**Breakout Session 5:**

**"T2T-omics" at scale: Improving our understanding of human genetic variation using AnVIL**

Professor Michael Schatz (Moderator)
*Bloomberg Distinguished Professor of Computer Science and Biology,*
*Johns Hopkins University*

# T2T Powered by AnVIL!



3202 samples from 26 populations

3202 samples x 30Gb = 96Tb input data

# T2T on AnVIL

# Core usage over 24 hours



Preview

1 hour    4 hours    **1 day**

instance/cpu/reserved_cores: 11,552.00

https://dockstore.org/workflows/github.com/schatzlab/t2t-variants/T2T_alignment

Samantha Zarate

# T2T Genomes Powered by AnVIL



https://anvil.terra.bio/#workspaces/anvil-datastorage/AnVIL_T2T

Sergey Aganezov  Stephanie Yan

Daniela Soto

Melanie Kirsche Samantha Zarate

**A complete reference genome improves analysis of human genetic variation**

# Science

$15
1 APRIL 2022
SPECIAL ISSUE
science.org

AAAS

# FILLING THE GAPS

Closing in on a complete human genome p. 42

---

# COMPLETING THE HUMAN GENOME

A fully sequenced human genome was announced more than 20 years ago. However, owing to technological limitations, some genomic regions remained unresolved. Here, *Science* and other journals present research by the Telomere-to-Telomere (T2T) Consortium, reporting on the endeavor to complete a comprehensive human reference genome.

FILTERS

6 RESULTS FOUND

**SPECIAL ISSUE RESEARCH ARTICLE**

## Segmental duplications and their variation in a complete human genome

BY MITCHELL R. VOLLGER, XAVI GUITART, PHILIP C. DISHUCK, LUDOVICA MERCURI, WILLIAM T. HARVEY, ARIEL GERSHMAN, MARK DIEKHANS, ARVIS SULOVARI, KATHERINE M. MUNSON, ALEXANDRA P. LEWIS, [...] EVAN E. EICHLER

SCIENCE • VOL. 376, NO. 6588 • 01 APR 2022

**SPECIAL ISSUE RESEARCH ARTICLE**

## Complete genomic and epigenetic maps of human centromeres

BY NICOLAS ALTEMOSE, GLENNIS A. LOGSDON, ANDREY V. BZIKADZE, PRAGYA SIDHWANI, SASHA A. LANGLEY, GINA V. CALDAS, SAVANNAH J. HOYT, LEV URALSKY, FEDOR D. RYABOV, COLIN J. SHEW, [...] KAREN H. MIGA  +48

SCIENCE • VOL. 376, NO. 6588 • 01 APR 2022

**SPECIAL ISSUE RESEARCH ARTICLE**

## From telomere to telomere: The transcriptional and epigenetic state of human repeat elements

BY SAVANNAH J. HOYT, JESSICA M. STORER, GABRIELLE A. HARTLEY, PATRICK G. S. GRADY, ARIEL GERSHMAN, LEONARDO G. DE LIMA, CHARLES LIMOUSE, REZA HALABIAN, LUKE WOJENSKI, MATIAS RODRIGUEZ, [...] RACHEL J. O

+16 authors • SCIENCE • VOL. 376, NO. 6588 • 01 APR 2022

**SPECIAL ISSUE RESEARCH ARTICLE**

## A complete reference genome improves analysis of human genetic variation

BY SERGEY AGANEZOV, STEPHANIE M. YAN, DANIELA C. SOTO, MELANIE KIRSCHE, SAMANTHA ZARATE, PAVEL AVDEYEV, DYLAN J. TAYLOR, KISHWAR SHAFIN, ALAINA SHUMATE, CHUNLIN XIAO, [...] MICHAEL C. SCHATZ  +22

SCIENCE • VOL. 376, NO. 6588 • 01 APR 2022

**SPECIAL ISSUE RESEARCH ARTICLE**

## Epigenetic patterns in a complete human genome

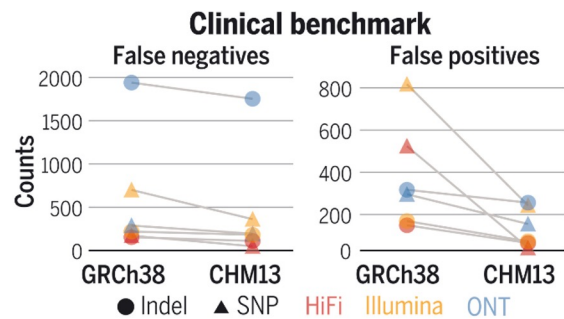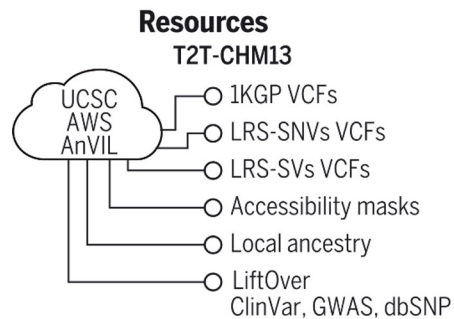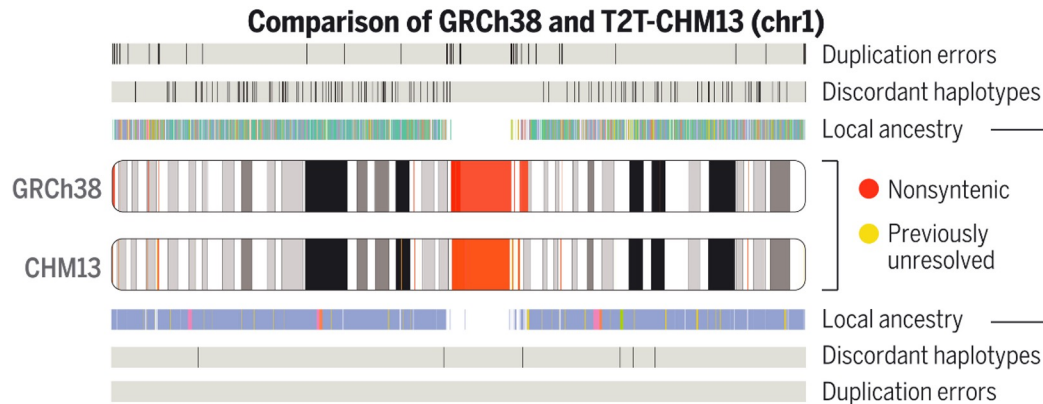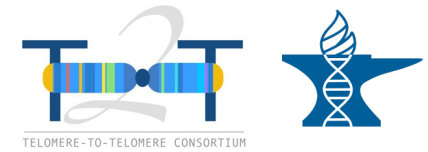BY ARIEL GERSHMAN, MICHAEL E. G. SAURIA, XAVI GUITART, MITCHELL R. VOLLGER, PAUL W. HOOK, SAVANNAH J. HOYT, MITEN JAIN, ALAINA SHUMATE, ROHAM RAZAGHI, SERGEY KOREN, [...] WINSTON TIMP  +9 authors

VOL. 376, NO. 6588 • 01 APR 2022

**SPECIAL ISSUE RESEARCH ARTICLE**

## The complete sequence of a human genome

BY SERGEY NURK, SERGEY KOREN, ARANG RHIE, MIKKO RAUTIAINEN, ANDREY V. BZIKADZE, ALLA MIKHEENKO, MITCHELL R. VOLLGER, NICOLAS ALTEMOSE, LEV URALSKY, ARIEL GERSHMAN, [...] ADAM M. PHILLIPPY  +89 au

SCIENCE • VOL. 376, NO. 6588 • 31 MAR 2022 : 44-53

# T2T-chrY: Human variation across 156 populations

**1000 Genomes Project (1KGP)**
3,202 samples from 26 populations



(Byrska-Bishop et al., Cell, 2022)

**Simons Genome Diversity Project (SGDP)**
279 open access samples from 130 populations



(Mallick et al., Nature, 2016)

**The complete sequence of a human Y chromosome**
Rhie *et al.* (2023) *Nature. https://doi.org/10.1038/s41586-023-06457-y*

Stephen Hwang          Dylan Taylor

# Hidden Variants in Breast Cancer Genes



Thanks to long reads we can now robustly detect entirely new types of variation

…

But how can we identify those variants with clinical & functional impact?

# CoLoRS: Consortium of Long Read Sequencing

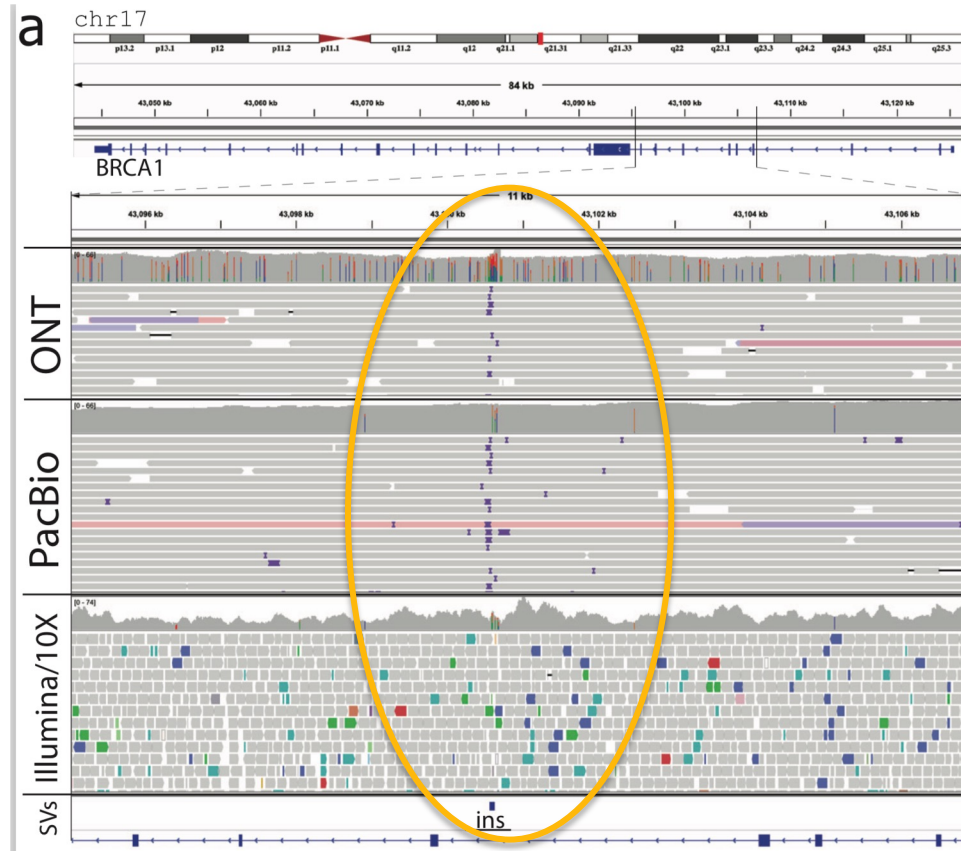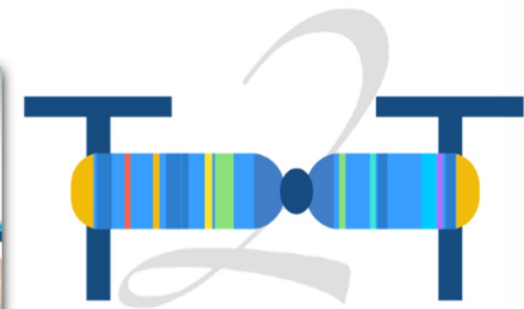| Organization | Number of samples | Samples | Coverage | Source | Orthogonal Data |
|---|---|---|---|---|---|
| Children's Mercy Research Institue | 1,071 | trios, 85% European | 571 are parents at 8-10x depth, 500 are individuals (probands, affected)20-30x depth | blood | WES (minimally),srWGS, many probands with some RNAseq/Iso-Seq |
| Human Genome Structural Variation Consortium (HGSVC) | 37 (goal 70) currently @ EBI | 1k (each population), healthy | >30-40x HiFi | cell lines | Comprehensive |
| Human PanGenome Reference Consortium (HPRC) | 127 (goal 350) | first 130 from 1000G, after that other populations, healthy | >30-40x HiFi | cell lines, future mix primary/cell lines | Illumina, Nanopore |
| University of Tokyo - Morishita Lab | 300 | HiFi genomes, all Japanese, healthy | 8x-20x HiFi | cell lines | Illumina, some Nanopore |
| HudsonAlpha Institute for Biotechnology (HAIB) | 80 | 50 probands (all affected), 30 parents, 60% European, 25% African American | 20x HiFi | blood | Illumina for nearly all |
| SolveRD | 100 (goal 510, 2022) | majority European, 100 trios, others singletons affected | 8-10x HiFi | largely blood | Illumina WES, occasionally genomes or array |
| Radboud UMC- Hoischen Lab | 5 trios CLR, 8 HiFi trios | probands with severe disease | 15-40x PacBio CLR, 30x HiFi | blood | Illumina WES, WGS, array, some bionano |
| University of Washington - Eichler Lab | Autism cohort (42, 12 families), quads & trios, goal 3x | families of autism with unsolved cases | >30x HiFi | largely blood, some cell lines | Illumina WES, arrays, half ONT |
| Amsterdam UMC - Holstege Lab | >100, goal 600 | Dutch population | 25x Hifi & PacBio CLR | Blood | WES & array data on all, |
| Kyushu University (Nagasaki lab) and National Center for Global Health and Medicine | 80 (goal 100) | HiFi genomes, all Japanese, healthy | 5 - 40x HiFi | Cell lines | Illumina |
| Chulalongkorn University | 250 (goal 300) | Patients with rare diseases and their parents. Thai ethnic. | 10 - 40x HiFi | Blood | Illumina, Nanopore |

Table reflects samples on 7/18/23, We expect these samples to grow with expansion of the projects above and by the addition of new collaborators.

**Open coalition of international researchers focused on cataloging all classes of variation using long-read whole genome sequencing.**

- The goal is to provide variant frequency data for public use and as a resource to the global scientific and clinical research community

- Complements existing databases such as gnomAD

- Develop state-of-the-art pipelines, execute at individual sites or within the AnVIL cloud platform

**>2195 samples and growing!**

**https://colorsdb.org/**

# Acknowledgements

## Schatz Lab

Enis Afghan

Ahmed Awan

Dannon Baker

Tyler Collins

Arun Das

John Davis

Sam Guerler

Katie Jenike

Sam Kovaka

Natalie Kucher

Qiuhui Li

Stephen Mosher

Matthew Nguyen

Bohan Ni

Alex Ostrovsky

Srividya
   Ramakrishnan

Michelle Savage

Michael Sauria

Margaret Starostik

Alex Sweeten

Jenn Vessio

Natalie Whitaker

### T2T, AnVIL, & Galaxy Teams

Miga, Phillippy, Eichler, Nekrutenko, Goecks, Tan, Leek, Morgan, Carey, Philippakis *et al.*

### CoLoRS

Eichler, Lake, Wenger, Korlach, Beck, Pastinen, Audano, Garimella, Schmutz, Chen *et al.*

## JHU

Battle Lab

Klein Lab

Genetic Resources Core

## Timp Lab

Carolina Montano

Jessica Hosea

Luke Morina

## Stanford

Montgomery Lab

Ashley Lab

## Mayo Clinic

Gloria Petersen / Sam Antwi

## University of Toronto

Steven Gallinger



CoLoRS

National Human Genome Research Institute

NIH National Cancer Institute

NSF

Bloomberg Professors

**National Institutes of Health**
*Office of Data Science Strategy*

HVD 21: Telomere-to-Telomere Consortium Analyses on the NHGRI AnVIL

HVD 22: Long Read Variant Frequency Database on AnVIL (CoLoRS)

# Thank you!
## schatz-lab.org