

総 括 研 究 報 告 書

課題番号： 28-12

課題名： 新規ヒトゲノム参照配列 GRCh38 および日本人基準配列を活用した
ゲノム変異診断主任研究者名（所属施設） 国立成育医療研究センター
（所属・職名） システム発生・再生医学研究部・
組織工学研究室長 岡村 浩司

（研究成果の要約） 小児に見られる異常の大部分はゲノムやエピゲノム変異が関わる疾患で、それらの診断に次世代シーケンサ(NGS)は欠かせない解析機器となっている。得られる配列長は従来のシーケンサと比べると短く 100 塩基程度であり、参照配列と呼ばれる既知のゲノム塩基配列と比較することで結果を得るため、参照配列の質が結果を大きく左右する。本センターでは 2009 年に公開された GRCh37 に基づくデータを一貫して使用しており、これまで特に問題は起きていないものの、新規参照配列 GRCh38 が発表されて 4 年以上が経過し、有用性が示され、かつ他機関との整合性を考慮すると移行する必要がある。このアップデートは NGS 登場後初めて、かつ複雑な作業である。センター全体にとって重要な事案ながら、これまでどの研究チームからも GRCh38 に関する報告はない。本研究では使用頻度の高いエクソーム解析を取り上げ、GRCh38 に対するマッピング、変異検出、アノテーションを行うパイプラインを再開発し、これまでと同等の変異検出ができることを確認し、特にミトコンドリアの変異検出に貢献することができた。また、これらゲノム配列データを深層学習することで、自動的に NGS データベースを構築する方法を編み出した。

1. 研究目的

大人とは異なり小児に見られる異常のほとんどはゲノムやエピゲノム変異が関わる疾患で、次世代シーケンサが欠かせない研究や診断のための機器となっている。次世代シーケンサが登場した当時、利用可能なヒトゲノム参照配列は NCBI Build 36 に基づく hg18 であったが、2009 年になって GRCh37 または hg19 と呼ばれる国際参照配列が発表され、次世代シーケンサやマイクロアレイなどのデータ解析に広く用いられてきた。現在でも特に断りがなければ、変

異等の染色体上の位置は GRCh37 に基づいた数値で表され、国内外の研究者、医師、カウンセラー間で問題なく連絡を取ることができる状況になっている。しかしながらこれらの参照配列は、今日に見られる次世代シーケンサの発展を視野に入れて用意されたものではなく、その後、マッピングと変異検出の正確さを上げるためにデコイ配列や alternate 配列が考案され、2013 年末に新規 GRCh38 が発表されるに至った。前バージョン GRCh37 の普及度合が高く、新規参照配列への切り換えはほとんど進んで

こなかったが、発表から4年以上が経ってその有用性が認められつつある。さらに1000人ゲノムプロジェクトによって日本人を含む各ヒト集団が持つ配列多様性が明らかになり、国際的に統一された単一参照配列の利用では問題が残ることも指摘され、国内でも主任研究者を含めいくつかの研究グループが日本人基準配列の決定に向け連絡を取り合っている。

本センターはこれまで一貫してGRCh37に基づいたデータ解析を行っており、新規参照配列への移行作業は全く行われていない。新しい配列の有用性は明らかで、かつ他機関との整合性の観点からも、新規参照配列に移行する必要があるが、この作業はデータベースを一つ入れ換えれば済むという単純なものではなく、解析パイプラインに組み込まれている個々のソフトウェアに対し設定と動作確認を行う必要があり、センター内の多くの部署で、また計算機の数だけいずれ必要になる作業でもあり、啓蒙も含め、イニシアチブを取って準備を進める必要がある。

まずは最も使用頻度の高いエクソーム解析を取り上げ、マッピング、変異検出、アノテーションを行う体制を整え、最終的には全ゲノム解析、DNAメチル化解析、クロマチン修飾解析、トランスクリプトーム解析、RNAメチル化解析等にも対応させたパイプラインソフトウェアの開発と整備を行う。また、近年、多くの状況で活用され始めた深層学習の技術を取り入れ、これらNGSデータを自動的に振り分けてデータベースを構築する仕組みについても検討を行う。

2. 研究組織

研究者	所属施設
岡村 浩司	国立成育医療研究センター
片桐 沙紀	お茶の水女子大学

3. 研究成果

次世代シーケンサから得られる配列データは従来のシーケンサから得られるものと比べると短く、解析方法は既知配列との比較が基本となる。通常は参照配列へのマッピングにより比較が行われるが、これはエクソーム解析、全ゲノム解析に限らず、DNAメチル化、クロマチン修飾のようなエピジェネティクス解析、さらにはトランスクリプトーム解析等においても同様である。小児や周産期疾患を中心とした研究を行っている本センターにおいてはエクソーム解析を行うことが最も多く、また所有する計算機クラスターがエクソーム解析を念頭に設計されたものであるため、初年度はエクソーム解析を最初のモデルケースとして取り上げることとした。

(1) GRCh38は2013年12月に発表されて以来、12のパッチがリリースされている。また、いくつかの研究機関が、研究者の便宜を計り、独自に改変した配列を公開しており、本研究課題においては欧州バイオインフォマティクス研究所(EBI)、カリフォルニア大学サンタクルーズ校(UCSC)、東北メディカル・メガバンク機構(ToMMo)から公開されているデータセットをダウンロード、比較検討することから研究を開始した。EBIに対し、UCSCではY, R, K, M, W, S, Bといった塩基が全てNで表記され、マッピングに利用しやすいことを確認した。また特に第5, 14, 19, 21, 22番染色体およびY染色体においてUCSCでは多くのNがT, C, A, Gの塩基で表記されていた。ToMMoの配列はUCSCの配列を基にしていると思われる。また国際参照配列に対し日本人ゲノムにおいて挿入により存在する配列を集約した配列をデコイ配列に加えたdecoyJRGを含んでおり、特に日本人のゲノム配列解析に有用であると判断され、ToMMo GRCh38を新規参照配列として利用することとした。

(2) GRCh38とはいえ、ミトコンドリアDNAの配列は統一されたものではないことが確認された。NC_012920が広く参照配列として利用されているようであったが、3107番目にNが

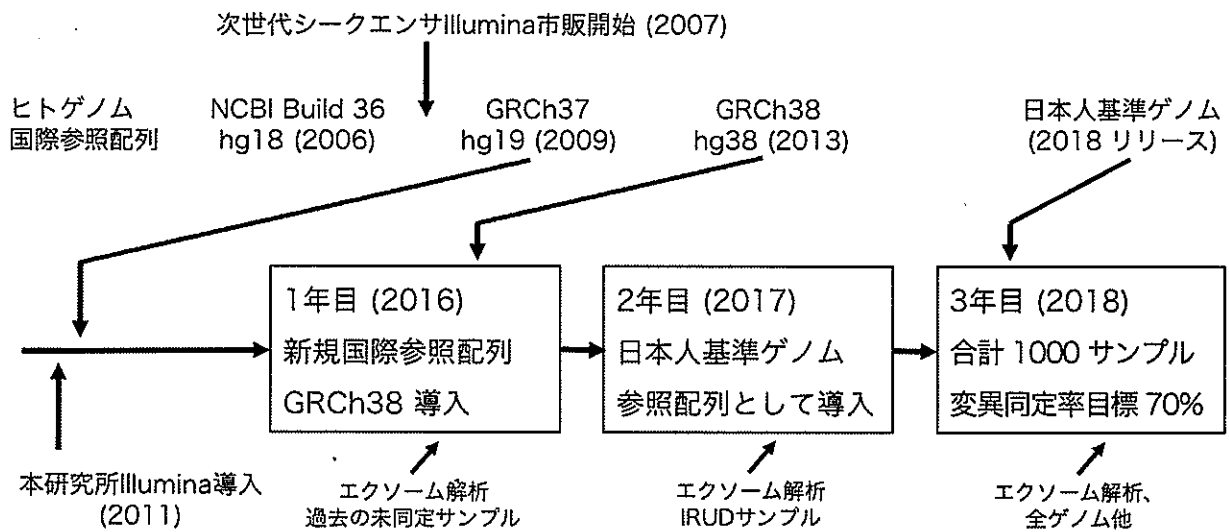


図 1. 参照配列 GRCh38 開発スキーム

入っていることが確認された。ToMMo が採用しているミトコンドリア DNA の配列は NC_012920 とは異なるが、主任研究者が別事業における「日本人由来細胞の同一性・ゲノム安定性評価法の確立」研究において決定した健康日本人 1 名 JOM005 のミトコンドリア DNA 配列 16,570 bp ときわめて類似していることが分かった。そこでミトコンドリア DNA の参照配列については、本研究所にて配列決定が行われた JOM005 を採用することとした。ToMMo の配列とは似ているものの、NC_012920 に対し 310 番目に C の挿入があるためそれ以降の座標がずれてしまうため、アノテーションデータを作り直す必要が生じた。そのため、dbSNP、ANNOVAR のデータベースを改変した。このツールを用いて得られたデータは、2017 年度に倫理審査申請を行った「ミトコンドリアの遺伝形式に関する研究」の基礎データとなった

(3) 全ゲノムに対し、エクソームにおいてはバイト設計領域周辺に限定した解析が必要となる。キャプチャーキットとしてはこれまで、Agilent SureSelect、Illumina TruSeq Exome Kit、Illumina Nextera Rapid Capture を採用したが、どれも GRCh38 に基づくデータが提供されていない。またそれぞれのキットに複数のバージョンが存在し、カバーしてい

る領域もそれぞれに若干異なり、データ処理を複雑にしている。そこで今回、これら既存のキャプチャーキットを全てカバーするような広い領域を解析対象として統一することとし、UCSC liftOver を利用して hg19 に基づくデータを GRCh38 に変換して作業を進めた。領域長は合わせて約 150 Mb となった。

(4) パイプラインは希少・未診断疾患イニシアチブ IRUD のために整備したスクリプトを基に、今回の作業においては、マッピング (BWA)、ヴァリントコール (Picard および GATK)、アノテーション (ANNOVAR) 等の各ツールはできる限り最新版を導入した。リードのトリミング、マッピング、リアラインメントによりサンプルごとの BAM 作成、さらにマルチサンプルコーリングによる VCF 作成までを行う 2 つのパイプラインを完成させ、目標のうち最も重要な部分を達成することができた。

(5) インハウス健康者コントロールとして用いている 2 サンプル、および主任研究者が筆頭著者として *Sci. Rep.* 誌にて発表した色素性乾皮症患者由来の 1 サンプル (Okamura *et al.* 2016)、合計 3 サンプルをモデルとして従来の hg19 および GRCh38 の両方で解析を行った。マッピング率は 96.66% から 96.54% とわずかながら下がり、ヴァリントコール数はそ

れぞれ 179,061 および 178,581 という結果が得られた。色素性乾皮症の原因である *XPA* 遺伝子のホモ変異は問題なく検出されており、診断に利用可能なパイプラインを完成することができた。

(6) アノテーションについては ANNOVAR および HGMD が GRCh38 に対応していることを確認したが、国内 5 拠点によるエクソーム解析から得られている日本人アレル頻度データ、ToMMo による健常日本人 2000 人全ゲノム解析から取得したアレル頻度データ等はどれも非対応で、位置情報を変換したデータベースを独自に作成する必要があり、作業を継続している。2016 年度および 2017 年度の日本分子生物学会年会において調査を行なったが、GRCh38 を利用した報告は見つけられず、国内において hg19 からの移行がほとんど進んでいない現状が分かった。しかしながら、次世代シーケンサから得られたデータが組み込まれている GRCh38 の優位性は明らかであり、日本人基準配列の活用も含め、本センター発の、エクソーム解析に限定しない一般的な NGS 解析環境を早急に完成させ、小児未診断疾患への応用とともに、広く有用性を示したいと考えている。なお、本課題により作成したヒトゲノム参照配列は GRCh38s1 として以下のサイトより一般公開している。

<http://epigenetics.nrichd.ncchd.go.jp/ge>

nomics/

(7) この一年間、ビッグデータと相まって社会のあらゆる分野で深層学習(ディープラーニング)が取り入れられ、大きな成果を上げている。医学や生命科学の分野においてもこういった人工知能の活用が望まれており、本課題のテーマとなっている参照配列について検討を進めた。その結果、参照配列に加えてこれまでに取得されデータベースに登録されている塩基配列データを用いて、NGS から得られる FASTQ データを全ゲノム、全エクソーム、RNA-seq、ChIP-seq 等に分類し自動的にデータベースを構築するフレームワークを提案することに成功した。この結果は日本分子生物学会年会の口頭発表に採択され(4P2T18-01)、さらに本研究所の定例セミナーにおいて発表を行った。病院の職員も含め、把握できる限りこれまでの定例セミナーの中で最も多い聴衆を集め、深層学習に対する期待と関心の高さがうかがわれた。

4. 研究内容の倫理面への配慮

本課題においては、倫理委員会により承認を得てシーケンシングされたサンプルのみを解析対象とした。