

# Microsoft Azure 上でのタンパク質間相互作用 予測システムの並列計算と性能評価

大上 雅史<sup>1,a)</sup> 山本 悠生<sup>1,2</sup> 秋山 泰<sup>1,2</sup>

**概要:** 多数のタンパク質のタンパク質間相互作用を計算機で予測するために、立体構造情報を活用して全てのタンパク質ペアに対して網羅的に計算を行う方法が提案されてきた。本研究では立体構造情報に基づくタンパク質間相互作用予測の並列版ソフトウェアである MEGADOCK を、パブリッククラウド計算環境である Microsoft Azure 上に移植し、並列計算性能を評価した。GPU 搭載型のバーチャルマシン 20 台を用いて 480 CPU コア・80 GPU による並列計算を行った結果、5 台利用時と比較して 89% の並列化効率が得られ、GPU による加速率も GPU 利用に伴う利用料金の増加率と比較して倍以上大きいことが示された。

**キーワード:** パブリッククラウド、Microsoft Azure、タンパク質間相互作用、MEGADOCK、GPU

## Parallel computing of protein-protein interaction prediction system MEGADOCK on Microsoft Azure

MASAHITO OHUE<sup>1,a)</sup> YUKI YAMAMOTO<sup>1,2</sup> YUTAKA AKIYAMA<sup>1,2</sup>

**Abstract:** To predict protein-protein interactions (PPIs) of a large number of proteins, a method of comprehensively calculating all pairs of proteins using tertiary structure information was proposed. In this study, MEGADOCK which was parallel computing software for structure-based PPI prediction software was ported on Microsoft Azure, which is a public cloud computing environment, and parallel computing performance was evaluated. As a result of parallel computation with 480 CPU cores and 80 GPUs using 20 GPU-equipped virtual machines, parallel efficiency of 89% was obtained compared with using 5 virtual machines. The acceleration ratio by GPUs was better than doubled of use fee which increased by using GPUs.

**Keywords:** Public cloud, Microsoft Azure, Protein-protein interaction, MEGADOCK, GPU

### 1. 序論

タンパク質間相互作用 (protein-protein interaction, PPI) は様々な細胞内プロセスや機能に関わる現象であり、疾病の理解や創薬標的の決定などにその情報が活用されている。PPI は一般に生化学実験によって決定されるものであるが、近年では計算機による予測手法がさかんに研究され

ており、既知配列情報を用いるもの、単体のタンパク質立体構造を用いるもの、既知複合体構造に対する相同性検索に基づくものなどが提案されている [1]。我々は単体のタンパク質立体構造を用いた手法として MEGADOCK [2] を提案しており、走化性 [3] やアポトーシス [4] などの PPI ネットワークの予測に応用してきた。

PPI ネットワークのように複数のタンパク質同士の相互作用関係を予測するためには、全てのタンパク質の組み合わせに対して予測計算を実施する必要がある。MEGADOCK では、このような「全対全」の予測を迅速に行うために、マルチノード並列化 [5] や、マルチ GPU ノード並列化 [6], [7] などが完了している。

<sup>1</sup> 東京工業大学 情報理工学 情報工学系  
Department of Computer Science, School of Computing,  
Tokyo Institute of Technology  
<sup>2</sup> 東京工業大学 情報生命博士教育院  
Education Academy of Computational Life Sciences, Tokyo  
Institute of Technology  
<sup>a)</sup> ohue@c.titech.ac.jp

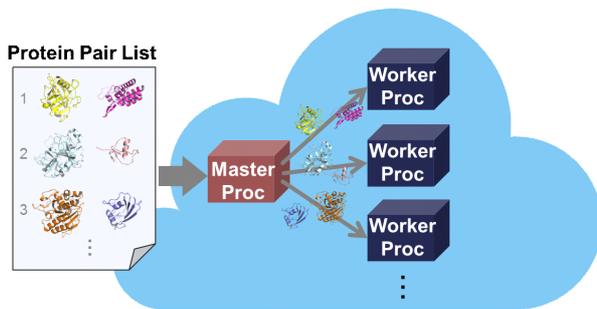


図 1 MEGADOCK の Azure VM を利用した並列計算の概要図。

一方でこれらはスーパーコンピュータ等の並列計算機を想定して実装されており、小規模な利用や、利用者の要求に応じた即時対応などは困難であった。こうした需要に応えられる計算リソースとしてパブリッククラウドが挙げられ、Amazon EC2 や Microsoft Azure (Azure) などのパブリッククラウドサービスが広く普及している [8]。パブリッククラウド環境で大規模なバイオインフォマティクス計算を行った事例もいくつか報告されており [9], [10]、パブリッククラウドへの注目が高まっている。

本研究では、PPI 予測ソフトウェアである MEGADOCK に対し、主要なパブリッククラウドサービスの 1 つである Azure の複数のバーチャルマシン上で並列実行可能にすることを目的とし、実装と並列計算性能の評価を行った。

## 2. MEGADOCK の Azure への実装

### 2.1 プロセス並列とスレッド並列

MEGADOCK の計算は、タンパク質ペアごとに独立して行うことができるため、タンパク質ペアのデータ並列と、1 つのタンパク質ペアの PPI 予測計算のスレッド並列のハイブリッド並列として実装されている [5]。MPIDP フレームワーク [5], [11] による Master-Worker 型の実装が行われており、Master プロセスによって PPI 予測のタスクが Worker プロセスに割り振られ、タスク分散が行われている。Master と Worker は MPI 通信を行い、Worker 同士の通信は行わず、Worker プロセス内は OpenMP (GPU を併用する場合には OpenMP と CUDA) によるスレッド並列で計算を実行する。本研究では並列実装については従来のものを踏襲し、Azure のバーチャルマシン (virtual machine, VM) 上で同様にハイブリッド並列計算を行えるようにした (図 1, 図 2)。

MPI プロセス数とスレッド数の割り振りについては Azure の VM の種類 (インスタンス) によって検討の余地があるが、本研究では各 VM に 4 プロセス、各プロセスに (コア数/4) のスレッド数を割り当てて実行した。また、本研究では GPU の利用を前提とし、4 つの GPU が搭載されているインスタンスを利用することで、各プロセスが 1 つずつ GPU を利用するようにした。

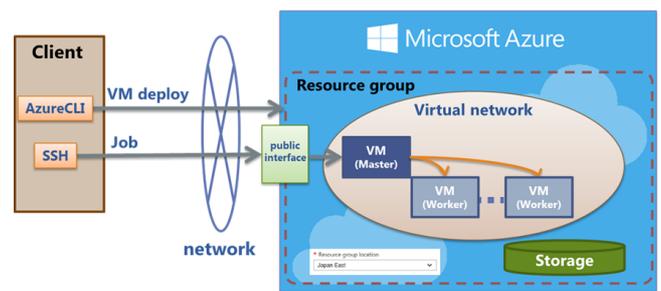


図 2 Azure 上のアーキテクチャの概要図。

表 1 使用した Azure インスタンスの詳細

CPU	Intel Xeon E5-2690v3 2.6GHz ×2 (24 コア)
GPU	NVIDIA Tesla K80 ×4 (物理的には 2 基)
RAM	224 GB
Disk	1.44 TB SSD
料金	NC24: 440.65 円/h, NC24r: 484.71 円/h
備考	NC24r は RDMA (Infiniband) で構成

### 2.2 インスタンスの種類

Azure ではインスタンスが約 70 種提供されており、それぞれハードウェアやネットワーク、ストレージなどが異なっている。目的の計算に応じて適切に選択する必要があるが、本研究では、代表的な GPU インスタンスである NC24 と NC24r を用いて検証を行った。インスタンスの詳細を表 1 に示す。NC24 および NC24r インスタンスでは GPU に Tesla K80 が搭載されているが、Tesla K80 は 1 つのボードに 2 つの GPU が含まれているため、各 VM に Tesla K80 が 2 基搭載され、利用可能な GPU が 4 つあることに注意されたい。

## 3. 性能評価実験

### 3.1 使用データセット

MEGADOCK の Azure 上での計算性能の評価には、Protein-Protein Docking Benchmark 1.0[12] の複合体構造データセットを用いた。このデータセットには、二量体の複合体構造が 59 個含まれており、二量体の各構成タンパク質は r と 1 という名前で区別されている。本研究では、この r と 1 の全ての組み合わせである  $59 \times 59 = 3,481$  通りのタンパク質ペアに対する予測計算を行い、計算が完了するまでの時間を計測した。

### 3.2 結果

以降、「 $n$  台の VM」を  $\#VM = n$  と表記する。Azure の NC24 および NC24r インスタンスで、 $\#VM = 5$ ,  $\#VM = 10$ ,  $\#VM = 20$ ,  $\#VM = 22$  について計算時間の計測を行った。1 分あたりに計算できたタンパク質ペア数の値を図 3 に示す。なお、Azure の VM の最大数は Microsoft 側で管理されており、我々のクォータ制限下では最大  $\#VM = 22$  まで利用が可能であった。

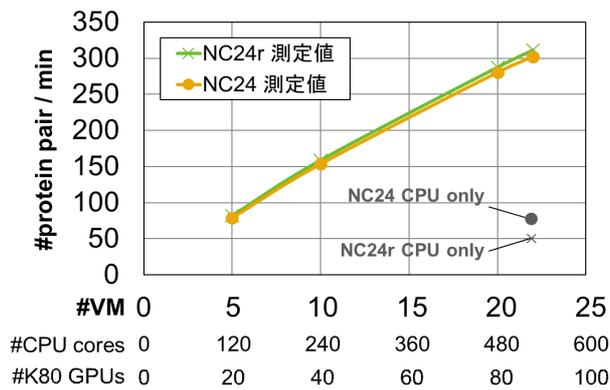


図 3 NC24, NC24r での計算時間の測定結果。

NC24 および NC24r のインスタンスの違いについては、本研究では差が見られなかった。NC24 インスタンスにおける #VM = 5 に対する #VM = 20 の並列化効率（強スケール値）は 89%であった。

GPU を用いずに CPU のみで計算した結果についても #VM = 22 に限って図 3 に示した。CPU のみでは NC24r より NC24 が若干高速であった。NC24r は RDMA ネットワークインターフェースが使えるインスタンスであり、NC24 より高い通信性能を有するが、MEGADOCK に関しては特に影響しないといえる。NC24 の方が NC24r より安価であるため、MEGADOCK では NC24 を利用する方が好ましい。

### 3.3 GPU の効果

NC24 について、GPU を用いた場合と用いない場合では、GPU を用いる方が約 3.8 倍高速であった。NC24 と同等の CPU 性能を持つインスタンスは現在のところ提供されていないが、同クロック周波数でコア数が 2/3 (16 コア) である A9 インスタンス (CPU: Intel Xeon E5-2670 2.6GHz、197.06 円/h) が存在する。A9 と比較すると、

- A9 (24 コア相当) = 197.06 円/h × 24/16 = 295.59 円/h
- 440.65 円/h (NC24) ÷ 295.59 円/h (A9) ≒ 1.5 倍

であるから、本研究の GPU による 3.8 倍の高速化は、利用料金の観点でも有利であるといえる。

## 4. 結論

我々が開発した PPI 予測ソフトウェア MEGADOCK を、パブリッククラウドサービスの 1 つである Microsoft Azure の VM 上で並列計算が行えるように実装し、GPU 搭載型の VM を用いてその性能を評価した。480 CPU コア・80 GPU 規模の並列計算を実施し、良好な並列化効率が得られ、GPU による加速も VM の利用料金に対して倍以上有利であることが示された。

本研究では Microsoft のクォータコア制限により、NC24 および NC24r インスタンスに関しては 22 台までしか VM

の同時利用ができなかったが、現在この制限を緩和する手続きを行っており、より大規模な並列実行での評価が可能になる見込みである。CPU インスタンスではより多くの VM 数での評価が進んでいるが [13]、GPU インスタンスでのさらなる評価の実施は今後の課題である。

謝辞 本研究は、JSPS 科研費(基盤研究(A)24240044、若手研究(B)15K16081、基盤研究(C)24118088)、JST CREST「EBD: 次世代の年ヨット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」、JST リサーチコンプレックス推進プログラム「世界に誇る社会システムと技術の革新で新産業を創る Wellbeing Research Campus “Tonomachi”」、Microsoft Business Investment Funding、リバネス研究費の支援を受けて行われた。

## 参考文献

- [1] Matsuzaki, Y. *et al.* (in press) Rigid-docking approaches to explore protein-protein interaction space. *Adv. Biochem. Eng./Biotechnol.*
- [2] Ohue, M. *et al.* (2014) MEGADOCK: An all-to-all protein-protein interaction prediction system using tertiary structure data. *Prot. Pept. Lett.*, 21(8): 766–778.
- [3] Matsuzaki, Y. *et al.* (2014) Protein-protein interaction network prediction by using rigid-body docking tools: application to bacterial chemotaxis. *Prot. Pept. Lett.*, 21(8): 790–798.
- [4] Ohue, M. *et al.* (2013) Highly precise protein-protein interaction prediction based on consensus between template-based and de novo docking methods. *BMC Proc.*, 7(Suppl 7): S6.
- [5] Matsuzaki, Y. *et al.* (2013) MEGADOCK 3.0: a high-performance protein-protein interaction prediction software using hybrid parallel computing for petascale supercomputing environments. *Source Code Biol. Med.*, 8(1): 18.
- [6] Shimoda, T. *et al.* (2013) MEGADOCK-GPU: acceleration of protein-protein docking calculation on GPUs. In *Proc. of ACM-BCB'13*, 883–889.
- [7] Ohue, M. *et al.* (2014) MEGADOCK 4.0: an ultra-high-performance protein-protein docking software for heterogeneous supercomputers. *Bioinformatics*, 30(22): 3281–3283.
- [8] Hashem, I.A.T. *et al.* (2015) The rise of “big data” on cloud computing: Review and open research issues. *Inf. Syst.*, 47: 98–115.
- [9] Ekanayake, J. *et al.* (2011) Cloud technologies for bioinformatics applications. *IEEE Trans. Parallel. Distrib. Syst.*, 22(6): 998–1011.
- [10] Shanahan, H.P. *et al.* (2014) Bioinformatics on the Cloud Computing Platform Azure. *PLOS ONE*, 9(7): e102642.
- [11] 青山 健人, 他. (2016) スーパーコンピュータ「京」上でのエクソーム解析パイプラインの開発. *情報処理学会論文誌 コンピューティングシステム (ACS)*, 9(2): 15–33.
- [12] Chen, R. *et al.* (2003) A Protein-Protein Docking Benchmark. *Proteins*, 52(1): 88–91.
- [13] 青山 健人, 他. (2017) コンテナ型仮想化による分散計算環境におけるタンパク質間相互作用予測システムの性能評価. *情報処理学会研究報告 バイオ情報学*, 2017-BIO-49(3): 1–8.