

## PCI Express 拡張 Box と仮想 GPU との計算性能比較

瀬戸口 幸寿<sup>†</sup>成見 哲<sup>‡</sup>

電気通信大学 情報・通信工学専攻

## 1 はじめに

近年,GPU(Graphics Processing Unit)が科学シミュレーションなどの汎用計算に用いられている. この技術はGPGPU(General-Purpose computing on GPU)と呼ばれており, HPCの分野に大きく貢献している. GPUは本来, 画像処理に特化したプロセッサであるが, 近年のGPUは大規模なデータに対し単純な演算を並列に行うことが可能である.

我々は, ソフトウェアを書き換えることなくネットワークを介して複数のGPUを使用できるようにするGPU仮想化ツールであるDS-CUDA[1]を提案している. GPUの仮想化により, 場所の制限なしに複数のGPUを使用するプログラムを簡単に記述することが出来る. しかし, クライアントノード-GPUサーバ間の通信がボトルネックになる.

本研究では複数のGPUを用いた場合に, DS-CUDAによる仮想GPUとネイティブなGPUとの性能を比較し, 性能モデルの構築を行った.

## 2 システム

本研究では, 以下の手法で複数のGPUを用いた計算性能の比較を行った.

- DS-CUDA:Infiniband ネットワークによる仮想GPUの利用
- ネイティブなCUDA:PCI Express 拡張BoxによるGPUの直接利用

## 2.1 GPU

本研究では両手法にて用いるGPUとして, ELSA GeForce GTX 780[2]を最大8台まで用いた.

## 2.2 CUDA, DS-CUDA

CUDA(Compute Unified Device Architecture)[3]はNVIDIA社が開発したGPGPUの為の統合開発環境であり, 拡張言語, コンパイラ,C言語のライブラリから構成さ

れている.

DS-CUDA(Distributed-Shared CUDA)[1]はネットワークを介してリモートGPUを用いる為のミドルウェアである.DS-CUDAのシステムはDS-CUDAコンパイラ,clientノード,(複数の)serverノードから構成される. DS-CUDAコンパイラはclientノード上で, client-server間通信の為に命令を付加してコンパイルを実行する. これによってDS-CUDAは,CUDAプログラムをローカルなGPU(仮想GPU)を用いて実行しているように見える(図1左). また, プログラムはネットワークを介したGPUの利用を気にせずに, 同一のCUDAソースコードを利用できる. 本研究では,8台のGPUに対してInfinibandネットワークを介した8つのserverノードを構築し実験を行った.

## 2.3 PCI Express 拡張 box

DS-CUDAと比較の為に,clientノードから直接複数のGPUを繋ぐ仕組みが必要である. 本研究ではPCI Express 拡張box[4]を利用した.PCI Express 拡張boxはPCI-Expressバスに挿すホストカード, GPUコンテナbox, 両者を繋ぐケーブルから構成される(図1右). 本研究では4台格納可能な拡張boxを2台用いて, 最大8台のGPUを接続する.

## 2.4 Claret

Claretは古石によって作られた融解塩の分子動力学(MD)シミュレーションプログラムである[5]. 本研究ではClaretを用いて性能比較, モデル構築を行う. オリジナルコードに対し, 複数GPUへの対応を行い, OpenGLによる画面描画機能は単純化の為無効化した. 粒子数を $n$ とすると1回のMDstepに対してCPU-GPU間の通信量は $O(n)$ ,GPUでの計算量は $O(n^2)$ である.

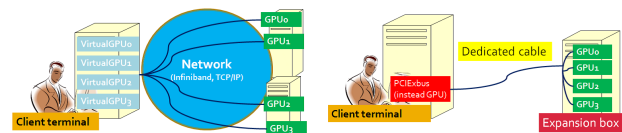


図1 (左)DS-CUDA使用時のGPUの接続

(右)PCI Express 拡張 Box 使用時のGPUの接続

Performance Comparison of GPGPU between PCI Expansion Box and Virtualization via DS-CUDA

<sup>†</sup> Setoguchi Yukitoshi

The University of Electro-Communications

<sup>‡</sup> Narumi Tetsu

The University of Electro-Communications

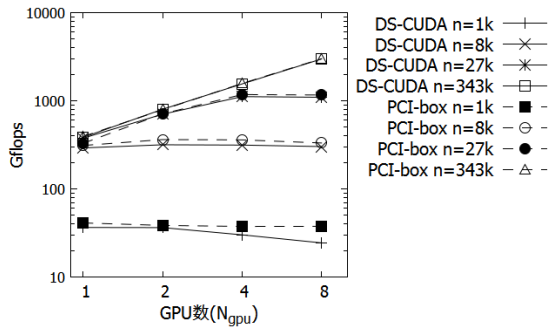


図2 GPU数に対する計算性能

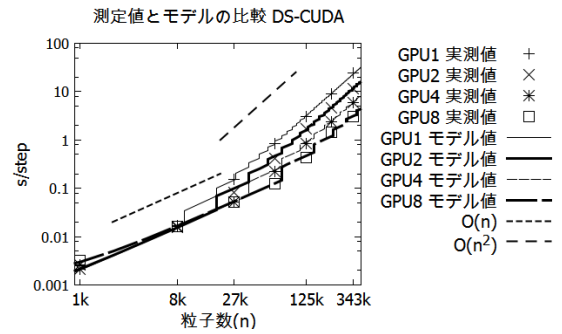


図3 1stepの計算時間の実測値とモデルの比較 (DS-CUDAの場合)

### 3 実験結果

粒子数, GPUの並列数, そしてGPUの利用手法を変えながらMDシミュレーションを実行した. 測定した1stepにおける計算時間を基に, GPUの並列数に対する計算性能(Gflops)を図2に示す. 粒子数が多い時にはGPU数に応じて性能が向上しているが, 粒子数が少ない時は逆に低下している.

### 4 モデルの構築

1stepの計算時間  $T$  をGPUで要する計算時間 ( $T_{gpu}$ ), CPUで要する計算時間 ( $T_{cpu}$ ), server-client間の通信時間 ( $T_{com}$ ) の和で以下のようにモデル化する.

$$T_{gpu} = \lceil \frac{n}{kcN_{gpu}} \rceil nkct_{gpu} \quad (1)$$

$$T_{cpu} = t_{cpu}n \quad (2)$$

$$T_{com} = \left( \frac{16nN_{gpu}}{B_{h2d}} + N_{gpu}L_{h2d} \right) + \left( \frac{12n}{B_{d2h}} + N_{gpu}L_{d2h} \right) + t_{kernel}N_{gpu} \quad (3)$$

- $n$  : 粒子数
- $N_{gpu}$  : GPU並列数
- $t$  : 測定で決定する各係数 (sec)
- $k$  : 測定で決定する整数 (今回は4)
- $c$  : GPUあたりのコア数 (今回は2304)
- $B$  : 転送スループット (byte/sec)
- $L$  : 転送レイテンシ (sec)

ここで,  $kcN_{gpu}$  は一度に実行出来るスレッド数,  $h2d$  は client(host) から server(device) へ,  $d2h$  は server から client への転送を表し,  $t_{kernel}$  はカーネル実行の遅延を表す. 測定の結果  $L, 1/B, t_{kernel}$  いずれも DS-CUDAの方が大きく, 小粒子, 並列大領域におけるPCI Express拡張boxとの差を決定づけていることが分かった. DS-CUDAのモデルと実測値の比較を図3に示す. 粒子数の少ない時は粒子数に比例して計算時間が上昇していくが, 粒子数が多い時は2乗で増えている. またモデルにおける

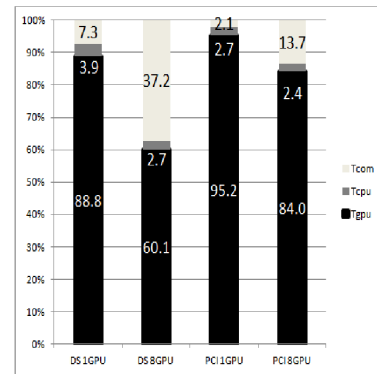


図4 1000粒子における  $T_{gpu}, T_{cpu}, T_{com}$  の比率

$T_{gpu}, T_{cpu}, T_{com}$  の比率を図4に示す. GPU数が増えると通信時間が増えていることが分かる.

### 5 おわりに

測定とモデル化の結果から, 両者共に小粒子数, GPU並列数大のときに, 通信時間の割合が大きくなることが確認出来た. 今後の展望として, ピーク時に対する効率に関するパラメータを導き, 最適な仮想GPUの運用の為の指標を提案したい.

### 参考文献

- [1] A. Kawai, K. Yasuoka, K. Yoshikawa, T. Narumi: "Distributed-Shared CUDA: Virtualization of Large-Scale GPU Systems for Programmability and Reliability", The Fourth International Conference on Future Computational Technologies and Applications, FUTURE COMPUTING 2012, Nice, France, 2012.
- [2] ELSA GeForce GTX 780: [http://www.elsa-jp.co.jp/products/products-top/graphicsboard/geforce/ultra\\_high\\_end/geforce\\_gtx780/](http://www.elsa-jp.co.jp/products/products-top/graphicsboard/geforce/ultra_high_end/geforce_gtx780/) (Last access:2015/05/13)
- [3] CUDA Zone: <https://developer.nvidia.com/cuda-zone> (Last access:2015/05/13)
- [4] ELSA VRIDGE X100 Dual16: [http://www.elsa-jp.co.jp/html/products/pes/vridge.x100\\_dual16/index.htm](http://www.elsa-jp.co.jp/html/products/pes/vridge.x100_dual16/index.htm) (Last access:2015/05/13)
- [5] Claret: <http://atlas.riken.go.jp/~koishi/claret.html> (Last access:2015/05/13)