

# GPUを用いたメタゲノム配列相同性解析ツールの MPI並列化と応用

藤井 智也<sup>1</sup> 角田 将典<sup>2</sup> 大上 雅史<sup>2</sup> 石田 貴士<sup>2,3</sup> 石原 和幸<sup>4</sup> 秋山 泰<sup>2,3,a)</sup>

**概要:** DNA シーケンサーの発展と共に、細菌叢中のゲノム DNA を断片化してランダムにシーケンシングするショットガンメタゲノムが、メタゲノム解析の一手法として注目されるようになった。それに伴い、DNA 断片配列データを短時間で処理できる配列相同性解析ツールの需要が増してきている。本研究では、Suzuki らによって開発された高速なショットガンメタゲノム向け配列相同性検索ツール GHOSTZ-GPU に対し、MPI ライブラリを用いたマルチノード並列化を行った。近年のクラスタ型計算機では、共有ディスクの他に各計算ノードにローカル SSD 領域を有するものが存在し、そのようなローカル SSD の活用による I/O の高速化を提案した。開発した GHOSTZ-MP は、4 ノード実行時を基準として 32 倍の計算資源である 128 ノード実行時と比較すると、ローカルストレージ領域を利用しない方法では約 9.7 倍、利用する方法では約 14.7 倍の処理速度を達成した。さらに、GHOSTZ-MP を口腔内から採取された歯肉縁下プラークのメタゲノム配列データに適用した。東京工業大学 TSUBAME 2.5 スーパーコンピュータ 128 ノードで約 30,000 reads/sec の速度で相同性検索が可能であることを示し、実際に Socransky の分類に基づく歯周病関連細菌の存在比率の解析も行った。

**キーワード:** ショットガンメタゲノム解析, MPI 並列化, GHOSTZ-MP

## 1. 導入

DNA シーケンサーの発展と共に、細菌叢中のゲノム DNA を断片化してランダムにシーケンシングするショットガンメタゲノムが、メタゲノム解析の一手法として注目されるようになった。最新のシーケンサーでは 1 日あたり 5000 億~6000 億塩基という膨大な量のゲノムが解読でき [1]、ショットガンメタゲノム解析の追い風にもなっている。しかし、ショットガンメタゲノム解析では単なるリードマッピングではなく、配列相同性検索を行うことが求められ、配列データの生成に比べて計算が追いつかないことが問題となっている。

このような背景の下で、BLAST [2] に代わる配列相同性検索ツールが多数提案されてきた。BLAT [3] や RAPSearch [4]、DIAMOND [5] などはその代表例であるが、Suzuki らによって開発された GHOSTX [6] および GHOSTZ [7] は、ショットガンメタゲノム解析に耐える感度を保ちつつ BLAST に比べて約 260 倍 (GHOSTZ)

の高速化を達成し、KEGG のメタゲノムアノテーションサービスである GhostKOALA のエンジンとしても活用されている [8]。GHOSTZ はクエリとデータベースの双方にハッシュテーブルを適用し、またデータベース中の似ているアラインメント同士の前処理 (クラスタリング) と三角不等式による類似度上界を用いることで高速化を実現している。また、GPU アクセラレータを保有している場合は GHOSTZ-GPU [9] によって、CPU クラスタマシンで使用する場合は GHOSTX を MPI 並列化した GHOST-MP [10] によって、それぞれさらなる計算時間の削減が可能となっている。

しかしながら、GHOSTZ-GPU は単一ノードでしか実行できず、GHOST-MP は GPU が利用できないという問題が存在する。GHOSTZ-GPU のマルチノード並列化が望まれるが、単純な並列化実装ではデータベースへのアクセス集中による性能低下も想定される。そこで本研究では、GHOSTZ-GPU のマルチノード並列化を、データベースを共有ファイルシステムではなくローカルストレージ領域に配置してアクセスすることで、I/O の集中を減らす実装を取り入れた GHOSTZ-MP を開発した。

この GHOSTZ-MP を用いて、本研究では実際の口腔内メタゲノム配列データの解析も行った。

<sup>1</sup> 東京工業大学 工学部 情報工学科  
<sup>2</sup> 東京工業大学 大学院情報理工学専攻 計算工学専攻  
<sup>3</sup> 東京工業大学 情報生命博士教育院  
<sup>4</sup> 東京歯科大学 大学院歯学研究科 歯学専攻  
a) akiyama@cs.titech.ac.jp

## 2. GHOSTZ-GPU の MPI 並列化

本研究では、東京工業大学 秋山研究室で開発されたマルチ GPU 対応の GHOSTZ-GPU に、MPI 並列化を適用する。MPI 並列版 GHOSTZ-GPU の実装においては、当研究室で開発された並列分散処理用負荷分散ツールである MPIDP を利用した。以下に MPIDP についての説明と、本研究での GHOSTZ-MP の MPI 並列実装内容について述べる。

### 2.1 MPIDP

MPIDP は、並列分散処理用負荷分散を目的とした MPI ライブラリによるマスター・ワーカー型の汎用負荷分散ツールである。汎用性と移植性に優れており、既存のプログラムに組み込むことで簡易 MPI プログラムとして仕立てることが可能で、GHOST-MP [10] や MEGADOCK 3.0 [11] などに利用されてきた。マスターはあらかじめ作成されているコマンドリストを読み込み、各行に記述されている処理の実行をワーカーに指示する。MPIDP 自体に耐ノード障害機能やログファイル生成などのオプションが実装されており、場合に応じて所望のプログラム上でこれらの機能を活用することが可能である。

### 2.2 ベースライン手法：MPIDP の単純適用

MPIDP を用いて GHOSTZ-GPU を単純に MPI 並列化した。この実装を以降「GHOSTZ-MP (単純並列実装)」と呼ぶ。なお、この処理には前処理としてクエリをあらかじめ分割しておく必要があることと、相同性検索用のデータベースは共有ファイルシステムに置かれることに注意されたい。

### 2.3 提案手法：ローカル SSD を利用した実装

ベースライン手法として挙げた MPIDP の単純適用では、実行ノード数が増えるにつれてデータベースファイルへの I/O 集中によって並列化性能が頭打ちになることが推測される (図 1)。一方、近年では TSUBAME や国立遺伝学研究所のスーパーコンピュータ [12] などに代表される、共有ファイルシステムの他にノード内にローカル SSD 領域を持つ並列計算機が登場しており、この領域を活用することで I/O 集中を避けることが可能である。本研究では、各計算ノード (ワーカーノード) のローカル SSD 領域に MPI の集団通信の B\_CAST 関数を用いてデータベースをコピーして配置し、各ノードは GHOSTZ-GPU を実行するとき各々のローカル SSD にあるデータベースファイルを参照させることで I/O の集中を防ぐ実装を提案した (図 2) 以降このローカル SSD を利用した実装を「GHOSTZ-MP (SSD 利用実装)」と呼ぶ。

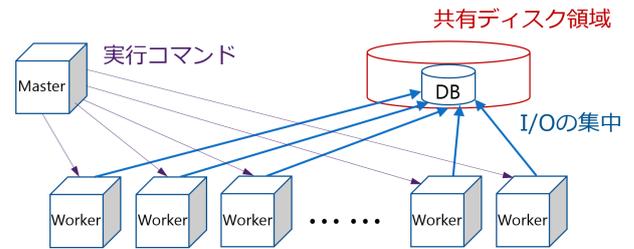


図 1 データベースファイルへの I/O 集中の概要図 (ベースライン手法)

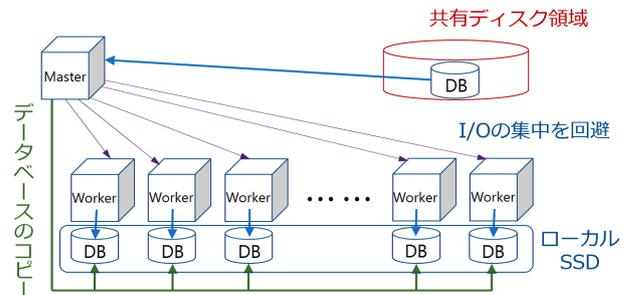


図 2 ローカルディスクを用いた I/O 分散の概要図 (提案手法)

TSUBAME 2.5 の Thin ノードのローカル SSD の容量は 50 GB 程度であり、KEGG GENES 2015.2 版を GHOSTZ のハッシュテーブルとクラスタリングによってインデックス化したデータベースの容量は合計で 29.9 GB である。提案手法ではワーカーノードが計算する部分のクエリとデータベースは各ワーカーノードのローカルストレージにコピーしているため、実行ノード数が  $n$  のとき、1 ワーカーノードあたりが計算するクエリの容量は (全体のクエリの容量) /  $(n-1)$  である。すなわち、

$$\begin{aligned} & \text{(ローカルストレージの容量)} \\ & \geq \frac{\text{(全体のクエリの容量)}}{n-1} + \text{(データベースの容量)} \end{aligned}$$

であるときにローカルストレージへのデータベース配置が可能となる。今回の例では、TSUBAME 2.5 の Thin ノードを考えた場合、入力クエリの大きさが少なくとも 20 GB 以内であれば少ないノードでも実行が可能であり、例えば  $n = 128$  ノード利用時には約 2.5 TB までのクエリ配列を扱うことが可能となる。

## 3. 評価実験

### 3.1 実験環境

GHOSTZ-MP の並列性能を確かめるため、以下の評価実験を行った。実行環境は TSUBAME 2.5 の Thin ノード (表 1) である。また、使用クエリとして舌背部のメタゲノムデータである SRS078182 (リード数 146,908,592 本 (18.9 GB)、最長リード長 95 塩基 (全体の 71.4%)) を用い、データベースには KEGG GENES 2015.2 版 (アミノ酸配列数 15,248,714 本 (6.2 GB)) を用いた。MPIDP のパラメータとして、計測においてはノード障害時のリトライ

表 1 TSUBAME 2.5 Thin ノードの環境

CPU	Xeon 5670 (2.93 GHz, 6 cores) ×2
Memory	54 GB
GPU	Tesla K20X (732 MHz, 2688 CUDA cores) ×3
SSD	50 GB

表 2 実行時間の測定結果. 値は 5 回計測したときの平均値 (sec) で, 括弧内は標準偏差, 斜体は実行速度倍率 (各実装の 4 ノード基準) である.

	Number of Nodes					
	4	8	16	32	64	128
単純実装	44,220 (3,350) <i>1.0</i>	19,055 (869) <i>2.3</i>	10,169 (1,140) <i>4.3</i>	5,805 (983) <i>7.6</i>	5,271 (149) <i>8.4</i>	4,537 (500) <i>9.7</i>
SSD 実装	44,349 (1,744) <i>1.0</i>	20,212 (1,498) <i>2.2</i>	10,368 (948) <i>4.3</i>	5,555 (253) <i>8.0</i>	3,800 (97) <i>11.7</i>	3,011 (79) <i>14.7</i>

機能を off に設定した. また, GHOSTZ-GPU のパラメータとして

- 使用 CPU コア数 (-a パラメータ) = 12
- 使用 GPU 枚数 (-g パラメータ) = 3
- クエリチャンクサイズ (-l パラメータ) = 33,554,432 (32 MB)

を用いた.

### 3.2 測定結果

ここでは GHOSTZ-MP (単純並列実装) をマルチノード実行させるときと GHOSTZ-MP (ローカル SSD を利用した実装) をマルチノード実行させるときのパフォーマンスをそれぞれ示す (表 2, 図 3, 図 4). 128 ノード時で単純並列実装に比べて, SSD 実装は 1.6 倍の速度向上を達成した. また, 並列化効率も単純並列実装より SSD 実装の方が優れていることが示された. また, 32 ノードまでは同程度のパフォーマンスであるが, 64 ノードからは提案手法の方がパフォーマンスが良い. これはデータベースのコピーを行わない単純実装について, 実行ノード数が 32 ノードまでならデータベースへの I/O 集中の影響はそれほど受けませんが, 64 ノードまで増えると I/O 集中の影響が無視できなくなるためだと考えられる.

## 4. 口腔内メタゲノム解析への応用

東京歯科大学にて採取された 6 人の被験者 (いずれも歯周病を患っていない健康者) の口腔内歯肉縁下プラークからシーケンシングされたメタゲノム配列データを対象に, GHOSTZ-MP による相同性検索を行った. シーケンシングは MiSeq によってペアエンドで行われ, それぞれの被験者は性別と年齢が分かっている (表 3).

GHOSTZ-MP の出力から Species レベルでの細菌種の存在比率を集計し, 歯周病の関連度に応じて口腔内細菌を分類した Socransky の分類 [13, 14] によって特徴付けられて

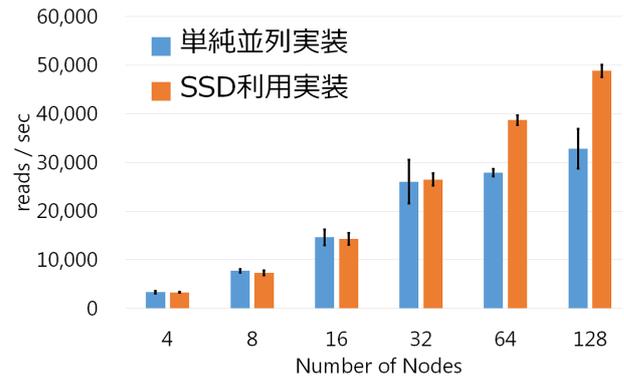


図 3 単純実装とローカル SSD を利用した実装でのパフォーマンス比較

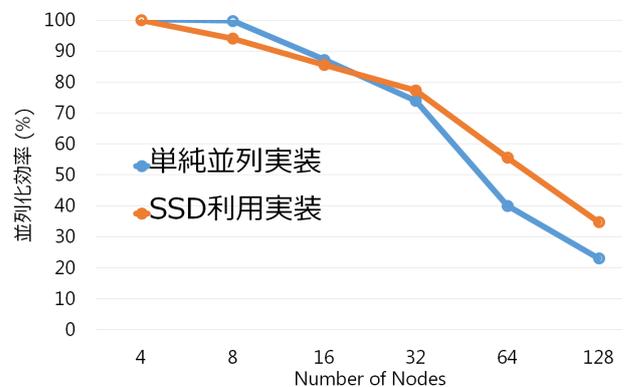


図 4 並列化効率比較 (SSD 利用実装の 4 ノードを基準, 強スケーリング)

表 3 各サンプルデータの詳細

サンプル id	被験者 S1	被験者 S2	被験者 S3
リード数	1,013,737	444,296	1,029,821
性別	男性	男性	男性
年齢	青年	壮年	壮年
サンプル id	被験者 S4	被験者 S5	被験者 S6
リード数	907,627	917,614	631,217
性別	女性	女性	男性
年齢	壮年	青年	青年

いる細菌種の存在比率を比較した (図 5). TSUBAME 128 ノードで GHOSTZ-MP を実行すると約で終了する規模の解析である. 解析の結果, 被験者 S3 と S5 は Socransky の分類中の Red Complex に属する菌種が極めて少なく, Yellow Complex に属する *Streptococcus* 属が多いという結果が示された. このことから被験者 S3 と S5 は他の被験者に比べて歯周病リスクが低いと推定される.

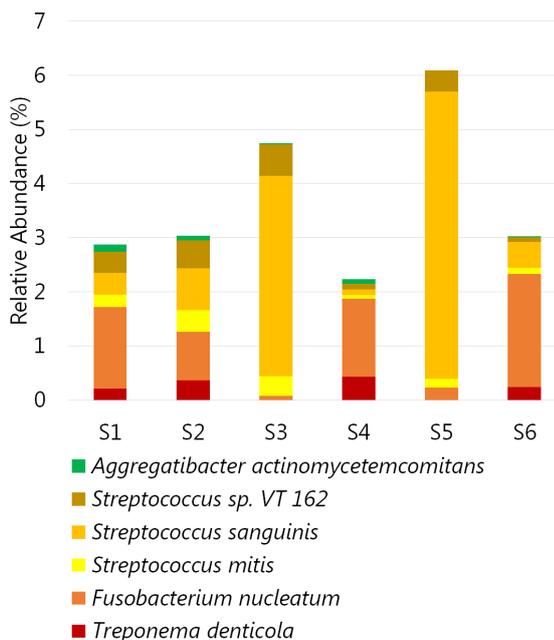


図 5 6 人の被験者サンプルにおける Socransky の分類による Species 階層細菌種の存在比率 (%)

## 5. 結論

本研究では配列相同性検索ソフトウェア GHOSTZ-GPU を MPI 並列化し、大規模に実行可能な GHOSTZ-MP を新たに開発した。検索用データベースを各ノードが保持するローカルストレージ領域に MPI 通信を用いて配置することで、データベースへの I/O の集中を避け、計算速度ならびに並列化効率の双方の改善に成功した。また、応用として口腔内メタゲノム解析を実施し、健常者間での細菌叢の違いを観察した。このショットガンメタゲノム解析は GHOSTZ-MP と計算ノード 128 台を用いてわずか 10 分以内に計算が完了するものであり、今後より多くのサンプル／リード配列が得られるようになっても現実的な時間内での解析が本研究によって実施可能となった。

今後の課題として、MPI プロセス間での通信の最適化による高速化と、口腔内メタゲノム応用で発見された相同遺伝子群の詳細な解析が挙げられる。

**謝辞** 本研究の一部は、文部科学省 HPCI 戦略プログラム 分野 1「予測する生命科学・医療および創薬基盤」および、JST CREST「EBD：次世代の年ヨッタバイト処理に向けたエクストリームビッグデータの基盤技術」の支援を受けて行われた。

## 参考文献

- [1] [http://genaport.genaris.com/GOC\\_sequencer\\_post.php?eid=00093](http://genaport.genaris.com/GOC_sequencer_post.php?eid=00093)
- [2] Altschul SF, et al. Basic local alignment search tool, *J*

- Mol Biol*, 215: 403–410 (1990)
- [3] Kent WJ. BLAT-the BLAST-like alignment tool, *Genome Res*, 12: 656–664 (2002)
- [4] Ye Y, et al. RAPSearch: a fast protein similarity search tool for short reads, *BMC Bioinform*, 12: 159 (2011)
- [5] Buchfink B, et al. Fast and sensitive protein alignment using DIAMOND, *Nat Methods*, 12: 59–60 (2015)
- [6] Suzuki S, et al. GHOSTX: an improved sequence homology search algorithm using a query suffix array and a database suffix array, *PLoS ONE*, 9: e103833 (2014)
- [7] Suzuki S, et al. Faster sequence homology searches by clustering subsequences, *Bioinformatics*, 31: 1183–1190 (2015)
- [8] Kanehisa M, et al. BlastKOALA and GhostKOALA: KEGG Tools for Functional Characterization of Genome and Metagenome Sequences, *J Mol Biol.* (in press)
- [9] Suzuki S, et al. GPU-Acceleration of Sequence Homology Searches with Database Subsequence Clustering. (submitted)
- [10] Kakuta M, et al. A massively parallel sequence similarity search for metagenomic sequencing data. (submitted)
- [11] Matsuzaki Y, et al. MEGADOCK 3.0: a high-performance protein-protein interaction prediction software using hybrid parallel computing for petascale supercomputing environments, *Source Code Biol Med*, 8: 18 (2013)
- [12] <https://sc.ddbj.nig.ac.jp/>
- [13] Socransky SS, et al. Microbial complexes in subgingival plaque, *J Clin Periodontol*, 25: 134–144 (1998)
- [14] Socransky SS and Haffajee AD. Dental biofilms: difficult therapeutic targets, *Periodontol*, 28: 12–55 (2000)