# The Human Pangenome Project:
Creating a Reference that Better Represents Human Global Genetic Diversity

TELOMERE-TO-TELOMERE CONSORTIUM

TOWARDS A COMPLETE REFERENCE OF HUMAN GENOME DIVERSITY
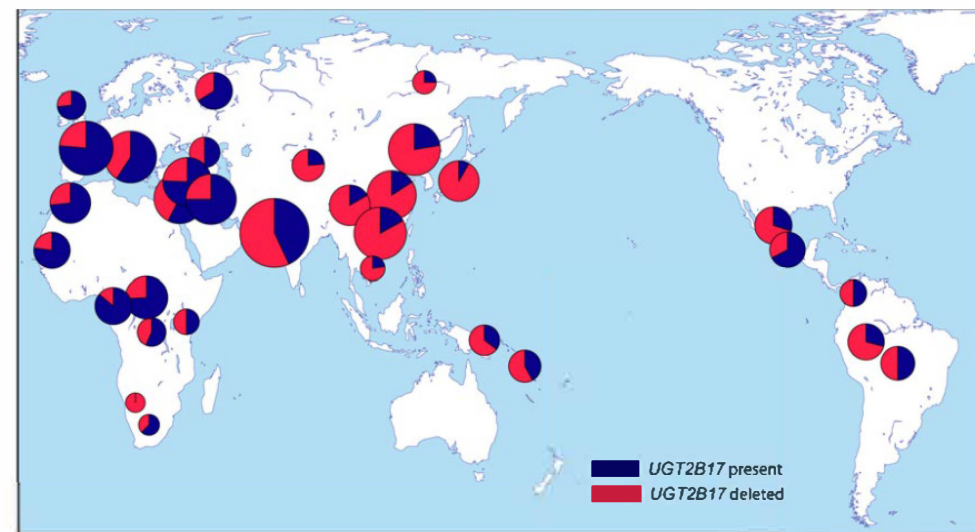
HUMAN PANGENOME

**Deanna Church**
**Presenting**

# The current human reference fails to faithfully represent a single individual genome much less the genomes of a global population.

RP11 hap1

*UGT2B17*

RP11 hap2



- **The human reference genome** is a foundational resource in human genetics and like most technology-driven resources, is overdue for an upgrade.

- The current structure is a **linear monoploid representation containing mixed haplotypes with too many gaps and errors. Additionally, the underlying sequence is predominately from a single individual.**

- **Mapping limitations of short reads and inherent reference biases** means we have missed more than 70% of structural variants in traditional whole-genome sequencing studies
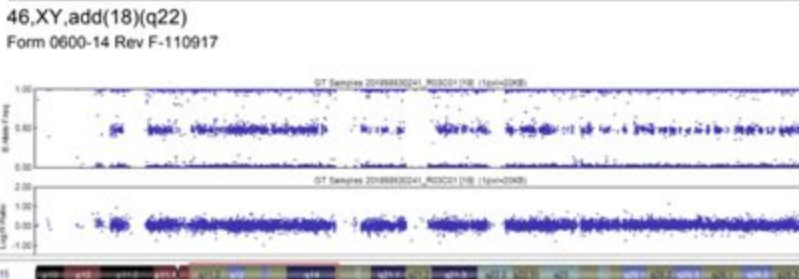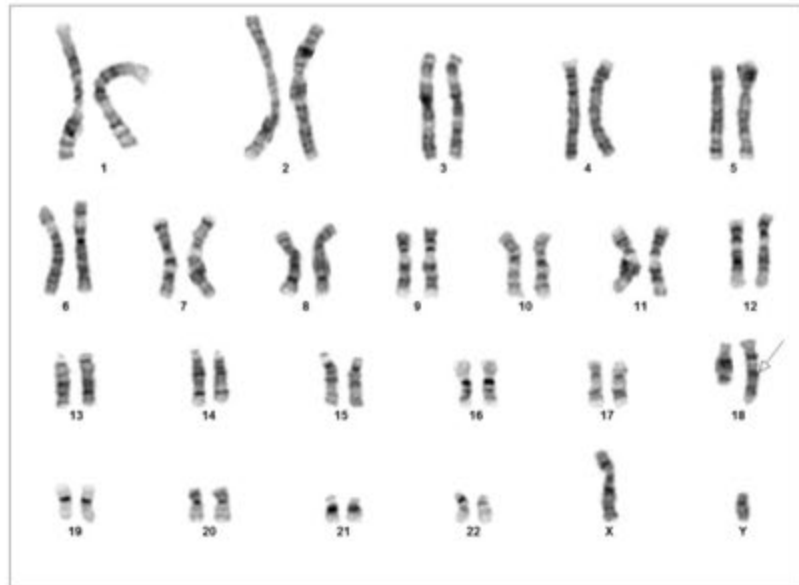
# Human Pangenome Reference Consortium

- Improve representation of **global genomic diversity** (>350 diverse diploid references)

- **Prioritizing quality**: we aim to release a complete (T2T) and comprehensive map of genome variation

- **Develop a new, non-linear reference data structure** and foster an innovative ecosystem of pangenomic tools

- Outreach, Education and Implementation

# Multi-Center Sequencing Technology and Production: Optimized for Efficiency and Quality
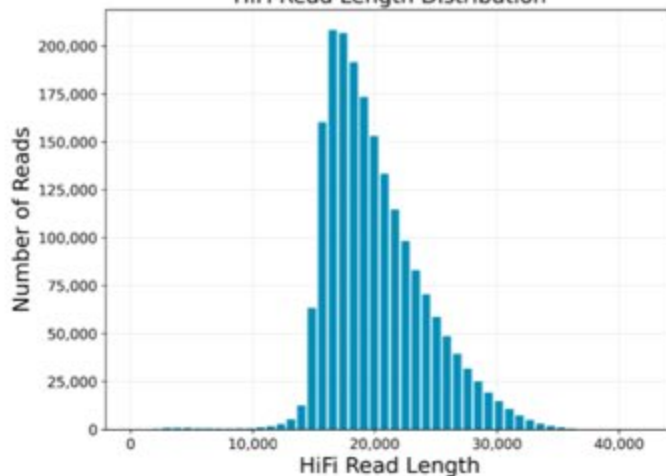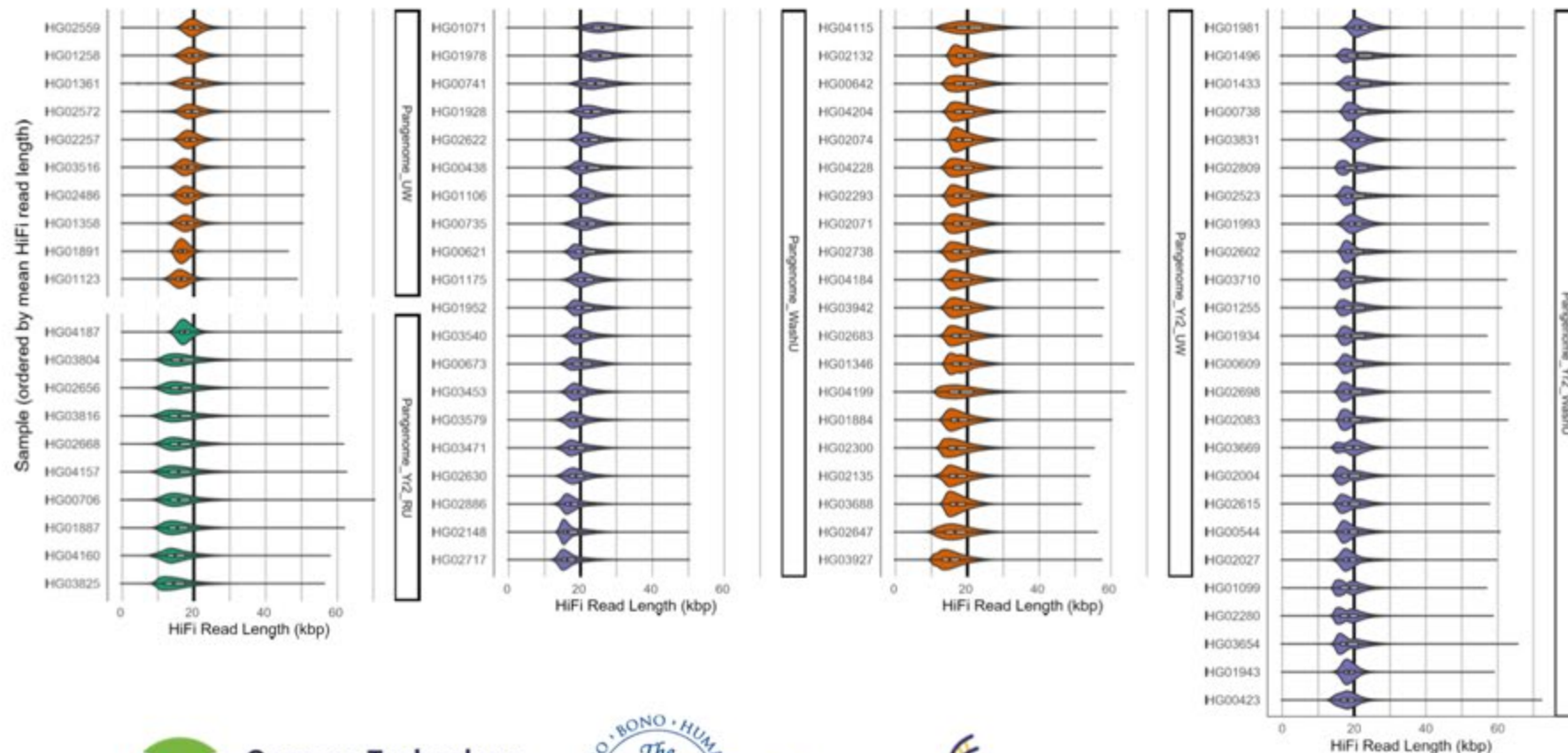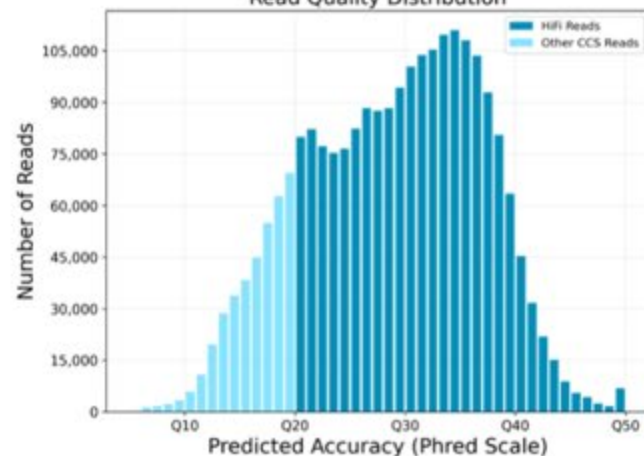
## Cell line stability/QC

Optimized and Consistent Long-Read HiFi Production

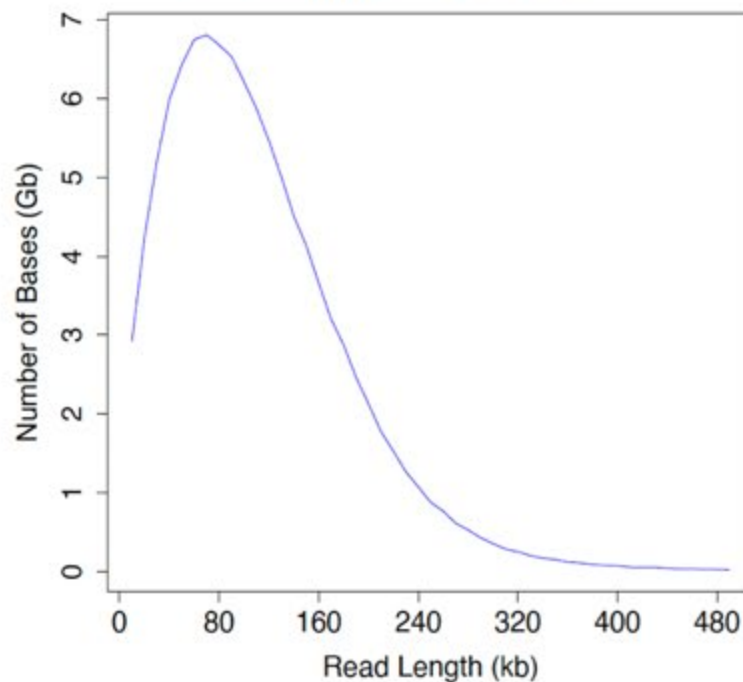# Forefront of ultra-long read sequencing innovation

**2020**

Read N50 ~76 kb
~9X coverage per flow cell
~3.5X coverage in 100 kb+

**2021**

Read N50 ~73 kb
~30X coverage per flow cell
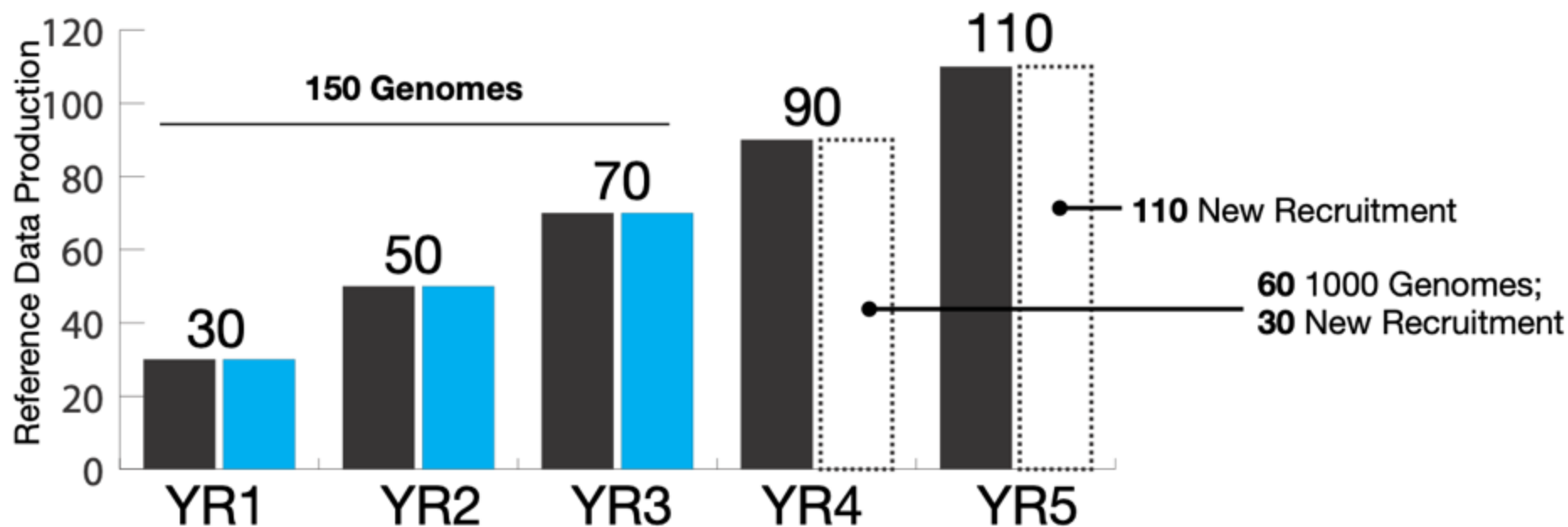~9X coverage in 100 kb+
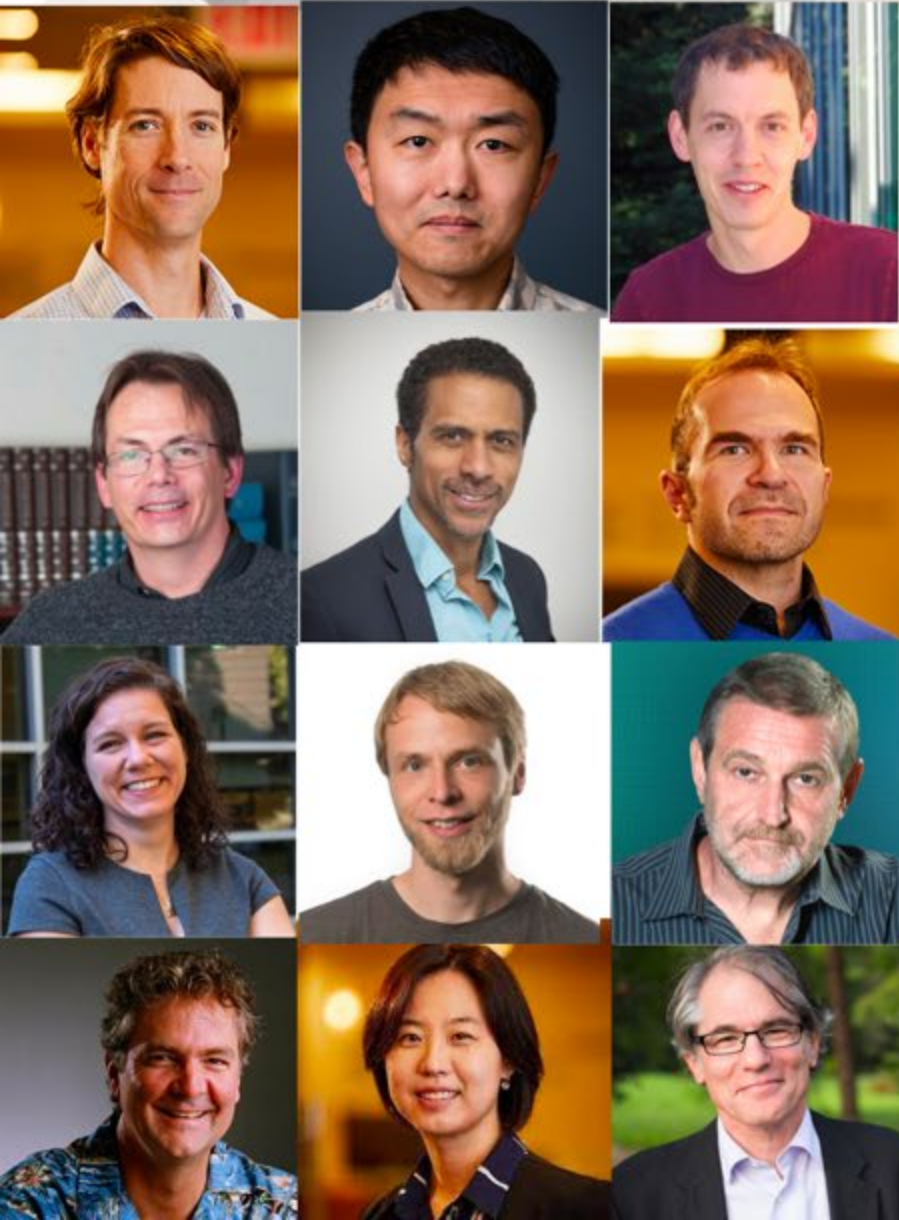
**3x increase in throughput**

# Overview

**Population Sampling and Representation (Phase 1)**: We are representing >99.9% common SNVs (1%) in the 1000 Genomes lines. **(Phase 2)** We are in position for perspective recruitment of remaining 150 individuals (BioMe)

**Sequence Technology and Production:** Highly efficient multi-center production effort, automated assembly and quality assessment

# Innovation in Long Read Assembly Methods



We have assembled the leaders in long-read assembly methods, with an emphasis in researchers involved in finishing and repeat assemblies.

# Science

$15
1 APRIL 2022
SPECIAL ISSUE
science.org

AAAS

## FILLING THE GAPS

Closing in on a complete human genome p. 42

---

# Automated assembly standards: High-quality References

We tested the current best practices in sequencing technologies and automated assembly algorithms on one human sample, HG002, an openly- consented Ashkenazi individual from the Personal Genome Project



**T2T-CHM13 3.055 Gb**

**QV 50-60**

Near complete chromosome scaffolds between HPRC-HG002 maternal and paternal assemblies.

# One genome is not enough....

# Population Sampling and Representation

**Goals:**

- Establish a framework to define diversity and prioritize cell lines from 350 globally diverse participants to build a new pangenome reference

- Aim to capture most common variants, defined as variants at >1% frequency, in human populations globally

**Challenges:**

- Incomplete understanding of the full spectrum of human genetic diversity

- Not one way to define diversity in this study: multiple well motivated approaches proposed for sample selection

# Phase 1 HPRC Sample Selection



✓ Cover genetic and geographic diversity

✓ Availability of low passage cell lines

✓ Availability of trios/parental data (YR1-2)

**Phase 1 Production: 150 total lines**

- Admixed American (AMR) — 37, 25%
- East Asian (EAS) — 30, 20%
- European (EUR) — 8, 5%
- African (AFR) — 45, 30%
- South Asian (SAS) — 30, 20%

Eimear Kenny, PI

**Mount Sinai's BioMe BioBank:**
>70,000 participants from 160 countries, with sequencing data for more than 40,000 participants.

**Washington University, St. Louis Recruitment Center:**
New recruitment from African American communities

# Phase 2: Prioritization of New Participants

- **Model 1:** Maximizing common variant diversity
Leveraging sequence data, iteratively select samples that maximize common variant coverage in an out-of-sample dataset

- **Model 2**: Maximizing genetic divergence
*Using sequence/array data, plot PCA and select participants to represent the continuum of genetic dissimilarity in principle component (PC) space*

- **Model 3**: Targeting underrepresented populations based on self-reporting/ geographical data
Using country-of-origin data, designate countries/subcontinental regions of interest to the project, and select participants from those regions

# Partnership with ELSI Team for Consent and Outreach



**Embedded Ethics Team**

- **Formal Review of Consent Language (compatibility with 1000 genomes consent)**

- **Design and review of outreach materials for prospective recruitment**

- **Alignment with internal ethics/genomic review board at BioMe**

- **Review of the use of external IRB ( BRANY )**

# Welcome to the Human Pangenome Project!

About Us ➔

The Human Pangenome Project is an International Alliance of genomics partners that aims to provide accessible high-quality genomes that represent the diversity of the human population. These genomes will be represented as a Pangenome Reference.

# Improving Population Sampling and Representation

**Continue reference production of 1000 genomes (200/350)**

**Move away from dependency on trios** (priority based on genomic diversity) Establish new assembly methods that do  not rely on parental data

**Launch new recruitment efforts and establish 100 new LCLs (NHGRI Collection: Human Pangenome)** outside of 1000 Genome Cell lines

**Establish new collaborations and international partnerships**

# The Human Pangenome



Alignments of high-quality assemblies were performed using three different methods, pioneered by our team:

Minigraph (Li et al., 2020),
Minigraph-Cactus (MC)
PanGenome Graph Builder (PGGB).

# HLA-A: Aligning Complex Pangenome Loci

# Applications



A pangenome approach captures annotation missing in a linear reference

A pangenome approach allows annotation at heterozygous structural variant sites

# Remaining challenges

- Sample selection and resource development

  - Ensuring this project is truly international in scope

  - Ensuring that we common (> 1% MAF) variation is included

  - Development of usable cell lines (iPSCs, LCLs, etc) for experimental work

- Pangenome implementation

  - Will one graph representation rule them all, or do we need different graphs for different applications?

- Pangenome adoption

  - In a world where many people still use GRCh37, how do we encourage adoption of this resource?

# Acknowledgements

**UNIVERSITY OF CALIFORNIA SANTA CRUZ** Genomics Institute

David Haussler, Miten Jain, Benedict Paten, Hugh Olsen, Karen Miga, Erik Garrison, Ed Green, Marina Haukness, Mark Akeson, Mark Akeson, Adam Novack

Adirna Fuller, Tony Tsung Yu Lu, Xian Chang, Trevor Pesout, Ryan Lorig-Roach, Charles Markello, Jean Monlong, Glenn Hickey, Jonas Sibbesen

Melissa Meredith, Kishwar Shafin, Jouni Siren, Jordan Eizenga, Beth Sheets, Julian Lucas, Brian Hannafious

**THE GENOME INSTITUTE** at Washington University

Ting Wang, Lucinda Fulton, Sarah Cody, Robert Fulton, Heather Lawson

Wen-Wei Liao, Nathan Stitziel, Milinn Kremitizki, Haley Abel, Eddie Belter

Derek Albracht, Chad Tomlinson, Allison Regier, Chris Markovic, Tina Lindsay

**EMBL-EBI**

Paul Flicek, Susan Fairley, Daniel Zerbino

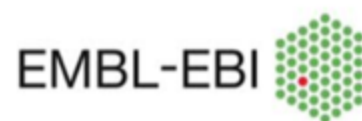**Icahn School of Medicine at Mount Sinai**

Eimear Kenny, Vimi Desai

**NIH** National Human Genome Research Institute

Adam Felsenfeld, Mike Smith, Carolyn Hutter, Taylorlyn Stephan, Heidi Sofia, Nicole Lockhart, Sara Currin

Adam Phillippy, Sergey Koren, Arang Rhie, Chirag Jain, Baergen Schultz, Kris Wetterstrand

**Genome Sciences UNIVERSITY OF WASHINGTON**

Evan Eichler, Katy Munson, Mitchell Vollger, Pete Audano, David Porubksy, Arvis Sulovari, Gene Myers

**CBG** Max Planck Institute of Molecular Cell Biology and Genetics

**Yale University**

Ira Hall, Wen-Wei Liao, Shuangjia Lu

**wellcome sanger institute**

Kerstin Howe

**UNIVERSITY OF CAMBRIDGE**

Richard Durbin

**HARVARD MEDICAL SCHOOL**

Heng Li, Shilpa Garg, Haoyu Cheng, Xiaowen Feng

**National Center for Biotechnology Information NCBI**

Valerie Schneider, Terence Murphy, Paul Kitts, Chunlin Xiao, Francoise Thibaud-Nissent

**AnVIL**

**CORIELL INSTITUTE FOR MEDICAL RESEARCH** DECODING THE GENOME

Alissa Resch, Brittany Kerr, Brittney Martinez, Ellen Kellly

**THE ROCKEFELLER UNIVERSITY** Science for the benefit of humanity

Erich Jarvis, Olivier Fedrigo, Giulio Formenti, Sadye Paez, Lauren Shalmiyev

Barbara Koenig (UCSF), Nanibaa' Garrison (UCLA), Bob Cook-Deegan (ASU), Alice Popejoy (UC Davis)

**TOWARDS A COMPLETE REFERENCE OF HUMAN GENOME DIVERSITY**

**HUMAN PANGENOME**

**NIST**

Justin Zook

**Company Partnerships**

**aws**

illumina, Google DeepVariant, ARIMA GENOMICS, Dovetail GENOMICS, PACBIO, NANOPORE Technologies, HudsonAlpha INSTITUTE FOR BIOTECHNOLOGY