# 密結合並列演算加速機構 TCA による GPU 対応 GASNet の実装と評価

佐藤 賢太<sup>1,a)</sup> 藤田 典久<sup>1,†1</sup> 塙 敏博<sup>2</sup> 松本 和也<sup>3,†2</sup> 朴 泰祐<sup>3,1</sup> Khaled Ibrahim<sup>4</sup>

概要:近年,GPUのような演算加速装置を用いたクラスタがHPC分野で多く用いられるようになってきている. 筑波大学計算科学研究センターでは、ノードを跨ぐ演算加速装置間での直接通信を実現するために、密結合並列演算加速機構TCA(Tightly Coupled Accelerators)を提唱している.このTCAの実装としてPEACH2(PCI Express Adaptive Communication Hub version 2)が開発されており、ノードを跨ぐGPU間での直接通信を行うことができる.しかしながら、TCA/PEACH2を利用するためには独自のAPIを用いる必要があり、プログラミングコストが高く、既存のアプリケーションの移植も容易ではないという問題がある.本研究では、PGAS言語を対象とした通信ライブラリであるGASNet に注目し、これをGPUを対象としたPEACH2に実装する.これによって、GASNetを介して各種のソフトウェアとの互換性が生じ、TCA/PEACH2 が広く利用できるようになると考えられる.既存のGASNetではGPUメモリを対象とした通信しか想定されておらず、GPUを対象とした拡張は現在開発段階にあり、本論文では現在進行中のGASNetのGPU対応の拡張についても触れる.TCA/PEACH2によるGASNetのプロトタイプ実装において、ノードを跨ぐGPU間の通信性能はTCA/PEACH2を直接使用した場合の性能と比較して、最小レイテンシの増大は15%程度に抑えられ、ソフトウェア支援によって最大バンド幅は1.2倍の性能向上を達成した.

キーワード:GPU クラスタ,GPU 間直接通信,通信ライブラリ

## Implementation and Evaluation of GPU-aware GASNet by Tightly Coupled Accelerators

Kenta Sato<sup>1,a)</sup> Norihisa Fujita<sup>1,†1</sup> Toshihiro Hanawa<sup>2</sup> Kazuya Matsumoto<sup>3,†2</sup> Taisuke Boku<sup>3,1</sup> Khaled Ibrahim<sup>4</sup>

Abstract: Recently, PC clusters equipped with GPU as accelerators are widely spread and operated. We have been proposing Tightly Coupled Accelerators (TCA) architecture to realize inter-node direct communication among GPUs, and we developed PCI Express Adaptive Communication Hub version 2 (PEACH2) as a prototype of TCA implementation. However, currently non-standard unique API is required to use TCA/PEACH2, so that the programming cost is expensive and porting of existing applications is not easy. On the other hand, GASNet library developed by Lawrence Berkeley National Laboratory provides low-level communication layer for Partitioned Global Address Space (PGAS) languages such as Unified Parallel C (UPC), Co-Array-Fortran and XcalableMP (XMP), and so on. GASNet assumes only CPU memory as communication target, and extension for GPU-aware GASNet is work in progress now. Beside of GPU-aware GASNet development on commodity network such as InfiniBand, we implement it on TCA/PEACH2 to provide general programming and system software porting on this hardware in this paper. We also mention currently planned features of GPU-aware GASNet. In the case of inter-node GPU communication using GASNet prototype implementation by TCA/PEACH2, the minimum latency is increased only 15 % from the case with native API, and the maximum bandwidth is increased by 1.2 times of native API thanks to software support.

Keywords: GPU cluster, direct communication between GPUs, Communication library

## 1. はじめに

近年,GPUのような高い演算性能とメモリバンド幅を 持つ演算加速装置を用いたクラスタが HPC 分野で広く用 いられるようになってきている.Top500 List[1]において も,2015年11月のリストで全体の13%のシステムが演 算加速装置としてGPUを利用している.一般に,GPUは PCIe(PCI Express)バスを介してCPUやノード内の他の GPUと接続されているが,PCIeの転送バンド幅はGPU のメモリバンド幅よりも低く,GPUアプリケーションにお いてボトルネックとなることも多い.また,ノードを跨ぐ GPU間の通信にはIB(InfiniBand)などのコモディティ ネットワークを介する必要があり,PCIeとIB間でのプロ トコル変換などが必要となるため,通信レイテンシが増大 し,強スケーリングの達成が困難となる.

このような問題を解決するために、筑波大学計算科学研究 センターでは、ノードを跨ぐ演算加速装置間での直接通信を 実現する密結合並列演算加速機構 TCA (Tightly Coupled Accelerators)を提唱しており、その実装として、PEACH2 (PCI Express Adaptive Communication Hub version 2)を 開発している [2]. TCA は演算加速装置間を直接的通信網 で結合するというコンセプトであり、PEACH2 は FPGA を 用いて PCIe を対象として実装した TCA のプロトタイプで ある. 以後、PEACH2 による TCA 実装を TCA/PEACH2 と記す. 現在、TCA/PEACH2 を利用するためには独自の API を用いる必要があるため、MPI と比べてプログラミン グコストが高く、既存のアプリケーションの移植が容易で はないという問題がある.

GASNet はプログラミング言語や通信インタフェースな どに依存しない低レベルな通信レイヤの提供を目指して開 発されている通信ライブラリである [3]. GASNet の API は UPC (Unified Parallel C) [4] や Co-array Fortran[5], XMP (XcalableMP) [6] などの PGAS (Partitioned Global Address Space) 言語において、ランタイムやコンパイラ が生成したコードなどが利用することを想定して設計され ており、MPI と比べて機能は乏しいが軽量なライブラリと なっている. しかし、既存の GASNet は CPU メモリを対 象とした通信しか想定されておらず、GPU メモリに対応 していない. このため、現在、GPU メモリを通信対象と する GASNet の GPU 対応拡張が進められている.

以上の背景のもと、本論文では、GPU 対応 GASNet を TCA/PEACH2 によって実装する.この実装は、現在

1	筑波大学大学院 システム情報工学研究科
	Graduate School of System and Information Engineering
	University of Tsukuba
2	東京大学 情報基盤センター

Information Technology Center, The University of Tokyo 3 筑波大学 計算科学研究センター

Center for Computational Sciences, University of Tsukuba

- <sup>4</sup> Lawrence Berkeley National Laboratory
- <sup>†1</sup> 現在, 筑波大学 計算科学研究センター Presently with Center for Computational Sciences, University of Tsukuba
- <sup>†2</sup> 現在,日本原子力研究開発機構 Presently with Japan Atomic Energy Agency
- <sup>a)</sup> ksato@hpcs.cs.tsukuba.ac.jp

開発中の GPU 対応 GASNet を基準としており,今後の 同システムをベースとするシステムへの適用を想定し ている.これによって,GASNet や上位レイヤの PGAS 言語などを介して各種ソフトウェアとの互換性が生じ, TCA/PEACH2 を広く利用できるようになると期待され る.本稿では,GPU 対応 GASNet の基本的な通信性能に ついて,TCA/PEACH2 と IB のそれぞれの実装で比較し, 評価を行う.また,TCA/PEACH2 による実装では、姫野 ベンチマークの性能について評価を行う.

## 2. TCA/PEACH2

本節では、本研究を理解するための TCA および PEACH2 についての概要を説明する. これらの詳細については [2] を参照されたい.

## 2.1 TCA

筑波大学計算科学研究センターでは、ノードを跨ぐ演算 加速装置間での直接通信を実現するために、密結合並列演 算加速機構 TCA (Tightly Coupled Accelerators)を提唱 している. TCA は、ノードを跨ぐ演算加速装置同士を密 に結合することで、演算加速装置間での直接通信を可能と し、通信レイテンシを削減することで強スケーリングを達 成するというものである.

## 2.2 PEACH2

TCA の実装として, PEACH2 (PCI Express Adaptive Communication Hub version2)が開発されており, ノー ドを跨ぐ CPU 間および GPU 間での直接通信ができる. 現状では, NVIDIA 社製 GPU に対応している. PEACH2 では, ノード間の接続に PCIe を用いており, PEACH2 が PCIe パケットのルーティングを行うことでノードを跨ぐ PCIe デバイス間での直接通信を実現している.

PEACH2 は 4 つの PCIe Gen2 x8 ポートを持っており, そのうち 1 ポートをホスト CPU との接続に使用し,残り 3 ポートを他のノードの PEACH2 との接続に使用する. ノード内の PCIe デバイスには,CPU 内部の PCIe スイッ チを介して接続されている.PEACH2 による GPU メモリ へのアクセスには,GDR (GPUDirect RDMA) [7] を利用 している.

PEACH2 はリモートノードに対するアクセスとして, RDMA Write にのみ対応おり, RDMA Read などが必要 な場合は,ソフトウェアで処理する必要がある. この点は アプリケーションやシステムソフトウェアの実装および性 能上,大きな影響を及ぼすので注意が必要である.

## 2.3 DMA 通信

PEACH2 には, DMAC (DMA Controller) が4 チャネ ル実装されている.これらの DMAC を利用した DMA 通 信は,ホスト上であらかじめ転送元アドレスや転送先アド レス,転送サイズなどを指定したディスクリプタを作成 しておき,使用する DMAC にこれを登録することで行わ れる.また, PEACH2 の DMAC は Chaining 機能を持っ ており、複数のディスクリプタを連結することで、連続し た DMA 処理を 1 回の DMA 要求で開始できる.ただし、 ディスクリプタの作成には時間がかかるため、できる限り 事前に作成しておき、以降は事前に作成しておいたディ スクリプタを DMAC に登録するだけに留めるのが望まし い.また、PEACH2 では通常の連続アドレス上のデータ のブロック転送の他、ブロックストライド転送も実行可能 である.したがって、定型的なブロックストライド転送を DMA Chaining を用いなくても実行可能である.

ディスクリプタを格納する領域として,FPGA に内蔵さ れたメモリ(内蔵メモリモード)と,ホスト CPUのメモ リ(ホストメモリモード)がある.内蔵メモリモードでは, ディスクリプタを読み出すコストが低いため,ホストメモ リモードよりも通信レイテンシが低くなる.しかしながら, メモリ容量には限りがあるため,最大で 1024 個のディス クリプタしか格納できない.また,ディスクリプタを作成 する度に,PEACH2 にアクセスする必要があるため,ディ スクリプタ作成にかかる時間が大きい.一方,ホストメモ リモードは通信レイテンシが若干増えるものの,ホストの メモリ容量の許す限りディスクリプタを格納することがで き,またディスクリプタの作成はホスト内で完結するため, 作成にかかる時間が小さい.このように,各モードで特性 が異なるため,通信パターンなどに応じて適切に選択する ことで最適化が可能である.

PEACH2の DMA 通信には、ハードウェアなどの制約 から、転送元や転送先のアドレスのアラインメントなどに 関する制限があり、任意のデータを対象として通信を行う ことはできないので注意が必要である.また、CPUメモ リを通信対象とする場合には、tcaMalloc()という専用の APIで確保された領域のみ通信対象とすることができる.

## 2.4 PIO 通信

PIO(Programmed I/O)通信は CPU の Store 命令に よってリモートノードのメモリに書き込みを行うというも のである. DMAC を使用しないため、ディスクリプタの作 成や DMAC の起動処理などが不要となり、非常に低いレ イテンシで通信を行うことができるが、バンド幅は低いた め、小さいデータの転送に適している.また、PIO 通信は 転送先のメモリとして *tcaMalloc()* で確保された CPU メ モリのみ対応しており、DMA 通信と同様にアラインメン トに関する制限もあるので注意が必要である.

#### 2.5 HA-PACS/TCA

筑波大学計算科学研究センターでは,TCA/PEACH2の 実験用システムとして,HA-PACS/TCA (Highly Accelerated Parallel Advanced System for Computational Sciences/TCA) が運用されている [8]. 表1にノード構成を 示す.本稿における今後の性能評価は,このクラスタを用 いて行う.



#### 2.6 PEACH2 の基本性能

今後の実装や性能評価の基準となる PEACH2 の基本的 な通信性能を図1および図2に示す. 図中の凡例における "PIO"は PIO 通信を示しており, "Internal" および "Host" は DMA 通信における DMA ディスクリプタの配置の違い で,前者は内蔵メモリモード,後者はホストメモリモード をそれぞれ示している.また,括弧内は転送元と転送先の メモリを示している.

図1より, PIO 通信の最小レイテンシは 0.9  $\mu$ s となって おり, DMA 通信の最小レイテンシは内蔵メモリモードの 場合はメモリの種別に依らず 2.0  $\mu$ s, ホストメモリモード の場合は CPU メモリで 2.2  $\mu$ s, GPU メモリで 3.0  $\mu$ s であ る.また, 図2より, PIO 通信の最大バンド幅は 0.1 GB/s, DMA 通信の最大バンド幅はモードに関わらず, CPU メモ リの場合は 3.5 GB/s, GPU メモリの場合は 2.9 GB/s で ある.

## 3. GASNet

本節では、本研究を理解するために GASNet やその API についての概要と、現在 LBNL で進められている GPU 対 応拡張について説明する. API の詳細については [3] を参 照されたい.

GASNet は図3に示すように Core APIと Extended API の2つの通信レイヤに分かれている. Extended APIには, Core APIのみを用いて実装されたリファレンス実装があ り, Core APIを実装することで GASNet の全ての APIが 利用可能となる. そして Extended APIの特定の APIの みを置き換えることが可能となっているため, ハードウェ



図 3 GASNet のソフトウェアスタック

アが対応している処理などは、リファレンス実装を使用せ ずにローカライズされた実装をすることで通信性能の最適 化が可能である.

## 3.1 Core API

Core APIは、GASNet の初期化や終了処理などの制御系 の APIと AM (Active Message)による通信 APIで構成さ れている。AM は RPC (Remote Procedure Call)の一種 で、あらかじめハンドラテーブルにハンドラ関数を登録し ておき、リモートノードから、ハンドラや引数、データを指 定し、ハンドラに対応する関数を呼び出すというものであ る。GASNet における AM は、データ転送の可否や転送先 のメモリ領域などに応じて Short, Medium, Long の3種類 が定義されている。以後、これらをそれぞれ AM\_Short(), AM\_Medium(), AM\_Long()と記す.

AM\_Short()では、データを転送はできず、ハンドラと引数 のみを転送できる。そして、AM\_Medium()とAM\_Long() ではAM\_Short()に加えてデータも転送することができる。 AM\_Medium()におけるデータの転送先はシステムが管理 するバッファとなるが、AM\_Long()ではユーザが指定した 領域を転送先として指定することができる。

## 3.2 Extended API

Extended API は put() や get() といった RMA (Remote Memory Access) 通信やバリア同期などの API で構成され ている.また,非公式 API として各種の Collective 通信や 非連続なデータ転送のための API も存在し,非連続データ の転送に関しては, Vector, Indexed, Strided の3種類が あり,それぞれに put() と get() が用意されている.これ らの API は現時点では非公式だが,現在開発が進められて いる GASNet-EX[9] では,公式 API になる予定である.

### 3.3 Segment

GASNet では、ライブラリ初期化時にユーザが指定した サイズのメモリを各ノードで確保する.このメモリ領域 を Segment と呼び、リモートノードからアクセスされる メモリは、この領域に収まっている必要がある.つまり、 AM\_Long()や put()における転送先と、get()における転 送元は Segment 内に限定される.

## 3.4 GPU 対応

前節で説明したように、GASNet ではリモートノードか らアクセスされるメモリに関しては Segment 内である必 要がある.そのため,Segment に GPU メモリを割り当て ることで、リモートノードの GPU メモリへのアクセスが 可能となる.ただし、Segment は単一の連続領域である必 要があるため, Segment に GPU メモリを割り当てた場合, 基本的にリモートノードの CPU メモリにアクセスするこ とはできなくなる. 例外として, AM\_Medium()に関して は,転送先はシステムが管理するバッファであり,Segment 内である必要はない. ここに CPU メモリを割り当てるこ とで、転送サイズなどに制限はあるものの、リモートノー ドの CPU メモリにデータを転送することができる.通常, GPU を利用するアプリケーションでは、主として GPU to GPU の通信行うため、この問題による大きな影響はない と考えられる.また、GASNet-EX[9]では、各ノードが複 数の Segment を持つことが可能となる予定であり、CPU メモリと GPU メモリをそれぞれ異なる Segment に割り当 てることで、このような問題は生じないと考えられる。一 方で,ローカルのメモリ領域に関しては特別な規定はなく, CPU メモリと GPU メモリの両方を扱うことができる.

## 4. 実装

本節では、本研究で開発した TCA/PEACH2 による GPU 対応 GASNet の実装について説明する.

## 4.1 概要

GASNetを実装するのに必要な機能として,各種の制御 メッセージを送受信する機能と,任意のデータに対応し たデータ転送機能が挙げられる. Core API で必要な AM は,必要に応じてデータの転送を行い,その後にハンドラ や引数などの情報を含むメッセージをリモートノードに送 信することで実装できる. Extended API の put() に関し てはデータ転送機能そのものであり,get()に関してもメッ セージの送受信やデータの転送の組み合わせで実装でき る.また,非公式 API のストライド転送機能は,PEACH2 にブロックストライド転送機能を利用することで,高い性 能が期待できる. Extended API の他の機能に関しては, PEACH2 との親和性が低く,直接実装しても高い性能が得 られるとは考えにくいため,リファレンス実装を利用する.

## 4.2 パケット通信

図4にパケット通信機構の構造を示す. 2.2 節で説明し たように, PEACH2 ではリモートノードのメモリ領域に



図 4 パケット通信機構

対する RDMA Write しか行うことができない. そのため, 本実装では,各ノードペア間にリングバッファを用意し, 送信バッファから受信バッファに対して RDMA Write を 行うことでパケットを転送する. これによって,ローカル の送信バッファとリモートノードの受信バッファには全く 同じデータが格納されている状態になる. その後,リモー トノードに対してリングバッファの更新を通知すること で,リモートノードはパケットの到着を検出することがで きる. リモートノードでは,到着したパケットの受信処理 が完了すると,送信元に受信完了を通知する. これによっ て,リングバッファのエントリが再利用可能な状態となる. この機構では,送信側が受信バッファの空き状況を把握す ることができるため,受信バッファに空きがない場合は送 信を停止することでフロー制御が可能となっている.

リングバッファに用いるメモリ領域は、ライブラリ初期 化時に tcaMalloc() で確保し、固定長単位で利用する.た だし、実際に転送されるパケットのサイズは可変となって おり、転送するメッセージ長に応じて変化する.2.6 節で の測定結果より、転送サイズが小さい場合は PIO 通信の 方が高速な転送が可能である.そのため、パケットの転送 は、パケットサイズが小さい場合は PIO 通信を利用し、大 きい場合は DMA 通信を利用する.

## 4.3 データ転送

2.3 節や 2.4 節で説明したように, PEACH2 の DMA 通 信や PIO 通信にはいくつかの制限があり, 任意のデータ を転送できるわけではない. このような制限はユーザが TCA の APIを用いて直接プログラムを開発する場合はそ れほど問題にはならないが,本研究のように通信ライブラ リを実装する場合は,上位レイヤのアプリケーションなど にハードウェアの制限を意識させないための一般化が必要 であり,任意のデータを柔軟に転送できるようにする必要 がある.

前節で説明したパケット通信によってデータを転送する 場合,データは送信バッファと受信バッファを経由するた め,任意のデータを転送することができる(図5(a)).しか し,この方法では2回のメモリコピーが行われるため高い 通信性能は期待できない.他方,3.3節で説明したように, GASNetではリモートノードからアクセスされる領域は Segment内である必要があり,転送元と転送先のどちらか 一方は必ずGPUメモリとなる.具体的には,AM\_Long() と put()の場合は転送先がSegment内となり,get()の場合 は転送元がSegment内となり,必ずGPUメモリとなる. 転送元と転送先の両方がSegment内の場合は,転送元から 転送先までDMA通信によって直接転送を行う(図5(b)). この場合,メモリコピーは発生しないため,高い通信性能が 得られると考えられる.また,転送先がSegment内の場合 は、送信バッファから転送先までは DMA 通信によって直 接転送でき(図5(c))、逆に転送元が Segment 内の場合は、 転送元から受信バッファまでは DMA 通信によって直接転 送できる(図5(d)). この場合も、メモリコピーは1回で 済むためパケット通信を利用する場合よりは高い通信性能 が得られると考えられる.以後、パケット通信によるデー タ転送を BtoB (Buffer to Buffer)モード,DMA 通信に よる直接転送を MtoM (Memory to Memory)モード、送 信バッファから転送先への転送に DMA 通信を使うものを BtoM (Buffer to Memory)モード、転送元から受信バッ ファへの転送に DMA 通信を使うものを MtoB (Memory to Buffer)モードと記す.

パケット通信では、パケットサイズが小さい場合は PIO 通信を利用するため、DMA 通信を用いるモードよりも高 速にデータを転送できる.そのため、転送するデータサイ ズが小さい場合は、その他のモードが利用可能な場合でも BtoB モードを利用して転送を行う.また、転送元や転送 先のアラインメントがずれている場合は、転送するデータ の先頭から数バイトのみを BtoB モードで転送してアライ ンメントの調整を行う.これによって、どのような場合で あっても MtoM, BtoM, MtoB のいずれかの転送モード が利用可能となり、BtoB モードは基本的に DMA 通信よ りも PIO 通信の方が高速な小さいデータの転送にのみ利 用されることになる.

MtoM モード以外のモードではメモリコピーが必要とな るが、メモリコピーと DMA 通信をオーバーラップするこ とで、メモリコピーにかかる時間を隠蔽することができる. また、MtoB モードで使用する受信バッファは受信側で転 送先へのメモリコピーなどを行う必要があり、フロー制御 が必要なため、BtoB モードのパケット通信で使用する受信 バッファをそのまま使用する. BtoM モードで使用する送 信バッファに関しても、BtoB モードで使用する送信バッ ファをそのまま使用することもできるが、敢えて専用の送 信バッファを2つ用意し、ダブルバッファリングの要領で 交互に使用する. この理由に関しては、次節で説明する.

## 4.4 ディスクリプタキャッシュ

2.3 節で説明したように、DMA 通信で必要なディスク リプタの作成には時間がかかるため、できる限り事前に作 成したものを再利用するのが望ましい.しかし,通信ライ ブラリを実装する場合は、事前にどのような通信を行うか を予測して作成しておくのは困難である.一方で姫野ベン チマークを始め, HPC アプリケーションでは転送元や転 送先の領域や転送サイズが変化しない固定パターンの通信 を何度も繰り返すというアプリケーションが数多く存在す る.このようなアプリケーションでは、事前に通信パター ンを予測できなくても、一旦作成したディスクリプタをそ のまま再利用できる可能性は高いと考えられる. もちろん, ユーザが直接 TCA/PEACH2 を使ってアプリケーション を記述する場合は、ディスクリプタの再利用は自然に行わ れるが、本研究では一般化された通信ライブラリを実装す るため、ディスクリプタの再利用性をライブラリ内で管理 する必要がある. そこで, 本実装では一旦作成したディス



クリプタをキャッシュする機構を備える.このようにディ スクリプタをキャッシュすることによって,ディスクリプ タ作成に時間がかかるという問題をある程度緩和できると 考えられる.

キャッシュのアルゴリズムに関しては、対象とするアプ リケーションの通信パターンによって最適なものが異なる. そこで、本実装では LRU と LFU のような一般的なキャッ シュアルゴリズムの他、キャッシュエントリが埋まるまで 先着順に割り当てるというアルゴリズムを実装しており, プログラム実行時に切り替えることで、アプリケーション の特性に応じた最適化を可能としている.また,2.6節での 測定結果より、ホストメモリモードのディスクリプタによ る DMA 通信の最小レイテンシは内蔵メモリモードよりも 大きいため、頻繁に発生する通信には内蔵メモリモードで ディスクリプタを作成するのが望ましい.一方で,内蔵メ モリモードは作成できるディスクリプタ数に限りがあり, ホストメモリモードよりも作成に時間がかかるため、低頻 度な通信での利用や頻繁な置換は避けるべきである.一旦 ホストメモリモードで作成した後、使用回数が一定回数以 上となったディスクリプタを内蔵メモリモードで再作成す れば、特に頻繁に発生する通信のみを内蔵メモリモードで 作成することができる.一般的に,HPC アプリケーション では、同じデータ通信パターンが多数反復されることが多 く、この機構は有効に動作すると考えられる.ただし、通 信パターンが少なく、全てのディスクリプタを内蔵メモリ モードで作成できる場合などに対応するため、この機構は 実行時に無効化することも可能としている. また, MtoM モードでは通信パターン毎にディスクリプタを1つ, BtoM モードでは送信バッファが2つあるため,通信パターン毎 にディスクリプタを2つ使用するが、BtoBモードや MtoB モードでは通信パターン毎にバッファのエントリ数分の ディスクリプタを使用するため、キャッシュしたディスク



図 6 Page-Locked Host Memory を経由したデータ転送

DeviceToHost

リプタの利用効率が悪く、内蔵メモリを大量に消費してし まう.そのため、内蔵メモリモードの利用は MtoM モード と BtoM モードに限定している.また、BtoM モードで送 信バッファに BtoB モードの送信バッファを使用せずに専 用のバッファ領域を用意するのはこれが理由である.

## 4.5 gdrcopy

MtoM 以外のデータ転送モードでは, Segment である GPU メモリと送信バッファや受信バッファとの間でメ モリコピーが必要となる.通常, GPU メモリのコピーに は CUDA[10] の API である cudaMemcpy() を利用するが, cudaMemcpy() によるメモリコピーはレイテンシが大きい という問題がある.そこで,本実装では, GPU メモリの コピーに gdrcopy[11] を利用する.

gdrcopy は NVIDIA によって提供されている GDR を用 いて GPU メモリのコピーを行うためのライブラリである. gdrcopy によるメモリコピーの特徴として,レイテンシが 非常に小さく, CPU to GPU のバンド幅が高いというもの がある.一方で,GPU to CPU のバンド幅は非常に低い ため,サイズが小さい場合にしか利用できない.しかし, GPU to CPU の転送は基本的に MtoB モードで転送され るため,GPU メモリのコピーはアラインメントの調整で のみ行われ,それほど大きな影響はないと考えられる.

## 4.6 Page-Locked Host Memory

2.6 節での測定結果より、PEACH2のDMA 通信でGPU メモリを転送した場合のバンド幅は512 KB をピークとして これ以上の転送サイズでは低下している. この現象は GDR に対応した IB HCA を用いて GPU メモリを転送した場合 でも同様に観測されるため、GDR を用いた GPU メモリの 読み出し性能の限界だと考えられる. 一方で, Page-Locked Host Memory[12] を転送先として, cudaMemcpy()を用い て GPU メモリのコピーを行うと, GDR を用いる場合よ りも高いバンド幅を得ることできる. そこで, MtoM モー ドや MtoB モードで大きいデータの転送を行う際に,一旦 Page-Locked Host Memory を経由することで, GDR によ る読み出し性能の限界を超えたバンド幅での転送が可能 になると考えられる. ただし, Page-Locked Host Memory は DMA 通信の通信対象とすることはできないため,図6 に示すように,送信バッファを経由する必要がある.した がって, MtoM モードは実際には BtoM モードでの転送 となり、MtoB モードは実際には BtoB モードでの転送と なる.

## 5. 性能評価

本節では、本研究で実装した GASNet の基本的な通信性





能と袖領域通信の性能,GASNetを用いて姫野ベンチマー クを実装したときの性能について評価を行う.また,HA-PACS/TCAに搭載されている IB HCA は QDR Dual-port だが,本研究では1ポートのみ使用する.

#### 5.1 基本通信性能

通信性能の評価は、GASNetの仕様上利用可能な GPU to GPU, CPU to GPU, GPU to CPU の3種類の転送パ ターンについて行う. GPU to GPU と CPU to GPU に関 しては、put() と AM\_Long() で利用可能であるが、put() を 利用する場合は、受信側で転送完了を検出するために GPU メモリのポーリングが必要となり、オーバヘッドが大きい. そのため、転送先が GPU メモリの場合は AM\_Long() によ る Ping-Pong 通信によって評価を行う. GPU to CPU に関 しては、get()の処理時間によって評価を行う. また、GPU to GPU に関しては、姫野ベンチマークで利用するため、 MPI でも MPI\_Send() と MPI\_Recv() によって Ping-Pong 通信を行い GASNet と比較する.

図7にレイテンシを、図8にバンド幅を示す. 凡例は,通 信に使用した実装の種類を示しており、"GASNet/TCA" はTCA/PEACH2によるGASNet実装、"GASNet/IBは IB"によるGASNet実装、"MPI/IB"はMPIによる通信 を示している.また、括弧内は転送元と転送先のメモリ種 別を示している.

GPU to GPU の最小レイテンシは, GASNet/TCA の 場合は 2.3 µs, GASNet/IB の場合は 1.9 µs, MPI/IB の 場合は 2.1 µs となっており, GASNet/IB が最も低いこと

がわかる. GASNet/TCA と MPI/IB でレイテンシを比較 すると,最小レイテンシは MPI/IB のほうが低いものの, 32 B ~ 64 KB では GASNet/TCA の方が低いことがわか る. 最大バンド幅に関しては, GASNet/TCA の場合は 3.4 GB/s, GASNet/IBの場合は2.9 GB/s, MPI/IBの場合は 3.8 GB/s となっており, MPI が最も高いことがわかる. GASNet/IB のバンド幅は1 MB をピークとして低下して いるが、これは PEACH2 の DMA 通信と同様に GDR に よる GPU メモリの読み出し性能の限界だと考えられる. 一方で,GASNet/TCA では,このようなバンド幅の低下 は見られず, PEACH2 における DMA 通信の最大バンド 幅の 1.2 倍の性能が得られている. これは, 4.6 節で説明 した Page-Locked Host Memory を経由した GPU メモリ の転送が有効に機能しているためだと考えられる.また, MPI/IB でもバンド幅の低下は発生しておらず、本研究と 同様の通信の最適化が行われていると考えられる.

CPU to GPUの最小レイテンシは、GASNet/TCAの場 合は 1.6  $\mu$ s, GASNet/IBの場合は 1.8  $\mu$ s となっており、 転送元が CPU メモリのため GPU to GPU よりも低いレイ テンシで通信が行えていることがわかる.GASNet/TCA の最小レイテンシは PEACH2の DMA 通信よりも低く、 PIO 通信と比較したオーバヘッドは 0.7  $\mu$ s である.最大 バンド幅に関しては、GASNet/TCAの場合は 3.4 GB/s、 GASNet/IBの場合は 3.3 GB/s となっている.

GPU to CPUの最小レイテンシは、GASNet/TCAの場 合は 3.5 µs, GASNet/IB の場合は 2.1 µs となっており, get()を行っているため 3 種類の転送パターンの中で最も レイテンシが大きい. IB はハードウェアで get()を処理す ることができるためレイテンシの増加はそれほど大きく ないが、TCA/PEACH2 ではソフトウェアで処理を行って いるため最小レイテンシが大きく、さらに、MtoB モード によって転送を行っているため、内蔵メモリモードが利用 できず、全体的にレイテンシが大きくなっていると考えら れる.最大バンド幅に関しては、GASNet/TCA の場合は 3.4 GB/s、GASNet/IB の場合は 3.0 GB/s となっており, GASNet/IB では GPU to GPU と同様のバンド幅の低下 が見られる.

#### 5.2 袖領域通信の性能

多次元ステンシル計算などでは、連続領域のデータ転送 だけでなく、ブロックストライドあるいはストライド配置 のデータ転送性能も重要となる. TCA/PEACH2では、こ のような需要を想定して、ハードウェア実装によるブロック ストライド転送機能を持っている. GASNet の Extended API においても非公式ではあるが、ストライド転送の API が用意されており、本研究ではこれに TCA/PEACH2 実 装を適用する最適化を行っている.

本節では、GPUメモリに確保した3次元配列の各面を 対象とした Ping-Pong 通信によって袖領域通信の性能評価 を行う.3次元配列を対象とした袖領域通信では、転送す る各面に応じて、ブロック転送、ブロックストライド転送 ストライド転送が行われる.GASNetではストライド転送 機能を利用してデータを転送した後、AM\_Short()を利用し 2016年八イパフォーマンスコンピューティングと計算科学シンポジウム High Performance Computing Symposium 2016



図 9 袖領域通信のレイテンシ

て転送完了を通知することで Ping-Pong 通信を行う.そし て, MPIでは MPI\_Type\_vector()を利用して, MPI\_Send() と MPLRecv() を行うことで Ping-Poing 通信を行う. ま た,配列の各要素は倍精度浮動小数点数(8bytes)とする. 図9にレイテンシを,図10バンド幅を示す.図より, GASNet/TCA のブロックストライド転送は, ブロック 転送とほぼ等しい性能が得られており,全体的に MPIよ りも大幅に高い性能が得られていることがわかる. GAS-Net/IBと比べた場合は、MPI ほどの性能差はないものの、 N = 4~64の範囲では GASNet/TCA の方が 10%以上小 さいレイテンシで転送できており, N = 32 では 30 %以上 小さいレイテンシで転送できている.これらの結果より, 小領域サイズにおける TCA/PEACH2 におけるブロック ストライド転送のハードウェア実装の有効性がわかる.ま た, N = 320 以上の大領域サイズにおいては, 前節での 測定同様に GASNet/IB ではバンド幅が大幅に低下してし まっているが、GASNet/TCA ではバンド幅の低下は見ら れず、本研究における通信の最適化がブロックストライド 転送においても有効に機能していることがわかる.一方 で、ストライド転送に関しては、GASNet/TCA の場合で N = 16以下, GASNet/IBの場合でN = 8以下の小領域サ イズでは MPI よりも低いレイテンシで転送できているが, これ以上のサイズでは MPI よりも低い性能となっている. これは、GASNet/TCAやGASNet/IBではPack/Unpack を行わずに細粒度の転送を行っているのに対して, MPI で は非連続データを Pack/Unpack して連続データとして転 送しているため, レイテンシは大きくなってしまうものの, ストライド転送でもブロックストライド転送と同程度の性 能が得られていると考えられる.

## 5.3 姫野ベンチマークの性能

本節では, 姫野ベンチマークを用いて TCA/PEACH2 に よる GPU 対応 GASNet の性能評価を行う.

姫野ベンチマークは非圧縮性流体解析コードの性能評価 のために開発されたもので、3次元ポアソン方程式をヤコ ビ反復法で解く場合の処理速度を測定するプログラムで ある.本稿では、表2に示す問題サイズを使用し、ヤコビ



HPCS2016

2016/6/7

法の反復回数を 1000 回に固定した際の処理時間によって 性能を評価する. GASNet による姫野ベンチマークは,オ リジナルの MPI 版 [13] を CUDA に移植したもの [14] に 対して NVIDIA Kepler アーキテクチャ向けの最適化を施 したもの [15] をベースとして実装している. 姫野ベンチ マークにおける通信は,袖領域の交換と収束判定のための Allreduce がある. 袖領域の交換には GASNet のストライ ド転送機能を使用している. 収束判定に関しては, GPU メ モリに存在する誤差の値を gdrcopy を使用して, CPU メモ リにコピーしてた後, Allreduce を行っている. Allreduce は AM\_Medium() を使用して Dissemination 法 [16] によっ て実装している.

表 2 姫野ベンき	チマークにおける問題	サイズ $(i  imes j  imes k)$
Small	Middle	Large
$64 \times 64 \times 128$	$128 \times 128 \times 256$	$256\times256\times512$

測定結果を図 11, 12, 13 にそれぞれ示す. 図中の凡例 における "Calc." は計算にかかった時間, "Halo exch." は 袖領域の交換にかかった時間, "Convergence" は収束判 定のための Allreduce にかかった時間をそれぞれ示してい る.図より,i方向で分割を行った場合に関しては、全て の問題サイズと分割パターンにおいて, GASNet/TCA, GASNet/IB, MPI でほぼ等しい性能となっていることが わかる.一方で,j方向で分割を行った場合に関しては, GASNet/TCA と GASNet/IB は MPI と比べて高い性能 が得られており、GASNet のストライド転送機能の有効性 がわかる.しかし、GASNet/TCAとGASNet/IBの間で の性能差はほぼなく、前節での測定結果と一致していない が、これは現在生じている PEACH2 の不具合への対処に 起因するオーバヘッドが原因である.具体的には、DMAC における転送完了通知の処理に不具合があり、DMA 通信 が完了して DMAC が完全に停止したことを正しく検出で きないという問題がある.現在は DMAC の制御レジスタ にある Peformance Counter の値をポーリングすることで DMACの停止を検出している. そのため, 通常の完了通知



図 11 姫野ベンチマークの測定結果 (Small)



図 12 姫野ベンチマークの測定結果 (Middle)



図 13 姫野ベンチマークの測定結果 (Large)

を利用する方法よりも余分な時間がかかっており、姫野ベ ンチマークの性能が低下してしまっている. PEACH2 は FPGA による実装を行っているため、今後の改良でこの点 を修正し、性能を改善する余地はあると考えている.

## 6. おわりに

本稿では、低レベル通信ライブラリである GASNet の GPU 対応拡張を TCA/PEACH2 を用いて実装し、基本的 な通信性能や袖領域通信、姫野ベンチマークの性能に関し ての評価を行った.その結果、GPU 間通信の最小レイテン シに関しては GASNet/IB や MPI よりも大きいが、比較的 小さいオーバヘッドで通信を行うことができた.最大バン ド幅に関しても、MPI よりは低いものの、ソフトウェア支 援によって PEACH2 の DMA 通信の 1.2 倍のバンド幅が 得られた.袖領域通信に関しては、ブロックストライド転 送において、PEACH2 のハードウェア機能を利用すること で、GASNet/IB や MPI と比べて高い性能が得られた.姫 野ベンチマークに関しては、MPI よりは高い性能が得られ たが、ハードウェア不具合への対処が原因で GASNet/IB と同程度の性能しか得られなかった. 今後の課題としては、ストライド転送の高速化が必要だ と考えられる.ストライド転送は PEACH2 のブロックス トライド転送でも実行は可能だが、高い性能は得らない ため、MPI 同様に Pack/Unpack による転送に対応する必 要があると考えられる.また、PEACH2 には DMAC が4 チャネル実装されているが、現在1チャネルしか利用しお らず、利用されていない DMAC を活用することで、バン ド幅の向上などが期待できる.

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.また, 本研究における HA-PACS/TCA の利用は,筑波大学計算 科学研究センター学際共同利用プログラム・平成 27 年度 課題「密結合演算加速機構アーキテクチャに向けたアプリ ケーションの開発と性能評価」による.

#### 参考文献

- [1] Top500 Supercomputer Sites, http://www.top500. org/.
- [2] 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 (ACS), Vol.6, No.3, pp.14-25, 2013.
- [3] GASNet Communication System, http://gasnet.lbl. gov/.
- [4] Unified Parallel C, http://upc.gwu.edu/.
- [5] Co-array Fortran, http://www.co-array.org/.
- [6] XcalableMP, http://www.xcalablemp.org/.
- [7] NVIDIA GPUDirect, https://developer.nvidia. com/gpudirect.
- [8] HA-PACS プロジェクト, http://www.ccs.tsukuba. ac.jp/research/research\_promotion/project/ ha-pacs.
- [9] GASNet-EX Collaboration, https://sites.google. com/a/lbl.gov/gasnet-ex-collaboration/
- [10] CUDA Toolkit, https://developer.nvidia.com/ cuda-toolkit.
- [11] NVIDIA gdrcopy, https://github.com/NVIDIA/ gdrcopy.
- [12] CUDA C Programming Guide, 3.4.2. Page-Locked Host Memory, http://docs.nvidia.com/ cuda/cuda-c-programming-guide/index.html# page-locked-host-memory
- [13] 姫野ベンチマーク, http://accc.riken.jp/2145.htm.
- [14] E. Phillips and M.fatica. : Implementing the Himeno benchmark with CUDA on GPU clusters, Parallel & Distributed Processing (IPDPS), 2010 IEEE International Symposium, pp.1-10, Apr. 2010.
- [15] Toshihiro Hanawa, Hisafumi Fujii, Norihisa Fujita, Tetsuya Odajima, Kazuya Matsumoto, Yuetsu Kodama, Taisuke Boku : Improving Strong-Scaling on GPU Cluster Based on Tightly Coupled Accelerators Architecture, IEEE Cluster 2015, Sep. 2015.
- [16] 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 による GPU 間直接通信にお ける Collective 通信の実装と性能評価, 情報処理学会論 文誌コンピューティングシステム (ACS), Vol.8, No.4, pp.36-49, 2015.