

# ハイブリッドクラウドを用いた ゲノム情報に基づく 構造多型パネルの構築とアノテーション



【課題番号】 jh220014

【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

JHPCN第14回シンポジウム 2022採択課題研究紹介 2022/7/7

# 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・松田 文彦

申請課題の全共同研究者

【京都大学】

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】

シーケンス技術の進展によりヒト全ゲノム配列情報とそれに付随して必要となる解析資源が爆発的に増えている。

申請者は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上で出てくる課題を解決することで

「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（令和2－3年度jh200047-NWH,jh210018-NWH）」を進めてきた。

当研究課題では、申請者が進めている日本人の長鎖型情報を活用し、令和4年度において、国外において先行研究が行った5,202検体と同程度の約5,000検体の短鎖型法の検体の解析を目標に構造多型のカタログの構築（図1）を行う。

図1 【研究目的】日本人のゲノム情報についてグラフゲノムによる解析による**構造多型（※）**のカタログ同定と疾患解析での活用

※ 構造多型：染色体上の複雑な配列の集団内の多様性を指す。SNVなどの1塩基の単純な集団内の多様性とは区別して記載する。IrWGSを用いることで徐々に解明されつつある。

海外の先行研究

RESEARCH ARTICLE SUMMARY

GENOMICS

Pangenomics enables genotyping of known structural variants in 5202 diverse genomes

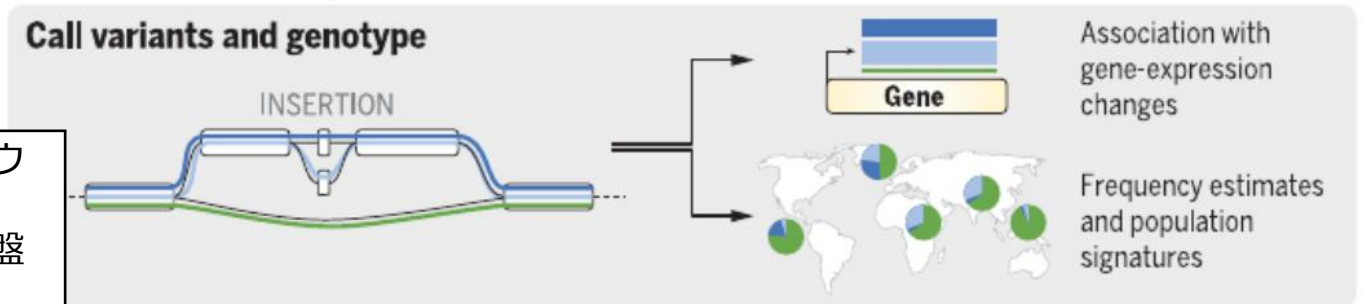
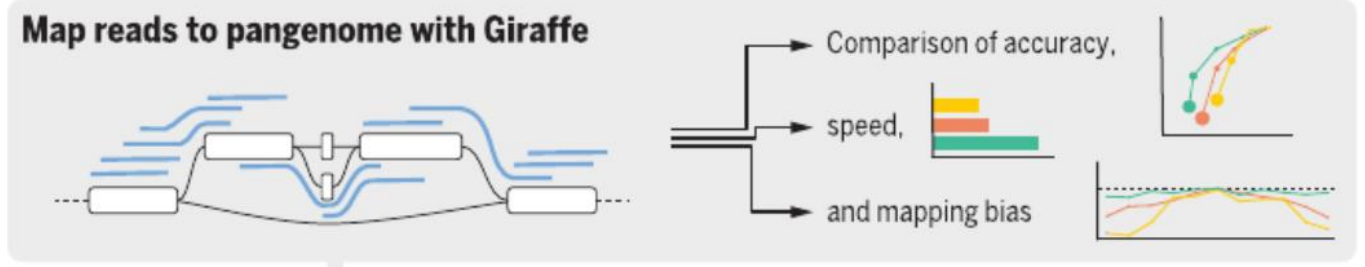
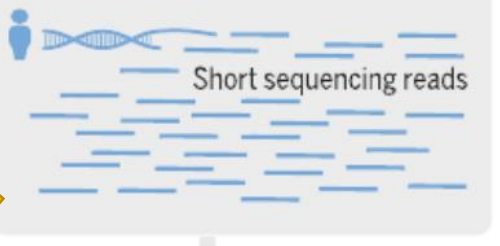
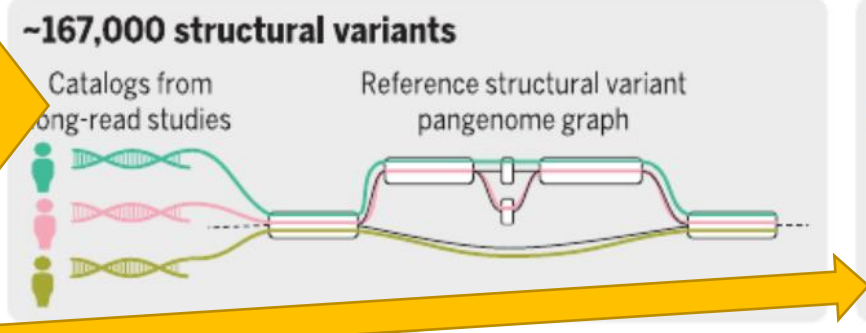
日本人のゲノム情報

約100人の日本人のIrWGSゲノム情報での解析

約5000人の日本人のsrWGSゲノム情報での解析

アノテーションと疾患解析での基盤としての活用

【課題1】ハイブリッドクラウド情報基盤での実装  
【課題2】アノテーション基盤情報の取得と拠点間転送



[https://github.com/graph-genome/graph\\_summarization](https://github.com/graph-genome/graph_summarization)

# 【研究計画】

**課題 1)** 構造多型のリファレンスパネルの構築とそのための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情報基盤の設計と運用 (図 2 参照)  
(長崎、松田、関谷、塙、深沢、村田)

**課題 2)** シークエンサから取得されたアノテーションに活用する情報を他の拠点に効率良く展開するための設計検討と実装 (図 3 参照) (大川、長崎、深沢、村田)

# 課題 1) 構造多型のリファレンスパネルの構築とそのための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情報基盤の設計と運用

## 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) 松田文彦 他

【役割】構造多型のリファレンスパネルの構築の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 (後スライド図 2 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SINET6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・埜 敏博

【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

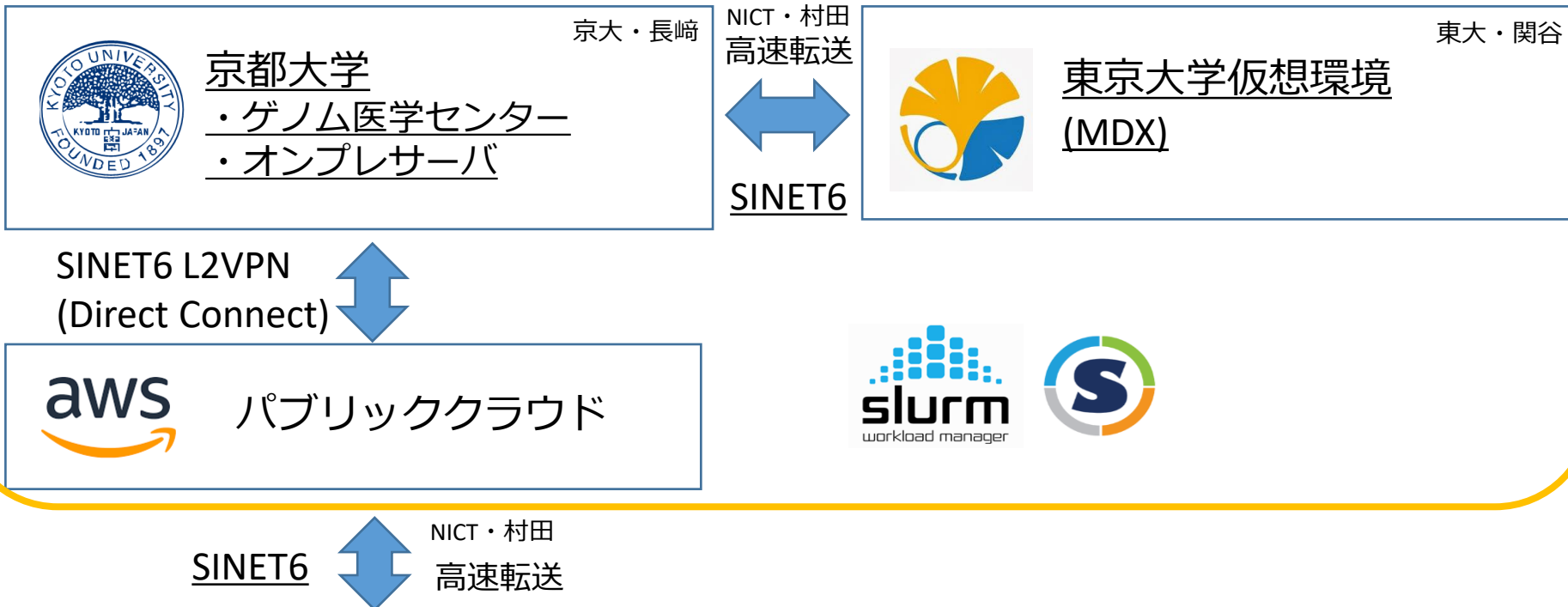
【担当者】村田健史

【役割】拠点間的高速データ転送技術提供と評価

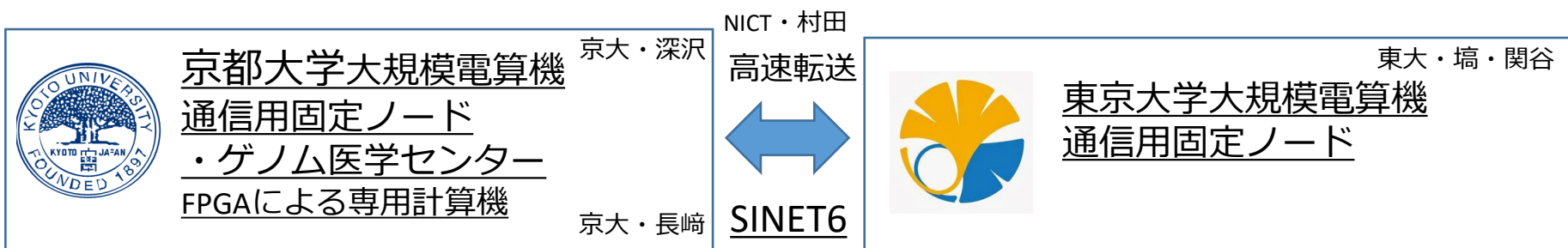
図2 課題1) 構造多型のリファレンスパネルの構築とそのための複数拠点間を効率的に運用できるハイブリッドクラウド情報基盤の設計と運用 長崎、松田、関谷、塙、深沢、村田

## 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



# 課題2) シークエンサから取得されたアノテーションに活用する情報を他の拠点に効率良く展開するための設計検討と実装

## 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】埴敏博・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

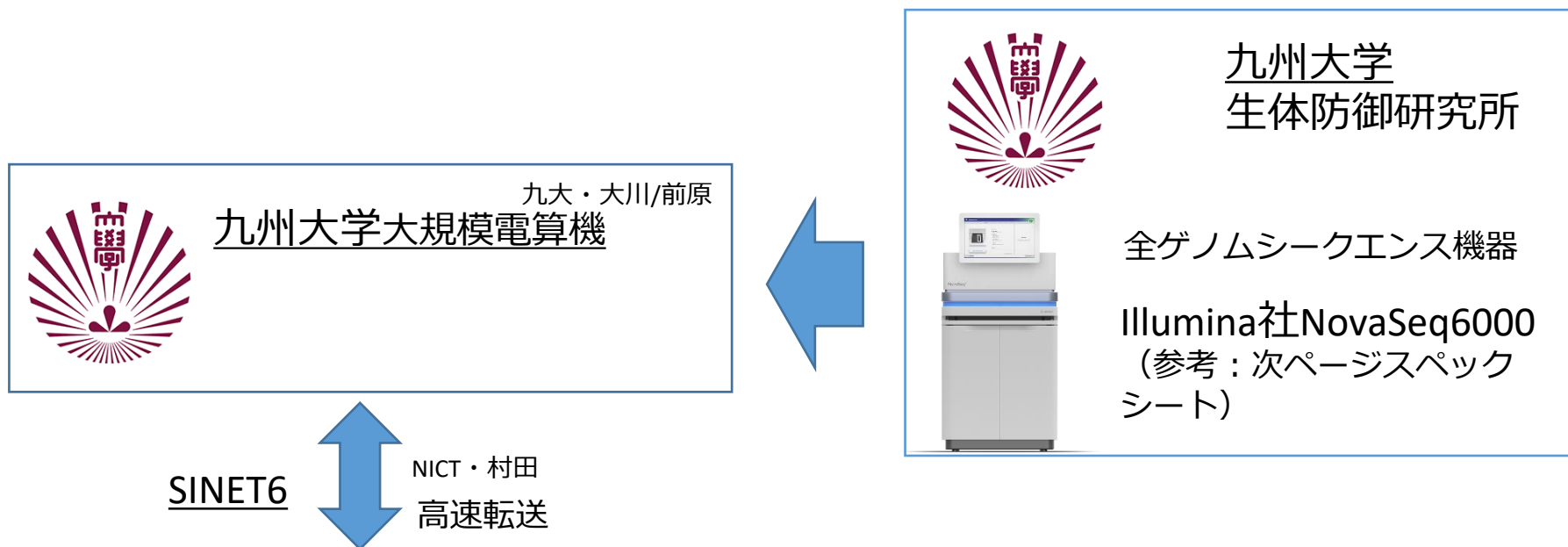
【担当者】大川恭行 前原一満

【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

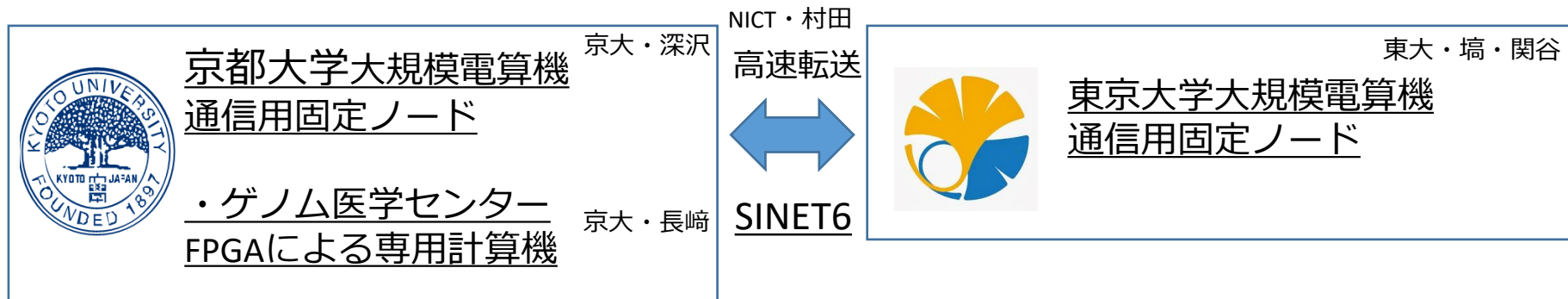
【役割2】解析結果の拠点1からの受け取りと結果評価



図3 課題2) シークエンサから取得されたアノテーションに活用する情報を他の拠点に効率良く展開するための設計検討と実装 (大川、長崎、深沢、村田)



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



【課題1】との連携

# 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
<b>Output per flow cell<sup>a,b</sup></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
<b>Quality scores<sup>c</sup></b>				
2 × 50 bp		≥ 85% ≥ Q30		
2 × 100 bp		≥ 80% ≥ Q30		
2 × 150 bp		≥ 75% ≥ Q30		
2 × 250 bp		≥ 75% ≥ Q30		
<b>Run time<sup>d</sup></b>				
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

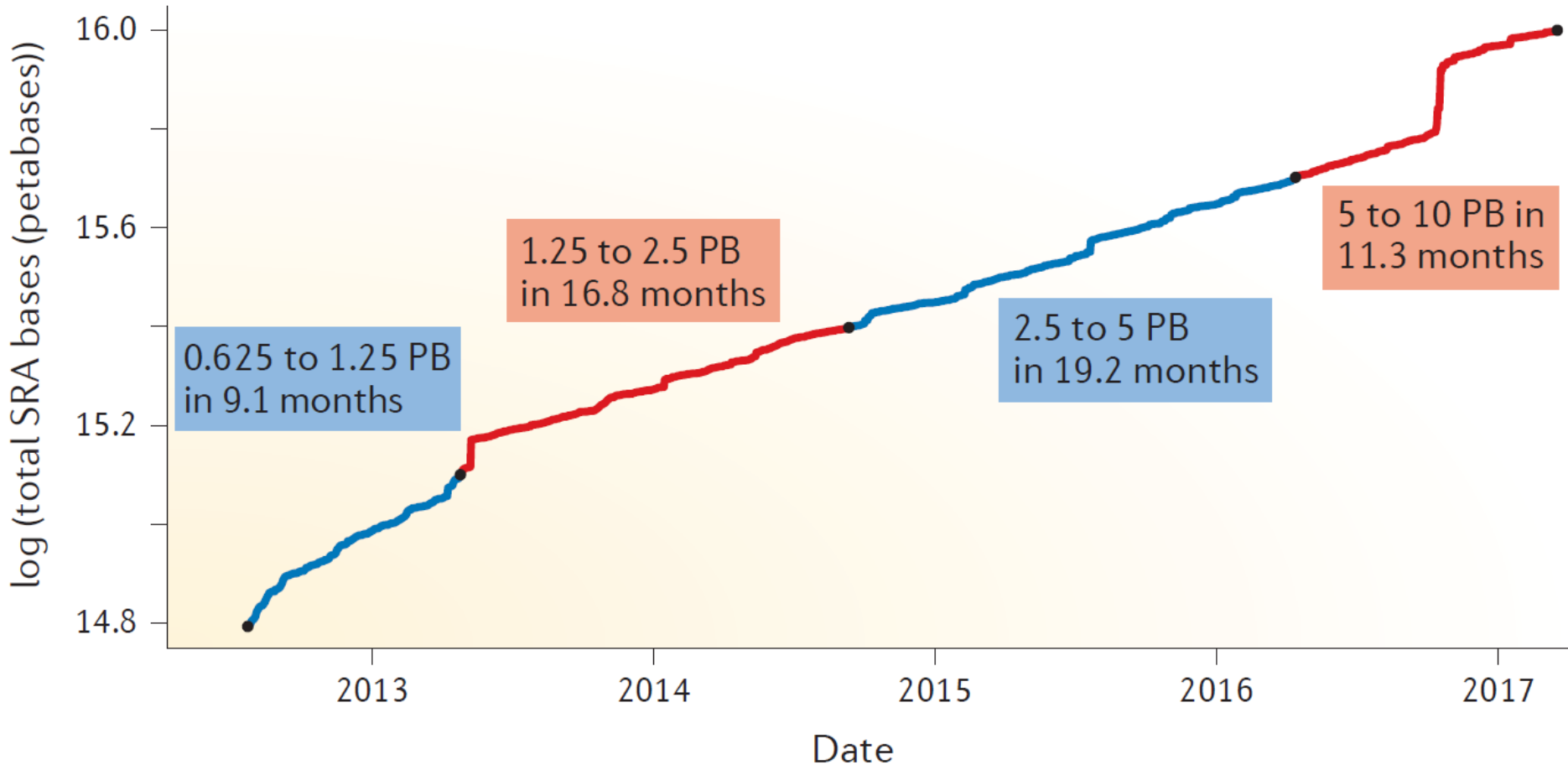
# 【参考関連論文】

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*, 29 (8):1282-1291, 2021.
8. T. Tanjo, Y. Kawai, K. Tokunaga, O. Ogasawara and M. Nagasaki. Practical guide for managing large-scale human genome data in research. *Journal of Human Genetics*, 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](https://special.nikkeibp.co.jp/atcl/NXT/20/aws1030_01)
3. 長崎 正朗, “Accelerating the pace of research in Kyoto University”, AWS Public Sector Summit Online, 2021/4/15-2021/4/16.
4. 長崎 正朗, 山口 泉, 川口 喬久, 寺岡 凌, 稲富 雄一, 深沢 圭一郎, 関谷 勇司, 塙 敏博, 大川 恭行, 王 妍雁, Olivier Gervais, Seik-Soon Khor, 植野 和子, 浅倉 章宏, 関谷 弥生, 人見 祐基, 小野 彰, 男澤 良子, 河合 洋介, 前原 一満, 南里 豪志, 村田 健史, 橋本 洋希, 丹生 智也, 小笠原 理, 山田 亮, 松田 文彦, 徳永 勝士, “ゲノム医科学における国内外のヒトゲノム解析の状況およびハイブリッドクラウド計算環境の構築と活用”, 第44回日本分子生物学会年会, 2021/12/1-2021/12/3.
5. 山口 泉, 川口 喬久, 寺岡 凌, 稲富 雄一, 深沢 圭一郎, 関谷 勇司, 塙 敏博, 大川 恭行, 王 妍雁, Olivier Gervais, Seik-Soon Khor, 植野 和子, 浅倉 章宏, 関谷 弥生, 人見 祐基, 小野 彰, 男澤 良子, 河合 洋介, 前原 一満, 南里 豪志, 村田 健史, 橋本 洋希, 徳永 勝士, 松田 文彦, 山田 亮, 長崎 正朗, “クラウドサーバとオンプレミスサーバを組み合わせたハイブリッドシステムの構築と活用”, 日本人類遺伝学会第66回大会, 2021/10/13-2021/10/16.
6. 山口 泉, 川口 喬久, 寺岡 凌, 稲富 雄一, 深沢 圭一郎, 関谷 勇司, 塙 敏博, 大川 恭行, 前原 一満, 南里 豪志, 村田 健史, 橋本 洋希, 松田 文彦, 山田 亮, 長崎 正朗, “クラウドサーバとオンプレミスサーバを組み合わせたハイブリッドシステムの構築と活用”, 第41回医療情報学連合大会, 2021/11/18-2021/11/21.

# ヒトゲノムシーケンス情報の増加の様子

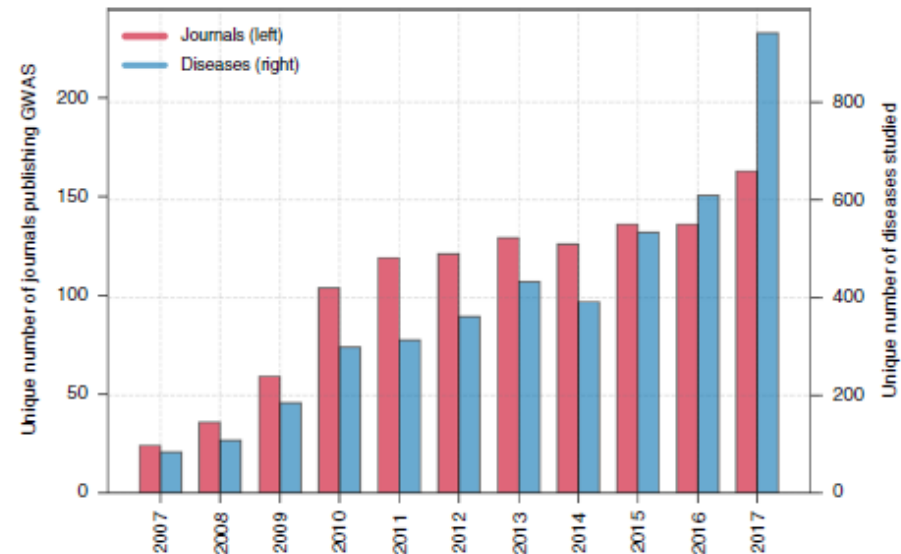
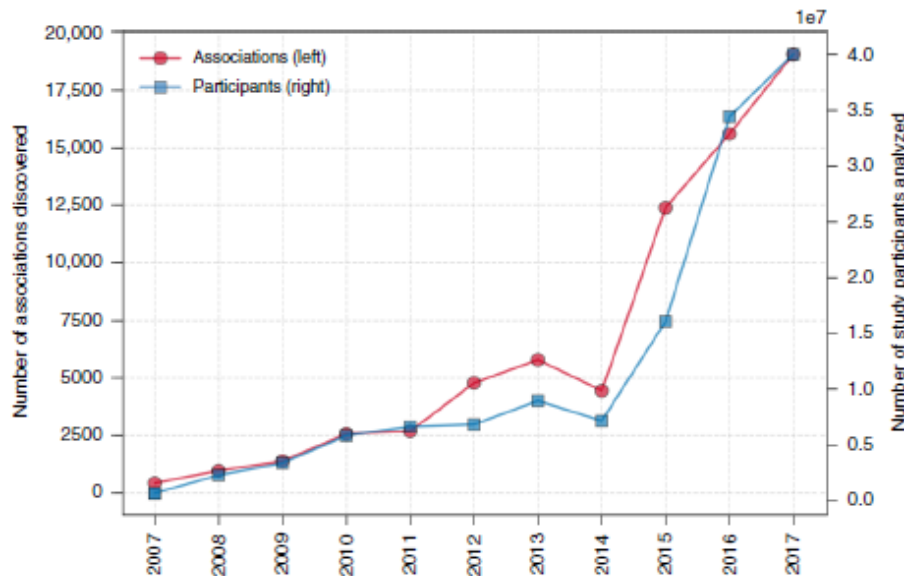
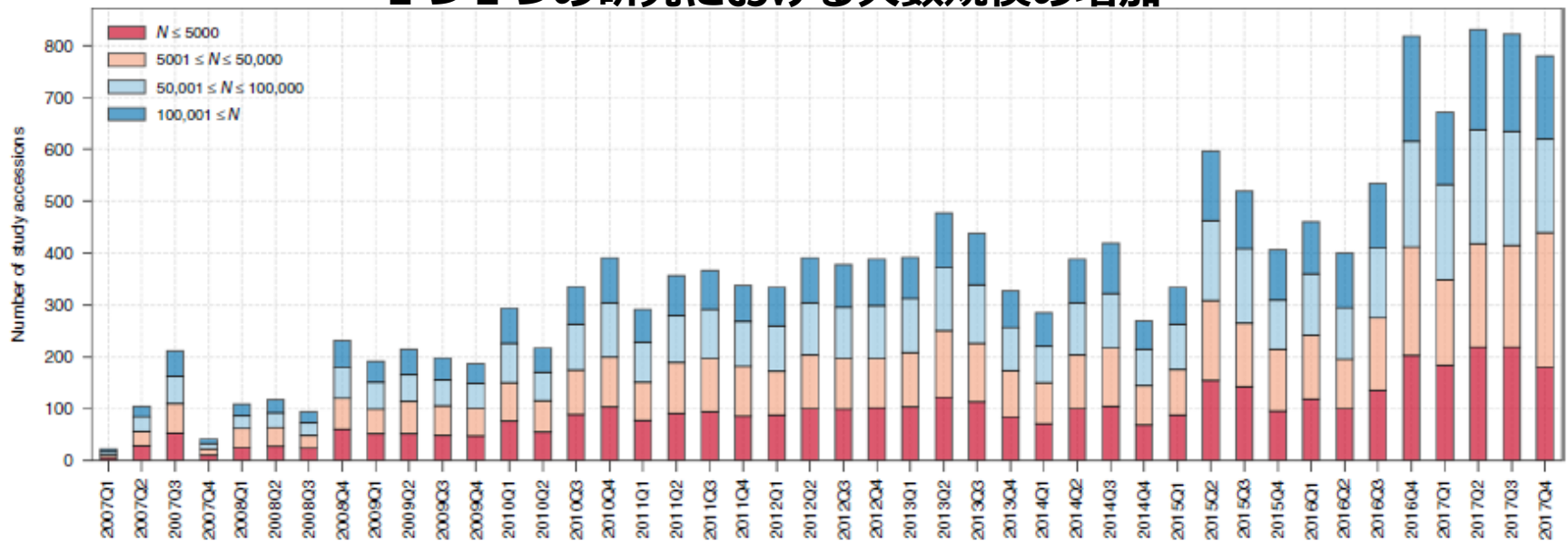


## 【図4】 ヒトゲノムシーケンス情報の増加の様子

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5】 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。**