

ハイブリッドクラウド構築と ゲノム情報解析の効率的な運用に関する研究



【課題番号】 jh210018

【研究代表者名】 長崎 正朗

京都大学学際融合教育研究推進センター スーパーグローバル
コース医学生命系ユニット
京都大学大学院医学研究科附属ゲノム医学センター

本研究課題のメンバ全体

研究課題代表者・副代表者

【京都大学】

長崎 正朗 山田 亮

申請課題の全共同研究者

【京都大学】

松田 文彦 山口 泉 川口 喬久 稲富 雄一

深沢 圭一郎 関谷 弥生 浅倉 章宏 寺岡 凌 男澤 良子

Olivier Gervais Wang Yen Yen 橋本 洋希

【情報通信研究機構】

村田 健史

【東京大学】

関谷 勇司

【東京大学情報基盤センター】

埴 敏博

【九州大学】

大川 恭行 前原 一満 南里 豪志

【研究目的】

ヒトゲノム情報についてシーケンス技術の開発により爆発的に出力される情報が増えてきている。これらの情報について、情報量の増加とともに適切な計算環境において計算を行うこと、また、それらの計算結果を複数拠点にバックアップを持つなどの運用が必要となる。

1つの拠点では、上の目的を達成することが困難な状況となっており、オンプレ、国内のスーパーコンピュータシステム、また、商用のクラウド環境の各々において、転送のコスト、費用、セキュリティなど総合的に勘案をして運用を行う必要がある。

昨年度は、有償で追加利用を行った東京大学および京都大学の計算資源も含め、**約5,000検体の全ゲノムリファレンスパネルの構築を進めた**。その中で、各拠点での計算機資源の特徴を考慮し、解析パイプラインの各ステップを試行錯誤しつつ実行を進めた。**本年度は、そこで得られた知見に基づき、有償で追加資源を確保することを並行しつつ、R3に利用可能となる合計約10,000検体の全ゲノムリファレンスパネルの構築を目標に進めていく。**

そこで、当研究においては、

1. 複数拠点間にわたる計算資源、ストレージを効率的に運用するにおいて出てくる課題を整理しつつ
2. 仮想環境や大規模電算機資源上でゲノム情報の解析パイプラインの実装を行い
3. 円滑に上の一部の情報（数百検体を予定）について試験的に拠点間転送と2のパイプラインを用いたデータ解析を行う

ことを目的とする。

特に、昨年度では十分に対応ができなかった拠点間のデータ転送について重点的に整備を進める。

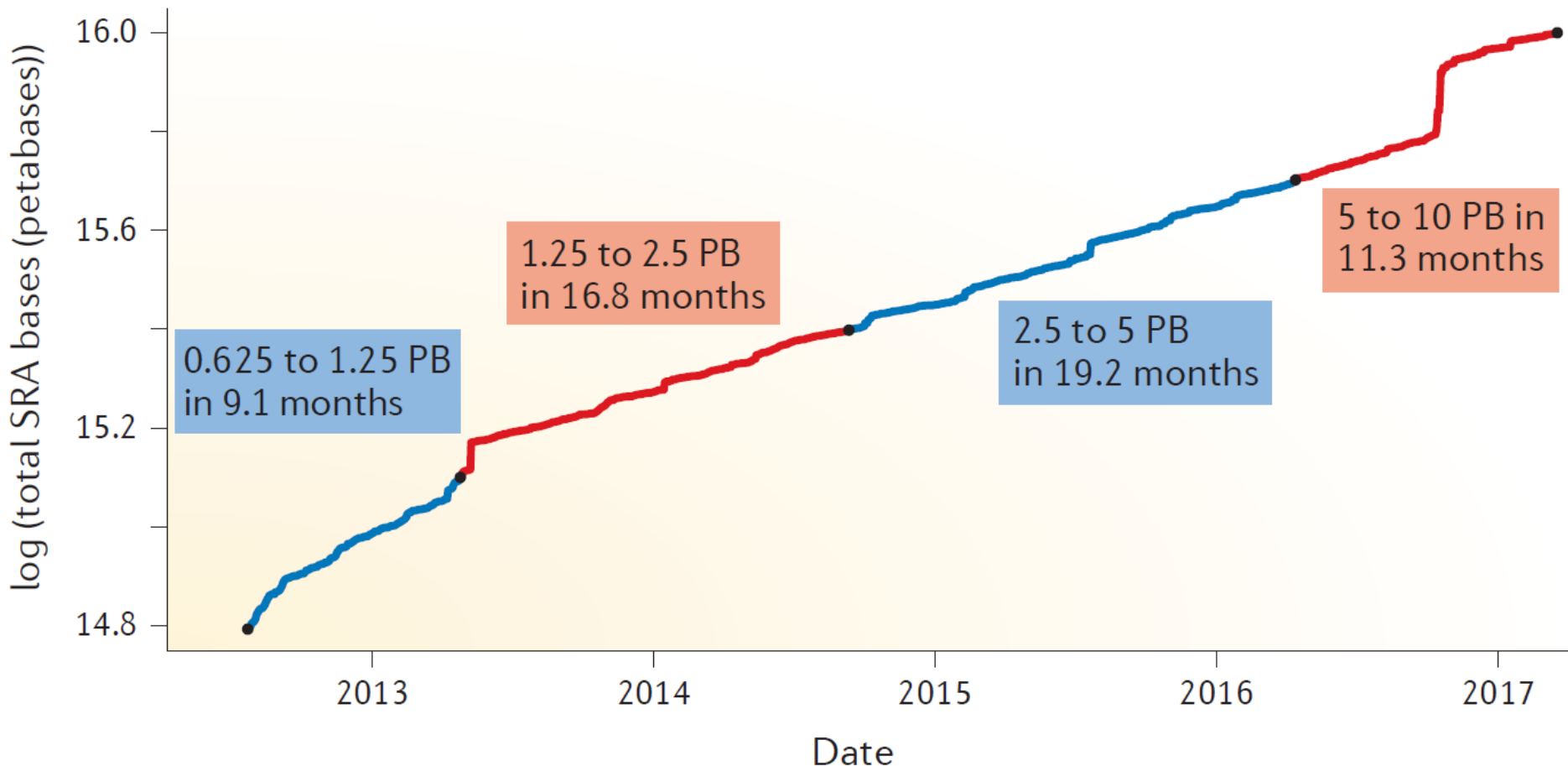
当拠点公募型共同研究として 実施する必要性

昨年度5,000検体の規模の解析を行ったが、海外の事例にあるように数万以上の規模が想定されている（次スライド参照）。当センターにおいても1万人規模の検体が収集されている。そのような規模において提案されてくるバイオインフォマティクス手法などによる解析が必須である。

また、それらの解析によって、シーケンスされた生データから新たな疾患のリスク要因が同定されることが想定されている。

そのため、今回の申請において、**各拠点でどのような解析を行うことで効率的に運用ができるか、また、将来的な情報量の増加に対応するか実際に設計・運用を行うことで検討を進めることが必要である。**

ヒトゲノムシーケンス情報の増加の様子

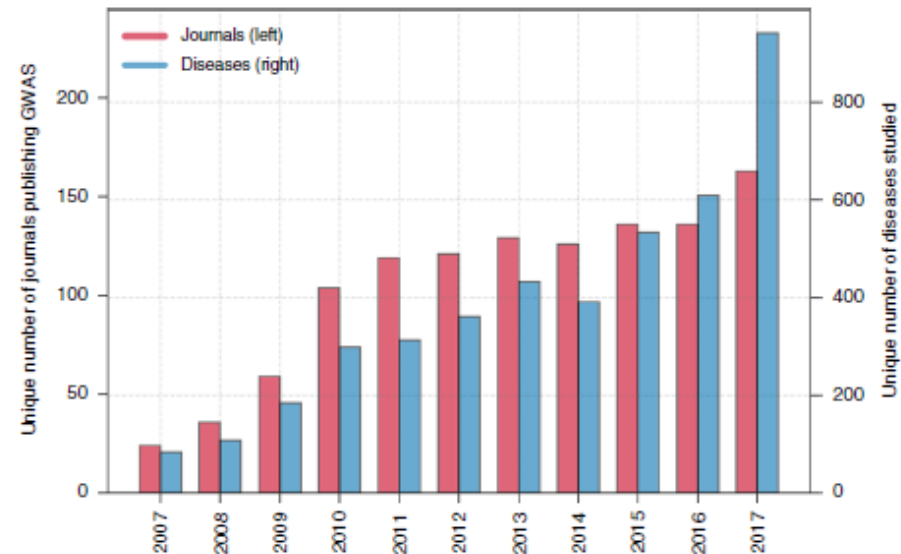
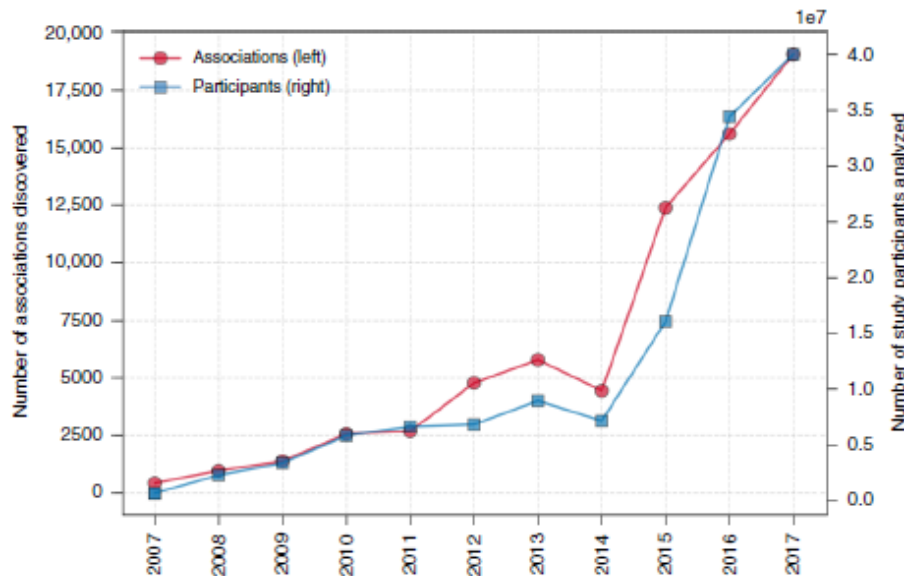
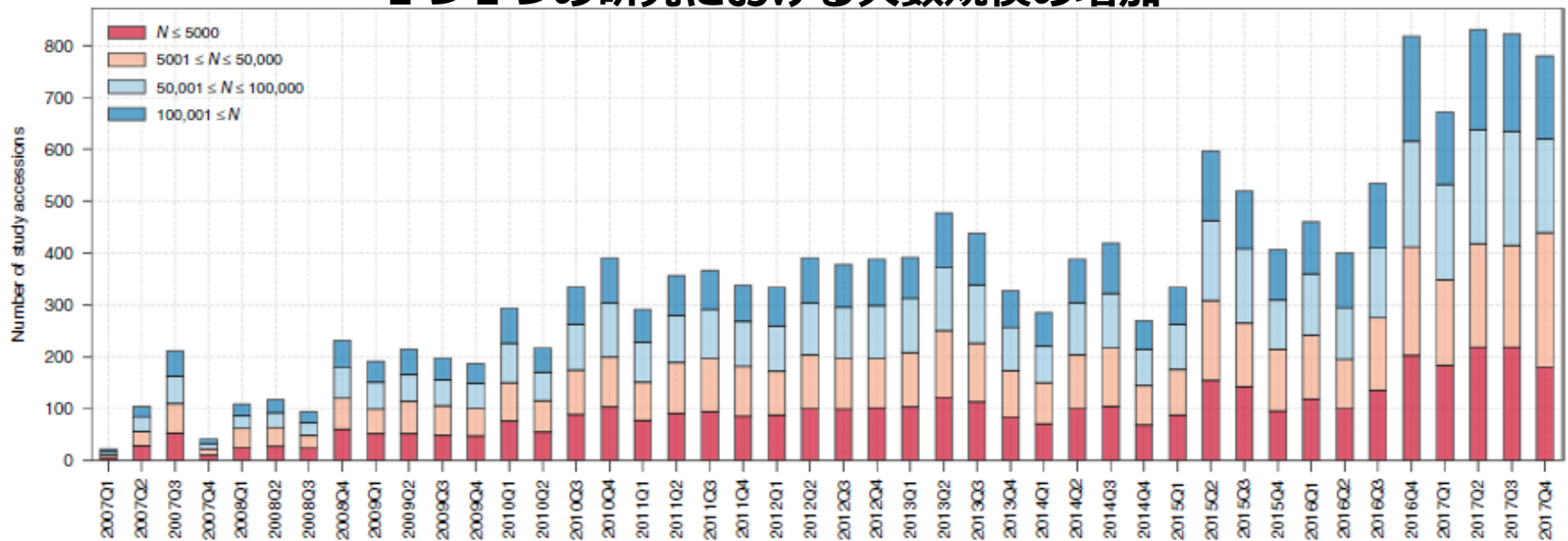


【図1】 ヒトゲノムシーケンス情報の増加の様子

米国のヒトゲノムシーケンス情報の国際データベース (Sequence Read Archive: SRA) に登録されているゲノム情報の総データサイズ (テラバイト) の推移 (2020年時点で12PB)

おおよそ1検体あたり30Gとして約40万人のゲノム情報が保存されていることになる。

1つ1つの研究における人数規模の増加



【図2】 1つ1つの研究における人数規模の増加の様子 研究プロジェクト毎の解析対象人数が大規模化するとともにまた数が増えていることがわかる。

【研究計画】

課題 1) 複数拠点間を効率的に運用できるハイブリッドクラウドシステムの設計と運用（長崎、山田、松田、関谷、深沢、村田）

課題 2) シークエンサから取得された情報および解析パイプラインの各ステップの入出力データの拠点間通信の最適化検討（大川、長崎、深沢、村田、塙、関谷）

課題 1) 複数拠点間を効率的に運用できるハイブリッドクラウドシステムの設計と運用

各拠点担当者とその役割

(拠点 1) 京都大学 ゲノム医学センター

【担当者】長崎正朗 全体統括 (分担) 山田亮、松田文彦 他

【役割】全ゲノム情報のハイブリッドクラウドにおける効率的な解析の設計
オンプレ、各電算資源間の効率的な解析パイプラインの構築 (後スライド参照)

(拠点 2) 京都大学 メディアセンター

【担当者】深沢圭一郎

【役割 1】京都大学と他拠点計算機資源との効率的なデータ分散及び拠点間転送支援

【役割 2】京都大学と他拠点とのSINET5を用いたパブリッククラウドへのVPN接続管理

(拠点 3) 東京大学

【担当者】関谷勇司

【役割】クラウド実装におけるアドバイス、また、試験環境の整備

(拠点 4) 情報通信研究機構

【担当者】村田健史

【役割】拠点間的高速データ転送技術提供と評価

課題 2) シークエンサから取得された情報および解析パイプラインの各ステップの入出力データの拠点間通信の最適化検討

各拠点担当者とその役割

(拠点 1) 京都大学 ゲノム医学センター

【担当者】長崎正朗 全体統括

【役割】拠点間的高速データ転送基盤整備と評価（拠点 1 担当）および全体評価

(拠点 2) 京都大学 メディアセンター

【担当者】深沢圭一郎

【役割】拠点間的高速データ転送基盤整備と評価（拠点 2 担当）

(拠点 3) 東京大学

【担当者】埜 敏博・関谷勇司

【役割】拠点間的高速データ転送基盤整備と評価（拠点 3 担当）

(拠点 4) 情報通信研究機構

【担当者】村田健史

【役割】拠点間的高速データ転送技術提供および運用支援

(拠点 5) 九州大学 生体防御医学研究所

【担当者】大川恭行 前原一満

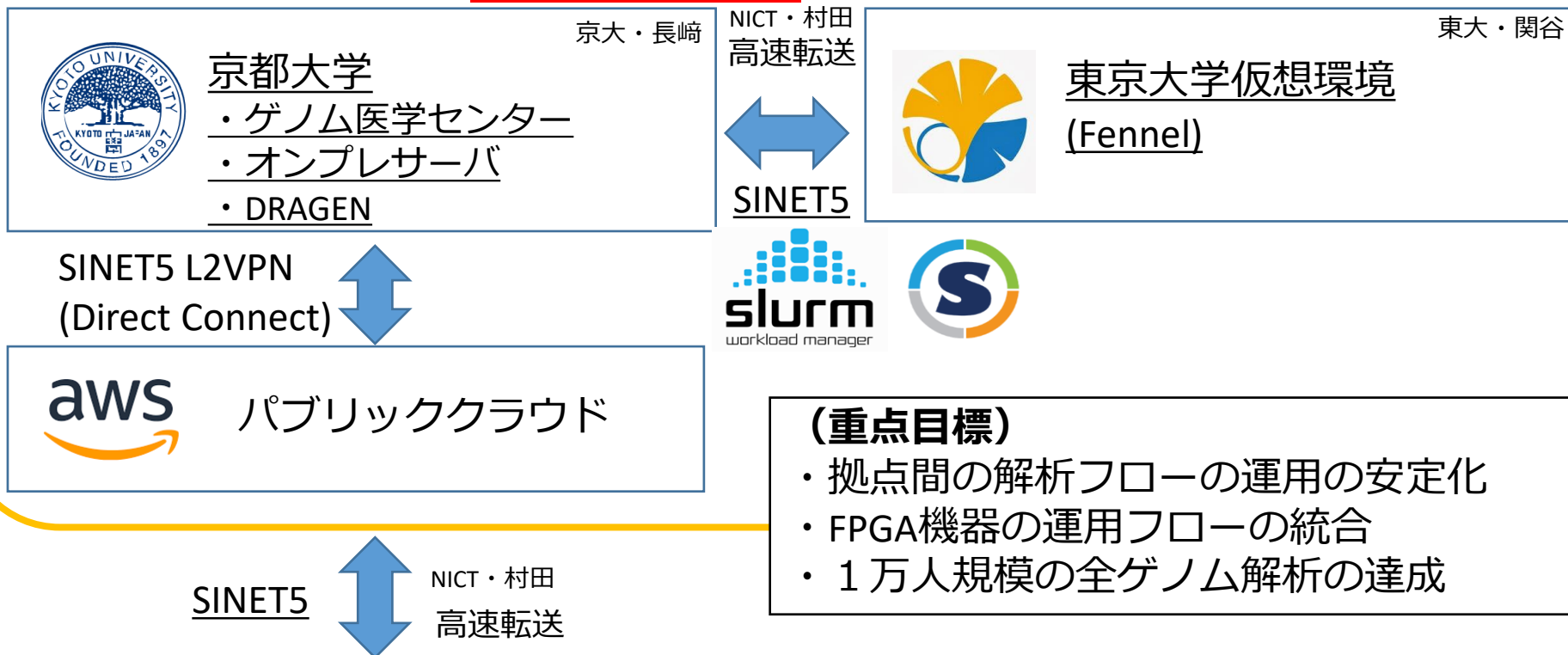
【役割 1】全ゲノムシークエンズデータの読み取りとその情報の拠点 1 への転送

【役割 2】解析結果の拠点1からの受け取りと結果評価

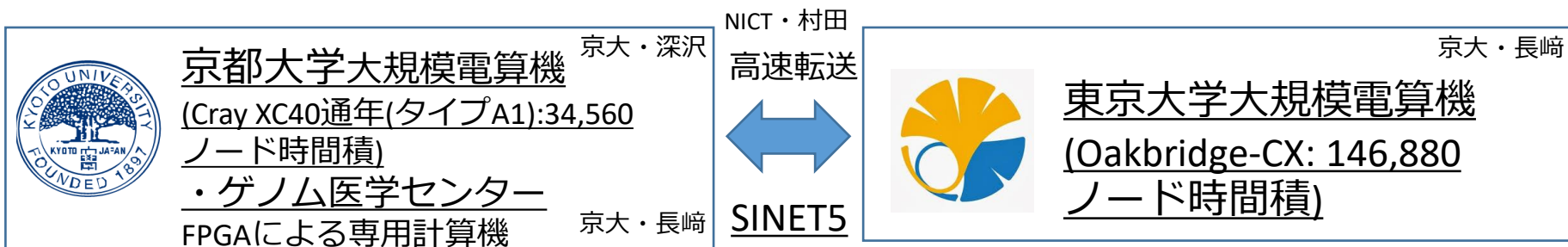
課題 1) 複数拠点間を効率的に運用できるハイブリッドクラウドシステムの設計と運用

システム全体構成と役割担当

ヒトゲノム情報解析でより汎用的な解析が求められる解析パイプラインの実装




ヒトゲノム情報解析で超高速な解析が求められる解析パイプラインの実装




ハイブリッドクラウド構築とゲノム情報解析の効率的な運用に関する研究

京大・長崎



全世界のTOP500 supercomputers **60%**で利用

→ 京大オンプレ & AWS & 東大仮想環境 (Fennel) での共通実装



Singularity v3を用いたパイプライン構築

→ コンテナの利用による再現性・再利用性・信頼性の担保

解析パイプラインをさらに実装予定

ヒトゲノム情報解析でより汎用的な解析が求められる解析パイプライン

パイプラインはPython / R / C++ / Javaなどさまざまなバイオインフォマティクス解析ソフトウェアのワークフローで構成

処理名称	パイプラインの概要	入力ファイル	出力ファイル
Genotyping	SNPアレイ (Japonica Array (CEL)ファイル) から約66万か所の遺伝型をクラスタリングによって決定するためのパイプライン	CEL	VCF/BED
Imputation	国際1000人ゲノムやそのほかの全ゲノムリファレンスパネルを用いることでSNPアレイでタイピングされた約66万か所のSNPから数千万のSNP情報を復元するパイプライン	VCF/BED	VCF/BED
GWAS	インピュテーション (1KGP / GRIFFIN Panel など) によって復元された数千万か所の変異情報について指定された条件でフィルタリングを行った後に疾患群と健常群などのcase/controlまたは検査情報などの連続量についての各SNPの偏りを統計手法により検定を行うパイプライン	VCF/BED	TXT
Annotation	GWASによってでてきた結果についてアノテーションを行うパイプライン	VCF/BED	TXT

NovaSeq6000 スペックシート

Table 1: NovaSeq 6000 System flow cell specifications

Flow cell type	SP	S1	S2	S4
Lanes per flow cell	2	2	2	4
Output per flow cell^{a,b}				
2 × 50 bp	65-80 Gb	134-167 Gb	333-417 Gb	N/A
2 × 100 bp	N/A	266-333 Gb	667-833 Gb	1600-2000 Gb
2 × 150 bp	200-250 Gb	400-500 Gb	1000-1250 Gb	2400-3000 Gb
2 × 250 bp	325-400 Gb	N/A	N/A	N/A
Single reads CPF	0.65-0.8 B	1.3-1.6 B	3.3-4.1 B	8-10 B
Paired-end reads CPF	1.3-1.6 B	2.6-3.2 B	6.6-8.2 B	16-20 B
Quality scores^c				
2 × 50 bp		≥ 85% ≥ Q30		
2 × 100 bp		≥ 80% ≥ Q30		
2 × 150 bp		≥ 75% ≥ Q30		
2 × 250 bp		≥ 75% ≥ Q30		
Run time^d				
2 × 50 bp	~13 hr	~13 hr	~16 hr	N/A
2 × 100 bp	N/A	~19 hr	~25 hr	~36 hr
2 × 150 bp	~25 hr	~25 hr	~36 hr	~44 hr
2 × 250 bp	~38 hr	N/A	N/A	N/A

S4のフローセルを用いて

1台で年間 約8,000人のヒト全ゲノム情報を取得できる。

データ量としては
546000G塩基

1塩基はA/T/G/Cの文字列とASCIIコードで表現される1文字

1検体あたり圧縮して30Gbぐらい

8000検体で解析前のデータで
250TB

【参考関連論文】

1. N. Nariai, K. Kojima, S. Saito, T. Mimori, Y. Sato, Y. Kawai, Y. Yamaguchi-Kabata, J. Yasuda and M. Nagasaki, HLA-VBSeq: accurate HLA typing at full resolution from whole-genome sequencing data, *BMC Genomics*, 16(2):S7, 2015.
2. M. Nagasaki *et al*, Rare variant discovery by deep whole-genome sequencing of 1,070 Japanese individuals, *Nature Communications*, 6:8018, 2015
3. Y Kawai, T Mimori, K Kojima, N Nariai, I Danjoh, R Saito, J Yasuda, M Yamamoto and M. Nagasaki, Japonica array: improved genotype imputation by designing a population-specific SNP array with 1070 Japanese individuals, *Journal of Human Genetics* 2015; 60: 581–587, 2015.
4. T. Hasegawa, K. Kojima, Y. Kawai, K. Misawa, T. Mimori, and M. Nagasaki, AP-SKAT: highly-efficient genome-wide rare variant association test, *BMC Genomics*, 17(1):745, 2016.
5. X. Jia, T. Horinouchi, Y. Hitomi, A. Shono, S.-S. Khor, Y. Omae, K. Kojima, Y. Kawai, M. Nagasaki, 17人略, K. Tokunaga, and K. Iijima, Strong Association of the HLA-DR/DQ Locus with Childhood Steroid-Sensitive Nephrotic Syndrome in the Japanese Population, *J. Am. Soc. Nephrol.*, vol. 29, no. 8, pp. 2189–2199, 2018.
6. Y.Y. Wang, T. Mimori, S. S. Khor, O. Gervais, Y. Kawai, Y. Hitomi, K. Tokunaga and M. Nagasaki, HLA-VBSeq v2: improved HLA calling accuracy with full-length Japanese class-I panel, *Hum Genome Var*, 6: 29, 2019.
7. O. Gervais, K. Ueno, Y. Kawai, Y. Hitomi, Y. Aiba, M. Ueta, M. Nakamura, K. Tokunaga and M. Nagasaki. Regional heritability mapping identifies several novel loci (STAT4, ULK4, and KCNH5) for primary biliary cholangitis in the Japanese population. *European Journal of Human Genetics*, 2021.
8. T. Tanjo, Y. Kawai, K. Tokunaga, O. Ogasawara and M. Nagasaki. Practical guide for managing large-scale human genome data in research. *J Hum Genet*, 66 (1) 39-52, 2021.

【参考発表】

- 1) 長崎 正朗, “ヒトゲノム情報統合解析に向けた京都大学ゲノム医学センターのハイブリッドクラウドシステム構築について”, AWS Summit Online Japan 2020 (2020/9/8-2020/9/30).
- 2) 長崎 正朗, AWS SUMMIT ONLINE JAPAN Report
https://special.nikkeibp.co.jp/atcl/NXT/20/aws1030_01/