

FMOプログラムABINIT-MPのGPU化と性能評価

GPU optimization and performance evaluation of FMO program ABINIT-MP

坂倉耕太¹⁾, 望月祐志²⁾³⁾, 中野達也⁴⁾, 大島聰史⁵⁾, 星野哲也⁶⁾, 片桐考洋⁷⁾

Kota Sakakura, Yuji Mochizuki, Tatsuya Nakano, Satoshi Oshima, Tetsuya Hoshino, Takahiro Katagiri

1)博(工) 大阪大学D3センター 特任准教授 (〒567-0047 大阪府茨木市美穂ヶ丘5-1)

2)理博 立教大学理学部 教授 (〒171-8501 東京都豊島区西池袋3-34-1)

3)理博 東京大学生産技術研究所 リサーチフェロー (〒153-8505 東京都目黒区駒場4-6-1)

4)博(理) 高度情報科学技術研究機構神戸センター (〒650-0047 神戸市中央区港島南町1-5-2)

5)博(工) 九州大学情報基盤研究開発センター 准教授 (〒819-0395 福岡市西区元岡744)

6)博(理) 名古屋大学情報基盤センター 准教授 (〒464-8601 名古屋市千種区不老町)

7)博(理) 名古屋大学情報基盤センター 教授 (〒464-8601 名古屋市千種区不老町)

We have been developing the ABINIT-MP program for fragment molecular orbital(FMO) calculations.

In this paper, We report on the acceleration and performance evaluation of ABINIT-MP on GPUs. The OpenACC directive was used for the two-electron integration of the high-cost part to speed up the process. By utilizing multiple processes with MPS, speedups of 1.5x for 8MPI+1GPU and 3.6x for 8MPI/4OpenMP+4GPU were achieved.

Key Words : Fragment molecular orbital, FMO, ABINIT-MP, GPU, Two electron repulsion integrals

1. はじめに

フラグメント分子軌道(FMO)法[1][2]は、タンパク質や核酸などの分子系をフラグメントに分割し、フラグメントとフラグメントペアに対する小規模な電子状態計算することで、巨大分子全体の電子状態を近似的に高速に解くことが可能な手法である。計算結果は、対象系の詳細な相互作用解析に適しており、理論創薬や材料化学の分野で広く利用されている。これまで、ABINIT-MPは、X86系計算機をはじめ、「富岳」や「不老」といった富士通A64FXや、ベクトル機であるNEC SX-Aurora TSUBASA環境など様々なアーキテクチャに対応してきた。昨今のHPC計算機環境は、CPU/GPU混成マシンが主流になりつつあり、スーパーコンピュータの性能ランキングTOP500上位の機種では、ほとんどがGPUを搭載したシステムによって占められている。今後大規模計算機資源を活用するためにはプログラムのGPU化対応は不可欠である。我々は、GPUの搭載が想定される富岳NEXT[3]を頂点とするHPCI計算機資源を活用したFMO計算のため環境整備を進めており、GPUによる高速化対応に注力している。

2. GPUによる高速化

ABINIT-MPはMPI/OpenMP並列計算に対応しており、ノード内のCPUコアをフルに使用し、高い並列性能での計算が可能である。GPUを付加することにより、さらなる性

能向上を目指すためには、MPIプロセス数、スレッド並列数、およびノードあたりのGPU搭載数の3つのパラメータのバランスを考慮したプログラム設計が必要である。ABINIT-MPでは、モノマー、ダイマーペア単位でのMPIプロセス並列が根本的な分割手法になるため、複数プロセスから1つのGPUを効率的に並列実行できるMPS (Multi-Process Service) の活用が望ましい。

ABINIT-MPの計算時間ホットスポットは、FOCK行列生成時に計算されるObara[4]のアルゴリズムで求められる3中心、4中心2電子積分計算（クーロン項、交換項、環境静電ポテンシャル）と、得られた積分値を、原子軌道番号を

```
do scf-iteration
!$acc data copyin(dc) copy(fc)
Subroutine eri_gpu_ssss
Subroutine eri_gpu_psss
Subroutine eri_gpu_psps
Subroutine eri_gpu_ppss
Subroutine eri_gpu_ppps
Subroutine eri_gpu_pppp
!$acc parallel num_workers(2) vector_length(16) async
!$acc loop gang worker private(a,c)
do n4_pppp = 1, nsize_pppp
  IJ_shell=list()
  KL_shell=list()
  do pq
    do rs
      calc ERI_pppp
```

図1 GPU化対象処理のOpenACC指示子修正イメージ

インデックス付きでFOCK行列に格納する加算処理であ

る。計算機アーキテクチャ、データ、近似手法によって差はあるものの、全体の計算時間の6割以上を占める。これらの処理部分に対して、2電子積分ループ構造の改変、OpenACC指示子によるGPU化、データ転送の最小化施策を適用し、1MPIプロセス比では数十倍のGPU高速化が可能となった[5][6]。図1のように、2電子積分計算は軌道角運動量タイプ(ss|ss)～(pp|pp)のサブルーチンによって構成される。これらの処理は相互に依存関係がなく、独立した処理が可能である。そこで、先述のMPSを活用し、各プロセスから同一GPUを効率的に処理できるよう、OpenACC指示子を用いて非同期計算ができるよう改修した。

3. 性能評価

(1) 性能測定環境、条件

NVIDIA社のA100搭載の東京大学情報基盤センターWisteria Aquarius1ノードを用いて測定した。測定データとしてTrpCage(残基数20)HF/6-31Gを採用した。TrpCageはデータサイズとしては小さいが、個々のフラグメントは大小様々なサイズを持ち、より実践的なデータであると言える。測定条件は、GPU高速化の比較としてCPUのみのMPI/OpenMP並列と、MPI/OpenMP+GPUを評価するため、(a) 10MPI/7OpenMP、(b) 1MPI+1GPU、(c) 8MPI+1GPU、(d) 8MPI/4OpenMP+4GPUの4条件で測定した。

(2) 測定結果、考察

上記4条件で測定した結果を表1に示す。two-electron Integrals(4center,3center)は、GPU化された2電子積分箇所であり、set_indexはGPU化に必要な前処理部分である。その他項目はGPU化しておらず、MPI/OpenMP並列で処理される。測定環境のMPSでは1GPUあたり最大CPU8コアが割り当て可能なため、+1GPU時は8MPI並列、+4GPU時は8MPI/4OpenMP並列とした。

表1 各並列条件における測定結果

MPI OpenMP	10MPI 7OpenMP	1MPI +1GPU	8MPI +1GPU	8MPI +4GPU
	(sec)	(sec)	(sec)	(sec)
two-electron Integrals(4center)	91.9	74.8	50.5	15.6
two-electron Integrals(3center)	7.5	6.8	5.0	4.3
dimer_es	2.6	104.5	12.8	4.7
dimer_oneint	1.1	58.3	7.2	1.9
monomer_esp	1.4	37.3	4.6	1.9
set_index	-	18.9	2.4	2.5
others	28.9	30.4	7.9	5.8
Total	133.3	330.9	90.5	36.8

まず、GPUを利用しない場合、計算時間は133.3秒であった。2電子積分部分が計算コストの大部分を占めた。次に、1MPI+1GPUでは、GPU化により、2電子積分部分は、2割程度高速化されているが、非GPU化部分に関しては、

並列処理されないため、全体の計算時間では2.5倍程度時間を要した。特にdimer_es部分が高コストになっており、次のGPU化ターゲットとして対応は必須である。MPSを利用した8MPI+1GPUでは、74.8秒から50.5秒と、さらなる2電子積分部分の高速化が確認できた。軌道角運動量タイプの積分処理においてMPIプロセス間で同期処理によって、効率よくGPUコアを利用できているものと考える。全体処理時間では90.5秒と当初の1.5倍の高速化となった。最後に、8MPI/4OpenMP+4MPIでは、複数GPUによる加速と、MPSによる効率的な処理により、さらに高速化され36.8秒となり、当初の3.6倍の性能向上であった。

4. まとめと今後の予定

本稿では、ABINIT-MPのGPU高速化への取り組みについて報告した。最計算コスト部分である3中心、4中心2電子積分計算部分のGPUを用いた高速化施策において、MPS活用が有効である結果が得られた。今後は、1) MPS、MIGの追加検証、2) dimer_es、1電子積分計算等の高コスト処理部分のGPU化検討、3) ポストHF法(MP2)のGPU化、の順で、高コスト部分はGPU、それ以外の処理は従来のMPI/OpenMPにて高速化整備を進めていく予定である。

謝辞

ABINIT-MPのGPU化対応は東京大学情報基盤センターのGPU移行推進プログラムの中で、エヌビディア合同会社の古家真之介氏のコーディネートで成されたものです。また、同センターの中島研吾先生、下川辺隆史先生、芝隼人先生(当時)のサポートにも感謝します。最後に、本研究開発は立教SFRの資金援助を受けていますことを記します。

参考文献

- [1] Fedorov, K, Kitaura ed., "The Fragment Molecular Orbital Method: Practical Applications to Large Molecular Systems", 2009, CRC.
- [2] Y. Mochizuki, S. Tanaka, K. Fukuzawa, ed., "Recent Advances of the Fragment Molecular Orbital Method", 2021, Springer.
- [3] <<https://cs-forum.github.io/roadmap-2023/>>
- [4] S. Obara, A. Saika, *J. Chem. Phys.*, 84, 3963 (1986)
DOI:10.1063/1.450106
- [5] Y. Mochizuki, T. Nakano, K. Sakakura, K. Okuwaki, et al., *J. Comput. Chem. Jpn.*, 23, 4(2024)
DOI: 10.2477/jccj.2024-0001
- [6] 坂倉耕太、望月祐志、中野達也、成瀬彰、大島聰史、星野哲也、片桐考洋、"FMOプログラムABINIT-MPのGPU化への対応"、計算工学講演会論文集、Vol.29,E-13-02,2024