装している [16].実装では、 TCA の Chaining DMA 機能を利用し、送信が必要なプロ セスにおいて DMA 通信の開始命令を1回だけ発行すれば 済むようにしている.図2はその実装の性能測定結果で ある.TCA による Broadcast 実装は MPI による実装より 256 KB ほどのサイズまでは高い性能を示している.なお、 サイズの大きいデータの通信において MPI に通信性能が 逆転されることは基本的に避けられない.これは、TCA の PEACH2 は PCIe Gen2 x8 技術を利用しているので4 GB/s が理論ピークバンド幅であるが、MPI はその倍の理 論ピーク性能を発揮する dual-rail InfiniBand QDR を利用 して通信を行われるためである*1.

¹ Dual-rail InfiniBand QDR の理論ピーク性能は 8 GB/s であ り、PCIe Gen3 x8 と同等の性能.

Collective 通信名	通信前			通信後					
Broadcast	Rank 0	Rank 1	Rank 2	Rank 3	Rank 0	Rank 1	Rank 2	Rank 3	
	x				x	x	x	x	
	Rank 0	Rank 1	Rank 2	Rank 3	Rank 0	Rank 1	Rank 2	Rank 3	
	x_0				x_0				
Scatter	x_1					x_1			
	x_2						x_2		
	x_3							x_3	
	Rank 0	Rank 1	Rank 2	Rank 3	Rank 0	Rank 1	Rank 2	Rank 3	
	x_0				x_0				
Gather		x_1			x_1				
			x_2		x_2				
				x_3	x_3				
Beduce	Rank 0	Rank 1	Rank 2	Rank 3	Rank 0	Rank 1	Rank 2	Rank 3	
neuuce	$x^{(0)}$	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$\sum_{j} x^{(j)}$				
Allgather	Rank 0	Rank 1	Rank 2	Rank 3	Rank 0	Rank 1	Rank 2	Rank 3	
	x_0				x_0	x_0	x_0	x_0	
		x_1			x_1	x_1	x_1	x_1	
			x_2		x_2	x_1	x_2	x_2	
				x_3	x_3	x_1	x_2	x_3	
Allreduce	Rank 0	Rank 1	Rank 2	Rank 3	Rank 0	Rank 1	Rank 2	Rank 3	_
Ameuuce	$x^{(0)}$	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$\sum_{i} x^{(j)}$	$\sum_{i} x^{(j)}$	$\sum_{i} x^{(j)}$	$\sum_{i} x^{(j)}$	

表 3 本稿で述べる Collective 通信(4 プロセス使用時)[15]





表 4 Broadcast(8 ブロセス使用時)								
$\setminus \operatorname{Rank}$	0	1	2	3	4	5	6	7
通信前	x							
Step 1	$x \to 4$							
Step 2	$x \to 2$				$x \to 6$			
Step 3	$x \to 1$		$x \to 3$		$x \to 5$		$x \to 7$	
通信後	x	x	x	x	x	x	x	x

3.2 Scatter

Scatter 操作は,Root プロセスが持つデータを他のプロ セスへ送る.この操作は,One-to-all personalized communication とも呼ばれる.Broadcast との相違点は,送られ るデータが各プロセスごとに異なることである.Scatter 操作では,各プロセスに合計 m サイズのデータを送る際 には,プロセスランク0には始めのm/pのサイズのデー タを送り,ランク1には次のサイズm/pのデータを送る ということを行い,最大ランクのランク *p* – 1 には最後の サイズ *m/p* のデータを送る.

Scatter の実装としては,Root プロセスが各プロセスに 対して順番にデータを送信するという単純なものを採用 している.Scatter に関しても Chaining DMA 機能を活用 し,Root プロセスが一回の DMA 通信の開始命令を発行 するだけで済むように実装している.図3に通信時間の測 定結果を示す.この Scatter 実装の最大通信性能は,TCA の GPU 間通信のバンド幅により制限される.各プロセス ごとに送られるデータ量 (m/b)が 64 KB ほどになると性 能が上がらなくなり,上限性能 2.8 GB/s で飽和する.そ の性能が上がらなくなるサイズ近辺で MPLScatter 実装に 性能が逆転されている.



3.3 Gather

Gather (Concatenation 操作とも呼ばれる)は Scatter と対になる Collective 通信で,Root プロセスへ他のプロ セスからデータを集める.Gather に関しても単純な実装 を採用しており,それぞれのプロセスが Root プロセスに 入力データを送信することで実現している.図4はその性 能の測定結果である.性能傾向として Scatter と同様であ るが,Gather 実装は上限性能として3GB/s ほど出ており Scatter よりはわずかながら高い性能を示す.

3.4 Reduce

Reduce (All-to-one reduction とも呼ばれる)は、各プ ロセスにあるデータ同士に加算や乗算などの結合則を満 たす演算を施し、その演算結果を Root プロセスに集める Collective 通信である. Reduce の実装は、始めに必要な データを Root プロセスに集め、その後に Reduction のた めの CUDA カーネルを一度だけ呼び出すことで実現して いる. 図5は Reduce 実装の通信時間を示す. Gather の 性能傾向とは異なり、小さいサイズにおいては TCA によ る現在の Reduce 実装は MPL Reduce と同程度の性能であ り速いとはいえない. 現在の実装では Reduction のための CUDA カーネルは実行終了までに最低でも 17 usec ほどの 時間を要し、その追加時間のために TCA の有利性が小さ くなってしまっている.

3.5 Allgather

Allgather (All-to-all broadcast とも呼ばれる)は、全プ ロセスが一斉に各プロセスへ Gather 操作を行う Collective 通信である。Allgather のアルゴリズムにも様々なものが あるが [16], [17], [18], [19].以前に我々は Allgather アルゴ リズムごとの性能の違いを調べた [5] が、本稿ではその中で もっとも高い性能を示した Recursive Doubling 法 [18] に よる Allgather についてのみ述べる。図 6(a) に Recursive Doubling 法の通信パターンを示す。このアルゴリズムは、 自ノードとデータを交換する通信相手ノードとの距離(ホッ プ数)を毎通信ステップごとに倍にしていく。通信ステッ



図 6 通信パターン. 図中の赤数字は、その通信ステップで送信する データを表す.

プ数は log₂ p で済むが毎回通信相手が異なり通信経路に おいて輻輳を起こしてしまう方法である. Allgather の場 合にはデータ通信量も毎ステップごとに倍になる. この Allgather 実装に関しては, 直前のステップで送られてい きたデータを他のプロセスに送る必要があるため Chaining DMA は使えず, ステップ数の分だけ待ち合わせをし通信 開始命令を発行しなければならない.

Allgather 実装の通信時間の測定結果を図7 (p = 2)と、 図8 (p = 4, 8, 16)に示す.この Allgather に関しては, ノードマッピングが性能に与える影響について調べ,その 中で最適なマッピングを用いている.例えば16 プロセス で128 KB データの Allgather を行う場合,図9(a)に示す 最適なノードマッピングを用いると134.8 μ sec で済むが, 図9(b)に示すマッピングを用いると174.9 μ sec かかって しまう.これは図9(a)のマッピングでは各ステップにおい てノード間通信のホップ数が均一なのに対して,図9(b)の マッピングではホップ数が均一なので輻輳が起きやす くなりそれが性能低下につながっていると考えられる.な お,最適なノードマッピングを用いることによりプロセス



300

ランク間の通信時間のばらつきも小さくなる.後述の CG 法実装ではこの図 9(a) のマッピングを用いている.この マッピングは他の Collective 通信では最適でない可能性が あるが,まだ詳細には評価できていない.

TCA を用いた Allgather 実装は 1200 KB までは MPI_Allgatherより高い性能を示している. これは, 153,600 要素の倍精度浮動小数点数型のベクトルデータを Allgather する際には TCA を用いた方が良いことを意味する. プロ セス数が 4, 8, 16 の場合においても以下のことがいえる.

- プロセス数4のとき,TCAのAllgather 実装は200 KBのサイズまではMPLAllgatherより速い(25,600 要素の倍精度データ).
- プロセス数 8 のとき, TCA の Allgather 実装は 180 KB のサイズまでは MPLAllgather より速い (23,040 要素の倍精度データ).
- プロセス数 16 のとき, TCA の Allgather 実装は 60 KB のサイズまでは MPI_Allgather より速い (7,680 要 素の倍精度データ).

プロセス数が増えると TCA と MPI の差が小さくなる が,TCA と MPI ではノード間接続に用いられているネッ トワークトポロジが異なるということが原因として挙げら れる.TCA の方は PEACH2 により 2 重リング状に接続さ れているが,このトポロジでは同時にデータ通信を行うプ ロセス数が増えると通信経路で輻輳が起こり通信が遅延し てしまう.それに対して MPI の方は,InfiniBand による フルバイセクションバンド幅の Fat Tree ネットワークに







図 9 ノードマッピングの例.番号は自身のプロセスランクを表す.

より接続されているため、プロセス数が増えても通信経路 での輻輳はほぼ起こらない.そのため、プロセス数の増加 が通信時間へ与える影響は MPI の方が TCA と比べて小さ くなっていると思われる.

情報処理学会研究報告 IPSJ SIG Technical Report



図 8 Allgather の通信時間 (p = 4, 8, 16)

表 5 Allredu	ce の通	自信時間	(単位は	$\mu sec)$
プロセス数	2	4	8	16
TCA	2.1	3.4	5.4	7.2
MPI	5.3	8.4	12.8	16.2

3.6 Allreduce

Allreduce (All-to-all reduction とも呼ばれる)は、全プ ロセスが一斉に各プロセスへ Reduce 操作を行う Collective 通信である.現時点において、実装できているのは TCA の PIO 通信を利用した CPU メモリ間の Allreduce だけで ある.Allreduce のアルゴリズムとしては Dissemination 法 [19], [20] を用いている^{*2}. Dissemination 法の通信パ ターンを図 6(b) に示す.この Allreduce のアルゴリズム は、 $\log_2 p$ 回の通信ステップが必要な方法で、Recursive Doubling 法と同様に通信先ノードとの距離を毎通信ステッ プごとに倍にしていく方法である.ただし Dissemination 法は、データを交換するのではなく全ノードの通信方向が 同じになるようにデータを流していく.

その CPU 間 Allreduce 実装の通信時間の測定結果を表 5 に示す.これは,倍精度スカラー変数(8 バイトのデー タ)を Allreduce した場合の結果であり,通信のレイテン シにより性能が決まっている.表に示した結果からもわか るように,TCA による実装は MPLAllreduce の半分以下 の通信時間で Allreduce を実現する.TCA の CPU 間 PIO 通信のレイテンシの短さは,Allreduce 通信において有効 に働いている.

4. CG 法実装への Collective 通信実装の利用

本研究では、TCA により実装した Collective 通信の利用 例として CG 法の実装を行っている [5]. CG 法は、対称正 定値行列を係数行列とする連立一次方程式を解くための反 復法である.本稿において、連立一次方程式はAx = b と 記す.ここでA は $N \times N$ の対称正定値行列であり、x およ び b は N 次元ベクトルである.本研究では、行列データの 格納形式は, Compressed Row Storage (CRS) 形式 (CSR: Compressed Sparse Row 形式とも呼ばれる)[21] を用いる. 浮動小数点演算は倍精度の実数に対して行う.なお, CG 法は前処理を行うことで収束性能を高められる可能性があ るが,本研究の実装では前処理は行っていない.

CG 法の逐次アルゴリズムを図 10 に示す [22], [23]. CG 法の主な演算は, 疎行列ベクトル積計算 (SpMV: Sparse Matrix-Vector multiply), 内積計算 (DOT product), ベク トル加算 (AXPY) である.図 10 のアルゴリズムは毎反 復において ($k \ge 2$), 1 回の SpMV(図 10 の行 11), 3 回の DOT(行 4, 12, 15*3), 3 回の AXPY(行 9, 13, 14) を行う. これらの行列とベクトルに対する 3 つの演算は基本的な演 算であり, CUDA による NVIDIA 社の数値計算ライブラリ でも提供されている.SpMV は cuSPARSE ライブラリ [24] に,DOT と AXPY は cuBLAS ライブラリ [25] にそれぞ れ cusparseDcsrmv, cublasDdot, cublasDaxpy ルーチンと して実装されている.本研究では,それらの cuSPARSE と cuBLAS ルーチンを利用し,これに TCA による複数ノー ドの GPU 間通信を加えて並列 CG 法を実装する.

本研究では、CG 法の並列化として行列 A を単純に一次 元分割する手法を用いる. CG 法を並列化するために、疎 行列 A を行方向にほぼ均等にプロセス数でデータ分割し、 かつベクトル x,b も同割合で均等に分割し各プロセスに初 期データとして持たせる. つまり、プロセス数を p と記述 し $n = \lfloor N/p \rfloor$ とするとき、各プロセスは $n \times N$ の A の部 分行列および n 次元の b と x の部分ベクトルを持つ (ただ し最大ランクのプロセスは $(N - (p - 1)n) \times N$ 行列およ び (N - (p - 1)n) 次元ベクトルを持つ). このようにデー タ分割を行うことにより、CG 法の並列アルゴリズムは図 11 のように記述できる.

各反復において,図 11 の並列アルゴリズムは,図 10 の 逐次アルゴリズムとは以下の点で異なる.

SpMV 計算(行 15)を行う前に、全プロセスが各プロセスに均等に分散されているベクトルデータ pl を集

 ^{*2} Allreduce アルゴリズムごとの性能の違いについても以前に調べ、 その中で Dissemination 法が最も良い性能を示している [5].

^{*3} ベクトルの 2-ノルムは内積計算を用いて計算する.

情報処理学会研究報告

IPSJ SIG Technical Report

1: r := b - Ax2: $norm0 := \operatorname{sqrt}(r^T r)$ 3: for $k := 1, 2, \cdots$ do $\rho := r^T r$ 4: if k = 1 then 5: 6: p:=r7: else 8: $\beta := \rho / \rho_{\text{prev}}$ 9: $p := \beta p + r$ end if 10: 11: q := Ap $\alpha := \rho/(p^T q)$ 12:13: $x := \alpha p + x$ 14: $r := -\alpha q + r$ $norm := \operatorname{sqrt}(r^T r)$ 15:16:if $norm/norm0 < \varepsilon$ then 17:break 18:end if 19: $\rho_{\text{prev}} := \rho$

20: end for

図 10 CG 法の逐次アルゴリズム

1: $x := \text{Allgather}(x_l)$ 2: $r_l := b_l - A_l x$ 3: $d_t := r_l^T r_l$ 4: $norm0 := \operatorname{sqrt}(\operatorname{AllreduceSum}(d_t))$ 5: for $k := 1, 2, \cdots$ do $\rho_t := r_l^T r_l$ 6: 7: $\rho := \text{AllreduceSum}(\rho_t)$ if k = 1 then 8: 9: $p_l := r_l$ 10:else 11: $\beta := \rho / \rho_{\text{prev}}$ 12: $p_l := \beta p_l + r_l$ 13:end if $p := \text{Allgather}(p_l)$ 14: $q_l := A_l p$ 15:16: $\alpha_t := \rho / (p_l^T q_l)$ 17: $\alpha := \text{AllreduceSum}(\alpha_t)$ 18: $x_l := \alpha p_l + x_l$ $r_l := -\alpha q_l + r_l$ 19:20: $d_t := r_l^T r_l$ 21: $norm := \operatorname{sqrt}(\operatorname{AllreduceSum}(d_t))$ 22:if $norm/norm0 < \varepsilon$ then 23:break 24:end if 25: $\rho_{\text{prev}} := \rho$ 26: end for 図 11 CG 法の並列アルゴリズム. 各変数の下付き文字"l"および

X II CG 伝の並列 アルコリズム. 谷変数の下村さ文子 で みよび "t"は、各プロセスごとにローカルに持つ部分データおよび 一時データであることをそれぞれ表す.

める必要がある (Allgather).

 各プロセスごとの DOT 計算(図 11 の行 6, 16, 20)の 後に、そのローカルなベクトル内積の総和を計算し、 全プロセスがその総和を持つ必要がある (Allreduce).

この Collective 通信に TCA により実装した Allgather と Allreduce を利用する. Allgather は各プロセスが持ってい るデータブロックを他のプロセスとやりとりし, Allreduce

表 6 性能評価に用いた疎行列の特性					
行列名	行数 (N)	非零要素数(nnz)	nnz/N		
nasa2910	2,910	174,296	59.9		
s1rmq4m1	$5,\!489$	281,111	51.2		
smt	25,710	3,753,184	146.0		
nd3k	9,000	$3,\!279,\!690$	364.4		
nd6k	18,000	$6,\!897,\!316$	383.2		
nd12k	36,000	$14,\!220,\!946$	395.0		
nd24k	72,000	28,715,634	398.8		

は倍精度の場合は 8 バイトという非常に少量のデータを 他プロセスとやりとりする.以上の特徴から,Allgather は TCA の GPU 間 DMA 通信による実装を利用し,Allreduce は TCA の CPU 間 PIO 通信による実装を利用する*4. TCA による通信を行うためには,DMA 通信の場合は DMA ディスクリプタを作成する必要があり,PIO 通信の場合は PIO 領域の準備が必要である.一度でも通信準備したも のは,同じものを使い回すことが可能である.それゆえ, Allgather と Allreduce 通信部分においては,初めて通信を 行う前に通信準備をしてしまい,それを再利用するように している.

4.1 CG 法実装の性能

本節では、TCA を用いた CG 法実装の性能測定結果を記 す.性能評価に用いる疎行列は、The University of Florida Sparse Matrix Collection[26] から取得した表 6 に記す実数 の対称正定値行列である.なお、実際の使用において、CG 法は解が収束するまで反復する必要があるが、本研究では 性能評価のために反復回数を 1000 回に固定している^{*5}.

図 12 に,各疎行列に対する TCA を用いた CG 法の実 行時間を示す.この図では,プロセス数を1,2,4,8,16 と 増やしたときの実行時間を示している.行列 nd3k, nd6k, nd12, nd24 に関しては,プロセス数を増やすことにより性 能向上を達成できている.行列 smt に関しては,4プロセ ス用いるときが最も高い性能を示す.残り2つの小さめの 行列 (nasa2910 と s1rmq4m1) に関しては,プロセス数 を増やすと性能は悪化する.

性能を更に分析するために,各処理ごとの合計実行時間の 内訳を見る.図13に3種の異なるサイズの行列(nasa2910, smt, nd24k)に対する各処理の内訳を示す.この内訳はプ ロセスランク番号0の結果で,比較のためにMPIを用いた 実装による結果も併記している.この図の結果から云える ことは, nasa2910のように行数n(および行あたりの非零 要素数)が小さすぎると,並列化をしても計算処理(SpMV,

^{*4} cublasDdot は計算したローカルな内積を CPU 側にも返すことが できるので、それにより返されたスカラー値を CPU 間 Allreduce を利用して内積を計算する.

^{*5} 本研究は CG 法における 1 反復当たりの処理時間の評価を目的 としており、収束するか否かは問題としない.性能評価における 反復当たりのバラ付きによる誤差をなくすための十分な反復回数 として 1000 回を選んだ.



図 12 TCA を用いた CG 法の実行時間

DOT, AXPY)の実行時間がほぼ一定で変わらず、データ 通信時間の分だけ遅くなってしまうということである. そ れに対して nd24k のように行列サイズが大きすぎると, 並 列化により性能は向上するが、TCA による Allgather の通 信時間が MPI のものより長くなってしまい、結果として TCA による実装の方が MPI による実装より性能が劣るも のになってしまう. それらの中間ぐらいの行列サイズ (smt のような 15.000 行から 35.000 行ほどのサイズ) の行列に 対しては、TCA を用いることで MPI を用いるよりも良い 性能を実現できている. 最後に, 図 14 に tlog プロファイ リングツールで時刻を基準とした各プロセスの振る舞いを 調べた結果を示す. これは, 行列 smt に関する CG 法実 装のプロファイリング結果である (p=8). この結果では TCA を用いることにより Allgather と Allreduce にかかる 通信時間が削減され、それが性能向上に貢献していること がわかる.

5. おわりに

本研究では、TCA による Collective 通信の実装をし、 HA-PACS/TCA GPU クラスタにおいてその実装の性能評 価を行なっている.本稿では Broadcast, Gather, Scatter, Reduce, Allgather, Allreduce 通信の実装とその性能を述 べた.TCA はノードを跨ぐ通信を低レイテンシで実現す るが、その特徴により小さいサイズの Collective 通信につ いては MPI による Collective 通信と比べて高速にその通 信処理を行うことが可能であった.

実装した Collective 通信を利用した CG 法の実装および その性能についても述べた. CG 法の並列アルゴリズムと しては, SpMV 計算に必要なデータを Allgather で集め, ベクトル内積を Allreduce を利用して実現するものを用い ている. TCA を用いた実装は疎行列のサイズ(行数)が 数千から数万の場合においては MPI を用いた実装よりも 高い性能を示した.

今後は Collective 通信実装の改善を行うと共にその性能 についてより詳細な解析を行う.また、本稿では TCA を 用いて CG 法を実装したが、並列化効果が大きいような大 きな行列に対しては、その通信にかかる時間が MPI を用



(c) hd24k 図 13 CG 法実装の各処理実行時間の内訳 (ランク 0)

いた実装より長くなってしまい必ずしも TCA が有効では なかった. どのような通信と計算を含んだアプリケーショ ンが TCA を用いるのに適しているのかを明らかにする必 要がある.

謝辞 本研究の一部は JST-CREST 研究領域「ポストペ タスケール高性能計算に資するシステムソフトウェア技術 の創出」,研究課題「ポストペタスケール時代に向けた演算 加速機構・通信機構統合環境の研究開発」による.

参考文献

- 塙 敏博,児玉祐悦,朴 泰祐,佐藤三久:Tightly Coupled Accelerators アーキテクチャに基づく GPU クラスタ の構築と性能予備評価,情報処理学会論文誌.コンピュー ティングシステム, Vol. 6, No. 3, pp. 14–25 (2013).
- [2] Hanawa, T., Kodama, Y., Boku, T. and Sato, M.: Tightly Coupled Accelerators Architecture for Minimizing Communication Latency among Accelerators, *Proc. IPDPSW 2013*, IEEE, pp. 1030–1039 (2013).
- [3] 朴 泰祐, 佐藤三久, 塙 敏博, 児玉祐悦, 高橋大介, 建部

情報処理学会研究報告 IPSJ SIG Technical Report



SpMV DOT AXPY Others Allgather Allreduce

修見,多田野寛人,蔵増嘉伸,吉川耕司,庄司光男:演算 加速装置に基づく超並列クラスタ HA-PACS による大規 模計算科学,情報処理学会研究報告, Vol. 2011-HPC-130, No. 21, pp. 1–7 (2011).

- [4] 藤井久史,藤田典久,塙 敏博,児玉祐悦,朴 泰祐,佐 藤三久,藏増嘉仲, Clark, M.: GPU 向け QCD ライブラ リ QUDA の TCA アーキテクチャ実装の性能評価,情報 処理学会研究報告, Vol. 2014, No. 43, pp. 1–9 (2014).
- [5] 松本和也,塙 敏博,児玉祐悦,藤井久史,朴 泰祐: 密結合並列演算加速機構 TCA を用いた GPU 間直接通信 による CG 法の実装と予備評価,情報処理学会研究報告, Vol. HPC-144, No. 12,情報処理学会, pp. 1–9 (2014).
- [6] 藤井久史,塙 敏博,児玉祐悦,朴 泰祐,佐藤三久:TCA アーキテクチャによる並列 GPU アプリケーションの性 能評価,情報処理学会研究報告, Vol. HPC-140, No. 37, pp. 1–6 (2013).
- [7] Sack, P. and Gropp, W.: Faster topology-aware collective algorithms through non-minimal communication, *ACM SIGPLAN Notices*, Vol. 47, No. 8, pp. 45–55 (2012).
- [8] Barnett, M., Shuler, L., van de Geijn, R., Gupta, S., Payne, D. G. and Watts, J.: Interprocessor collective communication library (InterCom), Proc. IEEE Scalable High Performance Computing Conference 1994, IEEE Computer Society, pp. 357–364 (1994).
- [9] Cevahir, A., Nukada, A. and Matsuoka, S.: High Performance Conjugate Gradient Solver on Multi-GPU Clusters using Hypergraph Partitioning, *Computer Science -Research and Development*, Vol. 25, No. 1-2, pp. 83–91 (2010).
- [10] Chen, C. and Taha, T. M.: A Communication Reduction Approach to Iteratively Solve Large Sparse Linear Systems on a GPGPU Cluster, *Cluster Computing* (2013).
- [11] Kodama, Y., Hanawa, T., Boku, T. and Sato, M.: PEACH2: An FPGA-based PCIe Network Device for Tightly Coupled Accelerators, *HEART 2014* (2014).
- [12] PCI-SIG: PCI Express Base Specification Revision 3.0 (2010).
- [13] NVIDIA: NVIDIA GPUDirect, (online), available from (https://developer.nvidia.com/gpudirect) (accessed Oct 10, 2014).
- [14] Panda, D. K.: MVAPICH2-GDR (MVAPICH2 with GPUDirect RDMA), The Ohio State University (online), available from (http://mvapich.cse.ohiostate.edu/overview/) (accessed Nov 7, 2014).

- [15] Chan, E., Heimlich, M., Purkayastha, A. and van De Geijn, R.: Collective communication: theory, practice, and experience, *Concurrency and Computation: Practice and Experience*, No. July, pp. 1749–1783 (2007).
- [16] Grama, A., Karypis, G., Kumar, V. and Gupta, A.: Introduction to Parallel Computing, Addison-Wesley, 2nd edition (2003).
- [17] Chen, J., Zhang, L., Zhang, Y. and Yuan, W.: Performance Evaluation of Allgather Algorithms on Terascale Linux Cluster with Fast Ethernet, *Proc. HPCASIA* '05, IEEE, pp. 437–442 (2005).
- [18] Kogge, P. M. and Stone, H. S.: A Parallel Algorithm for the Efficient Solution of a General Class of Recurrence Equations, *IEEE Transactions on Computers*, Vol. C-22, No. 8, pp. 786–793 (1973).
- [19] Bruck, J. and Ho, C.-T.: Efficient Algorithms for All-to-All Communications in Multiport Message-Passing Systems, *IEEE Transactions on Parallel and Distributed Systems*, Vol. 8, No. 11, pp. 1143–1156 (1997).
- [20] Hensgen, D., Finkel, R. and Manber, U.: Two Algorithms for Barrier Synchronization, *International Journal of Parallel Programming*, Vol. 17, No. 1, pp. 1–17 (1988).
- [21] Saad, Y.: Iterative Methods for Sparse Linear Systems, SIAM, 2nd edition (2003).
- [22] Golub, G. H. and Van Loan, C. F.: Matrix Computations, The John Hopkins University Press, 4th edition (2013).
- [23] Naumov, M.: Incomplete-LU and Cholesky Preconditioned Iterative Methods Using CUSPARSE and CUBLAS, Technical report, NVIDIA White Paper (2011).
- [24] NVIDIA: cuSPARSE Library, (online), available from (http://docs.nvidia.com/cuda/cusparse/index.html) (accessed Nov 7, 2014).
- [25] NVIDIA: cuBLAS Library, (online), available from (http://docs.nvidia.com/cuda/cublas/index.html) (accessed Nov 7, 2014).
- [26] Davis, T. A. and Hu, Y.: The University of Florida Sparse Matrix Collection, ACM Transactions on Mathematical Software, Vol. 38, No. 1, pp. 1:1–1:25 (online), available from (http://www.cise.ufl.edu/research/sparse/matrices) (2011).

図 14 CG 法実装における行列 smt に対する時刻を基準とした各プロセスの振る舞い(8 プ ロセス使用で,2 反復分の結果)