

タンパク質立体構造予測システム AlphaFold の TSUBAME3.0 上での高速化

藤田 隼斗^{1,a)} 野村 哲弘² 遠藤 敏夫^{2,1} 関嶋 政和¹

概要: タンパク質の立体構造の理解はタンパク質の機能の理解に役立ち、人体の働きを理解する基礎となる。ただし既知のタンパク質のアミノ酸配列に対して、実験的に立体構造が決定されているタンパク質は多いとは言えない。これらのギャップを埋めるために計算機によるアプローチが必要であり、そのアプローチの1つとして開発されたのが AlphaFold である。

AlphaFold は hhblits 等の既存のツールを使い、BFD などの巨大な遺伝子データベースから MSA を取得する。しかし、そのうちの1つの hhblits の実行に多大な時間がかかるという問題がある。この実行時間の一番の要因は I/O にあり、遺伝子データベースをどういったストレージに保管するかによって実行時間が大きく異なってくる。実際に TSUBAME3.0 の高速ストレージ領域にデフォルトの設定で保管した場合と Stripe 設定をして保管した場合とで実行時間が大きく異なっている。

そこで、本研究では TSUBAME3.0 上で選択できる遺伝子データベースの保管方法を変更し、実行時間を計測した。さらに、ボトルネックとなっている hhblits に対してプロファイリングを行い、その結果から並列化数のチューニングやデータベースの保存方法の変更、ソートの最適化を行った。また、AlphaFold の MSA 取得ツール間のデータ依存関係からツール実行を非同期に行うようにした。その結果、TSUBAME3.0 上で単一のアミノ酸配列から立体構造を予測する際の実行時間を平均して2分の1に短縮することに成功した。

1. はじめに

タンパク質のアミノ酸配列からそのタンパク質の立体構造を予測することは、「タンパク質のフォールディング(折り畳み)問題」として知られている。

ただし既知のタンパク質のアミノ酸配列に対して、実験的に立体構造が決定されているタンパク質は約1万分の1である [1][2][3]。実験的にタンパク質の立体構造を推定する方法はいくつか存在するが、どの方法も多大なコストや時間がかかる。このような状況と近年の情報技術の発達を受けて、計算機によるアプローチ台頭してきている。そのような中で第14回タンパク質構造予測精密評価(CASP14)[4]にて AlphaFold2[5]が登場し、他の手法を大きく上回る結果を残した。その精度は実験によるタンパク質の立体構造推定とほぼ同等である。

AlphaFold のソースコードは Apache License 2.0 でオー

ブンソース化されており、学習済みパラメータはモデルパラメータは2022年1月19日より CC BY 4.0 のライセンスで公開されている [6]。遺伝子データベースは2022年1月20日より PDB70 のライセンスが CC BY 4.0 に変更された [7] ことによって、AlphaFold は商用利用を含めて誰でも利用できる状態である。

ただし、AlphaFold の実行には合計で2TB以上にも及ぶ巨大な遺伝子データベースが必要であり、長いアミノ酸配列の予測には大きな VRAM を持った強力な GPU が必要である。さらに高速に実行するためには巨大な遺伝子データベースを HDD よりも高速な SSD やシステムメモリ上に置く必要があり、記憶域の面でも高性能な計算機が求められる。AlphaFold でも、ある程度の品質をトレードオフにしてより小さい遺伝子データベースで実行するオプションがリリースされている [8]。

AlphaFold は TSUBAME3.0 でも2021年8月より利用できるようになっている [9]。ただし、TSUBAME3.0 では hhblits の実行時間が遅く、また実行時間にむらがあった。特に I/O に難があるとされており、遺伝子データベースの保管設定が変更された [10]。これによって hhblits の実行時間は大幅に改善されるが、依然としてそのツールが占め

¹ 東京工業大学情報理工学院
School of Computing, Tokyo Institute of Technology

² 東京工業大学学術国際情報センター
Global Scientific Information and Computing Center, Tokyo
Institute of Technology

a) fujita.h.ai@m.titech.ac.jp

る実行時間は大きい。本研究では、hhblitsの高速化を主として行う。その他のツールについては非同期実行することによってAlphaFold全体の実行時間に関係することが少なくなるようにし、AlphaFold全体の実行時間が小さくなるようにする。

2. AlphaFold

AlphaFold[5]はDeepMindによって開発されたタンパク質立体構造予測システムである。

AlphaFoldの予測は大きく分けて2つステージで構成される。1つはMSA(Multiple sequence alignment)を構築するステージ、もう1つは実際に立体構造を推論するステージである。1つ目のMSA構築ステージでは、既存のツールであるjackhmmer, hhsearch, hmmsearch, hhblitsを用いてMSAとテンプレートを構築する。これらには遺伝子データベースが必要で、jackhmmer[11]はUniref90[12]とMGnify[13]を別々に実行し、hhsearch, hmmsearchはPDB70[14], hhblitsはBFDとUniclust30[15]を同時に使用して実行する。2つ目の推論ステージでは、Embeddingモジュール, Evoformerモジュール, Structureモジュールで構成され、立体構造の推論を行っている。

2つのステージで必要とするリソースは異なり、MSA構築ステージではCPUのみを用いて、推論ステージでは主にGPUを用いる。MSA構築ステージでは巨大なデータベースに対してアクセスするため、I/O性能も重要である。特にhhblitsがBFD(1.7TB)を使用しており、I/O性能によっては実行時間の大きな割合を占める。

2.1 TSUBAME3.0上のAlphaFold

TSUBAME3.0は東京工業大学のスーパーコンピュータである。TSUBAME3.0は540台の計算ノードを備え、総計12.15PFLOPSの演算能力を供給する。各計算ノードのスペックは、表1の通りである。

ストレージ環境として高速ストレージ領域があり、これはLustreファイルシステムで構成されている[17]。Lustreとは並列分散ファイルシステムであり、多くのHPCで利用されている。Lustreに適するのは、大きなファイル、シーケンシャルアクセス、プロセス毎で異なるファイルに書き込む並列アプリケーションである。逆に苦手とするのは、大量の小さなファイル、莫大なメタデータ操作である。Lustreでの実データはObject Storage Target(OST)と呼ばれるものに保管される。TSUBAME3.0では10本のHDD(10TB 7200rpm)を用いたRAID6で構成され、これが68ある。通常は1つのファイルは1つのOSTに保存されるが、File Stripingを使うことで複数のOSTに分散してファイルが保存される。

AlphaFoldの遺伝子データベースは高速ストレージ領域にStripeされて保存されている。Stripeの設定はstripe.count

```
>dummy_sequence
GWSTELEKHKREELKEFLKKEGITNVEIRIDNGRLEVRVEGGTERLK
RFLEELRQKLEKGYTVDIKIE
```

図1 68残基のダミーのアミノ酸配列

が68, stripe.sizeが1MiB(1048576B)となっている。

以降では特に断りがない場合は計算ノードを1台使用するものとする。

2.2 AlphaFoldのボトルネック

Stripeされていないデータベースに対してはダミー配列をStripeされているデータベースに対してはダミー配列とCASP14のT1050を用いてそれぞれAlphaFoldを実行する。

AlphaFoldの標準エラー出力から各MSA取得ツールの実行時間とそれを除いたAlphaFoldの推論の実行時間をまとめたものが表2である。実行はAlphaFoldv2.0.1を用いて、オプションに関してはTSUBAME3.0の例[18]にあるように実行した。(予測はモデル1つ)MSA取得ツールについては標準エラー出力にある実行時間をそのまま使う。AlphaFoldの推論部の実行時間については最初のログ出力と最後のログ出力を全体の実行時間とし、そこから各MSA取得ツールの実行時間を引いたものとした。

表2から実行時間の大部分をhhblitsが占めていることやStripeの有無というデータベースの保管方法の違いによってhhblitsの実行時間が大きく改善されていることが分かる。ただし依然としてhhblitsの実行時間の割合は大きい。

T1050は779残基とダミー配列に比べ配列が長く、実際のアミノ酸配列である。jackhmmer, hhsearchの実行時間は大きく変化せず、hhblitsやAlphaFoldの推論部の実行時間が増大している。これらの結果からTSUBAME3.0上におけるAlphaFoldのボトルネックはhhblitsであると考えられ、AlphaFoldを高速化するためには、hhblitsを高速化する必要がある。

2.3 hhblitsのボトルネック

hhblits[19]はHH-suiteに含まれるツールの1つであり、hhsearchもその1つである。

HHblits(HMM-HMM-based lightningfast iterative sequence search)は非常に大きな遺伝子データベースをに対して、高速にMSA検索を行うことができる。hhblitsは名前にある通り、隠れマルコフモデル(HMM)を用いる。

図2は簡略なhhblitsのワークフローである。hhblitsが非常に大きなデータベースに対して、高速にMSA検索を行うことが出来る訳はPrefilterにある。高速なPrefilterを用いることで、データベースを1~5%にフィルタリングをする。これによってPrefilter以降の重い計算をデータベー

表 1 計算ノードのスペック [16]

CPU	Intel Xeon E5-2680 V4 (Broadwell-EP, 2.4GHz) × 2 ソケット ソケットあたり 14 コア, ノードあたり合計 28 コア
GPU	NVIDIA TESLA P100 for NVlink-Optimized servers × 4 基
メモリ	256GB (DDR4-2400 32GB モジュール × 8 本)
SSD	Intel DC P3500 2TB (NVMe, PCI-E 3.0 x4, R2700/W1800)
ネットワーク	Intel Omni-Path Architecture HFI (100Gbps) × 4 本

表 2 各ブロックの実行時間 (秒)

	no Stripe		Stripe	
	dummy	dummy	T1050	
jackhmmer (uniref90)	345	339	432	
jackhmmer (mgy_clusters)	406	408	434	
hhsearch	230	608	464	
hhblits	16264	1589	8456	
alphafold(inference)	167	172	1478	

表 3 主な関数の実行時間 (秒)

	dummy	T1050
prefilter_db	1029	1882
alignment	6	3233
perform_realign	0	1088
mergeHitsToQuery	—	1906
main(全体)	1153	8337

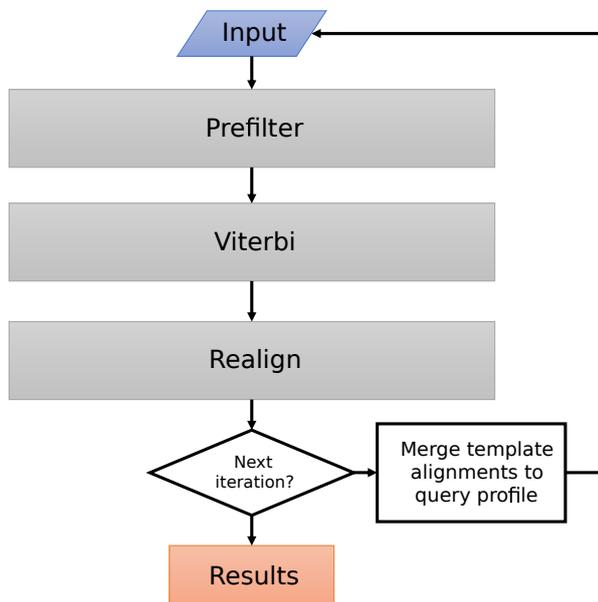


図 2 hhblits の簡略なワークフロー

```

$ module li
Currently Loaded Modulefiles:
 1) cuda/11.5.0           4) experimental
 2) openmpi/3.1.4-opa10.10 5) scorep/7.1
 3) gcc/10.2.0
  
```

図 3 Score-P を用いてビルドする際の環境

ス全体に対して行う必要がなくなる。

2章より, hhblits の高速化をする必要がある。そのため, hhblits のボトルネックを解析する。Score-P[20] を用いてボトルネックとなる関数を特定する。

Score-P でプロファイリングする際は Score-P を含めてビルドする必要があり, ビルド時の環境は図 3 を用いる。

Stripe されているデータベースに対して, ダミー配列と CASP14 の T1050 を入力配列として Score-P を含めてビルドした hhblits を実行する。

hhblits の本質的な処理が行われる関数は HHblits::run で

ある。この関数で直接呼ばれる関数の中で, Master Thread での実行時間が大きい関数をまとめたものが表 3 である。これらの 4 つ関数でどちらも全体の実行時間の 9 割ほどを占めている。実行されなかった関数は—で示した。

この 2 つの入力配列では関数の実行時間の傾向が大きく異なる。これは Prefilter を通過するアミノ酸配列の数が異なることに由来する。ダミー配列は配列に意味がなく, データベース内に似た配列は少なく Prefilter を通過する配列はごく少数である。そのため関数 prefilter_db 以外の実行時間が小さい。一方, T1050 はデータベース上に似た配列が数多く存在するために, 関数 prefilter_db 以外の実行時間も相対的に大きい。ダミー配列と T1050 の関数 prefilter_db の実行時間は表 3 では差が大きいのに見える。しかしこれは, AlphaFold での hhblits は 3 回のイテレーションを行う設定であるが, ダミー配列が 3 回のイテレーションをしていないためである。hhblits のログを解析すると, ダミー配列と T1050 の 1 回目の実行時間を見ると大きな差はない。また, 関数 prefilter_db は 2, 3 回目と呼ばれる際の実行時間は 1 回目と呼ばれる際の実行時間よりも小さい。

ボトルネックの 4 つの関数の中で関数 prefilter_db, alignment は I/O 性能に強く影響を受ける。両者の関数の I/O アクセスは mmap を介して行われている。また, これらの関数がアクセスするファイルは異なる。ソースコード解析とデバッグの結果, データベースが Uniclust30(86GB) の場合, 関数 prefilter_db では, 合計で 4GB ほどのファイルにアクセスし, 関数 alignment では少なくとも合計で 65GB ほどのファイルにアクセスする。

3. 手法

hhblits の高速化のために 3 つの手法を提案する。うち 2 つは TSUBAME3.0 上での高速化であるが, 1 つは TSUBAME3.0 に依らない hhblits の高速化である。

また、AlphaFold の MSA 取得の高速化のために各 MSA 取得を非同期に実行する手法を提案する。

3.1 hhblits の並列化数の適正化

hhblits は OpenMP でいくつかの関数が並列化されている。2.3 章で言及した 4 つの関数のうち関数 `prefilter_db`, `alignment`, `perform_realign` の 3 つの関数は実行時間的に大きな関数が並列化されている。

AlphaFold では hhblits の並列化数は 4 に設定されている。これを計算ノードの物理コア数の 28 に変更する。

hhblits は I/O に実行時間が大きく左右されるプログラムである。一般には I/O ボトルネックの場合に並列化数を上げることは、意味がない場合が多い。しかし TSUBAME3.0 上ではデータベースが Lustre ファイルシステム上で Stripe されて保管されているため、並列化数を上げることで I/O 性能が改善される場合がある。

3.2 hhblits のアクセスするデータベースの配置の変更

AlphaFold 中の hhblits は BFD(1.7TB) と Uniclust30(86GB) の 2 つのデータベースを使用している。関数 `prefilter_db` では BFD の 1.7TB のうち 18GB ほどのファイルのみにしかアクセスしていない。そのため、このファイルを高速ストレージから計算ノードの DRAM 上にコピーする。tmpfs ファイルシステムである `/dev/shm` にディレクトリを作成し、関数 `prefilter_db` がアクセスするファイルをコピーする。残りのファイルはシンボリックリンクを用いて高速ストレージにリンクする。

tmpfs は使用した分だけ DRAM 容量を消費するため、使用できる DRAM 容量が $235 - 18(\text{BFD}) - 4(\text{Uniclust30}) = 213\text{GB}$ に減少する。

また、DRAM 上に配置するためのコピー時間は 15 秒ほどである。

3.3 hhblits のソートの高速化

並列化が行われていない関数 `mergeHitsToQuery` を高速化する。この関数のコールツリーの子孫に関数 `Filter2` があり、これが関数 `QSortInt` を呼び出している。

関数 `QSortInt` はクイックソートで実装されたユーザ定義のソート関数である。入力配列が T1050 の場合は、この関数だけで 900 秒ほどの実行時間を占めている。この関数呼び出しを c++ 標準のソート関数 `std::stable_sort` に置き換える。

この手法は TSUBAME3.0 に依存しない変更である。hhblits は Github 上で `hh-suite` というリポジトリで公開されており、この変更の Pull Request を提出済みである [21]。

3.4 AlphaFold の MSA 取得の非同期実行

AlphaFold の MSA 取得は図 4 のように同期的に実

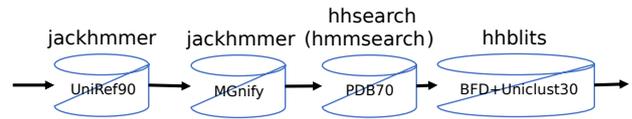


図 4 MSA 取得のパイプライン

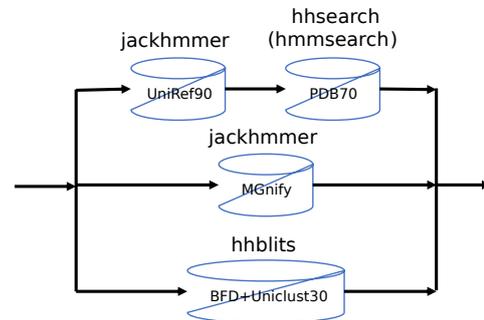


図 5 非同期化した MSA 取得のパイプライン

行している。各 MSA 取得ツール間の依存関係は `hhsearch(hmmsearch)` の実行が `jackhmmer(UniRef90)` の結果に依存しているのみであり、その他には依存関係は存在しない。これに着目して、図 5 のように MSA 取得を非同期に実行する。実装は `concurrent.futures.ThreadPoolExecutor` を用いて、実装した。

この手法によって、AlphaFold の MSA 取得全体の実行時間が一番大きいパスのみに依存するようになる。ただし、CPU、メモリ、I/O の必要量が増加する。また 3.1 章の手法を用いる場合、必要とされるコア数が最大 20 から 44 になる。

そのため、非同期化を行う前に比べて各 MSA 取得の実行時間が大きくなる可能性が考えられる。

4. 実験

4.1 実験環境

本実験では TSUBAME3.0 の計算ノードを 1 台用いる。AlphaFold は v2.1.1 を使用。実行の際は TSUBAME3.0 で提供されている `run_alphafold.sh` を用いて実行する。実行オプションは `-a 0,1,2,3 -max_template_date=2020-05-14` とする。

ただし、`run_alphafold.sh` は v2.1.0 での `preset`, `model_names` の廃止と `db_preset`, `model_preset` の追加という変更に対応していないため、可能な限り似た動作するように変更する。`run_alphafold.sh` では `preset` のデフォルトの値が `full_dbs` になっている。よって、`db_preset` は `full_dbs` を指定し、`model_preset` は指定せずデフォルトを使用する。(予測はモデル 5 つ)

また、AlphaFold では hhblits のログが捨てられているが、ファイルに出力するようにソースコードを変更する。

ベースラインは、以上の変更を加えた `run_alphafold.sh` を実行する。

提案手法の hhblits は 3.3 章で述べた通りに実装し、図 6

```
$ module li
Currently Loaded Modulefiles:
 1) cuda/11.5.0          3) gcc/10.2.0
 2) openmpi/3.1.4-opa10.10
```

図 6 hhblits をビルドする際の環境

表 4 実験に用いたデータセット

Target ID	残基数
T1024	408
T1029	125
T1032	284
T1041	242
T1047s1	232
T1047s2	365
T1053	580
T1061	949
T1087	93
T1095	665

の環境でビルドして使用する。

また、データベースには高速ストレージを使用しているが、その影響を同じにするために、同時に高速ストレージにアクセスすることがないようにした。さらに実験は3回ずつ行い、中央値を用いる。

4.2 実験データセット

実験に用いるデータセットとして、CASP14 で用いられたものから長さに偏りがないように適当なものを10個用いる(表4)。

5. 実験結果

ベースラインと提案手法の結果が図7である。提案手法では、MSA取得はjackhmmer(mgy_cluster), jackhmmer(uniref90) + hhsearch, hhblitsのうち一番時間がかかるツールだけに実行時間が依存するため、その他のツールの実行時間は表示していない。

AlphaFoldの実行時間が10配列のそれぞれの中央値で平均2分の1に短縮された。ベースラインの場合は、hhblitsが実行時間で1位か2位になっており、少なくとも4割程度を占めている。提案手法の場合はT1024, T1061を除いて、hhblitsが合計の実行時間に寄与していない。T1061でもhhblitsの実行時間はalphafold推論部の2分の1ほどになっており、hhblitsがボトルネックという状態にはなっていない。T1024については提案手法の段階でもhhblitsの実行時間が6000秒ほどある。

hhblitsのログから各配列の大まかなボトルネックを特定する。関数prefilter_dbの一部を含む、“Prefiltering database”から“HMMs passed 1st prefilter (gapless profile-profile alignment)”を含む行のログ間の実行時間に着目する。Score-Pによるプロファイリングの結果から、こ

表 5 手法1, 2の関数prefilter_dbの実行時間(秒)

	なし	手法1	手法2	手法1 + 2
T1024	1589	311	358	97

れは関数prefilter_dbの実行時間に近い値となる。T1024, T1061を除く配列は、この箇所が9割ほどの実行時間を占めている。そのため、これらの配列は関数prefilter_dbがボトルネックと言える。提案手法では、これらの配列のhhblits実行時間は隠蔽されるほど高速化されているが、これは手法1, 2の高速化の影響が大きい。

また、手法1, 2, 4によってhhblits以外のMSA取得ツールが扱える計算リソースが減少している。しかし、jackhmmer(mgy_cluster), jackhmmer(uniref90), hhsearchの両者の実行時間を見る限り、大きな差はなかった。

6. 各手法の寄与

提案手法の適用後でもhhblitsの割合が大きいT1024を用いて、各手法の寄与度を調査する。

まず手法4は、実験結果からhhblits以外のツールに影響がなかったため、単純に1番実行時間の小さいパスと2番目に小さいパスの実行時間が削減されると考えられる。hhblits以外のツールに絶対的に大きな実行時間の差は見られないため、hhblitsの実行時間が中程度の配列で効果的な手法となる。

手法3はhhblitsの関数mergeHitsToQueryに対する高速化である。関数mergeHitsToQueryに対してはこの手法以外で変更を行っていないため、関数mergeHitsToQueryのみに着目すれば良く、元の関数QSortIntの実行時間とおよそ等しい時間が短縮されている。この関数はデータベースにヒットする配列が多い入力配列で大きな実行時間を占める。そのような入力配列は全体の実行時間が比較して大きいいため、実行時間が大きい配列で効果的な手法となる。

関数prefilter_dbは手法1, 2で高速化される。表5に手法を適用した際の実行時間を示す。片方の手法を適用した時点で300秒台まで高速化されている。両方組み合わせることでさらに高速化されるが、T1024のhhblitsの実行時間から見ると5%以下の改善である。関数prefilter_dbが大きな割合を占める入力配列では、両方適用することでhhblitsが2倍以上高速化される可能性はあるが、手法4が存在するため、両方の適用は意味がない。ただし、高速ストレージのI/O状況によっては手法2のようにコピーすることは役立つ可能性がある。

手法1の関数prefilter_dbを除く関数alignmnet, perform_realignについて、関数alignmentはI/Oの影響を強く受けるが3000秒ほどかかっていた実行時間が1.3~2倍ほど高速化し、関数perform_realignは660秒から110秒に高速化した。またこれは手法3と同様に実行時間が大きい配列で効果的な手法となる。

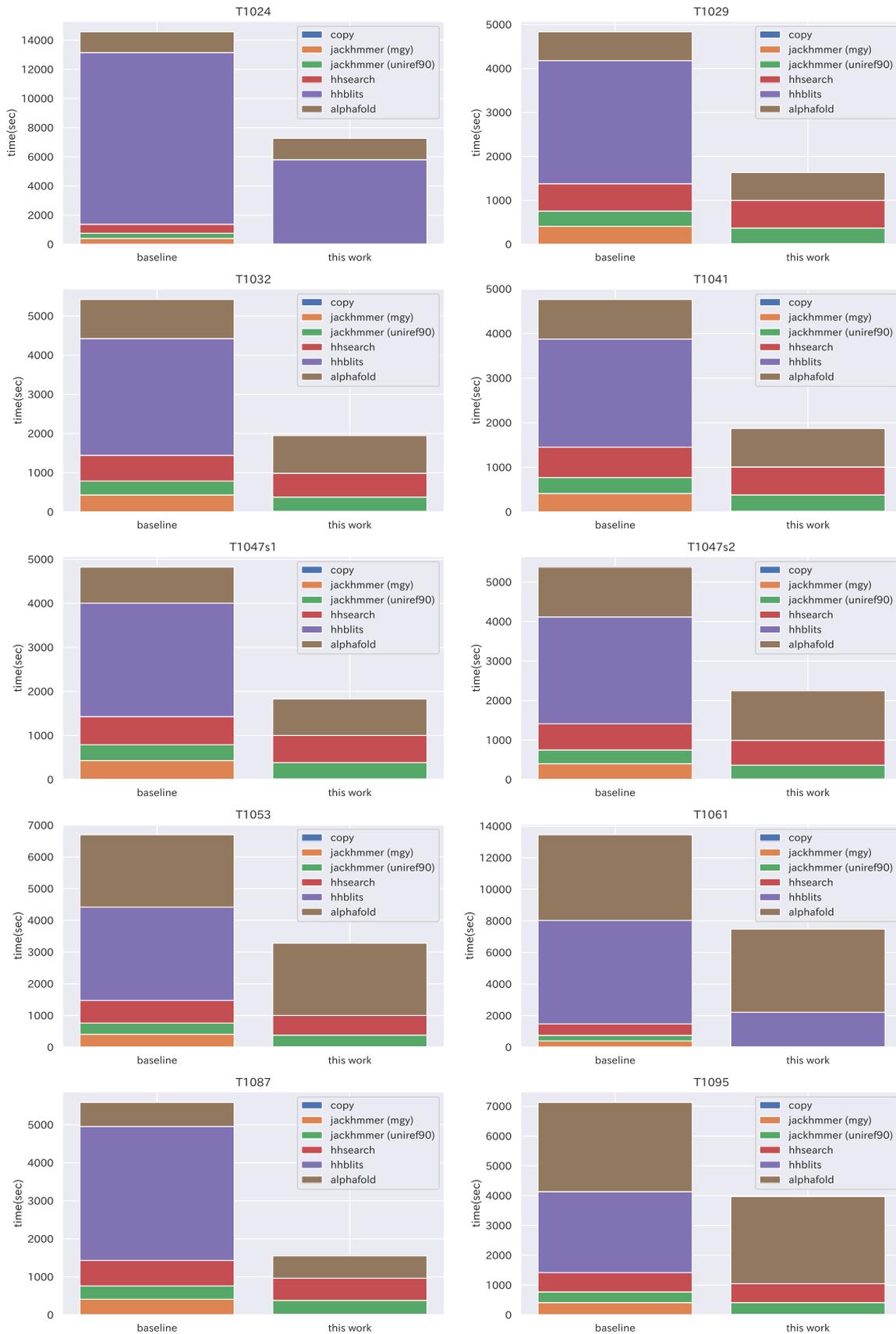


図 7 各配列の実行時間とその分布

I/O に強く影響を受ける 2 つの関数 (prefilter_db, alignment) では、手法 1 による高速化の寄与度が異なっている。

関数 prefilter_db では 5 倍ほど高速化されているが、関数 alignment では最大で 2 倍となっている。関数 alignment

については高速ストレージを使わずに計算ノードの SSD を用いることで 100 秒ほどまで高速化されるため、手法 1 適用後でも I/O ボトルネックである。手法 1 は共に I/O ボトルネックとみられる関数であっても、高速化の影響は関数によって異なる。

7. おわりに

本研究では、AlphaFold を高速に実行するために hhblits の高速化を行った。Score-P を用いたプロファイリングやソースコード解析を行い、ボトルネックとなった箇所を修正していった。修正箇所は並列化数の上昇、データベースの配置の変更、ソートの標準ライブラリ化である。「並列化数の上昇」はコア数が 4 以上かつ I/O 性能に余裕がある計算機に有用であり、「データベースの配置の変更」はリモートストレージがある計算機でローカルストレージとの I/O 性能に開きがある場合に有用である。最後の「ソートの標準ライブラリ化」は計算機による制限のない変更と考えられる。また、MSA 取得の非同期化を行い、AlphaFold の MSA 取得の高速化を図った。これにより、単一の入力配列に対する AlphaFold の実行時間が 2 分の 1 に短縮された。また多くの入力配列で hhblits の実行時間はボトルネックとは言えない段階まで高速化された。

謝辞 本研究は、東京工業大学のスーパーコンピュータ TSUBAME3.0 を利用して実施しました。また、本研究の一部は、国立研究開発法人日本医療研究開発機構 (AMED) 創薬等ライフサイエンス研究支援基盤事業 創薬等先端技術支援基盤プラットフォーム (BINDS) の課題番号 JP21am0101112 の支援を受けました。

参考文献

- [1] wwPDB Consortium. Protein Data Bank: the single global archive for 3D macromolecular structure data. *Nucleic Acids Res.* 47, D520–D528 (2018).
- [2] Mitchell, A. L. et al. MGnify: the microbiome analysis resource in 2020. *Nucleic Acids Res.* 48, D570–D578 (2020).
- [3] Steinegger, M., Mirdita, M. & Söding, J. Protein-level assembly increases protein sequence recovery from metagenomic samples manyfold. *Nat. Methods* 16, 603–606 (2019).
- [4] Moulton, J., Fidelis, K., Kryzhanovych, A., Schwede, T. & Topf, M. Critical assessment of techniques for protein structure prediction, fourteenth round. CASP 14 Abstract Book https://www.predictioncenter.org/casp14/doc/CASP14_Abstracts.pdf (2020).
- [5] Jumper, J., Evans, R., Pritzel, A. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* 596, 583–589 (2021). <https://doi.org/10.1038/s41586-021-03819-2>
- [6] Update license of AlphaFold parameters. · deepmind/alphafold@8173117 <https://github.com/deepmind/alphafold/commit/8173117130e6df8749ab7349722ec465666df548> (Accessed on 01/27/2022)
- [7] Martin Steinegger on Twitter <https://twitter.com/thesteinegger/status/1484165053402464259> (Accessed on 01/27/2022)
- [8] Release AlphaFold v2.0.1 · deepmind/alphafold <https://github.com/deepmind/alphafold/releases/tag/v2.0.1> (Accessed on 01/23/2022)
- [9] Alphafold2 公開のお知らせ — TSUBAME 計算サービス <https://www.t3.gsic.titech.ac.jp/node/454> (Accessed on 01/15/2022)
- [10] TSUBAME Computing Services on Twitter https://twitter.com/Titech_TSUBAME/status/1453594879007727617 (Accessed on 01/15/2022)
- [11] Johnson, L. S., Eddy, S. R. & Portugaly, E. Hidden Markov model speed heuristic and iterative HMM search procedure. *BMC Bioinformatics* 11, 431 (2010).
- [12] Suzek, B. E., Wang, Y., Huang, H., McGarvey, P. B. & Wu, C. H. UniRef clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics* 31, 926–932 (2015).
- [13] Mitchell, A. L. et al. MGnify: the microbiome analysis resource in 2020. *Nucleic Acids Res.* 48, D570–D578 (2020).
- [14] Steinegger, M. et al. HH-suite3 for fast remote homology detection and deep protein annotation. *BMC Bioinformatics* 20, 473 (2019).
- [15] Mirdita, M. et al. Uniclust databases of clustered and deeply annotated protein sequences and alignments. *Nucleic Acids Res.* 45, D170–D176 (2017).
- [16] TSUBAME3.0 ハードウェア・ソフトウェア仕様 <https://www.gsic.titech.ac.jp/sites/default/files/spec30j.pdf> (Accessed on 01/15/2022)
- [17] TSUBAME3 Lustre Seminar 公開用_20191028 https://www.t3.gsic.titech.ac.jp/sites/upload/TSUBAME3_Lustre_Seminar.pdf (Accessed on 01/15/2022)
- [18] 7. フリーウェア - TSUBAME3.0 利用の手引き <https://helpdesk.t3.gsic.titech.ac.jp/manuals/handbook.ja/freesoft/#alphafold> (Accessed on 01/15/2022)
- [19] Remmert, M., Biegert, A., Hauser, A. et al. hhblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat Methods* 9, 173–175 (2012). <https://doi.org/10.1038/nmeth.1818>
- [20] Knüpfner, Andreas, et al. “Score-P: A joint performance measurement run-time infrastructure for Periscope, Scalasca, TAU, and Vampir.” *Tools for High Performance Computing 2011*. Springer, Berlin, Heidelberg, 2012. 79–91.
- [21] Use std::sort instead of QSortInt by fuji8 · Pull Request #307 · soedinglab/hh-suite <https://github.com/soedinglab/hh-suite/pull/307> (Accessed on 01/27/2022)