

Breakout Session 5:

PII-secured AWS Computing Environment (PACE)- Some lessons learned working with PII in the STRIDES AWS Environment

Dr. Daniel Veltri

Health Scientist (Data Science), NIH/NIAID

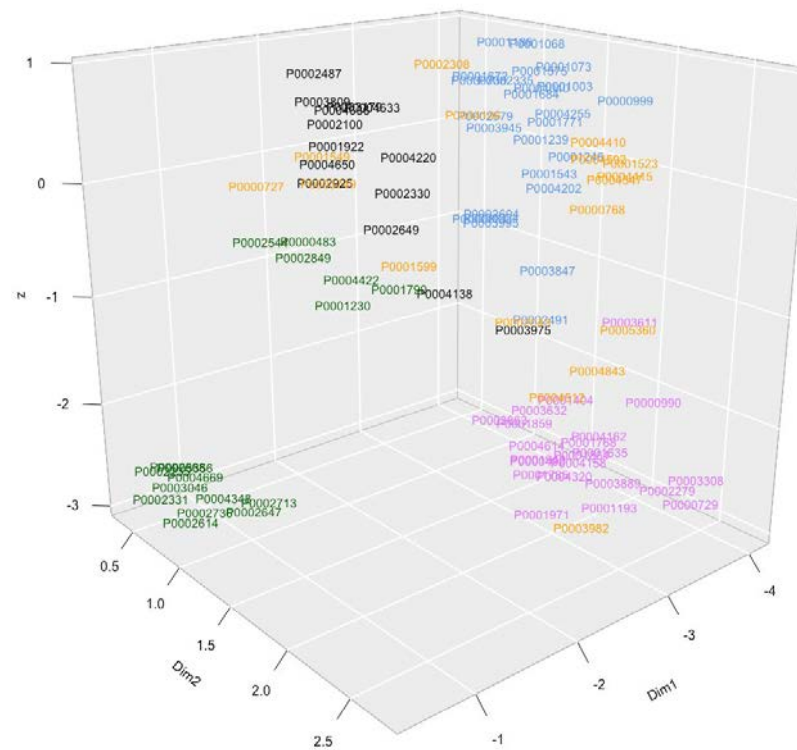
PII-secured AWS Computing Environment (PACE)

Some lessons learned working with PII in the STRIDES AWS Environment

Daniel Veltri, Ph.D.

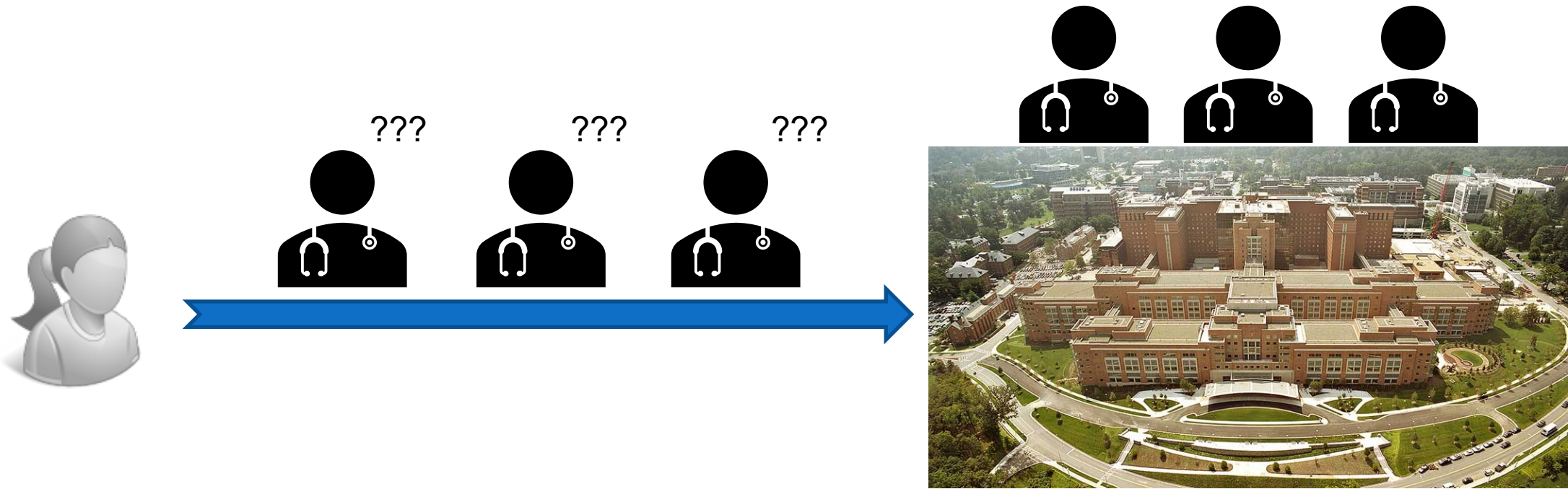
Bioinformatics and Computational Biosciences Branch (BCBB)

OCICB/OSMO/OD/NIAID/NIH



Leiden Clustering of 95 Juvenile Rare Disease Patients

Motivation: Diagnosing patients with rare diseases



NIH Clinical Center



NIAID Genomic Research Integration System (GRIS)

The screenshot shows the NIAID Genomic Research Integration System (GRIS) homepage. At the top left is the NIH logo and the text "Genomic Research Integration System". To the right are links for "System Availability", "Support", "Feedback", "Log In (Registered NIAID users)", and "Sign Up". Below these is the text "CRIMSON Data Updated M-F 12:00AM-6:00AM". The main heading is "Integrate Analyze Discover" in large yellow letters, with the subtitle "From recording to reporting, a comprehensive solution for clinical genomics research." Below this are two icons: a pedigree chart icon labeled "PhenoTips" and a magnifying glass over a DNA helix icon labeled "seqr". At the bottom, there are two descriptive sentences: "Browse, search, and record patient, pedigree, phenotypic, medical history, biospecimen, and genotype data." and "Perform Candidate Gene analysis".

NIAID GRIS Team



NIAID Central Sequencing Program



NIAID Genomic Research Integration System (GRIS)



The screenshot shows the NIAID Genomic Research Integration System (GRIS) website. At the top left is the NIH logo and the text "Genomic Research Integration System". To the right are links for "System Availability", "Support", "Feedback", "Log In (Registered NIAID users)", and "Sign Up". The main content area is titled "Clinical symptoms and physical findings" and is divided into two sections: "IMMUNE SYSTEM" and "INFECTIONS".

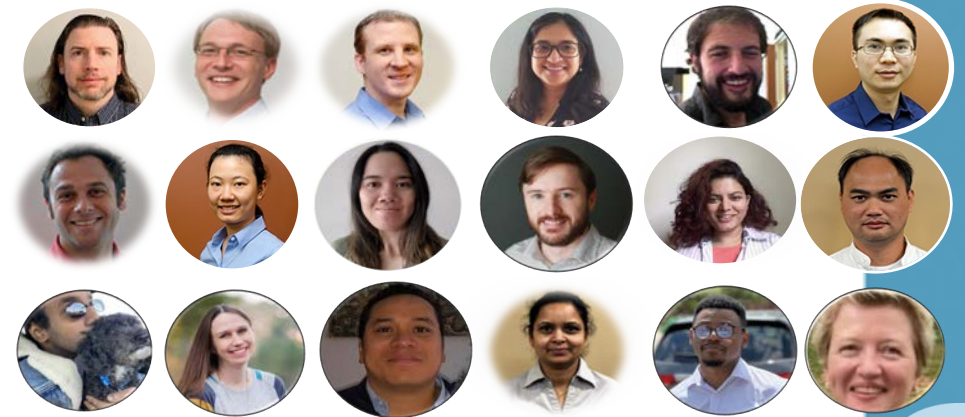
IMMUNE SYSTEM

- Immunologic hypersensitivity
 - Anaphylactic shock
 - Asthma
- Sinusitis
- Atopic dermatitis
- Pneumonia
- Allergic rhinitis
- Allergic conjunctivitis

INFECTIONS

- Recurrent infections
 - Recurrent otitis media

NIAID GRIS Team



NIAID Central Sequencing Program



How could Machine Learning (ML) help us?

- Promote faster diagnoses for patients
- Reduce time burden on researchers to manually curate symptoms (phenotypes)
- More consistent phenotype curation between labs extracted from a large collection of clinical notes

Working with the Lu Lab at NIH/NCI

- We assisted the Lu Lab in evaluating two of their new natural language processing deep learning models on GRIS patient clinical notes: PhenoTagger¹ and PhenoRerank²
- We would like a space to continue such collaborations and development of new tools, but this is complicated by the presence of PII - *major restrictions on most systems!*
- We wanted to leverage ODSS funding and STRIDES to create a secure but collaborative environment for clinical note evaluation
 - *We did it, but it took a long time to get authorized!*



Zhiyong Lu, Ph.D. FACMI FIAHSI
Deputy Director for Literature Search, NCBI
Senior Investigator, NLM

Challenges obtaining an Authority to Operate (ATO)

- At the time we started, we were the first NIAID project to try out STRIDES, and the first NIAID program to consider PII in the cloud
- STRIDES was still finalizing their AWS environment, security controls, and awaiting their own ATO
- Throughout, we have had to juggle cybersecurity requirements put on us from both CIT/STRIDES and NIAID
- While we inherit a majority of security controls from STRIDES, they do not allow themselves to be considered a “parent system”
 - This was a major conundrum for NIAID Cybersecurity
 - We ended up creating our own PACE security boundary
 - The privacy impact of our system as “Moderate”

ATO challenges continued ...

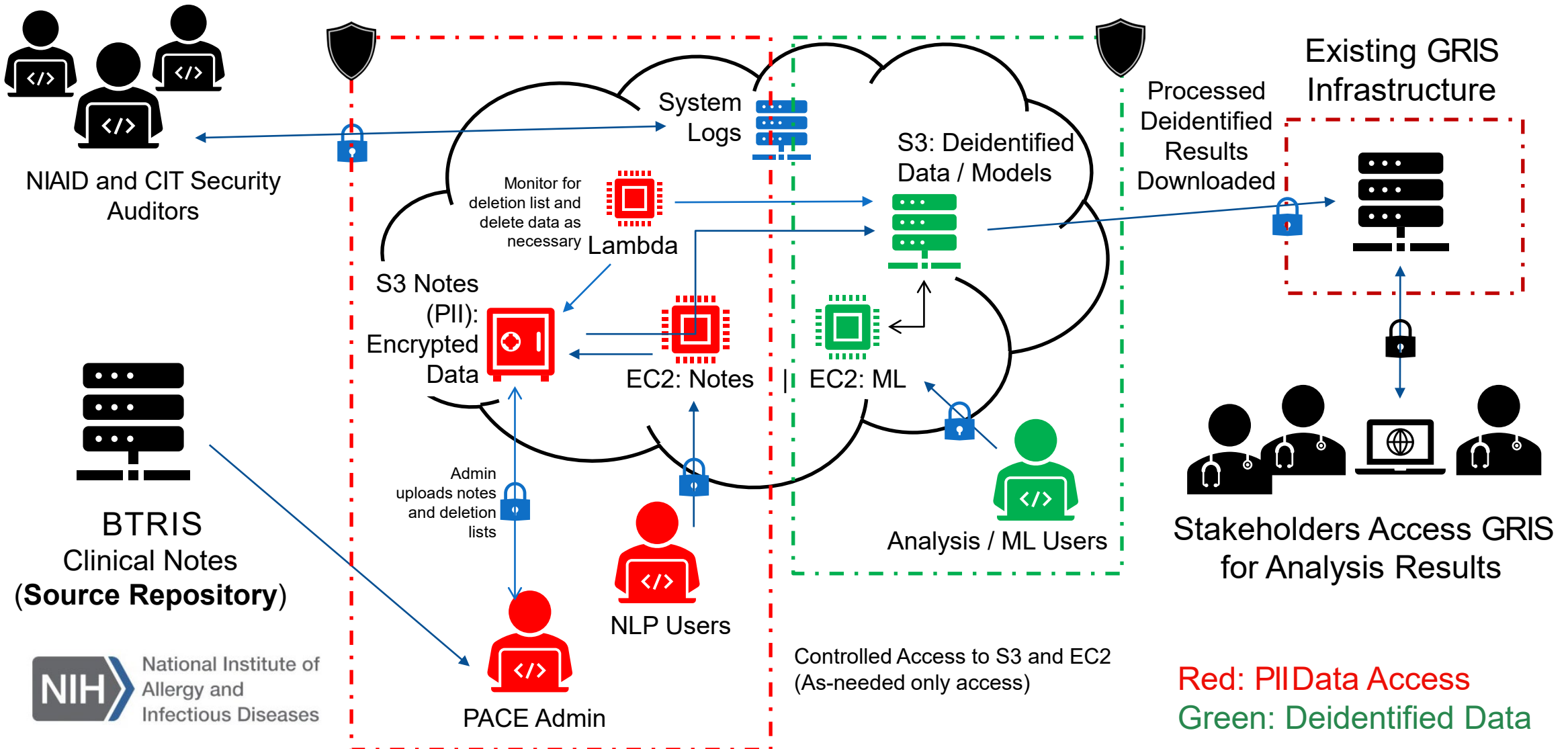
- NIAID required us to use specific software for security logging – the STRIDES team had to come up with a novel way to support it
- Our project also needed role-specific accounts to limit PII access – the STRIDES team had to come up with a novel implementation for this
- Midway, NIAID underwent major changes in cybersecurity leadership
- NIST released a major new cybersecurity controls revision after we had started on the old one

ATO challenges still continued ...

- STRIDES does not have an “official” recommended security scanning solution for AWS – NIAID’s current solutions were also not a great fit for us (we ended up getting approved to use AWS Inspector)
- In terms of setting up the environment roles to access resources – both *STRIDES* and the *AWS Enterprise Support teams* were *amazing help!*

We did manage to finally get ATO in 2023!

High level look at PACE architecture



Some lessons learned ...

- Dealing with PII? Avoid being a first-mover, wait for environments to be established and tested first. *It is complicated enough without PII!*
- Prior to starting – have a rock solid understanding of the *parent system* and the security boundaries you will use. *I should have asked more questions up front*
- Have a cloud expert in house - even for a simple system like this one, the STRIDES and AWS support teams can only do so much for you. *Ultimately, the design, implementation, and maintenance is on you*
- Does your group have other cloud environments or vendors? Think carefully how you can keep synchronized
 - Can you standardize the implementation, security scanning, and reporting procedures?
 - This could be a challenge for PACE moving forward in STRIDES as most of NIAID uses the Monarch Platform

Major Thanks!



Sandhya Xirasagar, Ph.D.
Initial PM for GRIS and huge
help to the start of this project



Giovanni Borjas and Metasebia Gizaw
Two incredible ISSOs that helped me navigate the
difficult ATO process



Andrew Kulak
Incredible support with AWS
environment configuration and
general questions for STRIDES

NIAID/OCICB: Dr. Darrell Hurt, Dr. Andrew Oler, Paul Suh (NIAID CISO), and Mike Tartakovsky

CIT/STRIDES: Nick Weber and James Davis

AWS: Tom Fonseca and Gargi Singh

NIAID Central Sequencing Team: Morgan Similuk, ScM and Dr. Rajarshi Ghosh

NLM Lu Lab: Drs. Zhiyong Lu, Ling Luo, and Shankai Yan, Lai Po-Ting, Chih-Hsuan Wei, and Robert Leaman

NIH BTRIS Team: Dr. Michael Ring and Andrea Beri

A big thanks to NIH/ODSS and STRIDES for supporting this project!

National Institute of Allergy and Infectious Diseases

Developing a Scalable and Reusable Framework for State-of-the-Art Structural Variant Calling of Whole Genome Sequencing Data



Andrew Oler, Ph.D.; Eric Karlins, M.S.; and Daniel Veltri, Ph.D.

Bioinformatics and Computational Biosciences Branch (BCBB)

OCICB/OSMO/OD/NIAID/NIH

NIAID



National Institute of
Allergy and
Infectious Diseases

NIH/ODSS Cloud Supplement Program PI Meeting
HVD21 Program Awardee | January 18th, 2024

Motivation: Looking for Structural Variants (SV) in genomes

- SVs comprise a range of large genomic rearrangements (> 50bp) and imbalances that can play significant roles in a variety of diseases and important for diagnosing patients
- This project uses NIH/ODSS funding and STRIDES to help run the Broad Institute's state-of-the-art GATK-SV pipeline on the Terra Platform to look for SVs in whole genome samples
- This was pilot initially started for Covid-19 patients but has now expanded to include general GRIS patients whenever whole genome sequencing data is available

Broad Institute GATK-SV Pipeline on Terra

Article

A structural variation reference for medical and population genetics

<https://doi.org/10.1038/s41586-020-2287-8>

Received: 2 March 2019

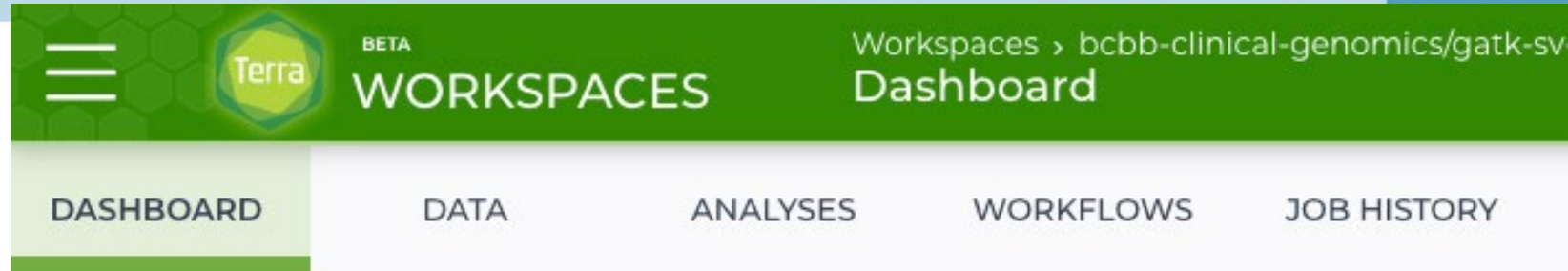
Accepted: 31 March 2020

Published online: 27 May 2020

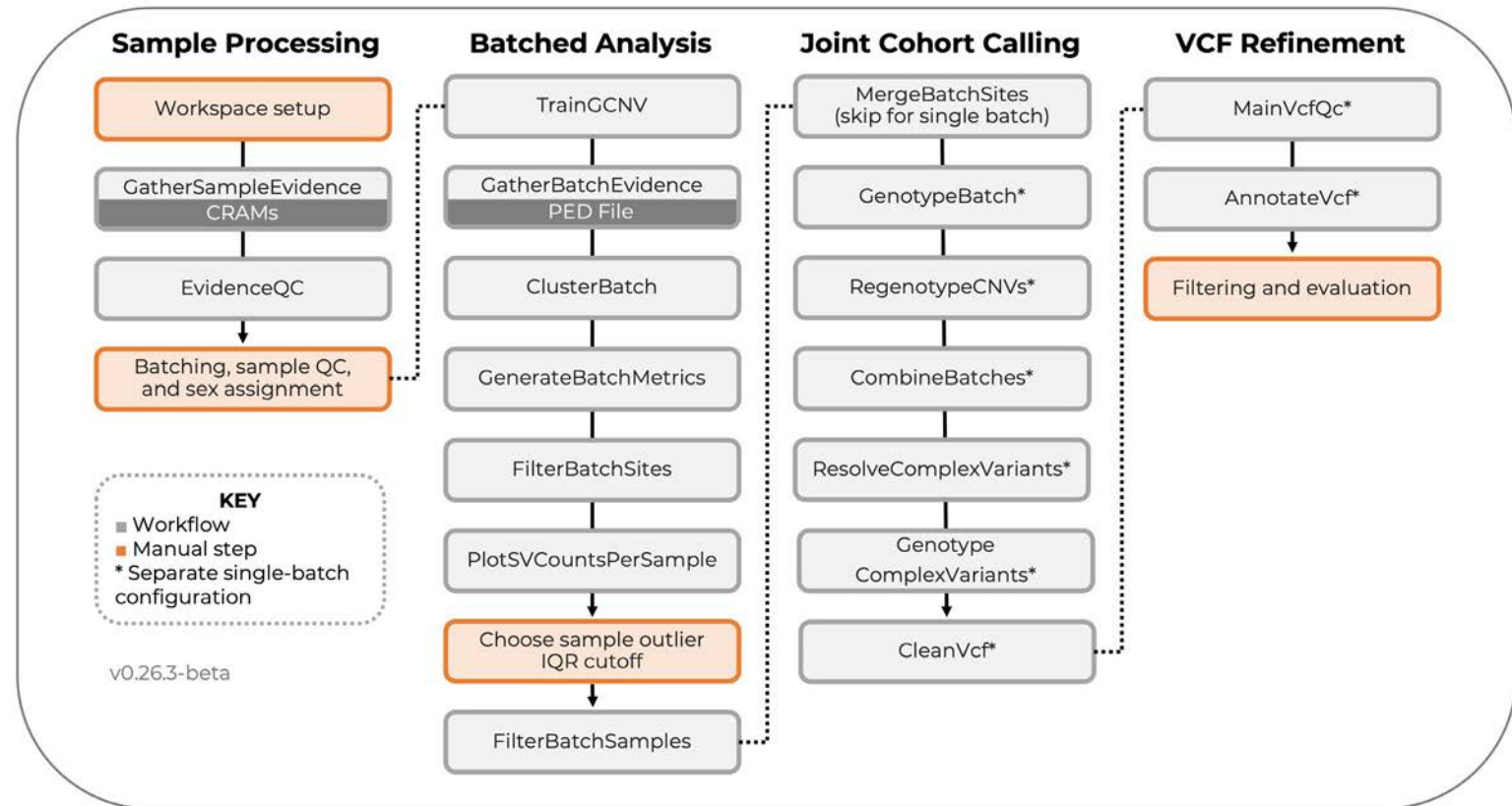
Open access

Check for updates

Ryan L. Collins^{1,2,3,106}, Harrison Brand^{1,2,4,106}, Konrad J. Karczewski^{1,5}, Xuefang Zhao^{1,2,4}, Jessica Alfoldi^{1,5}, Laurent C. Franciol^{1,5,6}, Amit V. Khera^{1,2}, Chelsea Lowther^{1,2,4}, Laura D. Gauthier^{1,7}, Harold Wang^{1,2}, Nicholas A. Watts^{1,5}, Matthew Solomonson^{1,5}, Anne O'Donnell-Luria^{1,5}, Alexander Baumann⁷, Ruchi Munshi⁷, Mark Walker^{1,7}, Christopher W. Whelan⁷, Yongqing Huang⁷, Ted Brookings⁷, Ted Sharpe⁷, Matthew R. Stone^{1,2}, Elise Valkanas^{1,2,3}, Jack Fu^{1,2,4}, Grace Tiao^{1,5}, Kristen M. Laricchia^{1,5}, Valentin Ruano-Rubio⁷, Christine Stevens¹, Namrata Gupta¹, Caroline Cusick¹, Lauren Margolin¹, Genome Aggregation Database Production Team⁸, Genome Aggregation Database Consortium⁸, Kent D. Taylor⁹, Henry J. Lin⁹, Stephen S. Rich⁹, Wendy S. Post¹⁰, Yii-Der Ida Chen⁹, Jerome I. Rotter⁹, Chad Nusbaum¹⁰⁶, Anthony Philippakis⁷, Eric Lander^{11,12}, Stacey Gabriel¹, Benjamin M. Neale^{1,2,3,13}, Sekar Kathiresan^{1,2,6,14}, Mark J. Daly^{1,2,3,13}, Eric Banks⁷, Daniel G. MacArthur^{1,2,5,6,106,107} & Michael E. Talkowski^{1,2,4,13,12}



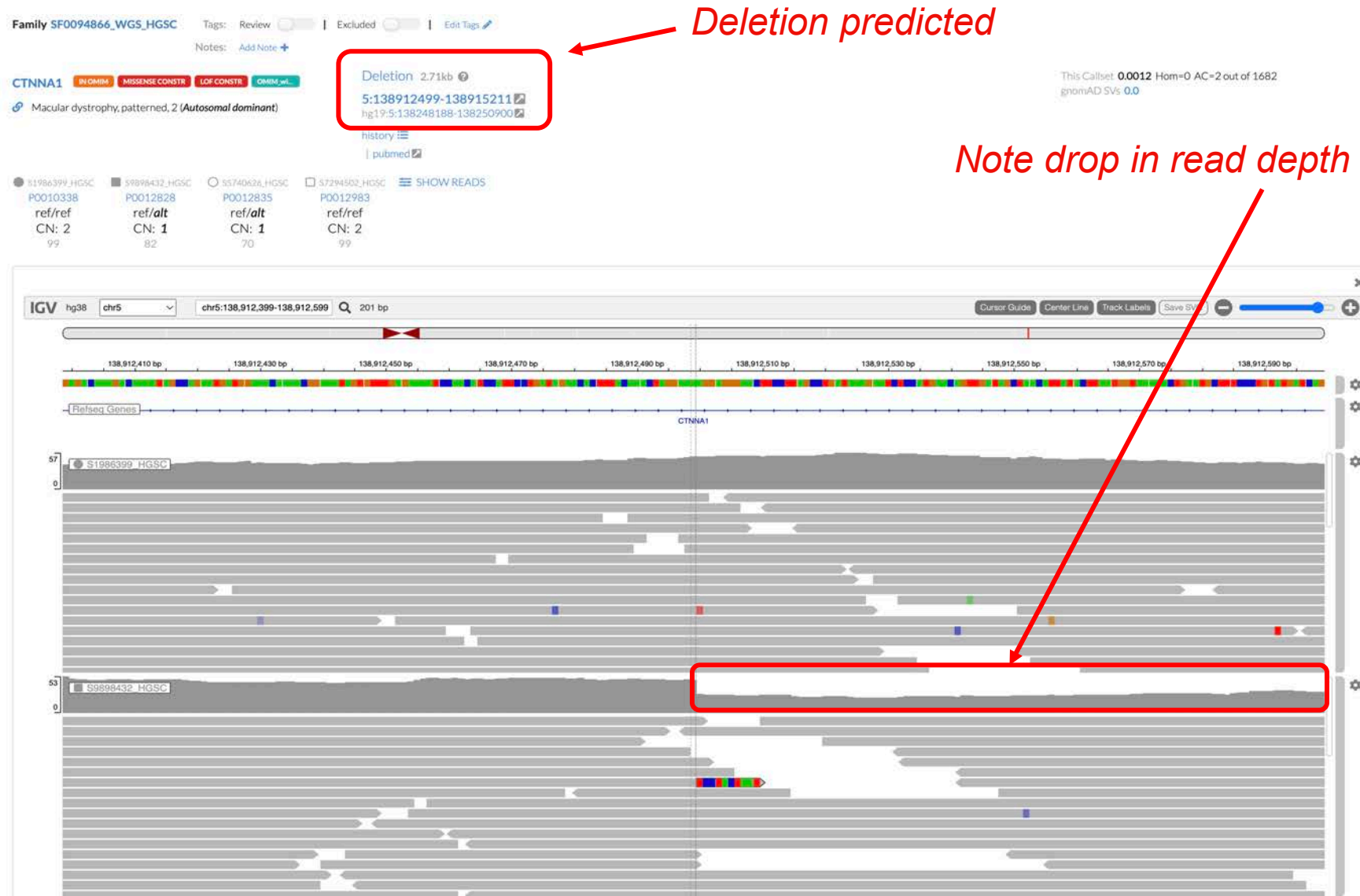
GATK-SV Cohort Mode in Terra



The same pipeline used for gnomAD-SV

<https://www.nature.com/articles/s41586-020-2287-8>

Example of an SV result in the Broad *seqr* program



Summary

- We have run the pipeline on >550 Covid-19 samples and >5000 GRIS patient samples
- Compute costs are in the \$3-5 range per sample
- Customized filtering of the results has been important to adjust sensitivity and specificity for the analysis – finalized protocols are still being developed
- This pipeline has assisted with the successful diagnosis of multiple patients enrolled in the NIAID Central Sequencing Program

Thanks

We would like to thank the NIAID GRIS Team for helping process and implement results and the NIAID Central Sequencing Program for all their help in testing and evaluation. Special thanks to NIH/ODSS and STRIDES for supporting this project.

Questions? Email: gris@mail.nih.gov