A Study of Performances Considering Communications on a GPU Cluster

Takahito Kawakami and Shinichi Yamagiwa School of Information Kochi University of Technology/JST PRESTO Kami, Kochi 782-8502 Japan

Abstract

Multiple processing nodes with GPUs connected by a high speed network organize a GPU cluster environment. It is expected to achieve very higher performance than the one of CPU-based cluster due to the drastic performance growth of the GPU. This paper focuses on the performance impact when the tasks are distributed by different assignments using MPI framework. Due to the difference between the local communication in a processing node and the remote one among the nodes, the performance is very likely to be changed Applying famous benchmarks available in the internet, this paper evaluates the performance impacts of different task assignments among the nodes on the 32-node Tesla-based GPU cluster.

1. はじめに

近年、PC(パーソナルコンピュータ)の低価格化・ 高性能化が顕著であり、PCを複数台ネットワークで 接続し、並列計算による高速化を図るクラスタ[1]が 用いられるようになっている。特に、本来グラフィッ クの計算を担当する GPU(Graphics Processing Unit) によって、CPU(Central Processing Unit) が行う計算 を肩代わりする GPGPU[2]を応用した GPU クラス タの利用が進んでいる。

一方で、GPU クラスタでは従来からの CPU クラ スタ環境と同様に、複数の GPU にタスクを分割し た場合、そのタスク間で通信が必要である。従って、 GPU クラスタにも通信のボトルネックが存在してい る[3][4]。これは、ノード内部の CPU や GPU のデー タ転送の速度に比べて、GPU クラスタのノード間を 接続するネットワークの速度が圧倒的に遅いことに 起因する。そのため、ノード間の通信システムが全 体の計算能力を左右する要因となってしまい、CPU や GPU に潜在する計算能力を引き出すことが難し い。そこで、ノード間の通信性能が計算能力に与え る影響を一般的なベンチマークソフトによる性能評 価を行うことで、通信性能とCPU・GPUの持つ潜在 能力の関係を明らかにする指標となる可能性がある。 本論文では、GPUクラスタにおいて、フリーで入 手可能なベンチマークソフトを使い、性能を評価す ることを目的とする。ベンチマークソフトとして、

姫野ベンチマーク・NPB(NAS Parallel Benchmarks)・ NAMDを使用し、計算の並列性の度合いや計算を並 列化する際のタスクの割り当て方を変化させ、ノー ド間の通信速度が与える影響を調査する。

2. GPU クラスタ

2.1. 構成

本論文で対象とする GPU クラスタを Figure 1 に 示す。このノードの構成を Figure 2 に示す。

GPU クラスタは全 32 ノードからなっており、1つ のノードには Intel 社製の CPU である Intel XeonR Processor E5645 2.4GHz と、NVIDIA 社製の GPUで ある Tesla M2050[5] が 2 台ずつ搭載されている。そ の他、メモリとして 12GB の DDR3、記憶装置とし て 80GB の SSD を有している。ノード間の通信に は、Infiniband QDR 40Gbps と Gbit Ethernet の 2 種 類を使用することができる。GPU クラスタを利用 した計算を行う場合は、Gbit Ethernet で接続された Headnode から各ノードに対してタスクをディスパッ チすることによって実行する。

2.2. GPU クラスタで使われるソフトウェア

2.2.1. MPI での並列化プログラム. 複数のノードから なる GPU クラスタにおいて並列計算を行う場合、複 数のノードで実行されるプロセス間においてネット ワークを利用した通信を行う必要がある。このプロセ ス間の通信に利用されるのが MPI(Message-Passing Interface)[6][7] である。

MPIとは、プロセス間のメッセージの送受信を定 義したものである。Message-Passingとあることから わかるように、Message-Passingモデルを採用してお り、送信側・受信側の両方のプロセスを操作すること



Figure 1. 本論文が使用する GPU クラスタ



Figure 2. GPU クラスタの構成

によって通信を行う。MPIそのものは、言語ではなく 仕様を定めたものとなっており、C・C++・Fortran-77・ Fortran-95 などの言語によって実装されている。

MPIを利用することで、プロセス間の通信に必要 なルーチンを手に入れることが可能になり、複雑な 通信操作に必要とされる手間を軽減できる。

この MPI を利用した並列計算を行う際の、ノード に対するタスクの割り当て方は machinefile で定義さ れる。machinefile の記入例を Figure 3 で示す。

通常は Figure 3 の (A) のように、使用するノード 名を1 つずつ記述する。これにより、各ノードのプ ロセッサコアに対して1 個ずつタスクがラウンドロ ビン (RR) によって割り当てられる。この方法を、1 コアずつのノード間 RR と呼ぶ。

また、(B)のように同じノード名を連続して記述 することで、同じノードのプロセッサコアに対する 連続したタスクの割り当てが RR によって行われる。



Figure 3. machinefile の記述例

この方法を、n コアずつのノード間 RR と呼ぶ。例 えばノード名を6回連続で記述すると、各ノードの プロセッサコアに対して、連続して6個のタスクが 割り当てられる6コアずつのノード間 RR となる。

今回使用する CPU には、6 つのプロセッサコアが 存在する。そのため、6 コアずつのノード間 RR を 行うことで、連続した 6 個のタスクが同じプロセッ サ内に割り当てられ、それらのタスクの通信は CPU 内で行われる。これをプロセッサ内通信という。さ らに、1 つのノードには 2 つの CPU が搭載されてい るため、12 コアずつのノード間 RR を行うことで、 12 個のタスク間の通信はノード内で行われる。これ をノード内通信という。プロセッサ内通信・ノード 内通信は、ノード間のネットワーク通信に比べて高 速である。

2.2.2. GPGPU(General-Purpose Computation on GPU). 本来、GPU や Video RAM(VRAM) を含むグ ラフィックアクセラレータは CPU によってコントロールされ、グラフィック処理を補助する目的で使用される [8]。この GPU を CPU が行うような汎用計算 に応用するのが GPGPU である。

GPUをGPGPUとして利用する場合、CPUはGPU による計算の実行を補助する目的で使用される。CPU はGPUの持つメモリに対してプログラムのダウン ロードを行う必要がある。また、GPUはプログラム の読み出しや結果の書き込みを VRAM に対して行 うので、CPUは VRAM とメインメモリ間でデータ のコピーを行うことで実行を補助する。

Figure 4 に示すような現在利用されている GPU は、 Stream Processor(SP) と呼ばれるプロセッサによって 構成されている。この SP が汎用計算用のプロセッサ として使用される。しかし、SP によってプログラム が実行されるため、プログラムをストリームベース



Figure 4. 近年の GPU の構成

に変更しなければならない。このストリームベース のプログラミングインタフェースは、NVIDIA 社に よって提案されている、GPGPU ためのランタイムで ある CUDA(Compute Unified Device Architecture)[9] によってサポートされる。CUDA 用のプログラミング ツールや API は、現在同社によって提供されている。

2.2.3. GPU プログラムの並列化. GPU によって実行 されるプログラムを並列化する場合、先に述べた MPI による並列化と GPGPU のためのランタイムである CUDA を利用して行う。CUDA によって GPU を汎 用計算に利用し、必要となるノード間のネットワー ク通信を MPI によって実現する。

この GPU プログラムの並列化を行う際に問題と なるのが、ノード間のネットワーク通信速度である。 GPU の処理速度に対してノード間の通信速度が圧倒 的に遅いため、GPU の処理能力を生かし切ることが できない可能性がある。つまり、通信速度に全体の 実行時間が支配されてしまうのである。この問題は、 CPU プログラムの並列化においても同様である。

2.3. 議論

先述したとおり GPU クラスタによる並列計算に おいて問題となるのが、ノード間のネットワーク通 信速度である。複数のノードによって実行される並 列計算において、避けることができない点であると 共に、計算速度を左右する大きな問題でもある。

この問題を解決する1つのアプローチとして、 machinefileによる並列計算プログラムの分配方法を 検討することが挙げられる。例えば、並列計算を行 うプログラム中で連続したタスク間の通信が多発す る場合、1コアずつのノード間RRを行うと、通信を 必要とするタスクは異なるノードに割り振られてし まう。そのため、低速なノード間通信が必要となり、 ノード間の通信速度に実行時間が左右される。逆に、 同じプロセッサ内・ノード内にタスクを割り振れば、

Table 1. 姫野ベンチマークの問題サイズ

S	$128 \times 64 \times 64$
Μ	$256\times128\times128$
L	$512 \times 256 \times 256$
XL	$1024 \times 512 \times 512$

高速なプロセッサ内通信・ノード内通信を利用する ことによる実行時間の短縮が見込まれる。

以上より、本論文では使用するプログラムやCPU・ GPUのプロセッサコア数を考慮したnコアずつのノー ド間RRを行うことで、高速なプロセス内通信・ノー ド内通信を促進し、低速なノード間ネットワーク通 信を減少させることで、インターネットなどで入手 可能なベンチマークを用いて、その性能に与えるイ ンパクトを調査し報告する。

3. ベンチマークソフト

本論文で対象とするベンチマークを以下に示す。

3.1. 姫野ベンチマーク

姫野ベンチマークとは、ポアッソン方程式解法をヤ コビの反復法で解く場合の主要ループによって計算 速度を測定するのものである[10]。実行時間が短いた め測定結果を直ちに求めることができる他、Table 1 に示すように複数の問題サイズが用意されている点 や、様々な環境下で測定された結果が公開されてい るなどの特徴がある。

3.2. NPB

NPB(NAS Parallel Benchmark) とは、NASA(National Aeronautics and Space Administration) によっ て開発されたベンチマークソフトである [11]。数値 流体力学 (CFD: Computational Fluid Dynamics) ア プリケーションから派生した、5 つのカーネルプロ グラムと3 つの疑似アプリケーションから構成され ている。問題のクラス (サイズ)を変更するオプショ ンが用意されており (A・B・C・D・S・W)、使用す る環境に合わせた性能の測定が可能である。

問題として、8つの中からCG・FT・IS・LUを使用 する。それぞれの詳しい説明がTable 2である[12]。

3.3. NAMD

NAMD[13] とは、大規模な生体分子システムのシ ミュレーションを行うアプリケーションである。イ リノイ大学によって開発されており、CUDA を利用 した GPU アクセラレーションが可能になっている。 本来はベンチマークとして開発された物では無いが、 CPU と GPU の計算能力の測定・比較を行うために 使用する。



Figure 5. 姫野ベンチマーク サイズ XL GCC

Table 2. 使用する NPB 問題の詳細

CG	大規模な正値対称な疎行列の最少固有値の近似値
	を、共役勾配法によって計算する
FT	FFT(Fast Fourie Transform) によって、3 次元偏微
	分方程式の解を計算する
IS	大規模な整数のソートを行う
LU	LU 分解を行うのではなく、5 × 5 の上下三角行
	列システムを SSOR(Symmetric Successive Over-
	Relaxation) 法で計算する

Table 3.3 種類の machinefile

1 コアずつの	ノード0から31まで1個ずつ順番
ノード間 RR	にタスクを割り当てる
6 コアずつの	ノード0から31まで6個ずつ順番
ノード間 RR	にタスクを割り当てる
12 コアずつの	ノード 0 から 31 まで 12 個ずつ順
ノード間 RR	番にタスクを割り当てる

4. 性能評価

姫野ベンチマーク・NPB・NAMD を使用した性能 評価を行う。

姫野ベンチマークは、[10] からソースコードをダ ウンロードし、GCC[14] と Intel Compiler[15] によ るコンパイルを行う。この際、問題サイズは XL を 選択する。そして、姫野ベンチマークを実行するに あたって、ノードへのタスクの割り当て方法として Table 3 に示す 3 種類の machinefile を利用する。

NPB は、8 つのベンチマークの中から CG・FT・ IS・LU を使用し、問題のクラスは C とする。ノー ドへのタスク割り当て方法として、姫野ベンチマー クと同様に、Table 3 に示す 3 種類の machinefile を 利用する。

NAMD は、実行環境に合わせた実行形式ファイル が用意されており、コンパイル環境に左右されない 普遍的な測定が可能になっている。今回は、CPUの みの実行を行うものと、CUDA による GPU での実 行をサポートしたものの2種類を使用する。さらに、 実行する問題サイズによる計算能力の差を測定する ため、[13] で公開されているサンプルの中から2種



Figure 6. 姫野ベンチマーク サイズ XL Intel

Table 4.2 種類の machinefile

1 コアずつの	ノード0から31まで1個ずつ順番
ノード間 RR	にタスクを割り当てる
2 コアずつの	ノード0から31まで2個ずつ順番
ノード間 RR	にタスクを割り当てる

類 (ATPase benchmark, STMV (virus) benchmark)を 使用する。

NAMD を GPU で実行するにあたって、CPU 同様 にタスクの割り当て方が計算能力に与える影響を調 べるために Table 4 に示す 2 種類の machinefile を使 用する。

4.1. 姫野ベンチマーク

Figure 5 と Figure 6 は、姫野ベンチマークの問題 サイズ XL を、GCC と Intel コンパイラによってコ ンパイルしたものを実行し、FLOPS 値を求めたもの である。

4.1.1. コンパイラによる比較. Figure 5 と Figure 6 に 関して、使用するコンパイラの違いという点で比較 すると、ピーク性能を発揮するプロセッサコア数に 違いがある。GCC の場合 256 個のプロセッサコアで 実行した際にピークとなっているのに対して、Intel コンパイラの場合 128 個のプロセッサコアで実行し た際にピークとなっている。さらに、両者がピーク として示す値には、差がほとんど存在しない。つま り、Intel コンパイラの方が少ないプロセッサコア数 で高い性能を引き出しており、プロセッサコア数 る一定を超えてしまうとコンパイラの最適化によっ て高速化されたタスクの処理時間以上に、ノード間 のネットワーク通信時間が長くなり、それがボトル ネックとなってプロセッサコア数が増えても性能が 下がる。

4.1.2. machinefile による比較. Figure 5 と Figure 6 を Table 4 に示した machinefile の違いという点で比



Figure 7. NPB CG クラス C

較すると、1 コアずつのノード間 RR によるタスク の割り当てを行った際の結果が優れている場合が多 い。これは、姫野ベンチマークがメモリバンド幅の 性能に左右されやすいベンチマークであることに起 因する。12 コアずつのノード間 RR のように1つの ノードに集中的にタスクを割り振ると、プロセッサ が持つプロセスあたりのメモリバンド幅が低下して しまう。そのため、各プロセスのメモリアクセス時 間の合計がノード間のネットワーク通信時間よりも 長くなり、複数のタスクが同じメモリを参照するこ とによるメモリバンド幅の圧迫によって全体性能が 低下する。

4.2. NPB

NPB に含まれるものの中から、クラス C の CG・ FT・IS・LU ベンチマークを選択して実行する。実行 結果として、実行時間と台数効果を求めている。

4.2.1. CG. Figure 7 は CG の実行結果である。CG ベンチマークに関しては、一部を除き1コアずつのノード間 RR による実行の際に結果が優れている。これは、選択的なプロセス間のリデュース操作が行われているためだと考えられる。同じノード内に連続してタスクを配置してしまうと、リデュース操作の送信者と受信者が、異なるノードに配置されている可能性が高くなる。そのため、外部のノードへの通信回数が増加してしまい、ネットワークでの通信時間が増加するため全体性能が低下する。

4.2.2. FT. Figure 8 は FT ベンチマークの実行結果で ある。FT ベンチマークに関しては、全体を通して 1 コアずつのノード間 RR による実行の際に結果が優 れている。特に、使用するプロセッサコア数が 16 個 のときは、最長で約 20 秒ほどの開きがある。これは、 FT ベンチマークで使用されている FFT おいて、バタ フライ通信を行っているためだと考えられる。バタ フライ通信を繰り返している。そのため、同じノード



Figure 8. NPB FT クラス C

内に連続してプロセスを割り当ててしまうと、ネットワーク通信に時間を必要とする。よって、1 コアず つのノード間 RR を使用した割り当ての方が、適度 にタスク間の間隔が取られ、良い結果となっている。

4.2.3. IS. Figure 9 は IS の実行結果である。IS ベン チマークに関しては、使用するプロセッサコア数に 対して問題サイズが小さいため、タスクの割り当て による実行時間の差は誤差と呼べる程度のものであ る。IS ベンチマーク中では、それぞれのプロセスの 計算結果を全体で送受信する全対全の通信を行って いる。そのため、タスクの割り当て方に関係なく通 信が発生してしまうため、実行結果に差が生まれに くい。

4.2.4. LU. Figure 10 は LU の実行結果である。LU ベンチマークに関しては CG ベンチマークと同様に、1 コアずつのノード間 RR での実行結果が優れており、16 個のプロセッサコアで実行した場合には約 20 秒ほどの差が確認される。ただし、こちらは全対全 通信を行っている分だけ、CG ベンチマークに比べる と実行方法の違いによる差が少ない。

4.3. NAMD

4.3.1. CPU と CPU+GPUの比較. Figure 11 と Figure 12 は、NAMD において ATPase benchmark(327,506 atoms)とNAMD STMV (virus) benchmark(1,066,628 atoms)を実行した結果である。それ ぞれの問題において CPU のみによる実行時間と、 CPU+GPU による実行時間を比較している。

ATPase benchmark においては、CPU による実行 と CPU+GPU による実行に大きな差は存在しない。 これは、問題サイズが大きくないため、両者ともに 問題に対する十分な計算能力を有しているためだと 考えられる。これに対して、NAMD STMV (virus) benchmark では両者の差は明確である。特に、計算 に使用するプロセッサコア数が少ない場合は、最高 で約 5.5 倍ほど CPU+GPU による実行の方が計算速



600

500

400

200

100

0

16

32



Figure 11. NAMD ATPase benchmark



Figure 13. NAMD ATPase benchmark GPU 比較

Figure 12. NAMD STMV (virus) benchmark

64 プロセッサコア数

CPU実行時間

CPU + GPU実行時間

CPU実行時間/CPU + GPU実行時間

128

6

5

4 効

1

0

256

率

倍



Figure 14. NAMD STMV (virus) benchmark GPU 比較

度が速いことがわかる。使用するプロセッサコア数 が増加するに伴って両者の差は減少していくが、こ れは問題サイズが小さい場合の理由と同様で、問題 サイズが CPU の計算能力に対して小さすぎてしま い、ネットワークの通信性能が全体性能に対して顕 著になるからである。

以上のことから、問題サイズが十分に大きいま たは、使用するプロセッサコア数が少ない場合は CPU+GPU による計算の実行に優位性があることが わかる。これは、GPU による計算の高速性を示すも のである。

4.3.2. GPU のタスク割り当てによる比較. Figure 13 と Figure 14 は、 Table 4 に示した 2 種類の machine file を使用して、CPU+GPU によって ATPase benchmark とNAMD STMV (virus) benchmark を実行し、実行時 間を比較している。ATPase benchmark・NAMD STMV (virus) benchmark 共に、実行時間に差は見られない という結果になった。これは、これらのベンチマー クを実行する際に、ノード間の通信を必要としてい ないことを示している。よって、ノード間のネット ワーク通信速度が計算能力に与える影響が少なく、 並列計算に適していることがわかる。

5. まとめ

本論文では、GPU クラスタにおいてインターネットで入手可能なベンチマークソフトによる並列計算 を行い、特にノード間のネットワーク通信に注目し、 複数の計算能力の測定を行った。

並列計算を行う場合は、ノードに対するタスクの 割り当て方が、計算速度に大きな影響を与えること を示すことができた。特にFT ベンチマークのよう なバタフライ通信を使用している場合、連続的なタ スクの割り当てが悪影響となり、計算速度が低下し てしまう。つまり、必ずしも同じノードに連続的に タスクを割り当てることが効果的とはいえず、使用 するプログラムのアクセスパターンを考慮した、タ スクの割り当て方を行うことが必要であることを示 した。

謝辞

本実験を行うにあたり、高知工科大学IACPのGPU クラスタを使用させていただいた。

References

- Rajkumar Buyya, "High Performance Cluster Computing: Architectures and Systems," *Prentice Hall*, 1999.
- [2] John D. Owens and David Luebke and Naga Govindaraju and Mark Harris, Jens Kruger and Aaron E. Lefohn and Timothy J. Purcell, "A Survey of General-Purpose Computation on Graphics Hardware," In Eurographics 2005, State of the Art Reports, pages 21-51, August 2005.
- [3] Vijay Karamcheti and Andrew A. Chien, "Software Overhead in Messaging Layers: Where Does the Time Go?" In ASPLOS-VI: Proceedings of the sixth international conference on Architectural conference on Architectural support for programming languages and operating systems, pages 51-60, 1994.
- [4] Shinichi Yamagiwa and Kevin Ferreira and Keiichi Aoki and Masaaki Ono and Koichi Wada and Leonel Sousa, "Maestro2: Experimental Evaluation of Communication Performance Improvement Techniques in the Link Layer," *Journal of Interconnection Networks* (*JOIN*), World Scientific Publishing, vol.7 no.2, pp. 295-318, June 2006.
- [5] David B. Kirk and Wen-mei W. Hwu, "Programming Massively Parallel Processors," MORGAN KAUF-MANN PUBLISHERS, 2010.
- [6] N. Doss William Gropp, E. Lusk and A. Skjellum, "A High-Performance, Portable Implementation of the MPI Message Passing Interface Standard," *In MPI Developers Conference*, 1995.

- [7] Message Passing Interface Forum, "MPI : A Message-Passing Interface Standard," September 4, 2009.
- [8] Pablo Lamilla Alvarez and Shinchi Yamagiwa and Masahiro Arai and Koichi Wada, "A Uniform Platform to Support Multigenerational GPUs for High Performance Stream-based Computing," *International Journal of Networking and Computing Vol 1, No 2*, 2011.
- [9] "NVIDIA Developer Zone," http://developer.nvidia. com/category/zone/cuda-zone.
- [10] "姫野ペンチマーク Homepage," http://accc.riken.jp/ HPC/HimenoBMT/.
- [11] "NAS Parallel Benchmarks Changes," http://www.nas. nasa.gov/Resources/Software/npb.html.
- [12] D. H. Bailey and E. Barszcz and J. T. Barton and D. S. Browning and R. L. Carter and L. Dagum and R. A. Fatoohi and P. O. Frederickson and T. A. Lasinski and R. S. Schreiber and H. D. Simon and V. Venkatakrishnan and S. K. Weeratunga, "THE NAS PARALLEL BENCHMARKS," *Intl. Journal of Supercomputer Applications, vol. 5, no. 3, pp.66-73,* Fall.1991.
- [13] "NAMD Scalable Molecular Dynamics," http://www. ks.uiuc.edu/Research/namd/.
- [14] "GCC, the GNU Compiler Collection," http://gcc.gnu. org/.
- [15] "Intel Compilers from Intel," http://software.intel.com/ en-us/articles/intel-compilers/.