

RECOT: 次世代シーケンサの比較ゲノムに向けた ゲノム座標の変換

伊澤 亜紀子^{†1} 瀬々 潤^{†2}

次世代シーケンサの台頭により多様な種のゲノムが読まれ、それらを利用した RNA-seq や ChIP-seq などの実験が行われている。この実験結果を種間比較するには、次世代シーケンサ用のブラウザは複数のゲノムが扱えず、比較ゲノムブラウザはリードを気軽に追加できない問題点がある。本研究ではこれらの問題点を解消するため、リードの位置や変異情報を異種のゲノム情報に対応させ、次世代シーケンサ用ブラウザで可視化可能になるよう変換する RECOT を開発した。

RECOT: A tool for the genome coordinate transformation of next-generation sequencing reads between related species

AKIKO IZAWA^{†1} and JUN SESE^{†2}

The whole-genome sequences of many non-model organisms have recently been observed. Based on the genome sequences, experiments using high-throughput sequencers such as RNA-seq and ChIP-seq have been performed, and the experiments have been compared between related species. Although these comparisons require changes of the genome coordinates of the reads between the species, current software tools are not suitable. In this paper, we introduce a set of programs, called REad COordinate Transformer (RECOT), which transform the coordinates of short reads from the genome of the originally observed species to that of another species after aligning the gene/genome sequences between the species.

^{†1} お茶の水女子大学
Ochanomizu University

^{†2} 東京工業大学
Tokyo Institute of Technology

1. はじめに

DNA を読む技術が急速に発展し、新型のシーケンサ (以下 NGS) が生まれたことで、多様な生物のゲノム配列が解読されている。これにより、今までゲノム既知の種限定で行われてきた RNA-seq¹⁾ や ChIP-seq²⁾ などの NGS を使った実験が多様な種で可能になった。これらの実験を比較することで、進化プロセスや生物メカニズムについての解明が期待できる。

一般に NGS を用いた実験では、ゲノムに配列断片 (以下リード) をアラインメントし、対応づける。近縁種同士の比較には、種間でリードをアラインメントし、対応づける必要がある。しかし現在、大量のリードを他種のゲノム配列に対応づけることに適したソフトウェアがない。そして、複数種の NGS による実験結果を同時に可視化し、種間で比較することは、現在公開されているソフトウェアでは容易ではない。このような問題点に対し本手法では、2 種で読まれた NGS のリードのアライメント結果を同時に表示し、比較可能にするためのツール群 REad COordinate Transformer (RECOT) を開発した。また、RECOT は NGS の計算結果で最も一般的な形式である SAM 形式を出力する。これにより、ユーザが使い慣れた可視化ソフトを選ぶことが可能である。

2. 提案手法

本研究では、リードが読まれた種を original species と呼び、リード位置の変換対象種を target species と呼ぶ。手法の概観を図 1 に示す。RECOT は 4 つの入力ファイルを入力とする。(1)NGS から読まれた original species のリード (2)original species のゲノム配列 (3) 遺伝子情報 (4)target species のゲノム配列である。出力は (1) の配列の (4) 上でのアラインメント結果である。

本手法では、2 つのアラインメント結果が必要になる。1 つは (1) と (2) のアラインメント結果。そして、(2) の遺伝子配列と制御領域の遺伝子周辺領域の配列を、(4) にアラインメントした結果である。前者は BWA や Bowtie、後者は GMAP などが利用可能である。遺伝子周辺領域の配列は提供されないことも多いので、ゲノム配列と遺伝子領域から、遺伝子周辺領域の配列を得るスクリプトも作成した。図 1 は、これらのアラインメントした例である。SAM 形式は塩基の対応を CIGAR で表す。図 1(B) では 3M3N12M となる。M はマッチ又はミスマッチ (以下マッチ) を示し、D は欠失、N はギャップを指す。

しかし、相同性のある遺伝子により、1 つの target species の領域に複数の遺伝子がアラ

インメントされ、遺伝子間の対応付けで混乱を招くかもしれない。このような問題を避けるため、同じ領域に複数の遺伝子がアラインメントされた場合、その複数の遺伝子から1つを選択するスクリプトを作成した。このスクリプトでは以下の2つの選択方法が用意されている。(1) ユーザ指定の遺伝子の対応表により、対応に優先順位をつける。(2) アラインメントスコアが最も高い遺伝子配列を優先する。

最後に2つのアラインメント結果から、リードを target species にアラインメントしていく。今までに計算した2つのアラインメント(1)-(2)と(2)-(3)の結果をもとに、(1)-(2)間のCIGARを(2)-(3)間のCIGARに対応づける。そのために(1) target species のゲノム配列にアラインメントされた original species の遺伝子と、その遺伝子の領域にアラインメントされたリードの対応関係を表す表を作成、(2) この対応表を使用し、リードのCIGARの書き換えをおこなう。図1(A)は、original species の遺伝子配列とリードの、図1(B)は original species の遺伝子配列と target species のゲノム配列のアラインメントを示す。また、図1(A)(B)の遺伝子配列は同一のものである。

3. 実行結果

NGSの実データに本手法を適用し、有用性を検証するため、*D. simulans*(ERR020078)と

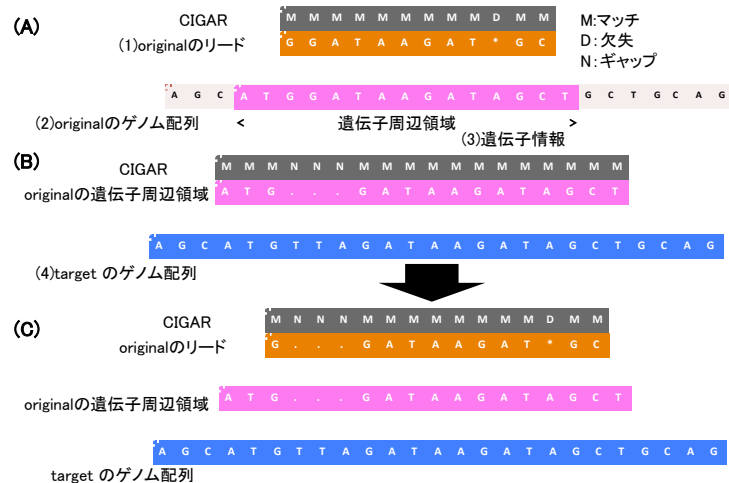


図1 提案手法の概観
Fig.1 Overview of Proposed Method

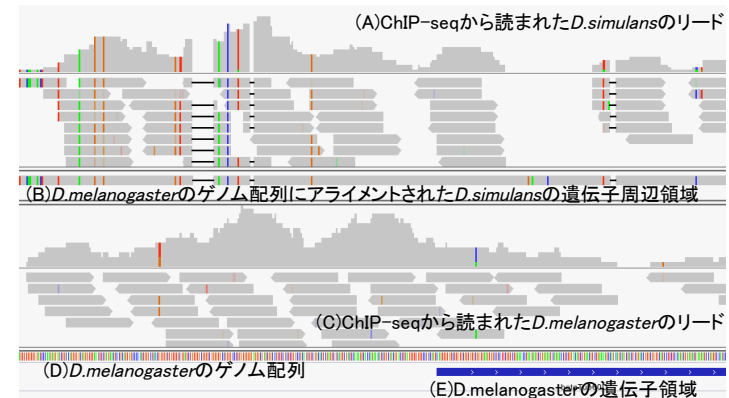


図2 実データによる実行結果
Fig.2 Result using a Real Dataset

キイロシヨウジョウバエ (*D. melanogaster*(ERR020066)) の ChIP-seq の比較をおこなった。これらは約 500 万年前に分岐した近縁種である。結果の可視化にはゲノムビューア Integrative Genomics Viewer(IGV)³⁾を使用した。図2(A)は、*D. simulans* のリードを RECOT により、*D. melanogaster* にアラインメントした結果である。図2(B)は、*D. melanogaster* の遺伝子周辺領域を *D. melanogaster* のゲノムにアラインメントした結果である。図2(C)は、*D. melanogaster* のリードを、*D. melanogaster* のゲノムにアラインメントした結果である。図2(D)は *D. melanogaster* のゲノム、図2(E)はその遺伝子である。2種のヒストグラムを比較すると、発現制御部位の変化を読み取ることができる。このように本手法は種間を比較するために有用である。

参考文献

- 1) Wang Z. *et al.*: RNA-Seq: a revolutionary tool for transcriptomics. *Nature Reviews Genetics*, Vol.10, pp.57-63 (2009).
- 2) Jothi *et al.*: Genome-wide identification of in vivo protein-DNA binding sites from ChIP-seq data. *Nucl Acids Res*, Vol.36, pp.5221-5231 (2008).
- 3) Robinson *et al.*: Integrative Genomics Viewer. *Nature Biotechnology*, Vol.29, pp. 24-26 (2011).