

BLAST アルゴリズム全体のハードウェア化による高速化の検討

石川淑[†] 田中飛鳥[†] 宮崎敏明[†]

会津大学大学院コンピュータ理工学研究科[†]

1. はじめに

Basic Local Alignment Search Tool (BLAST) は最も有名なシーケンスアライメントツールの一つである。シーケンスアライメントとはタンパク質 (または DNA) 配列データベース (DB) 内のシーケンスと検索対象となるタンパク質 (または DNA) 配列 (クエリシーケンス) を比較し、配列同士の類似部分検索を行うものである。シーケンスアライメントは、生物学上の進化や遺伝子系図を調べる上で重要であることから、バイオインフォマティクス分野では欠かせない情報となっている。BLAST のハードウェア化による高速化手法は多く提案されてきた [1, 2]。その殆どは、一部の処理をホスト PC で行うというものであり、通信ボトルネックが生じていた。本研究では、BLAST アルゴリズム全体をハードウェア化することによりホスト PC との通信ボトルネックを解消し、高速化を行ったので報告する。

2. BLAST アルゴリズム

BLAST アルゴリズムは、前処理、Seeding (ステップ 1)、Ungapped extension (ステップ 2)、Gapped extension (ステップ 3) の 3 つのステップと traceback 処理からなる。前処理では、検索に用いる隣接ワード (Neighborhood word) と呼ぶ文字列を生成する。隣接ワードはクエリシーケンスと置換行列 (Substitution Matrix) を用いて生成される。置換行列とは、タンパク質を構成する 20 種のアミノ酸同士の類似度を数値化した表であり、一般に 20×20 の行列形式で表現される。ここでは、比較する 2 つの配列の類似度を数値化したものをスコアと呼ぶ。置換行列は複数存在するが、本稿では、BLOSUM50 と呼ぶ置換行列を使用する。まず、クエリシーケンスを k 文字 (通常、タンパク質配列では k=3、本稿でも k=3 を用いる) のクエリワードに分割する。つぎに、そのクエリワードの 3 文字を他のアミノ酸の 20 文字と 1 文字ずつ比較しスコアを計算する。3 文字のスコアの合計が閾値 T (通常 T=12) 以上となったワードが隣接ワードとなる。Seeding (ステップ 1) では、前処理で生成した隣接ワードを使用して DB シーケンス内の seed を探す。seed とは、DB シーケンス上で隣接ワードが一致した位置のことである。ステップ 2、すなわち Ungapped extension ではステップ 1 で見つかった seed を拡張開始点とし、一致点の拡張を行う。拡張は合計スコアが最大値から閾値 S (ユーザが指定) 下

がるまで行う。拡張されたシーケンスペアのうち HSP (High scoring segment pair) と呼ぶ組み合わせだけが次のステップ 3 の入力となる。ステップ 3 の Gapped extension では Smith-Waterman アルゴリズムという動的計画法に基づく手法を使用している。当該手法では、計算処理に 2 次元データ行列 (以下、スコア行列と呼ぶ) を用いる。図 1 は、HSP をスコア行列の行と列に対応させ、各行列要素 (以下、セルと呼ぶ) のスコア計算を行った結果である。1 つのセルは、上、斜め上、そして左のセルのスコアを使用して計算された値と、数値 0 の 4 つの値を比較し、その最大値をスコアとして該当セルに保持する。図 1 中の矢印は、どの位置のスコアを使用したのか (経路) を示している。traceback 処理ではステップ 3 で求めたスコア行列上で、最も高いスコアを持つセルからスコアが 0 のセルまでの経路を辿ることにより、最適なシーケンスを出力する。図 1 の例では、“TALILA” と “TAL-LA” が最適なシーケンスとなる。ここで、シーケンスペア中の“-”はギャップと呼ばれ、文字列のずれを示す。

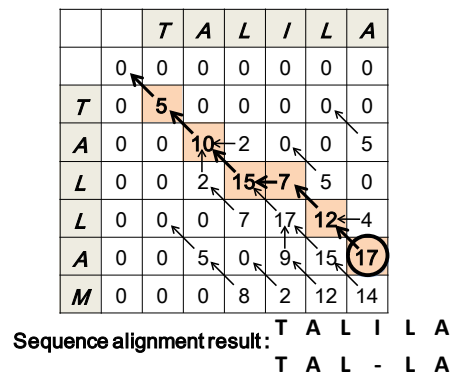


図 1 Gapped extension 処理と traceback 処理によって求められた計算結果と最適なシーケンス

3. 提案手法

従来、BLAST のハードウェア化は、前述のステップ 3、すなわち Gapped extension を中心に行われてきた。そのため、前処理である隣接ワードの生成や traceback 処理は、ホスト PC 上のソフトウェアで行う必要があり、PC とハードウェア間で通信ボトルネックが生じる可能性があった。本稿では、前述した traceback 処理を含む BLAST アルゴリズム全体をハードウェア化することを提案する。ここでは、BLASTP (タンパク質配列のシーケンスアライメントツール) を対象とする。図 2 に、Gapped extension 処理と traceback 処理を実行する提案回路を示す。Gapped extension 部内の計算には、Processing Element (PE) と呼ぶ簡単な処理を行うプロセッサを

An approach to accelerating BLAST algorithm using dedicated hardware

[†]Shizuka Ishikawa, [†]Asuka Tanaka, [†]Toshiaki Miyazaki
[†]Graduate School of Computer Science and Engineering,
 The University of Aizu

複数使用する. 図2のように複数のPE同士を線形接続して構成するPEアレイで並列処理を可能とし, ソフトウェア処理に比べ高速化を実現する. また, 各PEにはローカルメモリを接続し, 計算結果は, ホストPCに転送せずに, 当該メモリに格納する. ローカルメモリは $M \times N$ のサイズ ($M \geq$ クエリシーケンス長, $N \geq$ データベースシーケンス長)の2次元メモリとなっており, 各計算結果を図1で示したスコア行列に対応する位置に格納できる. ローカルメモリに格納される各計算結果は, 四つ組の情報(行列に対応したクエリ, DBシーケンス中の各1文字, どの位置のスコアを使用したかを示すparent(図1の矢印を数値で表したもの), 計算によって求められたスコア)からなる. PEアレイの計算が終了すると, クエリの終端文字を保持するPEから, 終了信号と最大スコア, 最大スコアが保存されているメモリアドレスが traceback部に送られ, 直ちに traceback処理が開始される. traceback部では, PEから届いた最大スコアのメモリアドレスを最初の読み出しアドレスとして, 2次元メモリにアクセスする. メモリから出力された四つ組データは traceback部に送られる. traceback部は, 四つ組データ中のparentを使用して次の読み出しアドレスを生成するとともに, 他のクエリ, DBの文字データは, そのまま出力する. 本操作を, ローカルメモリから読み出した四つ組データのスコアが0になるまで繰り返す. 上記の操作に於いて, traceback部から出力されたデータが最適なシーケンスとなる.

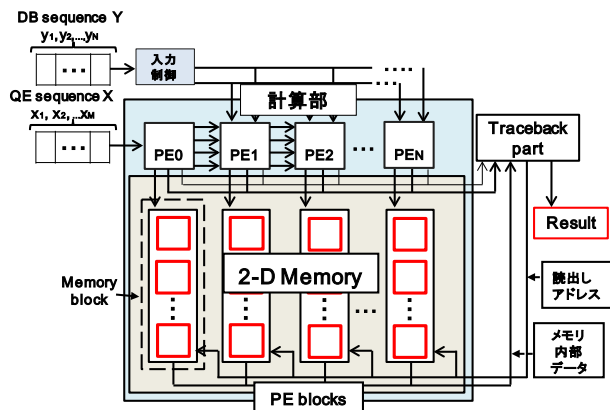


図2 Gapped extension処理と traceback処理を実行する提案回路

4. 全体構成

図3に提案回路の全体構成を示す. 提案回路は, Block 1からBlock 4の4つの部分から成る. 各部分は, 前述した隣接ワードを求める前処理(Block 1), ステップ 1, 2, 3(Block 2~4)に対応している. Block 4には, Gapped extension処理部に加え, traceback処理部も含まれる. 前処理部では, 入力されたクエリシーケンスを用いて隣接ワードメモリの内容を生成する. この隣接ワードメモリは, Block 2の処理で, 何度も参照される. Block 2では, DBシーケンスを3文字に分割し, 隣接ワードメモリに送り同一文字の一

致点(seed)を探す. Block 3の Ungapped extension処理では, パイプライン処理を効率的に行えるように文字列を左右方向に拡張するのではなく, 一方向へのみ拡張を行うようにしている. Block 3の処理で求められたHSPは時間調整用のFIFOを介して, 最終処理部であるBlock 4へ送られる. Block 4では, Smith-Watermanアルゴリズムに基づいた計算が行われ, 計算終了後, 最適なシーケンスを出力する.

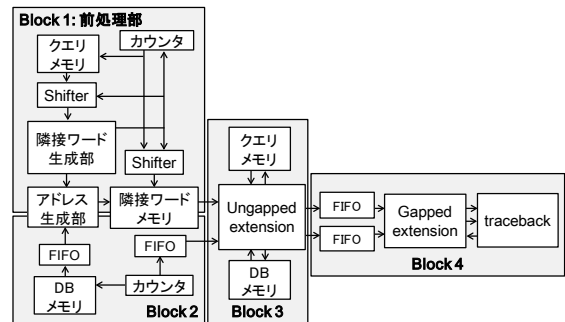


図3 BLAST処理を実行する提案回路の全体構成

5. 評価

動作確認のため, FPGAを用いて提案回路の実装を行った. 表1は, 図3に示した回路全体の実装結果であり, 表2は提案回路とソフトウェアBLASTの実行時間を比較したものである. 評価に用いたクエリシーケンス, DBシーケンス長はともに100である. 使用したFPGAはCyclone-IV E EP4CE115F29C7であり, 回路実装には同社の設計ツールQuartusII 10.1sp1を使用した. 提案回路の最大動作周波数は55.11 MHz, 必要となったクロックサイクル数は9383クロックであり, 実行時間は0.17 msであった. 比較対象として, ソフトウェアBLAST (Local BLAST 2.2.25)を使用した. 同一入力データに対して, その実行時間は135 msであった. よって, 表2に示すように, 提案回路はソフトウェアに比べて791倍の高速化を実現した.

表1 前処理部のFPGA回路規模

リソース	使用量	使用率
ロジックエレメント数	41,142	36 %
レジスタ数	3,656	3 %
メモリビット数	653,056	16 %

表2 ソフトウェアBLASTと提案回路の比較

	実行時間	高速化率
ソフトウェアBLAST	135.00 ms	1
提案回路	0.17 ms	791

6. おわりに

BLAST処理全体を実行するハードウェアアーキテクチャを提案し, ソフトウェアよりも790倍以上高速であることを示した. 今後は, PEに接続されているメモリをdual memoryとすることで, 計算処理と traceback処理を並列実行させ, さらなる高速化を目指す.

参考文献

- [1] K. Benkrid, Y. Liu, and A. Benkrid, "A Highly Parameterized and Efficient FPGA-Based Skeleton for Pairwise Biological Sequence Alignment," IEEE Trans. on VLSI Systems Vol. 17, No. 4, pp. 561-570, April 2009. [2] S. Kasap, K. Benkrid, and Y. Liu, "Design and Implementation of an FPGA-based Core for Gapped BLAST Sequence Alignment with the Two-hit Method," Engineering Letters, Vol. 16, Issue. 3, No. 25, (EL_16_3_25), August 2008.