# Integrating **structural variants** in genomic studies of rare and complex diseases with **long-read sequencing** and **pangenomes**

Jean Monlong

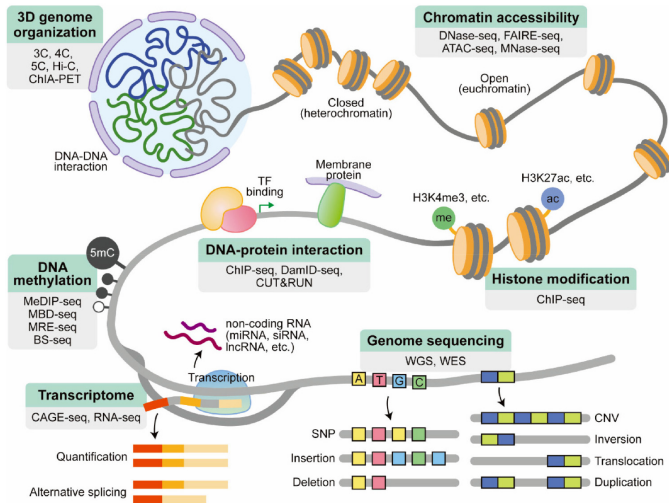JOBIM

11/07/2025

irsd
Institut de recherche en santé digestive

Inserm
La science pour la santé
From science to health

# Understanding functional impact of genomic variation



*Dev. Reprod. 2023; 27(1):9-24 DOI: 10.12717/DR.2023.27.1.9*

# Different types of genomic variants

Single-nucleotide polymorphisms
(**SNPs**)

Insertion-deletion polymorphisms
(**INDELs**)

Structural variants
(**SVs**)

GAT**C**AGC

GAT**CA**GC

GATCAGC

GAT**G**AGC

GAT - - GC

GATC AGC

CGC....300bp....GAT

# Different types of genomic variants

Single-nucleotide polymorphisms
(**SNPs**)

Insertion-deletion polymorphisms
(**INDELs**)

Structural variants
(**SVs**)

GAT**C**AGC

GAT**CA**GC

GATCAGC

GAT**G**AGC

GAT **- -** GC

GATC AGC

CGC....300bp....GAT

Single-nucleotide polymorphisms (**SNPs**)

Insertion-deletion polymorphisms (**INDELs**)

Structural variants (**SVs**)
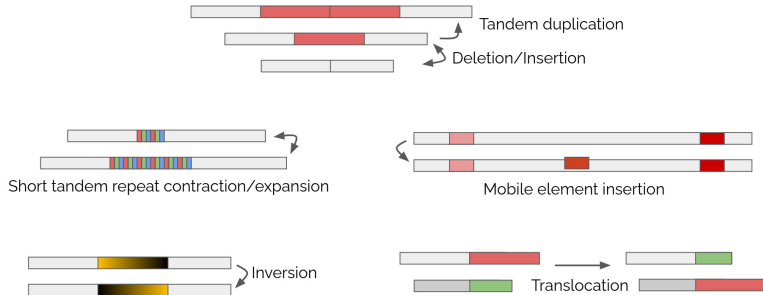
GAT**C**AGC
GAT**G**AGC

GAT**CA**GC
GAT **- -** GC

GATCAGC
GATC**∧**AGC

CGC....300bp....GAT

# Structural variants (SVs) come in diverse shapes and sizes
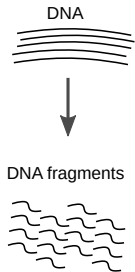
Variant size: from 50 bases to megabases.



Tandem duplication

Deletion/Insertion

Short tandem repeat contraction/expansion

Mobile element insertion

Inversion

Translocation

# Structural variants (SVs) come in diverse shapes and sizes

Variant size: from 50 bases to megabases.



Tandem duplication

Deletion/Insertion

Short tandem repeat contraction/expansion

Mobile element insertion

Inversion

Translocation

- ◆ High functional impact
- ◆ Involved in rare and common diseases, and cancers.
- ◆ **Hard to detect**

# Genome sequencing



DNA

DNA fragments
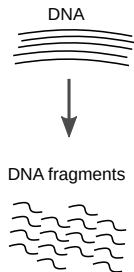
Sequencing machines
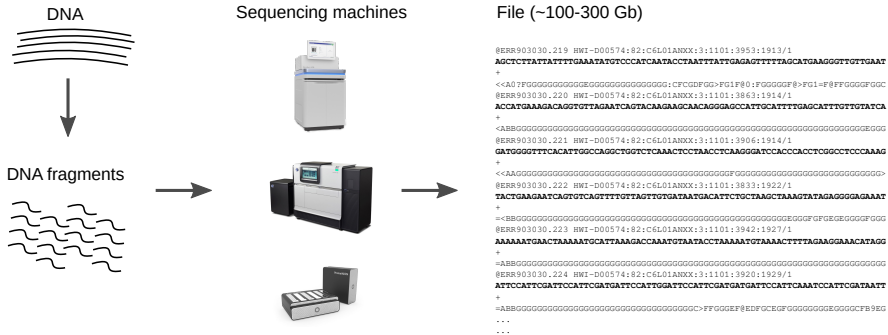
File (~100-300 Gb)

# Genome sequencing



DNA

Sequencing machines

File (~100-300 Gb)

DNA fragments

# Genome sequencing



DNA

Sequencing machines

File (~100-300 Gb)

DNA fragments

```
@ERR903030.219 HWI-D00574:82:C6L01ANXX:3:1101:3953:1913/1
AGCTCTTATTATTTTGAAATATGTCCCATCAATACCTAATTTATTGAGAGTTTTTAGCATGAAGGGTTGTTGAAT
+
<<A0?FGGGGGGGGGGGEGGGGGGGGGGGGGGGG:CFCGDFGG>FG1F@0:FGGGGGF@>FG1=F@FFGGGGFGGC
@ERR903030.220 HWI-D00574:82:C6L01ANXX:3:1101:3863:1914/1
ACCATGAAAGACAGGTGTTAGAATCAGTACAAGAAGCAACAGGGAGCCATTGCATTTTGAGCATTTGTTGTATCA
+
<ABBGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGEGGG
@ERR903030.221 HWI-D00574:82:C6L01ANXX:3:1101:3906:1914/1
GATGGGGTTTCACATTGGCCAGGCTGGTCTCAAACTCCTAACCTCAAGGGATCCACCCACCTCGGCCTCCCAAAG
+
<<AAGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGFGGGGGGGGGGGGGGGGGGGGGGGGGGGGG>
@ERR903030.222 HWI-D00574:82:C6L01ANXX:3:1101:3833:1922/1
TACTGAAGAATCAGTGTCAGTTTTGTTAGTTGTGATAATGACATTCTGCTAAGCTAAAGTATAGAGGGGAGAAAT
+
=<BBGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGEGGGFGFGEGEGGGGFGGG
@ERR903030.223 HWI-D00574:82:C6L01ANXX:3:1101:3942:1927/1
AAAAAATGAACTAAAAATGCATTAAAGACCAAATGTAATACCTAAAAATGTAAAACTTTTAGAAGGAAACATAGG
+
=ABBGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG
@ERR903030.224 HWI-D00574:82:C6L01ANXX:3:1101:3920:1929/1
ATTCCATTCGATTCCATTCGATGATTCCATTGGATTCCATTCGATGATGATTCCATTCAAATCCATTCGATAATT
+
=ABBGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGC>FFGGGEF@EDFGCEGFGGGGGGGGGEGGGGCFB9EG
...
...
```

## Sequencing reads

- **Short:** 150-250 bp (current tech)
- **Long:** 10,000s-100,000s bp (new tech. $$$)

Short-read sequencing, pangenomes, and complex diseases

Long-read sequencing and rare diseases

Pangenomes meet long-read sequencing

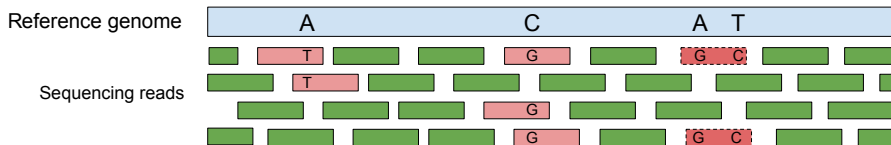# Short-read sequencing, pangenomes, and complex diseases

# Common variants associated with a complex disease


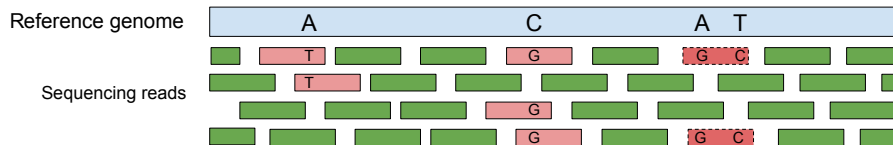
◆ Structural variant

## Goal
Genotype a comprehensive catalog of common variants across a large cohort.

# Aligning reads to a reference genome



**Assuming the reads are correctly placed**, small variants are identified as recurrent differences between reads and the reference genome.

**Assuming the reads are correctly placed**, small variants are identified as recurrent differences between reads and the reference genome.

Variants can be missed, resulting in **reference bias**.

# The challenges of structural variant detection

Around breakpoints, short sequencing reads are hard to map on the reference genome.

# The challenges of structural variant detection

Around breakpoints, short sequencing reads are hard to map on the reference genome.

# Pangenomics to the rescue. Which pangenomics?

Not "Genome-wide association studies" in French.

## Étude d'association pangénomique

文A **22 langues** ∨

Article  Discussion

Lire  Modifier  Modifier le code  Voir l'historique  Outils ∨

Une **étude d'association pangénomique** (en anglais *genome-wide association study*, GWAS) est une analyse de nombreuses variations génétiques chez de nombreux individus, afin d'étudier leurs corrélations avec des traits phénotypiques[1].

Ces études se concentrent généralement sur les associations entre les polymorphismes nucléotidiques (SNP) et des phénotypes tels que les maladies humaines majeures.

# Pangenomics to the rescue. Which pangenomics?

Not "Genome-wide association studies" in French.



Also not exactly the set of **genes** from all strains within a clade, like in microbial pangenome.

A pangenome represents a **collection of genomes** and the genetic variants among them.

A pangenome represents a **collection of genomes** and the genetic variants among them.

# Building a Human pangenome reference



- Human Pangenome Reference Consortium (**HPRC**)
- Latest sequencing technologies for 350 diverse individuals
- Pangenome containing a comprehensive catalog of (structural) variants
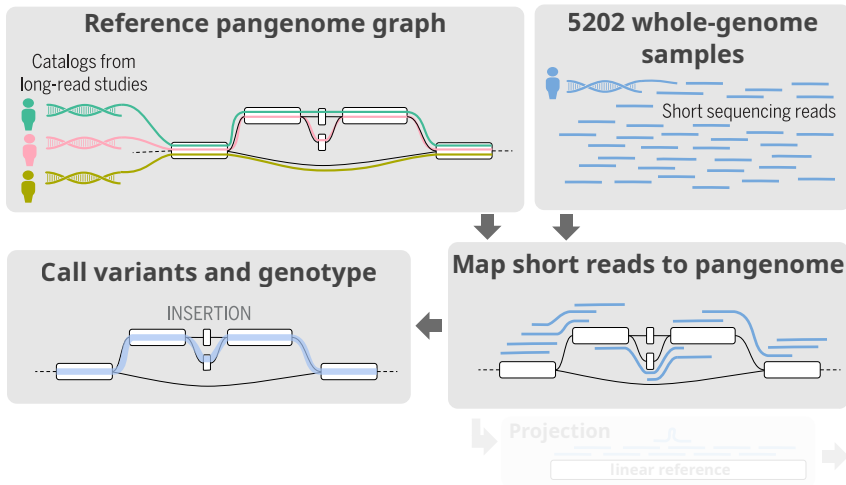
Liao*, Asri*, Ebler*, et al. Nature 2023
Hickey*, Monlong*, et al. Nat. Biotechnol. 2023
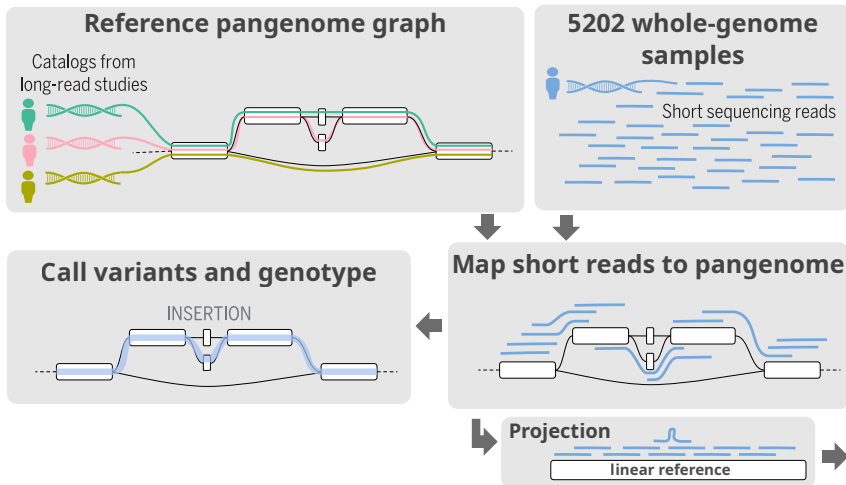
# Building a Human pangenome reference, a team effort



Check out the latest data at: `https://data.humanpangenome.org`

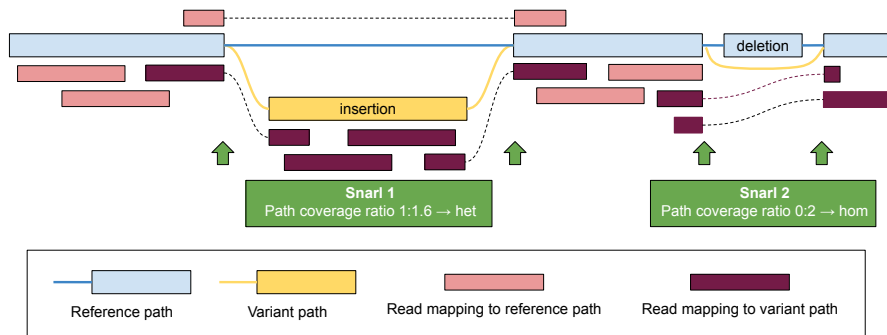# Short-read mapping and structural variant genotyping



Siren*, Monlong*, Chang*, Novak*, Eizenga*, et al. Science 2021

# Short-read mapping and structural variant genotyping



Siren*, Monlong*, Chang*, Novak*, Eizenga*, et al. Science 2021
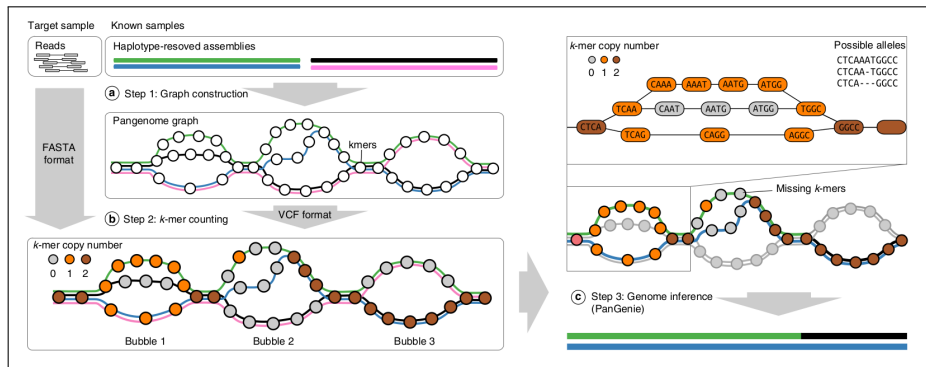
# Genotyping structural variation from pangenomic mapping



https://github.com/vgteam/vg

Hickey*, Heller*, Monlong*, et al. Genome Biology 2020

PanGenie uses k-mer and haplotype information to genotype SVs.



`https://github.com/eblerjana/PanGenie`

Ebler et al. Nature Genetics 2022

# Personalized pangenomes with haplotype sampling

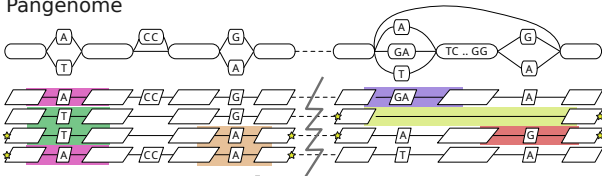With pangenomes becoming larger, analysis can suffer.

With pangenomes becoming larger, analysis can suffer.

One solution: k-mer-guided "down-sampling" of the full pangenome.



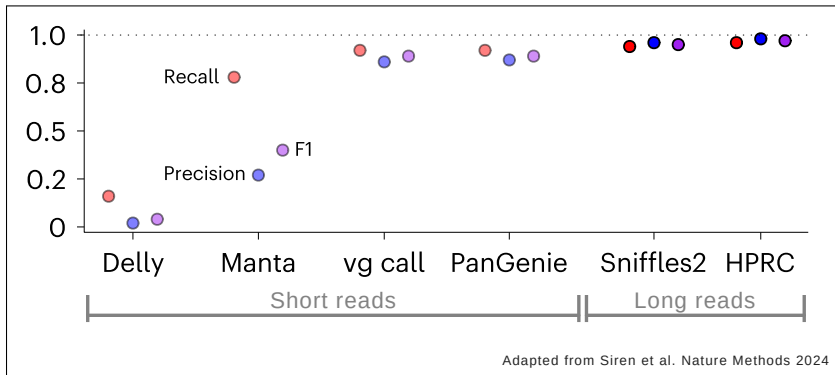K-mer counts from sequencing experiment

Pangenome

Personalized pangenome with only N haplotypes
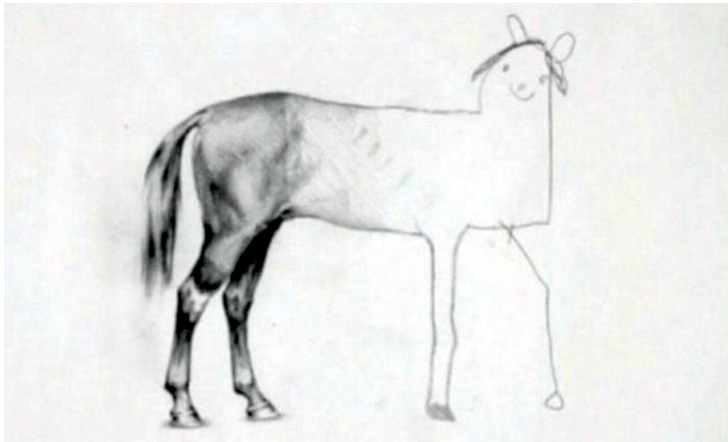
Sirén et al. Nature Methods 2024

# Structural variant genotyping performance



Adapted from Siren et al. Nature Methods 2024

*vg call and PanGenie using the "personalized pangenome" approach.
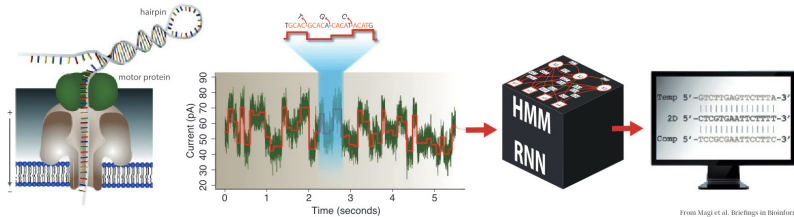
Construction    Complex variants    Annotation
Read mapping    Functional genomics    Association tests
Genotyping    Visualization    Multi-species
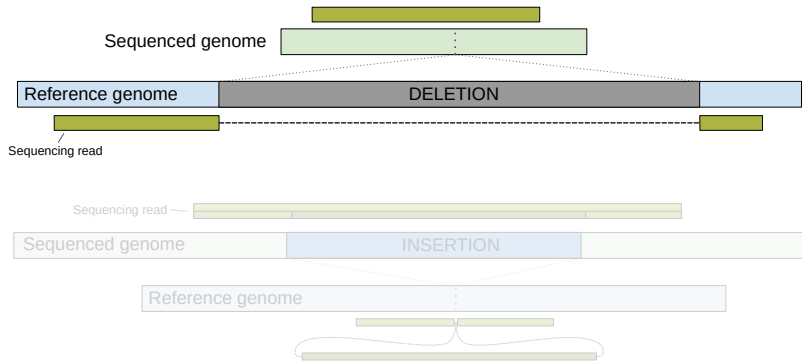
# Long-read sequencing and rare diseases
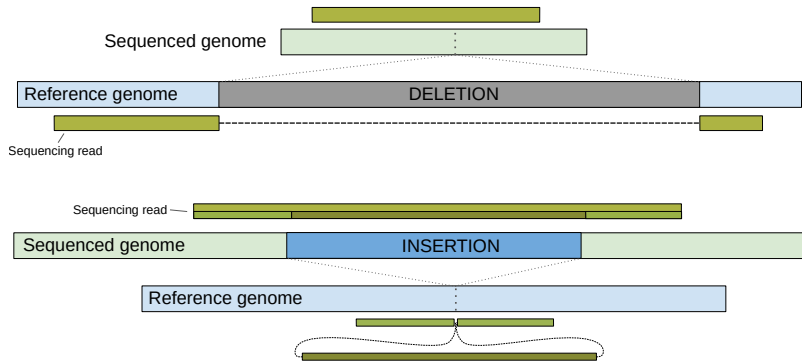
From Magi et al. Briefings in Bioinformatics

As the DNA (or RNA) fragment passes through the pore, the current changes and is decoded to predict nucleotides.

Reads length of 1,000s-100,000s of nucleotides.

Sequenced genome

Reference genome

DELETION

Sequencing read

Sequencing read

Sequenced genome

INSERTION

Reference genome

# Longer reads improve structural variant detection

# Oxford Nanopore is portable (space!) and fast

- Sequence as fast as possible
- Get a genomic diagnosis quick
- E.g. for newborns with suspicion of a rare genetic disease



Gorzynski et al. N. Engl. J. Med. 2022

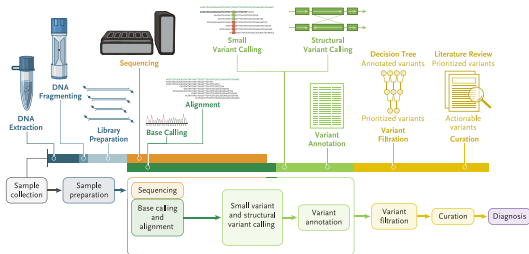Goenka*, Gorzynski*, Shafin*, et al. Nat. Biotechnol. 2022

- Sequence as fast as possible
- Get a genomic diagnosis quick
- E.g. for newborns with suspicion of a rare genetic disease



Ultrarapid Genome Sequencing Pipeline

"Fastest DNA sequencing technique": 5h2m
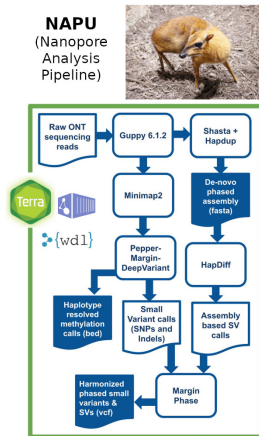


Gorzynski et al. N. Engl. J. Med. 2022
Goenka*, Gorzynski*, Shafin*, et al. Nat. Biotechnol. 2022

# Cost-efficient Nanopore pipeline

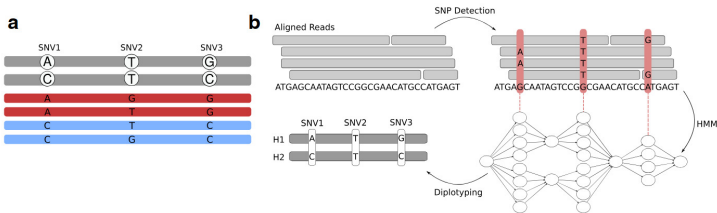- Only **one flow-cell** of Nanopore
- ∼30X coverage with 30 Kbp N50 reads

# Cost-efficient Nanopore pipeline

- Only **one flow-cell** of Nanopore
- ~30X coverage with 30 Kbp N50 reads
- Nanopore Analysis Pipeline (U?) to get haplotype resolved:
  1. small variants (SNPs/indels)
  2. structural variants
  3. *de novo* assembly
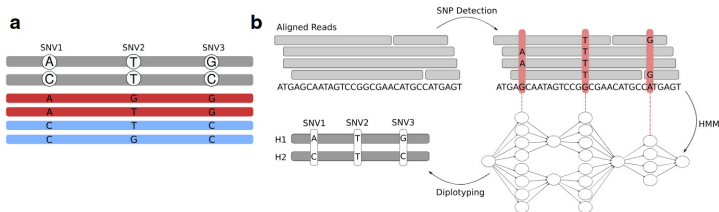  4. methylation marks



**NAPU**
(Nanopore
Analysis
Pipeline)

Kolmogorov*, Billingsley*, et al. Nature Methods 2023

Reads are **haplo-tagged** using information across heterozygous sites with Margin (Ebler*, Haukness*, Pesout*, et al. Genome Biology 2019).
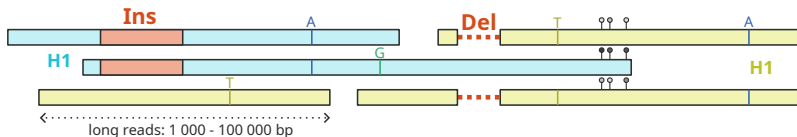
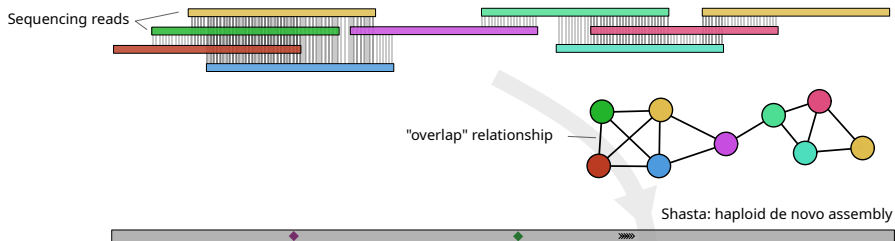# Under the hood: phased variants and methylation calls



Reads are **haplo-tagged** using information across heterozygous sites with Margin (Ebler*, Haukness*, Pesout*, et al. Genome Biology 2019).



Phased small variants (DeepVariant) and methylation calls (ModKit)

Reconstructs genomes without reference bias, hence better able to identify complex variants (e.g. combination of deletion/inversion)



Sequencing reads

"overlap" relationship

Shasta: haploid de novo assembly

Shafin*, Pesout*, Lorig-Roach*, Haukness*, Olsen*, et al. Nat. Biotechnol. 2020          Kolmogorov*, Billingsley*, et al. Nature Methods 2023

Reconstructs genomes without reference bias, hence better able to identify complex variants (e.g. combination of deletion/inversion)
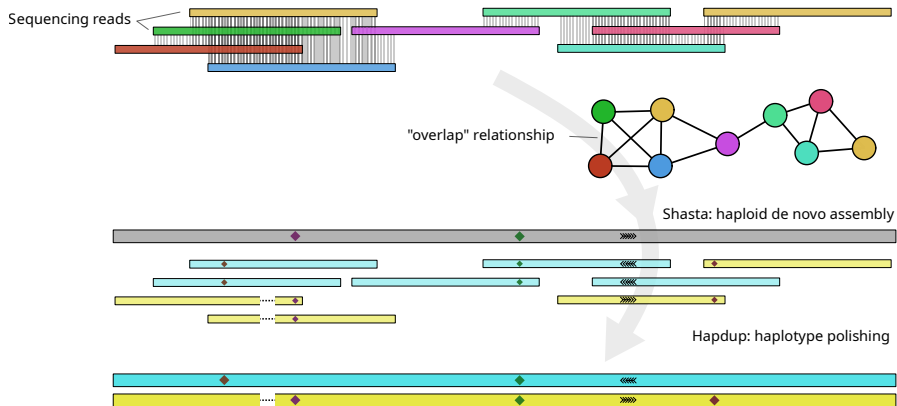


Shafin*, Pesout*, Lorig-Roach*, Haukness*, Olsen*, et al. Nat. Biotechnol. 2020    Kolmogorov*, Billingsley*, et al. Nature Methods 2023
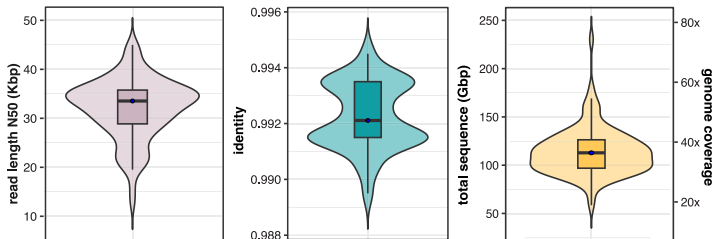
**Chan
Zuckerberg
Initiative**

Children's National.

GREGoR
consortium

42 probands and 56 unaffected family members, sequenced with one-flowcell of ONT long-read sequencing (R10).
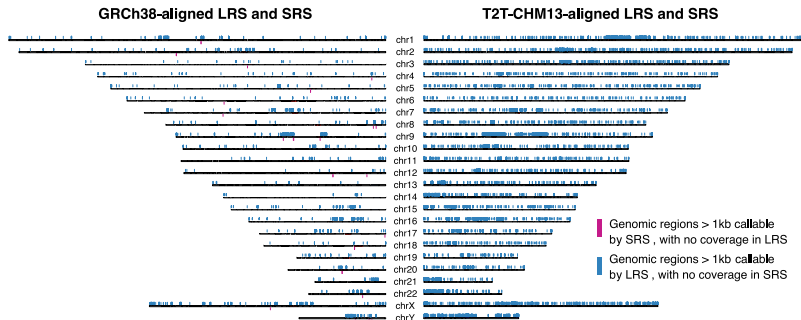


Negi et al. AJHG 2025

# Better coverage of confidently mapped reads

More of the CHM13-T2T genome covered with at least 10x.

◆ **93.99%** (LRS) vs. 88.27% (SRS)



**GRCh38-aligned LRS and SRS**          **T2T-CHM13-aligned LRS and SRS**

Genomic regions > 1kb callable by SRS , with no coverage in LRS

Genomic regions > 1kb callable by LRS , with no coverage in SRS

Negi et al. AJHG 2025

# Resolving compound heterozygous variants

In *LHCGR* gene, associated with Leydig cell hypoplasia:

- ◆ Coding SNV on haplotype 1 (left, blue reads)
- ◆ ~7 Kbp deletion of an exon on haplotype 2 (right, red reads)



Negi et al. AJHG 2025

Episignature: methylation pattern, across 10-100s of sites, associated with disease.
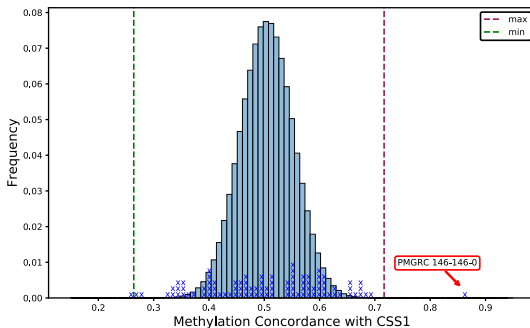
Episignature: methylation pattern, across 10-100s of sites, associated with disease.



One patient with suspected Coffin-Siris syndrome 1.

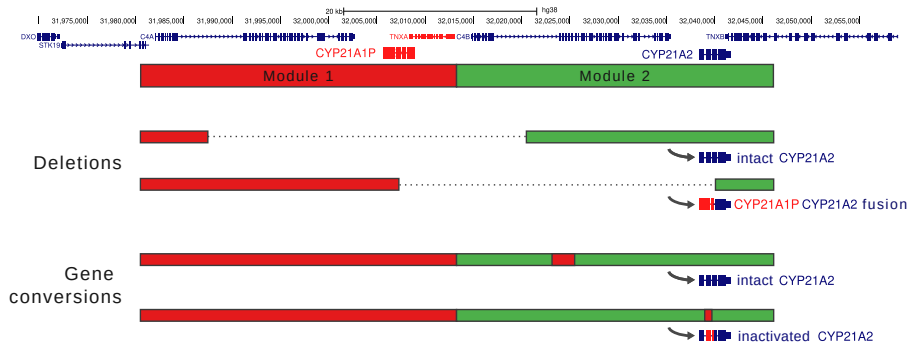Methylation across 106 differentially methylated CpG sites from Aref-Eshghi et al.

Negi et al. AJHG 2025

# Pangenomes meet long-read sequencing
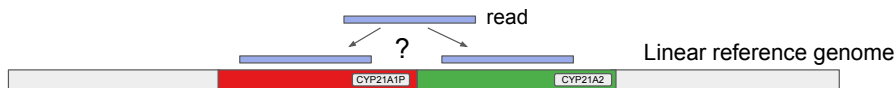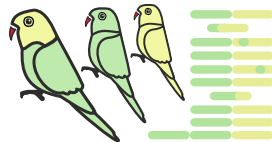
# Challenging RCCX modules in the HLA region

- Tandem-duplication of ∼30 Kbp genetic *module* (99% similar).
- CYP21A1P pseudogene and **CYP21A2 gene**.
- Variants cause congenital adrenal hyperplasia (recessive).

# **Parakit**: paralog toolkit using collapsed pangenomes

## Goal

Address multi-mapping confusion by mapping to a **collapsed pangenome** and by analyzing the alignment profile.





`https://github.com/jmonlong/parakit` Monlong et al. medRxiv 2025

# Parakit: paralog toolkit using collapsed pangenomes

## Goal

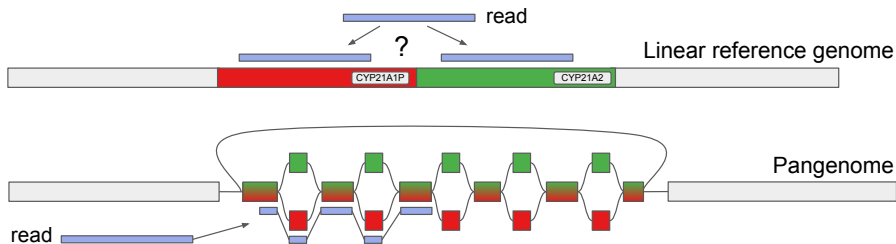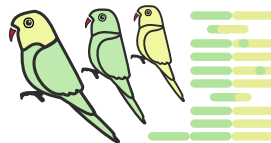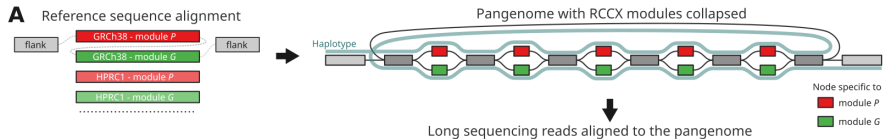Address multi-mapping confusion by mapping to a **collapsed pangenome** and by analyzing the alignment profile.



read

?

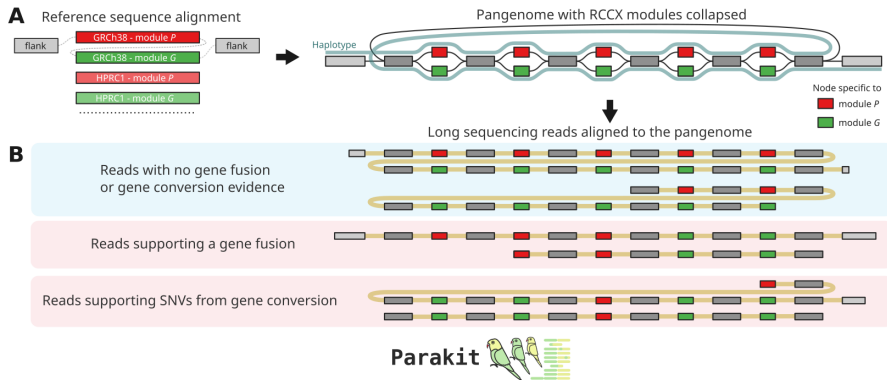Linear reference genome
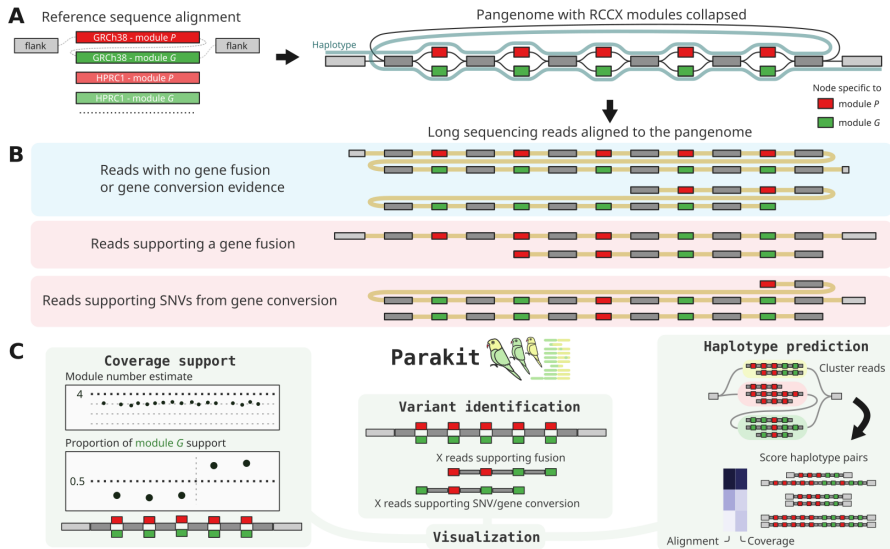
CYP21A1P

CYP21A2

Pangenome

read

https://github.com/jmonlong/parakit   Monlong et al. medRxiv 2025

**A** Reference sequence alignment

Pangenome with RCCX modules collapsed

Long sequencing reads aligned to the pangenome

Parakit

**A** Reference sequence alignment

Pangenome with RCCX modules collapsed

Long sequencing reads aligned to the pangenome

Node specific to
module P
module G

**B**

Reads with no gene fusion or gene conversion evidence

Reads supporting a gene fusion

Reads supporting SNVs from gene conversion

Parakit

https://github.com/jmonlong/parakit   Monlong et al. medRxiv 2025

https://github.com/jmonlong/parakit  Monlong et al. medRxiv 2025

# Example: patients with a gene fusion and pathogenic SNV
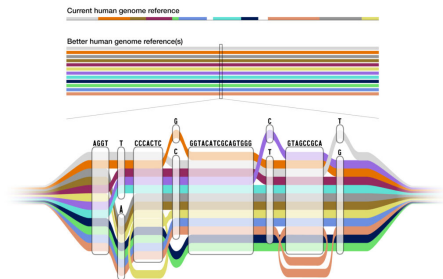


https://github.com/jmonlong/parakit  Monlong et al. medRxiv 2025

Two approaches to integrate structural variants into genomic studies:

Genotyping with **pangenomes** from **short-read sequencing** data, e.g. for **genome-wide association studies**.
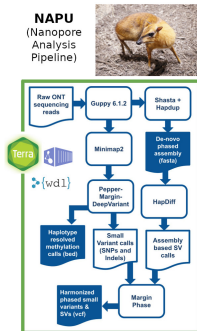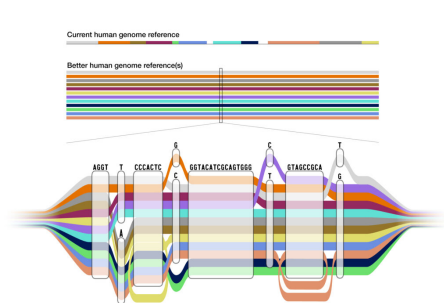
# Conclusions

Two approaches to integrate structural variants into genomic studies:

Genotyping with **pangenomes** from **short-read sequencing** data, e.g. for **genome-wide association studies**.

Cost-effective **long-read sequencing** using nanopore technologies to help solve undiagnosed **rare disease** cases.

**Methods, tools, benchmark, and analysis needed !**



Pangenomes with haplotype-resolved near-complete genomes.

Single-molecule long read sequencing (nanopore, PacBio).

**Methods, tools, benchmark, and analysis needed !**

Pangenomes with haplotype-resolved near-complete genomes.

Single-molecule long read sequencing (nanopore, PacBio).

- Complex variants
- Repeat-rich regions
- Association studies
- Functional genomics
- Epigenomics
- ...

# Acknowledgments

**Univ. California, Santa Cruz**

- **Benedict Paten**
- **Glenn Hickey**
- Jouni Sirén
- Adam Novak
- Xian Chang
- Jordan Eizenga
- **Shloka Negi**
- **Karen Miga**
- Brandy McNulty
- Melissa Meredith
- Paolo Carnevali
- Trevor Pesout
- Kishwar Shafin
- Mira Mastoras
- Mobin Asri

**INSERM IRSD**

- Sarah Djebali
- Matis Alias-Bagarre

**NIH**

- **Mikhail Kolmogorov**
- Cornelis Blauwendraat
- Kimberley Billingsley
- Pilar Alvarez Jerez

**Broad Institute**

- Anne O'Donnell-Luria
- Sarah Stenton
- Melanie O'Leary

**Univ. California, Irvine**
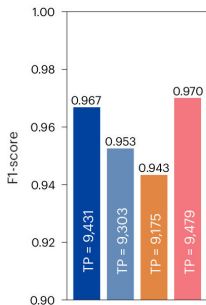
- Emmanuèle Délot
- Eric Vilain

**Children's National Research Institute**
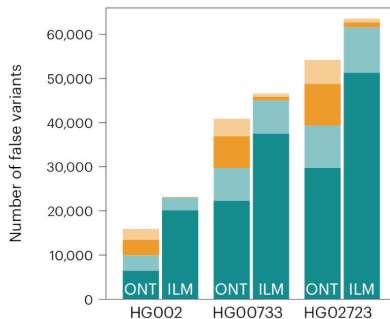
- Seth Berger
- Paolo Canigiula

# Better calls for both small and structural variants...



SV concordance with GIAB HG002 benchmark

Whole genome SNP performance, stratified by local context

Legend:
- Hapdup (ONT)
- Sniffles2 (ONT)
- CuteSV (ONT)
- Hifiasm (HiFi)

FP / FN Homopolymers
FP / FN Not in homopolymers

Kolmogorov*, Billingsley*, et al. Nature Methods 2023

# ...except for indels in homopolymers



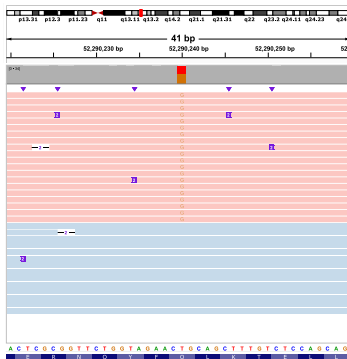Whole genome indel performance, stratified by local context

*Note: Results above are for the R9 chemistry. The new R10 chemistry has lower error rate and better (indel) calling performance.*

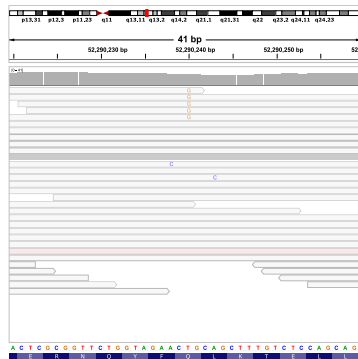Kolmogorov*, Billingsley*, et al. Nature Methods 2023

# Small variants found by long-reads only

Missense mutation in *KRT86* disease gene (monilethrix) invisible with short reads.

**chr12:52,290,220-52,290,259**



long-reads

short-reads

Negi et al. AJHG 2025

# Patient with complex neurodevelopmental phenotype

Variant of Uncertain Significance SNV in *ARID1B* gene (Coffin-Siris syndrome 1?).

◆ *De novo*, SRS and LRS, new splice site predicted *in silico* (SpliceAI).



Negi et al. AJHG 2025