

GPU を用いた McMurchie-Davidson アルゴリズムの高速化 A GPU Implementation of McMurchie-Davidson Algorithm

藤井 晴斗¹⁾ 伊藤 靖朗¹⁾ 横川 宜矢¹⁾ 鈴木 寛太¹⁾ 辻 聡樹¹⁾²⁾ 中野 浩嗣¹⁾ 笠置 明彦²⁾
Haruto Fujii Yasuaki Ito Nobuya Yokogawa Kanta Suzuki Satoki Tsuji Koji Nakano Akihiko Kasagi

1 はじめに

量子化学計算とは、分子のエネルギー等の特性を実際に行うことなく計算のみで導き出すことを目的とした、計算化学の一分野である。その中でも、二電子反発積分 (ERI) と呼ばれる計算のコストが特に高く、分子を M 個の基底関数で表す場合は $O(M^4)$ の計算量が必要となる。また、ERI の計算手法に関しては、McMurchie-Davidson 法 [1] や Obara-Saika 法 [2] など、様々な特性を持つアルゴリズムが提案されている。本研究では、McMurchie-Davidson 法を用いた計算に着目し、高速化を行った。

2 McMurchie-Davidson 法

本項では、ERI の計算手法の一つである McMurchie-Davidson (MD) 法について説明する。まず、ガウス型軌道の線形結合として表される基底関数を以下のように定義する。

$$\chi_\mu(r) = \sum_{e=1}^{K_\mu} d_{\mu e} G_{\mu e}(r) \quad (1)$$

ここで、 K_μ と $d_{\mu e}$ はそれぞれ線形結合の項数と係数であり、ガウス型軌道は以下のように定義される。

$$G_{\mu e}(r) = c_e (x-A_x)^{\mu_x} (y-A_y)^{\mu_y} (z-A_z)^{\mu_z} \exp(-\alpha_e |r-A|^2)$$

ここで、 $A = (A_x, A_y, A_z)$ はガウス軌道の中心、 c_e は正規化係数、指数部 μ_x, μ_y, μ_z は軌道の形状を決める非負整数、 α_e は軌道の大きさを決定する指数である。ERI は、4 つの基底関数 $\chi_\mu, \chi_\nu, \chi_\lambda, \chi_\sigma$ に対して以下のように定義される計算である。

$$(\mu\nu|\lambda\sigma) = \sum_{e=1}^{K_\mu} \sum_{f=1}^{K_\nu} \sum_{g=1}^{K_\lambda} \sum_{h=1}^{K_\sigma} d_{\mu e} d_{\nu f} d_{\lambda g} d_{\sigma h} [\mu_e \nu_f | \lambda_g \sigma_h] \quad (2)$$

MD 法では以下の式 (3)-(5) を用いて ERI を評価する。

$$[ab|cd] = \frac{2\pi^{5/2}}{pp'\sqrt{p+p'}} \times \sum_{i=0}^{a_x+b_x} \sum_{u=0}^{a_y+b_y} \sum_{v=0}^{a_z+b_z} \sum_{i'=0}^{c_x+d_x} \sum_{u'=0}^{c_y+d_y} \sum_{v'=0}^{c_z+d_z} E_i^{a_x, b_x} E_u^{a_y, b_y} E_v^{a_z, b_z} E_{i'}^{c_x, d_x} E_{u'}^{c_y, d_y} E_{v'}^{c_z, d_z} R_{i+i', u+u', v+v'}^0 \quad (3)$$

ここで、 $p = \alpha + \beta, p' = \gamma + \delta$ である。 $E_i^{a_x, b_x}$ の値は以下の漸化式で計算可能である。

$$\begin{cases} E_t^{i,j} = \frac{1}{2p} E_{t-1}^{i-1,j} - (A_x - B_x) \frac{q}{\alpha} E_t^{i-1,j} + (t+1) E_{t+1}^{i-1,j} \\ E_t^{i,j} = \frac{1}{2p} E_{t-1}^{i,j-1} + (A_x - B_x) \frac{q}{\beta} E_t^{i,j-1} + (t+1) E_{t+1}^{i,j-1} \\ E_t^{0,0} = \exp(-q(A_x - B_x)^2) \\ E_t^{i,j} = 0 \text{ if } t < 0 \text{ or } i+j < t \end{cases} \quad (4)$$

1) 広島大学大学院 先進理工系科学研究科

2) 富士通株式会社 コンピューティング研究所

ここで、 $q = \frac{\alpha\beta}{\alpha+\beta}$ であり、 $E_u^{a_y, b_y} E_v^{a_z, b_z}$ の計算にはそれぞれ y, z 座標の値を使い、 $E_{i'}^{c_x, d_x} E_{u'}^{c_y, d_y} E_{v'}^{c_z, d_z}$ の計算には G_c, G_d の 2 つのガウス型軌道を用いる。次に、 $R_{i+i', u+u', v+v'}^0$ の値は以下の漸化式で計算可能である。

$$\begin{cases} R_{t,u,v}^n = (t-1)R_{t-2,u,v}^{n+1} + (P_x - P'_x)R_{t-1,u,v}^{n+1} \\ R_{t,u,v}^n = (u-1)R_{t,u-2,v}^{n+1} + (P_y - P'_y)R_{t,u-1,v}^{n+1} \\ R_{t,u,v}^n = (v-1)R_{t,u,v-2}^{n+1} + (P_z - P'_z)R_{t,u,v-1}^{n+1} \\ R_{0,0,0}^n = (-2\omega)^n F_n(\omega|P - P'|^2) \\ R_{t,u,v}^n = 0 \text{ if } t < 0 \text{ or } u < 0 \text{ or } v < 0 \end{cases} \quad (5)$$

ここで、 $P = (P_x, P_y, P_z) = \frac{\alpha A + \beta B}{\alpha + \beta}, P' = (P'_x, P'_y, P'_z) = \frac{\gamma C + \delta D}{\gamma + \delta}, \omega = \frac{pp'}{p+p'}$ であり、 $R_{0,0,0}^n$ の値は Boys 関数

$$F_n(x) = \int_0^1 t^{2n} e^{-xt^2} dt$$

で計算される。式 (5) の計算の中では、Boys 関数の計算コストが支配的である。提案手法では、Boys 関数の評価に先行研究 [3] の手法を利用した。

3 提案アルゴリズム

本項では、MD 法を用いた ERI 計算に関する並列計算アルゴリズムを提案する。まず、 $s_x = a_x + b_x + c_x + d_x, s_y = a_y + b_y + c_y + d_y, s_z = a_z + b_z + c_z + d_z, K = s_x + s_y + s_z$ を定義する。

ここから、式 (3) の評価に必要な漸化式 R の値を全て事前計算する手法を説明する。式 (3) の評価には、 $0 \leq t \leq s_x, 0 \leq u \leq s_y, 0 \leq v \leq s_z$ を満たす $R_{t,u,v}^0$ の組が必要である。それらの各要素の値を得るために式 (5) を再帰的に計算すると、同じ項が多く現れ、冗長な計算が行われる。そこで、本手法では冗長な計算を回避することで無駄な計算を省略する。具体的には、再帰式 (5) に登場する各要素を、互いに依存関係がなく同時に計算可能な $K+1$ 個の計算単位 (以下、batch と呼ぶ) 毎に分割し、それらの値を順に評価することで再帰式 R を効率よく評価する。各 batch $k (0 \leq k \leq K)$ には、それぞれ

$$k = t + u + v, 0 \leq n \leq K - k \text{ を満たす全ての } R_{t,u,v}^n \text{ の組}$$

が含まれるとする。batch k に含まれる R の値の個数は $\frac{(k+1)(k+2)(K-k+1)}{2}$ 個、 R の値の総数は $\frac{(K+1)(K+2)(K+3)(K+4)}{24}$ 個となる。

ここで、式 (5) から、batch k に含まれるすべての値は batch $k-1, k-2$ のみに依存するため、batch k の各要素はそれぞれ独立に計算が可能である。図 1 に、 $K=3$ についての計算の流れを図示する。batch 0 の値は式 (5) の値を用いて計算され、batch k の値は batch $k-1, k-2$ の値を用いて計算可能であることが示されている。

ここからは、提案アルゴリズムを疑似コードとして説明する。

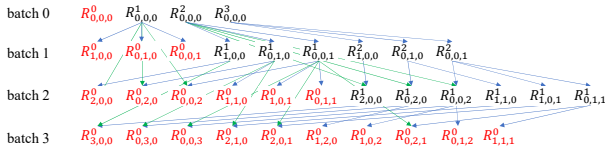


図 1 $K = 3$ の場合の例, 式 (3) に必要な値を赤で強調

4つのガウス型軌道についてのERIを計算する手法をAlgorithm 1に示す. このアルゴリズムの1-3行目では, 提案した事前計算を行う. その後, 得られた R の値を用いて, 5-8行目では式(3)に相当する計算を行う. $E_t^{a_x, b_x}$ のように, 式(3)内には E の計算も存在しているが, これは R と比較して計算量が小さいため, 式(4)の定義通りに再帰計算を行う.

Algorithm 1 ガウス型軌道についての並列 ERI 計算

Input: Four Gauss-type orbitals G_a, G_b, G_c, G_d

Output: ERI $[ab|cd]$

```

1: for  $k = 0$  to  $K$  do
2:   Compute  $R$  values in batch  $k$  in parallel
3: end for
4:  $s \leftarrow 0$  //  $s$  is a shared variable, and parallel addition is performed on it
5: for all  $t, u, v, t', u', v'$  do in parallel
6:    $s \leftarrow s + E_t^{a_x, b_x} E_u^{a_y, b_y} E_v^{a_z, b_z} E_{t'}^{c_x, d_x} E_{u'}^{c_y, d_y} E_{v'}^{c_z, d_z}$ 
7:    $(-1)^{t'+u'+v'} R_{t+u'+v'}^0$ 
8: end for
9: return  $\frac{2\pi^{5/2}}{\rho P \sqrt{\rho^{2+P}}} \cdot s$ 

```

4 提案 GPU 実装

本項では, 3章で提案したアルゴリズムを, NVIDIA社が提供するCUDA[4]を用いてGPUアーキテクチャ上に実装する手法を提案する. 本実装では, 1つのCUDA Blockに1組の基底関数に関するERI計算(式(2)に相当)を割り当てる. Algorithm 1については, Block内の各Threadを用いた並列計算を行う. 本実装では, 3つのShared Memory(Block内の各Threadが読み書き可能である高速なメモリ)の配列を繰り返し使用し, batch 0から順に計算を行うことで R の各要素を効率よく評価することが可能となっている. 3つの配列をそれぞれ sm_0, sm_1, sm_2 とすると, batch k の各要素を $sm_{(k \bmod 3)}$ に書き込む. これにより, 各batchの計算に必要な値を保持しつつ, 不要となった値を上書きすることでメモリ量を削減することが可能となる. また, 各batch計算の過程では, ERIの評価に必要な $R_{t,u,v}^0$ の各要素を保管するために別のShared Memoryに値を書き込む.

5 性能評価

本項では, 提案したGPU実装の性能評価を行う. 本研究では, GPUとしてNVIDIA A100 Tensor Core GPUを, CPUとしてAMD EPYC 7702 CPUをそれぞれ使用し, GPU/CPUのそれぞれにbatchベースの実装と再帰関数ベースの実装を行った. ここで, batchベースの実装は式(3)の漸化式 R の計算を提案アルゴリズムを用いて計算するものであり, 対して再帰関数ベースの実装は同様の計算を式(5)を用いて計算するものである.

事前実験として, 2種類の実装の実行時間を後述の対象分子について計測したところ, GPU/CPU共にすべてのケースでbatchベースの実装の方が高速であるという結果が得られ, batchベースの実装が R の冗長な計算を削減できていることが示された. よって, 本項ではbatchベースのGPU/CPU実装について性能評価を行う.

表1に, 単原子分子である水素原子に対してcorrelation consistent基底関数セット[5]を適用した場合のERI計算時間を示す. 結果より, ERI計算総数の少ないcc-pVDZを除いたすべての基底関数セットについてCPUに対して高速化を達成しており, 最大で101倍高速となった. また, ERI計算の総数や K の最大値が増加するにつれ全体的な実行時間は増加している.

表 1 提案手法を用いた単原子分子のERI計算時間 [sec]

基底関数	monatomic Hydrogen (H)				
	cc-pVDZ	cc-pVTZ	cc-pVQZ	cc-pV5Z	cc-pV6Z
M	5	15	35	70	126
$\max(K)$	4	8	12	16	20
CPU	0.001	0.135	9.680	766.249	48,479.200
GPU	0.002	0.008	0.183	8.713	476.624
高速化率	0.60	16.67	52.90	87.94	101.71

表2に, 様々な多原子分子でのERI計算時間を示す. 多原子分子としてベンゼン, ナフタレン, 酸化銅を使用し, 基底関数としてSTO-3Gと6-31G**をそれぞれ適用した. 結果として, CPUと比較し最大で24倍高速となった. ベンゼン及びナフタレンは, 分子数が異なるが, 原子配置が似たような構造を持っており, K の最大値は同じとなっている. よって, 実行時間は異なるものの, 高速化率は近い値となっている. また, 酸化銅はベンゼンやナフタレンと比較し, K の値が大きい基底関数の割合が高い. その上で, K の最大値が同じであっても高速化率が高くなっている為, batchベースの実装は K の値が大きいほど高い効果を発揮する実装だといえる.

表 2 提案手法を用いた多原子分子のERI計算時間 [sec]

基底関数	Benzene (C ₆ H ₆)		Naphthalene (C ₁₀ H ₈)		Copper oxide (CuO)	
	STO-3G	6-31G**	STO-3G	6-31G**	STO-3G	6-31G**
M	36	120	58	190	24	54
$\max(K)$	4	8	4	8	8	12
CPU	17.198	584.714	113.575	3972.26	14.847	274.386
GPU	2.385	44.736	15.722	290.210	0.841	11.141
高速化率	7.21	13.07	7.22	13.69	17.65	24.63

6 結論

本論文では, ERI計算に用いられるMDアルゴリズムの為の効率的なGPU実装を提案した. 提案手法では特にGPUに適した並列処理単位のbatchを導入し, 限られたサイズのShared Memoryを用いてERI計算を行うことが可能となった. NVIDIA A100 Tensor Core GPU上での実装は, AMD EPYC 7702 CPU上での実装と比較し, 単原子分子では最大101倍, 多原子分子では最大24倍の高速化を達成した.

参考文献

- [1] McMurchie, L. E., Davidson, E. R. (1978). One- and two-electron integrals over Cartesian Gaussian functions. *Journal of Computational Physics*, 26(2), 218-231.
- [2] Obara, S., Saika, A. (1986). Efficient recursive computation of molecular integrals over Cartesian Gaussian functions. *The Journal of Chemical Physics*, 84(7), 3963-3974.
- [3] Tsuji, S., Ito, Y., Nakano, K., Kasagi, A. (2023). Efficient GPU-Accelerated Bulk Evaluation of the Boys Function for Quantum Chemistry. In 2023 Eleventh International Symposium on Computing and Networking (CANDAR) (pp. 49-58).
- [4] NVIDIA Corporation. (2024). CUDA Programming Guide
- [5] Dunning, T. H. Jr. (1989). Gaussian basis sets for use in correlated molecular calculations. I. The atoms boron through neon and hydrogen. *The Journal of Chemical Physics*, 90(2), 1007-1023.