

C G T A C G T A  
A C G T A C G T

# What's new in the human (pan)genome

Adam M. Phillippy

CSHL Advanced Sequencing Technologies & Bioinformatics Analysis Course  
November 22, 2024

@aphillippy.bsky.social 



National Human Genome  
Research Institute

The **Forefront**  
of **Genomics**®

# A 20-year anniversary

articles

## Finishing the **euchromatic** sequence of the human genome

International Human Genome Sequencing Consortium\*

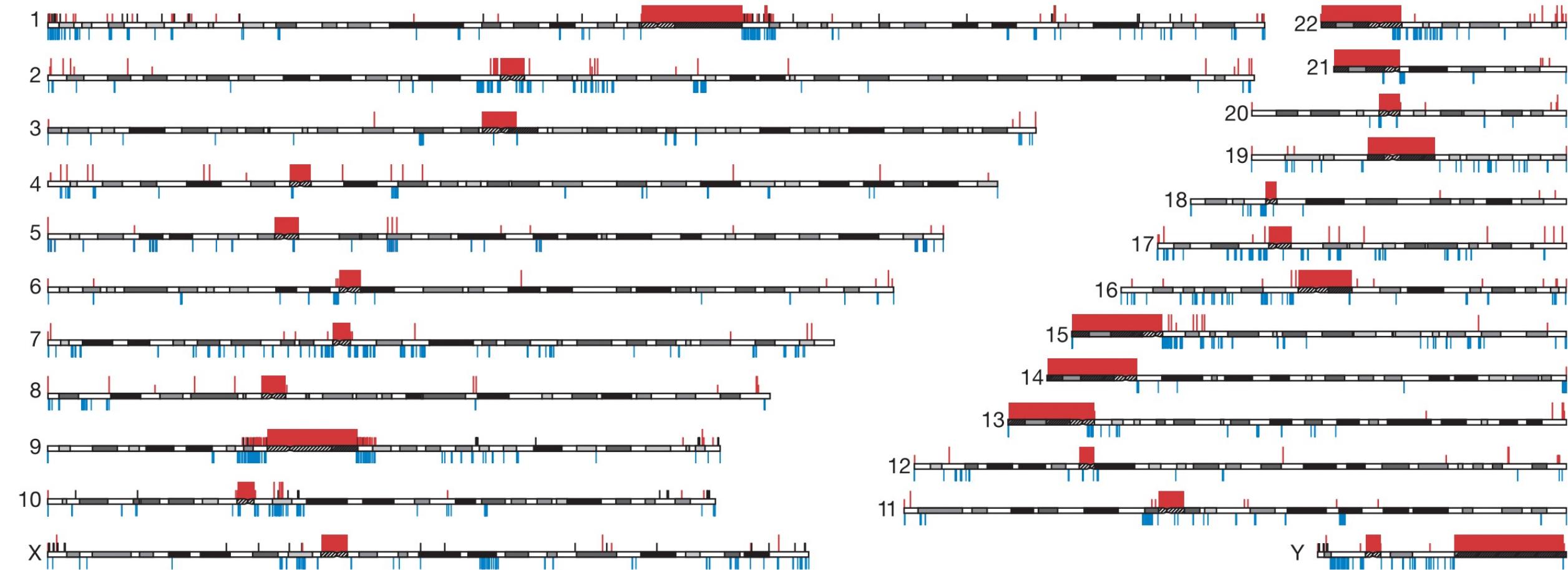
\* A list of authors and their affiliations appears in the Supplementary Information

---

The sequence of the human genome encodes the genetic instructions for human physiology, as well as rich information about human evolution. In 2001, the International Human Genome Sequencing Consortium reported a draft sequence of the euchromatic portion of the human genome. Since then, the international collaboration has worked to convert this draft into a genome sequence with high accuracy and nearly complete coverage. Here, we report the result of this finishing process. The current genome sequence (Build 35) contains 2.85 billion nucleotides interrupted by only 341 gaps. It covers ~99% of the euchromatic genome and is accurate to an error rate of ~1 event per 100,000 bases. Many of the remaining euchromatic gaps are associated with segmental duplications and will require focused work with new methods. The near-complete sequence, the first for a vertebrate, greatly improves the precision of biological analyses of the human genome including studies of gene number, birth and death. Notably, the human genome seems to encode only 20,000–25,000 protein-coding genes. The genome sequence reported here should serve as a firm foundation for biomedical research in the decades ahead.

# Where were the gaps?

Satellites  
Segmental duplications  
Gene arrays (e.g. rDNA)



# Finishing the human genome

**2003 (hg35)**

- 2.85 Gbp
- 341 gaps
- 1 error per 100,000 bases?
- \$5,000,000,000

**2023 (T2T-CHM13)**

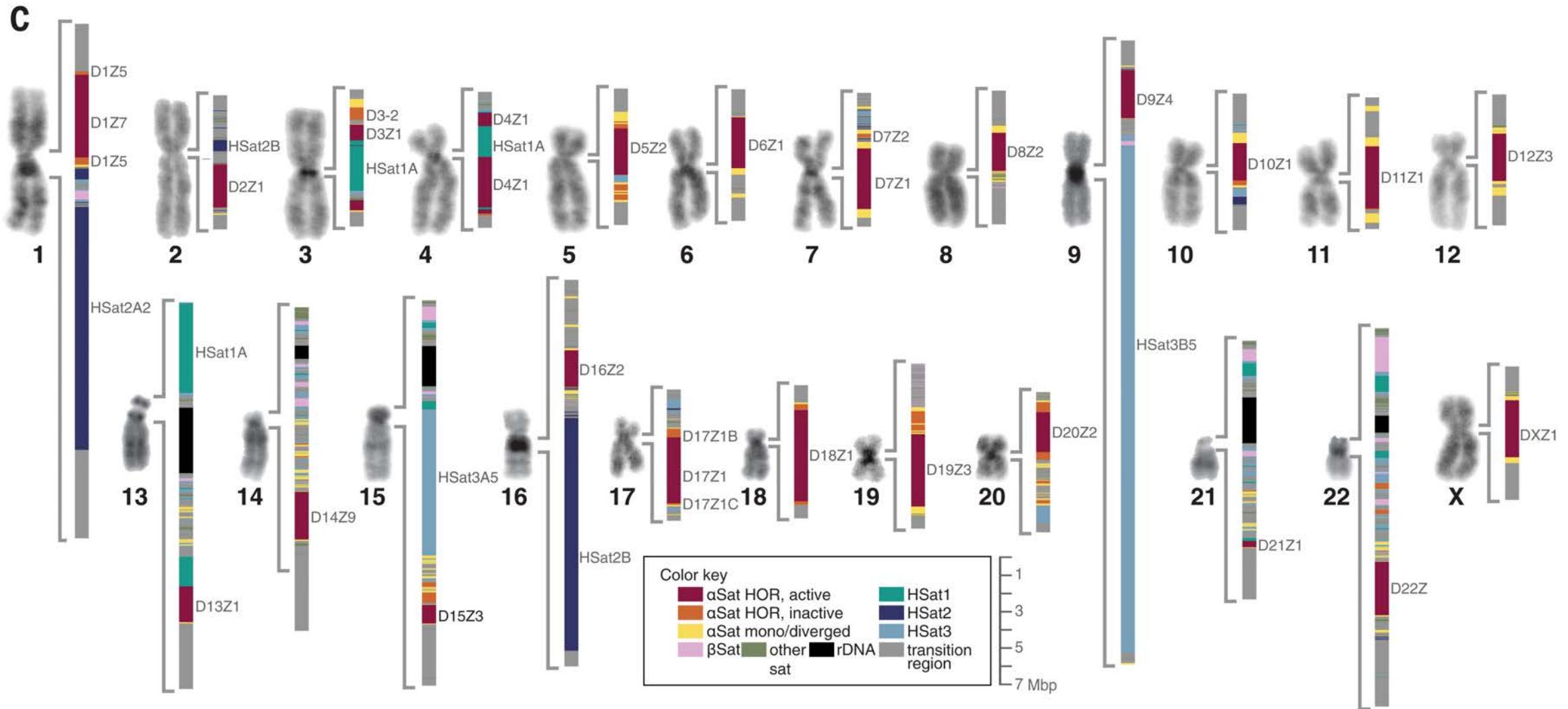
- 3.12 Gbp
- 0 gaps
- 1 error per 10,000,000 bases
- \$5,000



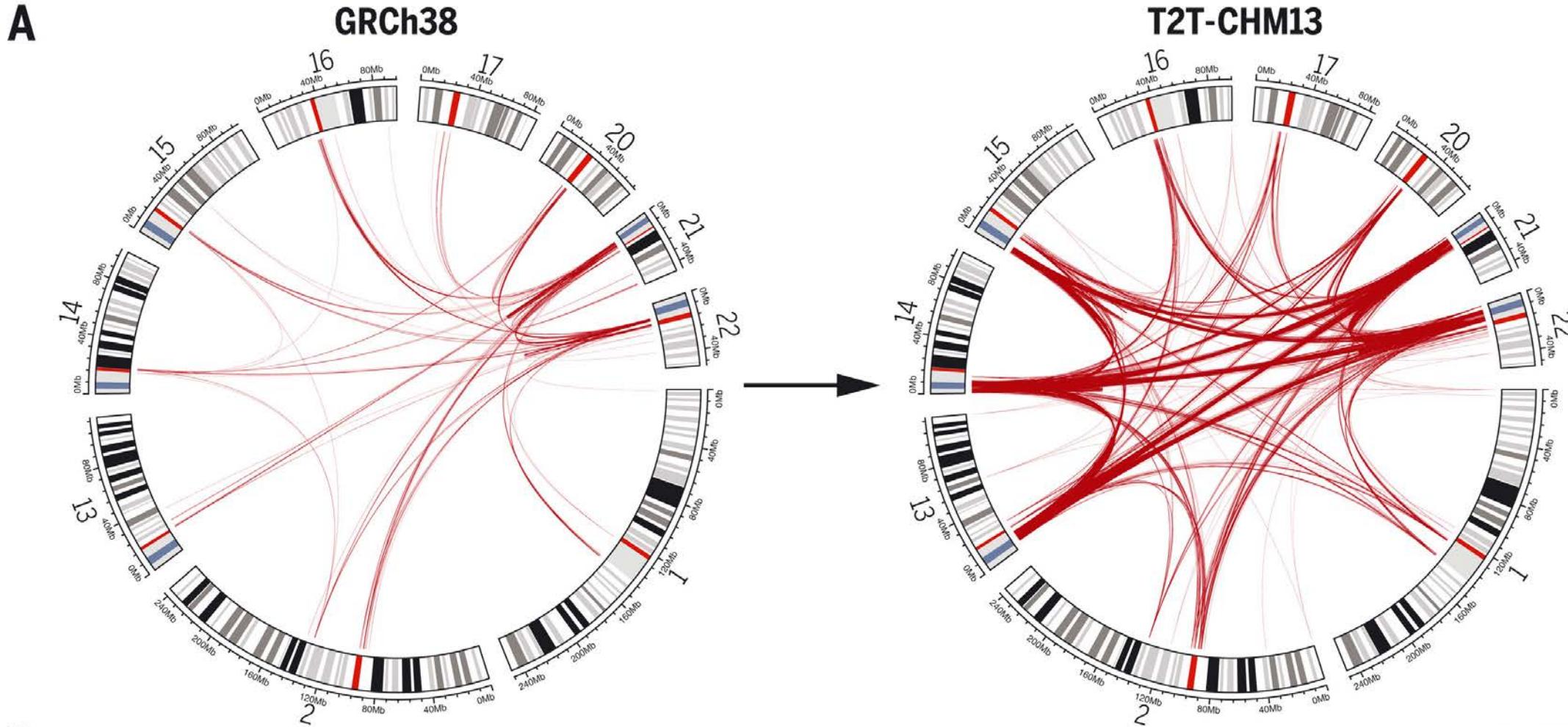
**MILLION-fold reduction in cost!**

Driven by decades of investment and technology development

# New satellites in CHM13

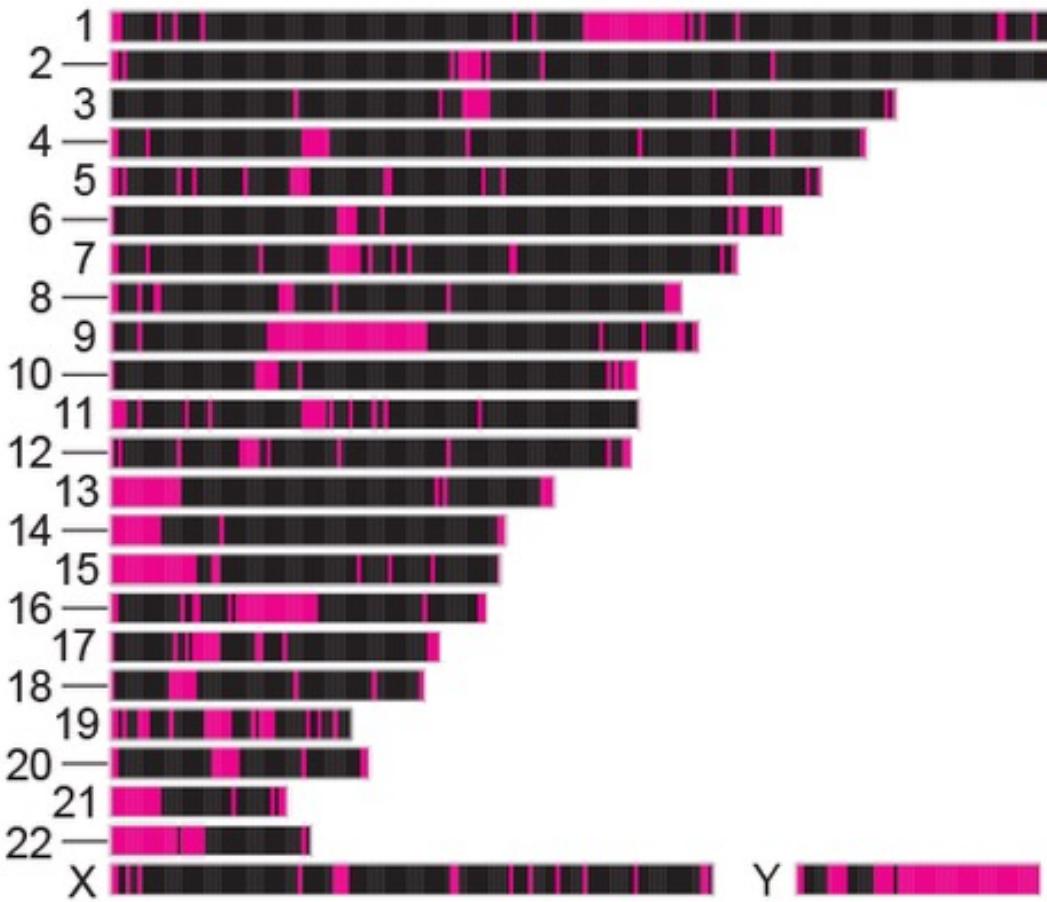


# New segdups in CHM13



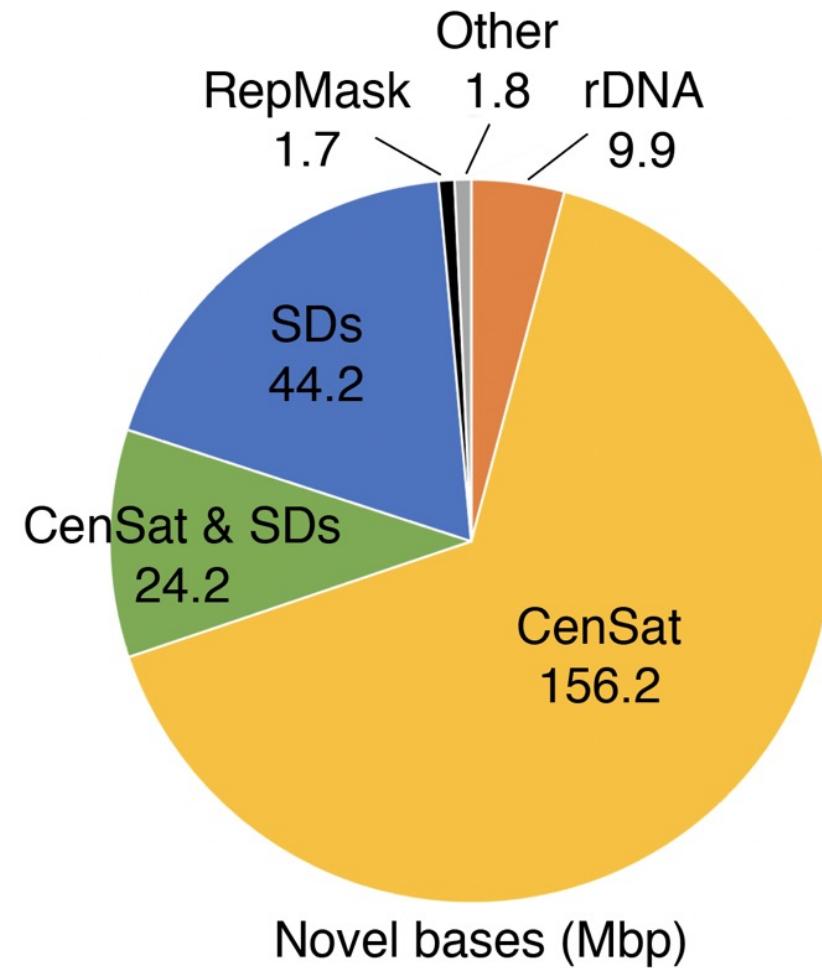
# Illuminating the dark regions of the human genome

# What's new? Repeats!



The complete sequence of a human genome.  
Nurk, Koren, Rhie, Rautiainen, et al. *Science* (2022)

The complete sequence of a human Y chromosome.  
Rhie, Nurk, et al. *Nature* (2023)

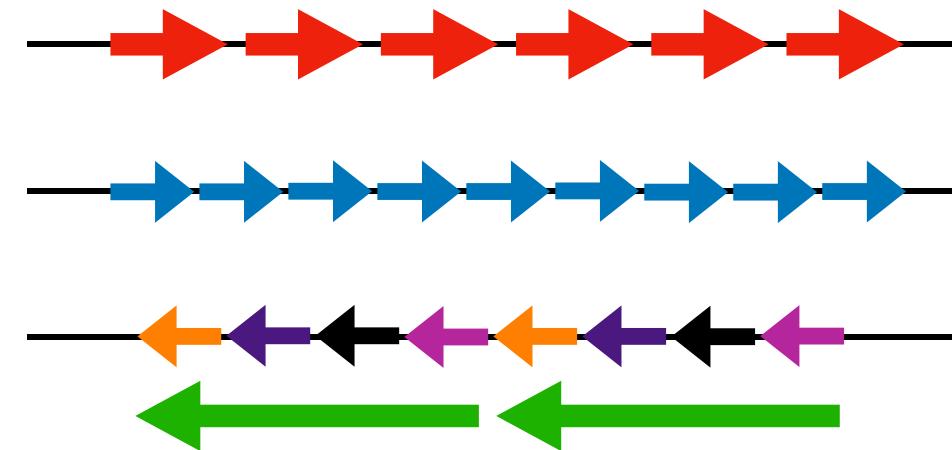
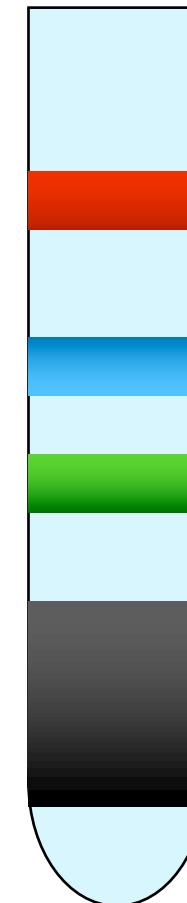


- >250 Mbp of new sequence
- ~2,000 new genes predicted



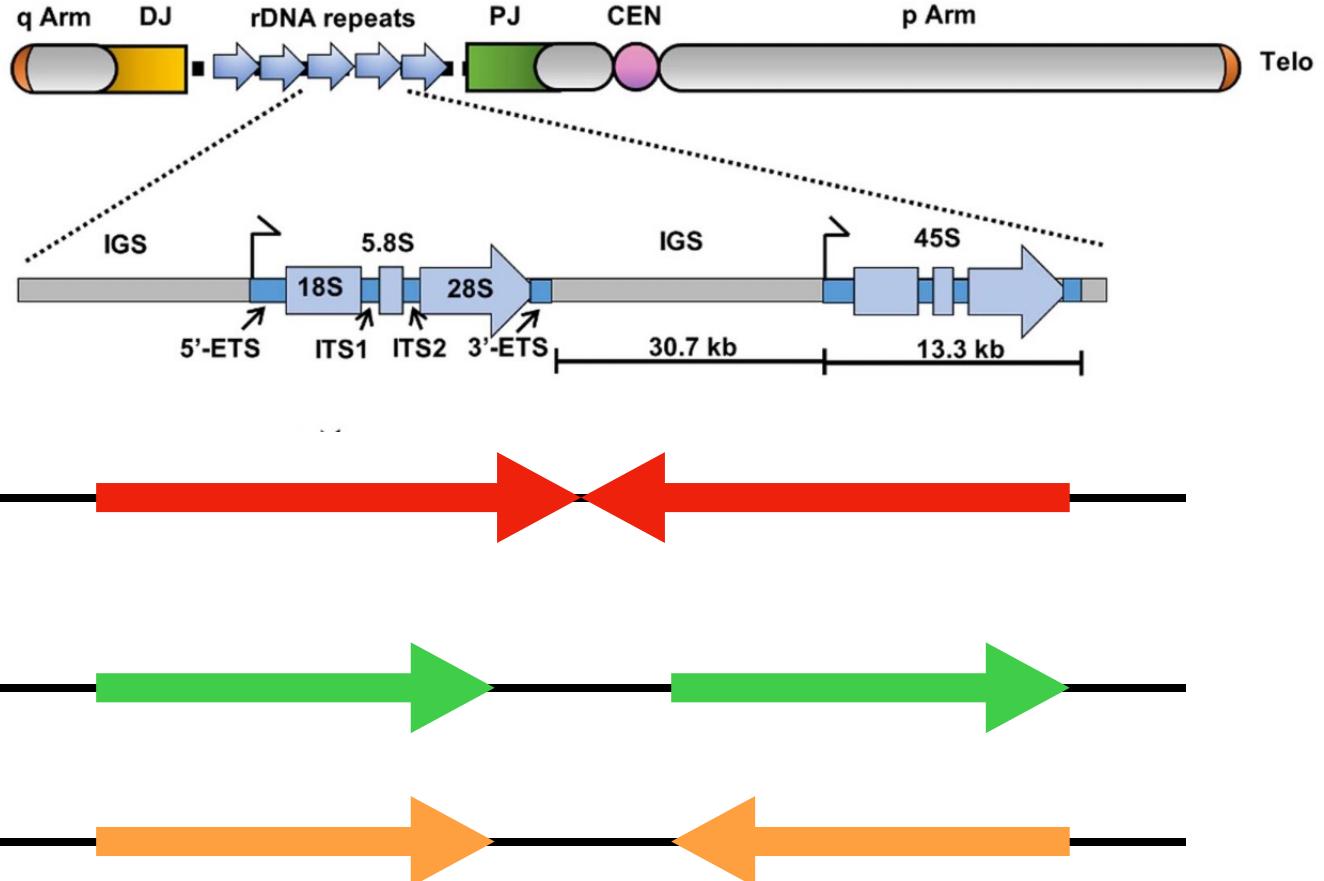
# Satellite DNA

- Consists of **long arrays of tandemly repeating sequences**
- Primarily found in **centromeres** and **acrocentric** short arms
- Uniform repeats result in **unique properties** compared to bulk DNA

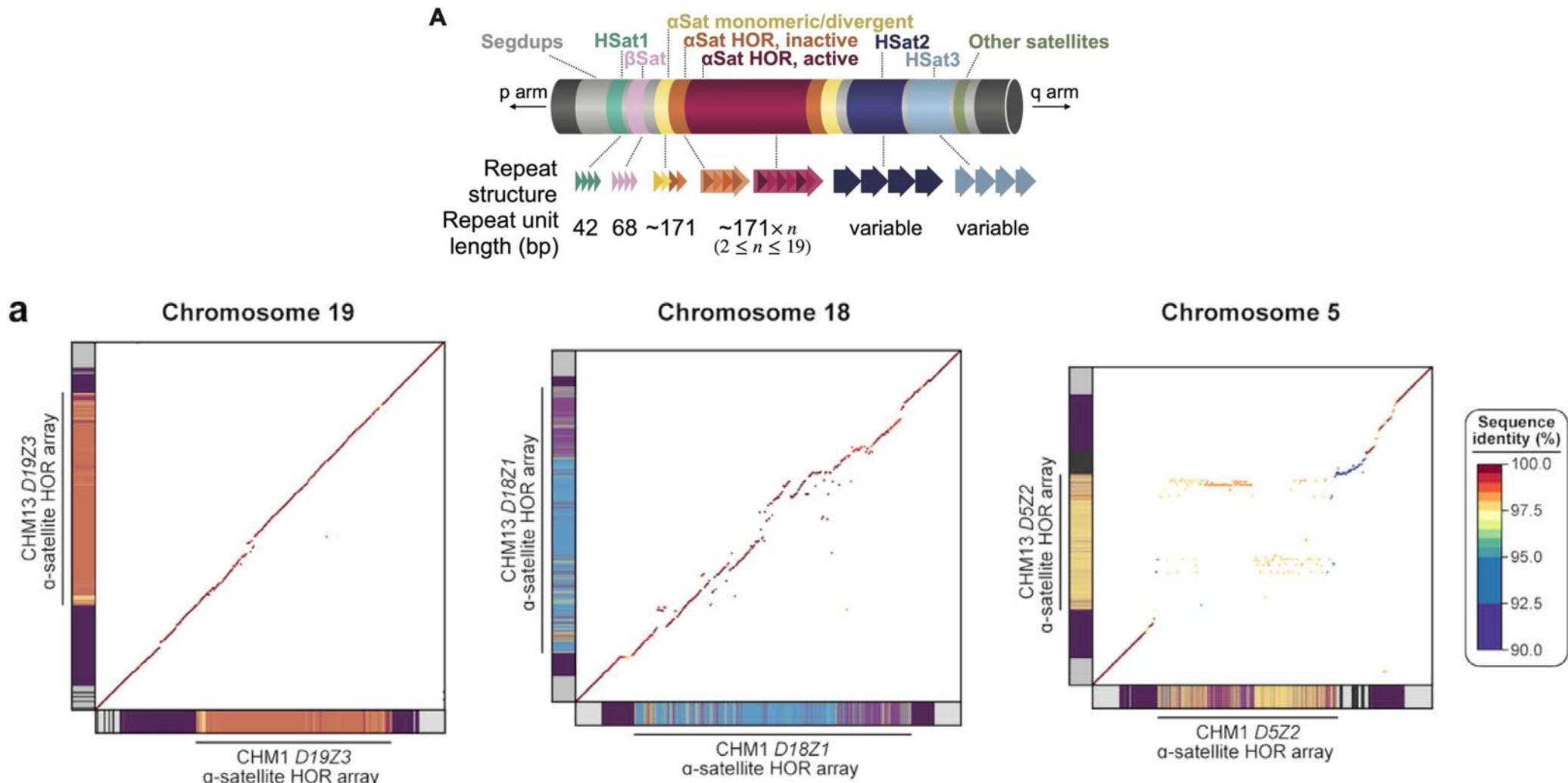


# Segmental duplications

- Gene arrays
- Palindromes
- Directed repeats
- Inverted repeats



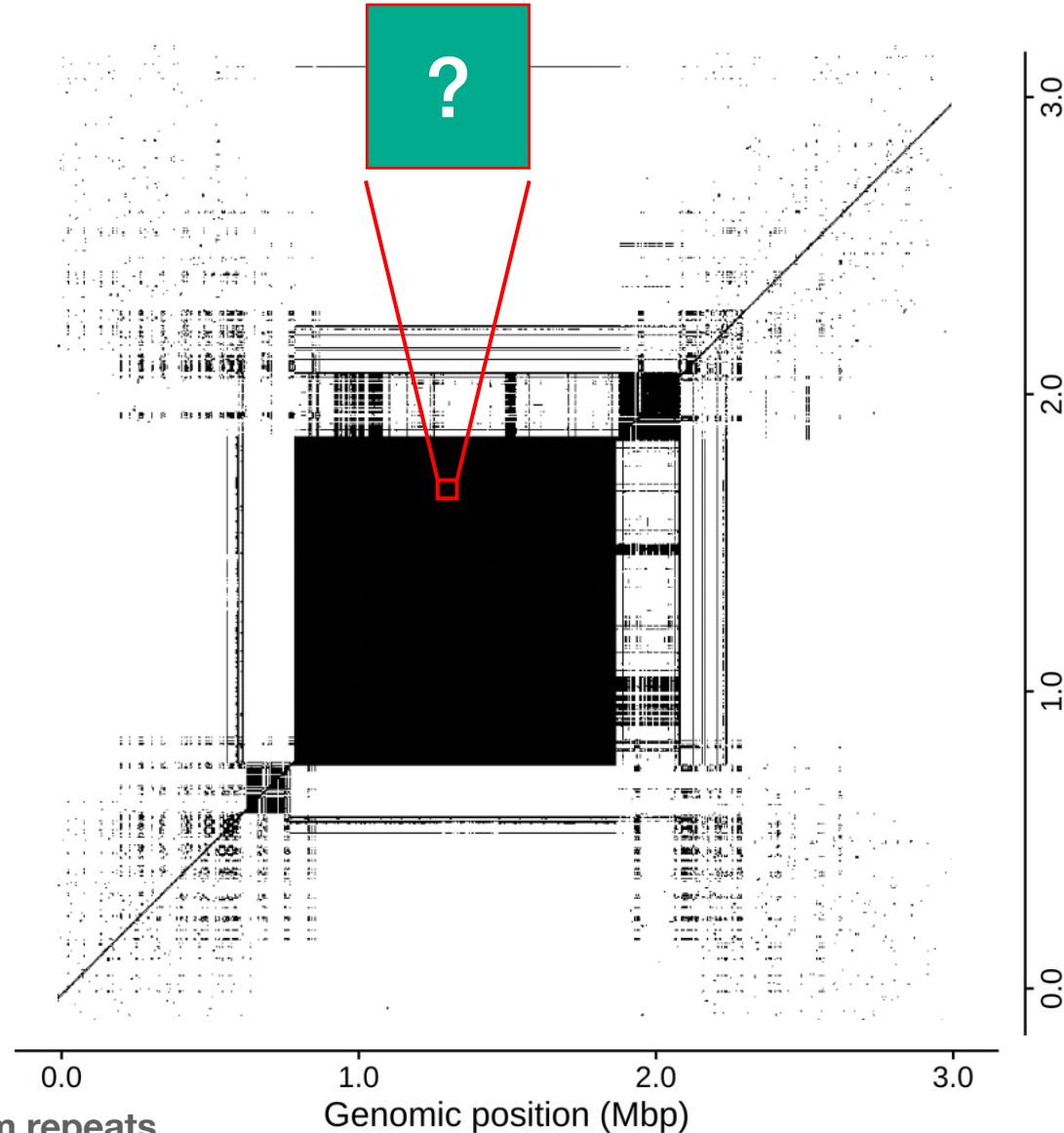
# Satellites are often unalignable



# A better dotplot

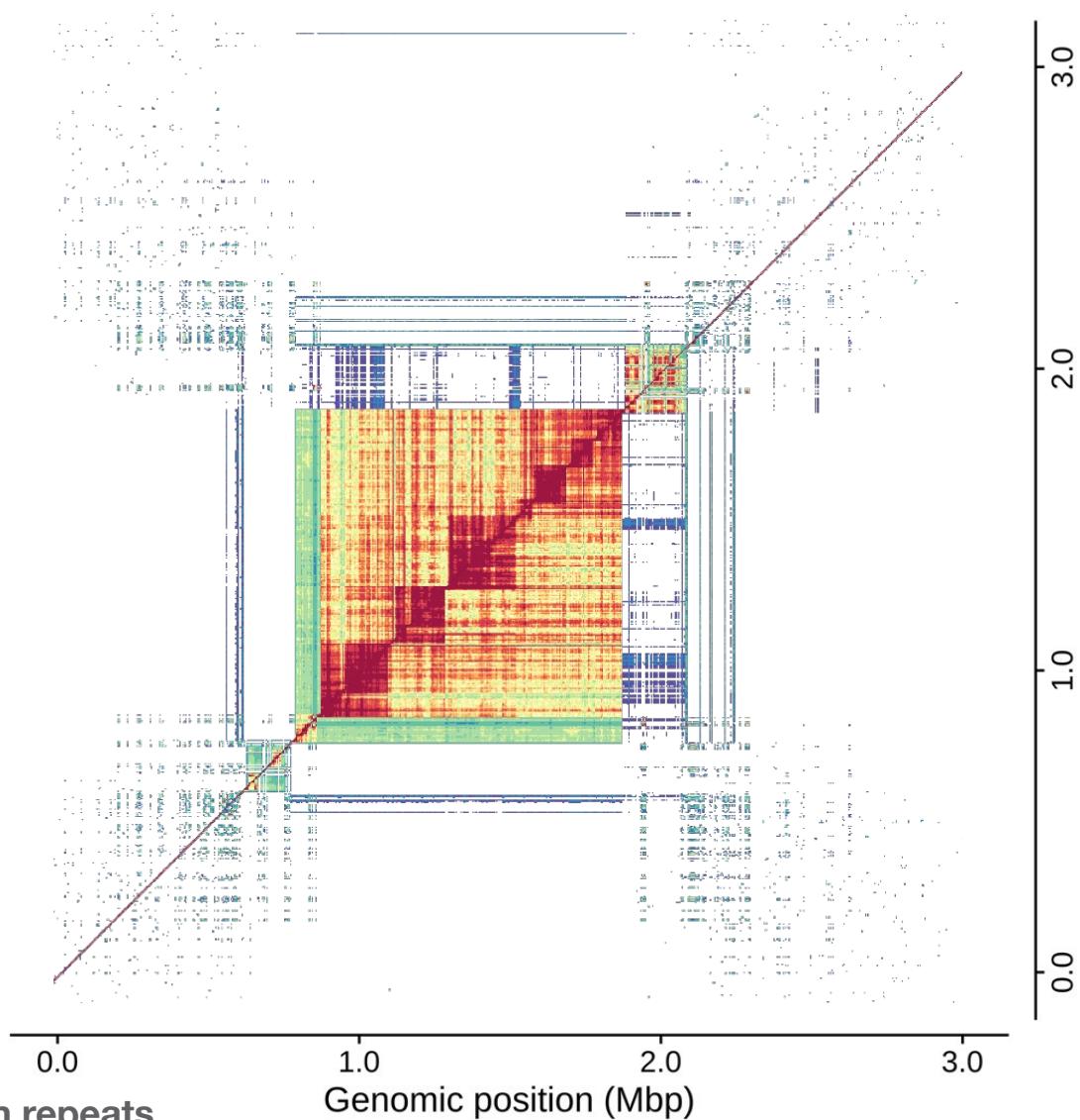
- Each pixel represents a pair of intervals of a fixed size
- Efficiently estimate the sequence identity for each pixel using min-hashing

E.g. 1 px = 1,000 x 1,000 bp

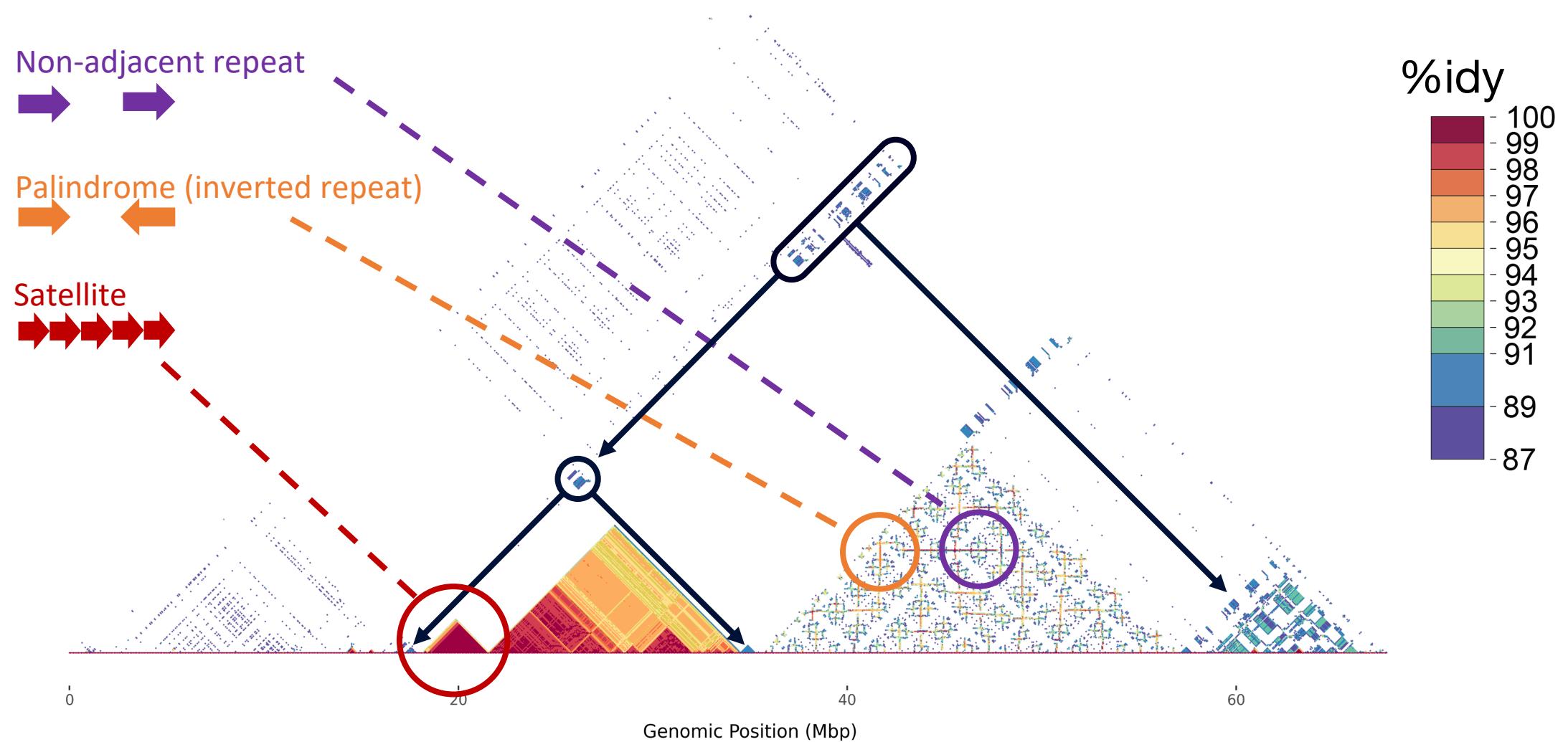


# A better dotplot

- Each pixel represents a pair of intervals of a fixed size
- Efficiently estimate the sequence identity for each pixel using min-hashing
- Reveals the structure and evolution of complex repeats



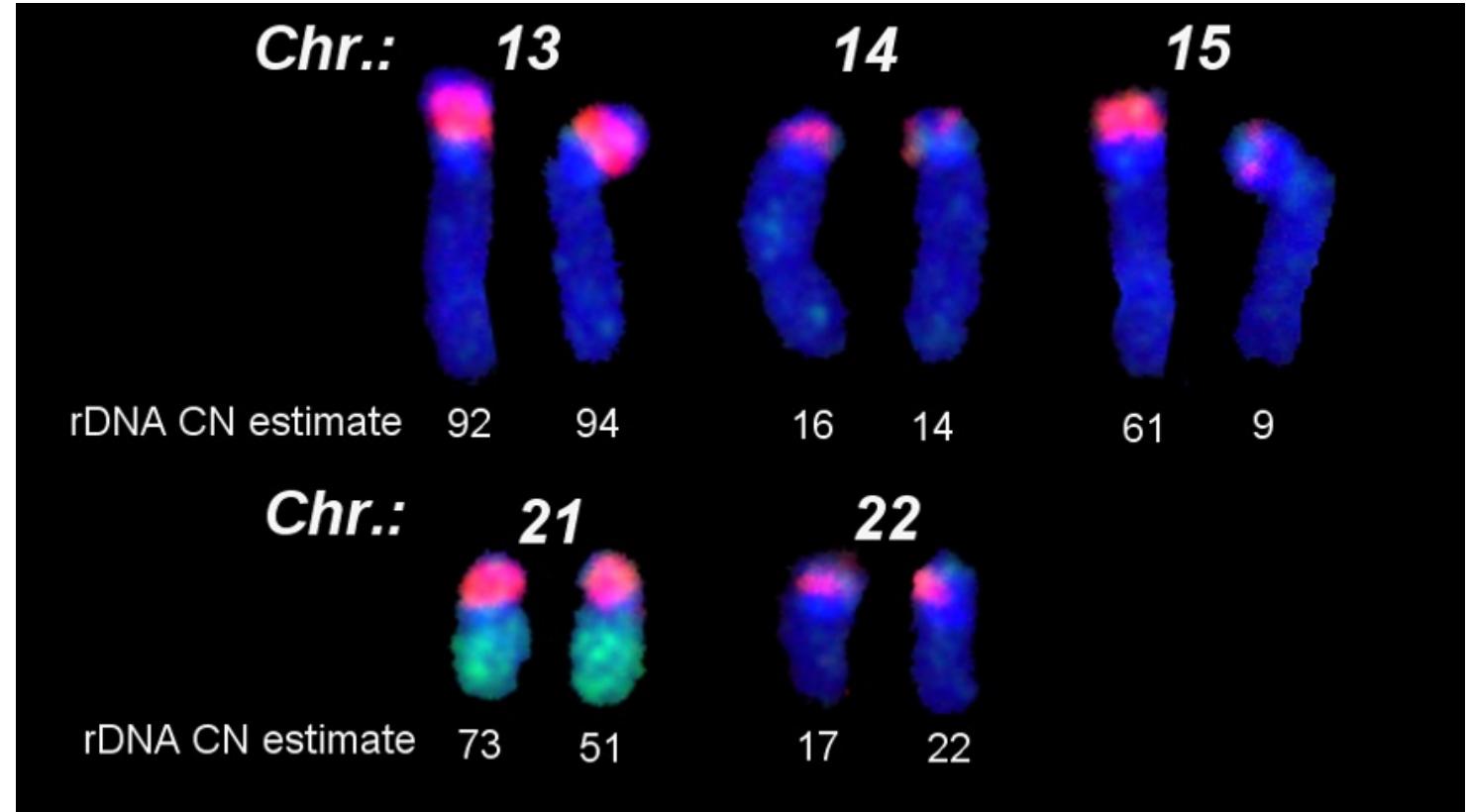
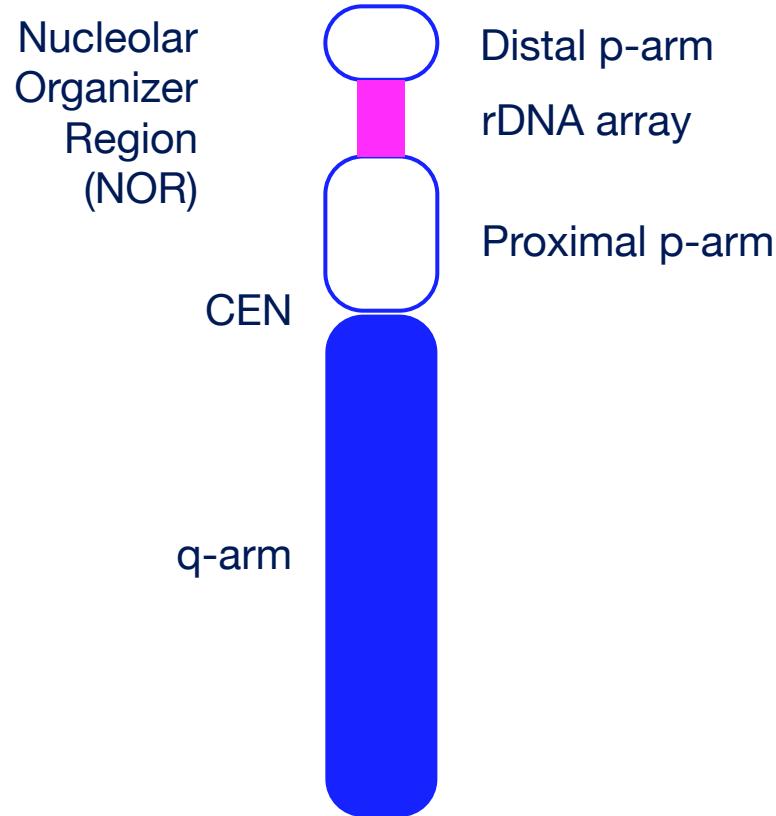
# StainedGlass & ModDotPlot



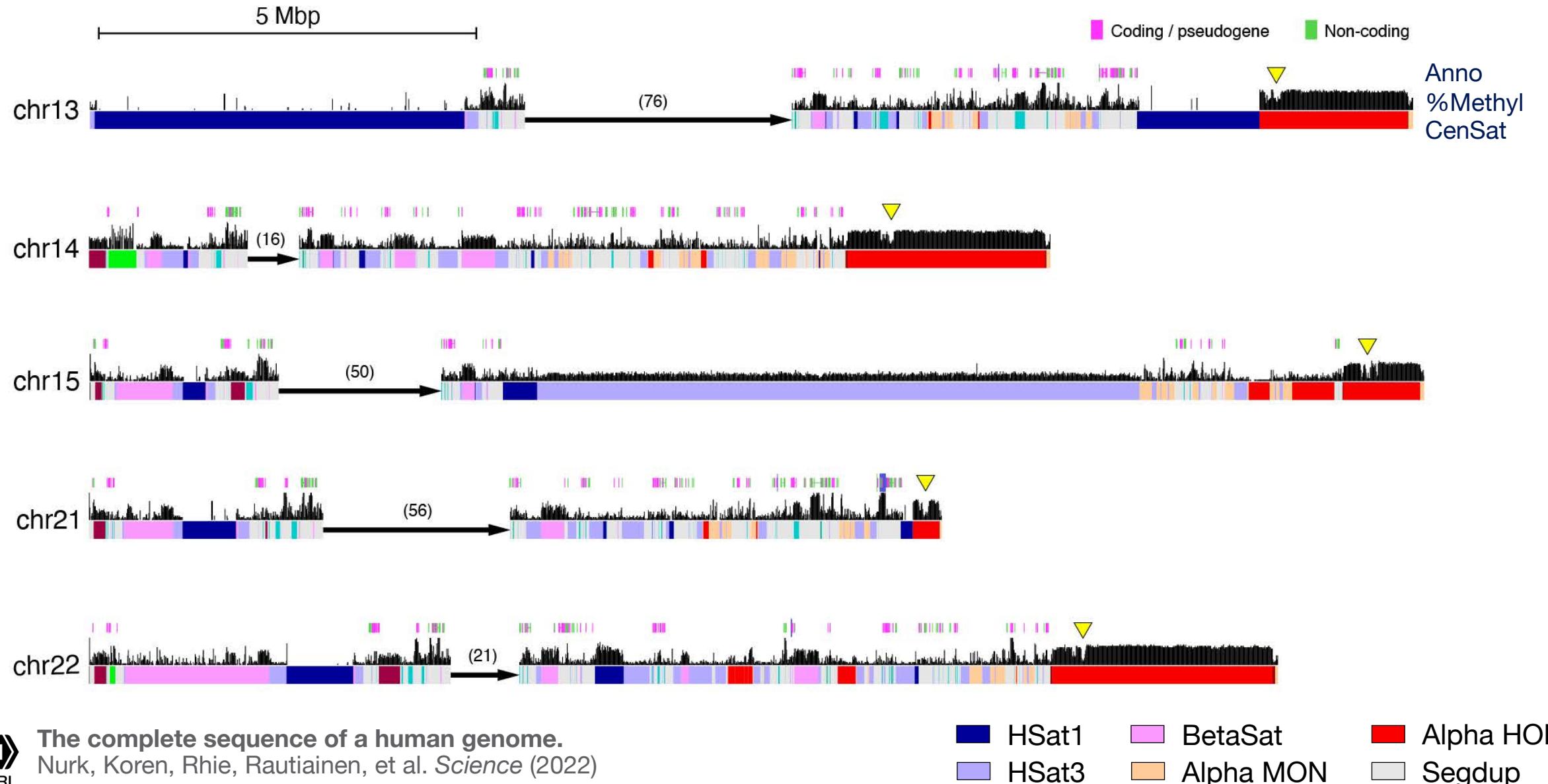
<https://github.com/mrvollger/StainedGlass>  
<https://github.com/marbl/ModDotPlot>

# The acrocentrics: Satellites and segdups everywhere

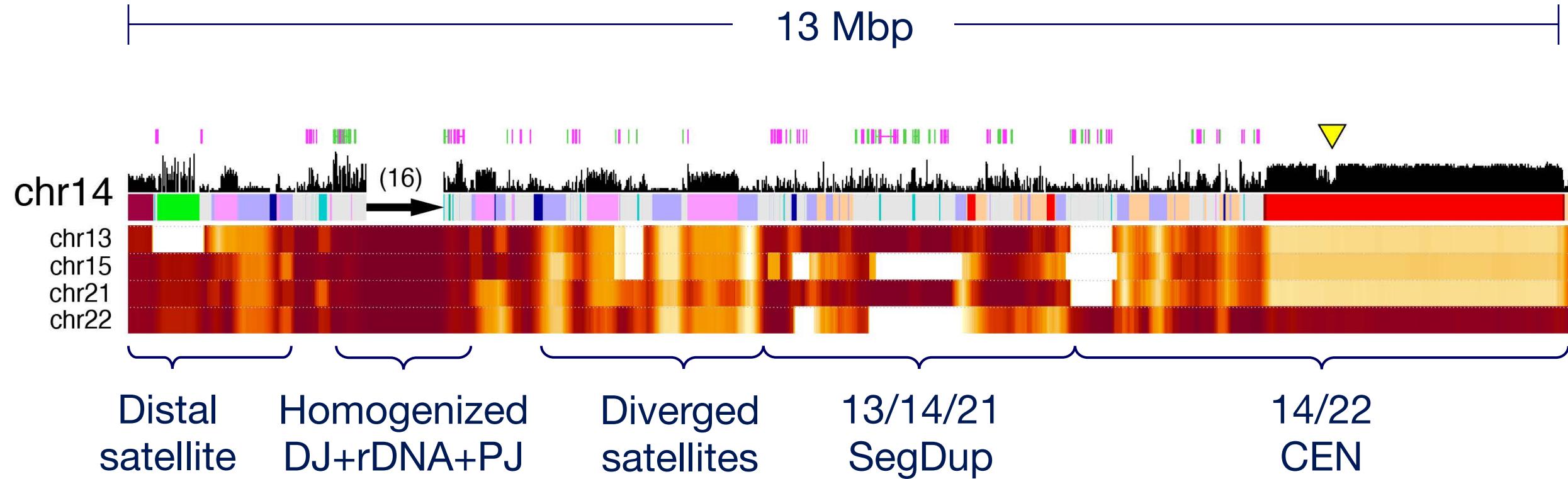
# The human acrocentrics



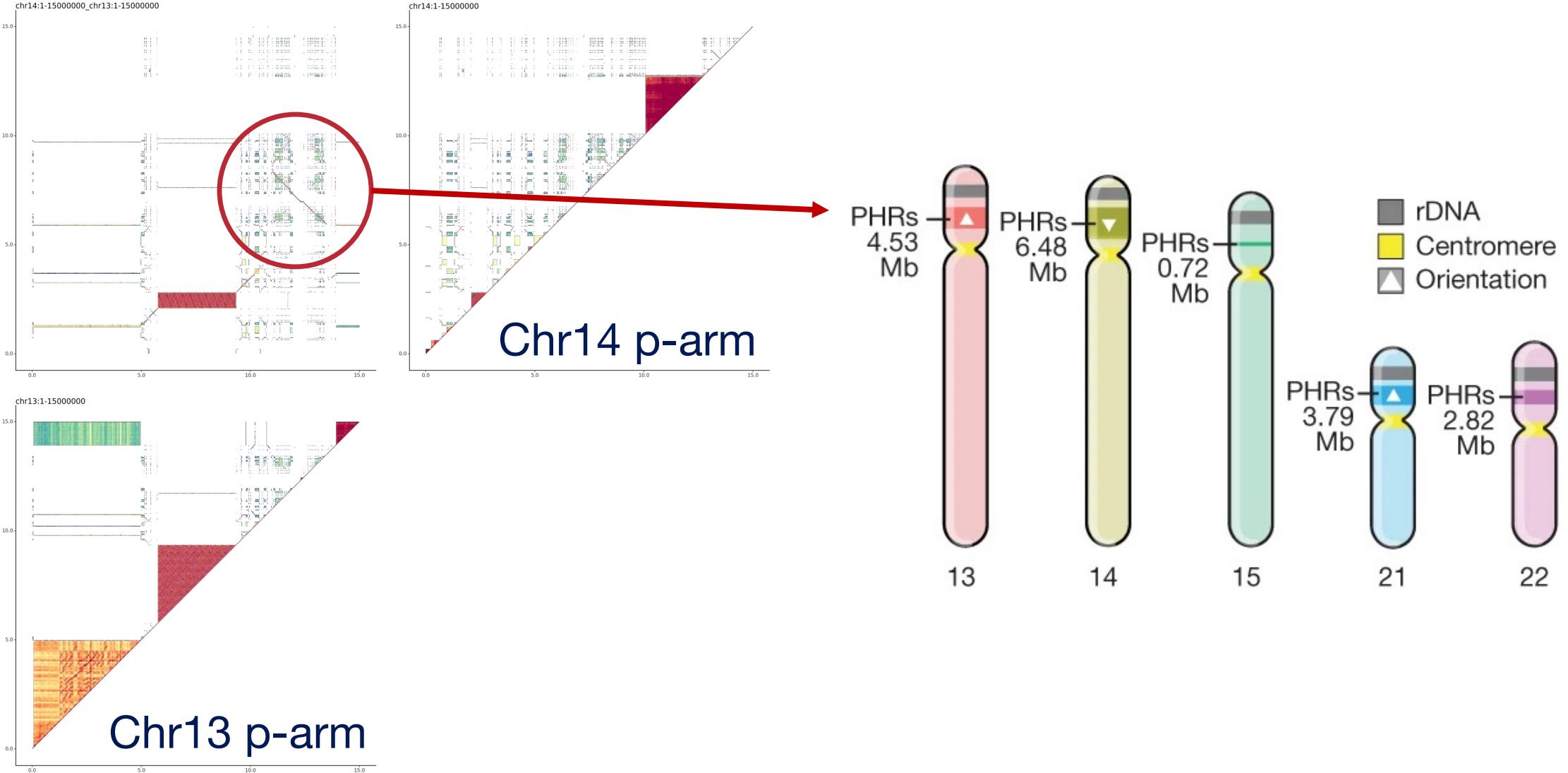
# Short arms of the acrocentrics



# High sequence similarity between across

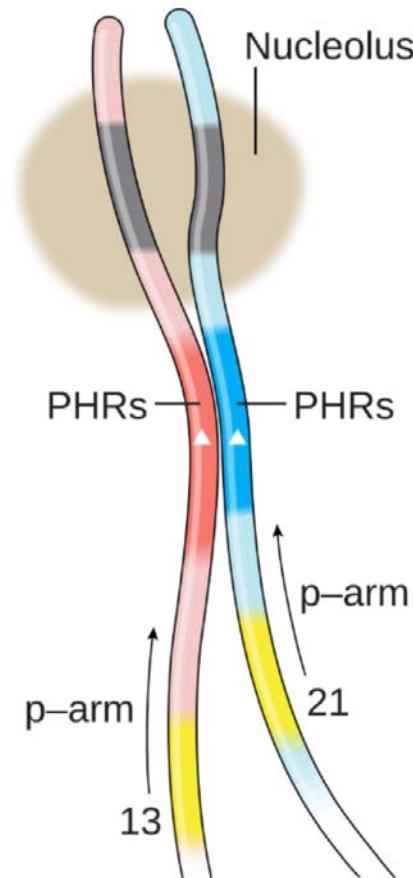


# Pseudo-homologous regions

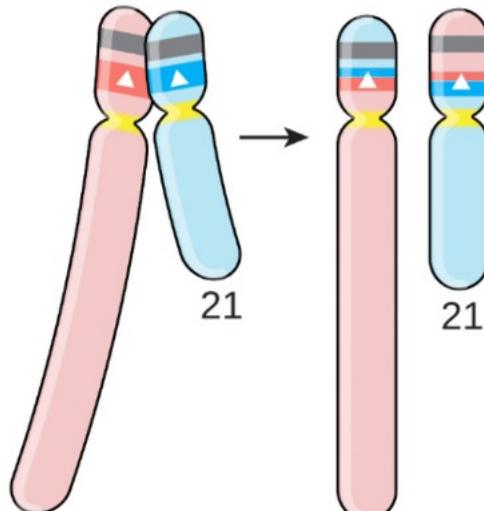


# Short arm recombination

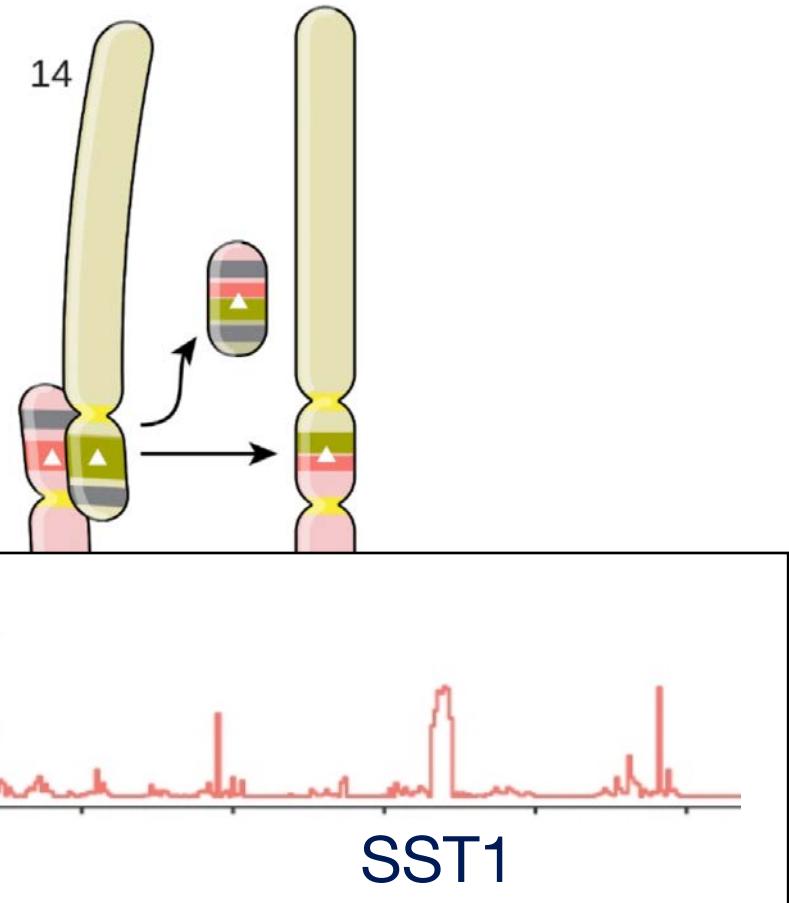
B. Physical Proximity



Heterologous  
C. Recombination

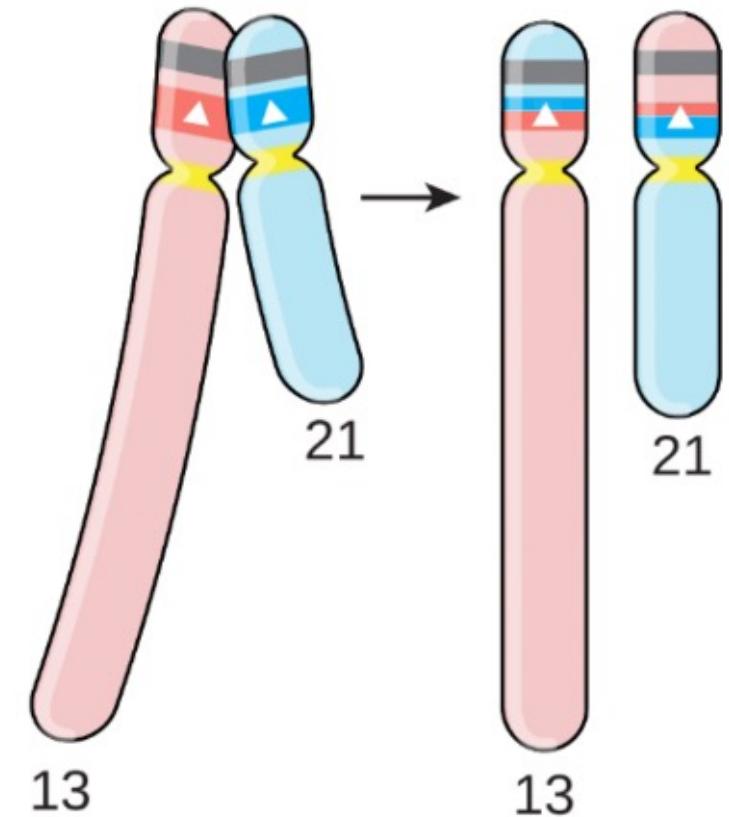


D. ROBs  
Robertsonian Translocations

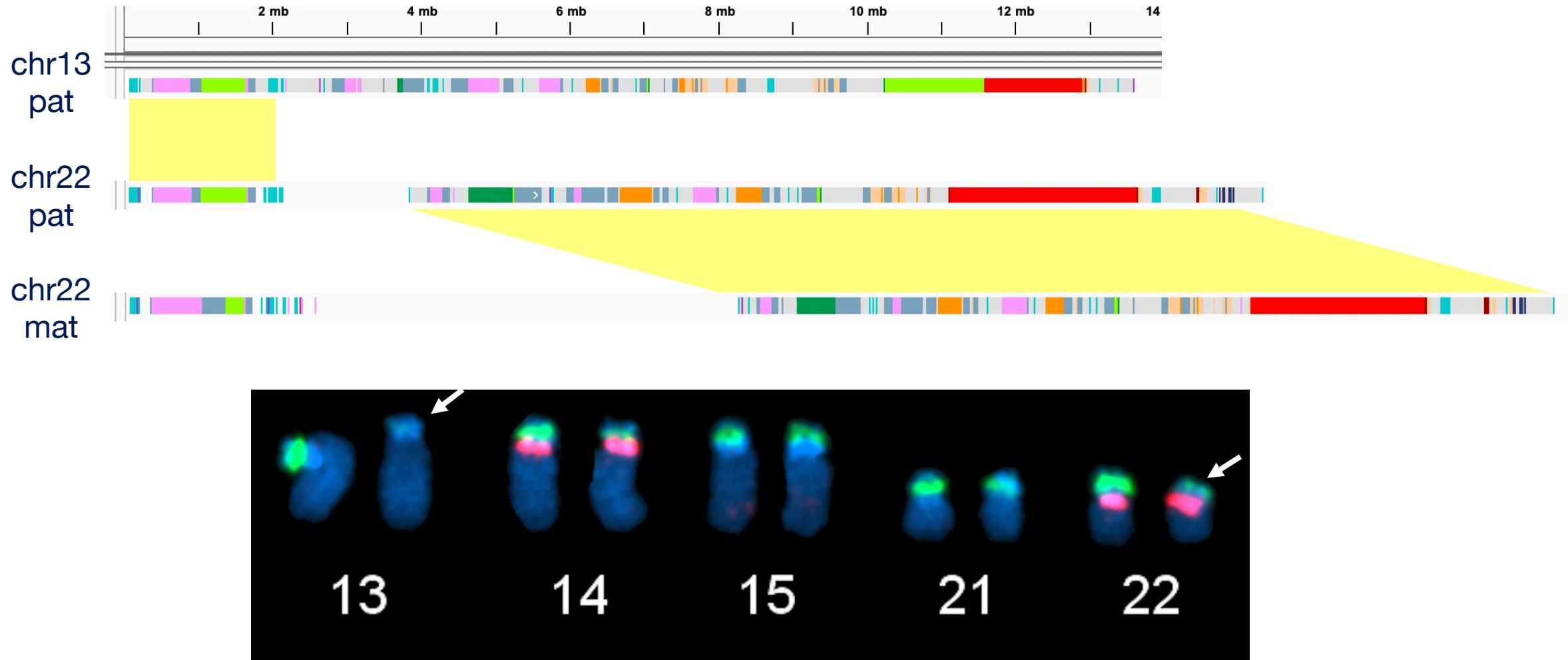


# Why are rDNAs on the acros?

- rDNAs seem to appear either:
  - In one array on a metacentric chr
  - In multiple arrays on short arms
- Why?
  - Permissive of crossovers
  - Maintains rDNA concerted evolution
- Conclusion
  - Human NORs not likely chr-specific

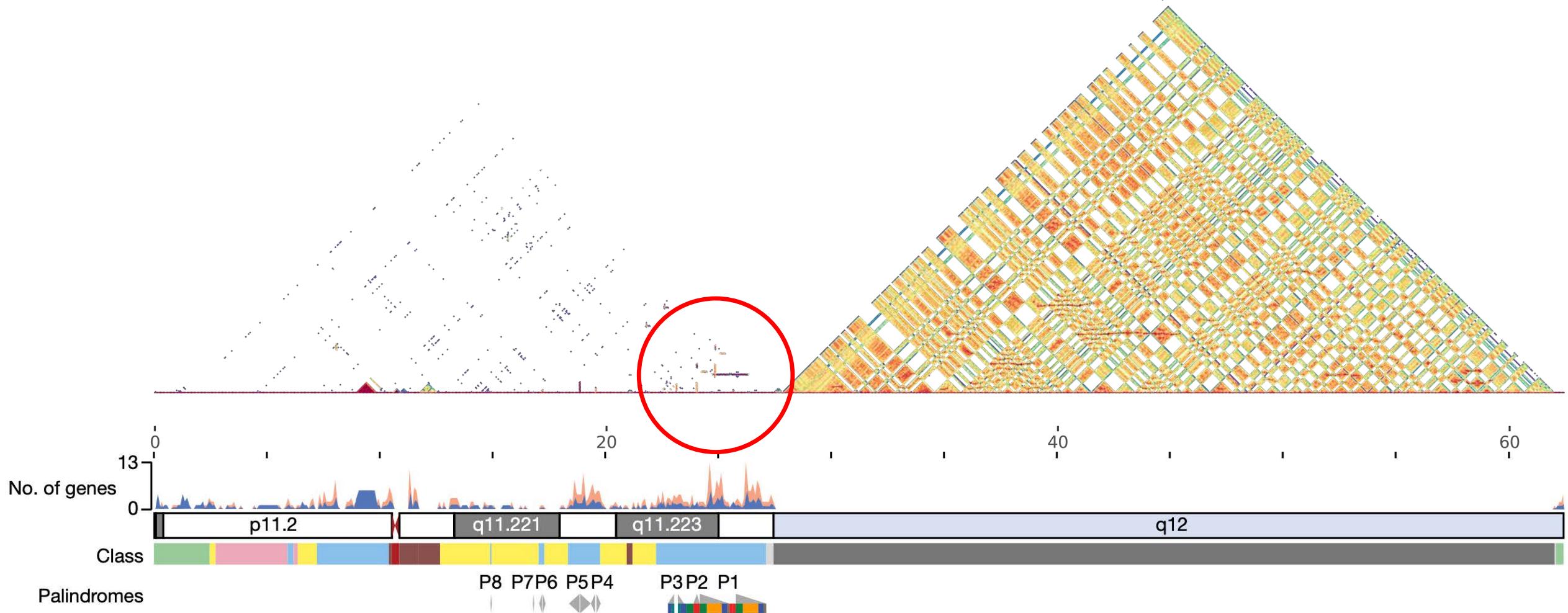


# HG002 distal copy between 13/22

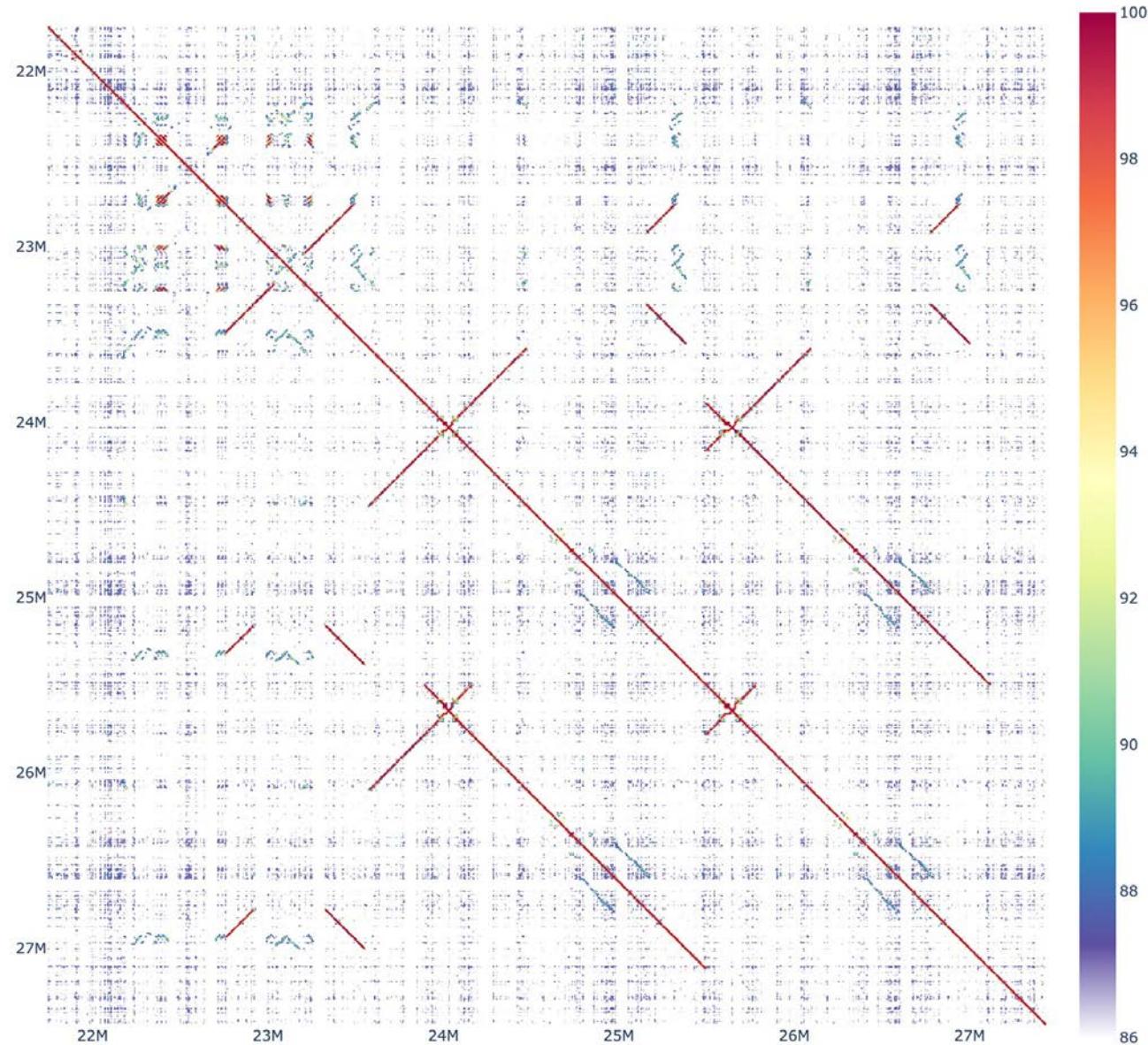


# The Y chromosome: A wannabe acrocentric

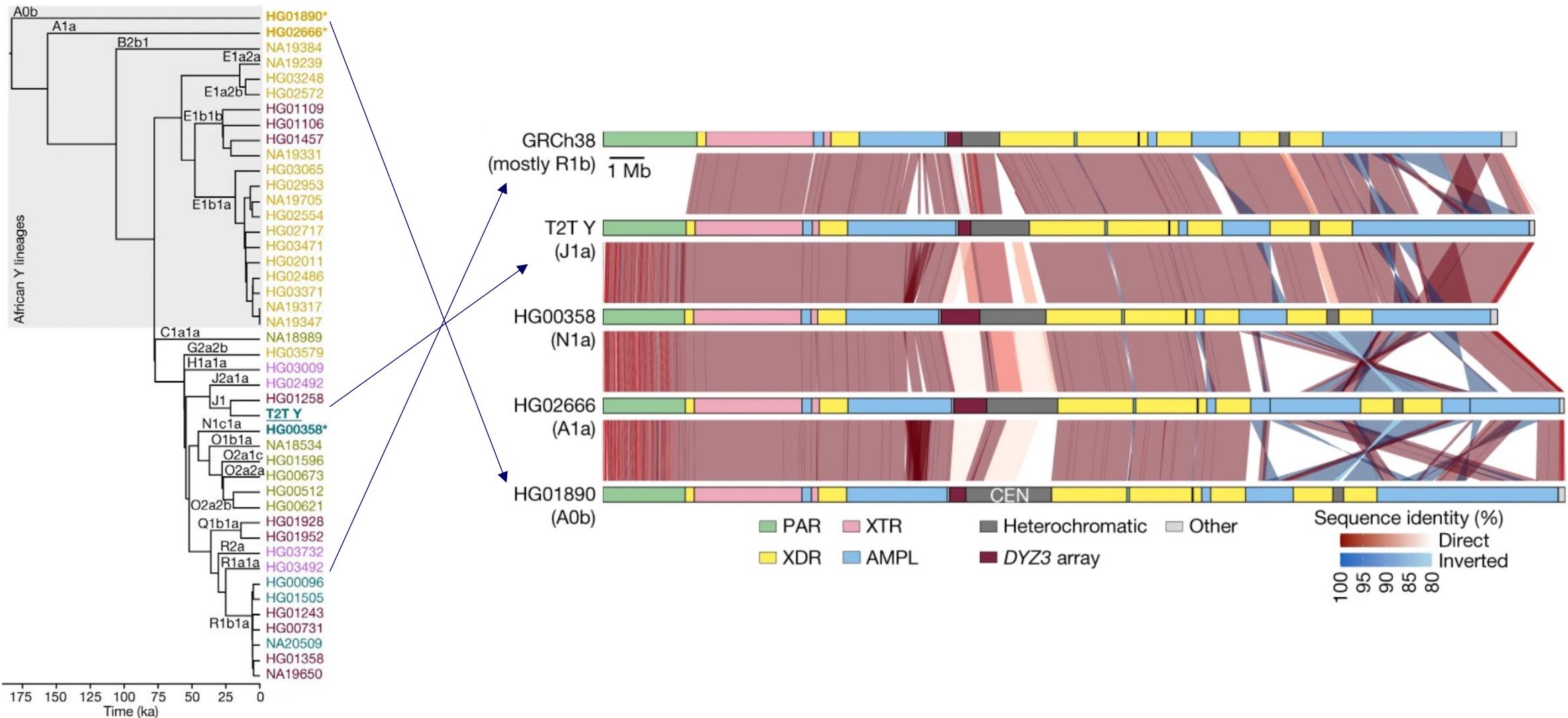
# Human HG002 T2T-Y



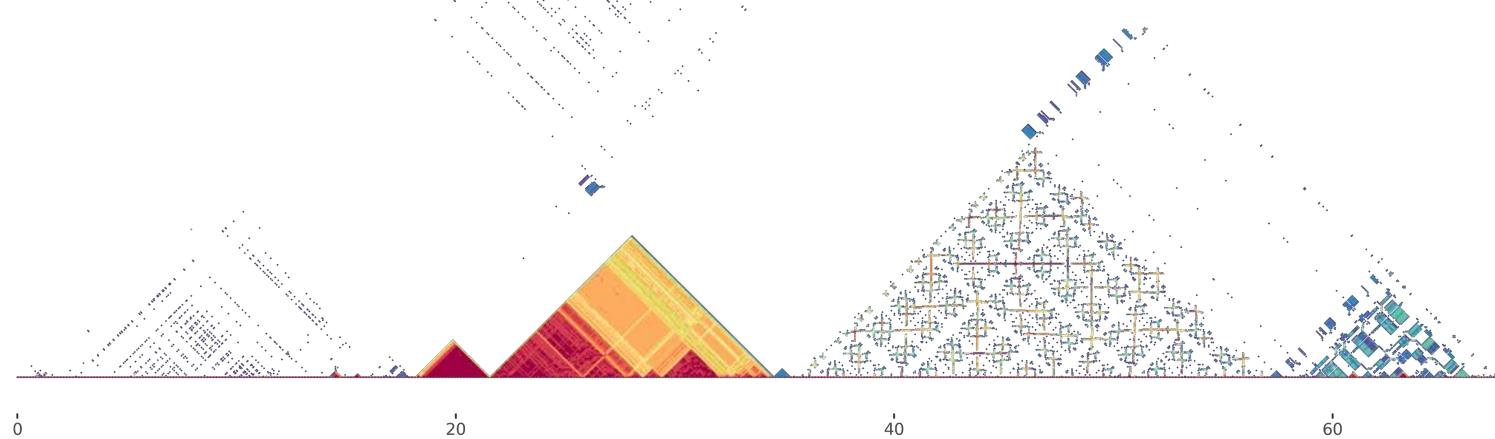
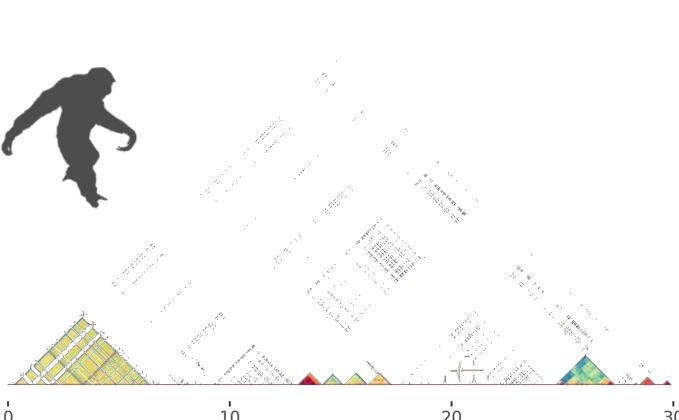
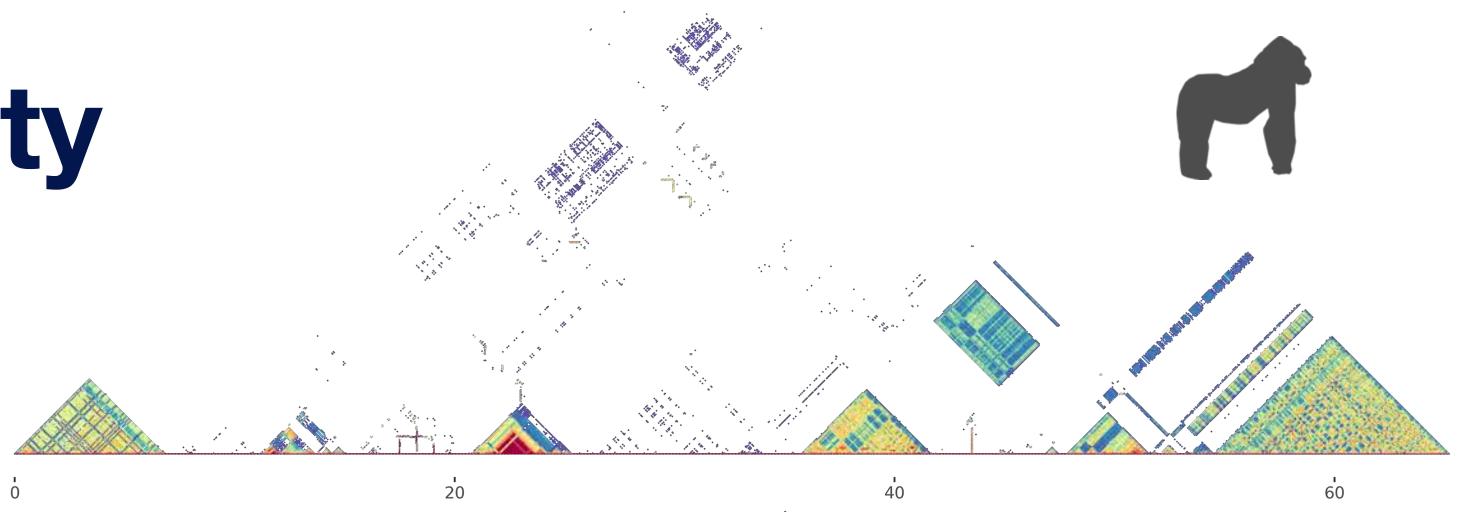
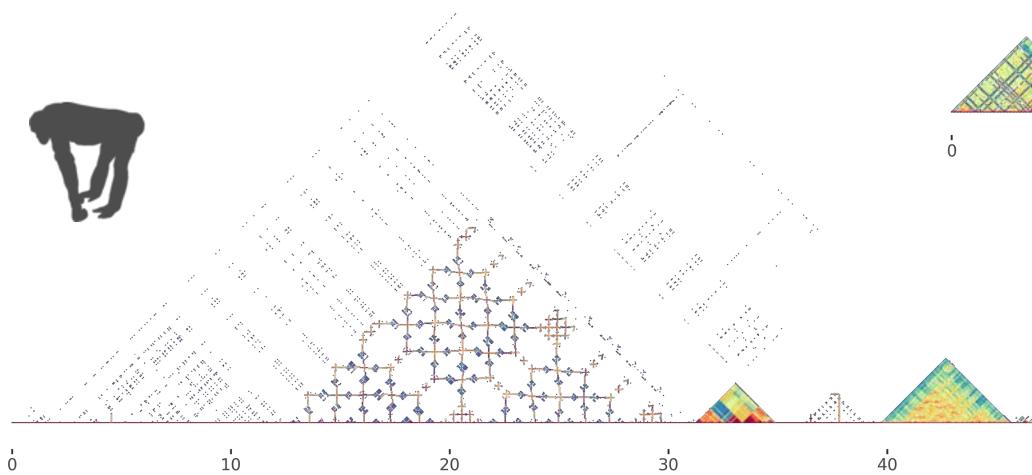
# T2T-Y palindromic region



# Human chrY variation

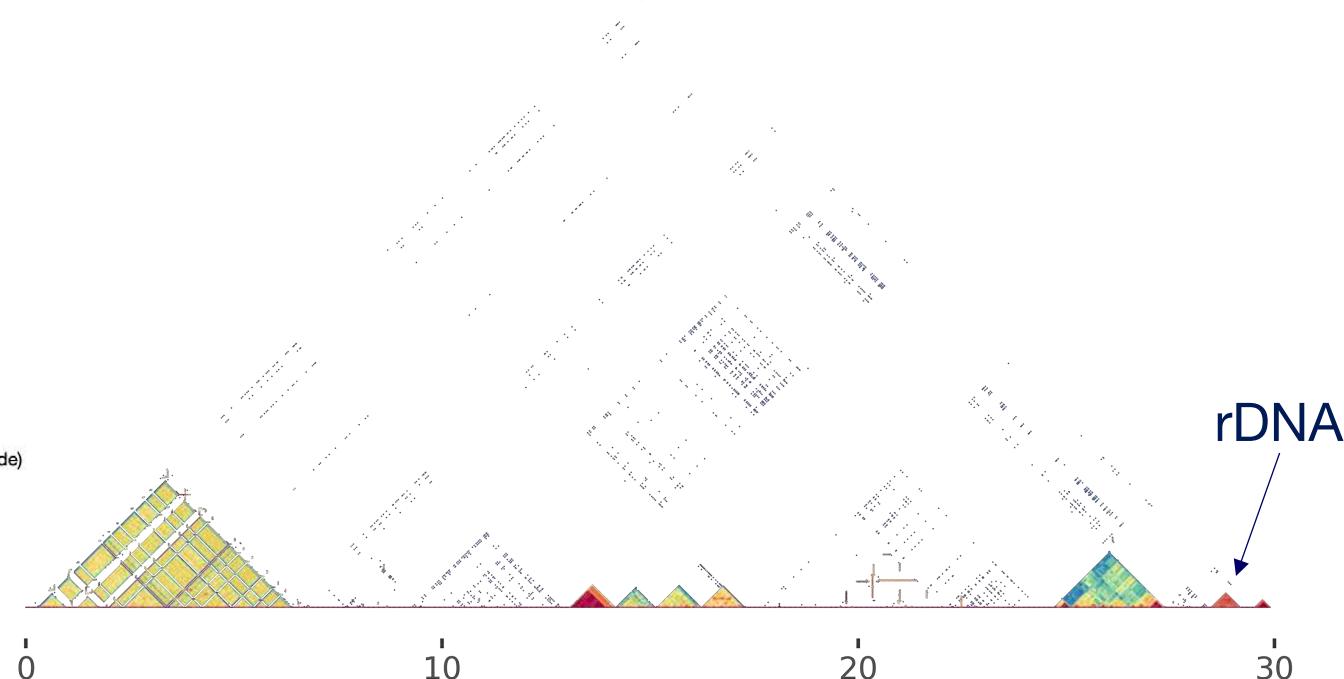
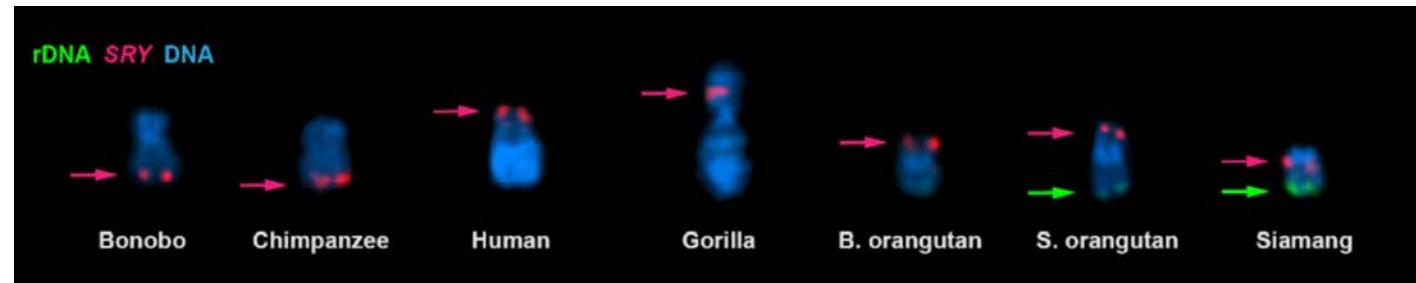
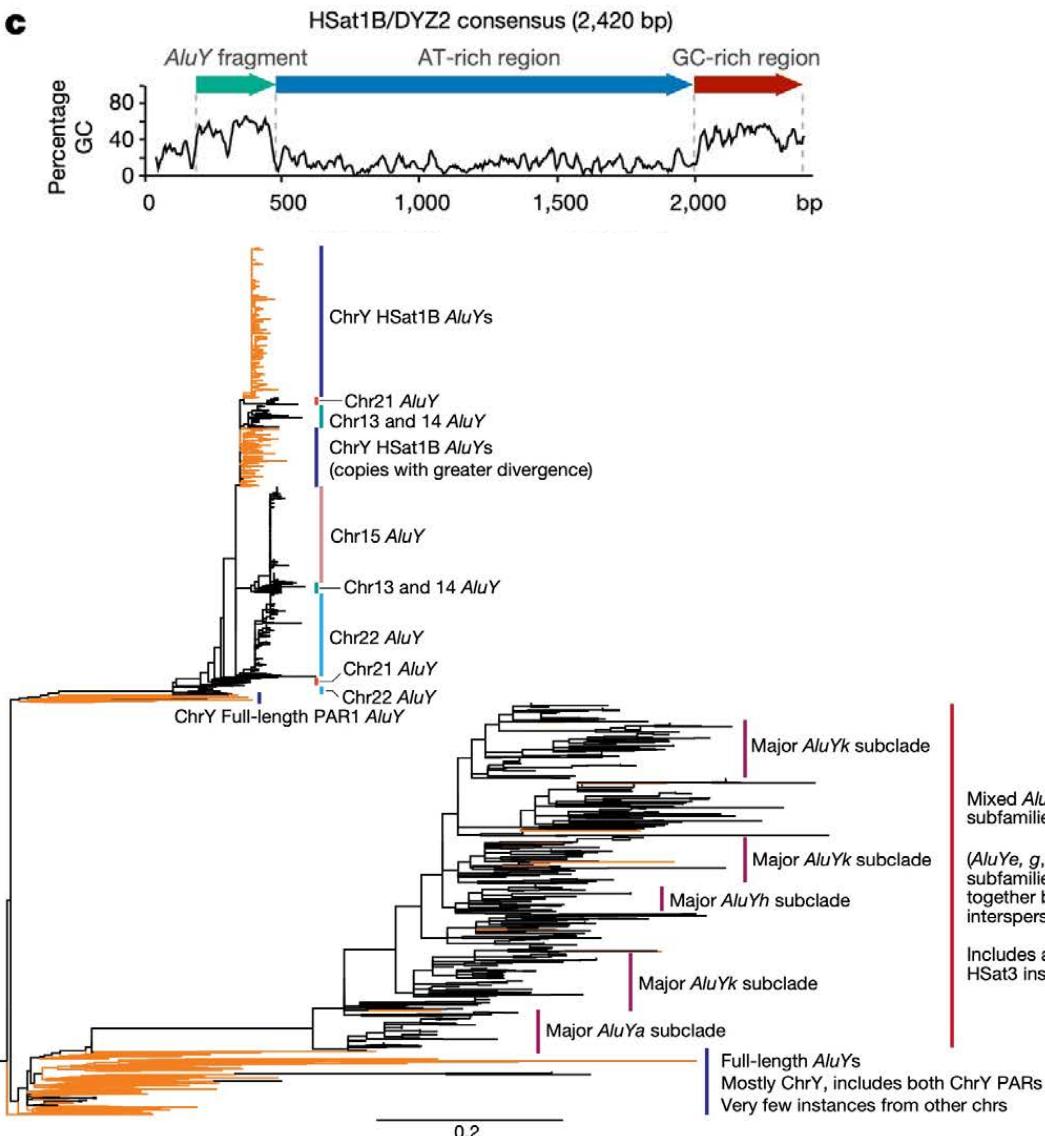


# Ape chrY diversity



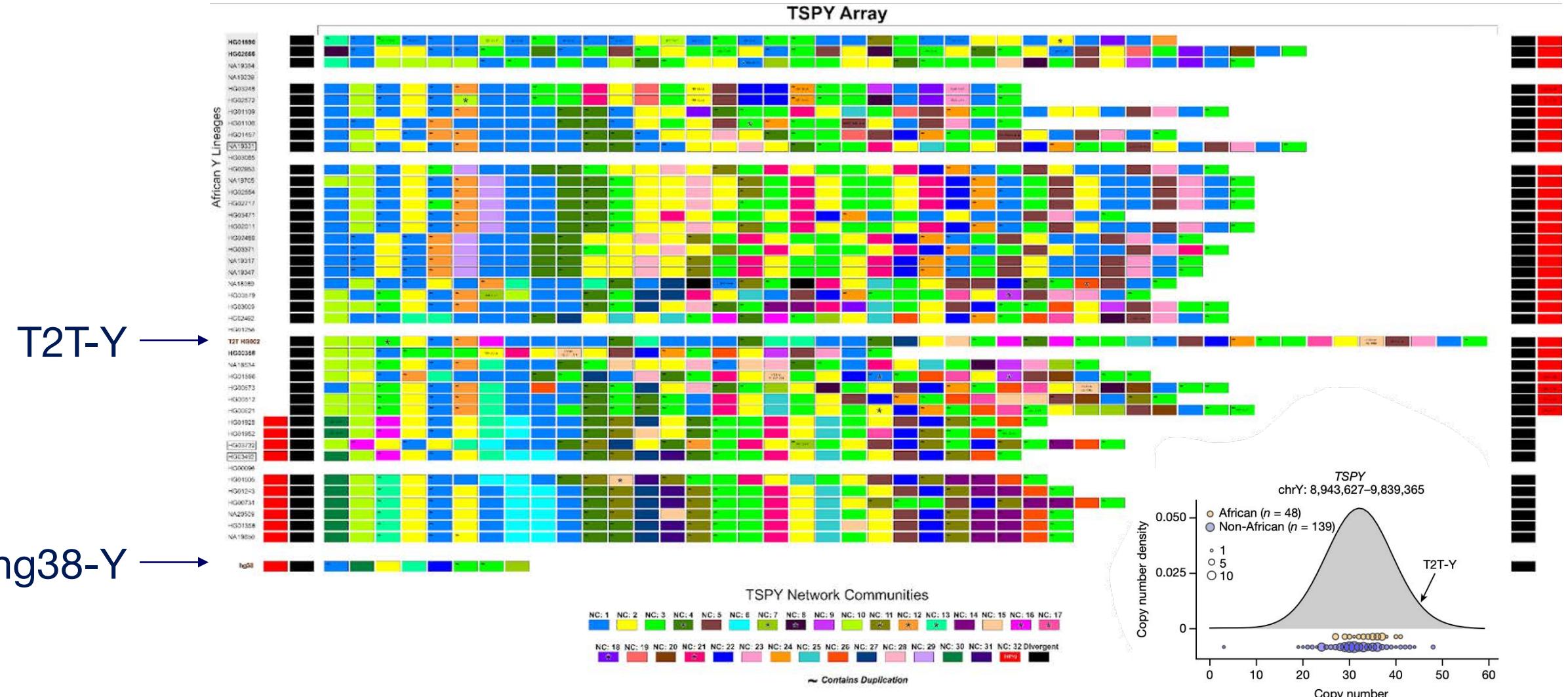
# Crosstalk between acros and chrY

c

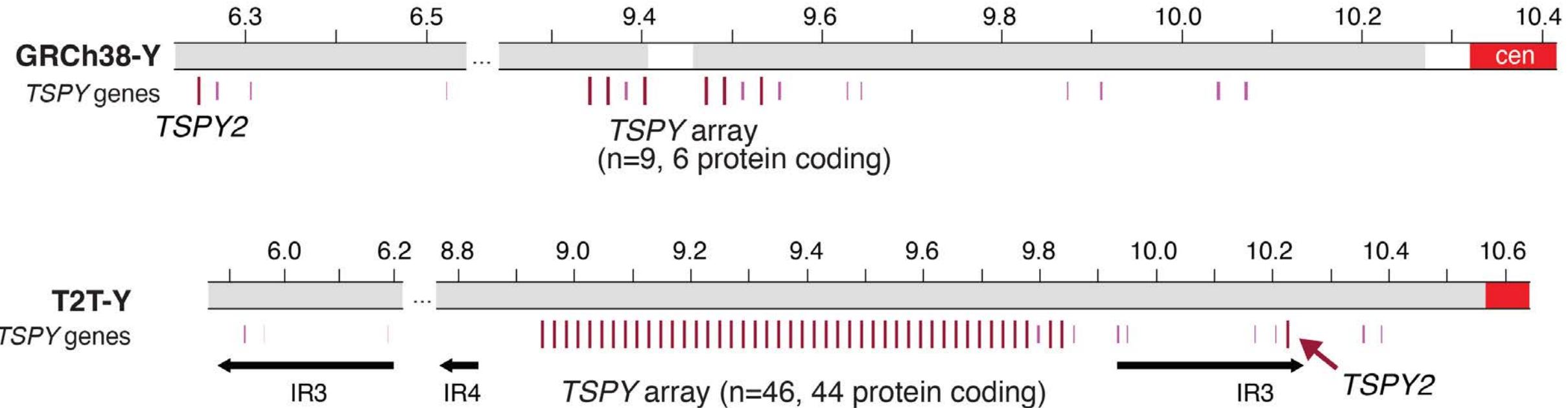


Mapping fails to characterize  
such dynamic regions properly

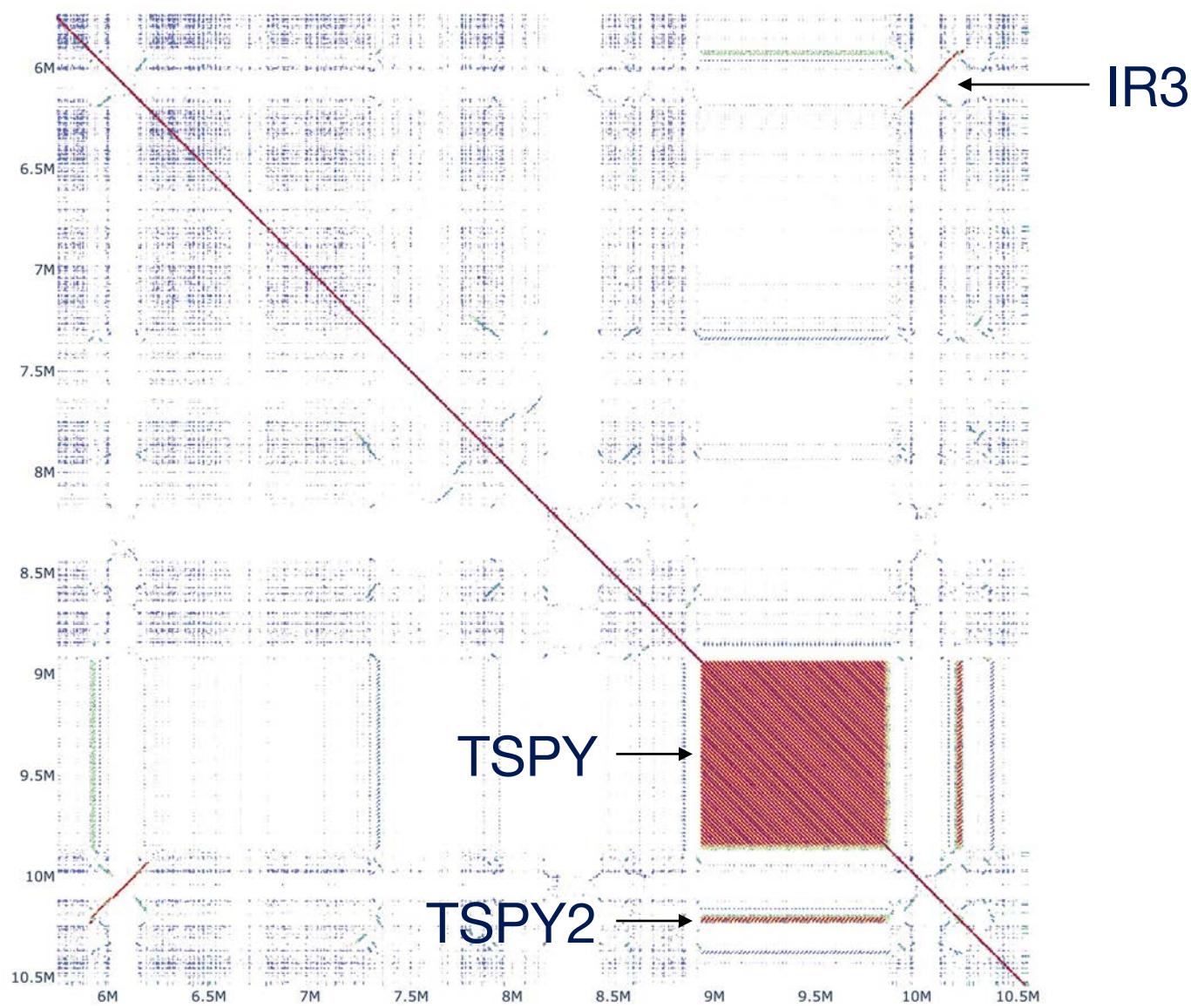
# TSPY gene array



# TSPY2 gene conversion



# TSPY region



# Mapped ONT reads clip at TSPY2

GRCh38 ref



HG002 reads

# Not just a problem on chrY

**Lost in the WASH. The functional human WASH complex I gene is on chromosome 20**

Daniel Cerdán-Vélez, Michael L.Tress

doi: <https://doi.org/10.1101/2023.06.14.544951>

This article is a preprint and has not been certified by peer review [what does this mean?].

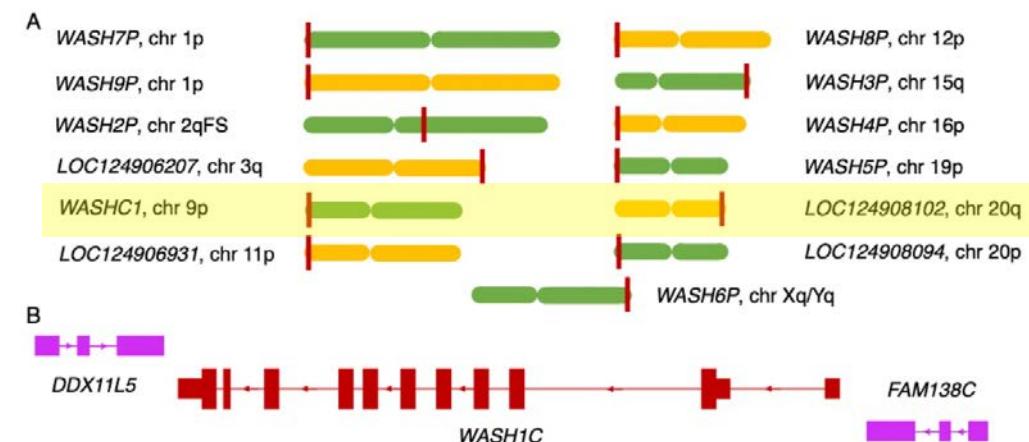
## Abstract

The WASH1 gene produces a protein that forms part of the developmentally important WASH complex. The WASH complex activates the Arp2/3 complex to initiate branched actin networks at the surface of endosomes. As a curiosity, the human reference gene set includes nine WASH1 genes. How many of these are pseudogenes and how many are *bona fide* coding genes is not clear.

Eight of the nine WASH1 genes reside in rearrangement and duplication-prone subtelomeric regions. Many of these subtelomeric regions had gaps in the GRCh38 human genome assembly, but the recently published T2T-CHM13 assembly from the Telomere to Telomere (T2T) Consortium has filled in the gaps. As a result, the T2T Consortium has added four new WASH1 paralogues in previously unannotated subtelomeric regions.

Here we show that one of these four novel WASH1 genes, *LOC124908094*, is the gene most likely to produce the functional WASH1 protein. We also demonstrate that the other twelve WASH1 genes derived from a single *WASH8P* pseudogene on chromosome 12. These 12 genes include *WASHC1*, the gene currently annotated as the functional WASH1 gene.

We propose *LOC124908094* should be annotated as a coding gene and all functional information relating to the *WASHC1* gene on chromosome 9 should be transferred to *LOC124908094*. The remaining WASH1 genes, including *WASHC1*, should be annotated as pseudogenes. This work confirms that the T2T assembly has added at least one functionally relevant coding gene to the human reference set. It remains to be seen whether other important coding genes are missing from the GRCh38 reference assembly.

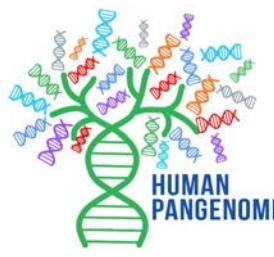


## How frequent are subtelomere translocations and what phenotypes are the mediating?

# The issues with mapping

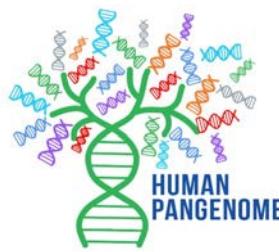
- **Reference bias**
  - Short reads, long repeats
  - A single, incomplete reference
  - Structural variants behave differently
- **Recombination between repeats**
  - Need *de novo* assemblies or specialized variant calling approaches for genotyping these atypical regions
- **Genes are not always where you think**
  - Gene conversion
  - Acrocentric and sex chromosome recombination
  - Subtelomere translocations

# What's the solution?

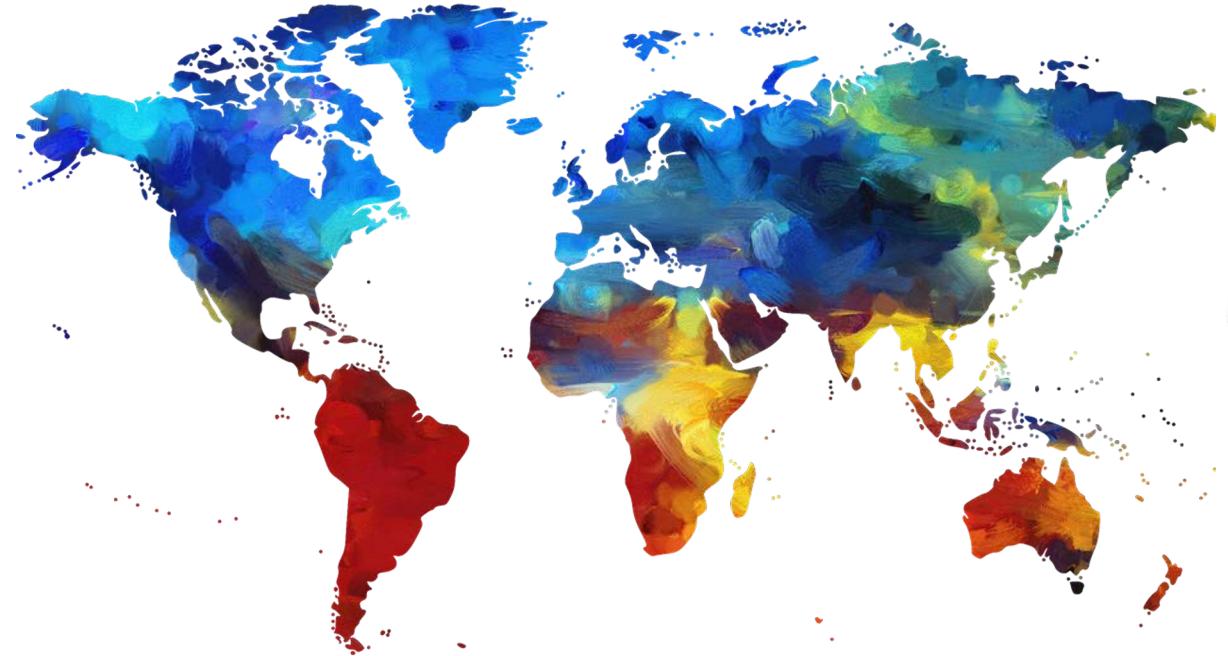


# The human pangenome

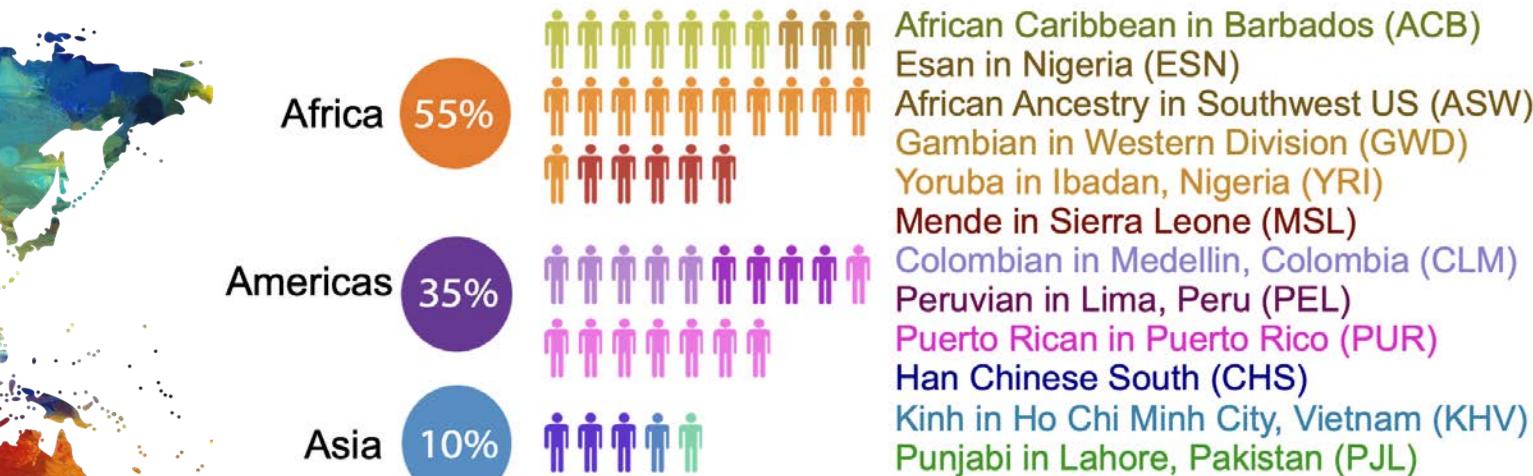
- The pangenome is just a **set of haplotypes**
- Replacing one haplotype with a set of haplotypes **reduces reference bias**
- This yields **immediate improvements** to variant calling accuracy and comprehensiveness
- Will ultimately enable **predictive models** of variation



# A draft human pangenome reference



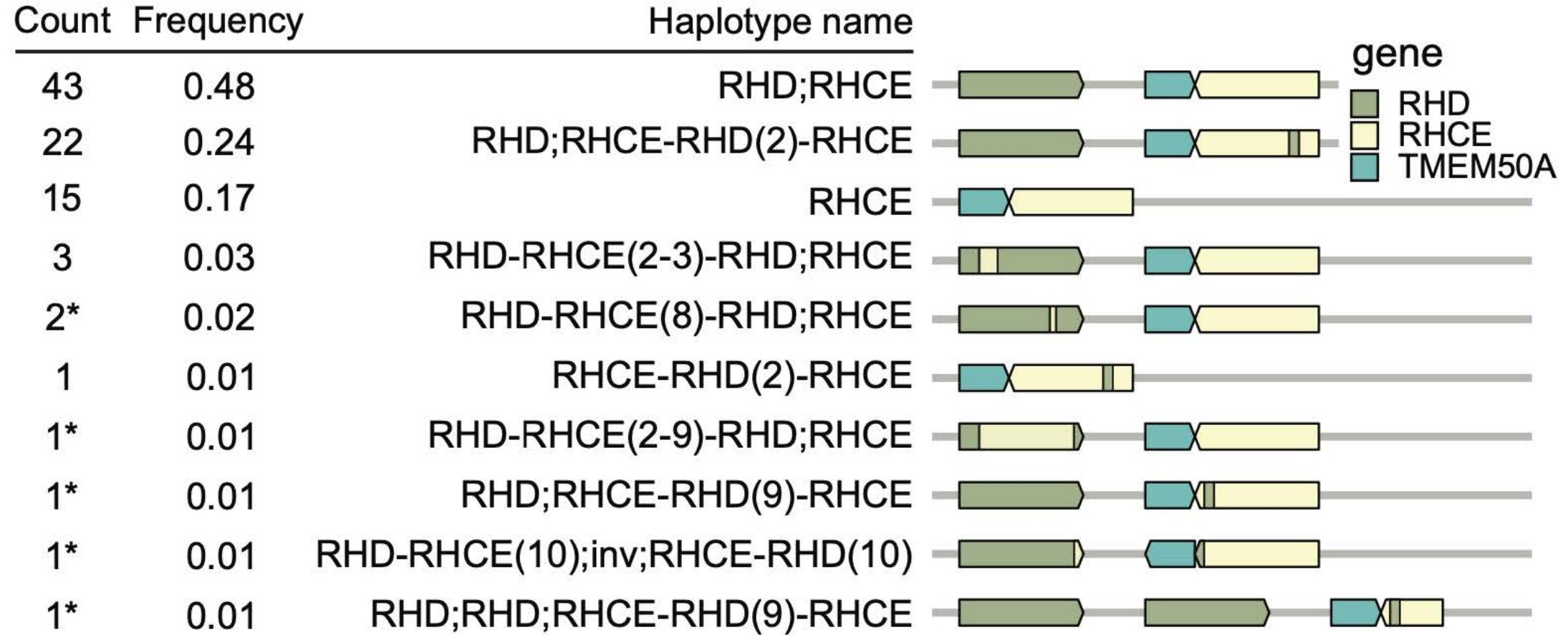
v1: HiFi assemblies (released)  
v2: HiFi+ONT assemblies  
v3: T2T assemblies



48 initial samples  
from 1KGP

- Immortalized 1000G cell lines
- Cover allelic and geographic diversity
- Availability of low passage cell lines
- Availability of trios/parental data
- Aiming for 350+ diploid genomes

# E.g. Rh blood group antigens



# You are your own best reference

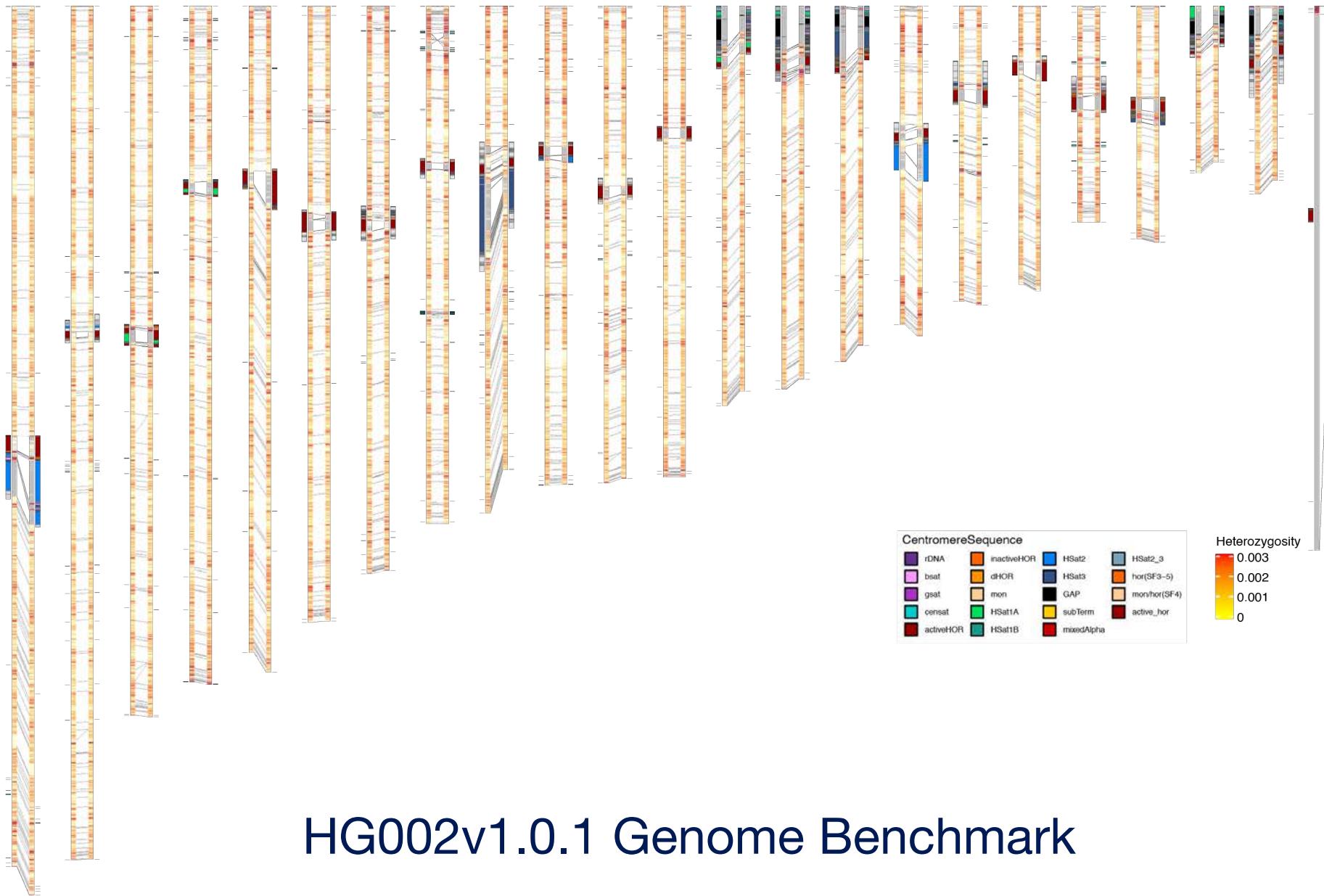
- **Goal:** Complete, diploid genomes for all
  - Variant calls are an incomplete genome assembly
  - *De novo* assembly is costly to scale
- Working towards “genome inference”
  - Use all information available (e.g. reads + pangenome)
  - Infer known alleles, assemble novel alleles
  - Genome assembly with priors

# The “Q100” project



- **NIST/GIAB standards are incomplete**
  - You are only told ~85% of what's in the bottle
- **A genome benchmark**
  - Vs. a variant benchmark
  - Complete, diploid, no errors
  - Q100: 1 error per 10 billion bases
- **T2T-HG002v1.0 released!**
  - Reliable regions approach perfect quality and phasing
  - Unreliable regions (<1%) are labeled (e.g. rDNAs)
  - Exceeds quality of current NIST benchmarks

chr1 chr2 chr3 chr4 chr5 chr6 chr7 chr8 chr9 chr10 chr11 chr12 chr13 chr14 chr15 chr16 chr17 chr18 chr19 chr20 chr21 chr22 chrX/Y



# HG002v1.0.1 Genome Benchmark

—

2000: Draft genomes  
2010: Reference genomes  
2020: T2T genomes  
2030: Genomes

# The future is complete genomes

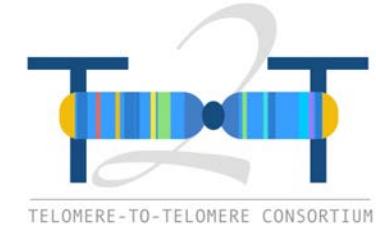
1. Routine assembly of T2T genomes
2. Comprehensive pangenome reference database
3. Cheap and routine *inference* of T2T genomes
4. ML-based annotation of personal, diploid genomes
5. Accurate functional/somatic view over lifespan

**Flip the model:** Don't bring your genome to the annotation, bring the annotation to your genome!

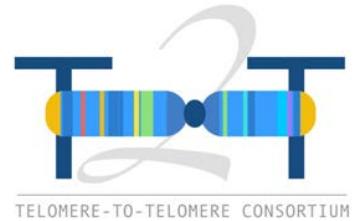
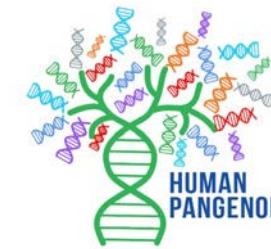
# Resources

# GitHub /MarBL

- **Verkko assembler**
  - <https://github.com/marbl/verkko>
- **T2T-CHM13 reference**
  - <https://github.com/marbl/chm13>
- **HG002 diploid genome benchmark**
  - <https://github.com/marbl/hg002>
- **Human Pangenome Reference Consortium**
  - <https://humanpangenome.org>
- **ModT2T**
  - <https://www.genomeark.org/>



# Team T2T, HPRC, (...and many more)



DNAexus



Google Health