

ゲノム科学におけるクラウド活用： 独占型スパコン構築から データ共有まで

*Cloud Utilization in Genome Science: From Building
Own Supercomputers to Sharing Data*

重信 秀治

Shuji Shigenobu

基礎生物学研究所 (NIBB)

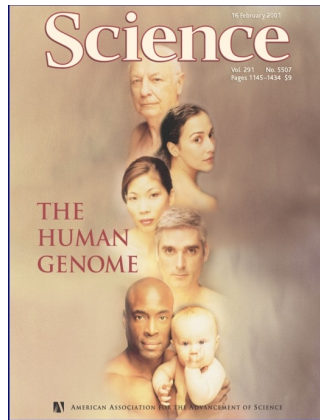
National Institute for Basic Biology, NINS

shige@nibb.ac.jp

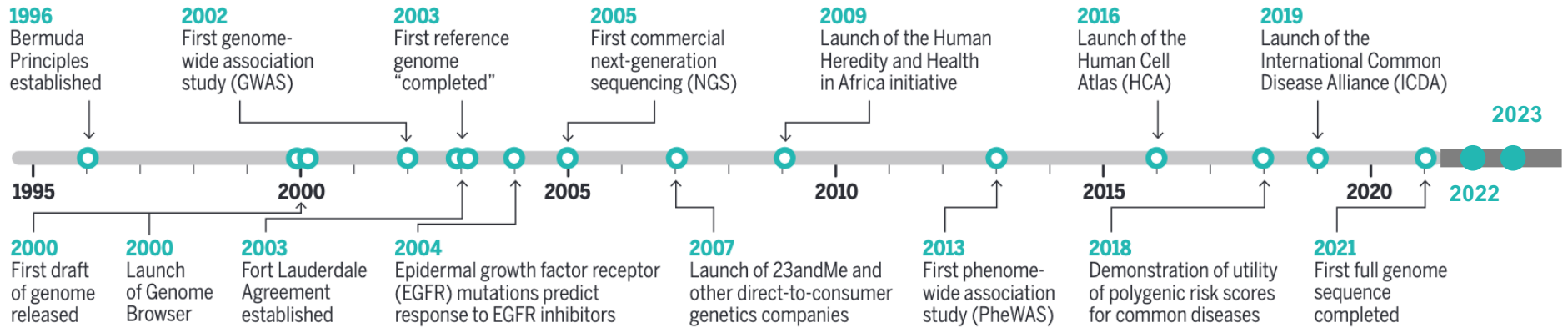
本日のトピック

ゲノム科学分野、特に基礎研究、におけるクラウド活用

- ▶ 計算資源(HPC)の確保
- ▶ データ共有による共同研究促進

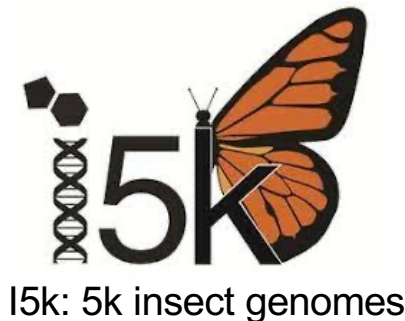


2000



(Rood et al. 2021 w/ modification)

急速に発展するゲノム科学



2022



2023

私が基礎生物学研究所で解読した 昆虫ゲノム（抜粋）



エンドウヒゲナガ
アブラムシ
細胞内共生・環境
適応・植物適応

(IAGC 2010 PLOS
Biol, Shigenobu et al.,
2000 Nature)



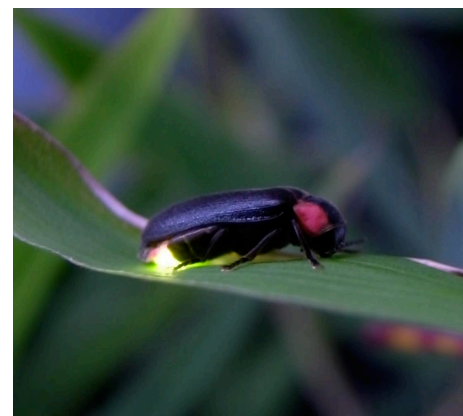
ヤマトシロアリ
社会性の進化

(Shigenobu et al.,
2022 PNAS)



カブトムシ
武器形質の進化と
性的二型

(Morita et al. 2023 Sci
Rep)



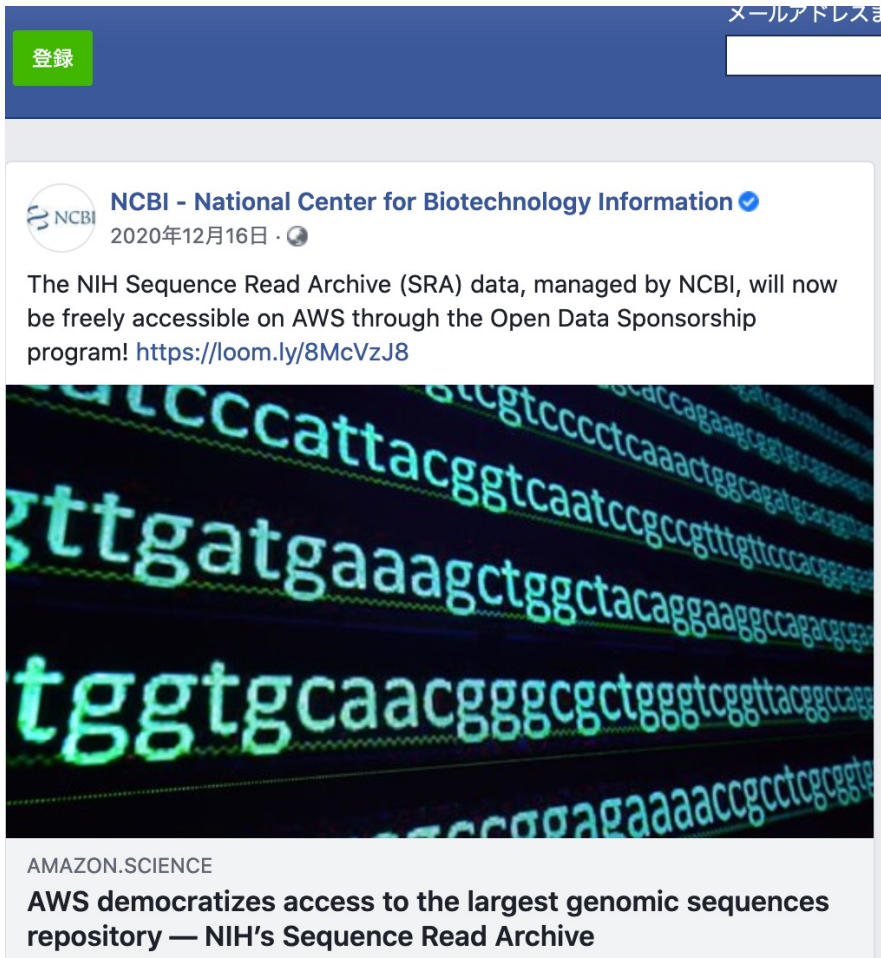
ヘイケボタル
発光機能の獲
得と進化

(Fallon et al., 2018
eLife)

ますますビッグデータ化するゲノム生物学
どうする コンピュータリソース？

ゲノム研究分野におけるクラウド化の動向

米国ではクラウド化の流れ



登録

メールアドレスを

NCBI - National Center for Biotechnology Information

2020年12月16日 ·

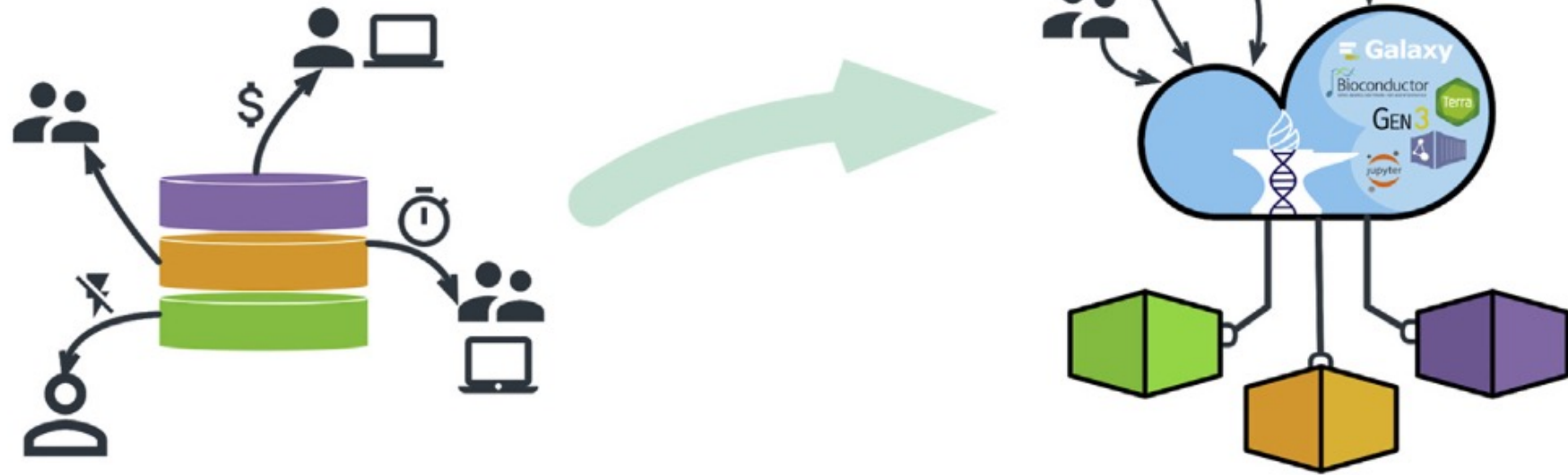
The NIH Sequence Read Archive (SRA) data, managed by NCBI, will now be freely accessible on AWS through the Open Data Sponsorship program! <https://loom.ly/8McVzJ8>

AMAZON.SCIENCE
AWS democratizes access to the largest genomic sequences repository — NIH's Sequence Read Archive

Facebook内ニュース 2020/12/16

- ▶ NCBIの次世代シーケンスデータ「SRA」はAmazon AWS上で提供される。
- ▶ SRA data:
 - ▶ 40 petabytes now
 - ▶ Will double every 1-1.5 year
- ▶ Saved in S3 storage service on AWS. Open access.
- ▶ NIHはSRA以外にもクラウド化推進の姿勢

AnVIL: the NHGRI Genomic Data Science Analysis, Visualization, and Informatics Lab-space



- AnVIL platform is NHGRI-supported data commons running on the Google Cloud Platform (GCP).
- Terra: analysis platform, Gen3: data search and artificial cohort creation; Dockstore : Docker-based genomic analysis tools and workflows.
- Support interactive analysis tools: Jupyter notebooks, Bioconductor, RStudio, and Galaxy.

<https://anvilproject.org/>

Schatz *et al.* *Cell Genom* (2022)

AWS Offers Genomics Analysis Platform



**Amazon Genomics
CLI**



**AWS HealthOmics
(Amazon Omics)**

日本のアカデミアにおけるクラウド動向

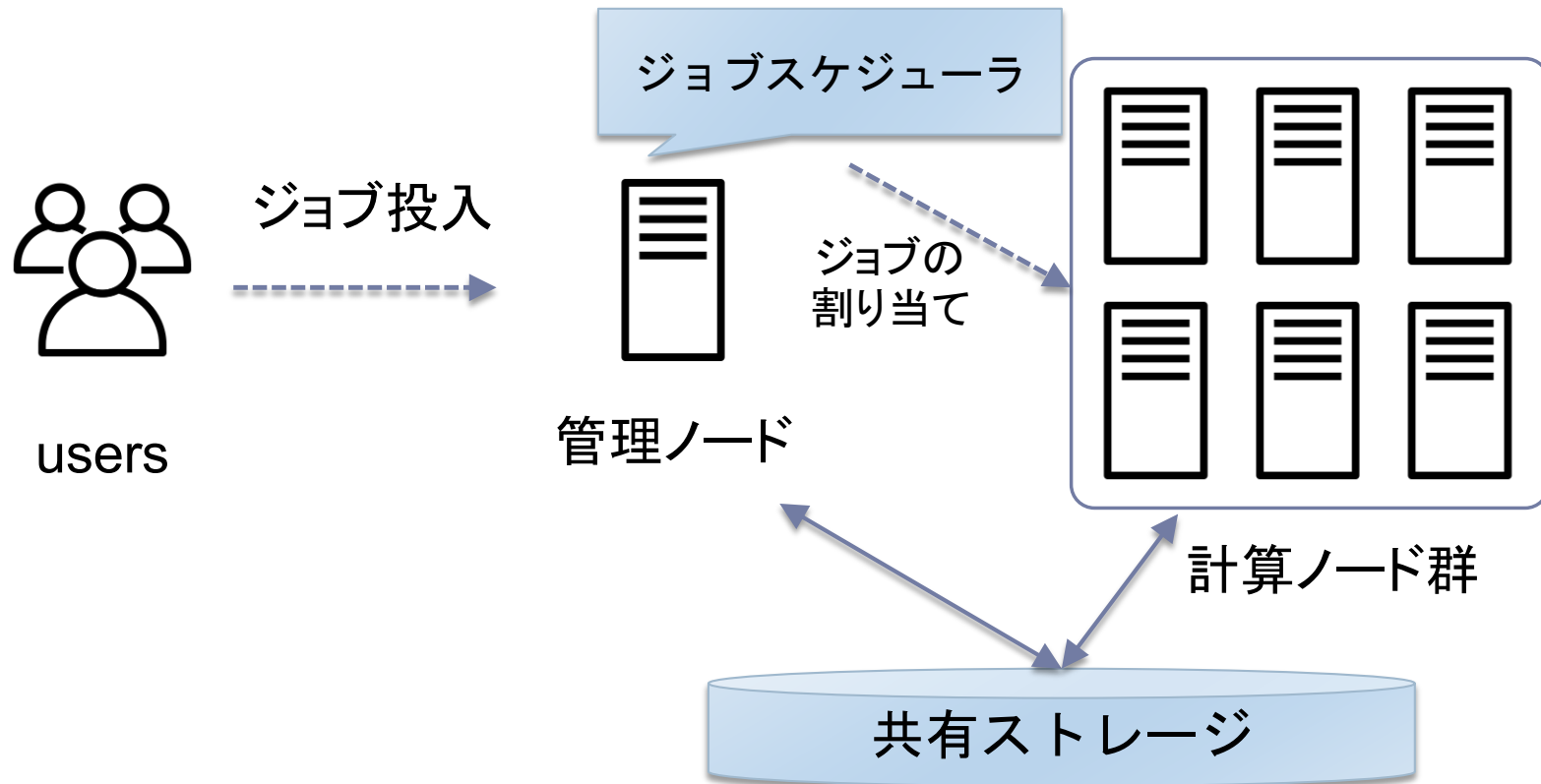
- ▶ 2018年、政府は「クラウドサービスの利用推進」を宣言
- ▶ 2023年5月、競争的研究費の直接経費からクラウド利用料の支出が可能であることが明確化
- ▶ 日本のゲノム科学界隈では、まだそれほどクラウド利用は浸透していない様子だが、先端的な利用例が報告されつつある。
 - ▶ (e.g.) Nagasaki, M. et al. Design and implementation of a hybrid cloud system for large-scale human genomic research. Hum Genome Var 10, 6 (2023).

ますますビッグデータ化するゲノム生物学
どうする コンピュータリソース？

1. HPC (High-Performance Computing)
2. ゲノム関連の情報とツールの共有

HPC (High-Performance Computing) on AWS

HPC環境の典型的な構成要素



多くの大学や研究機関では共有のHPC(クラスターコンピュータ・スパコン)を整備。

共用HPCの問題点

- ▶ 共用のオンプレミスHPCではジョブ待ちが頻繁。
- ▶ OSやライブラリのバージョンが固定。ソフトウェアが動かない場合、対処が困難。
- ▶ オンプレミスサーバは調達に時間がかかり、手続きも煩雑。メンテナンスなどの運用負担。

ゲノム科学の基礎研究に特有の事情

- ▶ 次世代シーケンスデータなどファイルサイズが大きくなりがち
- ▶ 基礎研究においては頻繁なトライ&エラーが重要
- ▶ ゲノム解析のワークフローでは多様なソフトウェアを組み合わせて利用する。それぞれのソフトウェアごとに必要なコンピュータリソースの特性が異なる。

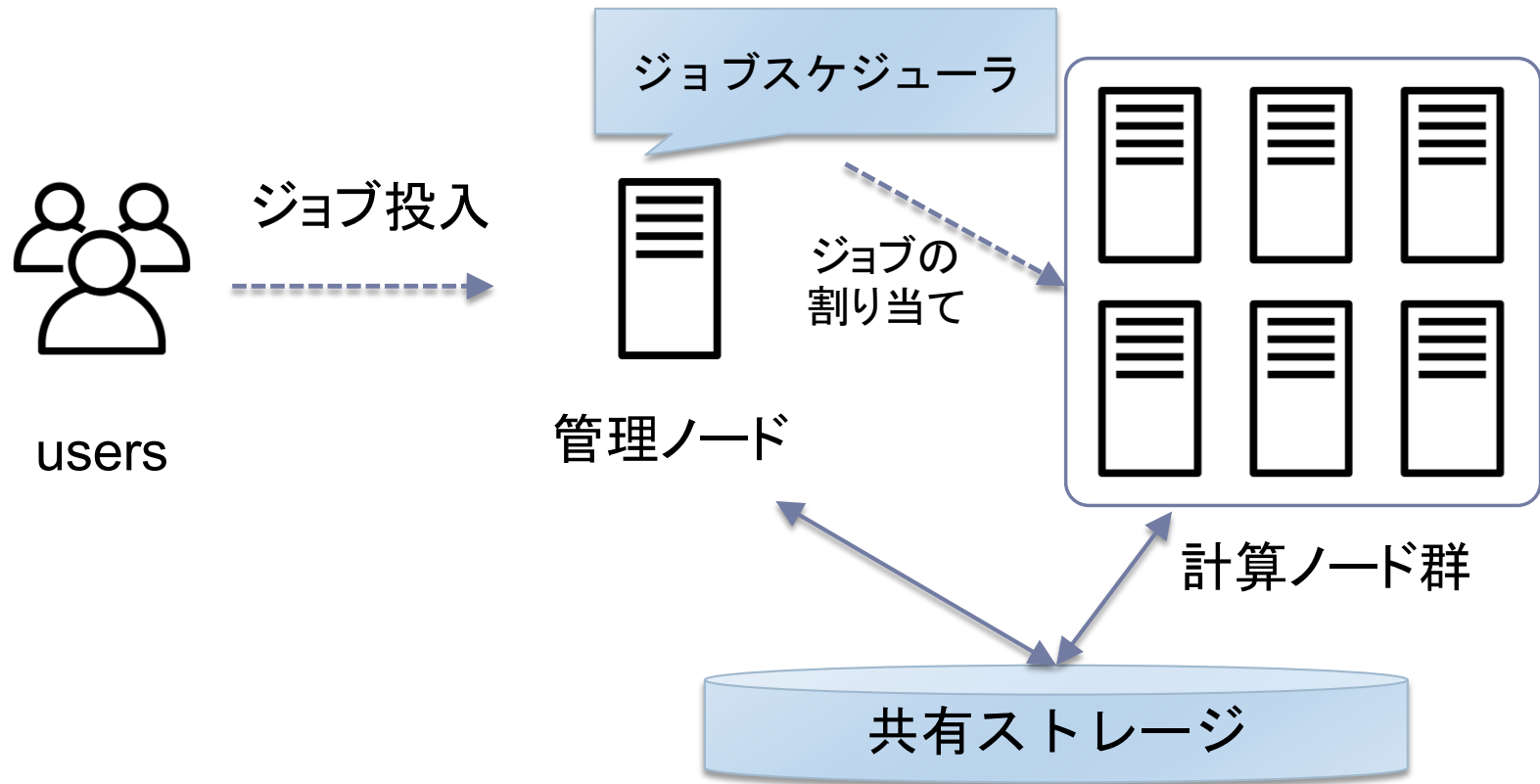
ゲノム科学の基礎研究に適した「HPC」はクラウドで実現可能か？

Yes



AWS Parallel Cluster

HPC環境の典型的な構成要素



AWSが提供するHPCフレームワーク



AWS Parallel Cluster



AWS ParallelClusterとは

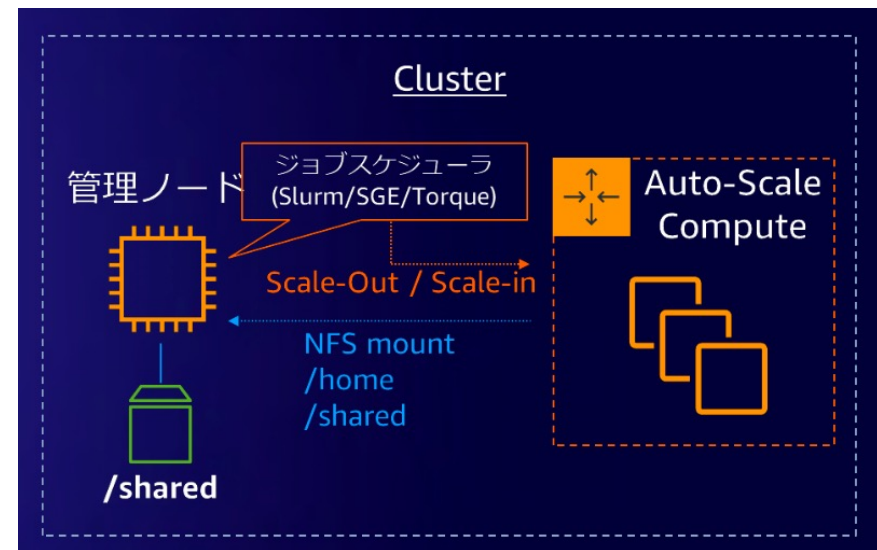
数コマンド操作で、ジョブ投入に応じて自動でスケールするクラスタをAWS上に構築可能なAWS公式のオープンソースソフトウェア

主な特徴

- ▶ HPC向けジョブスケジューラ(Slurm)と auto-scaling を連携した環境を作成
- ▶ 使用するOSや、インスタンスタイプ、ネットワーク環境、ストレージ構成などを柔軟にカスタマイズ可能
- ▶ OSS (<https://github.com/aws/aws-parallelcluster>)

利用方法

- ▶ ConfigファイルをYAMLフォーマットで記述。pcluster create コマンドを実行する。
- ▶ 管理ノードにログイン。実行したいジョブのシェルスクリプトを作成し、ジョブ投入。



pcluster CLI

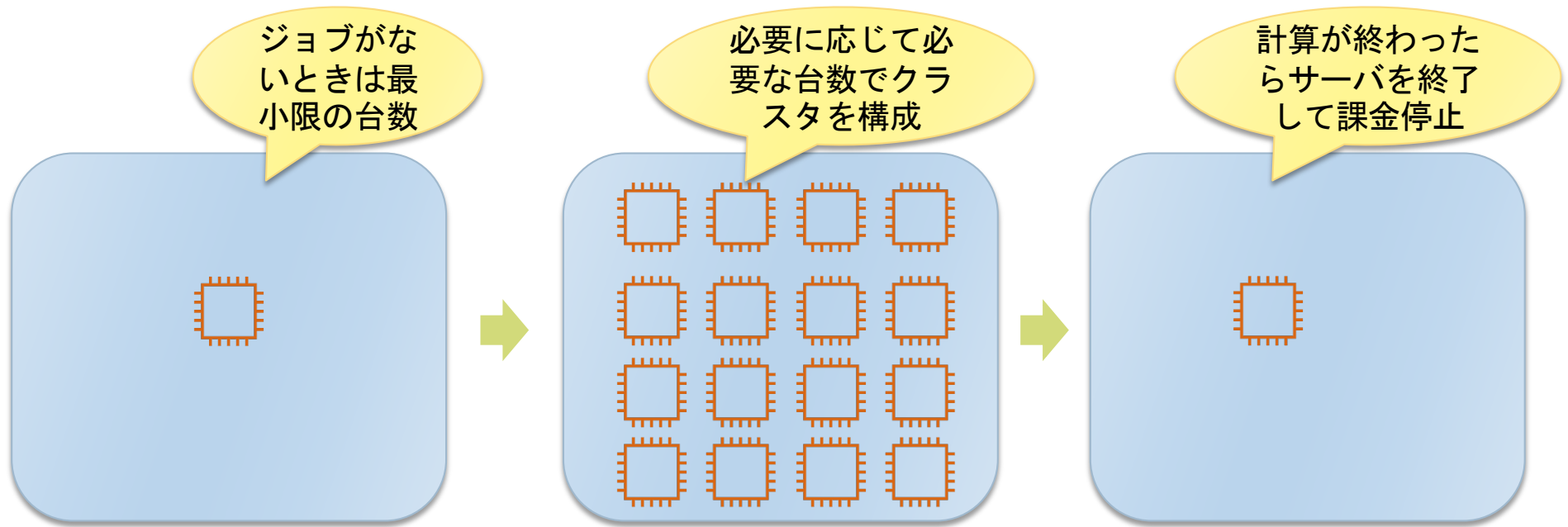
設定ファイル (YAML)



users

Amazon AWS summit 2021 の資料を一部改変

柔軟に自動スケールするクラスタ



技術的メリット






- ジョブ待ちの無いHPC環境を実現。
- 必要な時に必要な分の計算リソースを一気に確保
=> 解析時間短縮。

コスト的メリット

- 費用は「時間x台数」。無駄がない。
- スポットインスタンス(最大90%OFF)を使えば大幅なコスト節約。

多様なスペックのEC2インスタンスを選択可

高性能計算向けインスタンスタイプの例

高性能 CPU の選択肢		アクセラレータの選択肢		
				
Intel Xeon processor (x86_64 arch)	AMD EPYC processor* (x86_64 arch)	AWS Graviton Processor (64-bit Arm arch)	NVIDIA GPU	Xilinx FPGA
M6i インスタンス Ice Lake 最大時全コア 3.5 GHz 駆動	M6a インスタンス EPYC Milan 最大 3.3 GHz 駆動	C7g インスタンス 64bit Arm Neoverse V1ベース AWS Graviton3 CPU 搭載	P3 インスタンス V100 GPU 搭載 P4d インスタンス A100 GPU 搭載 P4de インスタンス* A100 (80GB版) GPU 搭載 G5 インスタンス A10G GPU 搭載	F1 インスタンス Virtex UltraScale+ VU9P 搭載
M5zn インスタンス Cascade Lake 最大全コア 4.5 GHz 駆動	Hpc6a インスタンス EPYC Milan HPC特化			

© 2022, Amazon Web Services, Inc. or its Affiliates. All rights reserved.

※ 2022年6月現在プレビュー提供中

Amazon AWS summit
2021 の資料を一部改変

ゲノム解析のワークフローでは多様なソフトウェアを組み合わせて利用する。それぞれのソフトウェアごとに必要なコンピュータリソースの特性が異なる。インスタンスタイプの選択肢が広いことはアドバンテージ。

- ▶ C, Mシリーズ(汎用・計算)はコア数とメモリのバランスが取れており、大小インスタンスを組み合わせることで幅広いアプリケーションに対応可能。
- ▶ I/Oの激しいジョブはローカルストレージを備えるdタイプのインスタンスが有用(例: M5d.8xlarge: 32 CPUs, 128 GB RAM, 2x600 SSD)
- ▶ ゲノム解析には高メモリを要求するソフトウェアも多い。R、X、ハイメモリシリーズが有用。
- ▶ 機械学習、AIにはGPUシリーズが有用もしくは必須(例: AlphaFold2)
- ▶ スポットインスタンスでコスト削減

Before



AWS Parallel Cluster

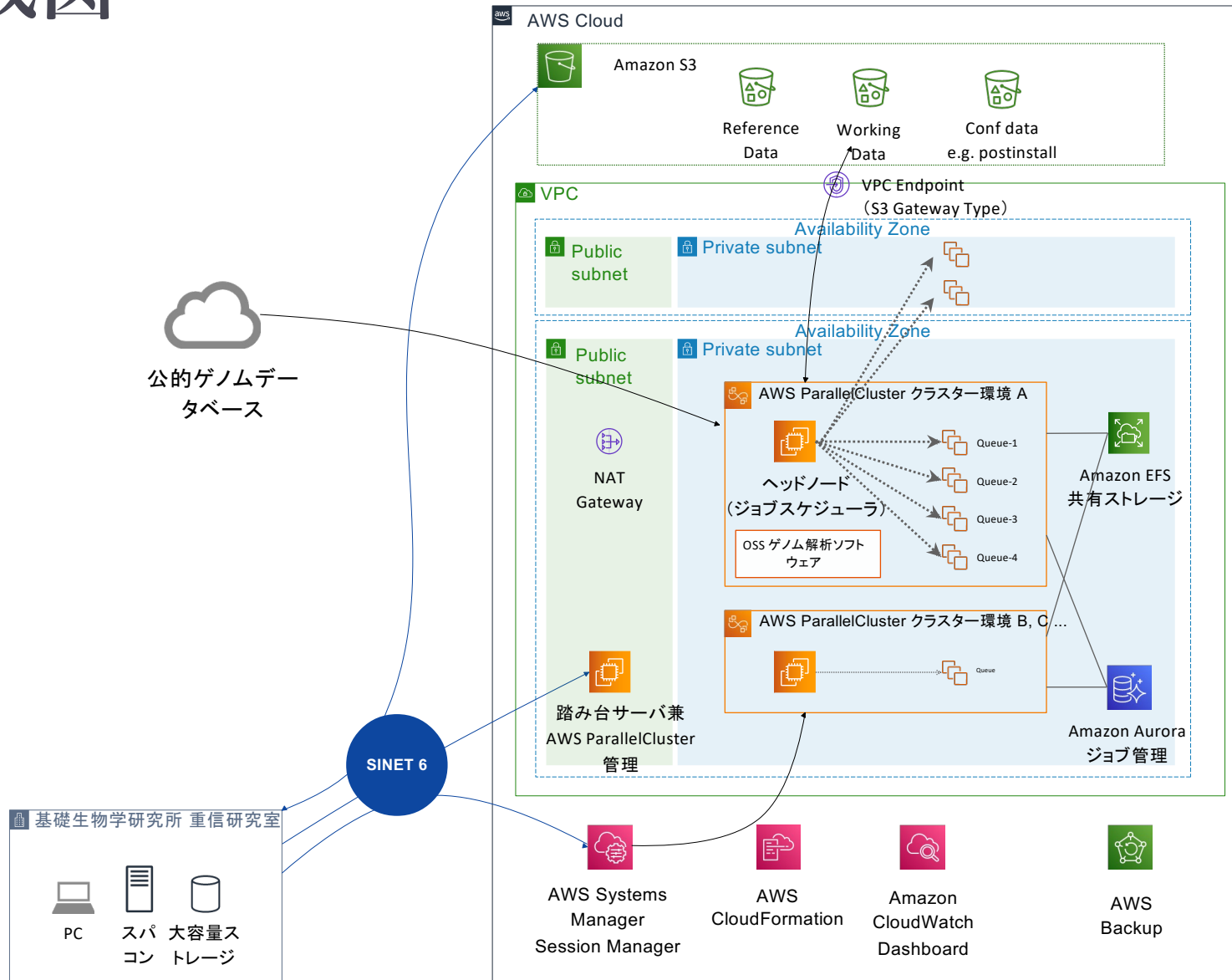


After

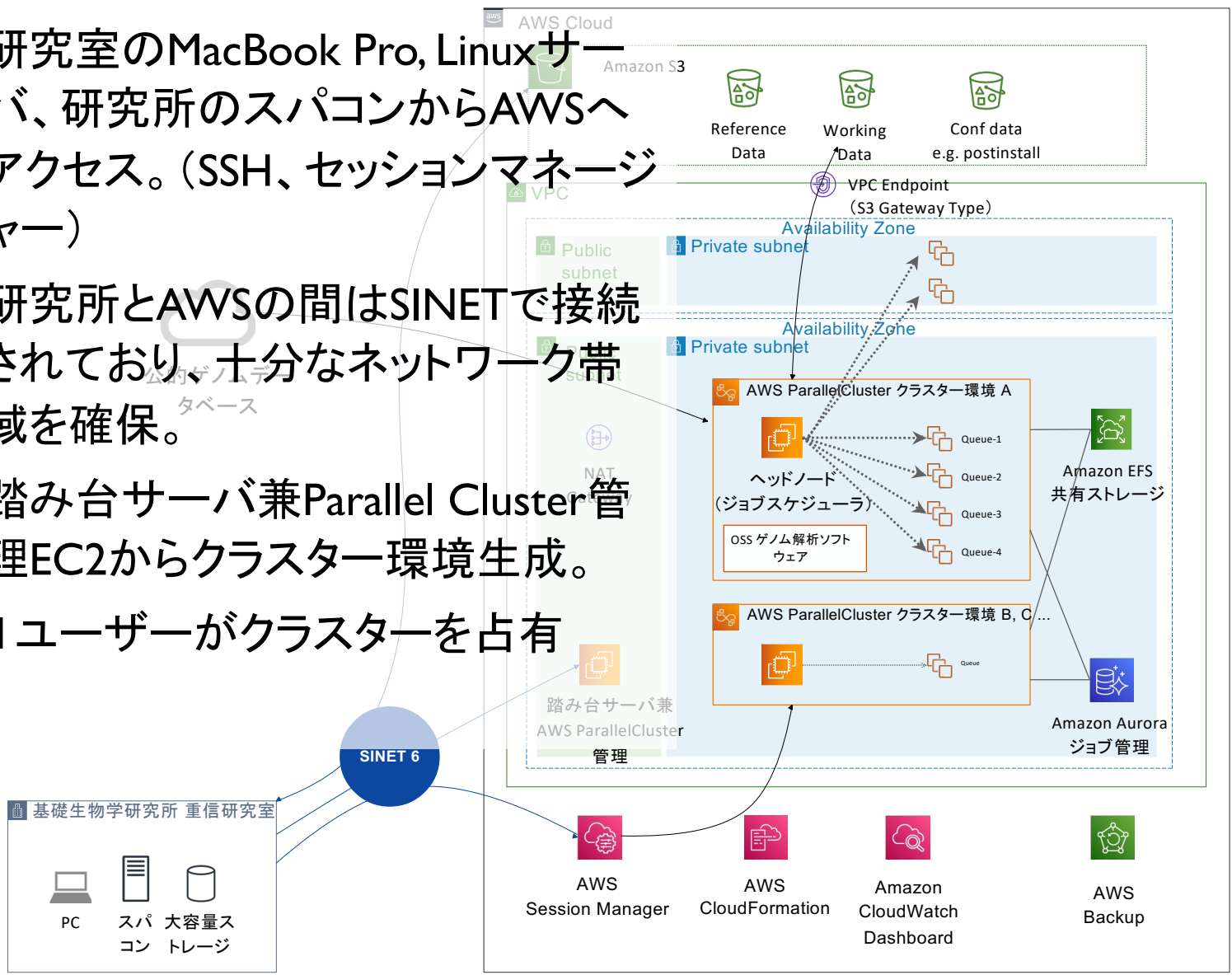
- ▶ 共用のオンプレミスHPCではジョブ待ちが頻繁。
 - ▶ OSやライブラリのバージョンが固定。ソフトウェアが動かない場合、対処が困難。
 - ▶ オンプレミスサーバは調達に時間がかかり、手続きも煩雑。メンテナンスなどの運用負担。
- ▶ HPCの計算機リソースを必要な時に必要な分を即時確保。
 - ▶ 理想的な計算環境を、root権限でスクラップ&ビルド。
 - ▶ 事務手続きやコンピュータメンテナンスから解放されて、研究そのものに集中。

当研究室での運用の実際

構成図



- ▶ 研究室のMacBook Pro, Linuxサーバ、研究所のスパコンからAWSへアクセス。(SSH、セッションマネージャー)
- ▶ 研究所とAWSの間はSINETで接続されており、十分なネットワーク帯域を確保。
- ▶ 踏み台サーバ兼Parallel Cluster管理EC2からクラスター環境生成。
- ▶ 1ユーザーがクラスターを占有



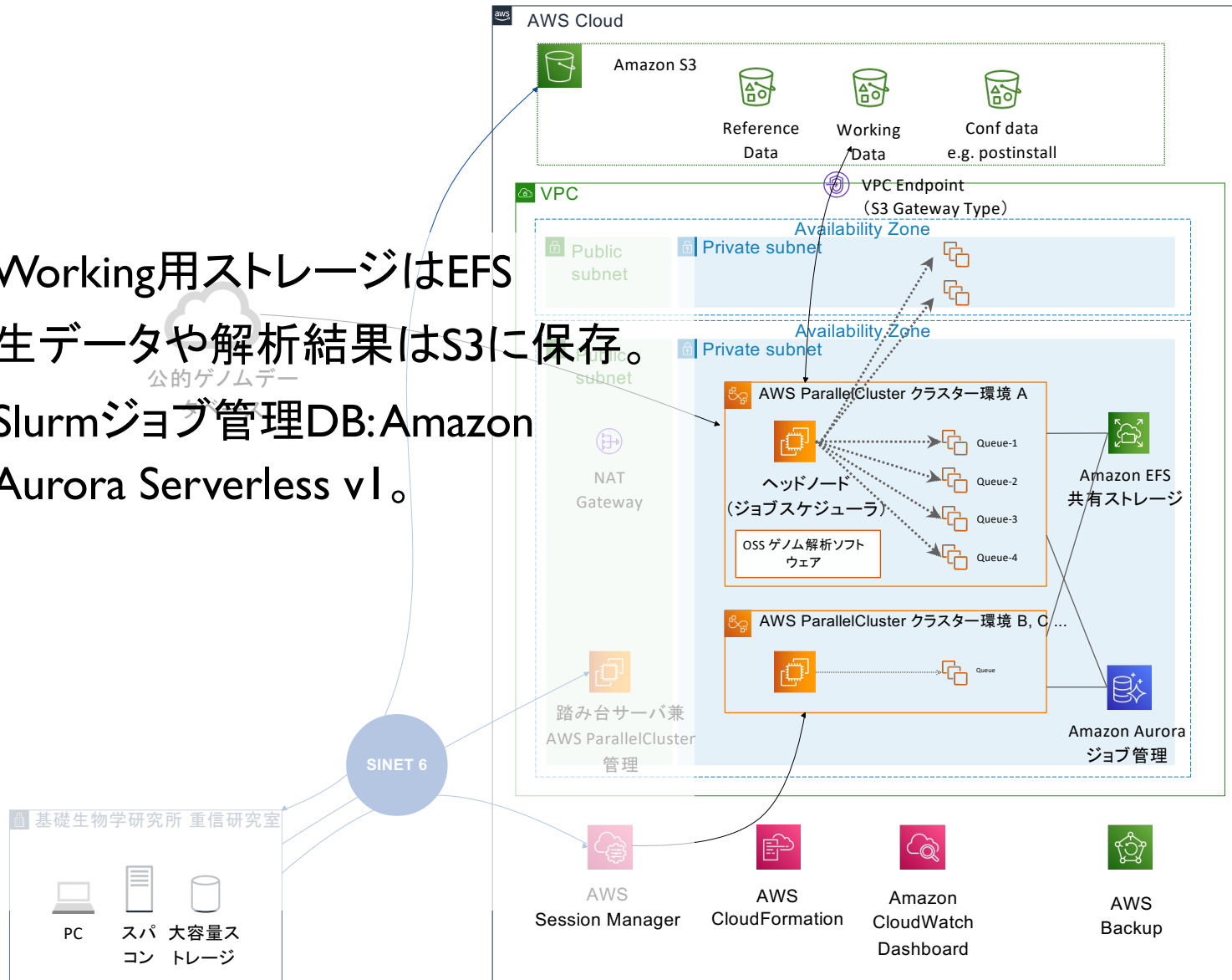
基礎生物学研究所 重信研究室

PC スパコン 大容量ストレージ

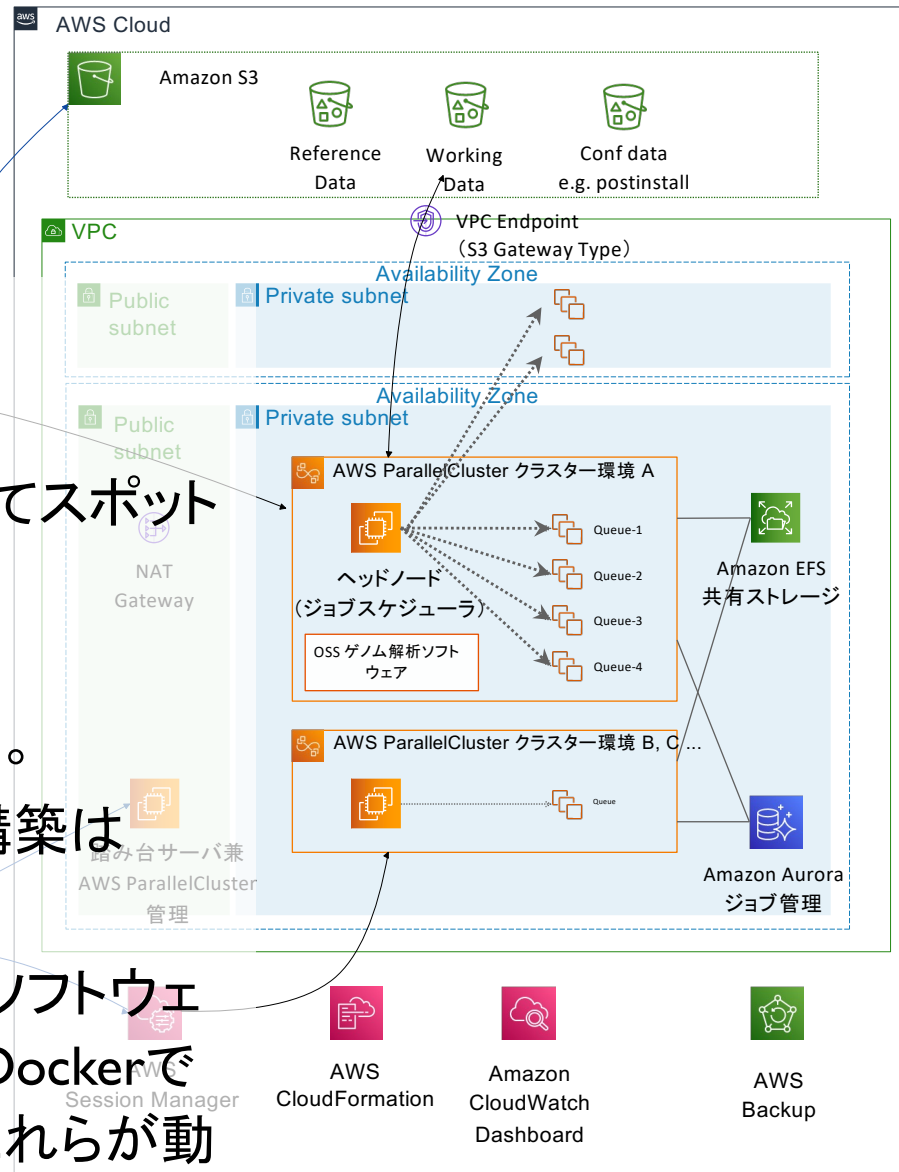
SINET 6

AWS Session Manager AWS CloudFormation Amazon CloudWatch Dashboard AWS Backup

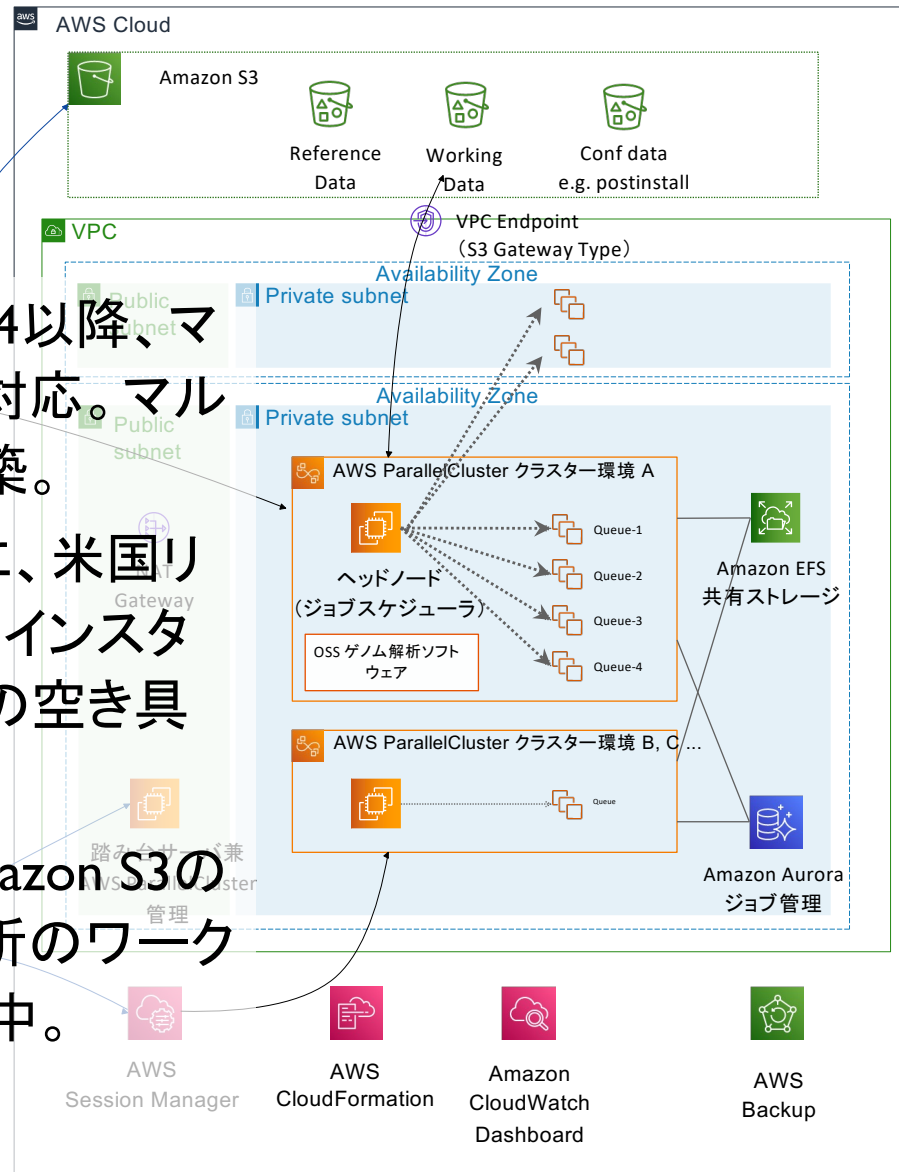
- ▶ Working用ストレージはEFS
- ▶ 生データや解析結果はS3に保存。
公的ゲノムデータ
- ▶ Slurmジョブ管理DB: Amazon Aurora Serverless v1。



- ▶ OS: Ubuntu 20.04 LTS
- ▶ CPU: intel系
公的ゲノムデータ
- ▶ 計算ノードは基本的に全てスポットインスタンス。
- ▶ AWSのリソース構築はCloudFormationで自動化。
- ▶ ヘッドノードの計算環境構築はAnsibleによって自動化。
- ▶ バイオインフォマティクスソフトウェアは、conda (Bioconda), Dockerで提供されることが多い。これらが動くように初期設定。



- ▶ AWS ParallelCluster がv3.4以降、マルチ Availability Zone に対応。マルチAZ対応のVPC環境構築。
- ▶ Tokyoリージョンをメインに、米国リージョンも併用。(スポットインスタンスのコンピュータードの空き具合と価格次第)
- ▶ 最近、Mountpoint for Amazon S3のサービス開始。ゲノム解析のワークフローに有用、導入試行中。



例：ヘビゲノムのBUSCO解析

シマヘビゲノムのBUSCO5の解析をParallelClusterで、5パターンで並列解析。
BUSCOとはゲノムアセンブリの完全性をユニバーサル遺伝子セットの網羅率で評価するソフトウェア。

- ▶ BUSCO解析をデータベースを5パターンで一気に入実施。
- ▶ BUSCOはconda環境でインストール済。(Dockerの利用も可能)
- ▶ インスタンスはm5d.8xlarge, l2x, l6x を選択。
- ▶ バッチジョブスケジューラ slurmにジョブを投げると、自動的に仮想マシンが立ち上がる。
- ▶ 経費節約のためにスポットインスタンスを利用した。
- ▶ I/Oの激しいジョブなので、インラインSSDを持つインスタンスタイプを選択。インラインSSDをマウントする/scratchで作業し、計算終了後、rsyncでEFS下のホームディレクトリに計算結果をコピー。
- ▶ 5パターンの解析は、それぞれ1h46m, 1h35m, 3h57m, 4h04m, 5h10mで計算終了。

```
#!/bin/bash

#SBATCH -p medium-m5d
#SBATCH -N 1
#SBATCH -n 16

INPUT=shimahebi_genome.fasta
MODE=genome
LINEAGE=metazoan #ここをvertebrataなど
異なるデータベースで合計5パターン
NCPUs=16

busco -i $INPUT \
      -l $LINEAGE \
      --out_path /scratch \
      -o $OUTF \
      -m $MODE \
      -c $NCPUs \

rsync -avh /scratch/$OUTF ./
```

AWS ParallelCluster 個人的使用感

- ▶ オンプレミスのバッチジョブとほぼ同じ感覚で使用できる。
- ▶ オンプレミスの共用サーバと異なり、ユーザーが一人でクラスターを占有できる。ジョブの待ち時間がゼロ。スパコン級環境独占の圧倒的快適さ。
- ▶ 計算機資源(インスタンス、ストレージ)は事実上無限大。
- ▶ 仮想マシン上に理想的な計算環境を自らがroot権限でスクラップ & ビルドすることができる、高い自由度により、フットワーク軽く、研究が進められる。
- ▶ 高性能なEC2インスタンスは結構高価。しかしスポットインスタンスを利用することでコストを数分の1に抑えられる。
- ▶ 最近、スポットインスタンスが以前より混んできた。対策として、マルチAZ対応や海外リージョン利用。
- ▶ 実行するソフトウェアの特性とAWSの料金体系を理解してインスタンスの選択やバッチジョブスクリプトを最適化する必要がある。
- ▶ AWSの知識が必須。ユーザーがAWS力を高めていく必要がある。

PI（研究室主宰者）目線でクラウドHPCの 管理運営面でのメリット

- ▶ 研究時間捻出：マシンのメンテナンスなどの雑用から解放され、研究そのものに集中することができる。
- ▶ 柔軟なコスト管理：使用した分のみ支払い。=>研究費に制限のある小さな研究室でも気軽にHPCが利用できる。
- ▶ 調達手続き不要：高額なコンピュータをオンプレで購入する場合、入札等で導入まで数ヶ月以上かかる上、手続きが面倒。
- ▶ セキュリティ：自分で計算機を管理するより確実・安全。（もちろん十分な知識と対策は必要）。
- ▶ オンプレの共用スパコンを否定するものではない。効果的なハイブリッド活用。

Before



AWS Parallel Cluster



After

- ▶ 共用のオンプレミスHPCではジョブ待ちが頻繁。
- ▶ OSやライブラリのバージョンが固定。ソフトウェアが動かない場合、対処が困難。
- ▶ オンプレミスサーバは調達に時間がかかり、手続きも煩雑。メンテナンスなどの運用負担。
- ▶ HPCの計算機リソースを必要な時に必要な分を即時確保。
- ▶ 理想的な計算環境を、root権限でスクラップ&ビルド。
- ▶ 事務手続きやコンピュータメンテナンスから解放されて、研究そのものに集中。

- HPCはクラウドで独占型スパコンを実現可能。快適。
- ゲノム科学や基礎科学に特有のニーズにも柔軟に対応。

ゲノム情報をクラウドを介して共有する

HPCで得られたゲノム解析結果をどのように共有するか？

- ▶ 「ゲノムブラウザ」を介して、共同研究者とデータを共有する。
- ▶ 「ゲノムブラウザ」の構築にはクラウドが最適。

ゲノムブラウザの例

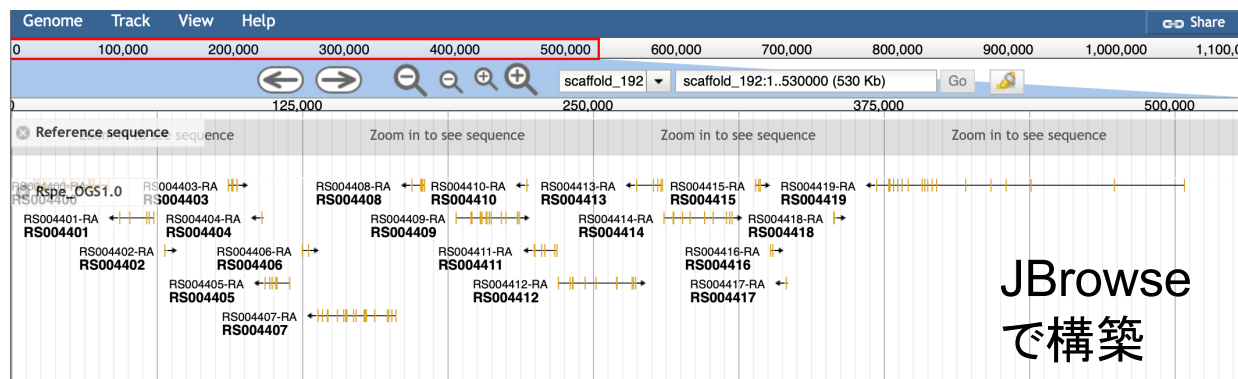


2022年1月19日

高度な社会性を持つシロアリのゲノム情報を解読
～遺伝子重複が社会性進化の原動力であることを明らかに～



Shigenobu et al., 2022 PNAS



公開ゲノムブラウザ <http://www.termite.nibb.info/retsp/>
(AWS EC2, Route 53 等を利用)

ゲノムブラウザ

- ▶ ゲノムブラウザ (Genome Browser) は、ゲノムの情報を視覚的に閲覧・検索するためのツールやウェブサイト。
- ▶ 遺伝子構造に関する情報や、その遺伝子に関連するアノテーション情報 (機能や発現データなど) を閲覧。
- ▶ 公開ゲノムブラウザの例としては、UCSC Genome Browser や Ensembl など。これらのブラウザはウェブベースでアクセスできる。
- ▶ 自前で構築するためのツールとしては、JBrowse2 と IGV (Integrative Genome Viewer) が有名。

理想のゲノムブラウザはGoogle Mapのゲノム版

Google Map

- ▶ Basic geological information
- ▶ Zoom in/out : multiple scale
- ▶ Locate any objects (buildings, railways, road)
- ▶ Many attributes
- ▶ Searchable
- ▶ Overlap other information
- ▶ Link to outside data

Ideal Genome Browser


- ▶ Gene structure information
- ▶ Zoom in/out : multiple scale
- ▶ Locate any objects (gene, repeat, SNP)
- ▶ Many attributes
- ▶ Searchable
- ▶ Overlap other omics data
- ▶ Link to outside data

SOFTWARE

Open Access



JBrowse: a dynamic web platform for genome visualization and analysis

Robert Buels¹, Eric Yao¹, Colin M. Diesh², Richard D. Hayes^{3,6}, Monica Munoz-Torres³, Gregg Helt^{3,4}, David M. Goodstein^{3,6}, Christine G. Elsik², Suzanna E. Lewis³, Lincoln Stein^{5,7} and Ian H. Holmes^{1,3*} 

Abstract

Background: JBrowse is a fast and full-featured genome browser built with JavaScript and HTML5. It is easily embedded into websites or apps but can also be served as a standalone web page.

Results: Overall improvements to speed and scalability are accompanied by specific enhancements that support complex interactive queries on large track sets. Analysis functions can readily be added using the plugin framework; most visual aspects of tracks can also be customized, along with clicks, mouseovers, menus, and popup boxes. JBrowse can also be used to browse local annotation files offline and to generate high-resolution figures for publication.

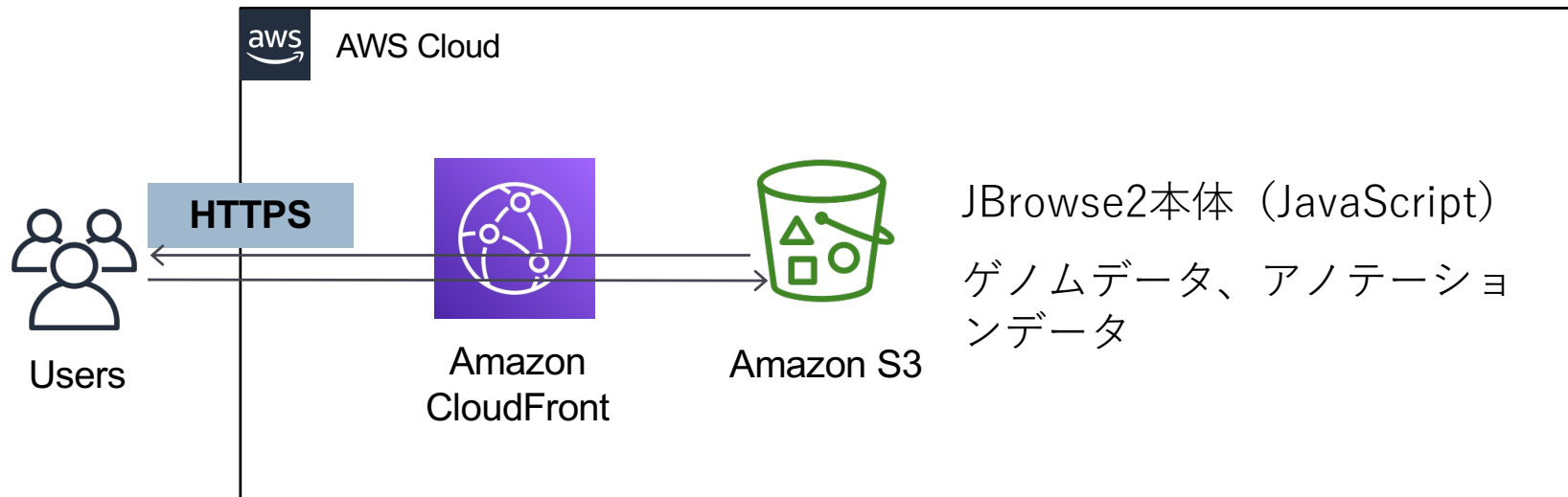
Conclusions: JBrowse is a mature web application suitable for genome visualization and analysis.

Keywords: Genome, Browser, Bioinformatics

<http://jbrowse.org/>

新規ゲノムブラウザをJBrowse2を用いて AWS上にサーバレスに構築する

- ▶ JBrowse2はJavaScriptで書かれている。リレーショナルデータベースは不要。
- ▶ S3 + CloudFront のサーバレスアーキテクチャーで実装



ゲノムブラウザをAWS上にサーバレスに構築するメリット

- ▶ サーバレスなのでサーバの管理が不要。
- ▶ コストはデータ保存料とデータ転送料のみで安価。
- ▶ CloudFrontの機能により、高負荷耐性、世界中へ高速配信。
- ▶ 最近では大学や研究機関のネットワークセキュリティが強化され、内部ネットワークに公開サーバを設置することがNGであることがほとんど。

例) 当研究室で解読したアブラムシゲノムのゲノムブラウザ

The screenshot displays the JBrowse genome browser interface. The main window shows a genomic track for the contig `ptg0000181` with a zoomed-in view of a 525Kbp region. The track includes various gene models such as `AL4ok.liftover`, `LOC100574563`, `LOC103309468`, `ACYP153291`, `LOC107884783`, `LOC100158757`, `ACYP1087002`, `LOC100574649`, `LOC103309466`, `Homeotic protein distal-less`, `protease filzig`, `Ribosome quality control complex subunit NEMF homolog`, `Homeobox domain-containing protein`, `Uncharacterized protein`, `Actin-like protein 6A`, `C2H2-type domain-containing protein`, and `Uncharacterized protein`. A detailed view of the `Homeotic protein distal-less` match is shown on the right.

Feature details

HOMEOTIC PROTEIN DIS... - MATCH

Core details

Position	ptg0000181:7,958,391..7,977,847 (-)
Name	Homeotic protein distal-less
Length	19,457
Type	match

Attributes

source	exonerate:protein2genome:local
phase	0
id	P20009
target	P20009 63 172
gap	M30 D3 M96 D19125 M203
identity	77.98
similarity	80.73
dbxref	EMBL:P20009 Uniprot:DLL_DROME tax
alias	Dll

[SHOW FEATURE SEQUENCE](#) ?

SUBFEATURES

MATCH_PART

Core details

Position	ptg0000181:7,958,391..7,958,593 (-)
Length	203
Type	match_part

Attributes

source	exonerate:protein2genome:local
phase	0

HPCで得られたゲノム解析結果をどのように共有するか？

- ▶ 「ゲノムブラウザ」を介して、共同研究者とデータを共有する。
- ▶ 「ゲノムブラウザ」の構築にはクラウドが最適。
- ▶ 配列検索ツール「BLAST」のニーズが高い。
- ▶ BLASTをサーバレス化した。



Blast-help > Blast Searches at a Cloud Provider > ElasticBlast

ElasticBLAST

Overview

Requirements

Quickstart

Tutorials

IAM Policy

Budget

Commands

Configuration variables

Exit codes

Tips for GCP

Limiting search by
taxonomy

Support

Troubleshooting

Known issues on AWS

Known issues on GCP

ElasticBLAST, version 1.1.0

ElasticBLAST is a cloud-based tool to perform your BLAST searches faster and make you more effective.

ElasticBLAST is ideal for users who have a large number (thousands or more) of queries to BLAST or who prefer to use cloud infrastructure for their searches. It can run BLAST searches that cannot be done on [NCBI WebBLAST](#) and runs them more quickly than stand-alone [BLAST+](#).

ElasticBLAST speeds up your work by distributing your searches across multiple cloud instances. The ability to scale resources in this way allows large numbers of queries to be searched in a shorter time that you normally could with BLAST+.

The National Center for Biotechnology Information ([NCBI](#)), part of the National Library of Medicine at the NIH, develops and maintains ElasticBLAST.

ElasticBLAST status: beta

Platforms available: AWS, GCP (account required)

Camacho *et al.* *BMC Bioinformatics* (2023) 24:117
<https://doi.org/10.1186/s12859-023-05245-9>

BMC Bioinformatics

SOFTWARE

Open Access

ElasticBLAST: accelerating sequence search via cloud computing



Christiam Camacho, Grzegorz M. Boratyn, Victor Joukov, Roberto Vera Alvarez and Thomas L. Madden*

Our Light-weight Serverless BLAST

BLAST search: *Cladonema pacificum*

BLAST search condition

Database: protein

Program: blastp -- protein > protein

Query: >test
MRVITTAIVPLISLLVSEAVWMTLGTQLVSSDPAKKSMLCYVTRQFNSRQOEVCRNPDLMHVHAGAKYGVHCEHRHOFNRRMNCSTIRESGLFESVLSKGCRE
AAFVHAVTAAGVAHSVTDACSKGRIESDCDCRNLGSRSSKGTWWSGCNSNIKFGVWFSKQFTEARERGDLLRQINWRHNSRAGRKALEELVWRKYKCHGLSGSCSMKT
CWMQANFRQIGDHLKVKYDSAVEMTTKVNRRGKRLPKPYSHFKKPSDKDLIYFETSPNYCDKNVTVGLSGTRGCRNYTNGIDGCELLCCGRGHNIQQAQITRNC

Format: 7: table with comments

E-value: 0.001

submit

BLAST Result:

Result will appear here.

© 2023 Copyright: [National Institute for Basic Biology](#). All rights reserved.

入力画面

BLAST search: *Cladonema pacificum*

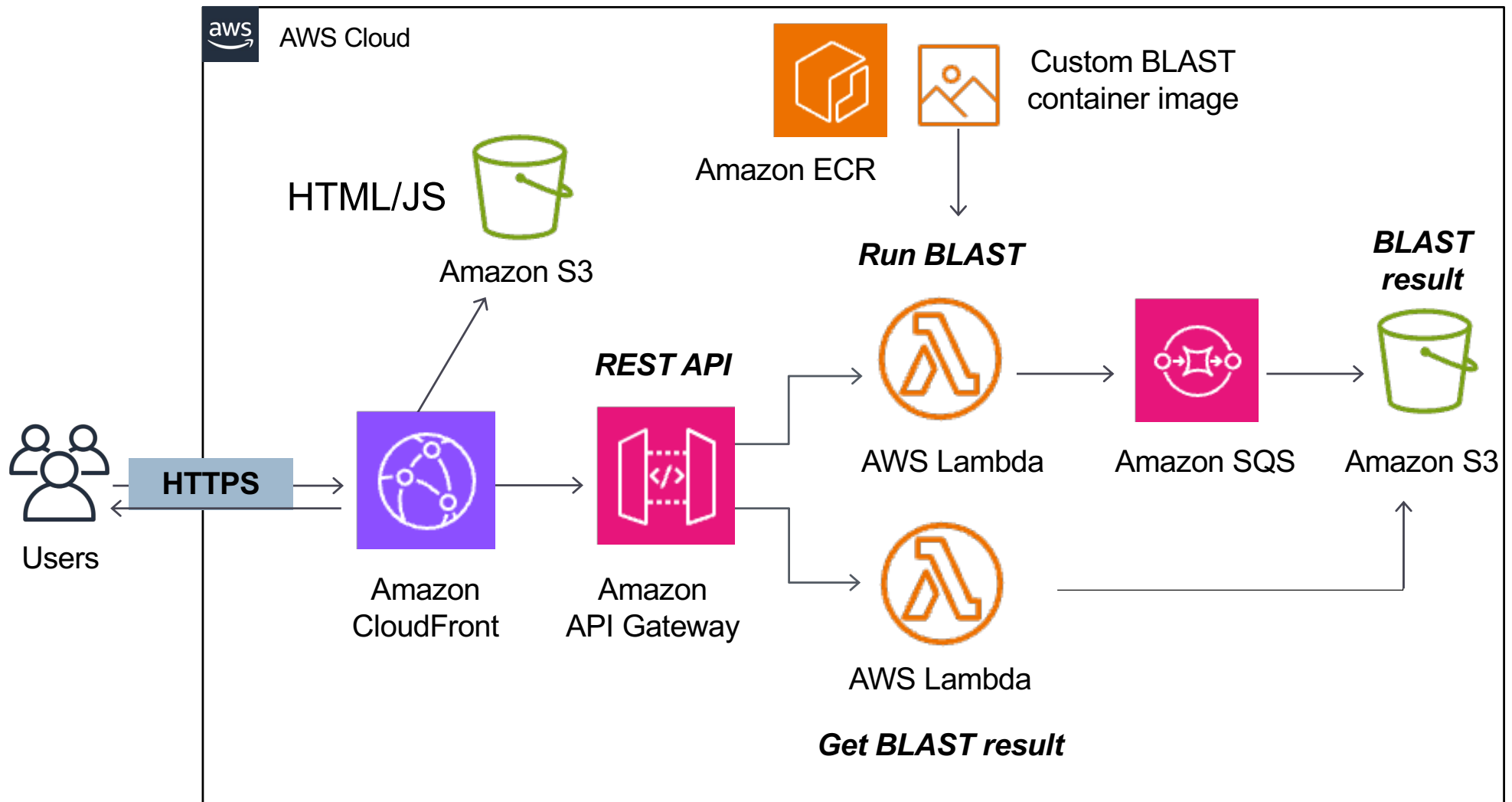
BLAST search condition

BLAST Result:

```
# BLASTP 2.14.0+
# Query: test
# Database: /var/task/blastdb/Cladonema_pacificum.fasta
# Fields: query acc.ver, subject acc.ver, % identity, alignment length, mismatches, gap opens, q. start, q. end, s. start, s. end, e-value
# 32 hits found
test Cp2_g9527.t1 42.241 348 186 5 7 341 10 355 4.59e-100 300
test Cp2_g9527.t2 44.444 297 156 2 55 342 1 297 2.39e-95 286
test Cp2_g28444.t1 39.118 340 192 5 8 340 5 336 1.16e-80 252
test Cp2_g16735.t1 39.118 340 192 5 8 340 5 336 1.20e-80 252
test Cp2_g8900.t1 37.391 345 203 8 4 340 3 342 7.62e-80 248
test Cp2_g16735.t2 39.130 345 197 6 1 340 18 354 6.40e-79 248
test Cp2_g28444.t2 39.130 345 197 6 1 340 18 354 6.97e-79 248
test Cp2_g8900.t2 41.812 287 158 6 62 340 1 286 6.33e-76 236
test Cp2_g21666.t1 31.157 337 210 7 21 340 6 337 3.62e-57 189
test Cp2_g9436.t1 32.424 330 198 9 24 339 14 332 1.43e-53 180
test Cp2_g25205.t2 39.527 296 130 7 76 341 717 993 2.30e-53 189
test Cp2_g13701.t1 30.769 364 210 10 2 339 8 355 7.35e-51 173
test Cp2_g1833.t1 30.769 364 210 10 2 339 8 355 1.10e-50 173
test Cp2_g13701.t3 30.769 364 210 10 2 339 8 355 2.37e-50 173
test Cp2_g13701.t2 30.769 364 210 10 2 339 8 355 9.17e-50 174
test Cp2_g8512.t2 29.345 351 225 9 3 342 2 340 2.91e-47 164
test Cp2_g8512.t3 31.834 289 179 8 65 342 48 329 8.38e-47 162
test Cp2_g25205.t1 38.554 249 128 5 74 302 799 1042 2.15e-41 155
test Cp2_g21666.t2 36.318 201 115 3 151 340 25 223 1.84e-40 142
test Cp2_g16261.t2 29.070 308 203 6 41 335 43 350 3.14e-35 132
test Cp2_g16261.t1 29.866 298 198 5 49 335 6 303 3.63e-35 131
test Cp2_g8656.t1 29.448 326 171 11 65 338 58 376 2.07e-30 119
test Cp2_g17158.t1 26.648 364 236 11 4 340 1 360 3.57e-30 119
test Cp2_g8656.t2 29.538 325 170 12 66 338 1 318 4.68e-30 117
test Cp2_g17158.t2 27.301 326 211 9 38 340 19 341 4.70e-28 112
test Cp2_g19192.t1 28.283 297 183 12 68 341 42 331 7.82e-28 112
test Cp2_g25970.t2 28.283 297 183 12 68 341 42 331 8.83e-28 111
test Cp2_g19192.t2 28.521 284 177 11 77 341 23 299 1.06e-27 110
test Cp2_g25970.t1 28.531 284 177 11 77 341 23 299 1.23e-27 110
```

結果画面

Our light-weight serverless BLAST



BLASTはNCBIよりUNIX CUI環境で動くバイナリが提供されている。これをDockerコンテナ化した上でAWS Lambda化。フロントエンドはJavaScriptで実装。ユーザーがウェブインターフェイスからquery配列を入力すると、非同期にBLAST実装Lambdaが起動し、計算が終わったら、結果が画面に表示される。

HPCで得られたゲノム解析結果をどのように共有するか？

- ▶ 「ゲノムブラウザ」を介して、共同研究者とデータを共有する。
- ▶ 「ゲノムブラウザ」の構築にはクラウドが最適。
- ▶ 配列検索ツール「BLAST」のニーズが高い。
- ▶ BLASTをサーバレス化した。

ゲノム解析後のゲノム情報の共有やツール提供にはクラウドが有効。開発後のメンテナンス面を考慮すると、サーバレスアーキテクチャを目指すべきである。

まとめ

ゲノム科学、特に基礎研究、におけるクラウド活用

▶ 計算資源(HPC)の確保

- ▶ HPCはクラウドで独占的スパコンを実現可能。快適。
- ▶ ゲノム科学や基礎科学に特有のニーズにも柔軟に対応。

▶ データ共有による共同研究促進

- ▶ ゲノムブラウザ等ゲノム情報の共有やツール提供にはクラウドが有効。
- ▶ サーバレスアーキテクチャを目指すべき。



謝辞：基礎生物学研究所・TSBセンター、進化ゲノミクス研究室
科研費・倉谷新学術ほか
クラスメソッド、アマゾンAWS