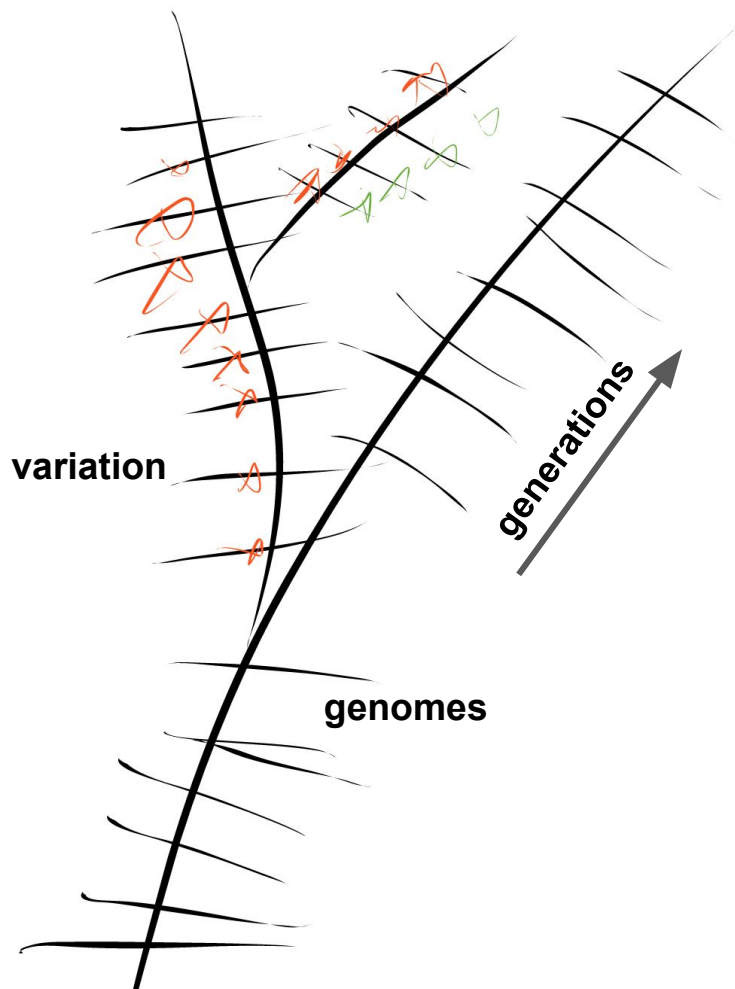


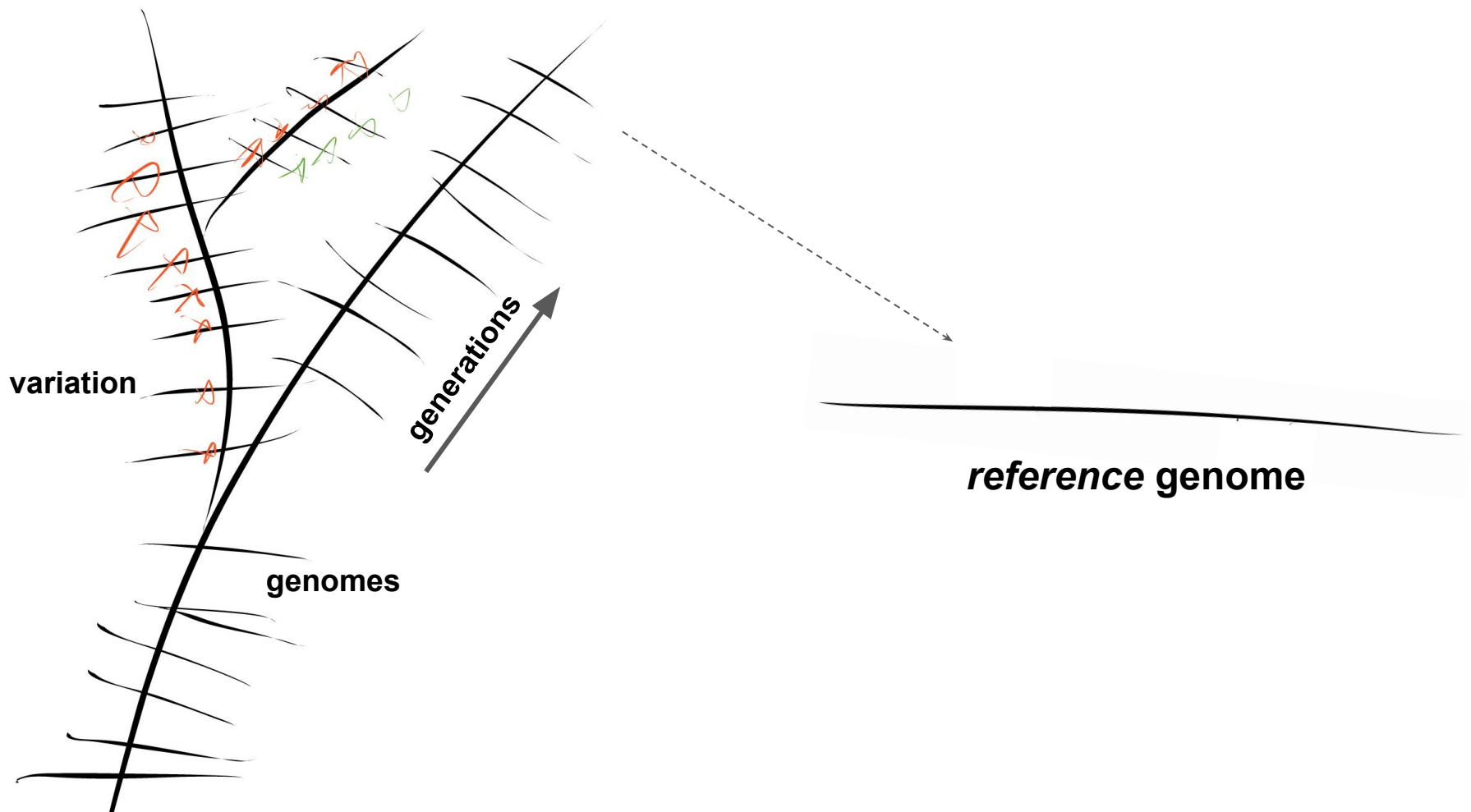
Building, understanding, and using pangenomes

Erik Garrison

University of Tennessee (UTHSC), Memphis

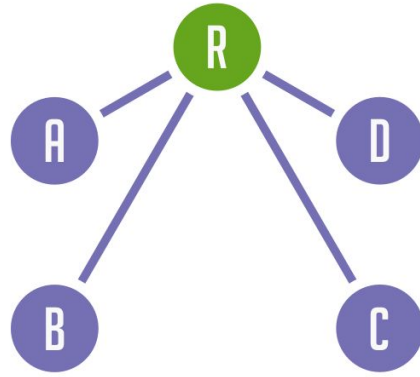
@Workshop on Genomics, Český Krumlov
January 13, 2024





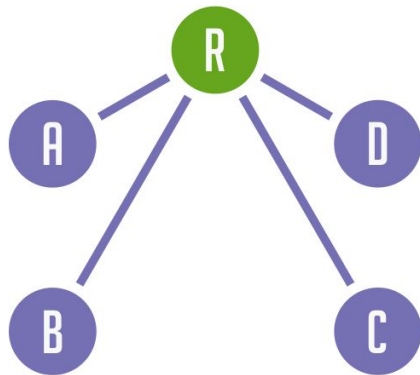
Genomic

Reference model

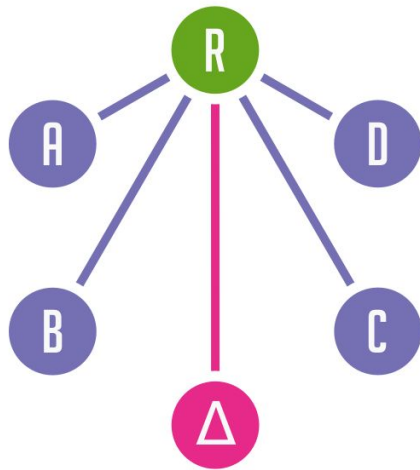


Genomic

Reference model



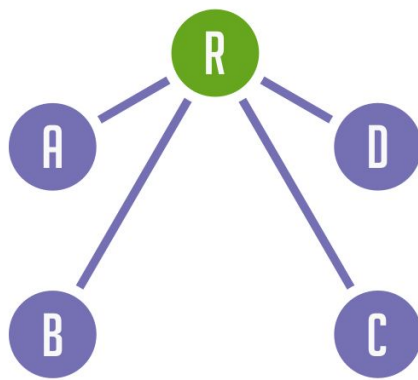
Extending the model



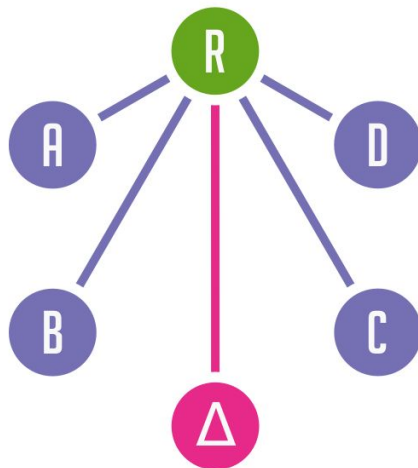
Δ : new genome; R: reference genome.

Figure from [Eizenga et al., 2020](#).

Reference model



Extending the model



alignment and
variant calling

Variation (*VCF*)

Reads
(*FASTQ*)

[illegible]

GAACAACATAAATATGTTATCAATCATCATTTTACTGTGACATAAATAAGTTCTAAATCACTGCAGAGTGTAAATGG
CAATAGACTTCCCACTAAACAAGACCATCTCGAAAAGTTTGGTTCATTTAGAGAAAAAATTTTAAACCTGAGAC
CATGAGATATGATTTTGTAAATAGTTTAAATGAAGAAACACAAATATAAATCTACTATCTTCCGATGATTATA
AAGAGTTCATATGATGATTAATTTTATTTTCCGAATTTGATATCCCAAGTAGTATCACTACAGCTTCTAA
TCTCATCTCAATATCTTCAATTTAAATTAATGACATAATAGTCTGTATATAAAACAACAGCTCTATAGCTCGCTATTC
AGACCAGTAATAAGAGTTTAAAGGCTTGTGATGACCAATGAAGTTTCTTATGGATTTTAAAGAAAAATTTTATAAAAA
TATGTGAGGTATTCAATGAATACATTTTAAATTTGCGAAGCATTTTGCAGAGATCTAGACATGATGAAAGATTTAA
ATTGCGAGCCCTTAAATAGTTTGTAAATAGTTGAGATCAAGATTTAGTACATACAGAGCAAGTAGAGCAAGTTAGG
TTATTAGAGATTTAGTCAGTATCTGTTGTATGTAAGCTCTGAGGAGAGCAAGCTGATCTTTCTGCTACATCTGCG
AAGTAGCTATGTTTGGACAGTAGGAATTTAATACCCCTCTCCCATCTTTCTCTGTTGTTGTCAAACTGTGACAA
CTCTACGTCAGATAGCTCAGGGCAAAAATGATAAAGTTCAAGTTAAGAGAGCTCTGCAAGTGTCTCAAGTCTCTCTCTGG
TGAAGAGGAGAGAGGTTGTGTTTAAATTAATGAATCTGGGATTTCAAAAATCTTACCATGCGCTGCTGCTGCCCTCATTA
GTAGAGGCTGTTATTTAATTAATGTCAGATACCAAGCTTTAGTAGCTCTCTGAGGTAAAGAGATGAAATTAATGTTG
TTTATCTATGTTCTACATTTAGCTAGCTGATTTATTAATTAATCAAGATTTATCTGAAGACTGAGTCACTCTGAGGAA
AAAAAAGAAAAATTTCCCTAAACAAGTTAGAGAGGATGATCGAGACACATCTTATCTCAAGGCTCTCAATCAACAGC
AATGCTTACCAACCTCTATTCAAAATTTTGGCCGAGAGTTCTGATAGACCACAGCAGAGTTGACACATTTAAT
CTACTCTTCTCAGTCTCTGATGATGTTCTCTCCAAATCTACCAAGATCTCAAAATATTTCAGGAATCTTCTCGACAG
AGAAACAGGTTGTGATGATACCACTTTGCTCCAAATCGAGGAGAGATATATATGGAAGATCTATTGATGACCTA
TAATATAGTTTGGAGTTGTTATCACTGAGAGTTTCTCCGAGGATACCAACAGAGCAGATGAGAGCTATTGTCTATT
GTTATGATCTCTTTTACGCGCTGTGAGGGCAGTCTACATCTATTGTTATTAACCTGAGACCAGGAGCCAGTGAAT
GACCTTCAGGCTCTCATTTGTCAAAAAATCAAAATGTGAGGCTTGTGCACTAGAGAAACAGATGTTCAAGAAACCG
TGCCCTGGTTCTGTAAATATCTTGAGATGTTCTGTGCTGAAGTGTGATGAGATGAATAGGATGATTTGAGAGGTGACAG
CTGTCGAGCTCTCATACAGCTCTGCTGCTGCTCAGGCGCTCTGCTGCTGCTCCCACTTTGGCGACATCTGAGAG
CCCTTCAGCCACCATCTGCTACTGTGGAGACCCCTTCTGGGCTGGCCAGGCGCGAGCCGCTCTCCATGTTGAGGGA
GTTGTGGAGGAGAGGCGGACGGACCCGGGGCTGGCGACGCGCTTGGCGGCGCAGCTGGAGTTCTGGGTTGGGCGCTGG
GCTTGGCGGCGCCCGCATCTCGAGACGCGGACGCGCTCTCAGCCCGCAGGCAATGAGAGGCTTGAACCCGGGCGACGA
GCTGGCGGAGGTTGACTCTGCTCCCGCAGCAGCTCGCAGCTCAAGCGGCTGCTGCTCAATTTCTCAGCGGGCTCTGAGCTGCC
TTGGCGGGGAGGTTGCTGGGAGCTGAGCAGCCGATGCTGAGCTCCGCTCCGCTCCGCTCAGGCTCTGAGGCTCTGAGCTGCC
CTCCGCTGAGGAGCTGCGCCGCTCTCAGCGGCGCCAGAGCTGAGGAGCTGAGGAGCTGAGGAGCTGAGGAGCTGAGGAG
CGCGGAGCTGCGAGGAGCTCCGCTCAGCTCTGCTGGGAGTCAATCGGGAGAGCAGCTGGGCTCTGAGGCTTGGT
GGGAGCTGGGAGAACCTTTTATCTAGCTCAGGAGTTGATAATACCACTGGGAGCTGCTGATCTAGCTCAGGTTGTT
AAACACCACTACGACAGCTCTGTTCTAGCTAGTGTGTGGAACGACCAAGCCACATCTGTATCTAGCTACTCTGGT
GGGCTTTTGGAGAACCTTTGTTCTCAGCTCTGTAGCAGTAATCTTGGTGGAGACATGGAGAACCTTTGGGTTGAGTCT
AGGATTTTGAAGCACTGAGTCAAGAGGCTCTGCTGCTACCAATGAGAGTCTGCTCTACCACTAGCAGGATGAGTGTGGTGGGAG
AGTAAAGAGCAATAAAGCAGGCTGCTGAGCAGCAGTGGGCAACCCGCTGGGTTCCCTTCCACATCTGGAGAGGTTGTT
TCTTTCTACTGTTTGAATAAATCTTGCTGCTGCTCACTCTTGGGGTCCACATGCTCTTTTGAAGCTGACACACTGCTG
GAAGTCTGCAAGCTTCACTCTGAAGCCAGGAGCAGGAGCCAGCAGGAGGACCAACAACTCAGAGAGCCGCGCT
TAAAGAGCTGACAACTCAGCTGGAAGCTCTGACGTTCACTCTGAGGACGAGAGCAGGAGCAAGAACTCAGAGAGAGAGAA
ACTCTGGAACATCTCGAGATCAGAGAGCACTCAGTCTGAGCGCCCAATTAAGAGTGTAACTACCTCAGCGAGAGGTC
CTGGGCTCATCTTGAAGTCAAGTGAAGCAGAGCAAGCCCACTTTTGGACAGATTTGGAATAAATTTAACAATCAAT
ATCTCTAAGGAATCAAACTACAGATTAATAATAGTAACAGGTCAGTAAGTAATTAATAAAGACATGATGACCA
AGATGGGCAAGAGTTTGTGTTGTCAGAAACCAAGTTTGGCTAAGTACATGACCGAGAAATTTAATCTGAGGATGAT
TATGTTGGGAAAGCAAGTGTGTTTATCTGTTGATTTAATAAATGATGACCGAGAAATTTAATCTGAGAGATAC
TCTAGTAAGCTCTGAGTCAAGCAATTTATTAAGTCTGATTTAGTACGATGCTGTTTGGTGAAGAGTCTGAGTCTG
TGAAGACATCCCAAGTGTCTATGAGAGGCTGGAATGATTGGAAGGAGTGTGGAATGATGATTAATGAGTCACTCA
GCTGTTTGAATAAAGAGCTCCAATTAACATACCAAAAGATACATTTAAACAAAGAGCTTGAAGAGAAAGAGTCC
ATCAAAACAATACACAACTATCTACACATTTTATAAATCTCAAAATAGCAATTAATAACATGATTTTAAAGGAG
CAAGCAATGTAGCAATGATTTTAAAGGCTGAGGAGTGAATGAATGAAGAGTCAAGAGTCAAGTCTGAGGAGCAG
GAAGAGGAGAGTGCAGTGTGGGAGAGGATACATAGTACAGCAGTAGAGAGGATCTTTTCTTCTTCTTCTTCTTCT
CTTATTTCTTCTAGTCTGATTTCTTTAGCTATGAGTATCTTTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCTTCT
ATCATACAGTATGTTTGGTGTGATGAAGTCTCTTAAATAGCTTTTATTTATTTTATTTTATTTTATTTTAAAGTTCAAG
GGTACATGTCGAGATTTGTTATATAGTAAATCTGTGTCATGTGTTTGGTGAACAATTTATCTTCAACCAAGGATTT
AAGATAGTACCACTAATGATTTGTTCTGTCATGCTCTCTCCCATCTCAACTCTTCACTCTCAATAGGCCCGAC
TGTGTGCTGTTCCCTCATGTGTGTCATGTGTTCTCATCACTAGCTCAGTCACTCAATAGTGAAGACATGCAATGTGGG
TTTCTGTTCTGCTGATTTAGCTCAGGATATTTGCTCAGGCTCATTTCTGCAAGACATATTTGCTTCTTCTTCTTCT
TATGTTGTCATGATTTGTTGGTGTGATGATGACCCGATTTCTTATGACGAGTCTATCATTTGATGAGCATTAGGTTGAT
TCCATGCTCTGCTCTAATCATTTTAAACAGTCTCTGAGTGAATAGGGGAGGCTGGTGTAAAGAAATGATCTGTTTAA
TGGAGAACATACATACATGATTAATAAAGCATATATATATAGACACATATATGATTTATCATATTCATATTAATCA
CAAGCTCTGAGAAATTTCCGCAAAAAGCATGTGGAGGAGAGAGAGGGGTTATGTTGATGATTTGTAGGCTCGGA
ATAAGATGATCTCATATTTGCTTCATTTTGTATAAATTTTTTTAAAGCATATGATGATCTAGAGTATGAGAGAA
ATTCTGGGAAGATTTGGAGAACTGTGGGAGTGTGAGAGTGGTCTACCTTTGCTTGTGATTTATGTTATTTT
TCATTTAATGTACTTTGTAATCAAAATTTATTTGATATTTTACTTCATTTTAAATGGCATACGAGCTTTAATTT
TATACTAAATAGCTTTGCTCCAGCGGAGAGGAGTGGGAGAGAGATGAGAGAGGCGGATCTCAAGACAGCAGC
ATTCTGCTCTGCTCTGGAAGCCAGGACGACATCACTCTCCGCGAGGCGGACGACGCTGCGCTGCGAGCGGCG
GAGGCGCGGAGCGGAGCGGAGCGGAGCGGAGCGGAGCGGAGCGGAGCGGAGCGGAGCGGAGCGGAGCGGAGCGG
GAGGAGCAGCAGCCGCGCAGCGCAGCTGGTGGCTGGCAGGCGGCGGCTGCTGGCGGAGGAGTCAATCAATCCCG
GGCTCTCTGCTCTCAACCCGCGCCCTGGCTCAGCTCCGCGCGCGCTGCTCTGGGAGCGCGGGGCGGAGGAGC
ATACACCCAGCTCTGCTGCTGCGCGCGCAGCTCAGCAGCCCTACCCAGGAGCATATCACTGCGTGGGCGCTGCGC
GGGCTCGGGAGCGCAAGGCTGCGCGCTGGGCGAGCGCTGAGCTCAGAGAGCAGAGAGCGGGCTCTCCGCGTGGCCCA



We cannot update
a linear reference
sequence

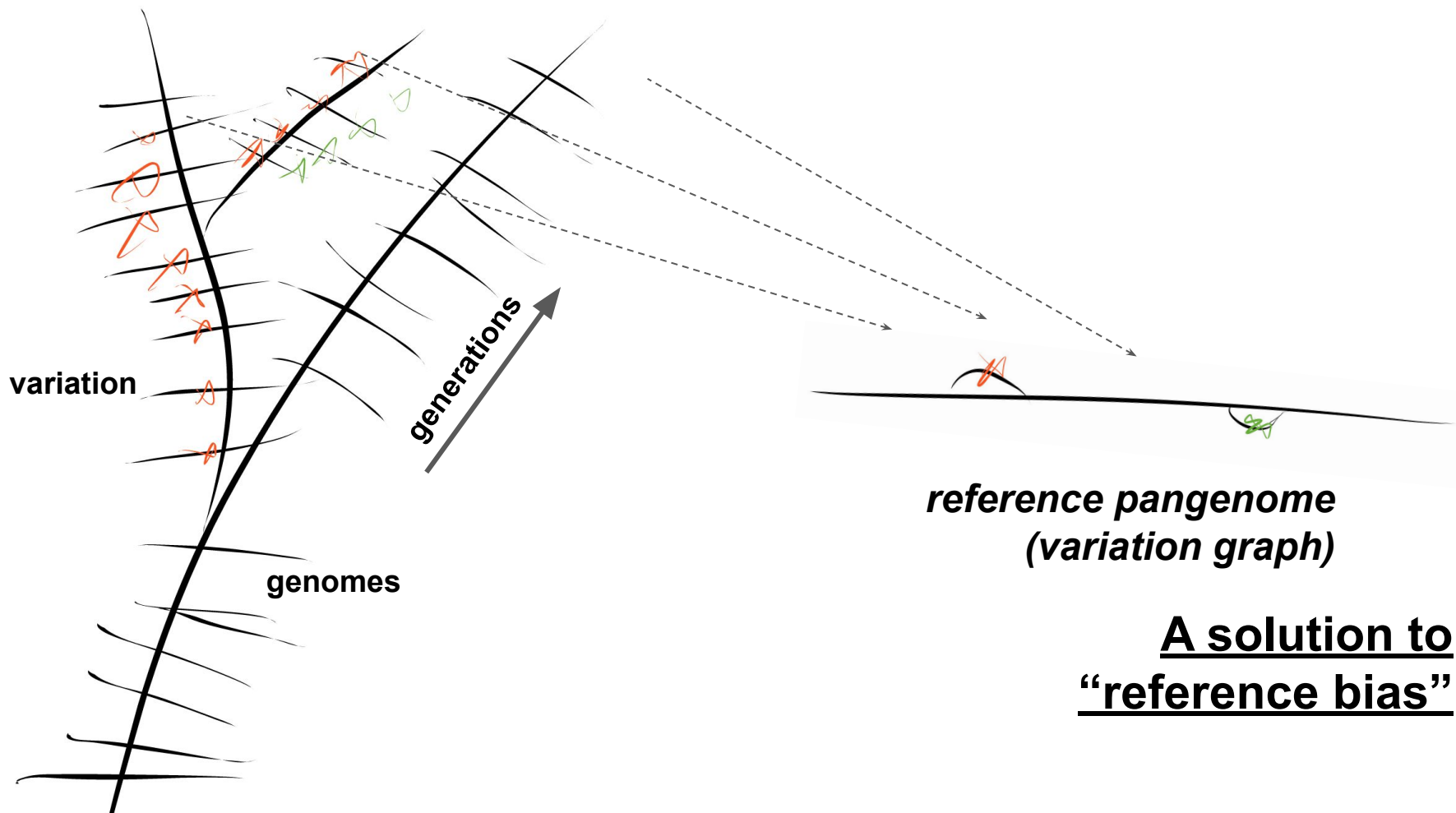
```
##source=mutatrix population genome simulator
##seed=1373927256
##reference=chr0.fa
##hasindex=true
##command=linearmatrix -S sample -p 2 -n 100 chr0.fa
##filters="AC > 0"

##INFO=ID=TYPE,Number=A,Type=String,Description="Type of each allele (snp, ins, del, mnp, complex)">
##INFO=ID=NA,Number=1,Type=Integer,Description="Number of alternate alleles">
##INFO=ID=LEN,Number=A,Type=Integer,Description="Length of each alternate allele">
##INFO=ID=MT,COSAT,Number=0,Type=Flag,Description="Generated at a sequence repeat loci">
##FORMAT=ID=GT,Number=1,Type=String,Description="Genotype">
##INFO=ID=AF,Number=A,Type=Float,Description="Total number of alternate alleles in called genotypes">
##INFO=ID=AF,Number=A,Type=Float,Description="Estimated allele frequency in the range [0,1]">
##INFO=ID=NS,Number=1,Type=Integer,Description="Number of samples with data">
##INFO=ID=AN,Number=1,Type=Integer,Description="Total number of alleles in called genotypes">
```

CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	sample001	sample002			
chr0	1252	-	C	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	1/0	0/0	0/0	1/0
chr0	3646	-	T	TC	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=ins	GT	0/0	0/1	0/0	0/1	1/0
chr0	6283	-	C	T	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	0/1	0/0	0/1	1/0
chr0	7412	-	C	T	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	1/0	0/0	0/1	1/0
chr0	7935	-	C	T	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	1/0	0/0	0/1	1/0
chr0	8131	-	T	C	99	-	AC=2;AF=0.5;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/1	0/0	0/1	1/0
chr0	8682	-	AA	TG	99	-	AC=1;AF=0.25;AN=4;LEN=2;NA=1;NS=2;TYPE=mnp	GT	0/0	0/1	0/0	0/1	1/0
chr0	10926	-	C	T	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	0/1	0/0	0/1	1/0
chr0	11921	-	G	GTT	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=ins	GT	0/1	0/0	0/0	0/1	1/0
chr0	12955	-	T	C	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	1/0	0/0	0/1	1/0
chr0	13808	-	T	TG	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=ins	GT	1/0	0/0	0/0	0/1	1/0
chr0	15271	-	A	G	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/0	0/0	0/1	1/0
chr0	15487	-	A	C	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	1/0	0/0	0/0	0/1	1/0
chr0	16486	-	C	G	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	1/0	0/0	0/0	0/1	1/0
chr0	16561	-	T	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/0	0/0	0/1	1/0
chr0	16748	-	GTT	G	99	-	AC=1;AF=0.25;AN=4;LEN=2;NA=2;NS=2;TYPE=del	GT	0/0	0/1	0/0	0/1	1/0
chr0	17697	-	G	C	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	0/1	0/0	0/1	1/0
chr0	19568	-	C	G	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/0	0/0	0/1	1/0
chr0	20750	-	T	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	1/0	0/0	0/1	1/0
chr0	21532	-	C	C	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	1/0	0/0	0/0	0/1	1/0
chr0	22291	-	C	T	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	0/1	0/0	0/1	1/0
chr0	23193	-	C	T	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	0/1	0/0	0/1	1/0
chr0	23954	-	CTAA	TTAA	99	-	AC=1;AF=0.25;AN=4;LEN=4;NA=2;NS=2;TYPE=ins	GT	0/0	0/1	0/0	0/1	1/0
chr0	24467	-	C	T	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	0/1	0/0	0/1	1/0
chr0	26100	-	G	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/0	0/0	0/1	1/0
chr0	29654	-	T	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	1/0	0/0	0/0	0/1	1/0
chr0	30670	-	T	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/0	0/0	0/1	1/0
chr0	31790	-	A	G	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/0	0/0	0/1	1/0
chr0	32792	-	T	C	99	-	AC=3;AF=0.75;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	1/1	0/1	0/0	0/1	1/0
chr0	33376	-	CC	C	99	-	AC=2;AF=0.5;AN=4;LEN=1;NA=1;NS=2;TYPE=del	GT	1/0	0/1	0/0	0/1	1/0
chr0	33483	-	C	T	99	-	AC=2;AF=0.5;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/1	0/0	0/1	1/0
chr0	33802	-	A	G	99	-	AC=2;AF=0.5;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	1/0	0/0	0/1	1/0
chr0	34450	-	C	T	99	-	AC=4;AF=1;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	1/1	1/1	0/0	0/1	1/0
chr0	34716	-	G	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	1/0	0/0	0/1	1/0
chr0	35484	-	G	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	1/0	0/0	0/1	1/0
chr0	36547	-	G	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	0/1	0/0	0/1	1/0
chr0	38015	-	T	A	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	1/0	0/0	0/1	1/0
chr0	38281	-	T	C	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/0	0/0	0/1	1/0
chr0	38467	-	T	C	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/0	0/0	0/1	1/0
chr0	40581	-	A	G	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	0/1	0/0	0/1	1/0
chr0	40601	-	A	T	99	-	AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/0	0/1	0/0	0/1	1/0
chr0	41968	-	G	A	99	-	AF=1;AF=0.75;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	0/1	0/1	0/0	0/1	1/0

Genome (FASTA)

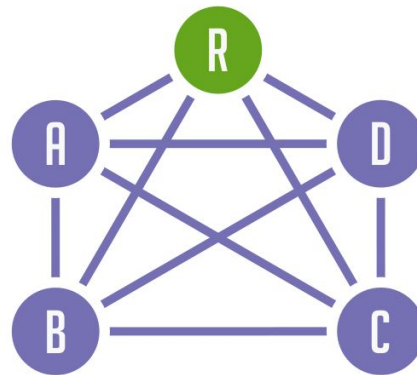
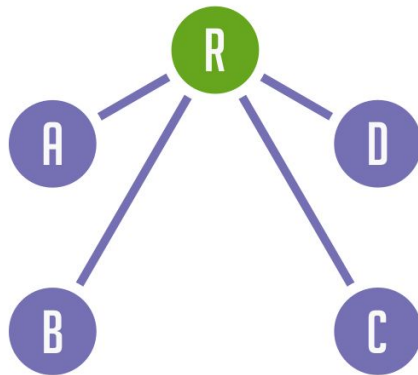
Variation (VCF)



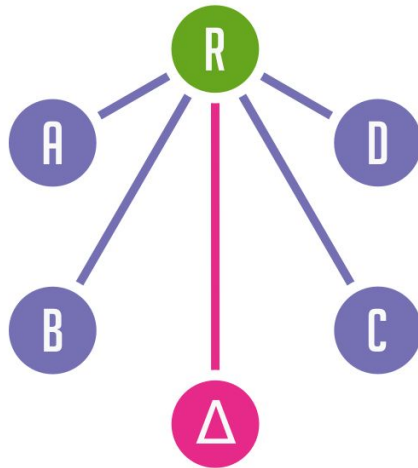
Genomic

Pangenomic

Reference model



Extending the model



Δ : new genome;

R: reference genome.

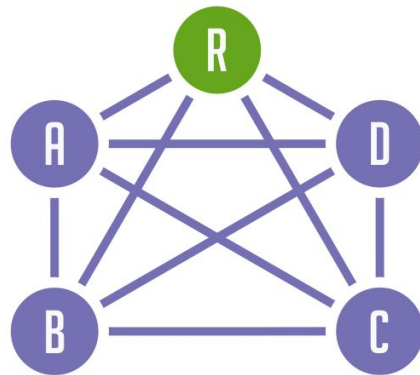
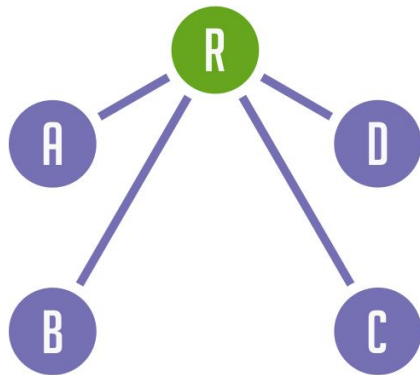
Figure from

[Eizenga et al., 2020.](#)

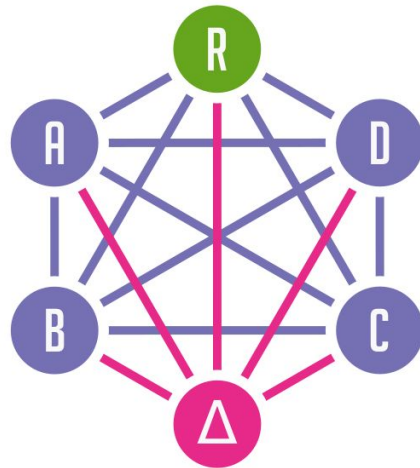
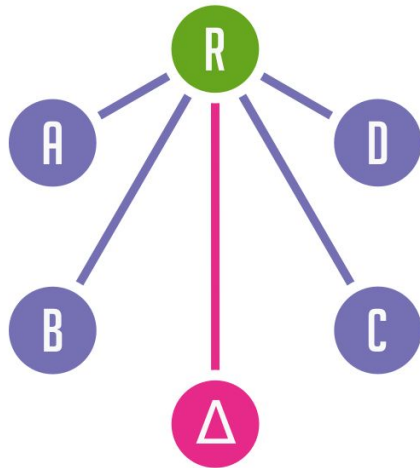
Genomic

Pangenomic

Reference model



Extending the model

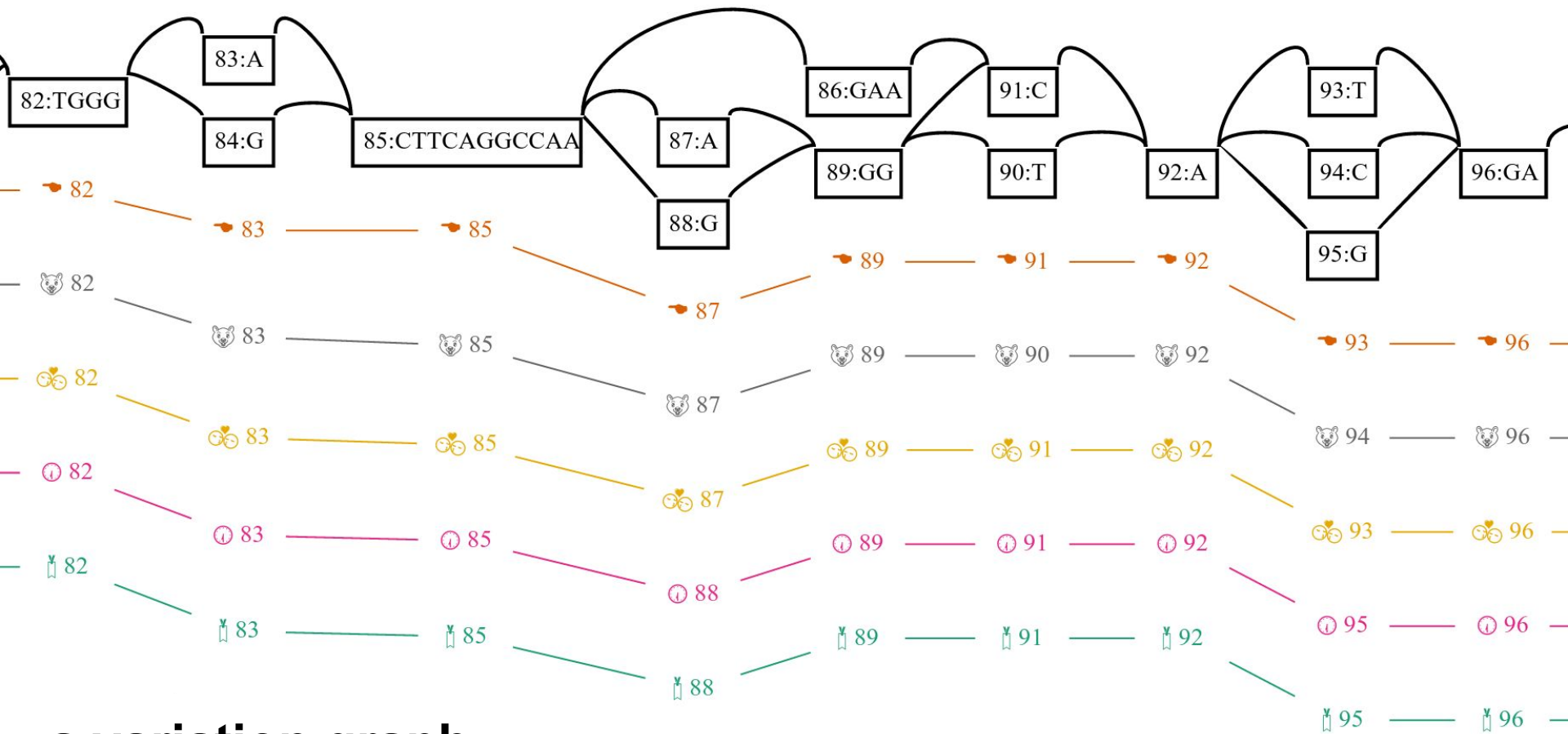


Δ : new genome;

R: reference genome.

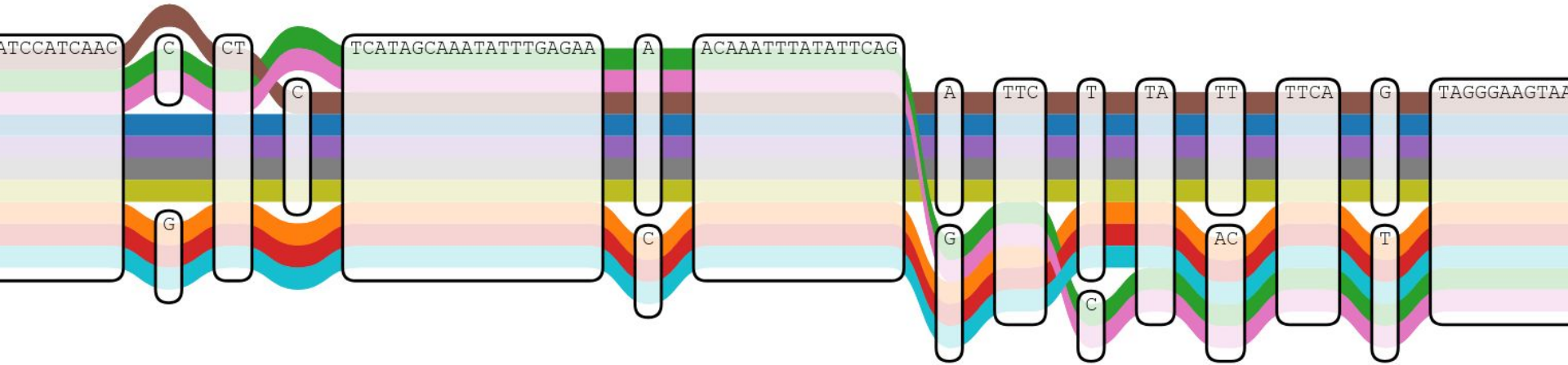
Figure from

[Eizenga et al., 2020.](#)

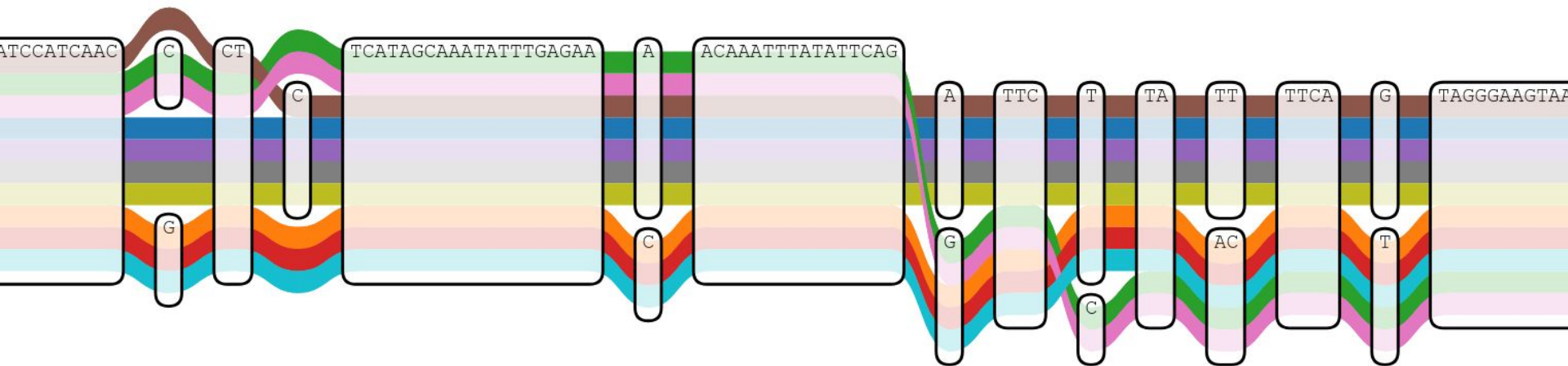


a variation graph

Variation graphs answer a key problem in bioinformatics:

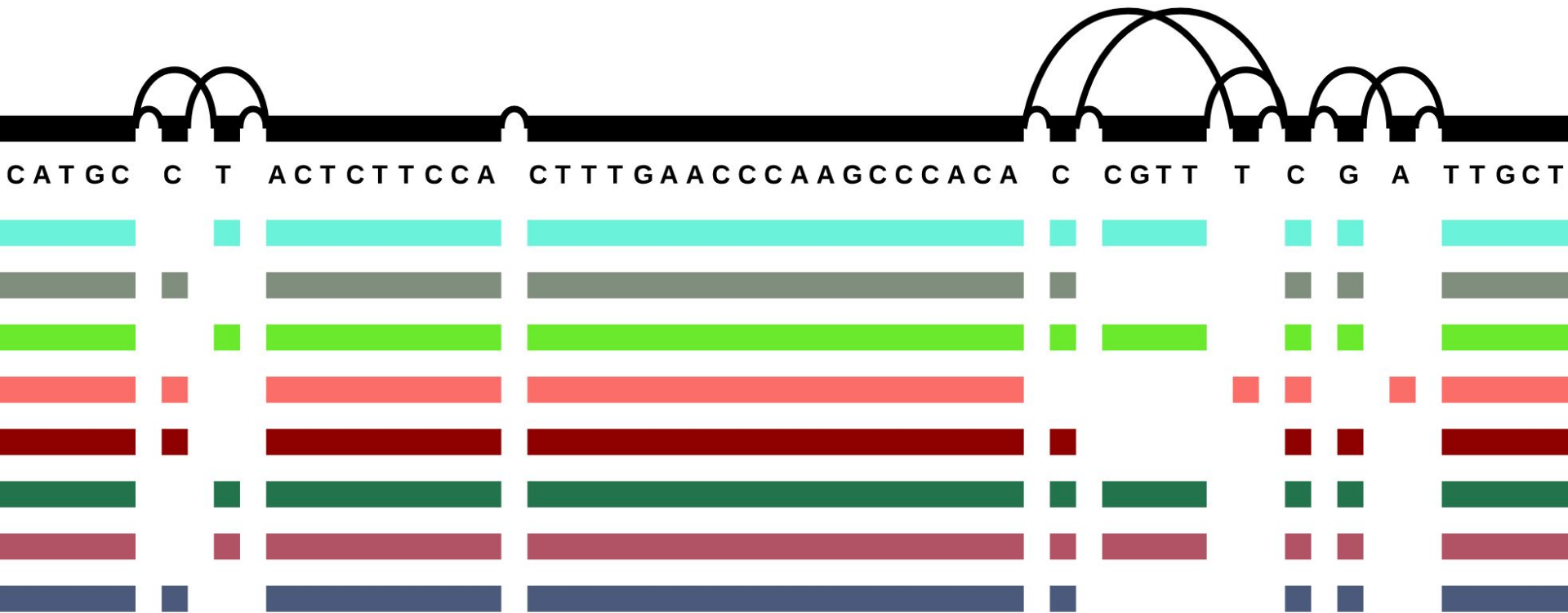


Variation graphs answer a key problem in bioinformatics:

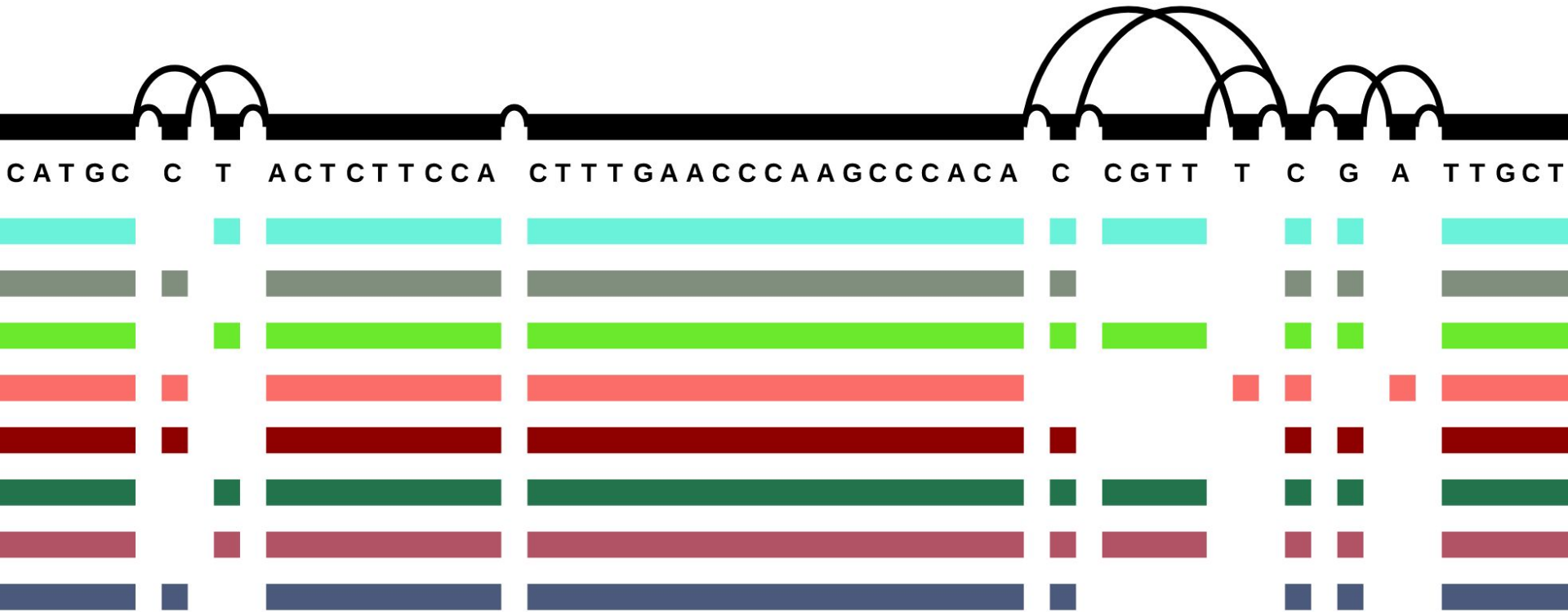


How to represent both sequences and any kind of variation between them.

variation graphs are pangenome models



variation graphs are pangenome models



... which give us a simple way to project many genomes into vector spaces.

New ideas often have a long history

This all seems cool and “new” but ideas are rarely that.

Pangenomes and variation graphs have a long[†] history.

([†]for genomics)

St. Agatha pipetting a
biosample into a
nanopore sequencer
c. 1420

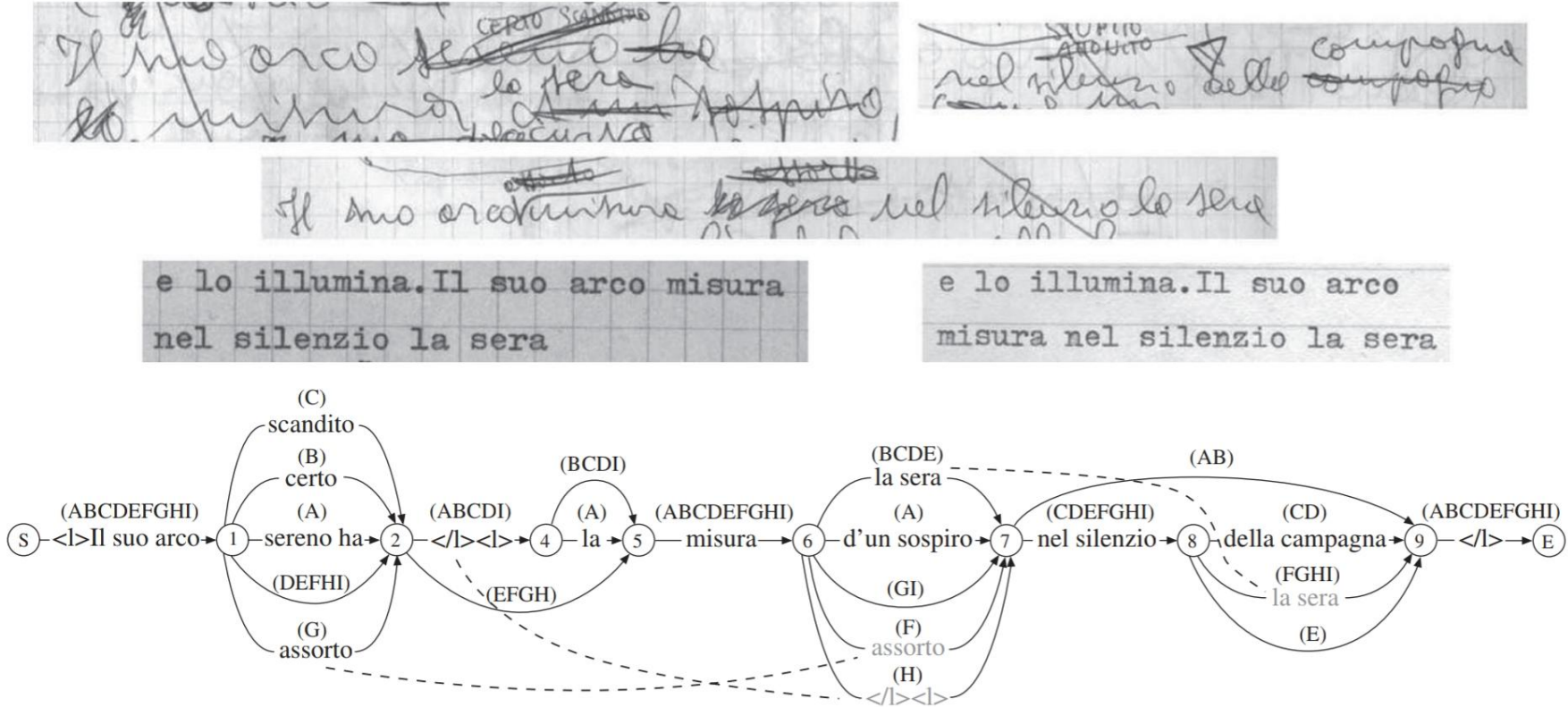


pipette

biosample

nanopore
sequencer

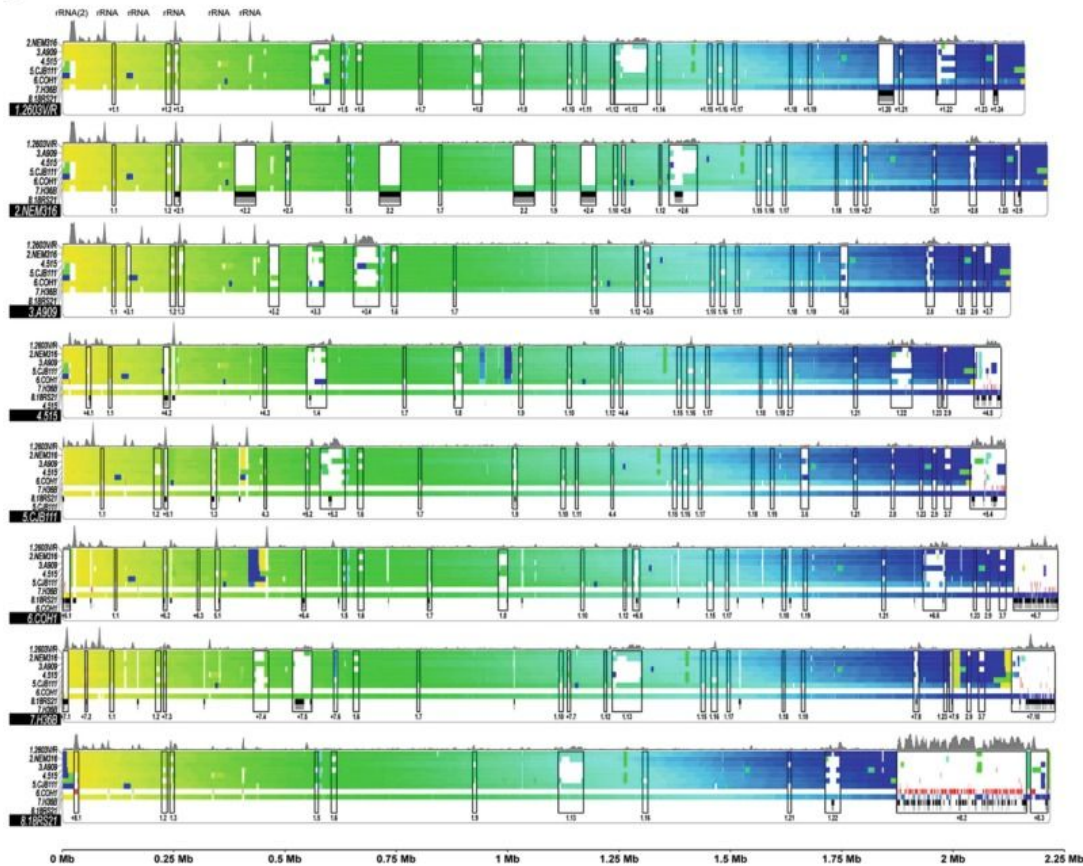




variation graphs: not a new idea!

Fig. 5. A variant graph.

nine versions Valerio Magrelli's poem "Campagna Romana" (1981)



Group B Streptococcus assemblies from 2002

<https://doi.org/10.1007/978-3-030-38281-0>

Pangenome: not a new concept

First collections of multiple genomes from the same species demonstrated substantial differences.

This was unexpected and required new theory to understand.

A single reference is not enough to explain genomic diversity. Even many genomes may not be enough.

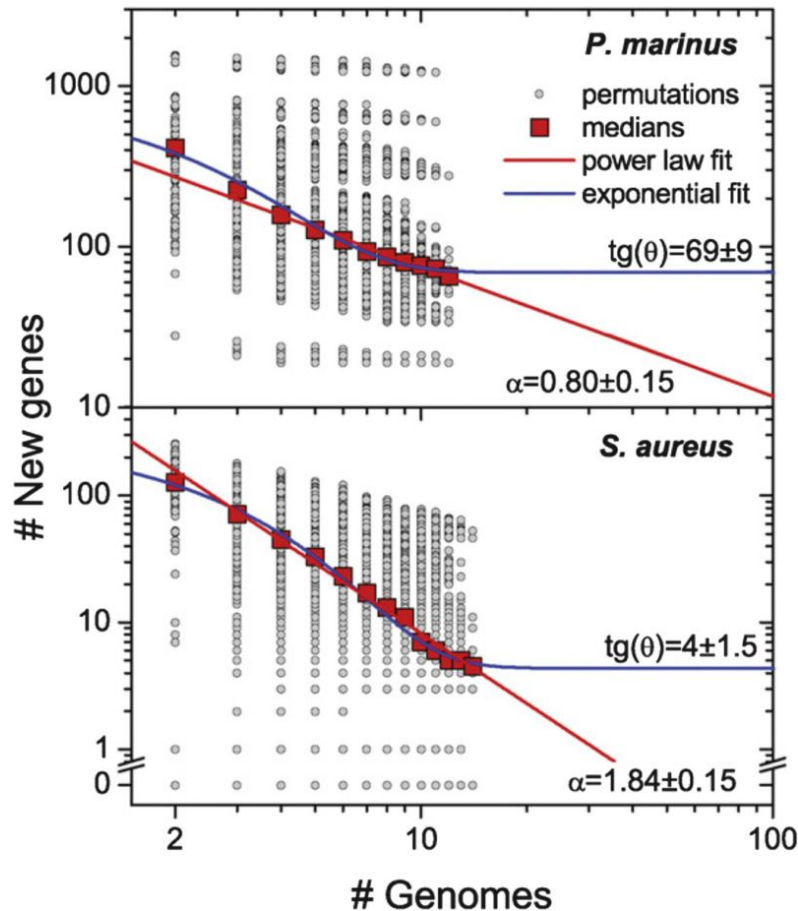
Some genes are shared among all individuals: these are “core”, while others are not—we call them “accessory”.

Lessons from language modeling: Heaps' law

A pangenome is:

Closed: our observations of new genes with new genomes diminish.

Open: we continue to see new genes as we add more genomes.



The exponent α determines whether the pangenome is open ($\alpha \leq 1$) or closed ($\alpha > 1$). The top panel shows data for an open pangenome species, *P. marinus*; the bottom panel for a closed pangenome species, *S. aureus*

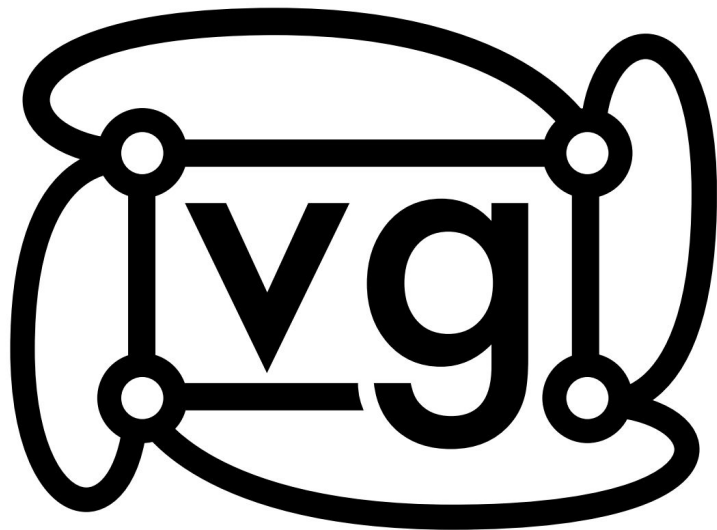


Pangenome research timeline

2000-2010s: counting genes

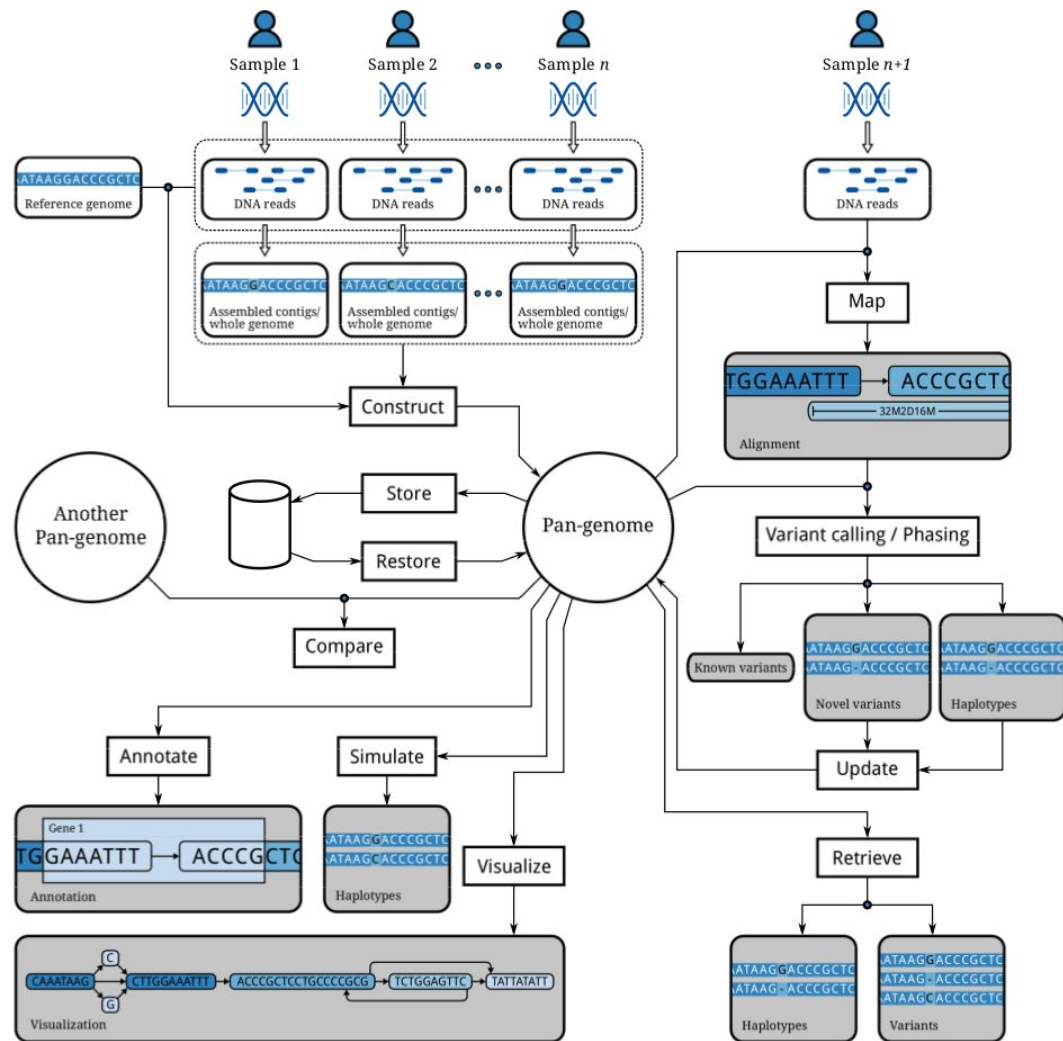
~2015: let's take it to the sequence level (genome graphs)

2020s: complete assemblies (T2T pangenomes)



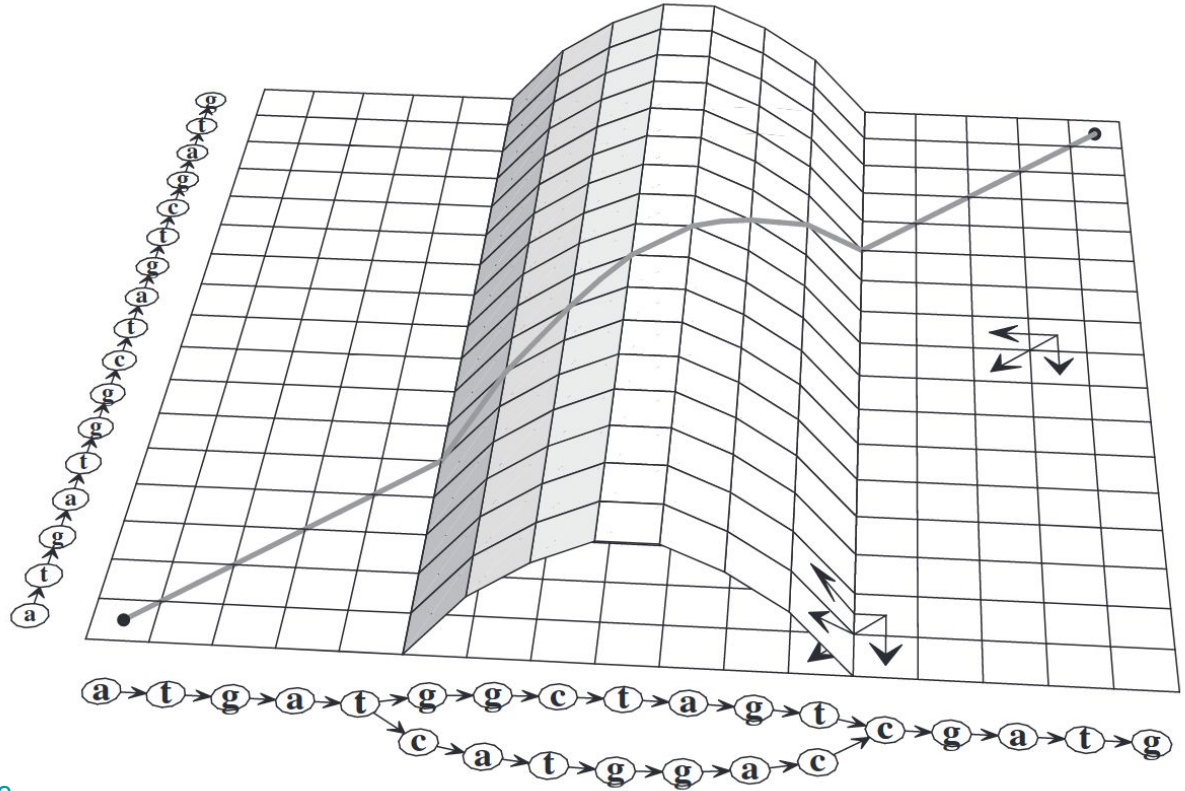
github.com/vgteam/vg

“Computational pan-genomics: status, promises and challenges” <https://doi.org/10.1093%2Fbib%2Fbbw089>



Wait! You can
align sequences
to graphs?

yup... we can generalize most
standard bioinformatic
algorithms to graphs, as in
Partial Order Alignment →



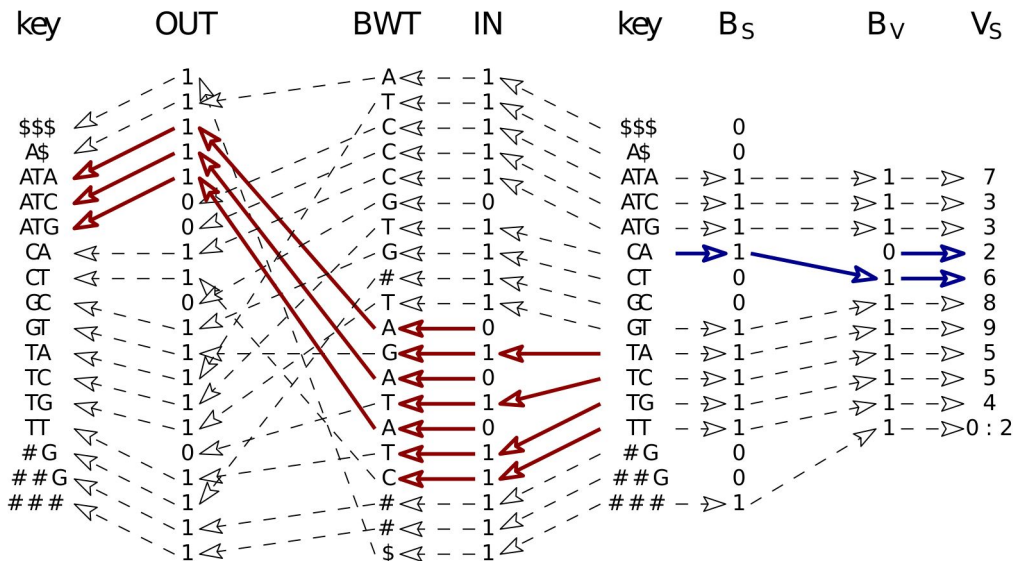
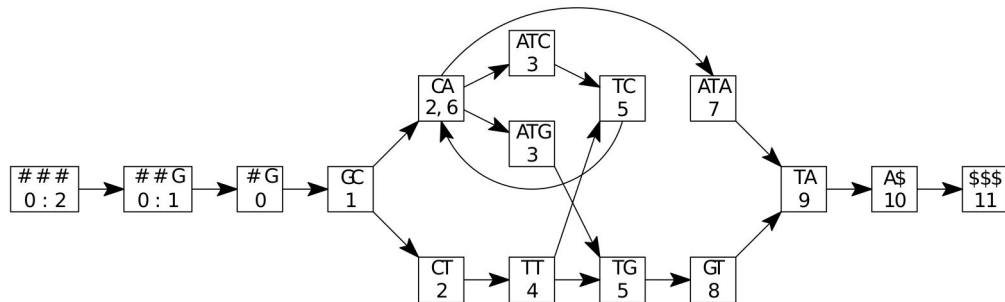
And the FM-index?

Jouni Sirén generalized the FM-index to work on a transformation of the variation graph (technically a de Bruijn graph with $k=256$).

GCSA2 →

We use it to find MEMs just as in bwa mem.

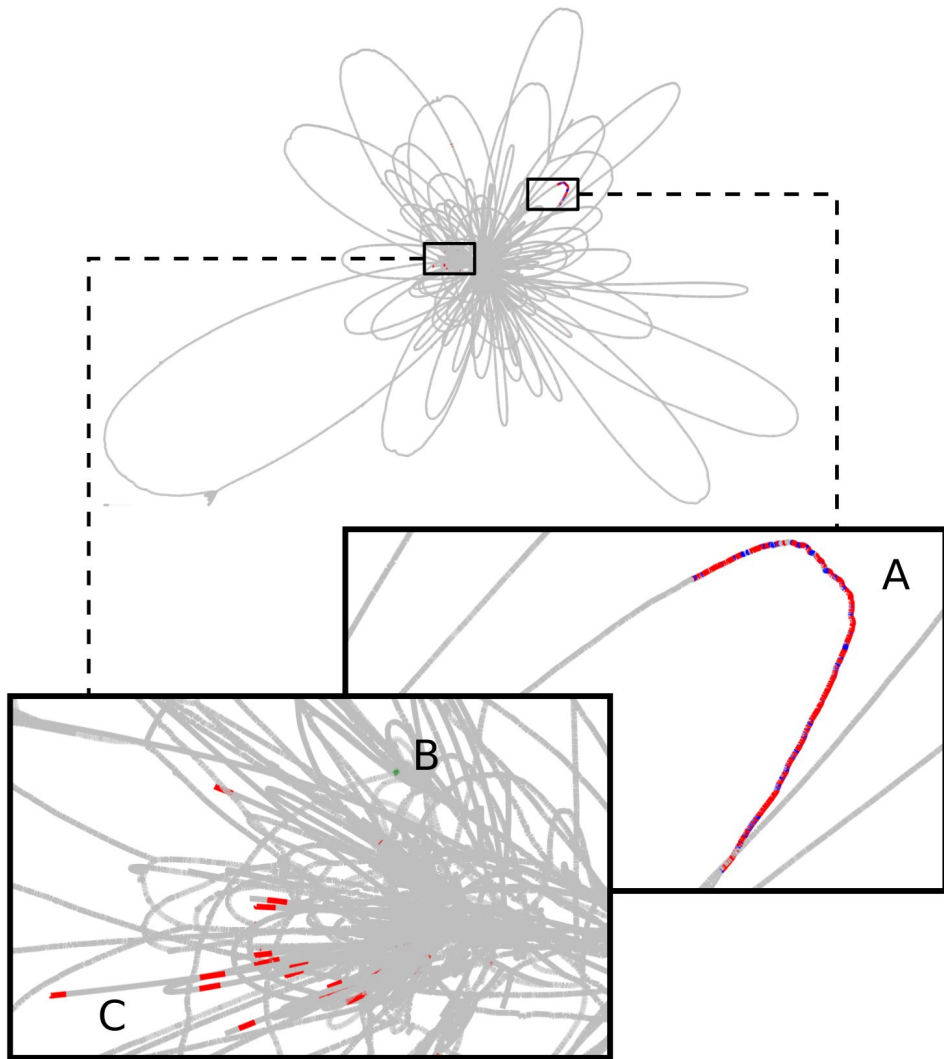
This seeds alignment to the graph.





and long reads too!

a pacbio read vs. a yeast graph:



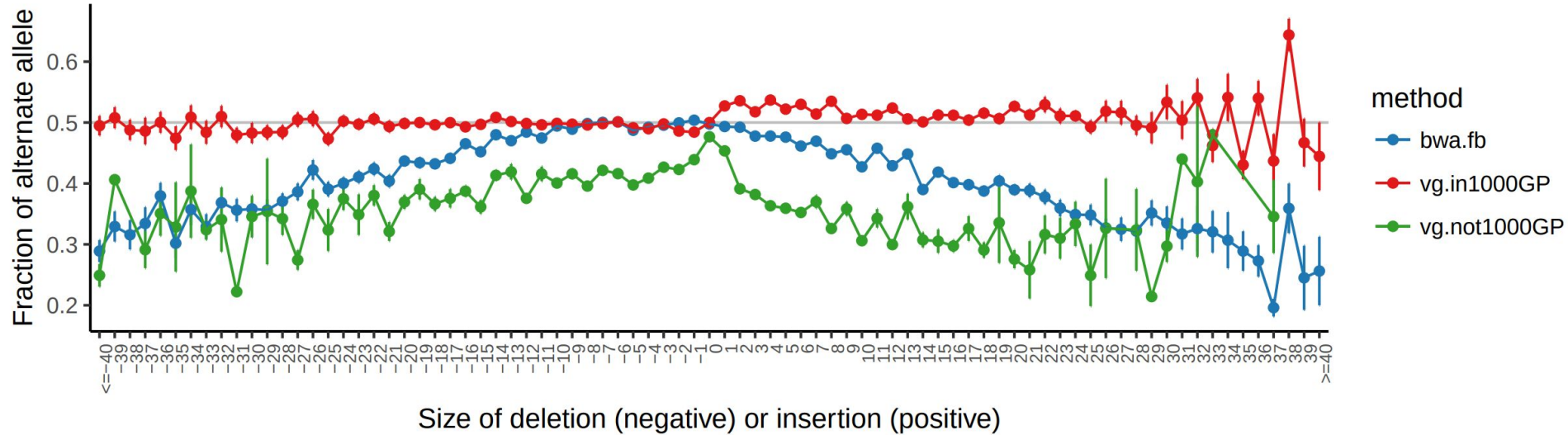
Letter | Published: 20 August 2018

Variation graph toolkit improves read mapping by representing genetic variation in the reference

Erik Garrison , Jouni Sirén, Adam M Novak, Glenn Hickey, Jordan M Eizenga, Eric T Dawson, William Jones, Shilpa Garg, Charles Markello, Michael F Lin, Benedict Paten & Richard Durbin 

Nature Biotechnology **36**, 875–879 (2018) | [Download Citation](#) 

vg resolves reference bias at known indels in *HG002*



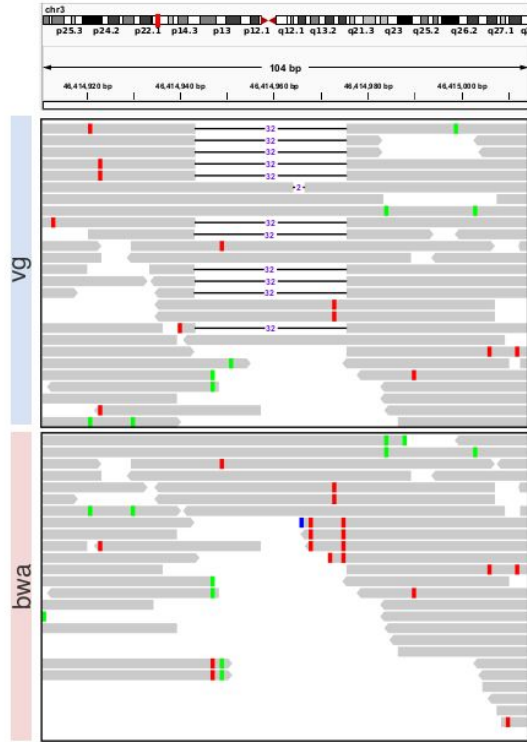
50x 2x150bp Illumina sequencing of *HG002*

Research | [Open Access](#) | [Published: 17 September 2020](#)

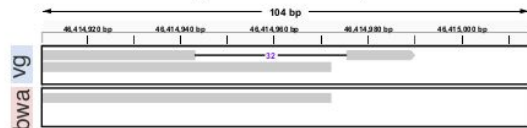
Removing reference bias and improving indel calling in ancient DNA data analysis by mapping to a sequence variation graph

[Rui Martiniano](#), [Erik Garrison](#), [Eppie R. Jones](#), [Andrea Manica](#) & [Richard Durbin](#) 

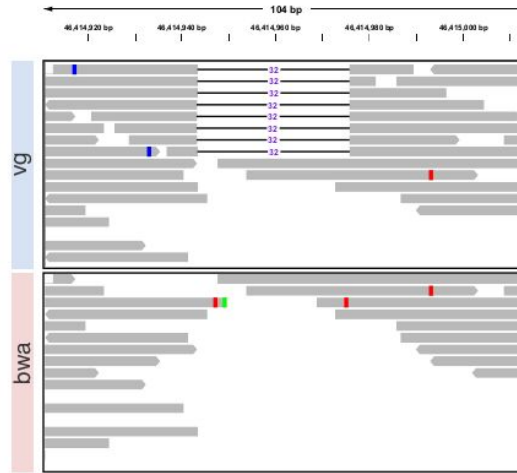
Yamnaya (Early Bronze Age Kazakhstan)



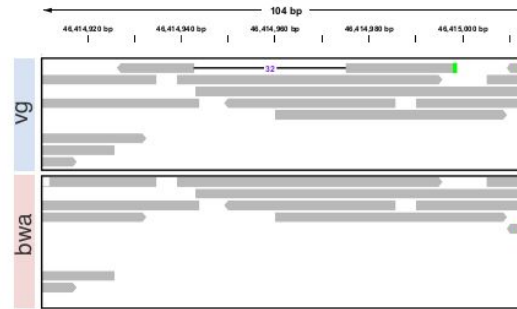
6DT3 (Roman Britain)



12880A (Iron Age Britain)



15577A (Anglo-Saxon Britain)



**Using variation graphs
to observe CCR5-delta
in ancient samples**

Rui Martiniano

 OPEN ACCESS  PEER-REVIEWED

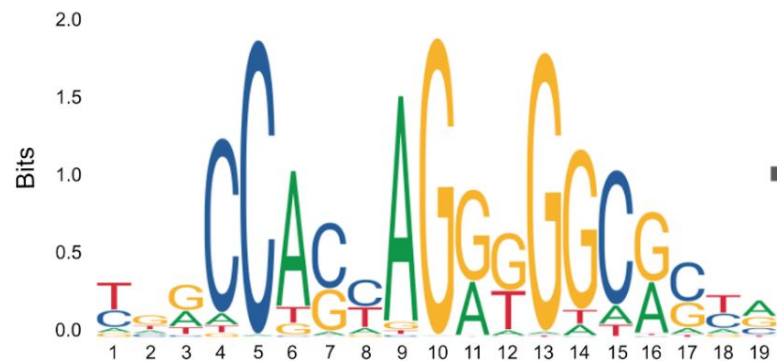
RESEARCH ARTICLE

GRAFIMO: Variant and haplotype aware motif scanning on pangenome graphs

Manuel Tognon, Vincenzo Bonnici, Erik Garrison, Rosalba Giugno , Luca Pinello 

A

CTCF Motif (MA139.1)

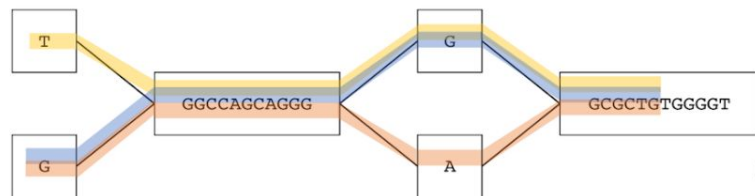


Processed motif PWM

S1	G	G	G	C	C	A	G	C	A	G	G	G	G	G	C	G	C	T	G	
S2	T	G	G	C	C	A	G	C	A	G	G	G	G	G	C	G	C	T	G	
S3	G	G	G	C	C	A	G	C	A	G	G	G	G	A	G	C	G	C	T	G
A	754	812	858	716	545	943	686	773	955	504	873	713	581	716	771	883	750	781	890	
C	893	832	736	982	993	620	945	928	605	51	491	613	51	577	975	620	938	902	852	
G	775	924	931	662	51	761	905	734	699	993	952	941	992	979	536	942	898	772	886	
T	901	822	792	677	424	758	573	870	616	459	555	875	545	737	732	611	680	888	721	

B

Pangenome variation graph (VG)


Reference motif candidate: **GGGCCAGCAGGGGGCGCTG**Haplotype motif candidate: **TGGCCAGCAGGGGGCGCTG**Haplotype motif candidate: **GGGCCAGCAGGGAGCGCTG**

Retrieved motif occurrences and haplotype frequencies

Sequence	Log-odds score	P-value	q-value	Reference	Haplotype frequency
GGGCCAGCAGGGGGCGCTG	28.22	7.51e ⁻¹²	3.86e ⁻⁶	non ref.	32
TGGCCAGCAGGGGGCGCTG	26.16	3.12e ⁻¹⁰	7.72e ⁻⁶	ref.	5063
GGGCCAGCAGGGAGCGCTG	19.43	1.71e ⁻⁷	1.73e ⁻⁴	non ref.	1

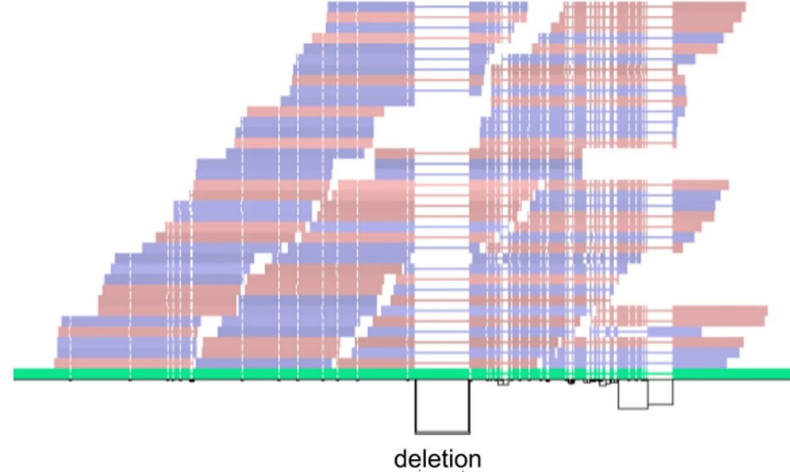
Method | [Open Access](#) | [Published: 12 February 2020](#)

Genotyping structural variants in pangenome graphs using the vg toolkit

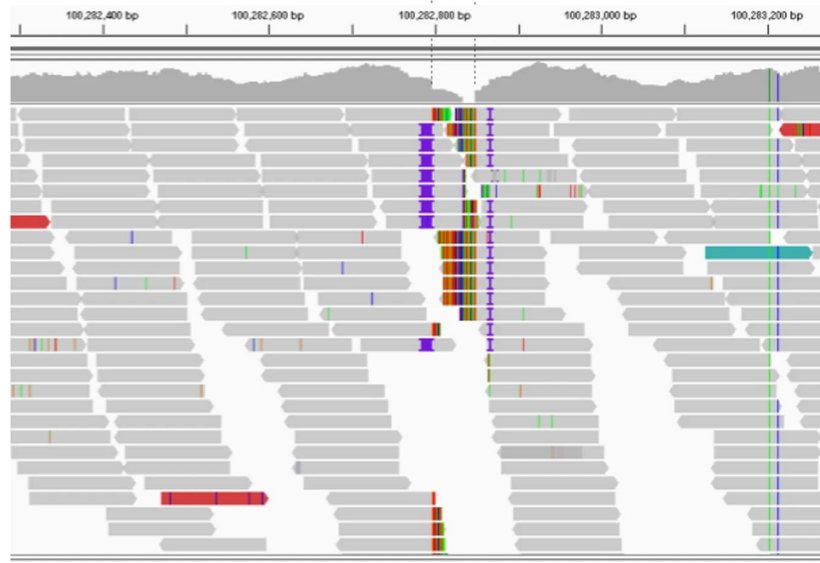
[Glenn Hickey](#), [David Heller](#), [Jean Monlong](#), [Jonas A. Sibbesen](#), [Jouni Sirén](#), [Jordan Eizenga](#), [Eric T. Dawson](#), [Erik Garrison](#), [Adam M. Novak](#) & [Benedict Paten](#) 

**Exonic deletion in the
HGSVC dataset correctly
genotyped by vg**

a)



b)



HOME > SCIENCE > VOL. 374, NO. 6574 > PANGENOMICS ENABLES GENOTYPING OF KNOWN STRUCTURAL VARIANTS IN 5202 DIVERSE GENOMES



RESEARCH ARTICLE | GENOMICS

Pangenomics enables genotyping of known structural variants in 5202 diverse genomes

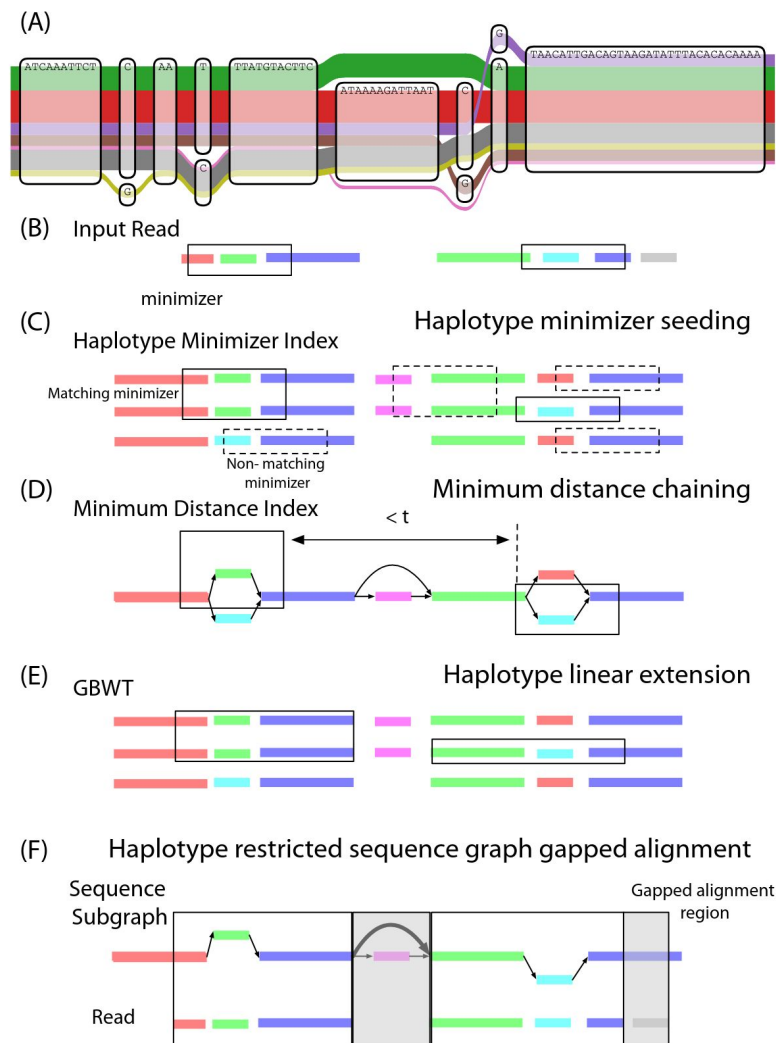
JOUNI SIRÉN , JEAN MONLONG , XIAN CHANG , ADAM M. NOVAK , JORDAN M. EIZENGA , CHARLES MARKELLO , JONAS A. SIBBESEN ,
GLENN HICKEY , PI-CHUAN CHANG , ANDREW CARROLL , NAMRATA GUPTA , STACEY GABRIEL, THOMAS W. BLACKWELL, AAKROSH RATAN ,
KENT D. TAYLOR , STEPHEN S. RICH , JEROME I. ROTTER , DAVID HAUSSLER , ERIK GARRISON, AND BENEDICT PATEN

fewer

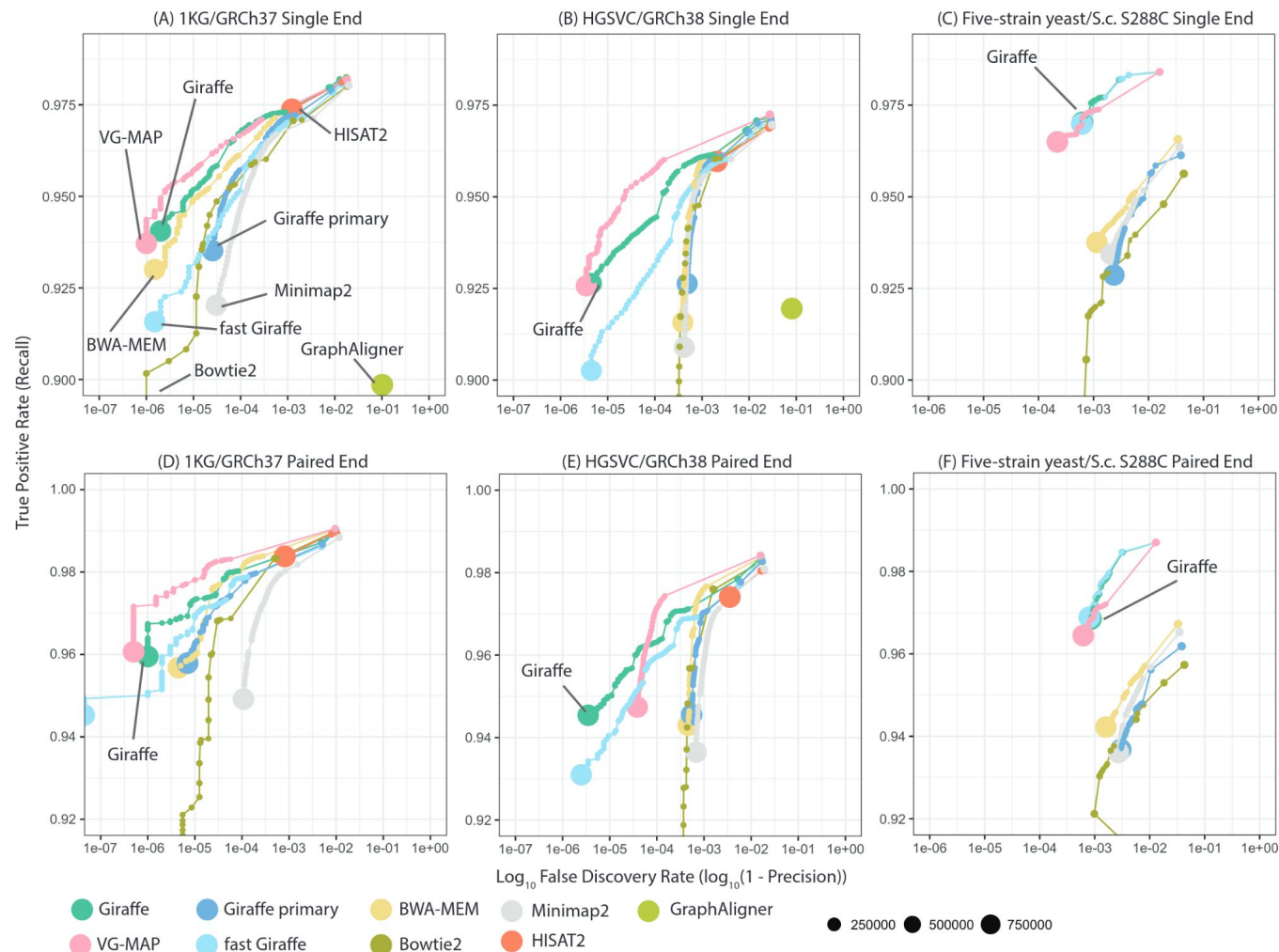
[Authors Info &](#)

[Affiliations](#)

vg giraffe: approach

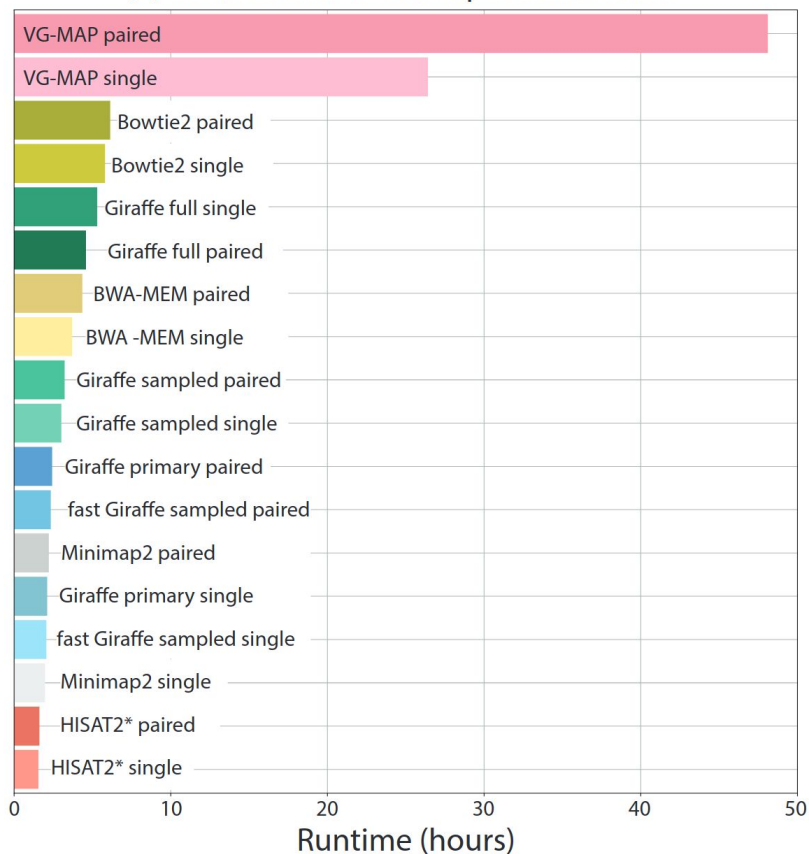


vg giraffe
is accurate
enough

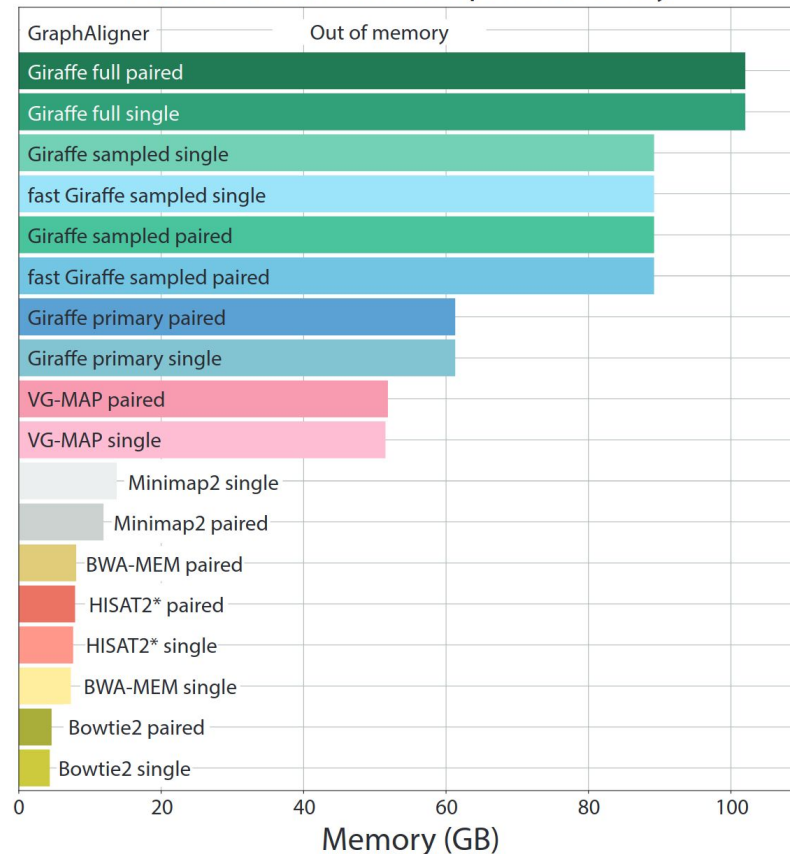


vg giraffe is very fast

(A) 1KG/GRCh37 NovaSeq 6000 Runtime

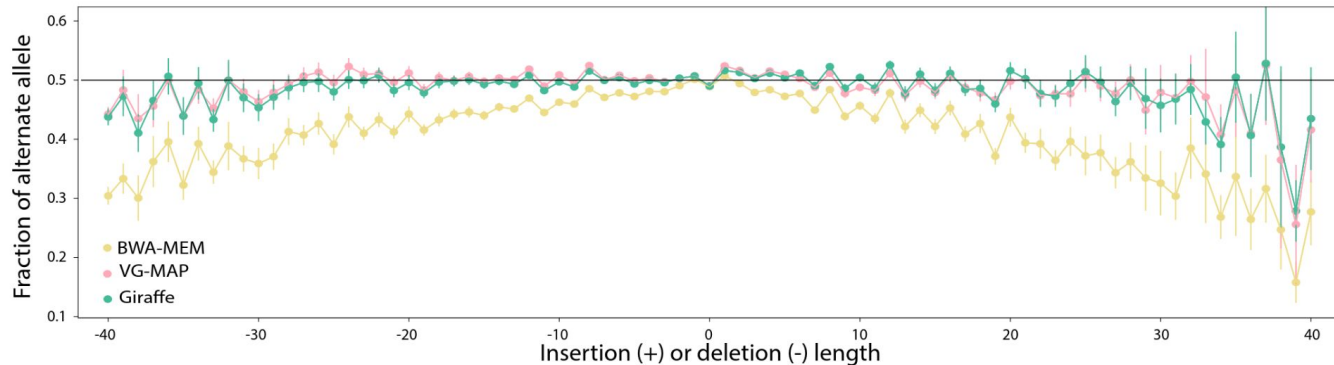


(C) 1KG/GRCh37 NovaSeq 6000 Memory

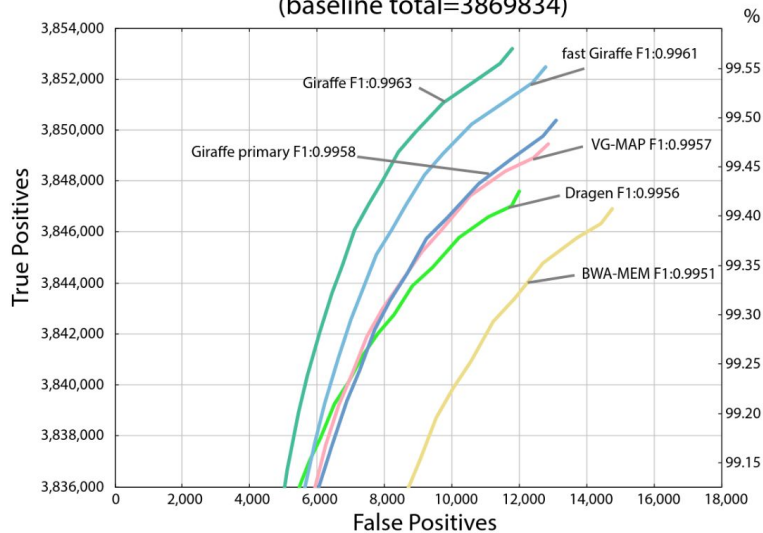


vg giraffe improves variant calling

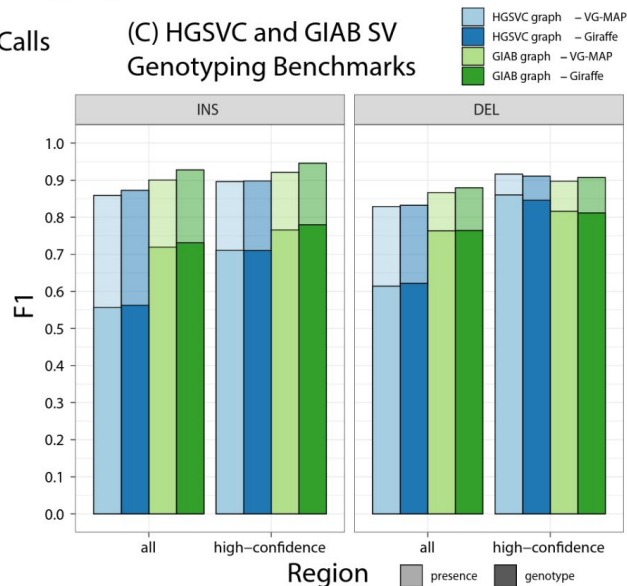
(A) Allele Balance - NovoSeq 6000 reads mapped to 1KG/GRCh37



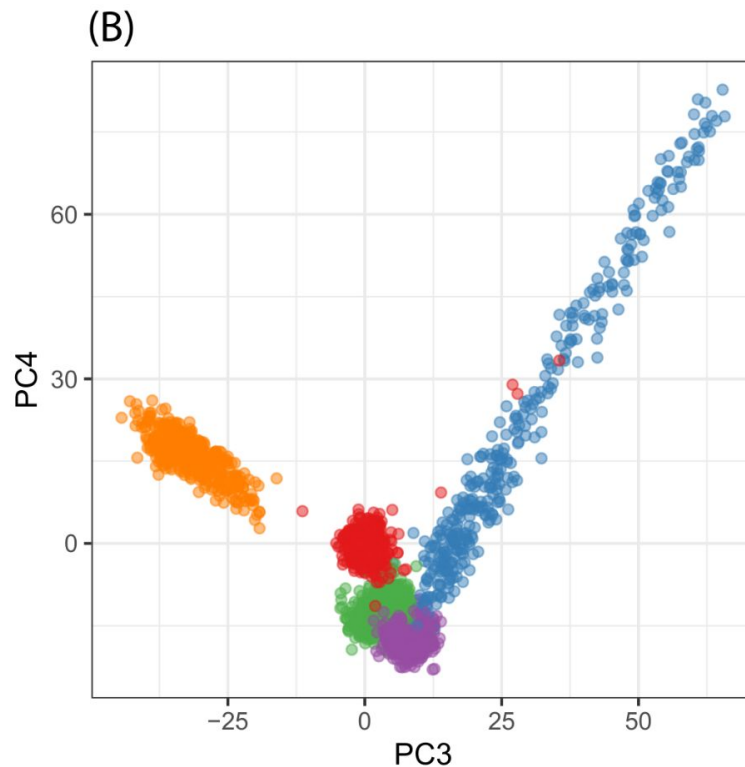
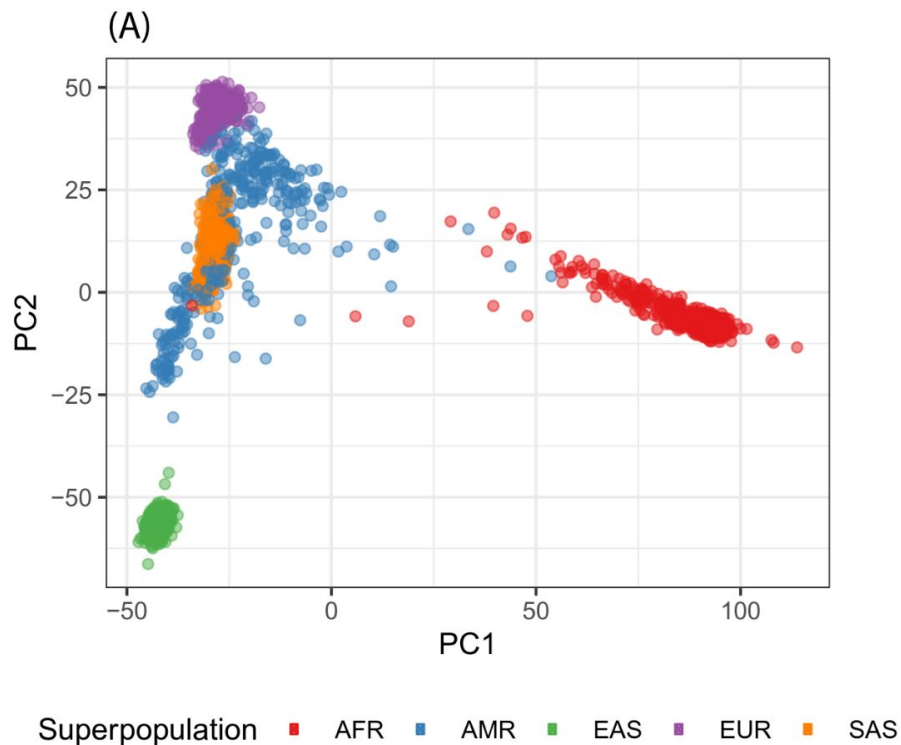
(B) HG002 v4.1 WGS High Confidence Regions Dragen Genotype Calls (baseline total=3869834)



(C) HGSVC and GIAB SV Genotyping Benchmarks



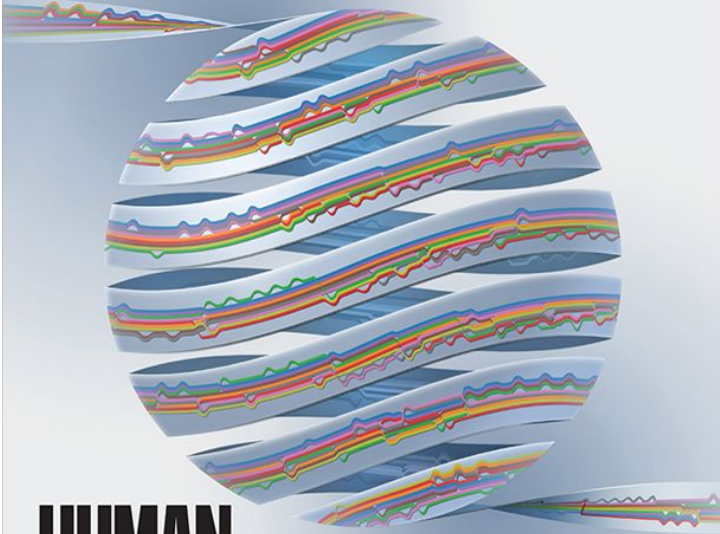
vg giraffe lets us scale: PCA from SVs in 5k genomes



The human pangenome project

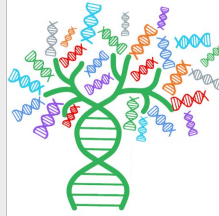
The international journal of science / 11 May 2023

nature



HUMAN PANGENOME

Data from 47 individuals combine to create
reference resource that reflects human diversity

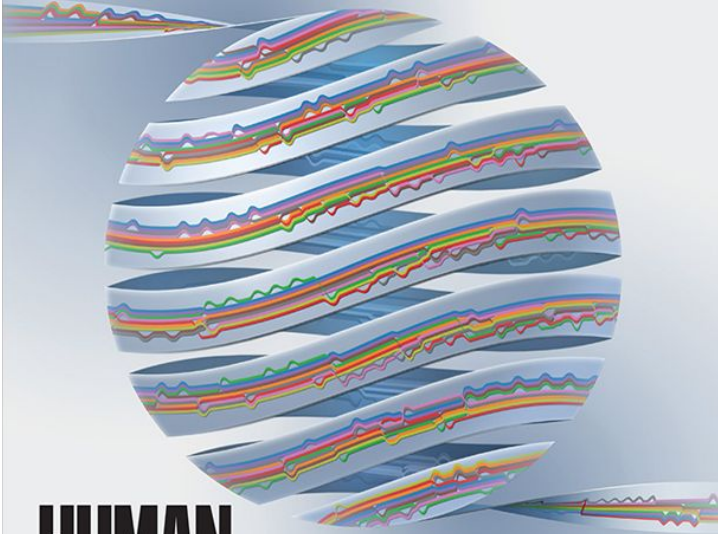


Human Pangenome Reference Consortium

- 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
- **First draft is available!**

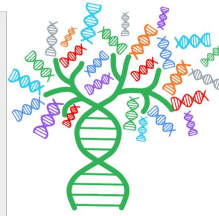
The international journal of science / 11 May 2023

nature



HUMAN PANGENOME

Data from 47 individuals combine to create
reference resource that reflects human diversity



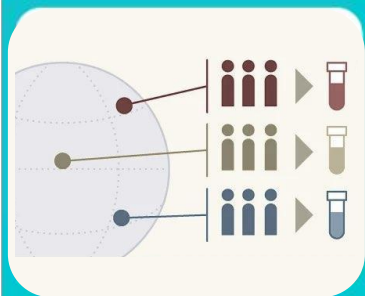
The Human Pangenome

Composed of three As:

- **Assemblies**
 - Haplotype resolved, soon T2T, but also 37, 38, T2T-CHM13.
- **Alignment**
 - Provides canonical homology information
- **Annotations**
 - Genes, etc. Should be consistent with alignment



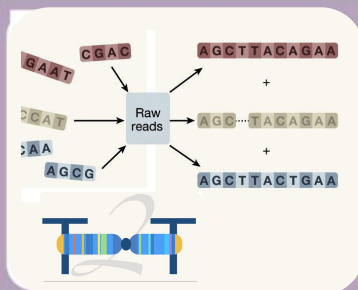
Population sampling and representation



Technology Production



Phased/Finished T2T Assemblies



Pangenome and new workflows/tooling

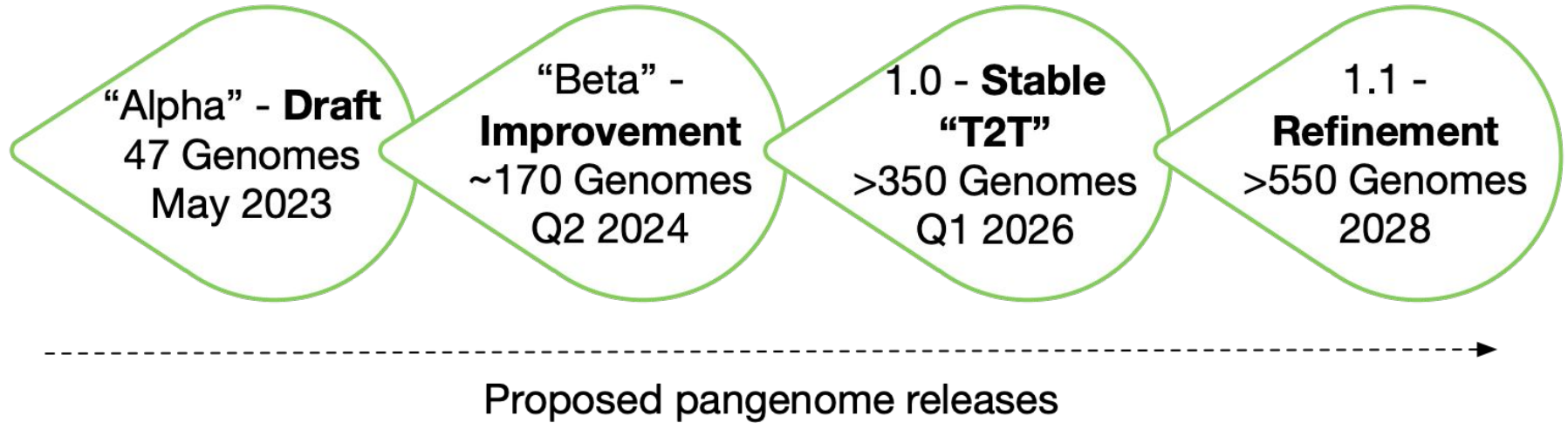


International Pangenome Project



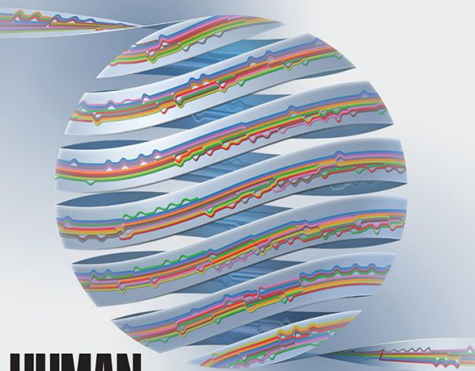
Embedded Ethics and Policy:
Inter-disciplinary ethics working group/oversight committee

Human Pangenome Timeline



Building a draft human pangenome

nature



HUMAN
PANGENOME

Data from 47 individuals combine to create
reference resource that reflects human diversity

Article

A draft human pangenome reference

<https://doi.org/10.1038/s41586-023-05896-x>

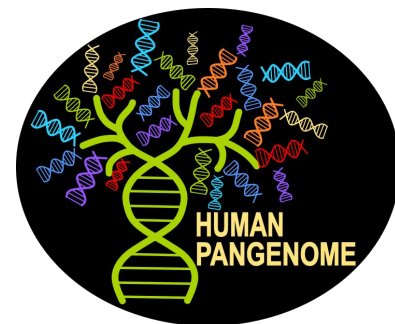
Received: 9 July 2022

Accepted: 28 February 2023

Published online: 10 May 2023

Open access

 Check for updates

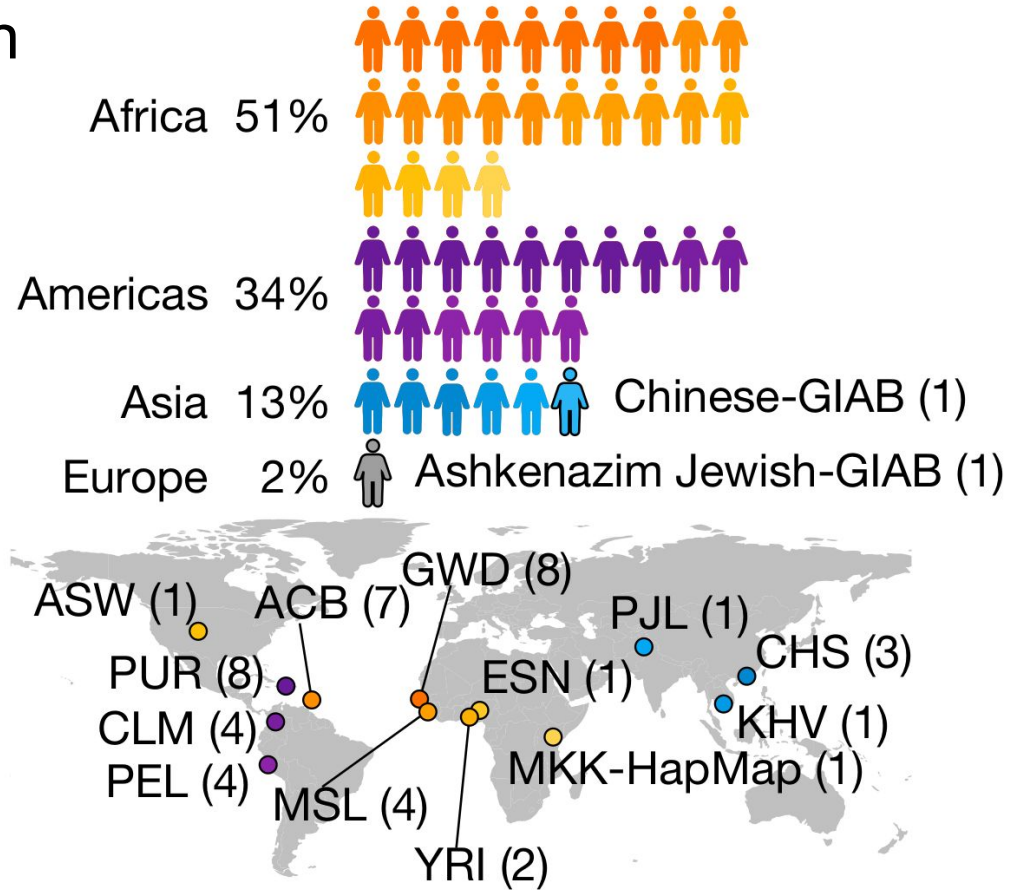


Wen-Wei Liao^{1,2,3,60}, Mobin Asri^{4,60}, Jana Ebler^{5,6,60}, Daniel Doerr^{5,6}, Marina Haukness⁴, Glenn Hickey⁴, Shuangjia Lu^{1,2}, Julian K. Lucas⁴, Jean Monlong⁴, Haley J. Abel⁷, Silvia Buonaiuto⁸, Xian H. Chang⁴, Haoyu Cheng^{9,10}, Justin Chu⁹, Vincenza Colonna^{8,11}, Jordan M. Eizenga⁴, Xiaowen Feng^{9,10}, Christian Fischer¹¹, Robert S. Fulton^{12,13}, Shilpa Garg¹⁴, Cristian Groza¹⁵, Andrea Guarracino^{11,16}, William T. Harvey¹⁷, Simon Heumos^{18,19}, Kerstin Howe²⁰, Miten Jain²¹, Tsung-Yu Lu²², Charles Markello⁴, Fergal J. Martin²³, Matthew W. Mitchell²⁴, Katherine M. Munson¹⁷, Moses Njagi Mwaniki²⁵, Adam M. Novak⁴, Hugh E. Olsen⁴, Trevor Pesout⁴, David Porubsky¹⁷, Pjotr Prins¹¹, Jonas A. Sibbesen²⁶, Jouni Sirén⁴, Chad Tomlinson¹², Flavia Villani¹¹, Mitchell R. Vollger^{17,27}, Lucinda L. Antonacci-Fulton¹², Gunjan Baid²⁸, Carl A. Baker¹⁷, Anastasiya Belyaeva²⁸, Konstantinos Billis²³, Andrew Carroll²⁸, Pi-Chuan Chang²⁸, Sarah Cody¹², Daniel E. Cook²⁸, Robert M. Cook-Deegan²⁹, Omar E. Cornejo³⁰, Mark Diekhans⁴, Peter Ebert^{5,8,31}, Susan Fairley²³, Olivier Fedrigo³², Adam L. Felsenfeld³³, Giulio Formenti³², Adam Frankish²³, Yan Gao³⁴, Nanibaa' A. Garrison^{35,36,37}, Carlos Garcia Giron²³, Richard E. Green^{38,39}, Leanne Haggerty²³, Kendra Hoekzema¹⁷, Thibaut Hourlier²³, Hanlee P. Ji⁴⁰, Eimear E. Kenny⁴¹, Barbara A. Koenig⁴², Alexey Kolesnikov²⁸, Jan O. Korbel^{23,43}, Jennifer Kordosky¹⁷, Sergey Koren⁴⁴, HoJoon Lee⁴⁰, Alexandra P. Lewis¹⁷, Hugo Magalhães^{5,6}, Santiago Marco-Sola^{45,46}, Pierre Marijon^{5,6}, Ann McCartney⁴⁴, Jennifer McDaniel⁴⁷, Jacquelyn Mountcastle³², Maria Nattestad²⁸, Sergey Nurk⁴⁴, Nathan D. Olson⁴⁷, Alice B. Popejoy⁴⁸, Daniela Puiu⁴⁹, Mikko Rautiainen⁴⁴, Allison A. Regier¹², Arang Rhie⁴⁴, Samuel Sacco³⁰, Ashley D. Sanders⁵⁰, Valerie A. Schneider⁵¹, Baergen I. Schultz³³, Kishwar Shafin²⁸, Michael W. Smith³³, Heidi J. Sofia³³, Ahmad N. Abou Tayoun^{52,53}, Françoise Thibaud-Nissen⁵¹, Francesca Floriana Tricomi²³, Justin Wagner⁴⁷, Brian Walenz⁴⁴, Jonathan M. D. Wood²⁰, Aleksey V. Zimin^{49,54}, Guillaume Bourque^{55,56,57}, Mark J. P. Chaisson²², Paul Flicek²³, Adam M. Phillippy⁴⁴, Justin M. Zook⁴⁷, Evan E. Eichler^{17,58}, David Haussler^{4,58}, Ting Wang^{12,13}, Erich D. Jarvis^{32,58,59}, Karen H. Miga⁴, Erik Garrison¹¹, Tobias Marschall^{5,6}, Ira M. Hall^{1,2}, Heng Li^{9,10} & Benedict Paten⁴

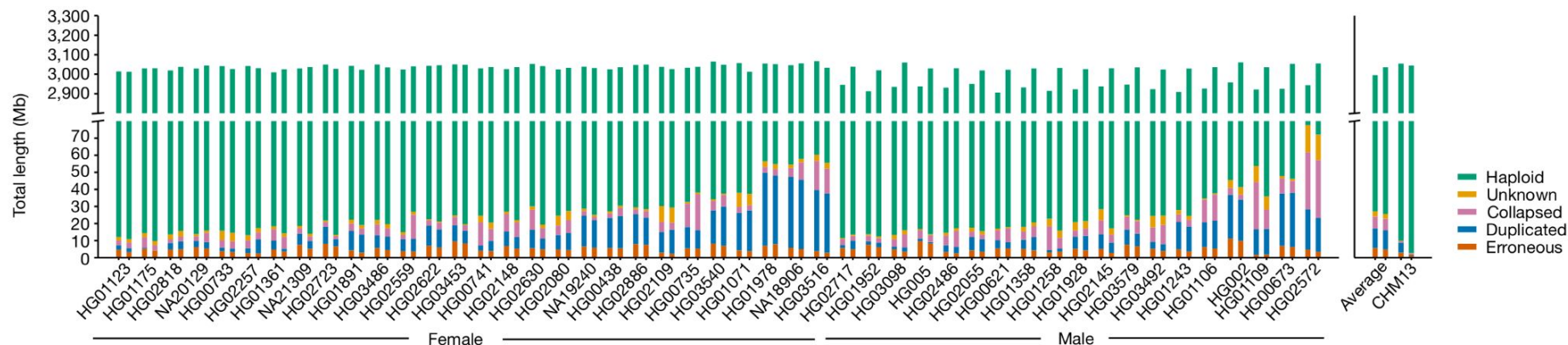
Draft pangenome composition

Sample selection was constrained by:

- trio status in Coriell biobank (-Europeans)
- low cell line passage count (--Europeans)
- genetic diversity (+++Africans)
- drift (+Asians, ++Americas)



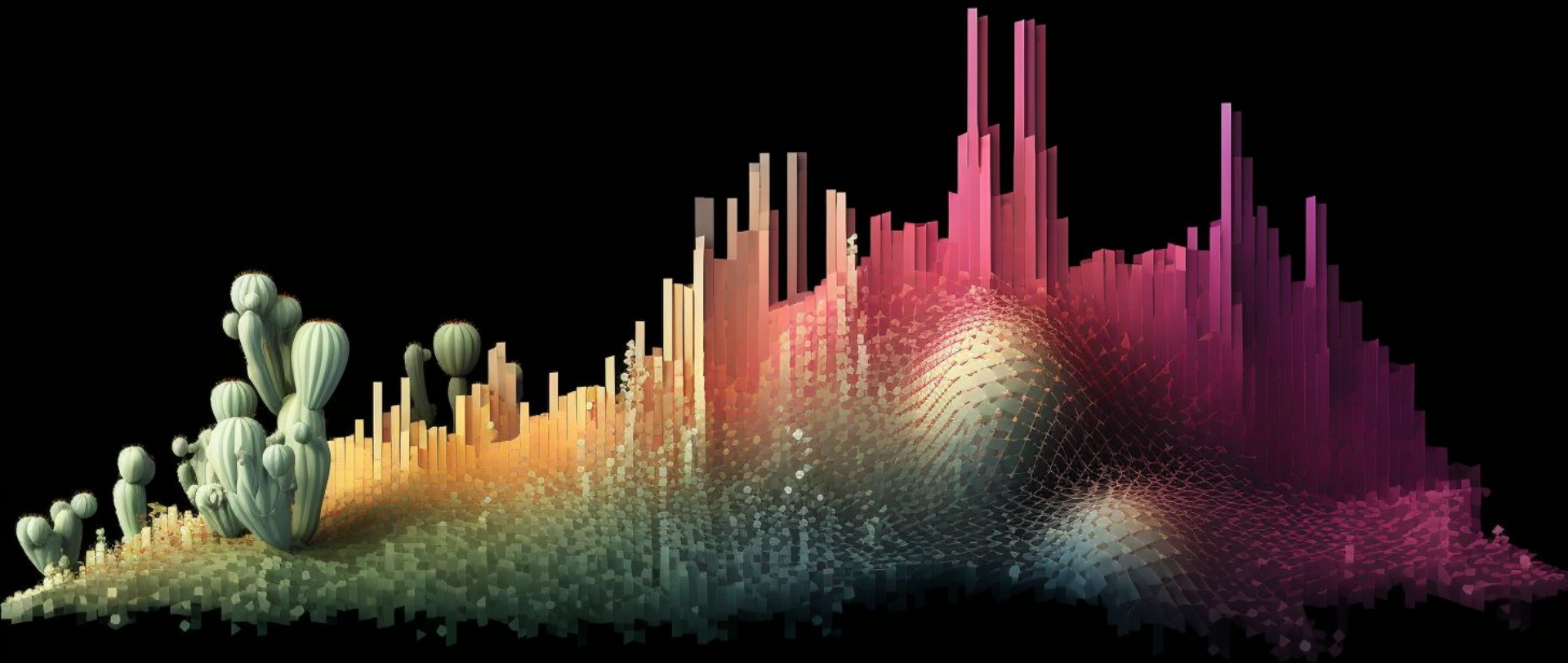
Amazing assemblies approach reference quality



Haplotype-resolved assemblies from trio-hifiasm.

They are really good, according to realignment of reads to the assemblies and model of assembly completeness—nearly as good as T2T-CHM13!

Mobin Asri



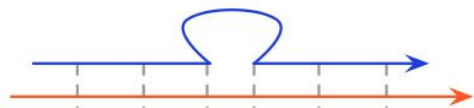
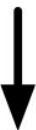
Then we made 5 pangenome (reference) graphs...

Minigraph

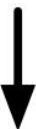
Graph 1
(asm 1):



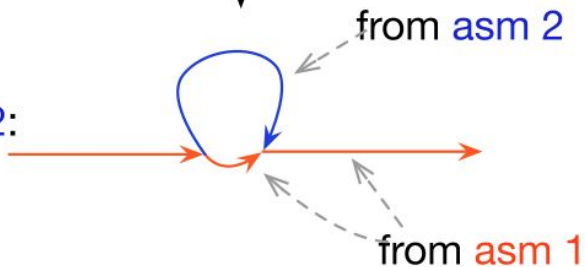
Align asm 2
to graph 1



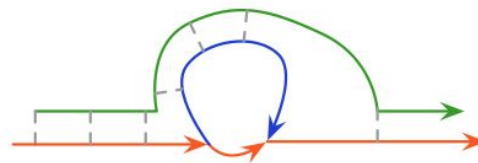
Construct
graph 1/2



Coarse
graph 1/2:



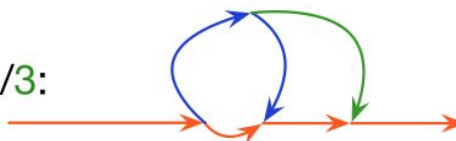
Align asm 3
to graph 1/2



Construct
graph 1/2/3

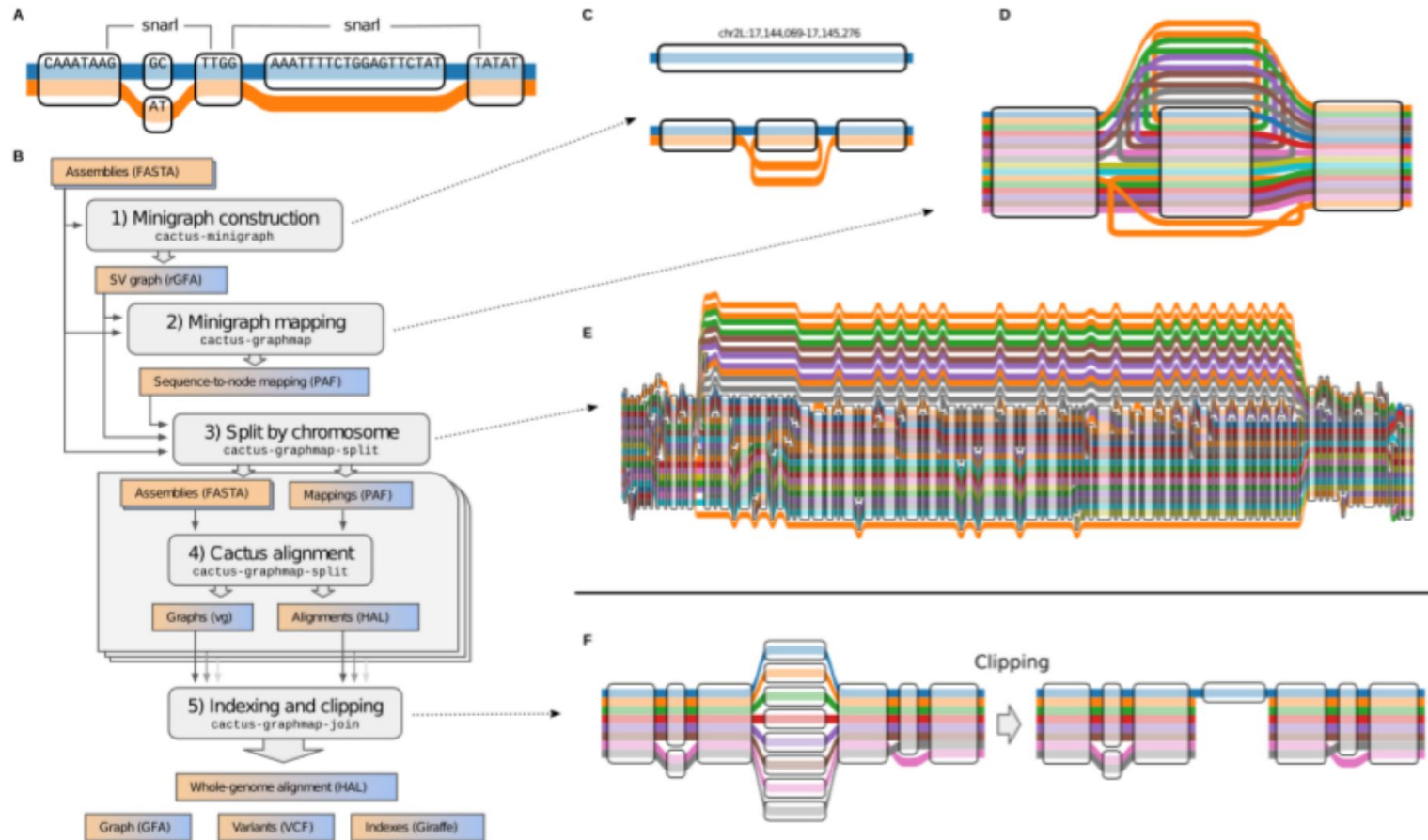


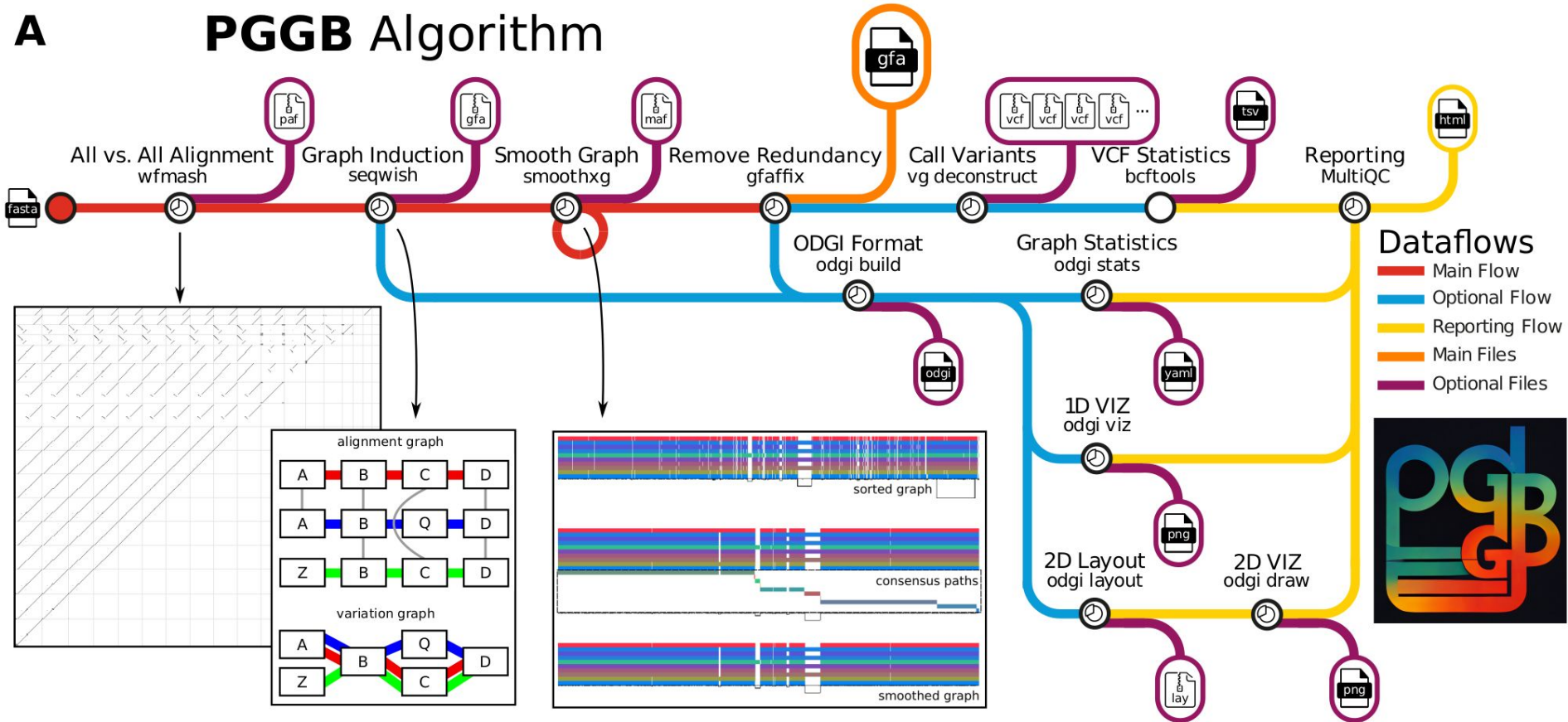
Coarse
graph 1/2/3:



Minigraph-Cactus

<https://doi.org/10.1038/s41587-023-01793-w>





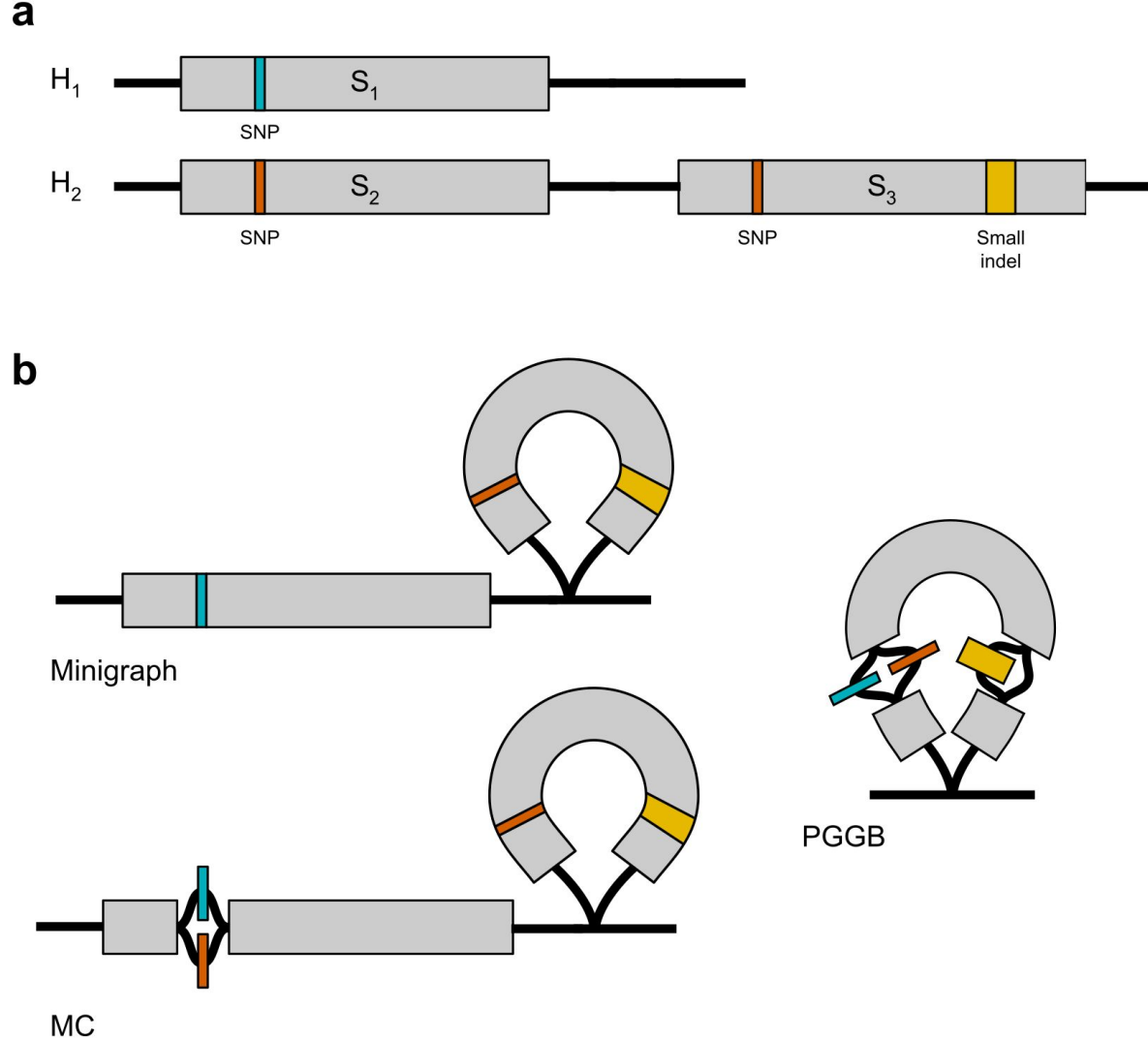
Key conceptual differences between HPRC pangenome construction methods

minigraph: just SVs, no complex stuff, one reference.

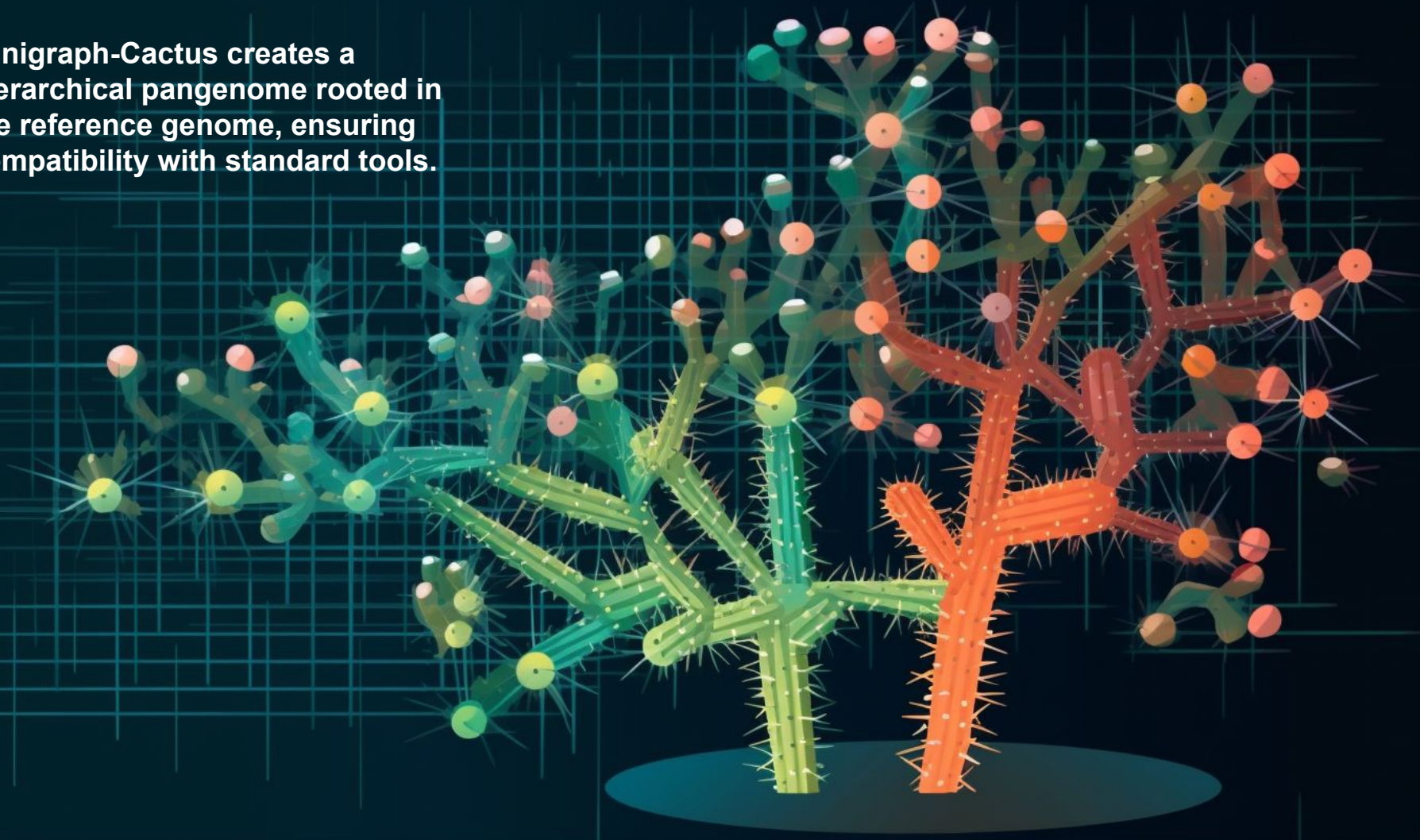
minigraph-cactus: add SNPs, clean up the breakpoints, useful for alignment, one reference.

pggb: everything-vs-everything, hard to align to, useful for studying evolution and pangenome structure at all scales, all genomes are references.

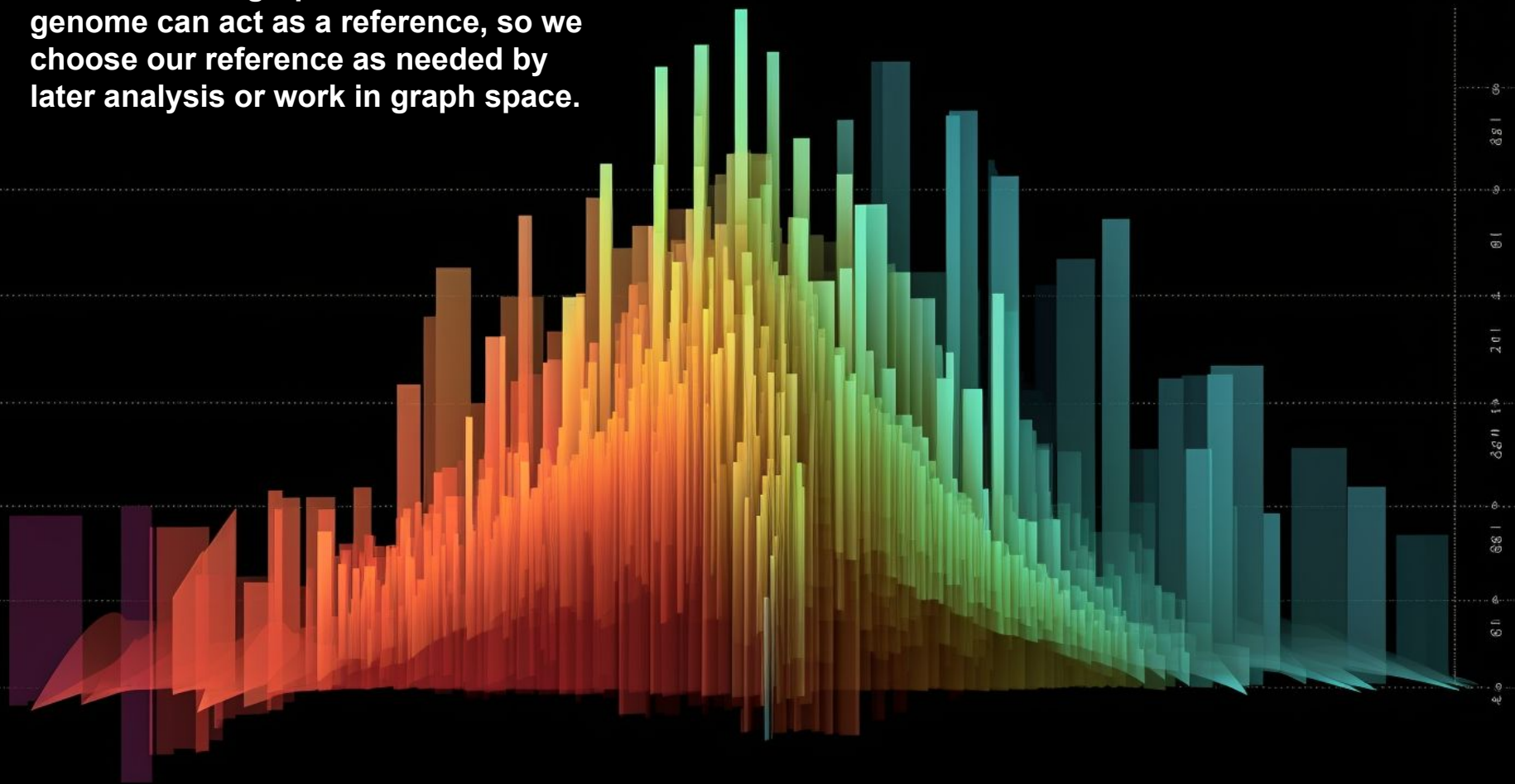
“Collapse” in high-copy repeats →



minigraph-Cactus creates a hierarchical pangenome rooted in the reference genome, ensuring compatibility with standard tools.

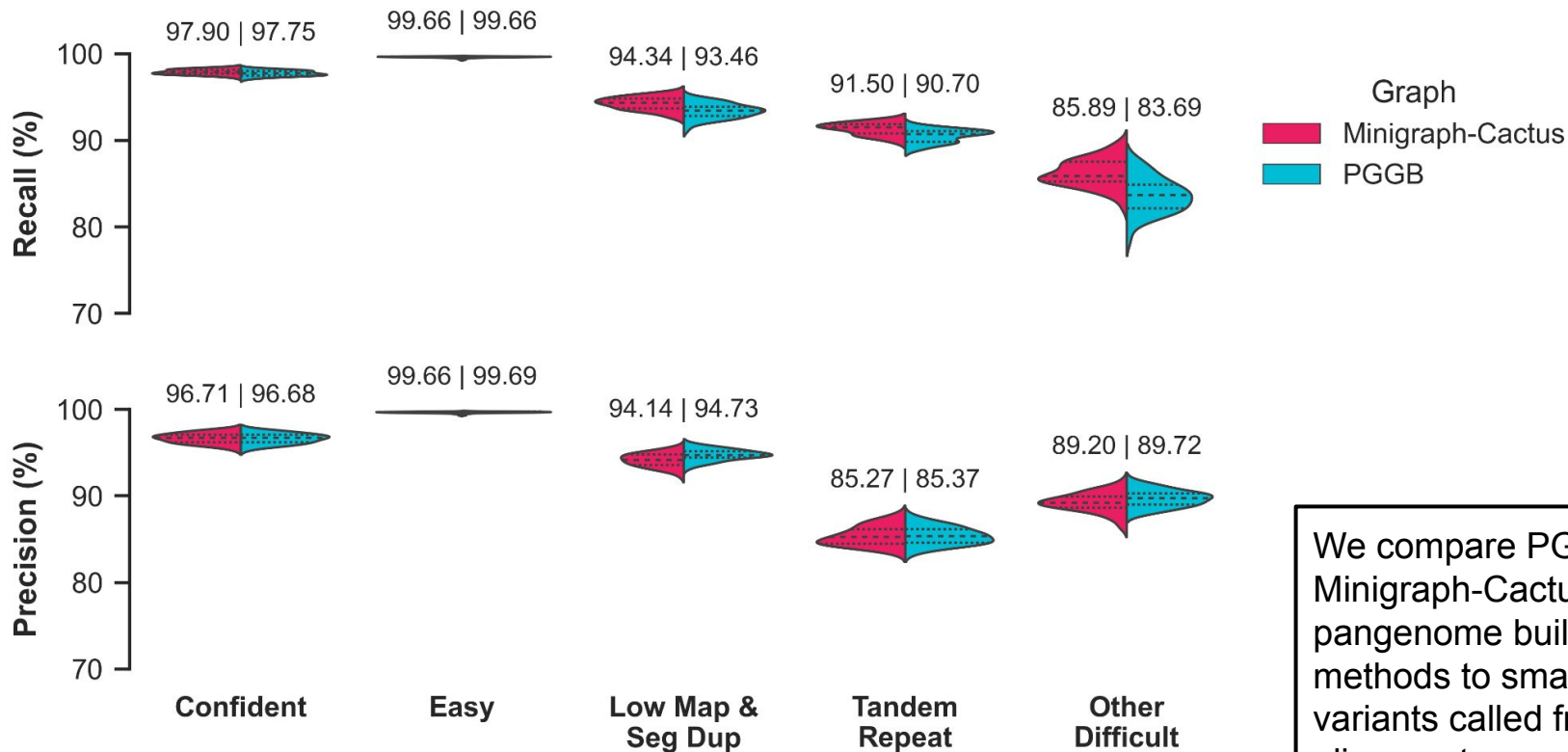


PGGB creates graphs in which each genome can act as a reference, so we choose our reference as needed by later analysis or work in graph space.



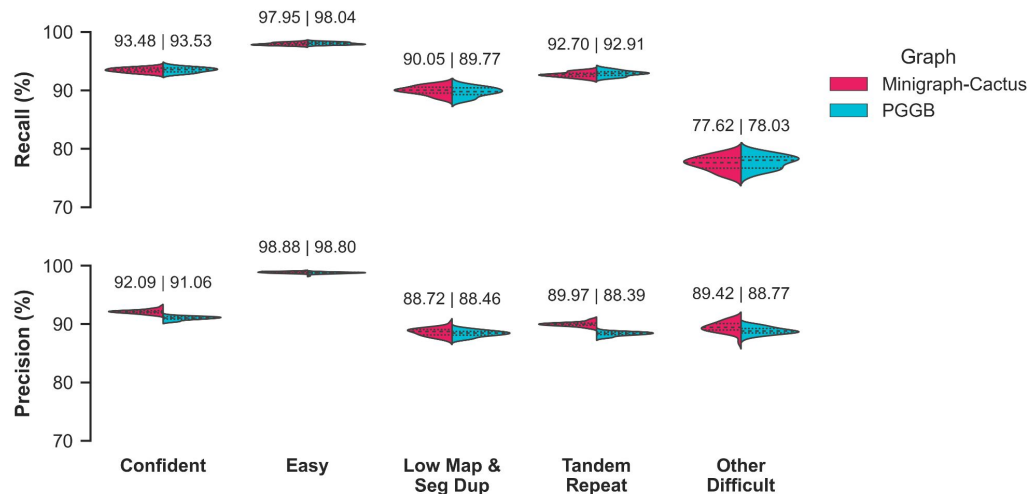
Multiple graph building methods show consistent quality

* variants extracted from graphs with vg deconstruct

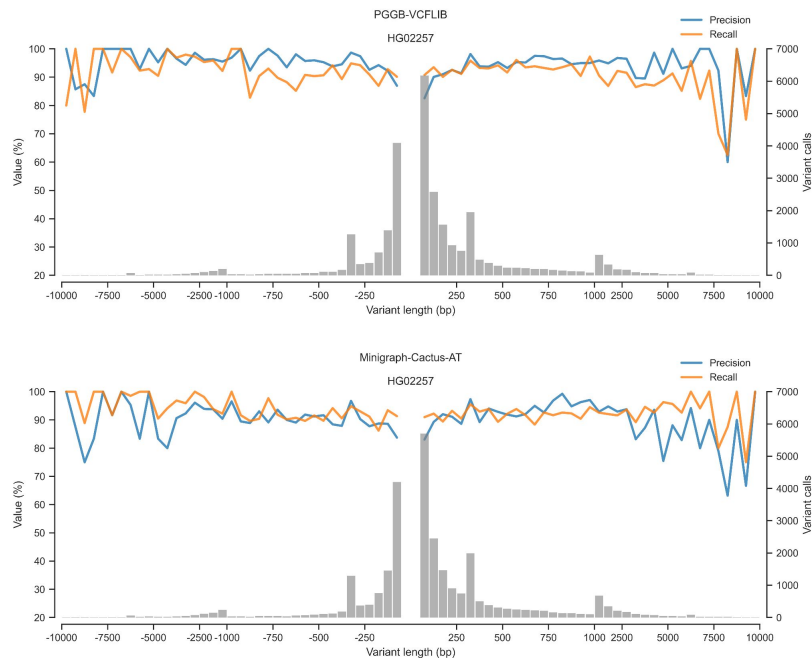


We compare PGGB and Minigraph-Cactus pangenome building methods to small variants called from HiFi alignments.

The graphs accurately characterize structural variants

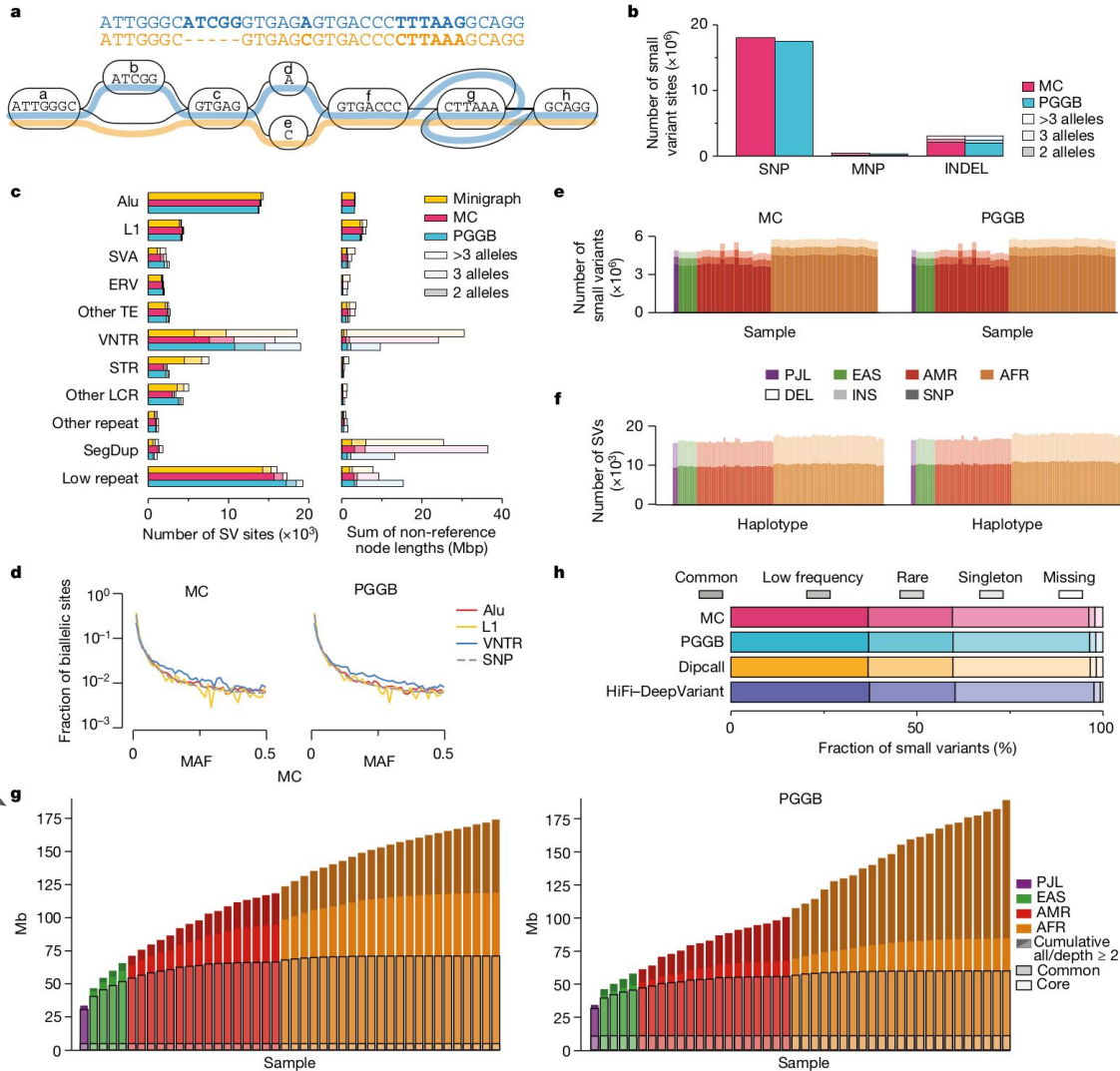


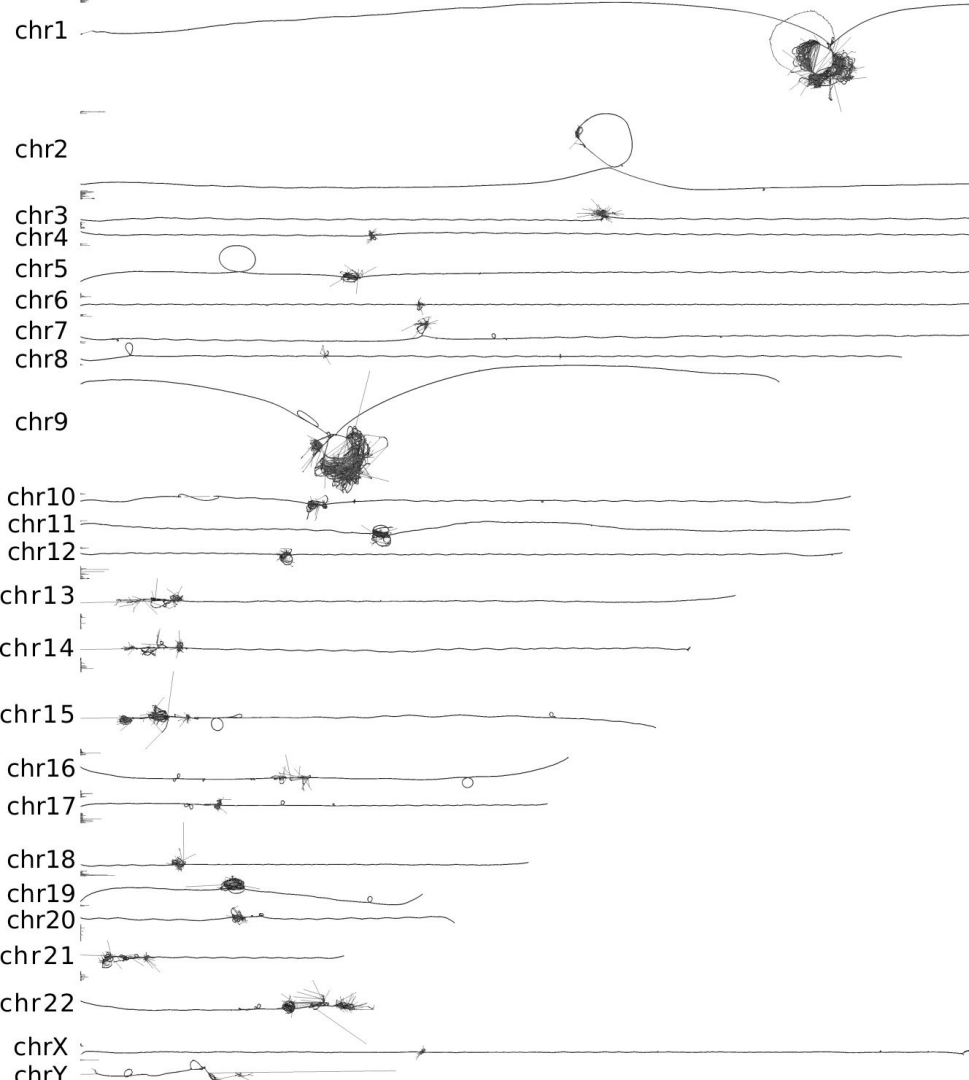
Comparing to consensus calls made by many reference-based SV callers.





MC and PGGB show very similar pictures of the pangenome

- T2T-CHM13 adds **~200MB of heterochromatin** to reference.
- Draft pangenome adds **~100MB of polymorphic euchromatin** (and a lot more heterochromatin),
- **0.6-4.4 Mb of additional genic sequences per haplotype** compared to GRCh38 (38 gene CNVs/haplotype).





PGGB: all chromosomes, layout with path-guided SGD












THE PREPRINT SERVER FOR BIOLOGY

<https://doi.org/10.1101/2023.09.22.558964>

New Results [Follow this preprint](#)

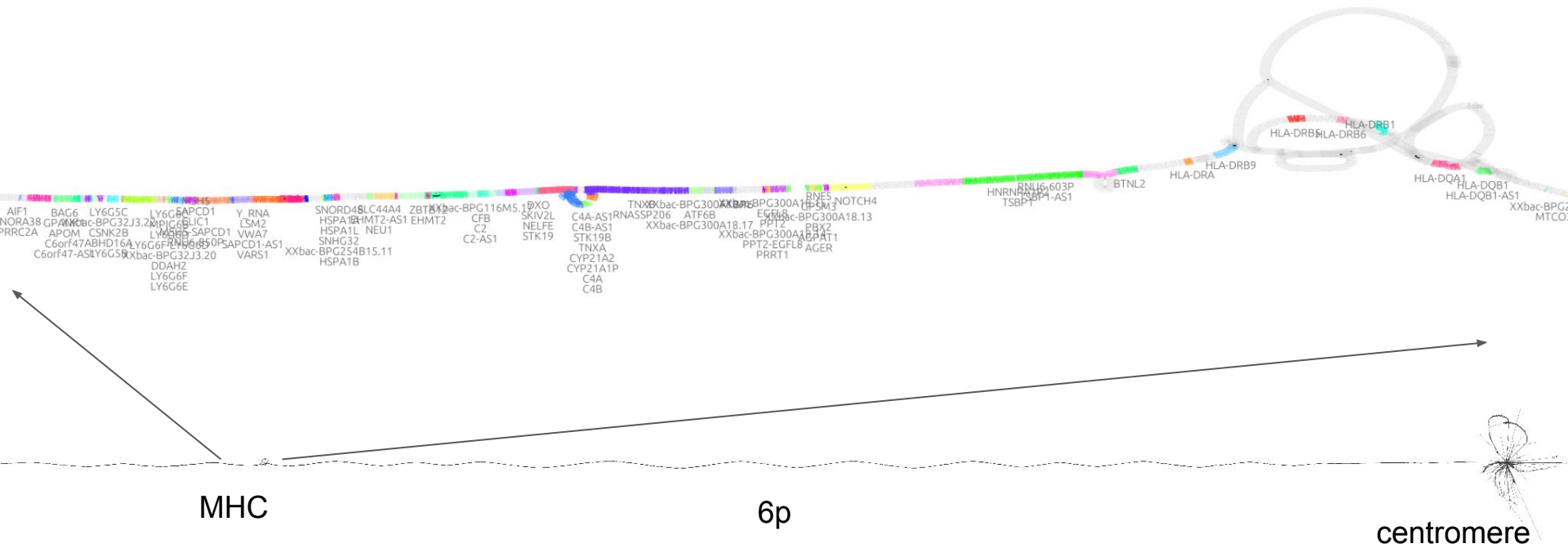
Pangenome graph layout by Path-Guided Stochastic Gradient Descent

 Simon Heumos,  Andrea Guarracino,  Jan-Niklas M. Schmelzle,  Jiajie Li,  Zhiru Zhang,
 Jörg Hagmann,  Sven Nahnsen,  Pjotr Prins,  Erik Garrison

Simon Heumos

C4A/B in pggb graph

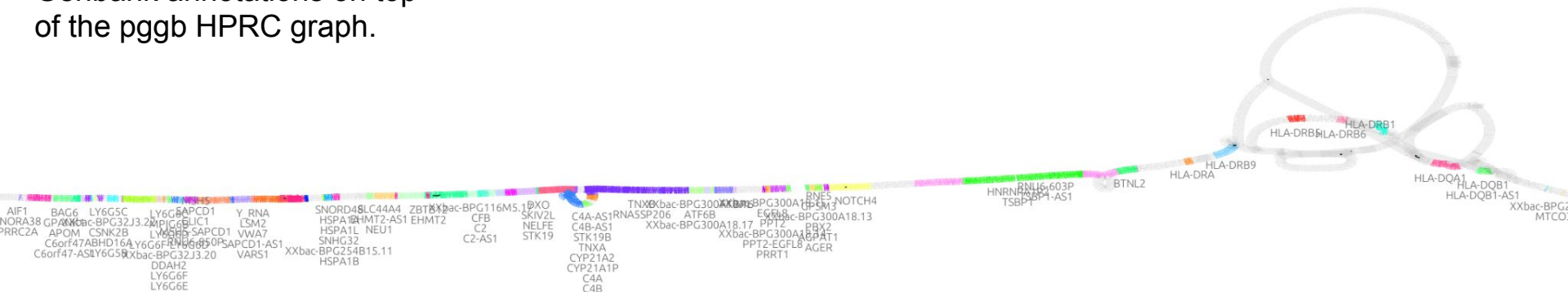
MHC class II



C4A/B in pggb graph

Genbank annotations on top
of the pggp HPRC graph.

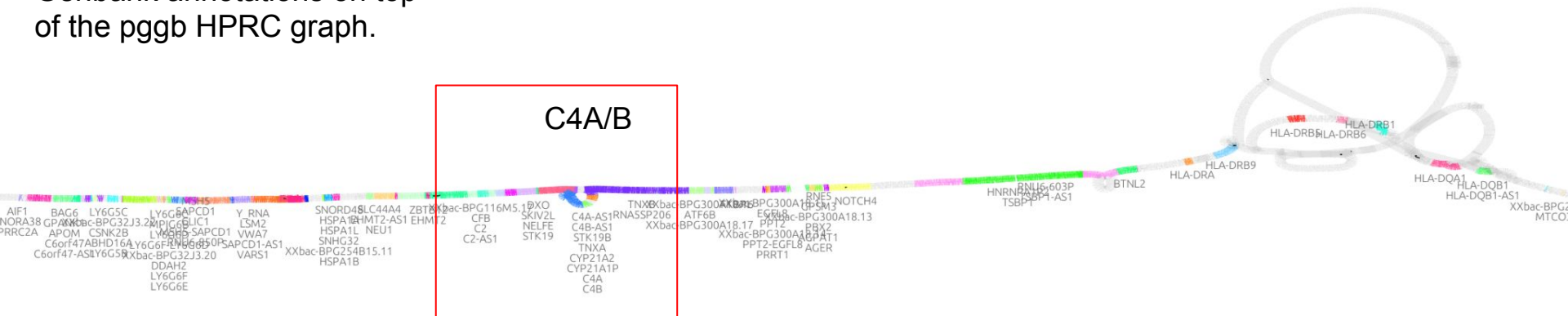
MHC class II



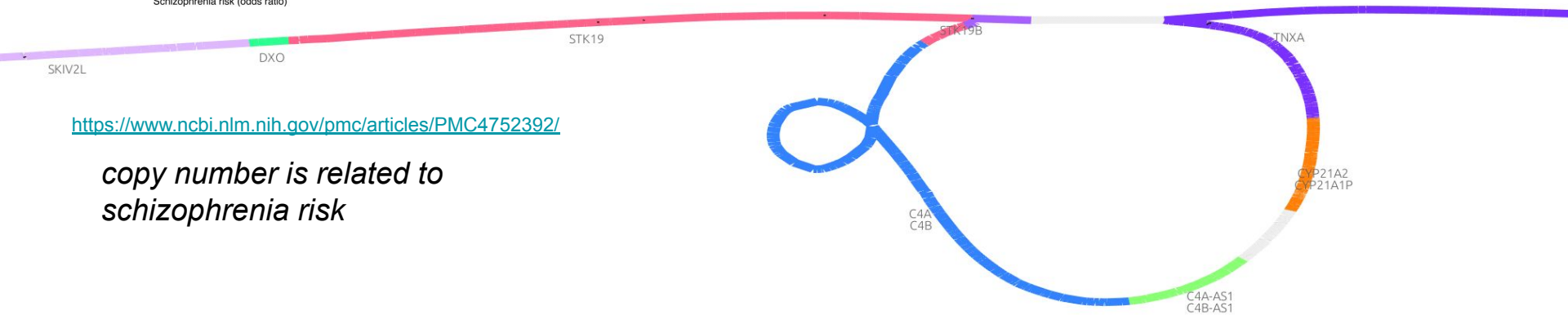
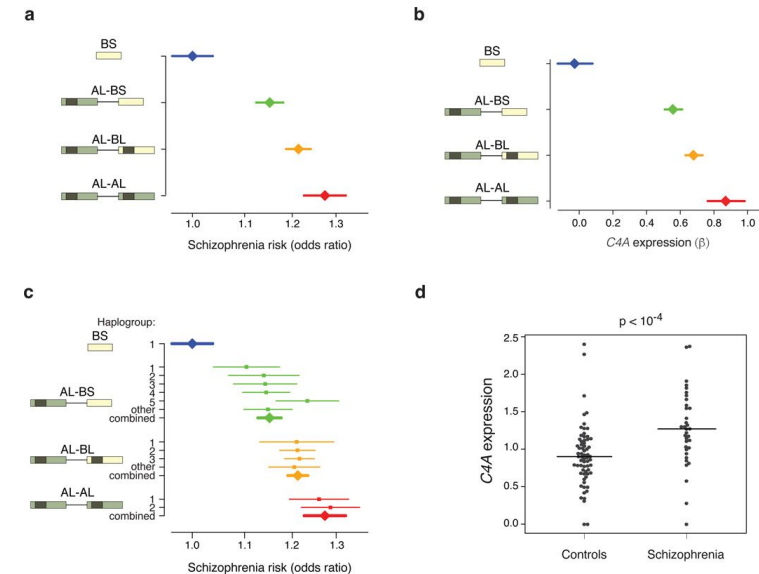
C4A/B in pggb graph

Genbank annotations on top
of the pggp HPRC graph.

MHC class II



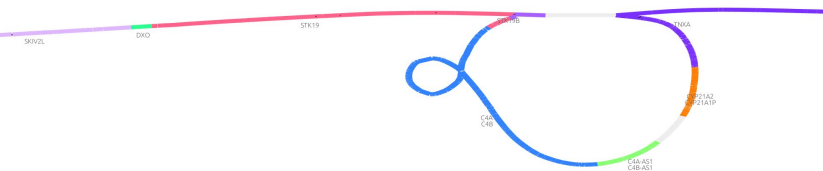
C4A/B in pggb graph



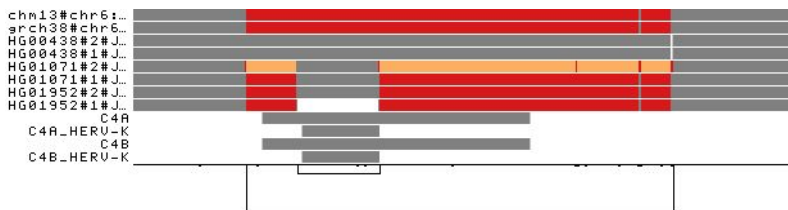
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4752392/>

*copy number is related to
schizophrenia risk*

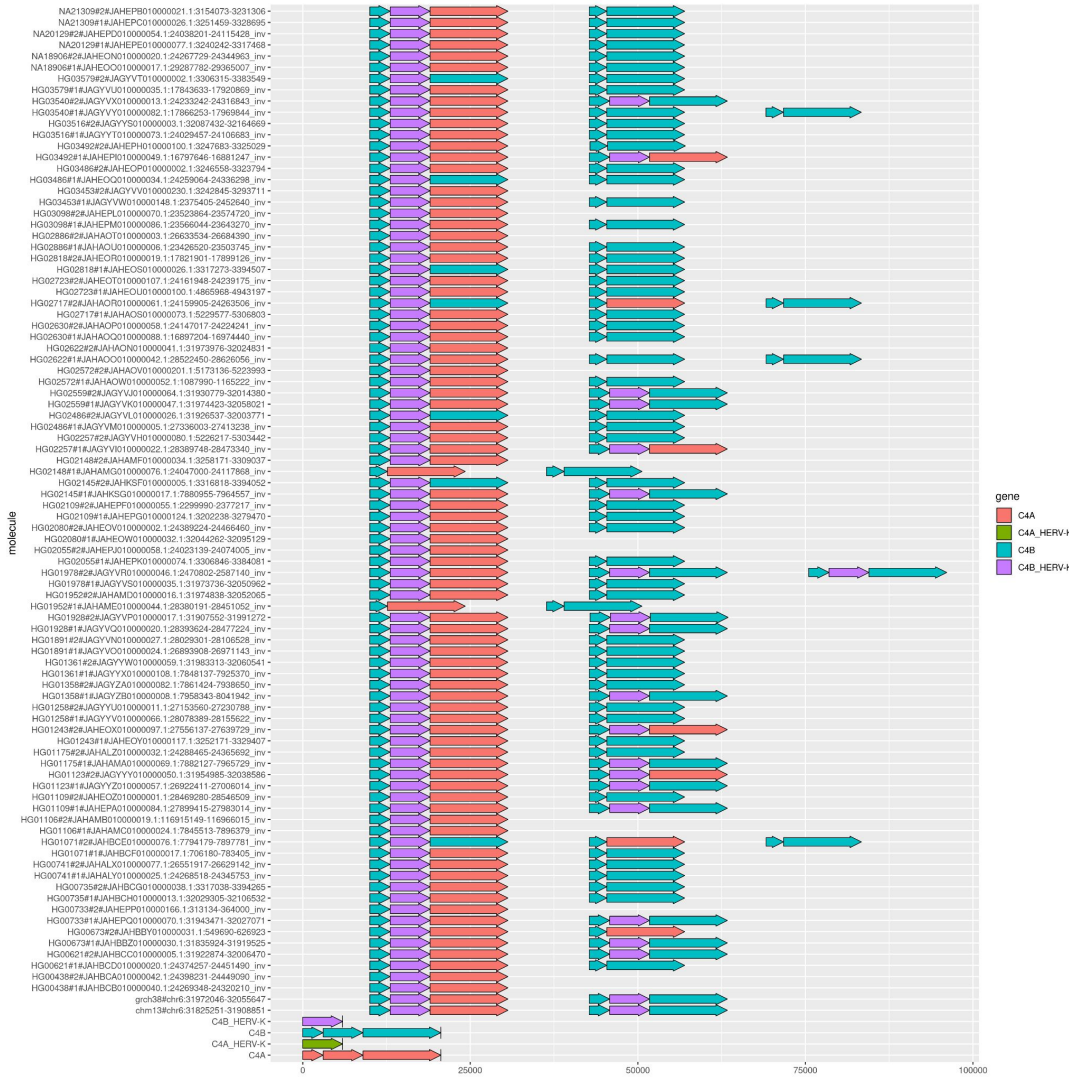
We learn that genome evolution is often nonlinear



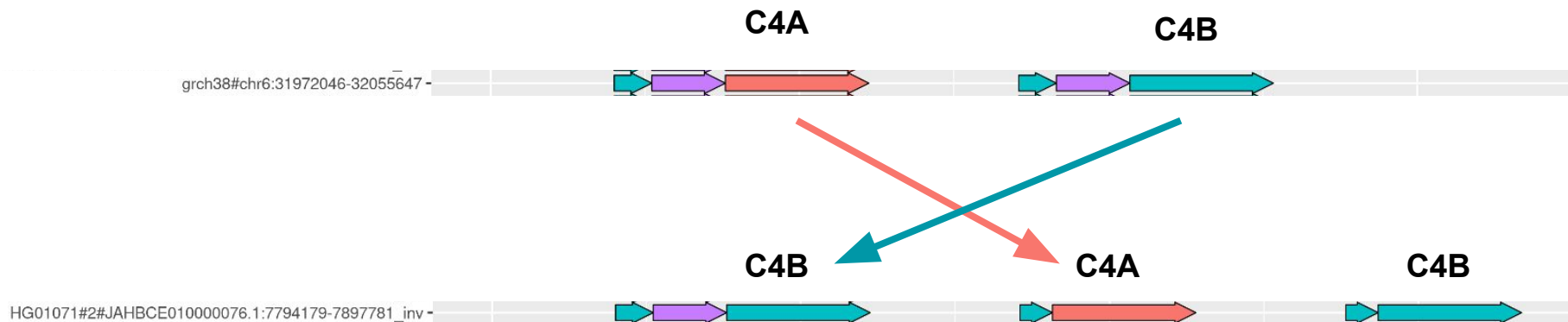
complement component 4 locus



Large SVs predominantly occur at VNTRs which are simply loops in our pggg graphs.



We learn that genome
evolution is often
nonlinear

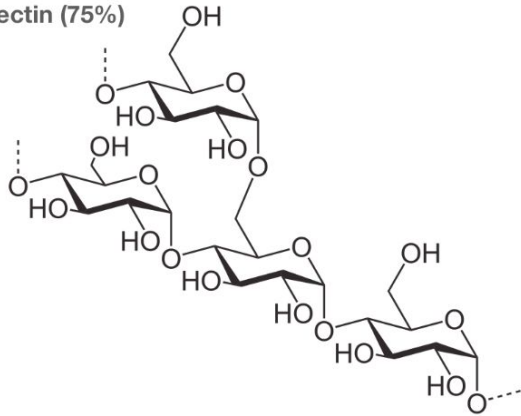


**The human
pangenome exposes
selection at the
amylase locus**

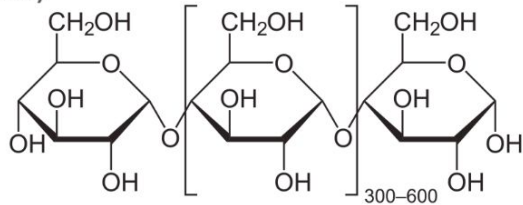
Amylase digests starch into sugar

Starch

amylopectin (75%)

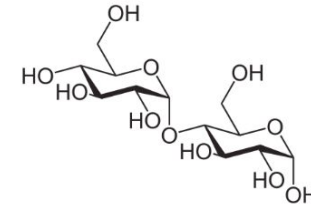


amylose (25%)

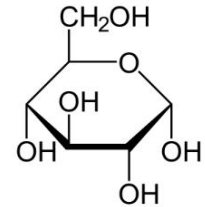


human salivary
 α amylase

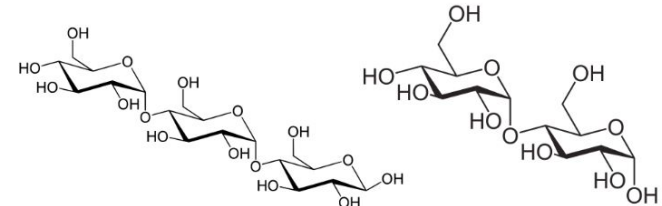
Sugar



Maltose



Glucose



Maltotriose

Maltose

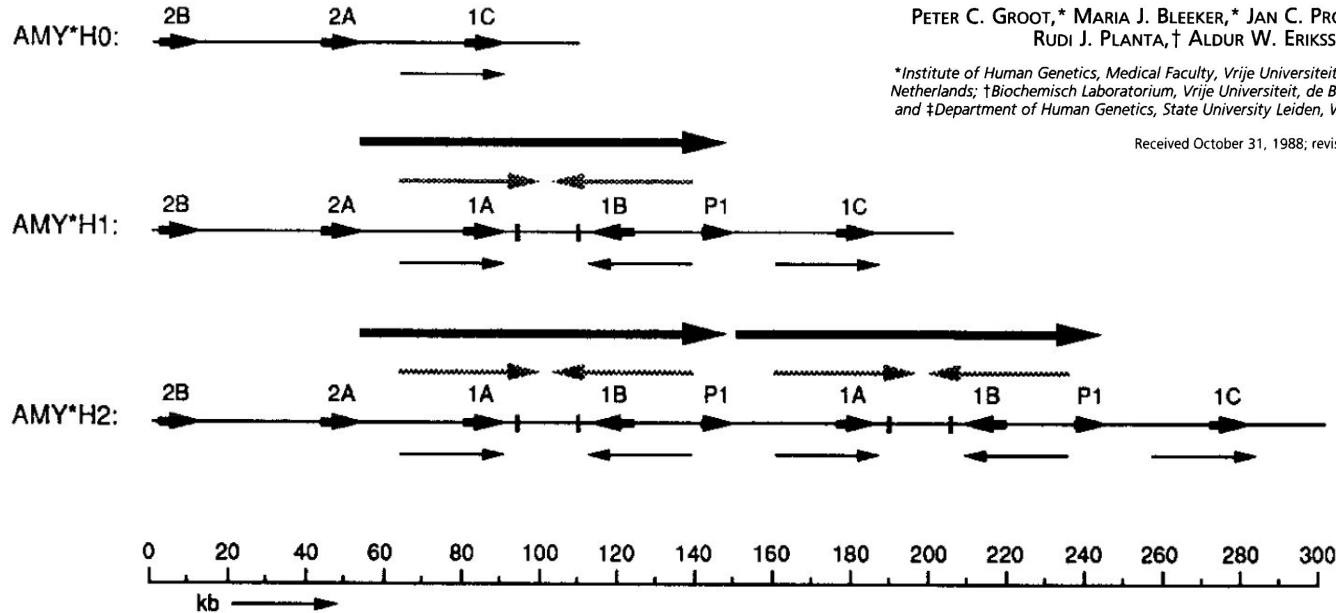
Amylase is a multi-copy gene family

The Human α -Amylase Multigene Family Consists of Haplotypes with Variable Numbers of Genes

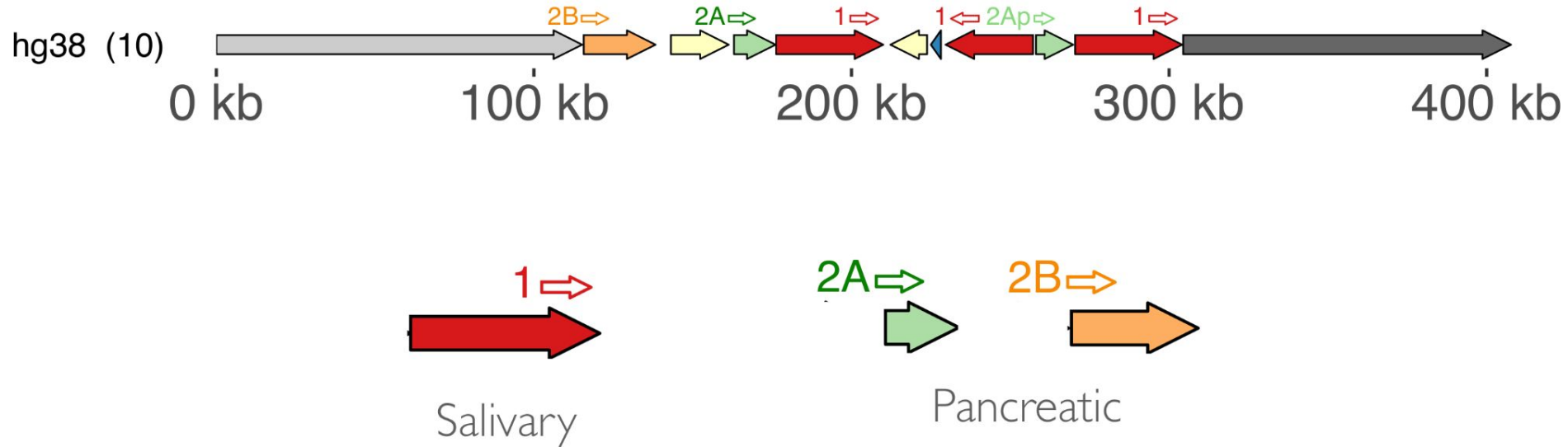
PETER C. GROOT,* MARIA J. BLEEKER,* JAN C. PRONK,* FRÉ ARWERT,* WILLEM H. MAGER,†
RUDI J. PLANTA,† ALDUR W. ERIKSSON,* AND RUNE R. FRANTS‡

**Institute of Human Genetics, Medical Faculty, Vrije Universiteit, van der Boechorststraat 7, 1081 BT, Amsterdam, The Netherlands; †Biochemisch Laboratorium, Vrije Universiteit, de Boelelaan 1083, 1081 HV, Amsterdam, The Netherlands; and ‡Department of Human Genetics, State University Leiden, Wassenaarseweg 72, 2333 AL, Leiden, The Netherlands*

Received October 31, 1988; revised February 10, 1989



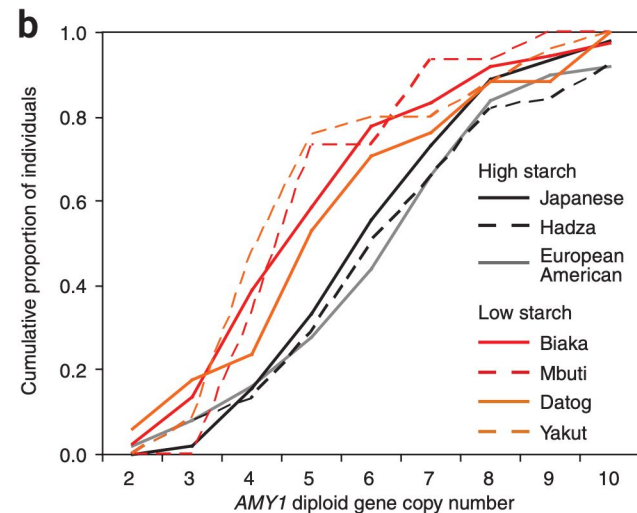
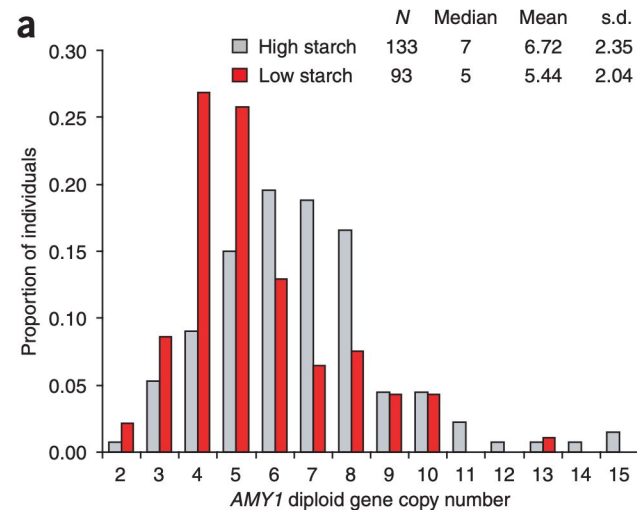
Amylase is a multi-copy gene family



Across human populations, diet correlates with amylase copy number

Diet and the evolution of human amylase gene copy number variation

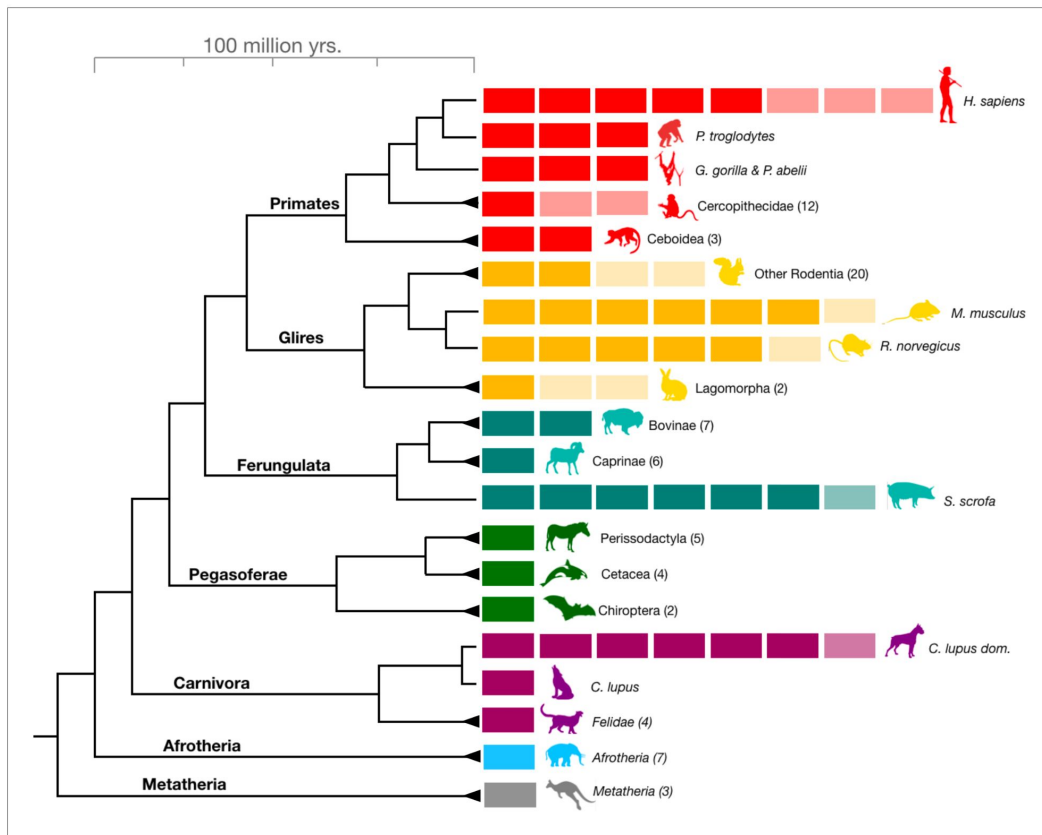
George H Perry^{1,2}, Nathaniel J Dominy³, Katrina G Claw^{1,4}, Arthur S Lee², Heike Fiegler⁵, Richard Redon⁵, John Werner⁴, Fernando A Villanea³, Joanna L Mountain⁶, Rajeev Misra⁴, Nigel P Carter⁵, Charles Lee^{2,7,8} & Anne C Stone^{1,8}



Across mammals, amylase copy correlates with diet

Independent amylase gene copy number bursts correlate with dietary preferences in mammals

Petar Pajic^{1,2}, Pavlos Pavlidis³, Kirsten Dean¹, Lubov Neznanova², Rose-Anne Romano², Danielle Garneau⁴, Erin Daugherty⁵, Anja Globig⁶, Stefan Ruhl^{2*}, Omer Gokcumen^{1*}



But, no evidence for selection in humans!

FADS1 and the Timing of Human Adaptation to Agriculture

Sara Mathieson¹ and Iain Mathieson^{*2}

¹Department of Computer Science, Swarthmore College, Swarthmore, PA

²Department of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA

***Corresponding author:** E-mail: mathi@pennmedicine.upenn.edu.

Associate editor: Evelyn Heyer

Abstract

Variation at the *FADS1*/*FADS2* gene cluster is functionally associated with differences in lipid metabolism and is often hypothesized to reflect adaptation to an agricultural diet. Here, we test the evidence for this relationship using both modern and ancient DNA data. We show that almost all the inhabitants of Europe carried the ancestral allele until the derived allele was introduced ~8,500 years ago by Early Neolithic farming populations. However, we also show that it was not under strong selection in these populations. We find that this allele, and other proposed agricultural adaptations at *LCT*/*MCM6* and *SLC22A4*, were not strongly selected until much later, perhaps as late as the Bronze Age. Similarly, increased copy number variation at the salivary amylase gene *AMY1* is not linked to the development of agriculture although, in this case, the putative adaptation precedes the agricultural transition. Our analysis shows that selection at the *FADS* locus was not tightly linked to the initial introduction of agriculture and the Neolithic transition. Further, it suggests that the strongest signals of recent human adaptation in Europe did not coincide with the Neolithic transition but with more recent changes in environment, diet, or efficiency of selection due to increases in effective population size.

Key words: Human evolution, selection, ancient DNA, agriculture, diet.

Selective sweep on human amylase genes postdates the split with Neanderthals

Charlotte E. Inchley¹, Cynthia D. A. Larbey¹, Nzar A. A. Shwan^{2,3}, Luca Pagani^{1,4}, Lauri Saag⁴, Tiago Antão⁵, Guy Jacobs⁶, Georgi Hudjashov^{4,7}, Ene Metspalu⁴, Mario Mitt^{8,9}, Christina A. Eichstaedt^{1,10}, Boris Malyarchuk¹¹, Miroslava Derenko¹¹, Joseph Wee¹², Syafiq Abdullah¹³, François-Xavier Ricaut¹⁴, Maru Mormina¹⁵, Reedik Mägi⁸, Richard Villems^{4,16,17}, Mait Metspalu⁴, Martin K. Jones¹, John A. L. Armour² & Toomas Kivisild^{2,4}

Discussion

In this study we have analysed genetic regions surrounding the human *AMY* cluster for evidence of natural selection and we have found: that human populations within and outside Africa are characterized by unusually low genetic diversity in the flanks of amylase locus relative to other genetic loci genome-wide; a young coalescent date postdating the human-Neanderthal population split; a significant Tajima's D signal in Africans; and the lack of strong signal of recent positive selection in human population groups we studied. These results are generally in line with Middle Pleistocene⁹ rather than Holocene¹ selection at the *AMY* locus although the significantly

And, conflicting GWAS results!

LETTERS

nature
genetics

Low copy number of the salivary amylase gene predisposes to obesity

Mario Falchi^{1,39,40}, Julia Sarah El-Sayed Moustafa^{1,39}, Petros Takousis¹, Francesco Pesce^{1,2}, Amélie Bonnefond³⁻⁶, Johanna C Andersson-Assarsson^{1,7,8}, Peter H Sudmant⁹, Rajkumar Dorajoo^{1,10}, Mashael Nedham Al-Shafai^{1,11}, Leonardo Bottolo¹², Erdal Ozdemir¹, Hon-Cheong So^{1,3}, Robert W Davies¹⁴, Alexandre Patrice^{6,15,16}, Robert Dent¹⁷, Massimo Mangino¹⁸, Pirro G Hysi¹⁸, Aurélie Dechaume^{3,4,6}, Marlène Huyvaert^{3,4,6}, Jane Skinner¹⁹, Marie Pigeyre^{4,6,15,16}, Robert Caiazzo^{4,6,15,16}, Violeta Raverdy^{6,15,16}, Emmanuel Vaillant^{3,4,6}, Sarah Field²⁰, Beverley Balkau^{21,22}, Michel Marre^{23,24}, Sophie Visvikis-Siest²⁵, Jacques Weill²⁶, Odile Poulain-Godefroy^{3,4,6}, Peter Jacobson^{7,8}, Lars Sjöström^{7,8}, Christopher J Hammond¹⁸, Panos Deloukas^{20,27,28}, Pak Chung Sham¹³, Ruth McPherson^{29,30}, Jeannette Lee³¹, E Shyong Tai³¹⁻³³, Robert Sladek³⁴⁻³⁶, Lena M S Carlsson^{7,8}, Andrew Walley^{1,37}, Evan E Eichler^{9,38}, Francois Pattou^{4,6,15,16}, Timothy D Spector^{18,40} & Philippe Froguel^{1,3-6,40}

2014

Low AMY associated with obesity

LETTERS

nature
genetics

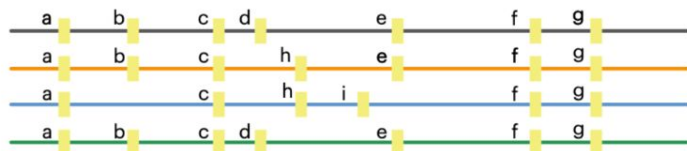
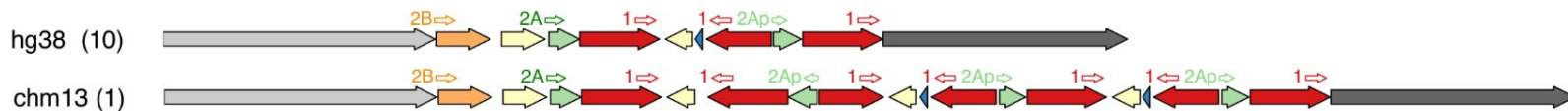
Structural forms of the human amylase locus and their relationships to SNPs, haplotypes and obesity

Christina L Usher¹, Robert E Handsaker¹⁻³, Tõnu Esko^{1,2,4-6}, Marcus A Tuke⁷, Michael N Weedon⁷, Alex R Hastie⁸, Han Cao⁸, Jennifer E Moon^{1,2,4,5}, Seva Kashin^{2,3}, Christian Fuchsberger^{9,10}, Andres Metspalu^{6,11}, Carlos N Pato¹², Michele T Pato¹², Mark I McCarthy¹³⁻¹⁵, Michael Boehnke^{9,10}, David M Altshuler^{1,2,16}, Timothy M Frayling⁷, Joel N Hirschhorn^{1,2,4,5} & Steven A McCarroll¹⁻³

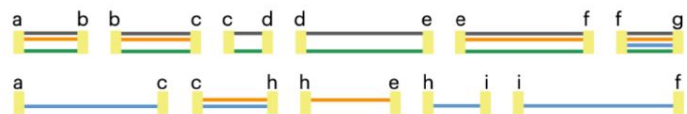
2015

Low AMY **NOT** associated with obesity

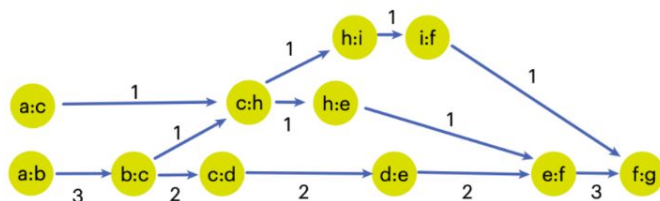
Human amylase copy number diversity



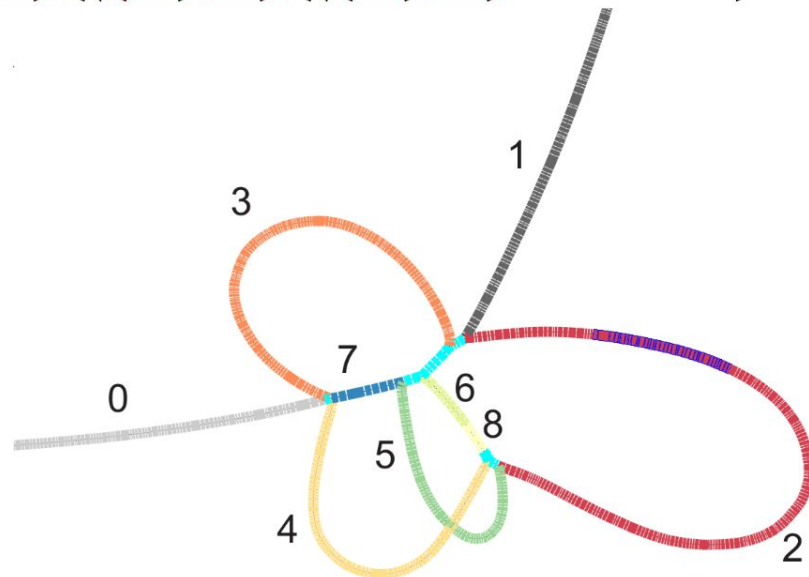
Sequence with
minimizer anchors



Graph vertex:
A set of sequences
with shared minimizer
anchors at both ends

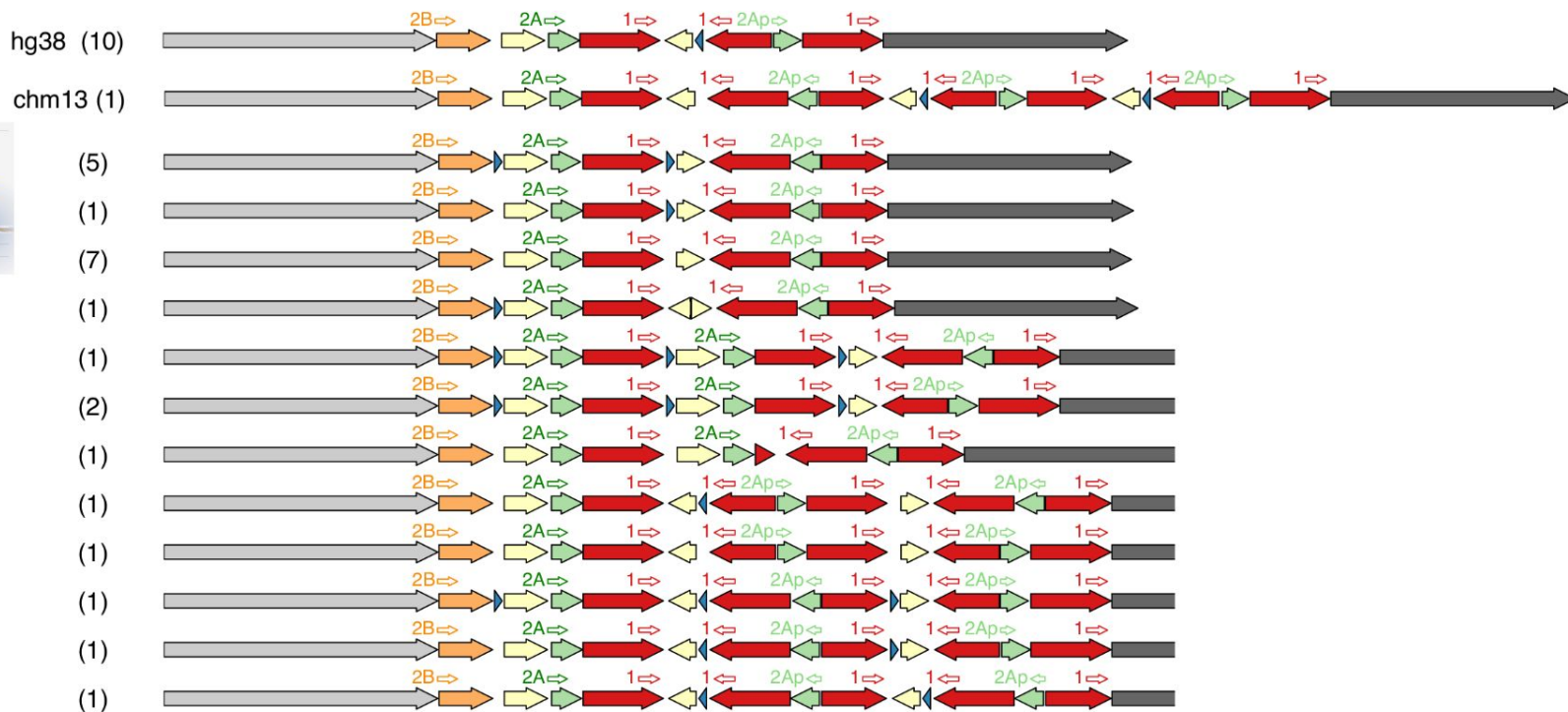


Induced
MAP-graph



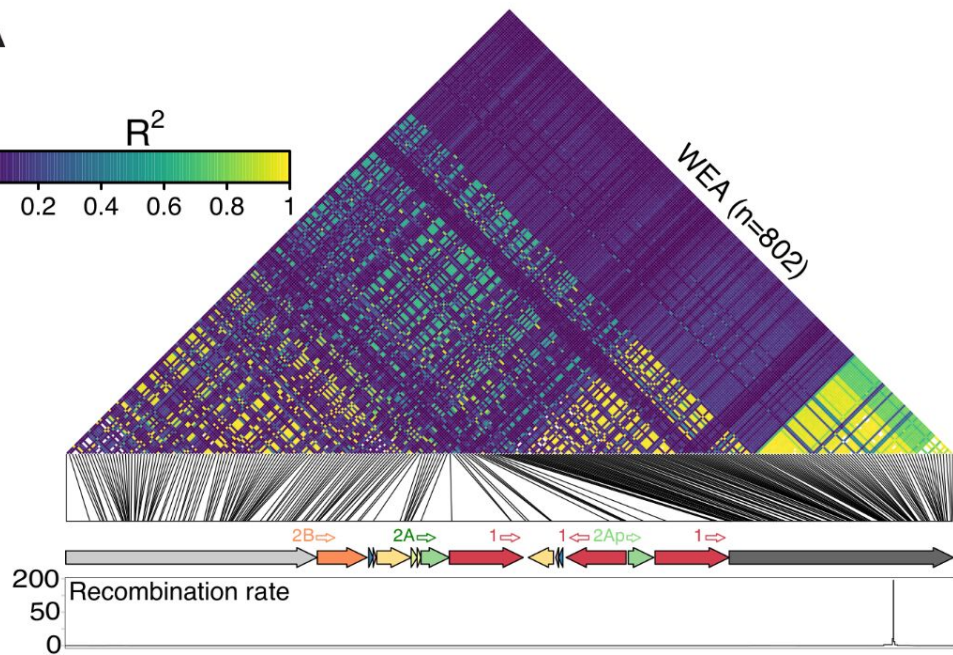
MAP graph (PGR-TK)

<https://doi.org/10.1038/s41592-023-01914-y>

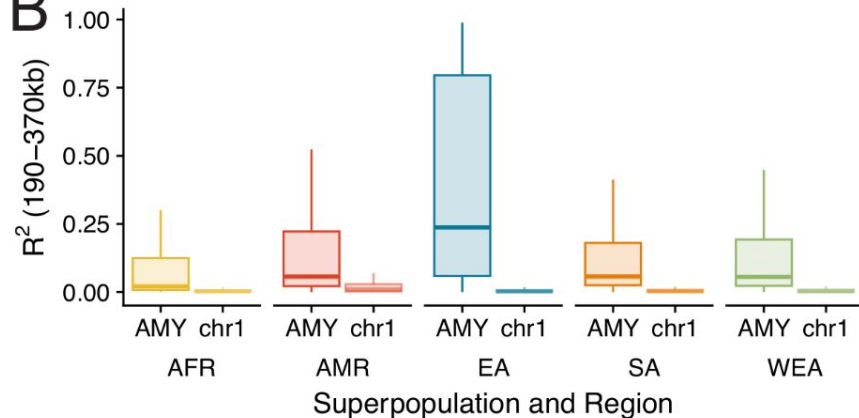


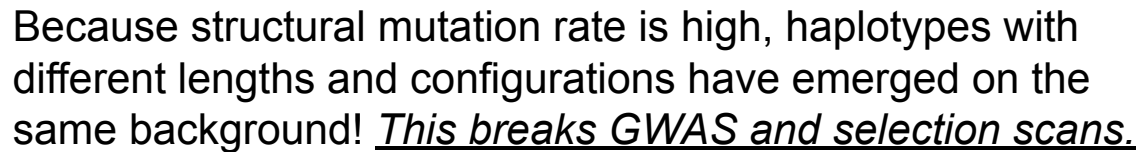
Strong linkage disequilibrium block across AMY locus

A

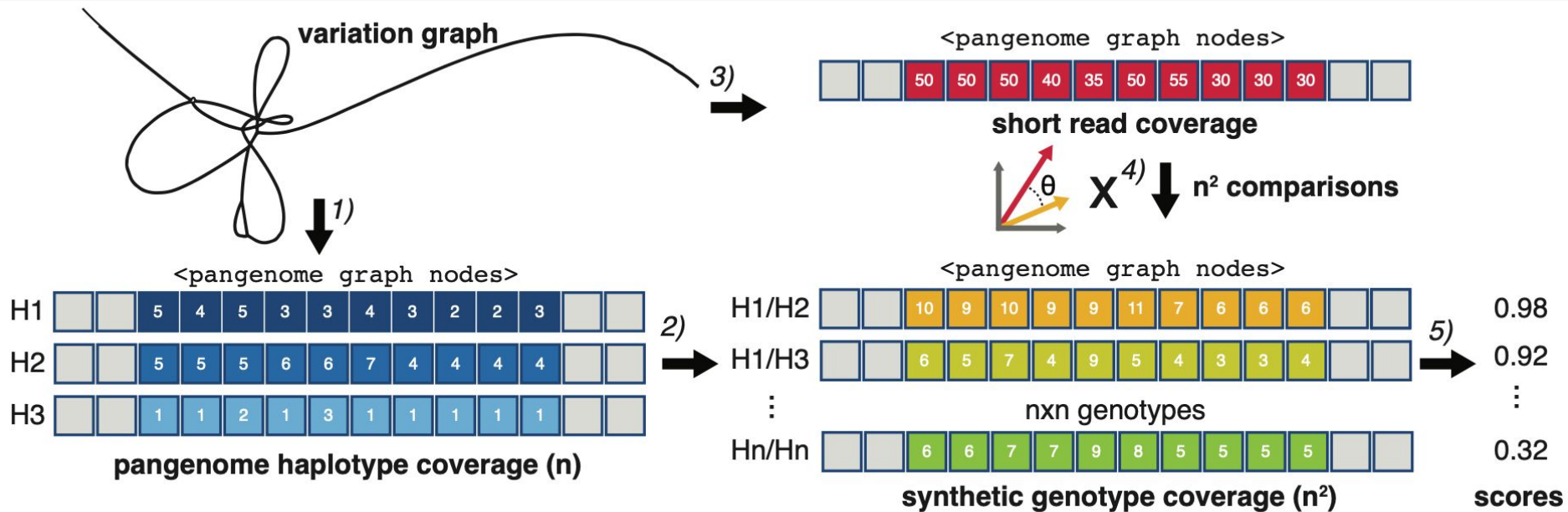


B



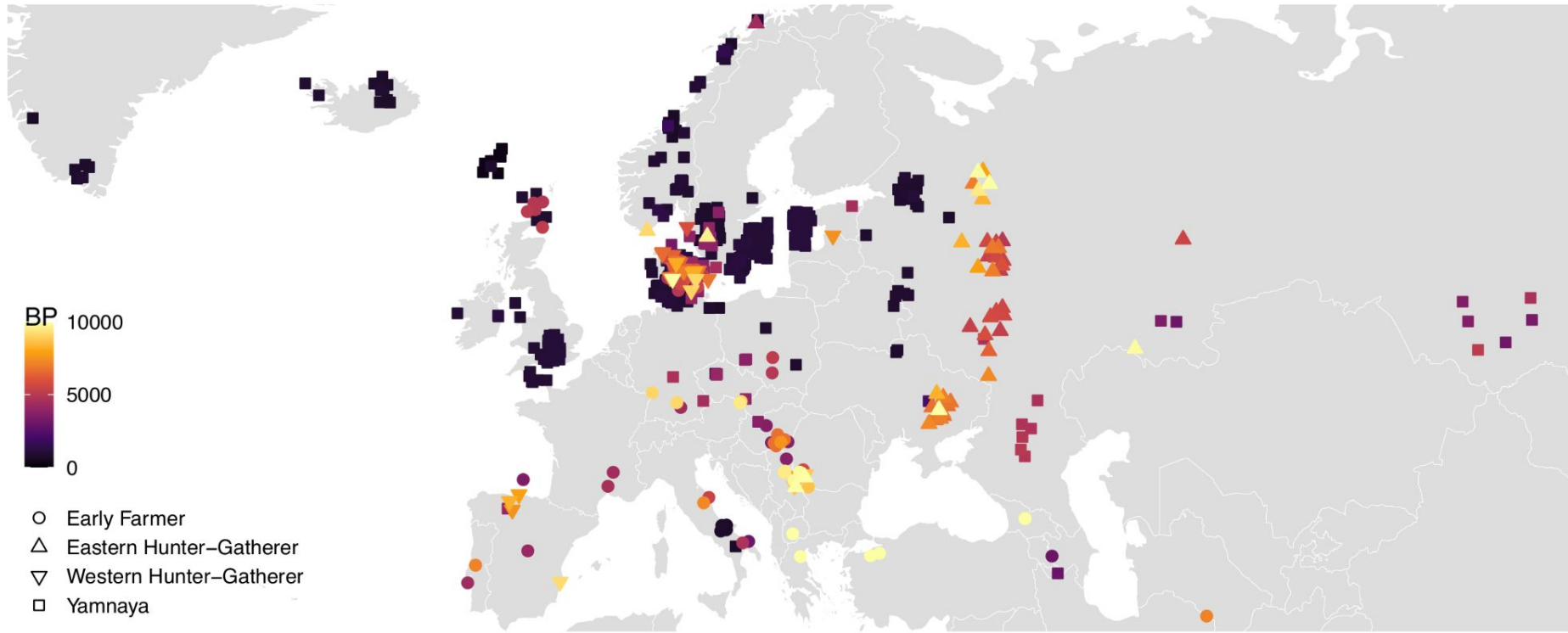


Deconvolving haplotypes from short reads

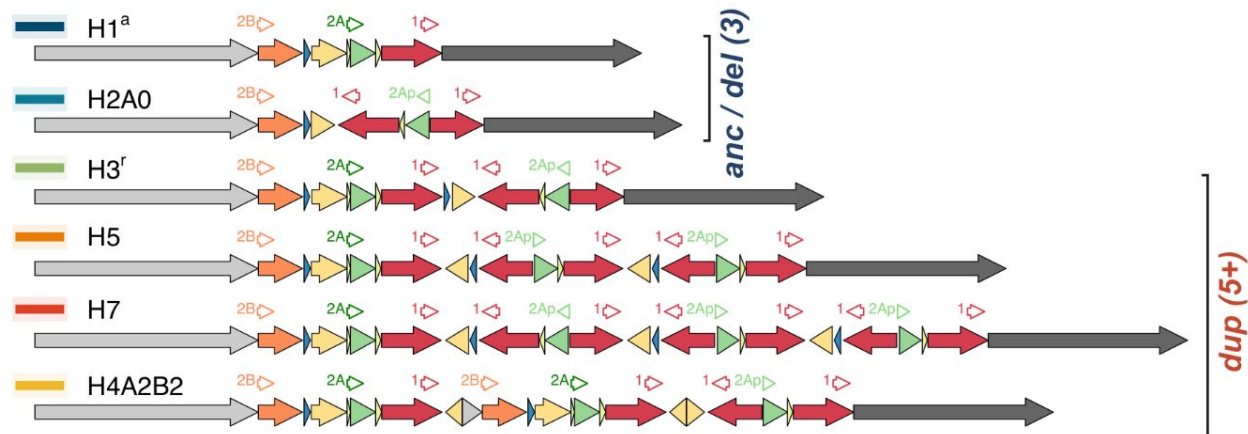
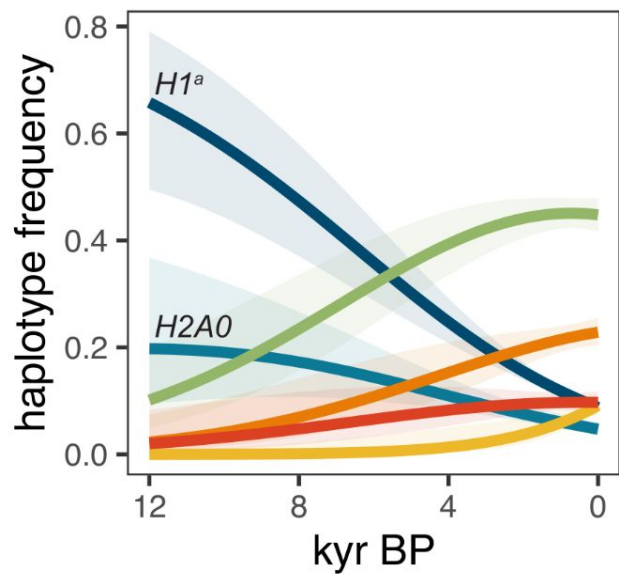


cross-validation with ddPCR, copy number, and hold-one-out experiments shows ~95% accuracy

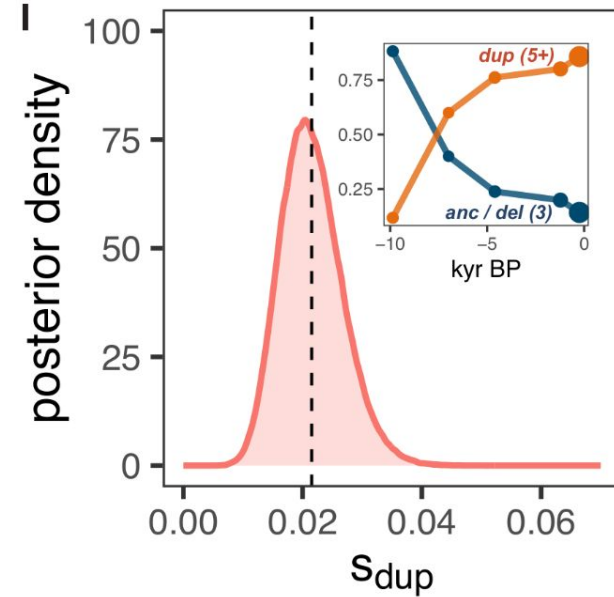
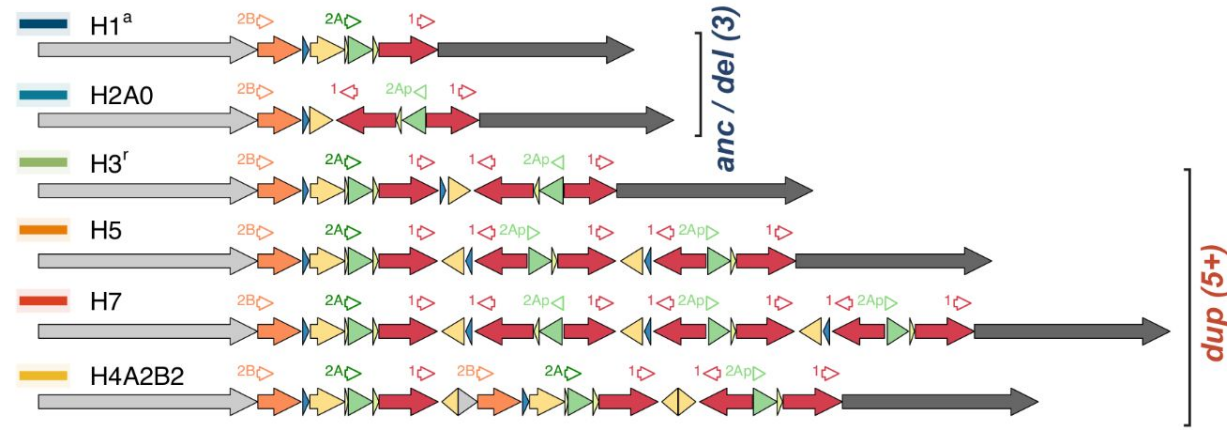
Recent evolution from 534 ancient European genomes



Recent evolution of human amylase copy number diversity



Evidence for selection of high-copy amylase haplotypes



There *is* selection at amylase in humans for haplotypes with more AMY1 copies.

Selection coefficient of 0.02 is equivalent to selection at lactase!

**The human
pangenome explains
acrocentric evolution**



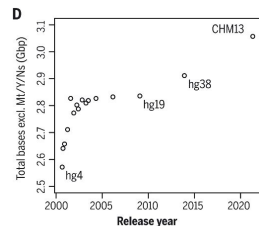
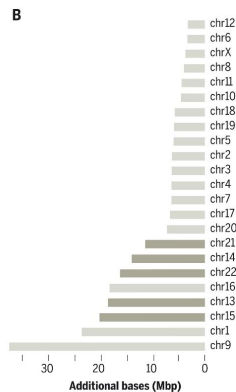
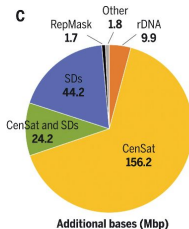
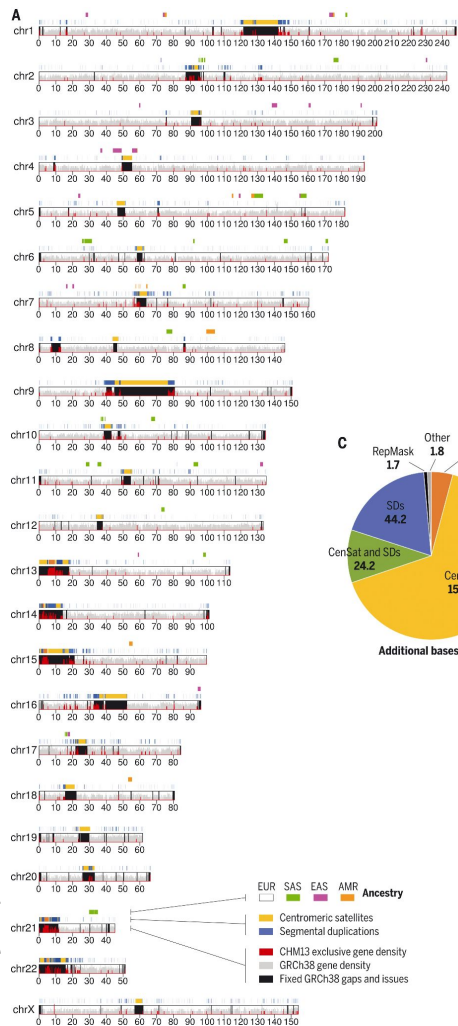
<https://doi.org/10.1126/science.abj6987>

The complete sequence of a human genome

SERGEY NURK [ID](#), SERGEY KOREN [ID](#), ARANG RHIE [ID](#), MIKKO RAUTIAINEN [ID](#), ANDREY V. BZIKADZE [ID](#), ALLA MIKHEENKO, MITCHELL R. VOLLGER [ID](#),
NICOLAS ALTEMOSE [ID](#), LEV URALSKY [ID](#), ARIEL GERSHMAN [ID](#), SERGEY AGANEZOV [ID](#), SAVANNAH J. HOYT [ID](#), MARK DIEKHANS [ID](#), GLENNIS A. LOGSDON [ID](#),
MICHAEL ALONGE [ID](#), STYLIANOS E. ANTONARAKIS [ID](#), MATTHEW BORCHERS [ID](#), GERARD G. BOUFFARD [ID](#), SHELISE Y. BROOKS, GINA V. CALDAS, NAE-CHYUN CHEN
[ID](#), HAORYU CHENG [ID](#), CHEN-SHAN CHIN [ID](#), WILLIAM CHOW [ID](#), LEONARDO G. DE LIMA [ID](#), PHILIP C. DISHUCK [ID](#), RICHARD DURBIN [ID](#), TATIANA DVORKINA,
IAN T. FIDDES [ID](#), GIULIO FORMENTI [ID](#), ROBERT S. FULTON, ARKARACHAI FUNGTAMMASAN [ID](#), ERIK GARRISON [ID](#), PATRICK G. S. GRADY [ID](#), TINA A. GRAVES-LINDSAY
[ID](#), IRA M. HALL [ID](#), NANCY F. HANSEN [ID](#), GABRIELLE A. HARTLEY, MARINA HAUKNES [ID](#), KERSTIN HOWE [ID](#), MICHAEL W. HUNKAPILLER, CHIRAG JAIN, MITEN JAIN
[ID](#), ERICH D. JARVIS [ID](#), PETER KERPEDJIEV, MELANIE KIRSCH [ID](#), MIKHAIL KOLMOGOROV [ID](#), JONAS KORLACH [ID](#), MILINN KREMITZKI [ID](#), HENG LI [ID](#),
VALERIE V. MADURO [ID](#), TOBIAS MARSHALL [ID](#), ANN M. MCCARTNEY, JENNIFER MCDANIEL [ID](#), DANNY E. MILLER [ID](#), JAMES C. MULLIKIN [ID](#), EUGENE W. MYERS [ID](#),
NATHAN D. OLSON [ID](#), BENEDICT PATEN [ID](#), PAUL PELUSO, PAVEL A. PEVZNER [ID](#), DAVID PORUBSKY [ID](#), TAMARA POTAPOVA [ID](#), EVGENY I. ROGAEV,
JEFFREY A. ROSENFELD [ID](#), STEVEN L. SALZBERG [ID](#), VALERIE A. SCHNEIDER, FRITZ J. SEDLAZECK [ID](#), KISHWAR SHAFIN [ID](#), COLIN J. SHEW, ALAINA SHUMATE [ID](#),
YING SIMS, ARIAN F. A. SMIT [ID](#), DANIELA C. SOTO [ID](#), IVAN SOVIĆ [ID](#), JESSICA M. STORER [ID](#), AARON STREETS [ID](#), BETH A. SULLIVAN [ID](#),
FRANÇOISE THIBAUD-NISSEN [ID](#), JAMES TORRANCE [ID](#), JUSTIN WAGNER, BRIAN P. WALENZ [ID](#), AARON WENGER [ID](#), JONATHAN M. D. WOOD [ID](#), CHUNLIN XIAO [ID](#),
STEPHANIE M. YAN [ID](#), ALICE C. YOUNG [ID](#), SAMANTHA ZARATE [ID](#), URVASHI SURTI, RAJIV C. MCCOY [ID](#), MEGAN Y. DENNIS [ID](#), IVAN A. ALEXANDROV [ID](#),
JENNIFER L. GERTON [ID](#), RACHEL J. O'NEILL [ID](#), WINSTON TIMP [ID](#), JUSTIN M. ZOOK [ID](#), MICHAEL C. SCHATZ [ID](#), EVAN E. EICHLER [ID](#), KAREN H. MIGA [ID](#),
AND ADAM M. PHILLIPPY [ID](#) [fewer](#) [Authors Info & Affiliations](#)

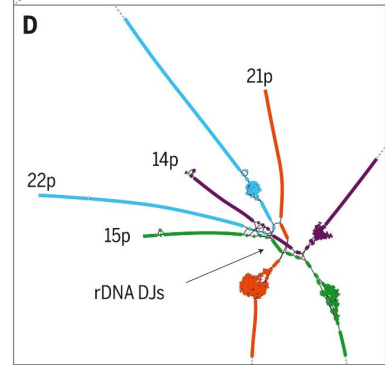
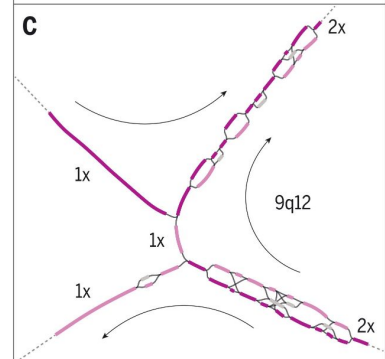
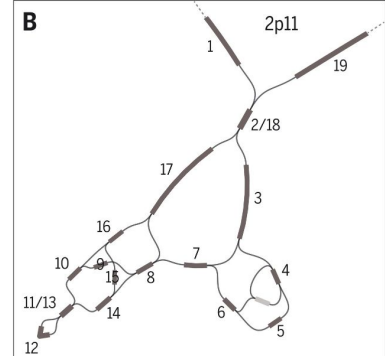
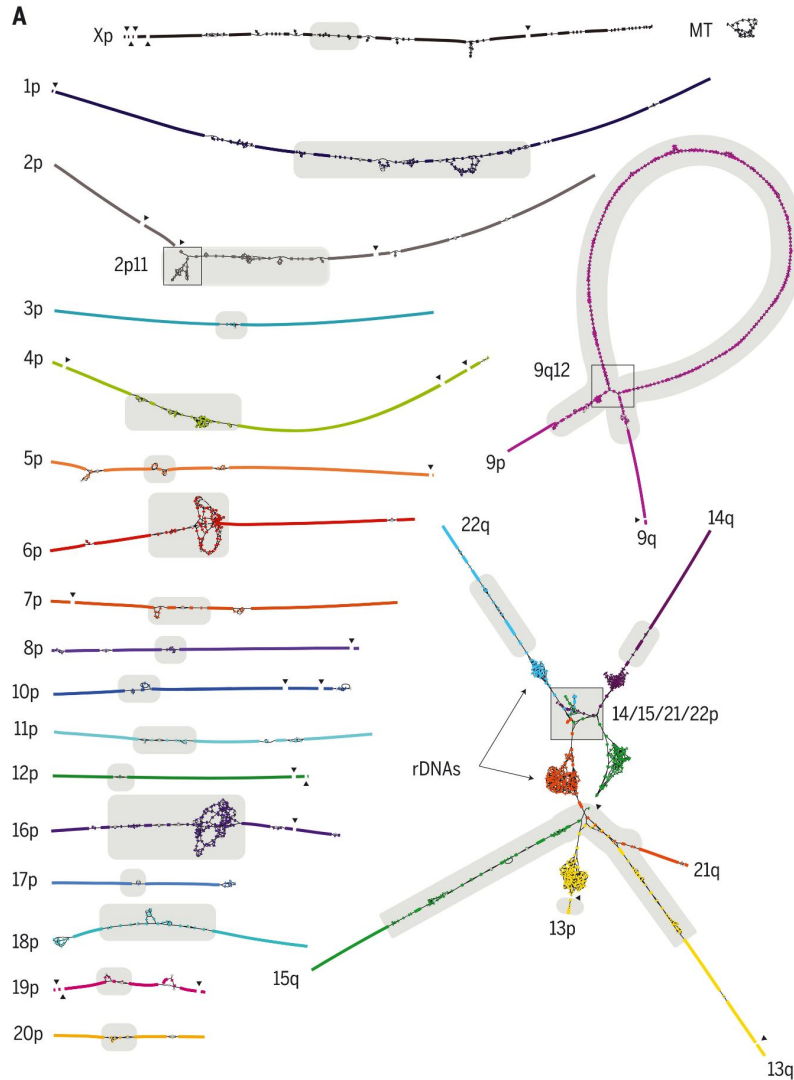
T2T-CHM13 fills 8% of the reference which was incomplete

All of the acrocentric
p-arms were
assembled for the
first time!

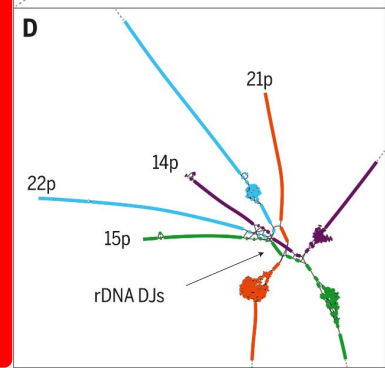
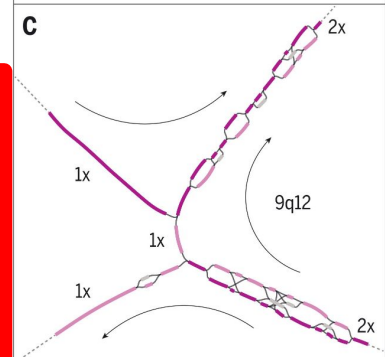
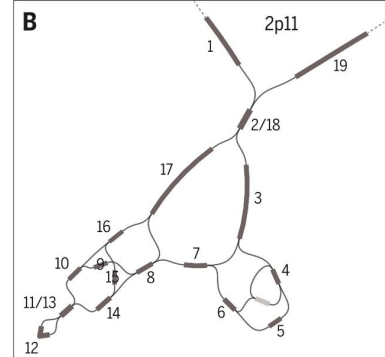
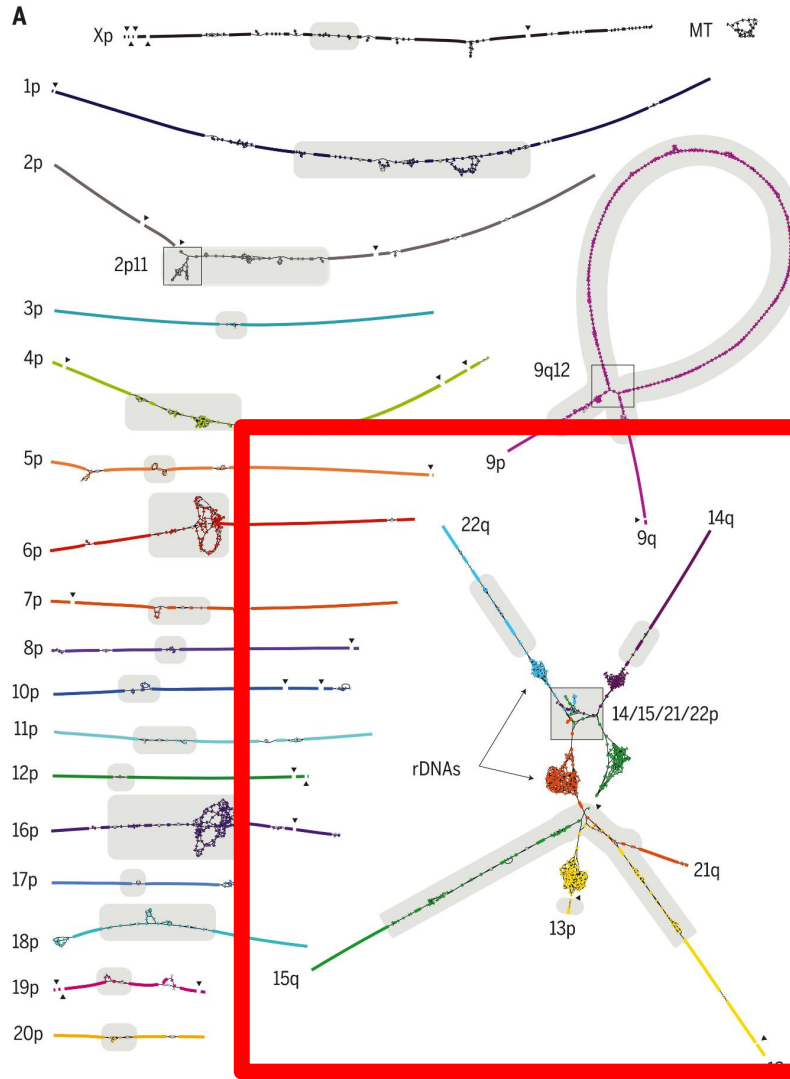


acrocentric: a chromosome where the centromere is almost at one end. In humans the short arms of the acrocentrics are the location of ribosomal DNA and organize the nucleoli.

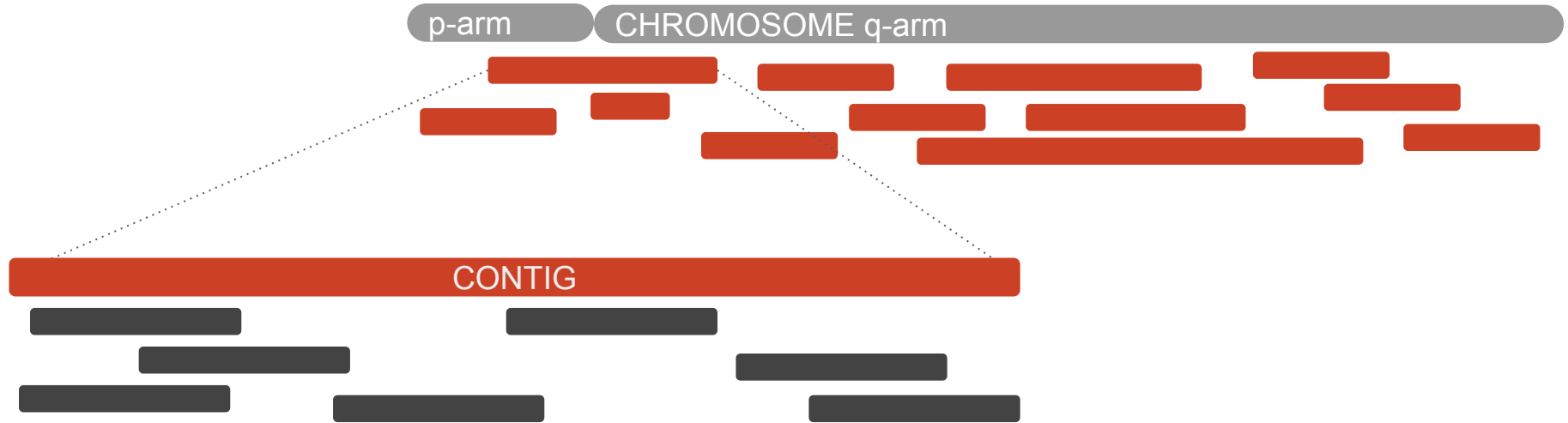
Revealing new mysteries...



Revealing new mysteries...

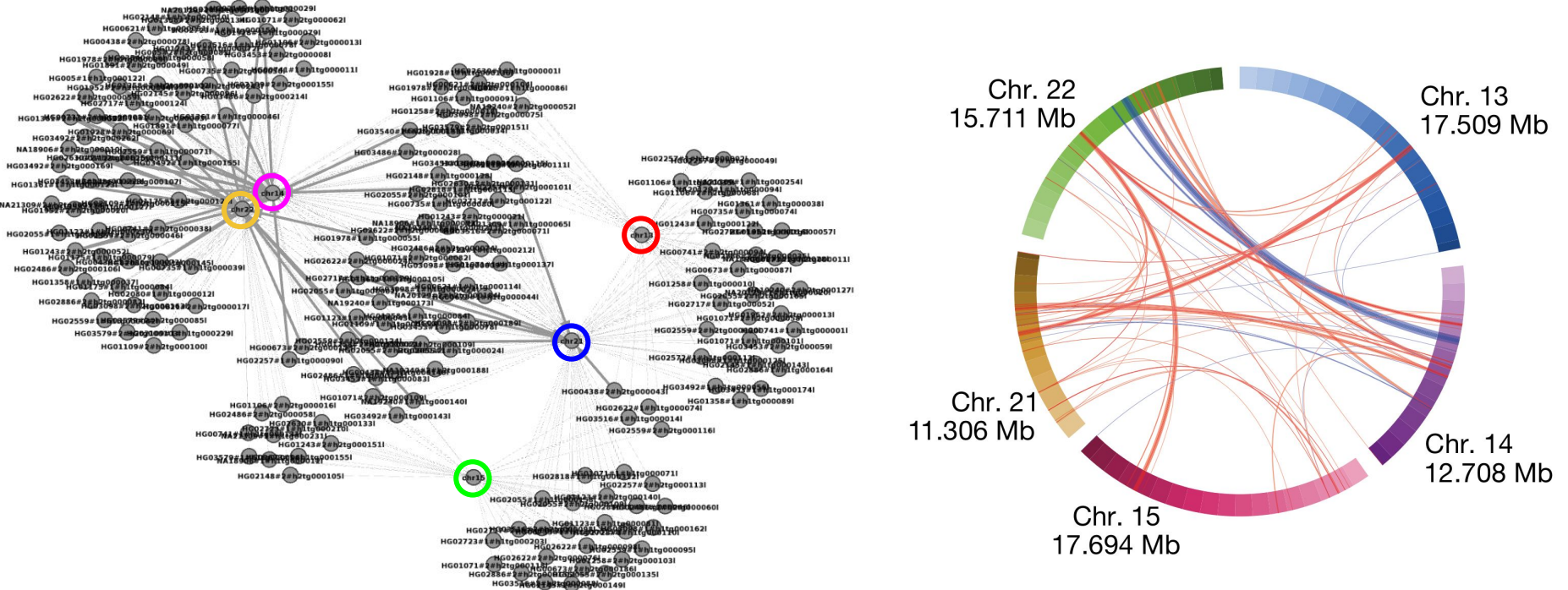


Assembly jargon: *Contigs*



```
contigs_are_assembled_from_sequencing_reads
ntigs_ar          d_from_sequenc
gs_are_c          cing_reads
contig            _assemble    uencing_r
```


HPRC acrocentric “misjoins”

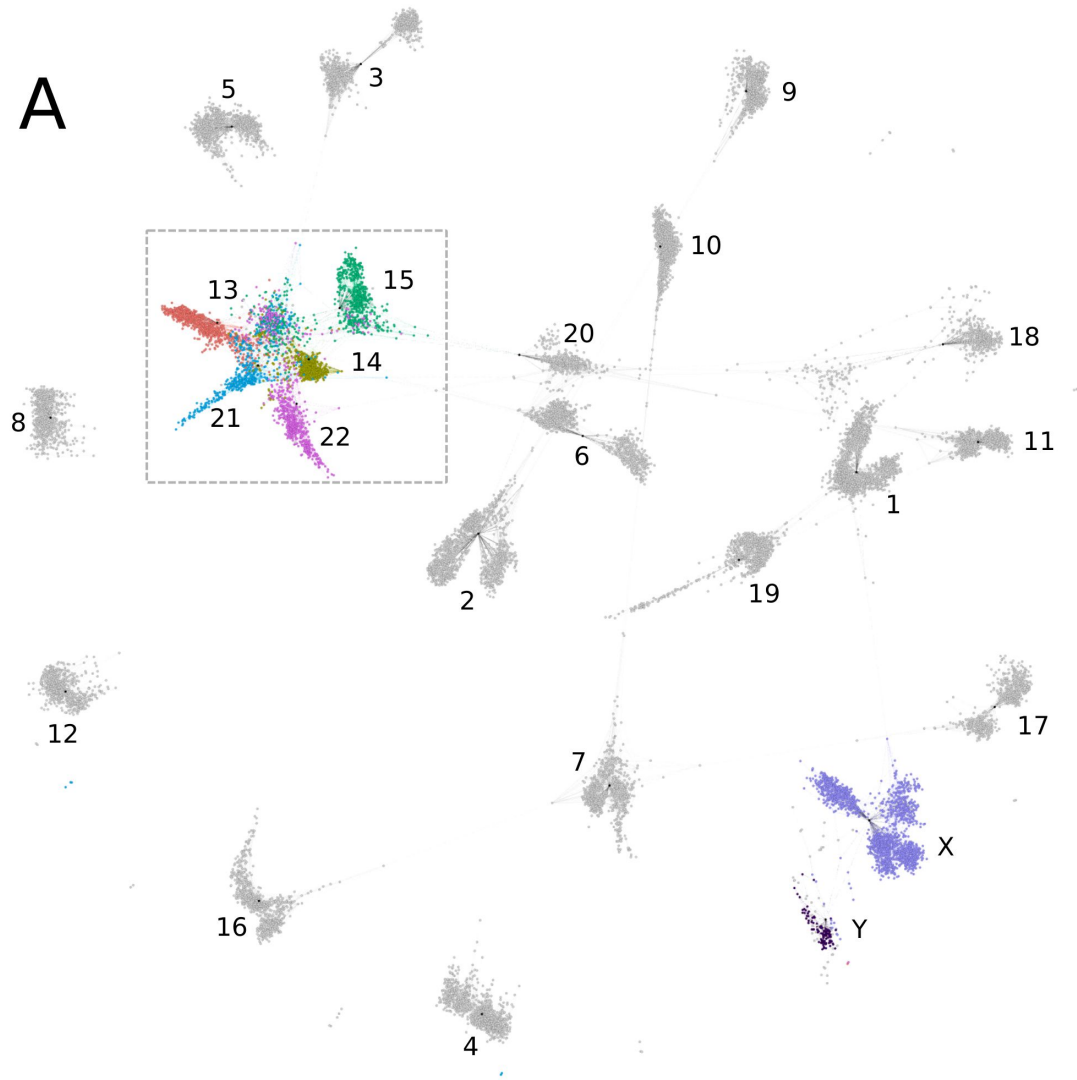


Alignment graph of the misjoins. Every node is a contig and every edge represents the number of mapping between nodes. Alignment graph obtained with [pafnet](#) and visualized with [gephi](#). Color code: **chr13**, **chr14**, **chr15**, **chr21**, **chr22**.

From HPRC main paper.

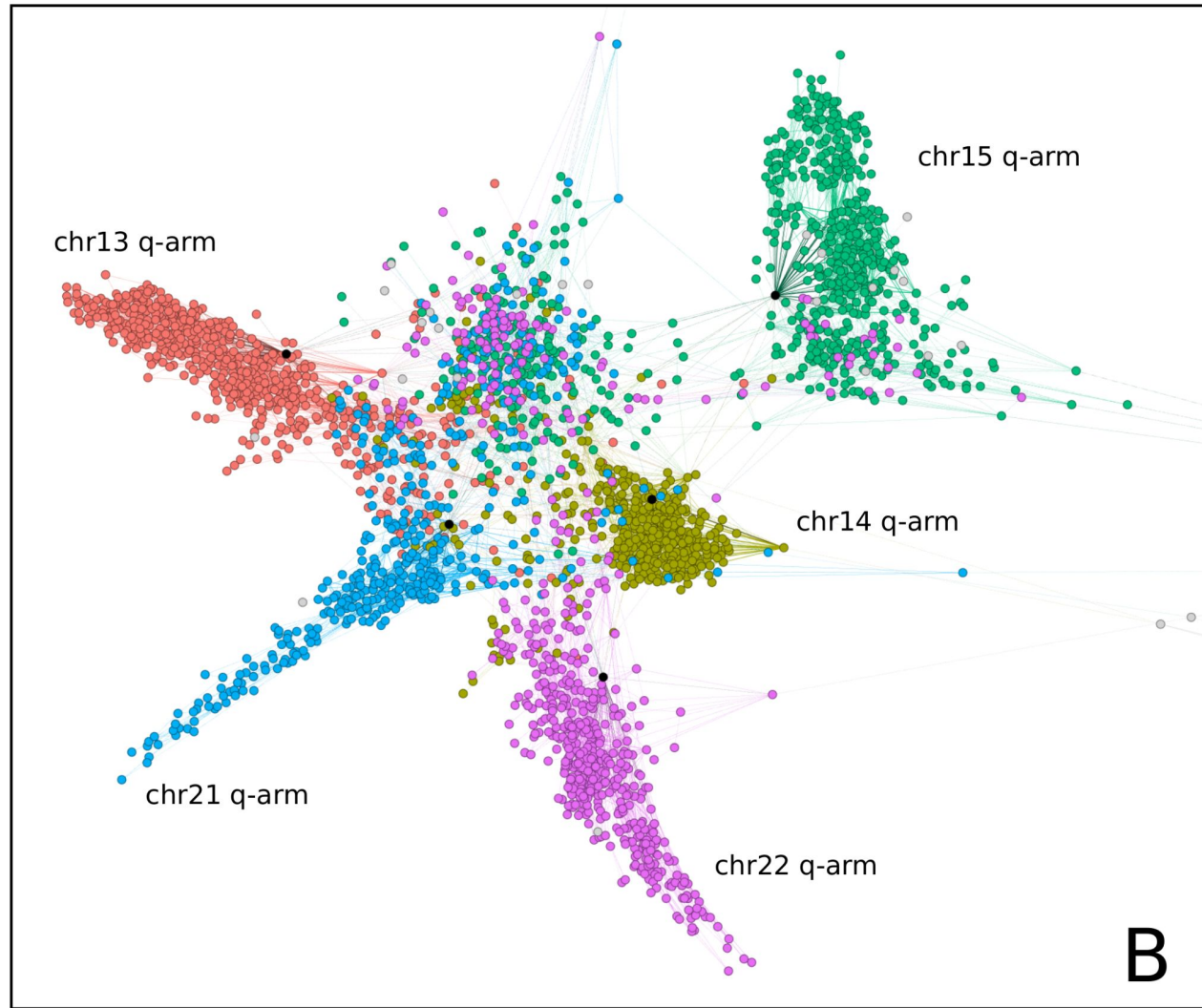
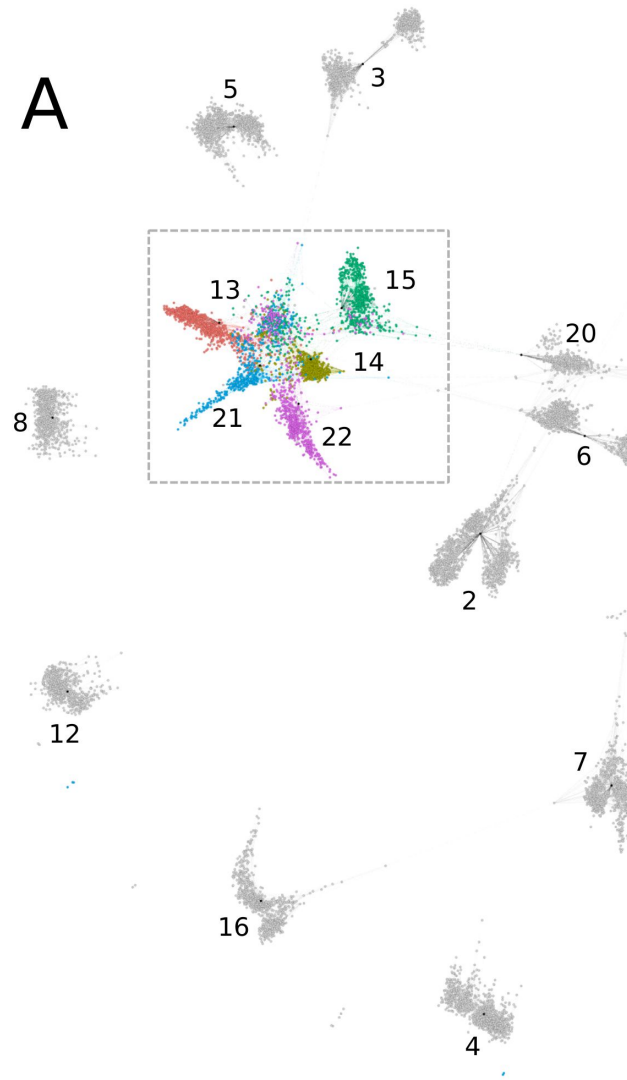
* With the exception of assembly errors in one haplotype (HG02080 paternal).

A

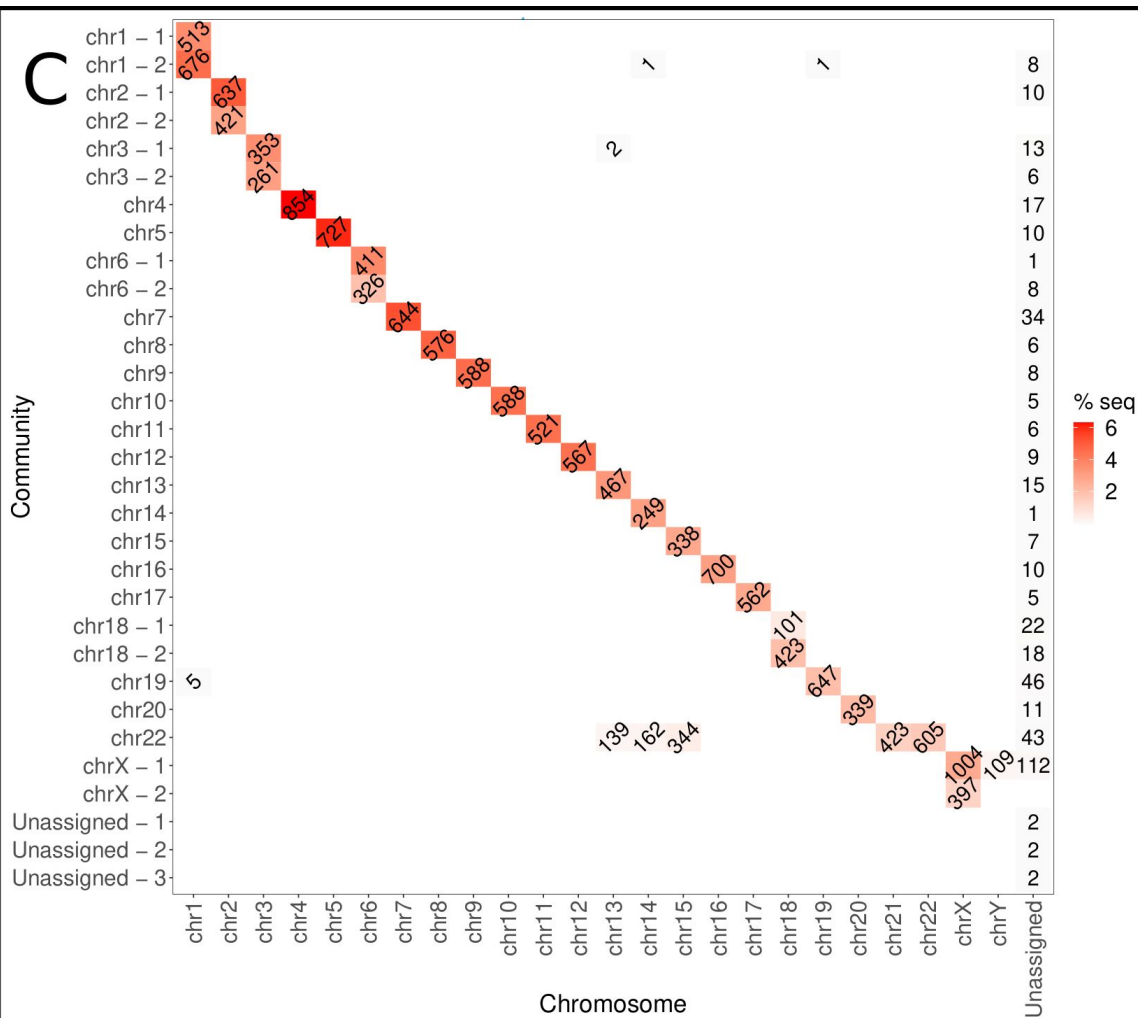


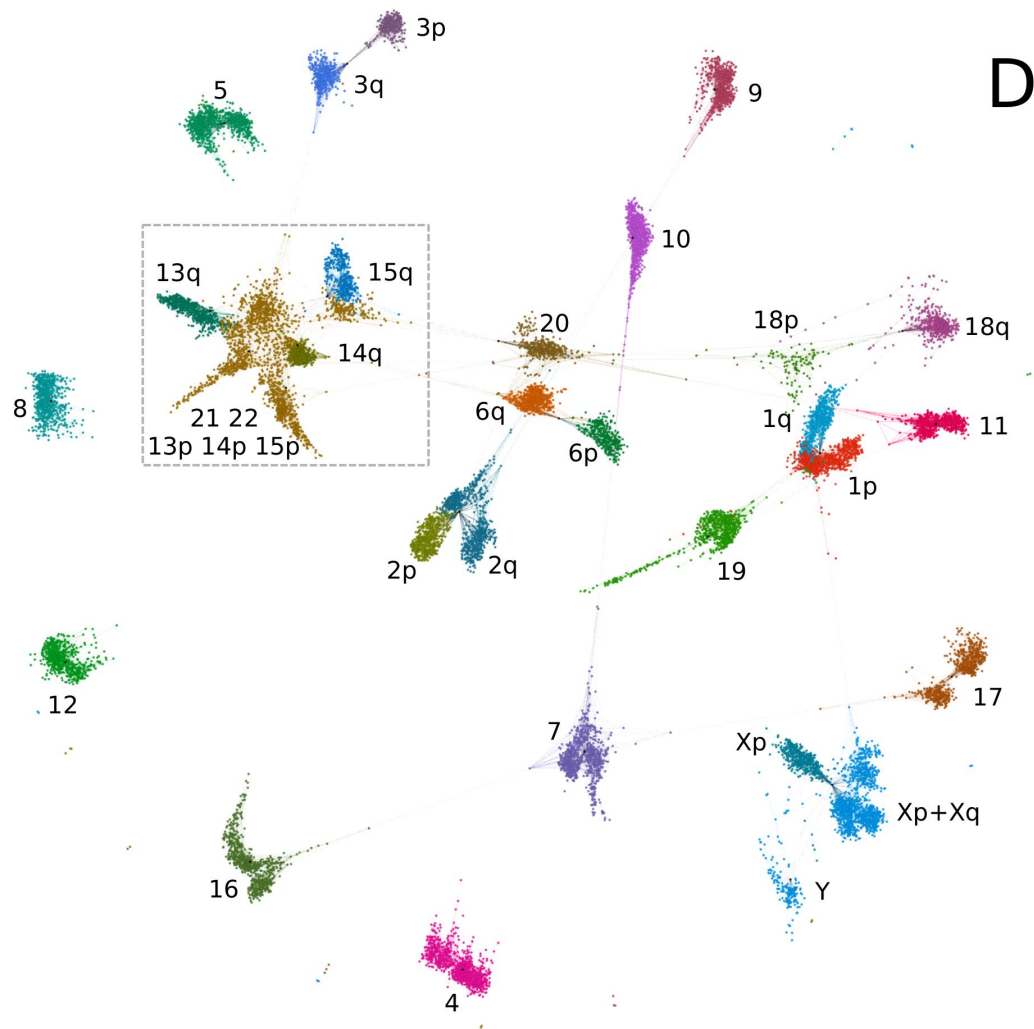
Chromosome communities in the HPRC

An all-vs-all mapping graph for
the HPRC contigs >1mbp.

A**B**

Leiden community detection



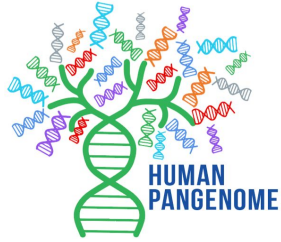


Leiden community detection

Labeling the layout with
community assignments.

We decided to take a closer look, focusing on the best assemblies in these regions.

Workflow



HPRC assemblies

https://github.com/human-pangenomics/HPP_Year1_Assemblies

We decided to take a closer look, focusing on the best assemblies in these regions.

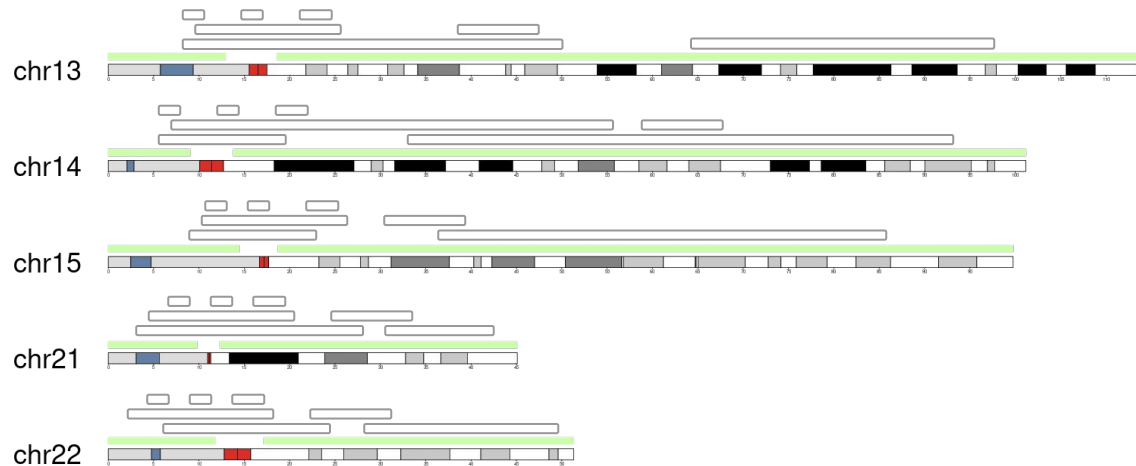
Workflow



HPRC assemblies

Mapping against
the whole CHM13

https://github.com/human-pangenomics/HPP_Year1_Assemblies



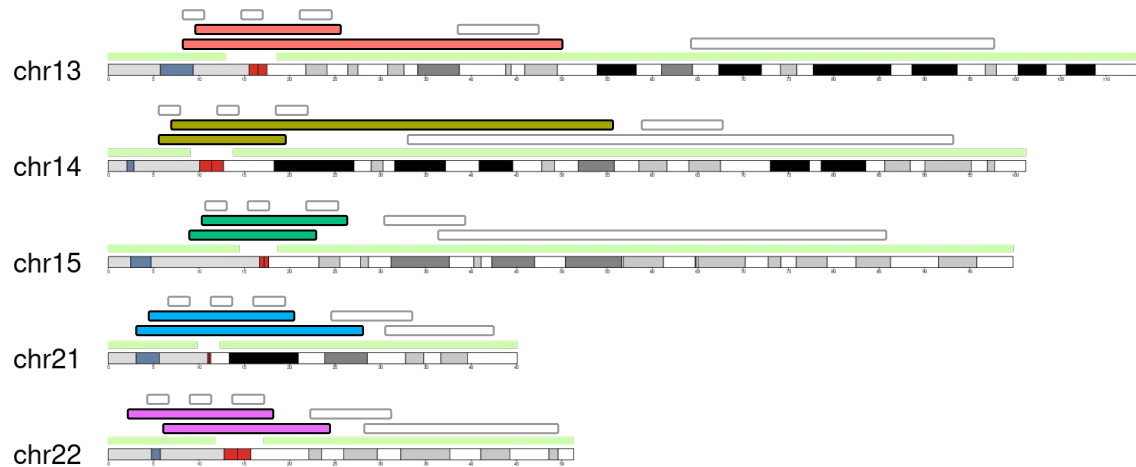
Workflow



HPRC assemblies

https://github.com/human-pangenomics/HPP_Year1_Assemblies

Mapping against
the whole CHM13



Acrocentric contigs covering (+/- 1Mbp) both the p and q arms (pq-contigs)

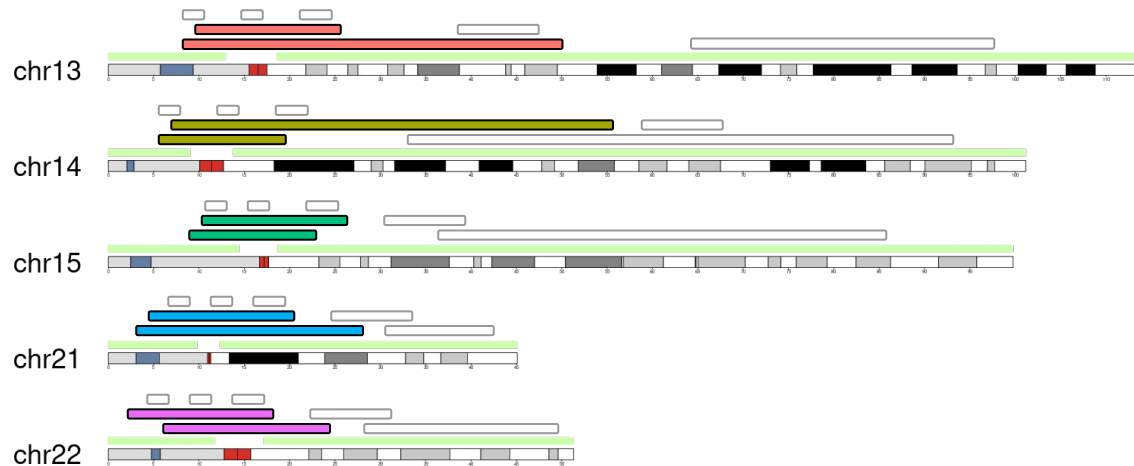
Workflow



HPRC assemblies

https://github.com/human-pangenomics/HPP_Year1_Assemblies

Mapping against
the whole CHM13

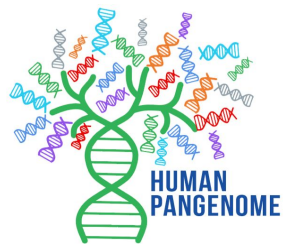


Acrocentric contigs covering (\pm 1Mbp) both the p and q arms (pq-contigs)

↓ + HG002 contigs \geq 300kbps
which map to acrocentrics

PanGenome Graph
Builder (PGGB)

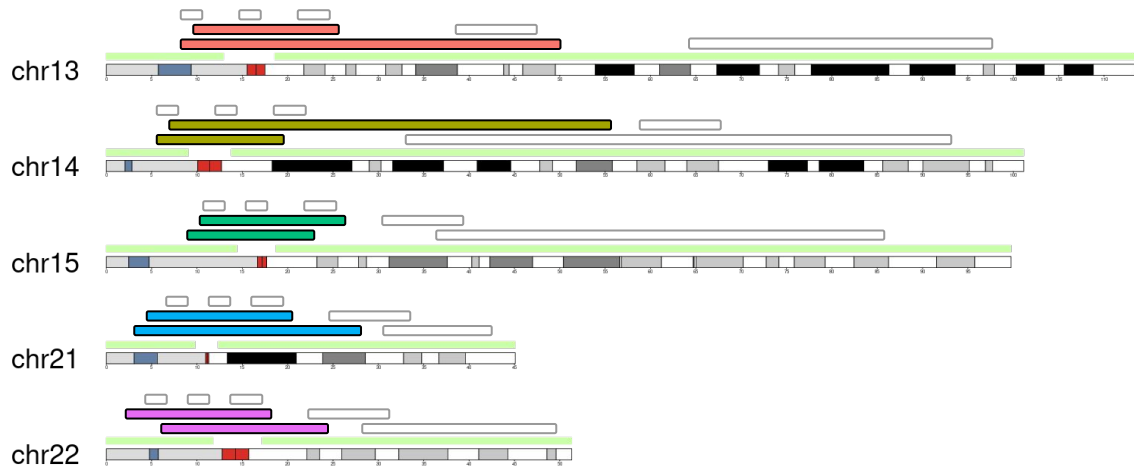
Workflow



HPRC assemblies

https://github.com/human-pangenomics/HPP_Year1_Assemblies

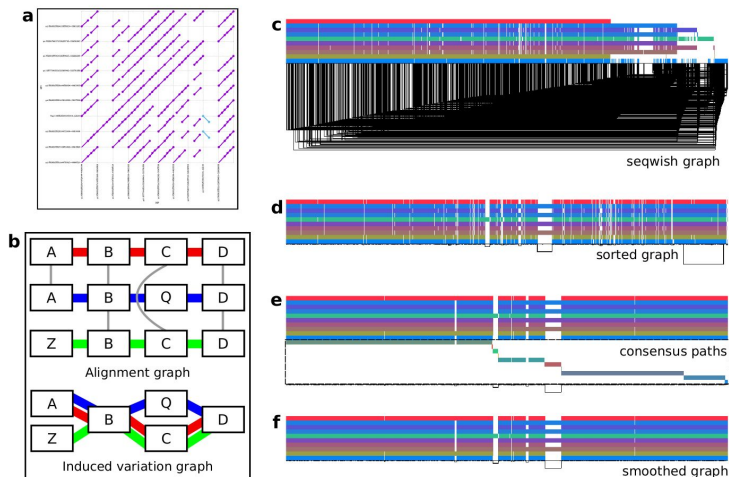
Mapping against
the whole CHM13



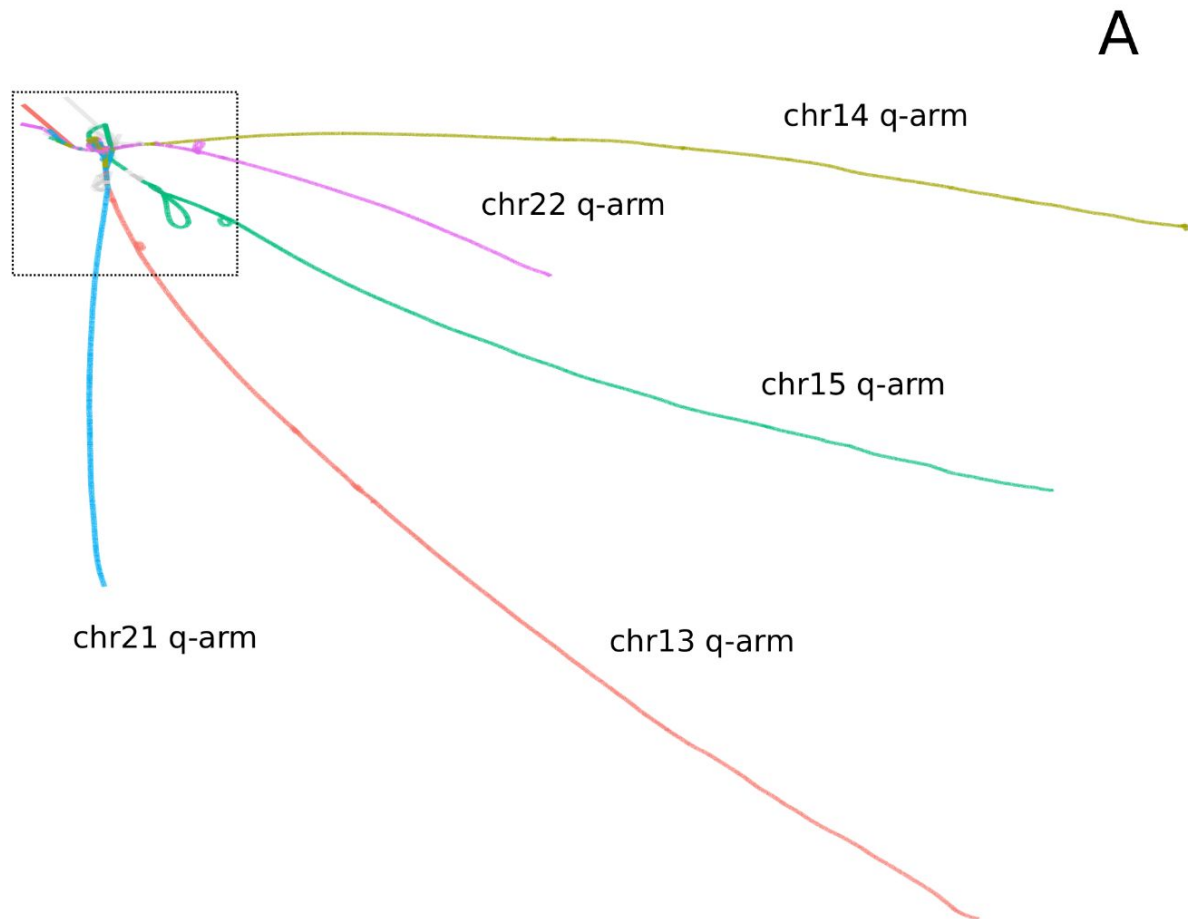
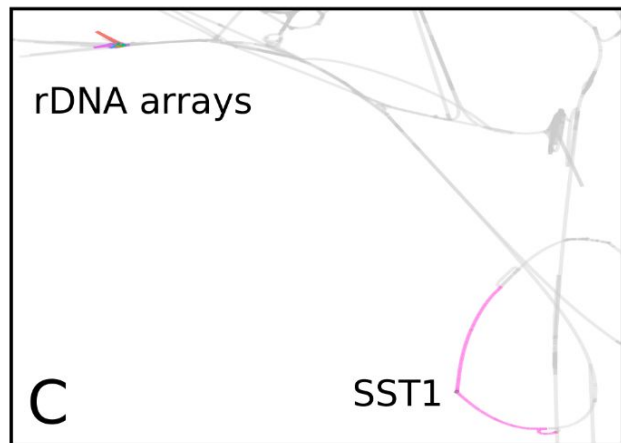
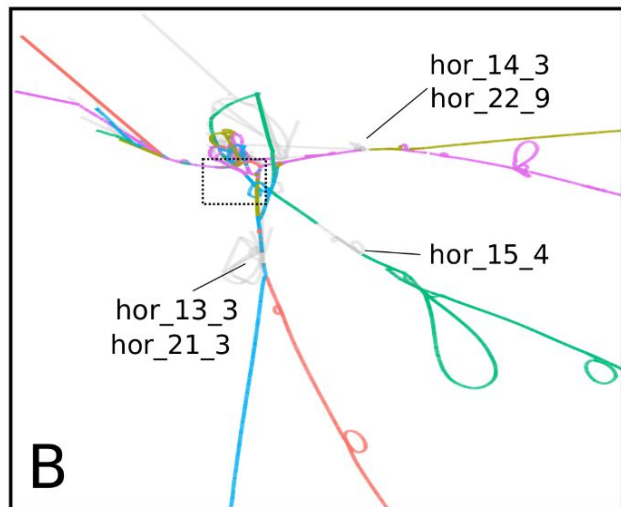
Acrocentric contigs covering (\pm 1Mbp) both the p and q arms (pq-contigs)

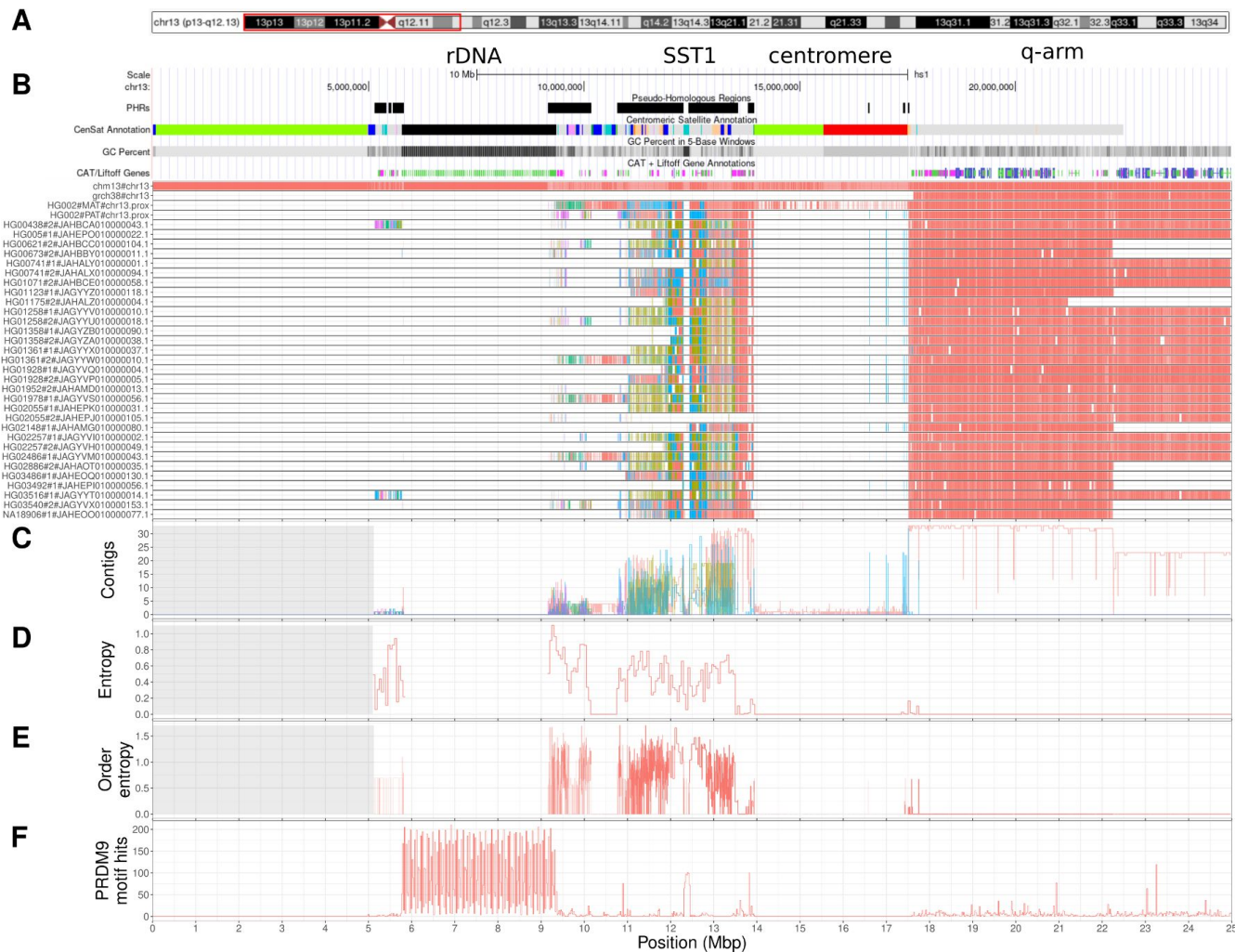
+ HG002 contigs \geq 300kbp
which map to acrocentrics

PanGenome Graph
Builder (PGGB)



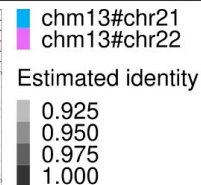
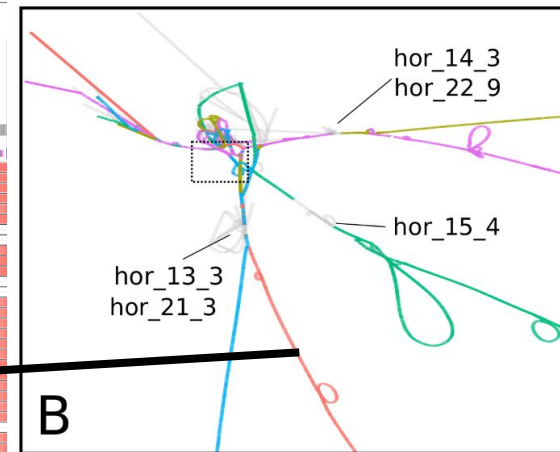
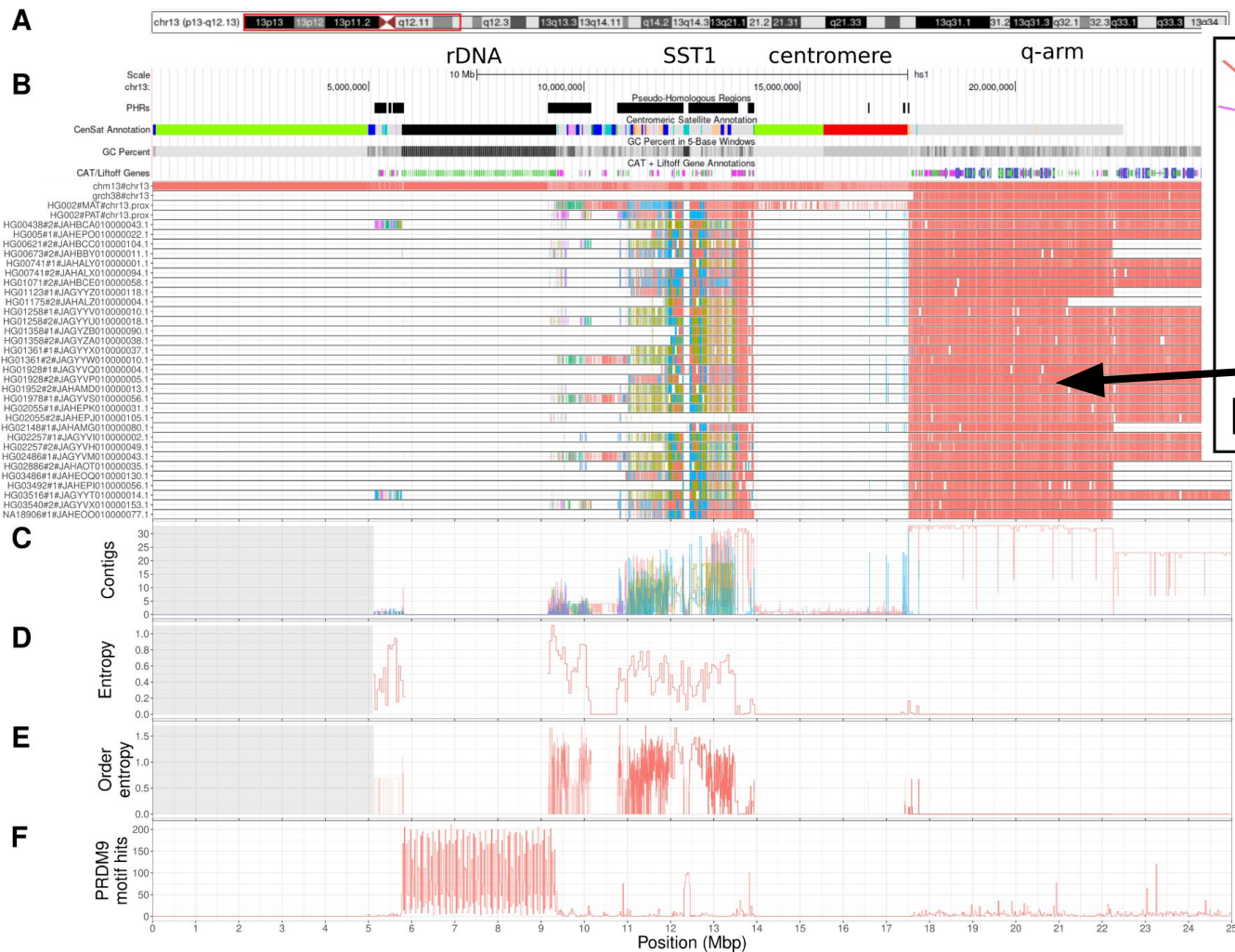
<https://github.com/pangenome/pggb>



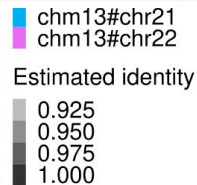
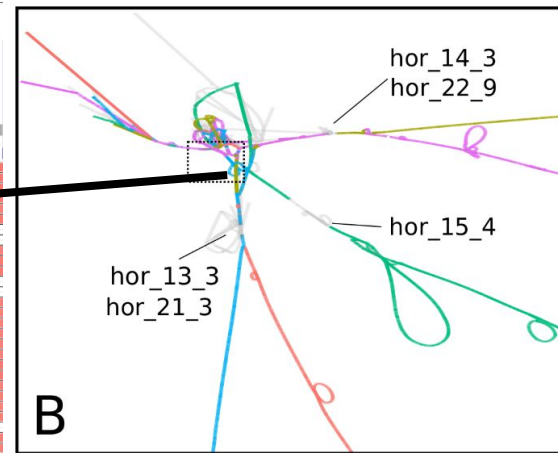
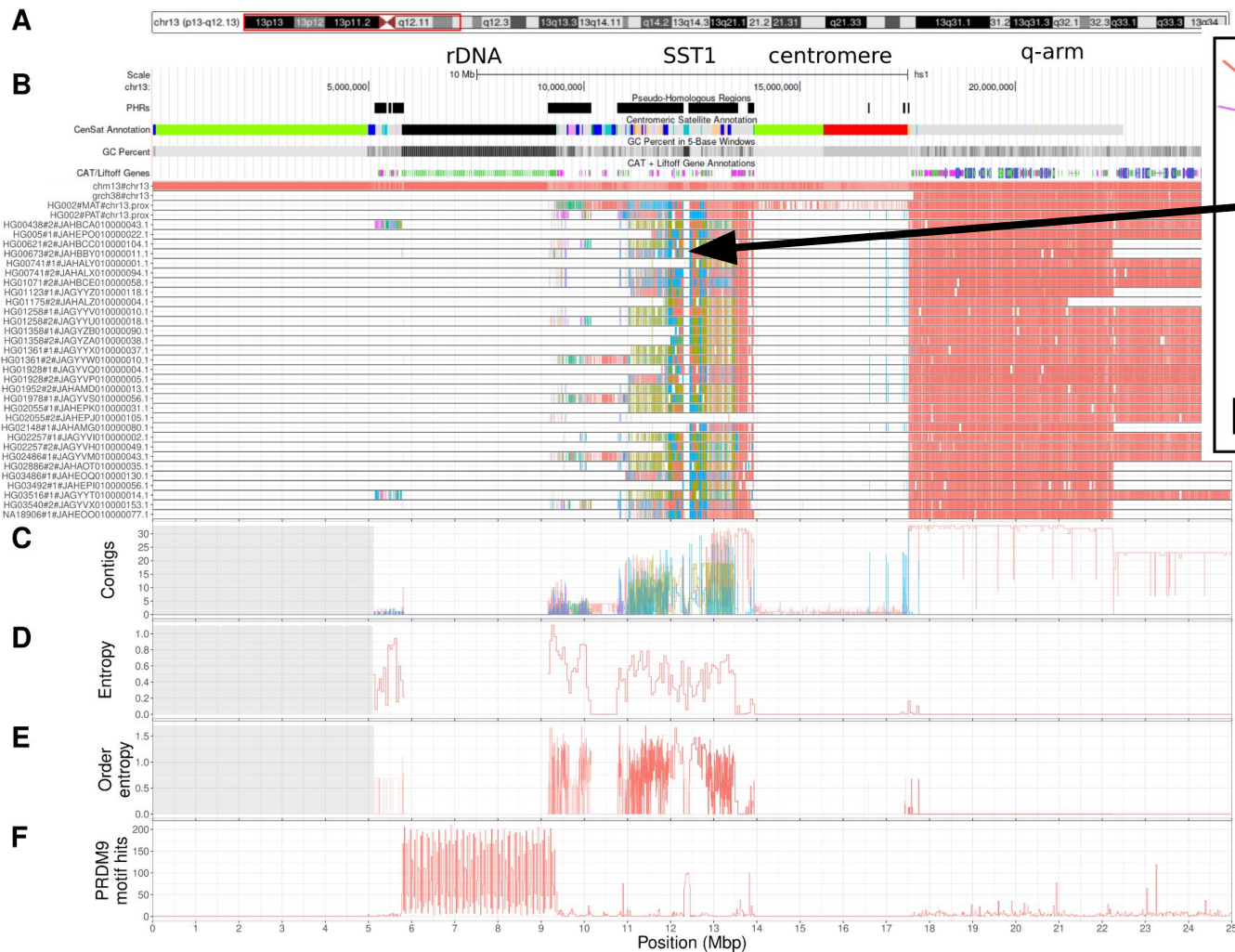


Untangling the pangenome graph

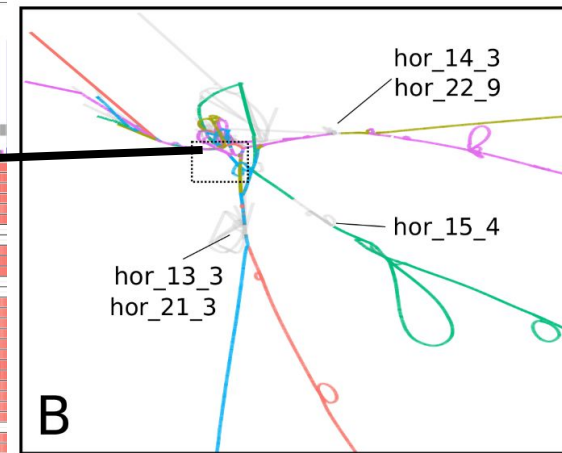
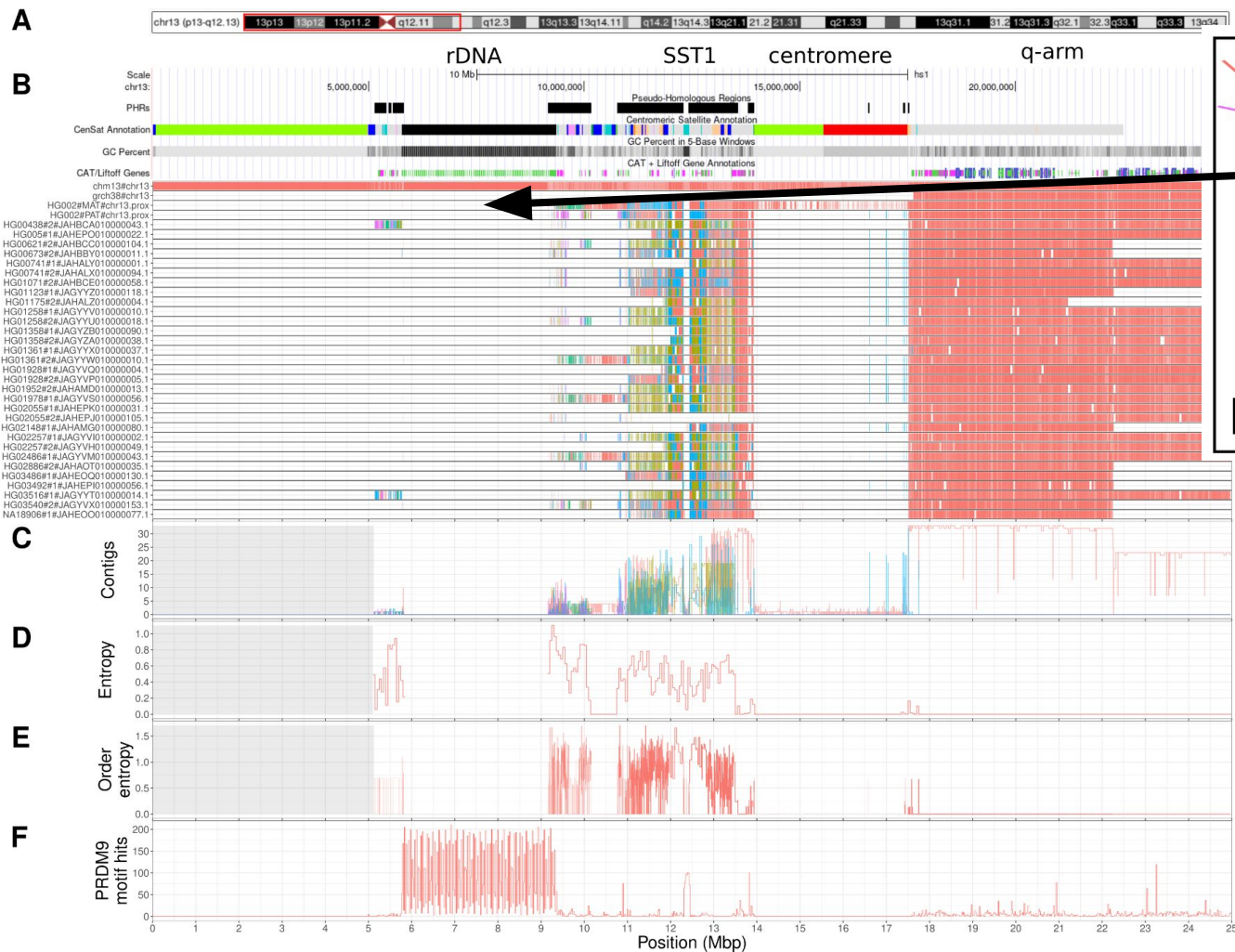
We look into the graph from the perspective of **chromosome 13**. Full information from pangenome plus reference annotations.



We look into the graph from the perspective of **chromosome 13**. Full information from pangenome plus reference annotations.



We look into the graph from the perspective of **chromosome 13**. Full information from pangenome plus reference annotations.



■ chm13#chr21
■ chm13#chr22

Estimated identity

0.925
0.950
0.975
1.000






We look into the graph from the perspective of **chromosome 13**. Full information from pangenome plus reference annotations.

A**B**

Pseudo
Homologous
Regions

A



 chm13#chr13
 chm13#chr14
 chm13#chr15
 chm13#chr21
 chm13#chr22

0.925
0.950
0.975
1.000

A



 chm13#chr13
 chm13#chr14
 chm13#chr15
 chm13#chr21
 chm13#chr22

0.925
0.950
0.975
1.000

A

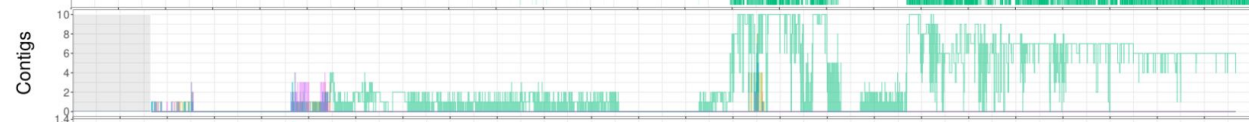
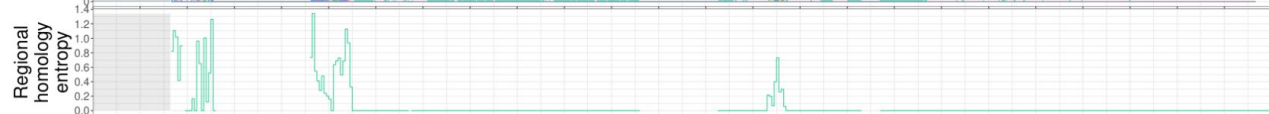
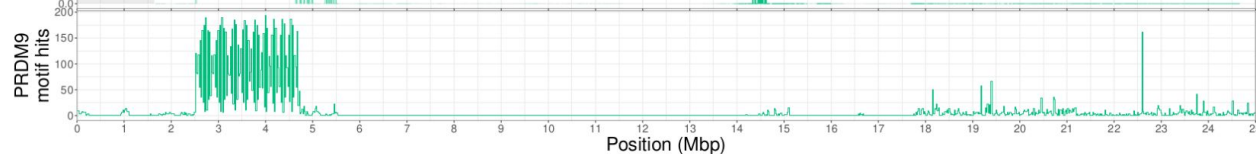
chromosome 15

Target

- chm13#chr13
- chm13#chr14
- chm13#chr15
- chm13#chr21
- chm13#chr22

Estimated identity

- 0.925
- 0.950
- 0.975
- 1.000

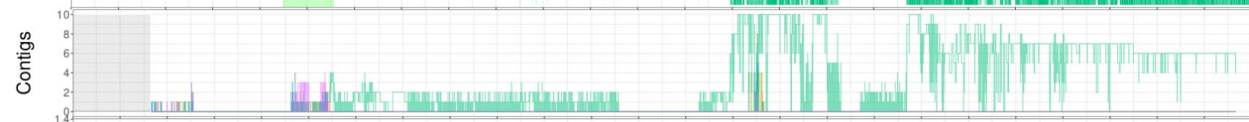
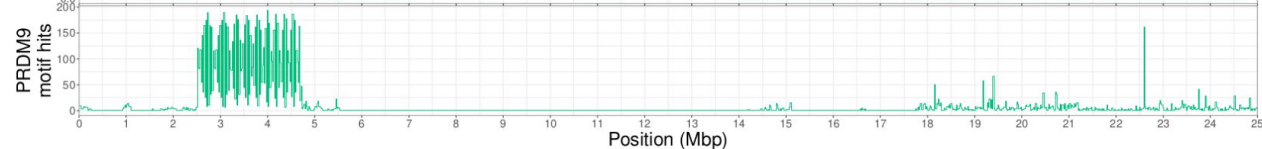
B**C****D****E**

A**chromosome 15****Target**

chm13#chr13
chm13#chr14
chm13#chr15
chm13#chr21
chm13#chr22

Estimated identity

0.925
0.950
0.975
1.000

B**C****D****E**



chromosome 21

Target

- chm13#chr13
- chm13#chr14
- chm13#chr15
- chm13#chr21
- chm13#chr22

Estimated identity

- 0.925
- 0.950
- 0.975
- 1.000

chromosome 21

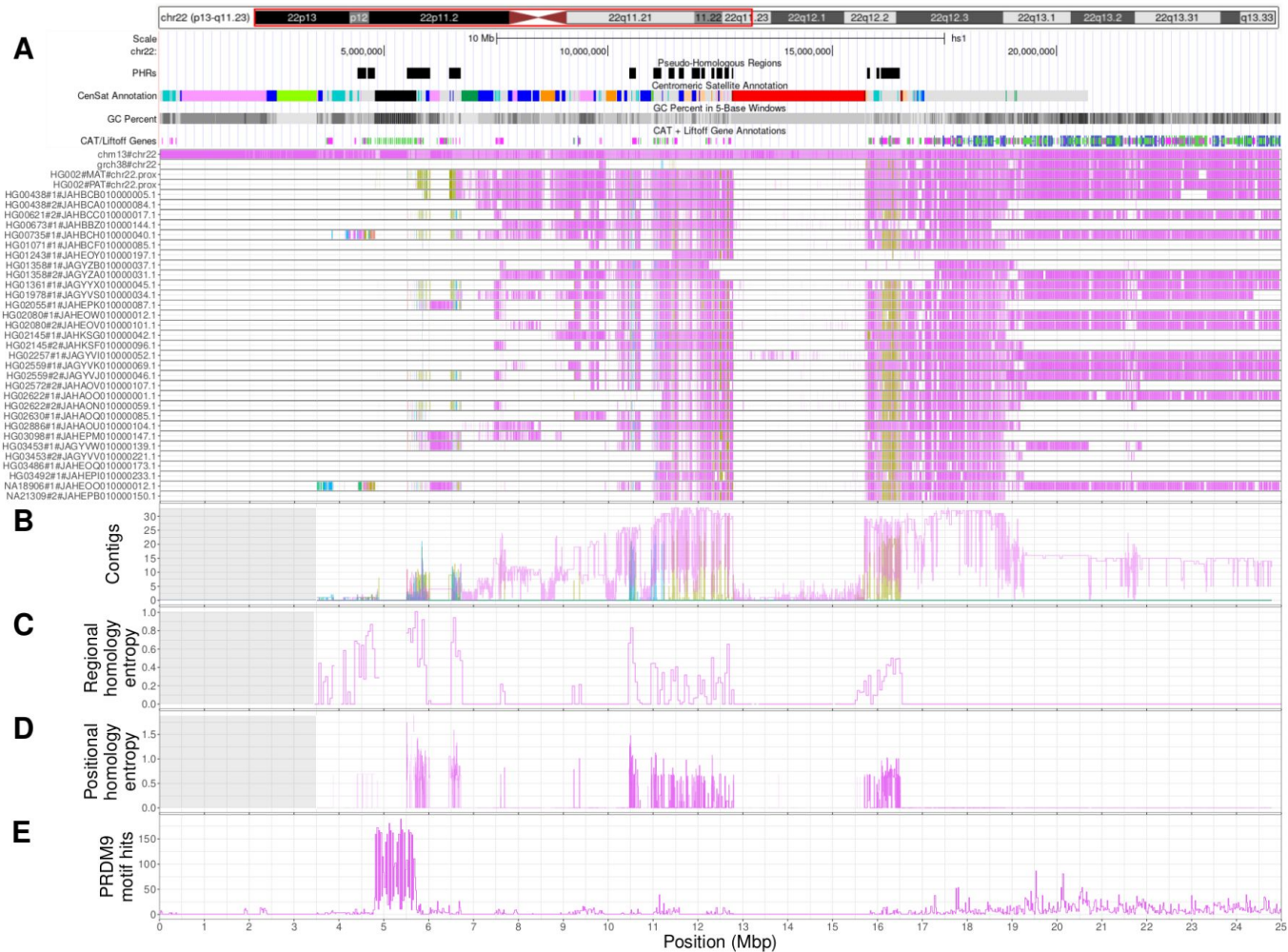


Target

- chm13#chr13
- chm13#chr14
- chm13#chr15
- chm13#chr21
- chm13#chr22

Estimated identity

- 0.925
- 0.950
- 0.975
- 1.000



chromosome 22

Target

chm13#chr13
chm13#chr14
chm13#chr15
chm13#chr21
chm13#chr22

Estimated identity

0.925
0.950
0.975
1.000



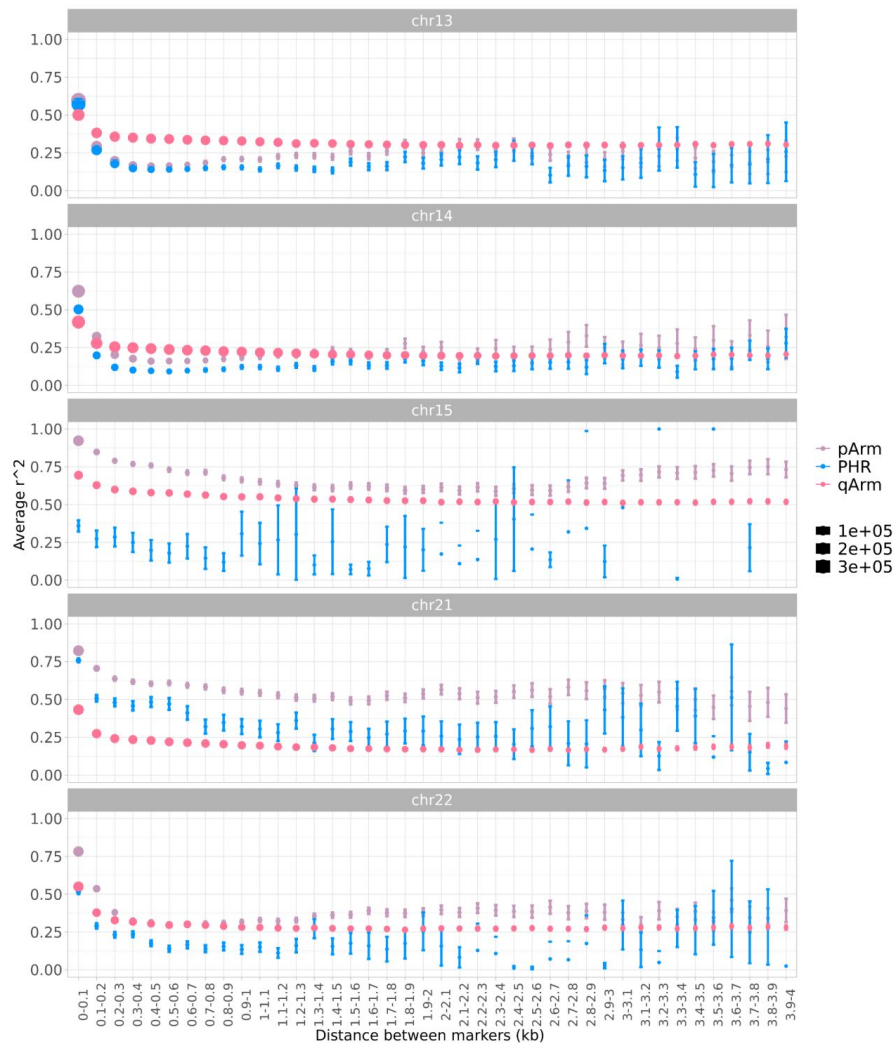
chromosome 22

Target

- chm13#chr13
- chm13#chr14
- chm13#chr15
- chm13#chr21
- chm13#chr22

Estimated identity

- 0.925
- 0.950
- 0.975
- 1.000

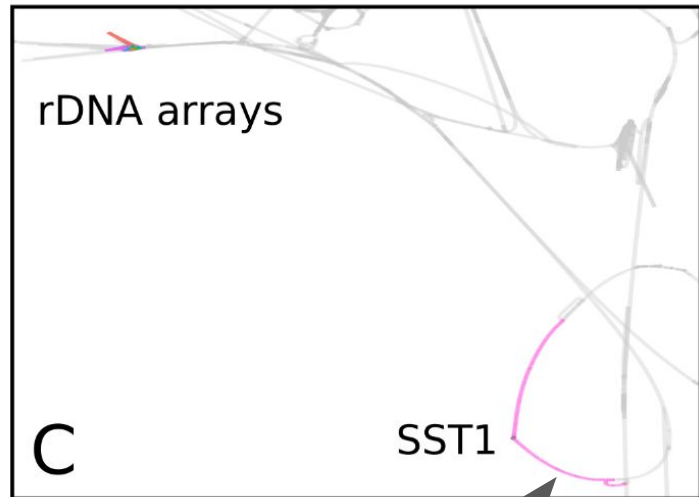


LD decays faster in pseudo-homologous regions than elsewhere in the p-arms or in the q-arms.

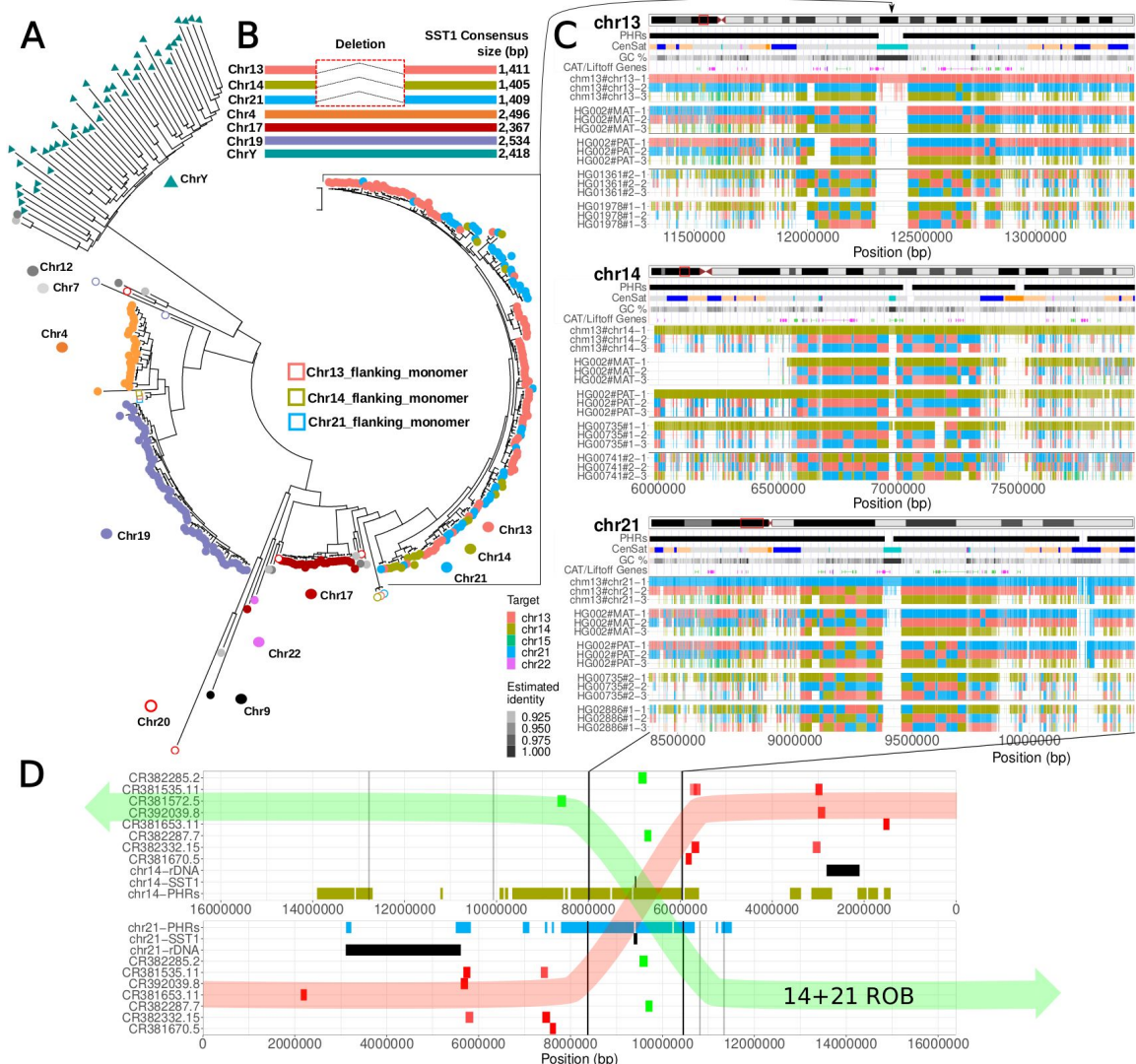
This pattern is consistent with higher recombination rates and/or effective population size in these regions.

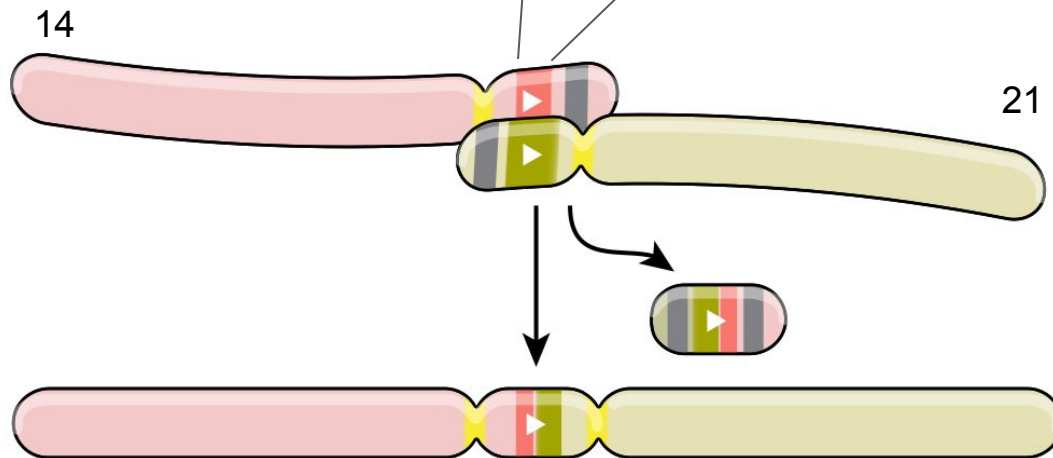
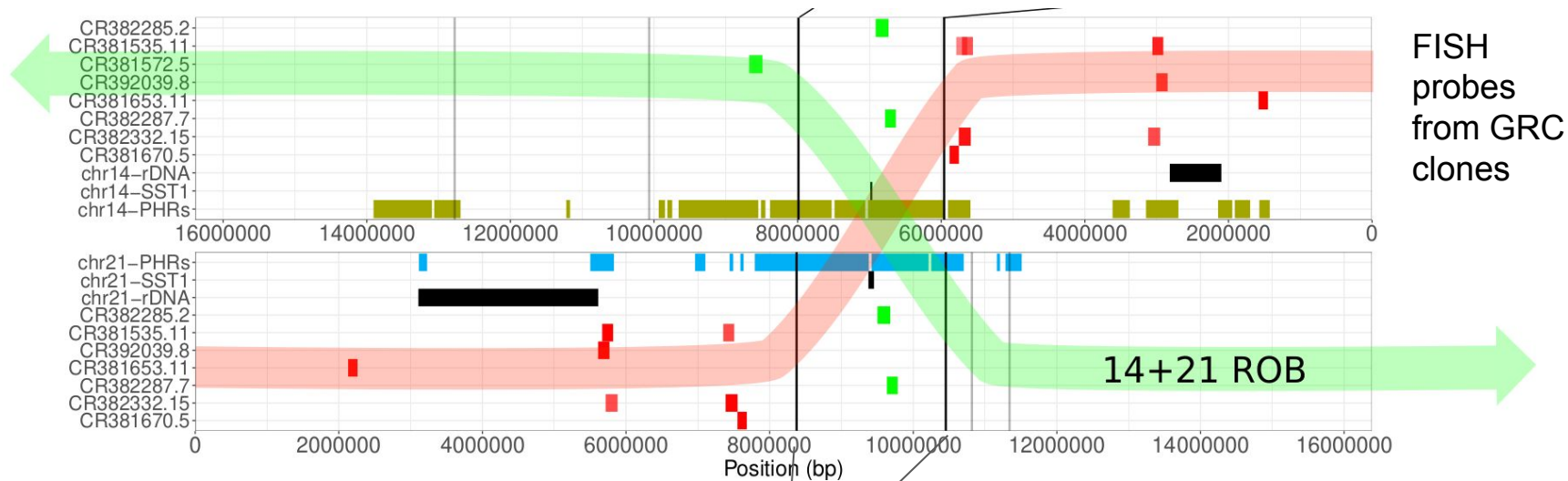
chromosome 13/14/21

Pseudo-homologous regions



The SST1-linked PHR is the site of the most intense signals of recombination between heterologous chromosomes.





Recombination between heterologous chromosomes

The high level of homology of the acrocentric chromosomes is likely due to **recombination between heterologous chromosomes!**

High-quality *de novo* assemblies and pangenomic approaches thus shed light on the most difficult regions of the human genomes.

This answers questions that arose in the early era of cytogenetics, ~50 years ago.

Volume 16 Number 4 1988

Nucleic Acids Research

Homologous alpha satellite sequences on human acrocentric chromosomes with selectivity for chromosomes 13, 14 and 21: implications for recombination between nonhomologues and Robertsonian translocations

K.H.Choo*, B.Vissel, R.Brown, R.G.Filby and E.Earle

ABSTRACT

We report a new subfamily of alpha satellite DNA (pTRA-2) which is found on all the human acrocentric chromosomes. The alphoid nature of the cloned DNA was established by partial sequencing. Southern analysis of restriction enzyme-digested DNA fragments from mouse/human hybrid cells containing only human chromosome 21 showed that the predominant higher-order repeating unit for pTRA-2 is a 3.9 kb structure. Analysis of a "consensus" *in situ* hybridisation profile derived from 13 normal individuals revealed the localisation of 73% of all centromeric autoradiographic grains over the five acrocentric chromosomes, with the following distribution: 20.4%, 21.5%, 17.1%, 7.3% and 6.5% on chromosomes 13, 14, 21, 15 and 22 respectively. An average of 1.4% of grains was found on the centromere of each of the remaining 19 nonacrocentric chromosomes. These results indicate the presence of a common subfamily of alpha satellite DNA on the five acrocentric chromosomes and suggest an evolutionary process consistent with recombination exchange of sequences between the nonhomologues. The results further suggests that such exchanges are more selective for chromosomes 13, 14 and 21 than for chromosomes 15 and 22. The possible role of centromeric alpha satellite DNA in the aetiology of 13q14q and 14q21q Robertsonian translocations involving the common and nonrandom association of chromosomes 13 and 14, and 14 and 21 is discussed.

[Chroo et al., 1988.](#)

HUMAN MEIOSIS I. THE HUMAN PACHYTENE KARYOTYPE ANALYZED BY THREE DIMENSIONAL RECONSTRUCTION OF THE SYNAPTONEMAL COMPLEX

by
PREBEN BACH HOLM
and
SØREN WILKEN RASMUSSEN

Department of Physiology, Carlsberg Laboratory
Gamle Carlsberg Vej 10, DK-2500 Copenhagen, Valby

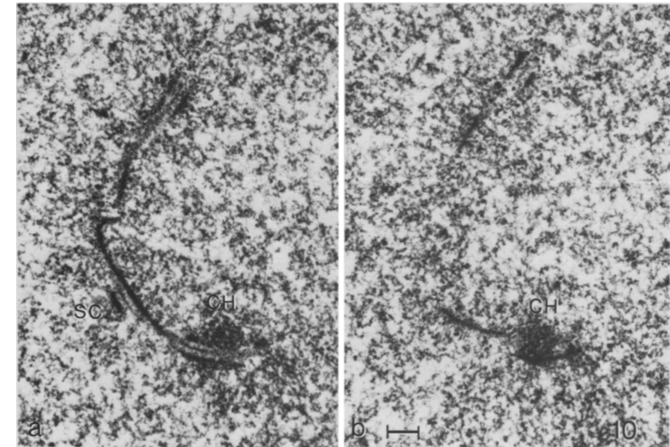
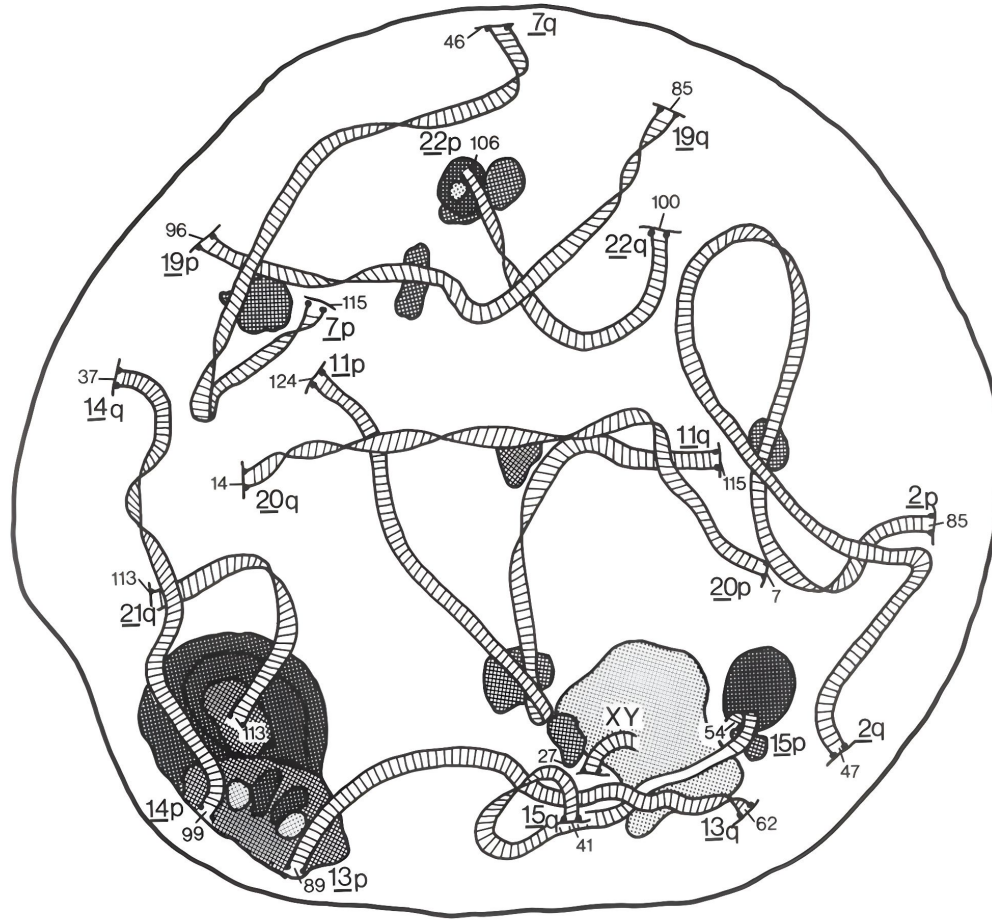
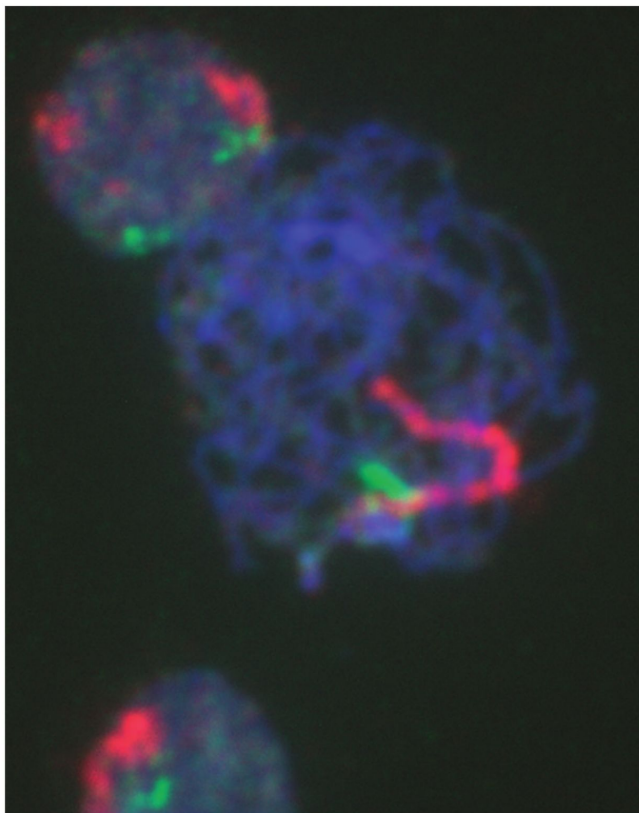




Figure 10. Two consecutive sections through the centromeric heterochromatin of a bivalent at early pachytene. The synaptonemal complex (SC) passes unaltered through the centromeric heterochromatin (CH). (Bar = 0.2 μ m)



TRANSACTIONS OF THE 70TH ANNUAL MEETING OF THE PACIFIC COAST OBSTETRICIANS AND GYNECOLOGICAL SOCIETY | VOLUME 190, ISSUE 6, P1781-1785, JUNE 01, 2004

FISHing for acrocentric associations between chromosomes 14 and 21 in human oogenesis

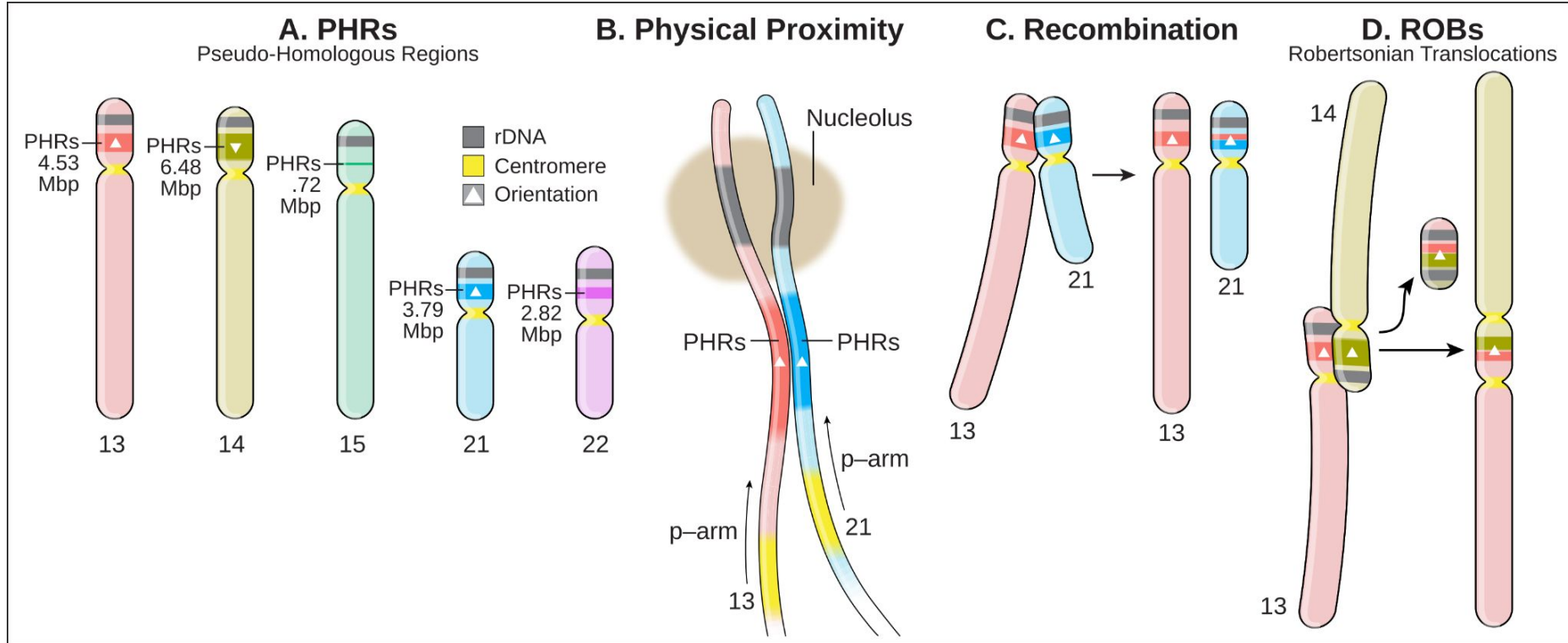
Edith Y Cheng, MD   Theresa Nalulai-Cecchini, BA

observed 300 cells to see a
single putative 14:21 synapse!

Figure The pachytene nucleus with 2 hybridization signals is oriented in a linear fashion. The *red signal* represents chromosome 14, and the *green signal* represents chromosome 21.

<https://doi.org/10.1016/j.ajog.2004.02.062>

Pseudo-homologous regions (PHRs)



Recombination between heterologous human acrocentric chromosomes

<https://doi.org/10.1038/s41586-023-05976-y>

Received: 15 August 2022

Accepted: 17 March 2023

Published online: 10 May 2023

Open access

 Check for updates

Andrea Guarracino^{1,2}, Silvia Buonaiuto³, Leonardo Gomes de Lima⁴, Tamara Potapova⁴, Arang Rhie⁵, Sergey Koren⁵, Boris Rubinstein⁴, Christian Fischer¹, Human Pangenome Reference Consortium⁶, Jennifer L. Gerton⁴, Adam M. Phillippy⁵, Vincenza Colonna^{1,3} & Erik Garrison¹✉

The short arms of the human acrocentric chromosomes 13, 14, 15, 21 and 22 (SAACs) share large homologous regions, including ribosomal DNA repeats and extended segmental duplications^{1,2}. Although the resolution of these regions in the first complete assembly of a human genome—the Telomere-to-Telomere Consortium’s CHM13 assembly (T2T-CHM13)—provided a model of their homology³, it remained unclear whether these patterns were ancestral or maintained by ongoing recombination exchange. Here we show that acrocentric chromosomes contain pseudo-homologous regions (PHRs) indicative of recombination between non-homologous sequences. Utilizing an all-to-all comparison of the human pangenome from the Human Pangenome Reference Consortium⁴ (HPRC), we find that contigs from all of the SAACs form a community. A variation graph⁵ constructed from centromere-spanning acrocentric contigs indicates the presence of regions in which most contigs appear nearly identical between heterologous acrocentric chromosomes in T2T-CHM13. Except on chromosome 15, we observe faster decay of linkage disequilibrium in the pseudo-homologous regions than in the corresponding short and long arms, indicating higher rates of recombination^{6,7}. The pseudo-homologous regions include sequences that have previously been shown to lie at the breakpoint of Robertsonian translocations⁸, and their arrangement is compatible with crossover in inverted duplications on chromosomes 13, 14 and 21. The ubiquity of signals of recombination between heterologous acrocentric chromosomes seen in the HPRC draft pangenome suggests that these shared sequences form the basis for recurrent Robertsonian translocations, providing sequence and population-based confirmation of hypotheses first developed from cytogenetic studies 50 years ago⁹.



to you, and...

Thanks!

Andrea Guarracino (pggb, wfmash, seqwish, odgi, chromosome communities)

Simon Heumos (pggb, odgi)

Flavia Villani (pggb, applications to mouse, popgen)

Njagi Mwaniki (wfmash, WFA applications)

Santiago Marco-Sola (WFA, wfmash)

Pjotr Prins (vcflib, vcfwave)

Richard Durbin (PhD guidance)

Nicole Soranzo (support)

Benedict Paten (vgteam)

Hao Chen (rat, mouse)

Zhigui Bao (plant applications)

Lorenzo Tattini (yeast pangenomes)

Enza Colonna (applications to mouse, popgen)

Nadia Pisanti (algorithms)

Luca Pinello (applications)

Peter Sudmant (primate pangenomes)

Robert Williams (guidance)

HPRC pangenomes working group and many others

funders:

NLnet

NSF

NIH (NIDA)

Amylase project!

Alessandro Raveane (Human Technopole)

Davide Bolognini (Human Technopole)

Peter Sudmant (Berkeley)

Joana Rocha (Berkeley)

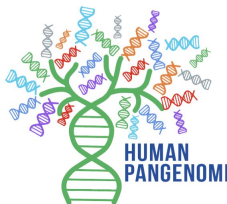
Andrea Guarracino (UTHSC)

Alma Halgren (Berkeley)

Jason Chin (GeneDX)

Nicholas Lou (Berkeley)

TOWARDS A
COMPLETE
REFERENCE OF
HUMAN GENOME
DIVERSITY



UNIVERSITY OF CALIFORNIA
SANTA CRUZ

Genomics
Institute



CORIELL INSTITUTE
FOR MEDICAL RESEARCH
DECODING THE GENOME

EMBL-EBI



HARVARD
MEDICAL SCHOOL



Icahn School of Medicine
at Mount Sinai



the
sanger
institute



Yale University



TOWARDS A
COMPLETE
REFERENCE OF
HUMAN GENOME
DIVERSITY



We would like to acknowledge the National Genome Research Institute (NHGRI) for funding the following grants which are in support of creating the human pangenome reference: 1U41HG010972, 1U01HG010971, 1U01HG010961, 1U01HG010973, 1U01HG010963, and the Human Pangenome Reference Consortium (<https://humanpangenome.org/>)



illumina

Google Health



Cantata Bio



Global Alliance
for Genomics & Health

Collaborate. Innovate. Accelerate.

NIST



National Human Genome
Research Institute

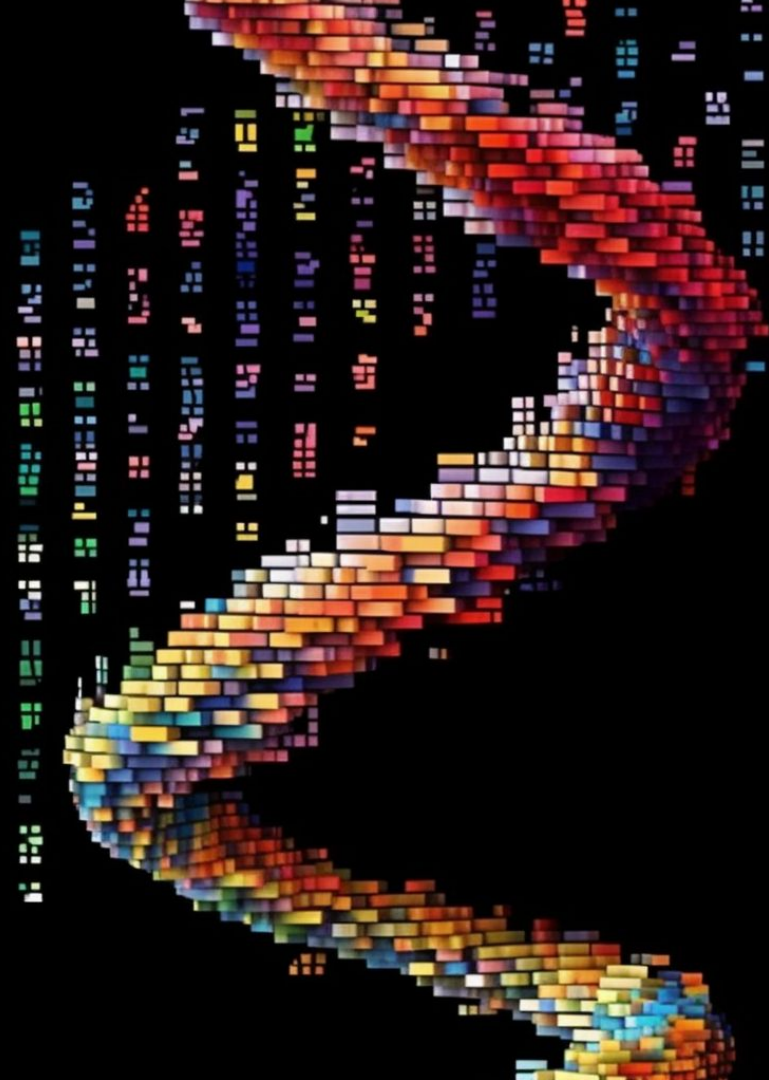
Practical!

Let's build some pangenome variation graphs with **pggb**!

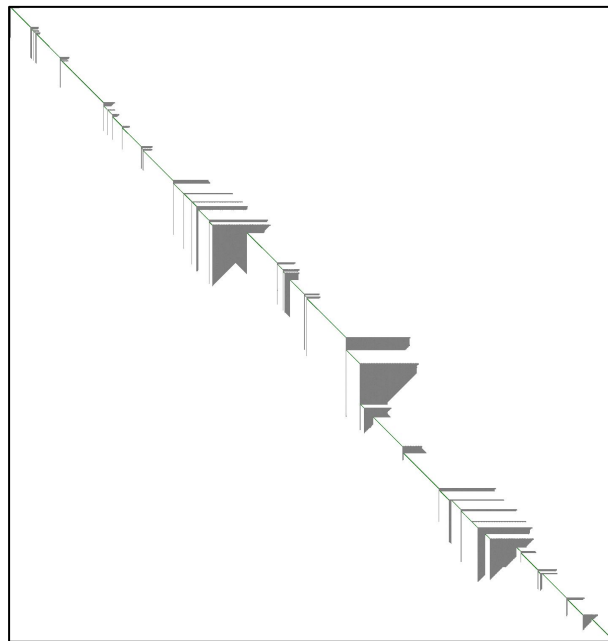
First: a deeper dive into how the method works.

Then: we'll work through small examples to learn how to drive it.



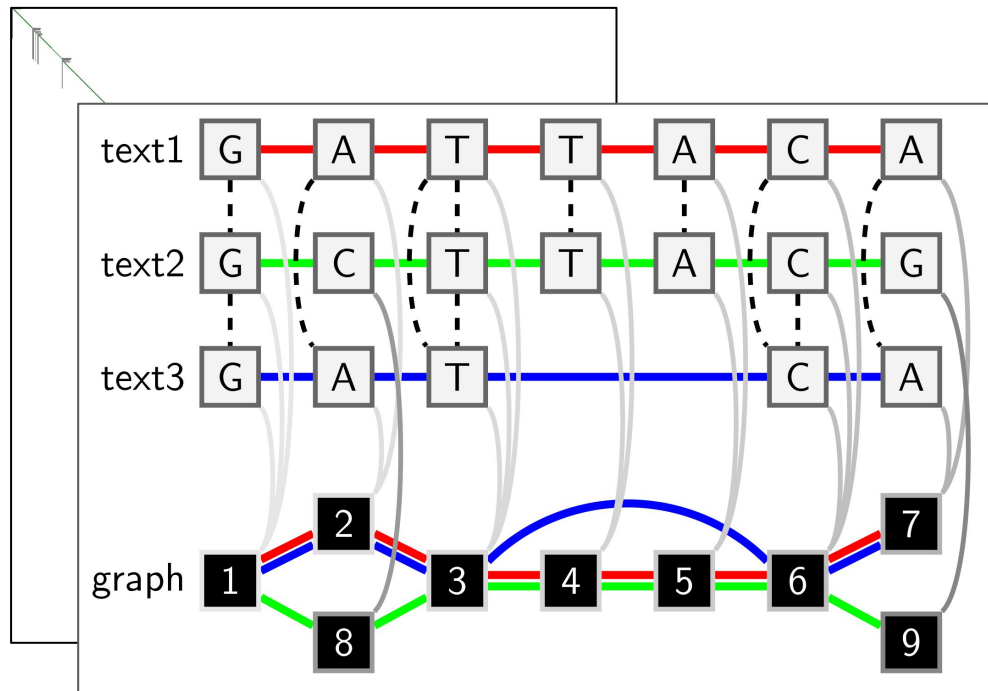


PanGenome Graph Builder



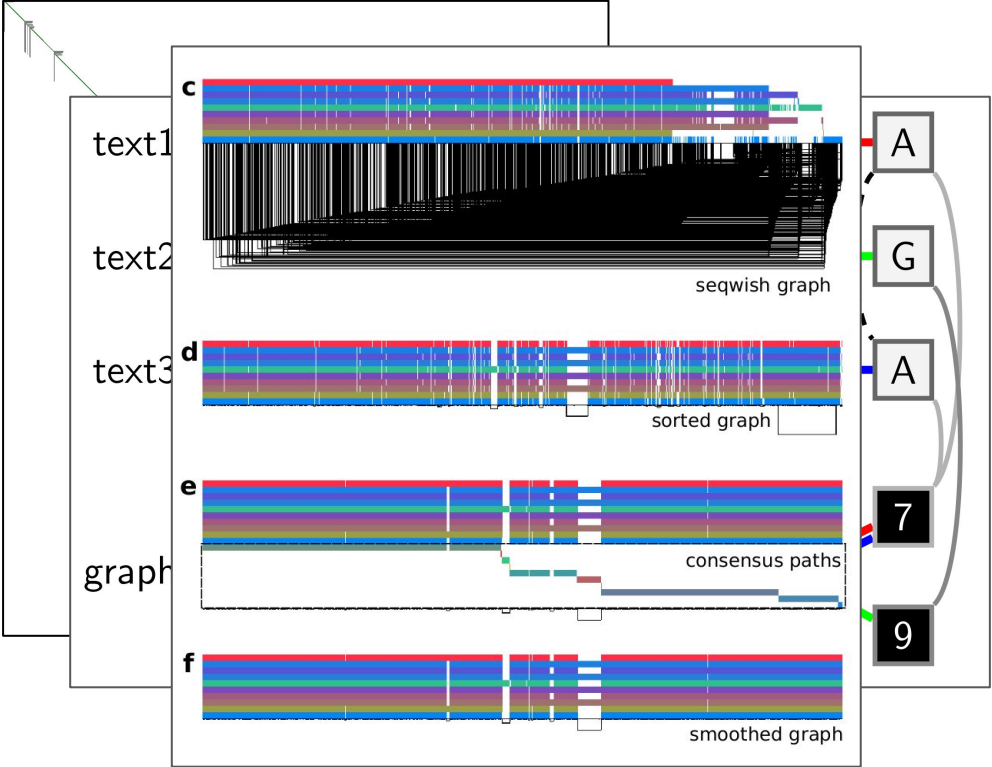
wfmash (biWFA)

PanGenome Graph Builder



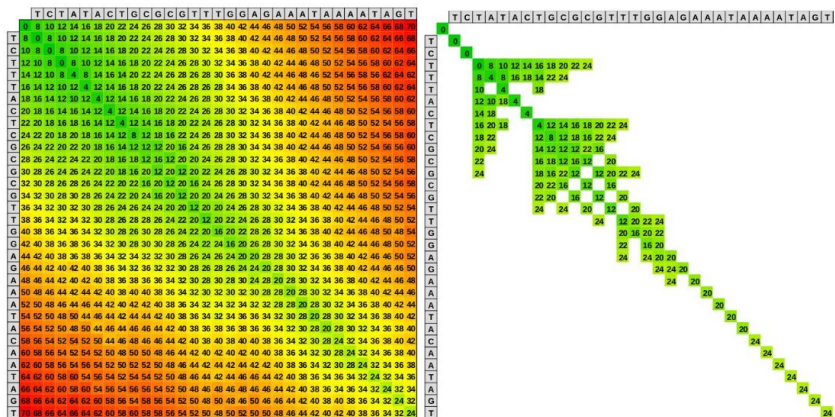
seqwish (unbiased graph builder)

PanGenome Graph Builder



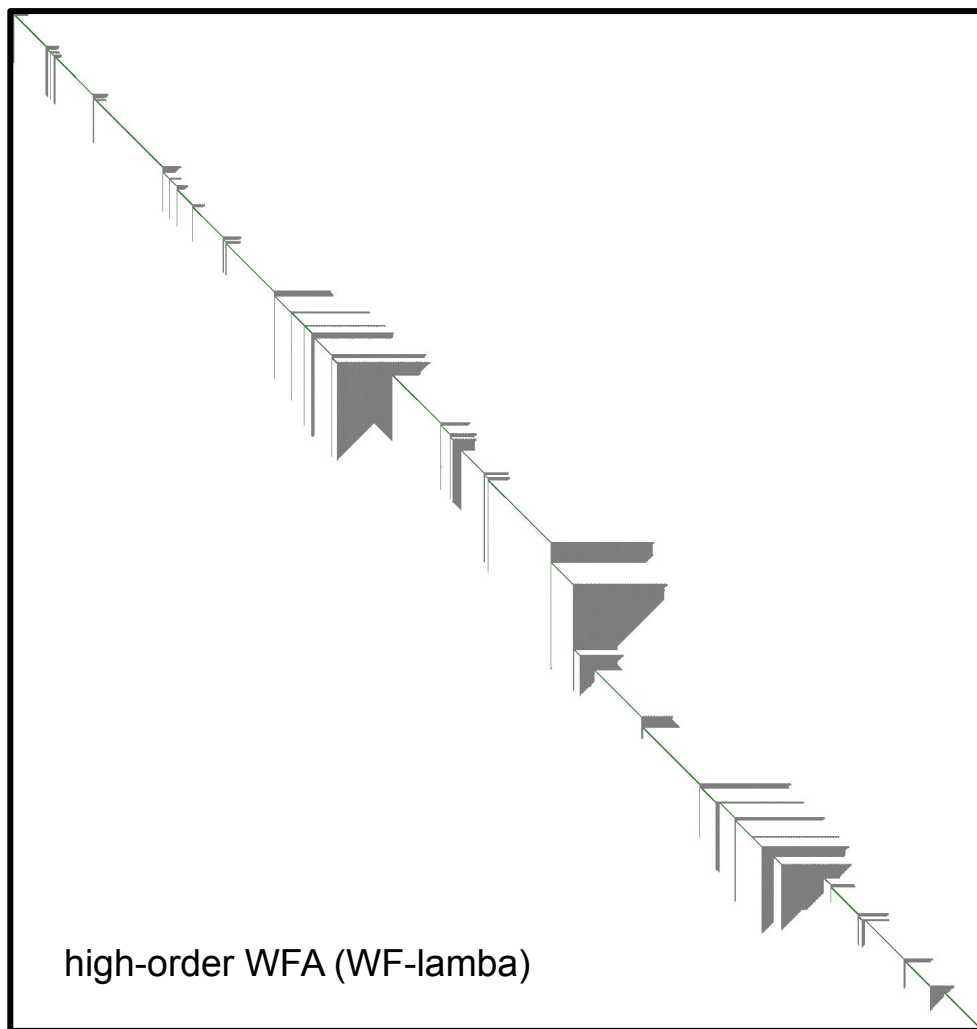
smoothxg (graph normalization)

wfmash makes initial alignments

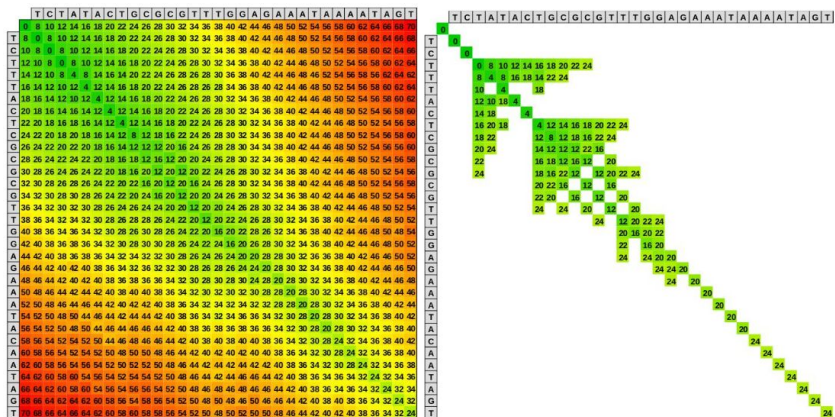


Needleman-Wunsch

the wavefront
algorithm (WFA)



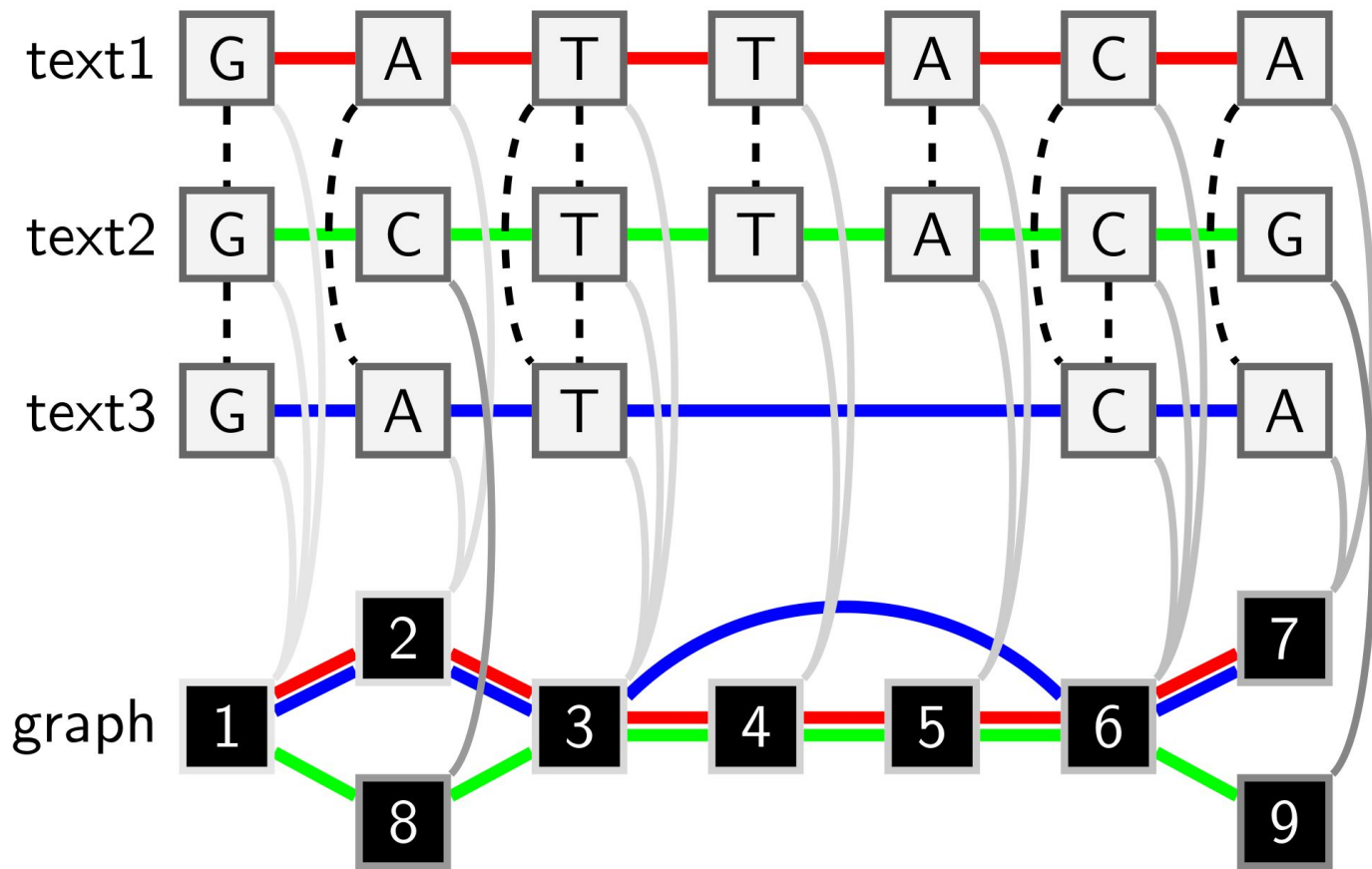
wfmash makes initial alignments



Needleman-Wunsch
the wavefront
algorithm (WFA)

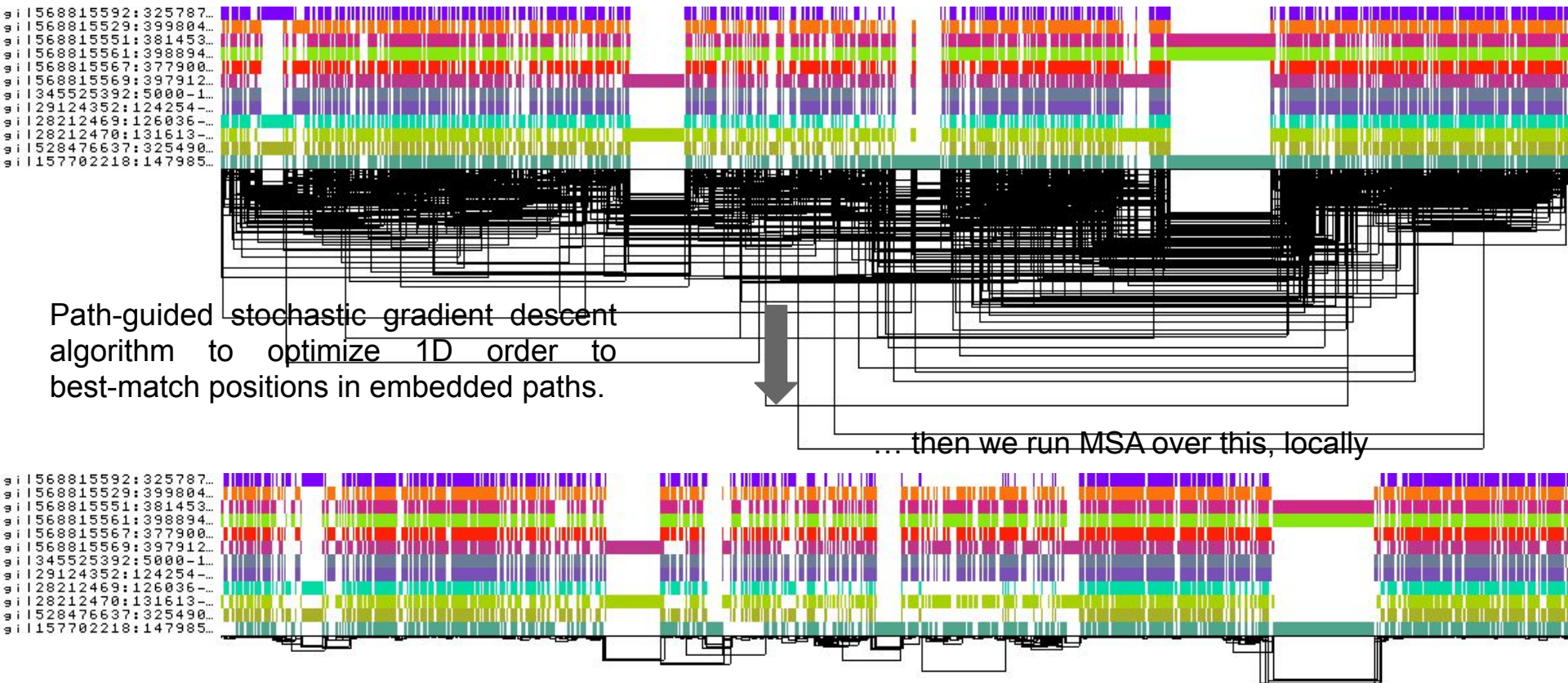
high-order *bidirectional* WFA (BiWFL)

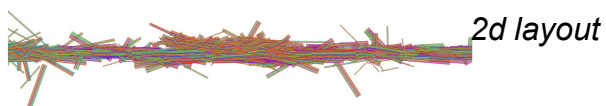
seqwish
builds the
graph



smoothxg organizes & normalizes the graph

Pangenome graph with 12 ALT sequences of the HLA-DRB1 gene from the GRCh38 reference genome.



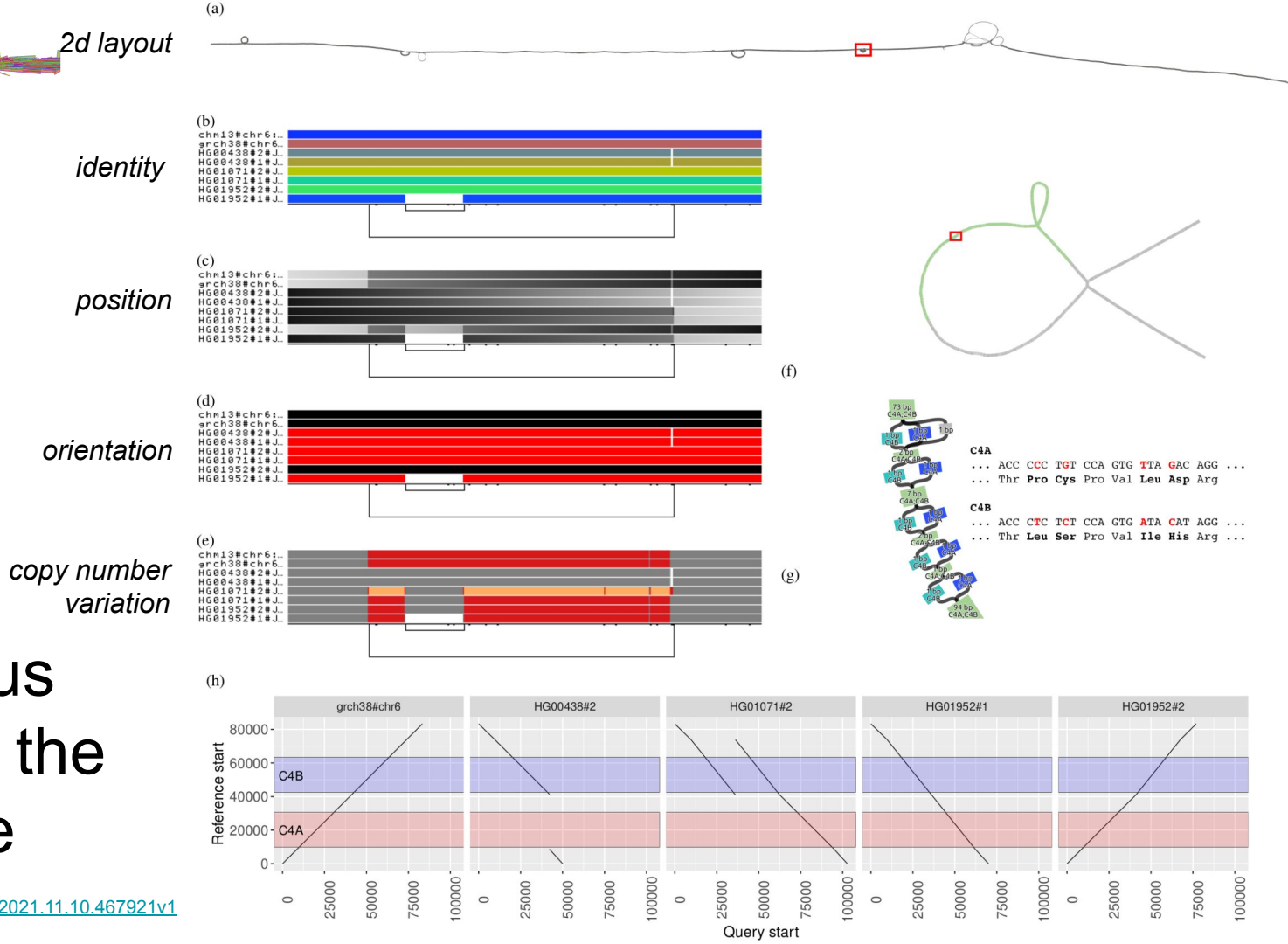


ODGI is meant to be a basic toolkit for interacting with pangenome graphs.

It uses the embedded genomes as references.

odgi helps us understand the pangenome

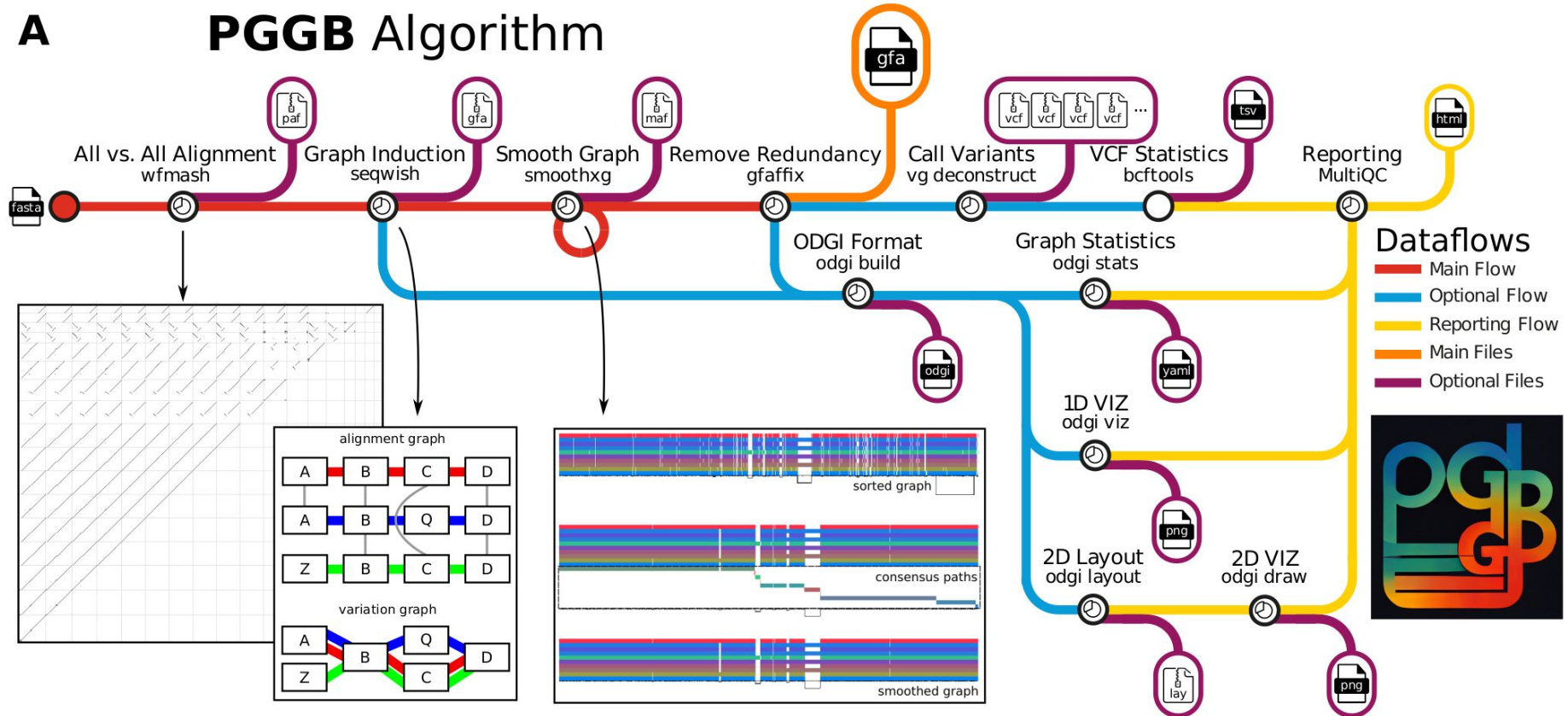
<https://www.biorxiv.org/content/10.1101/2021.11.10.467921v1>



Putting it all together!

A

PGGB Algorithm

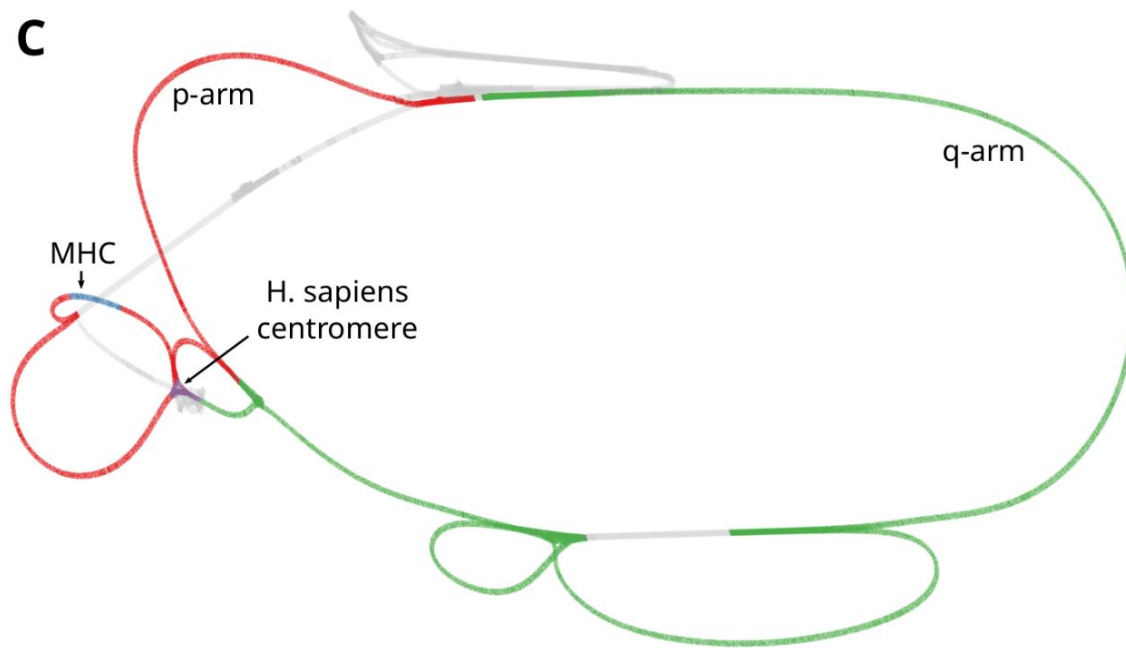


T2T Primate Pangenome: Chromosome 6

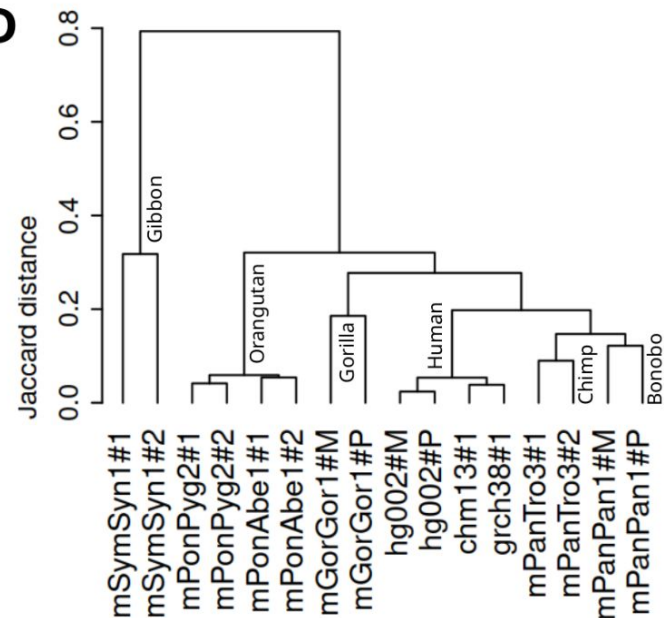
B



C



D

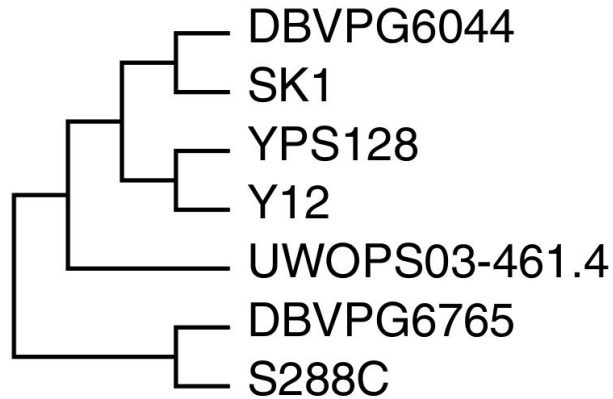


Test material today

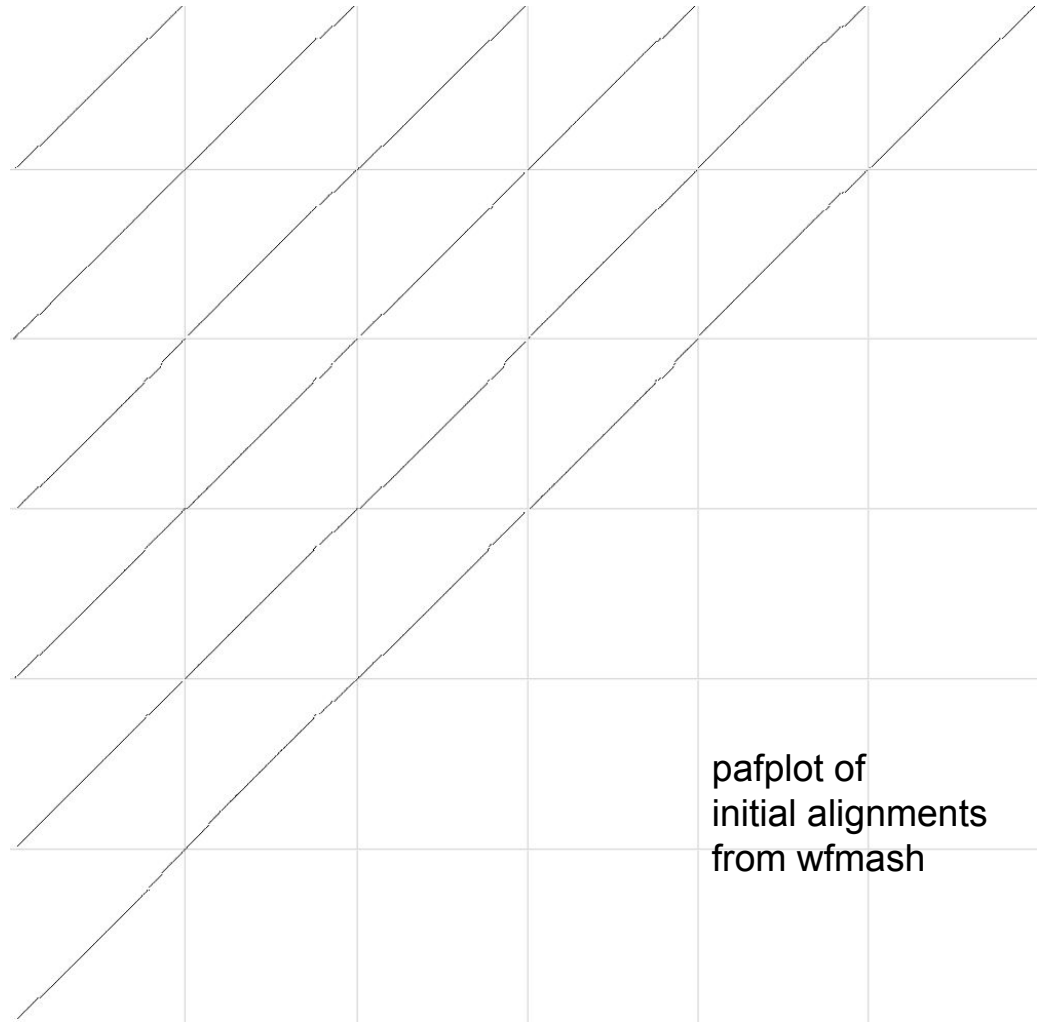
1. A few genes from HLA-D (MHC class II) in humans — getting started
 - a. <https://github.com/pangenome/pggb-workshop/tree/evomics2024>
2. Yeast chromosome 6 — scaling up
 - a. ~/workshop_materials/pangenomics/cerevisiae.chrV.fa
 - b. you will want to apply samtools faidx to this... pggb will warn you
3. Whole yeast chromosomes — looking at chromosome variation
 - a. ~/workshop_materials/pangenomics/cerevisiae.pan.fa.gz

Example: yeast chromosome 6

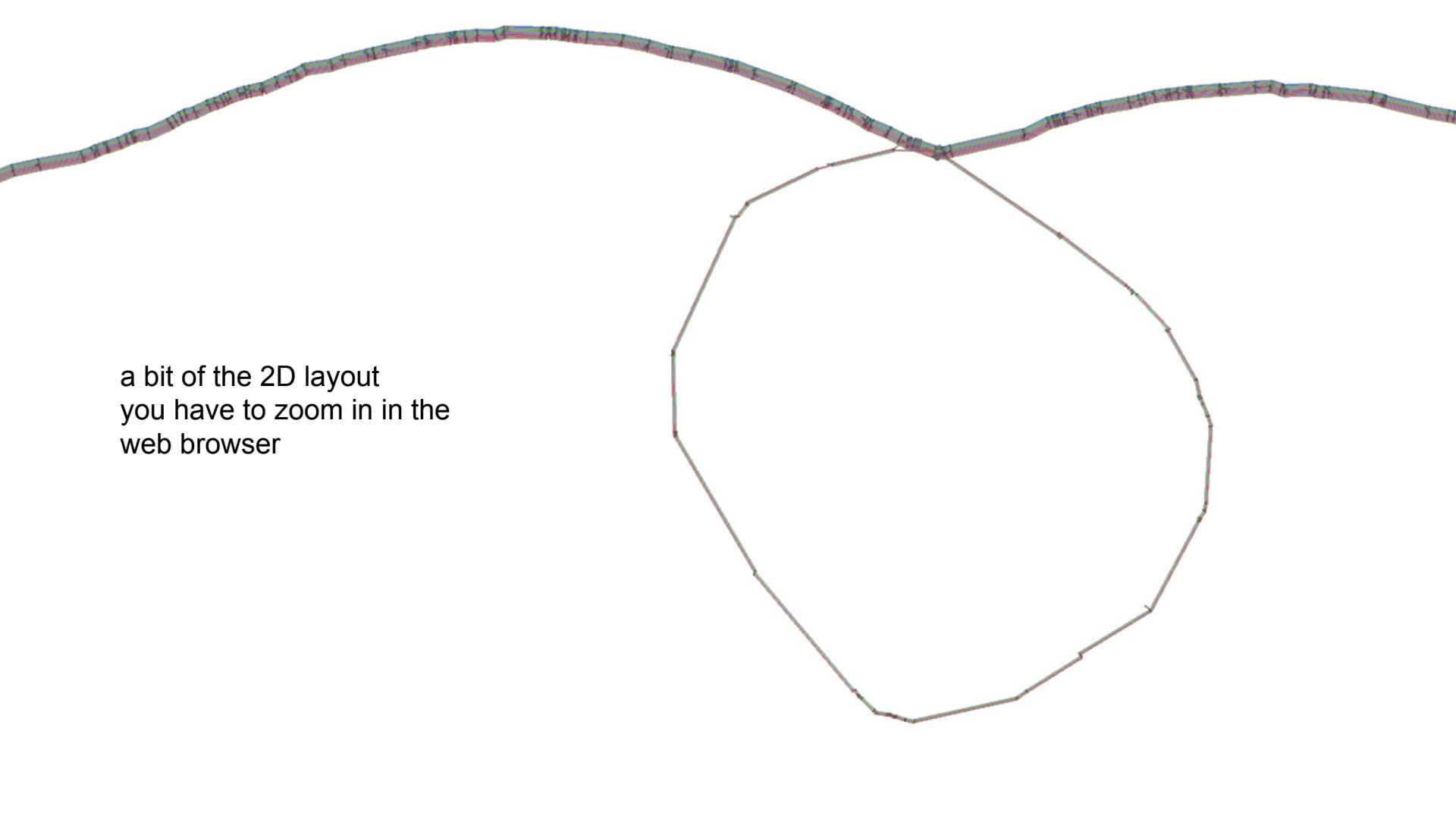
Yue, JX., Li, J., Aigrain, L. et al. Contrasting evolutionary genome dynamics between domesticated and wild yeasts. Nat Genet 49, 913–924 (2017). <https://doi.org/10.1038/ng.3847>



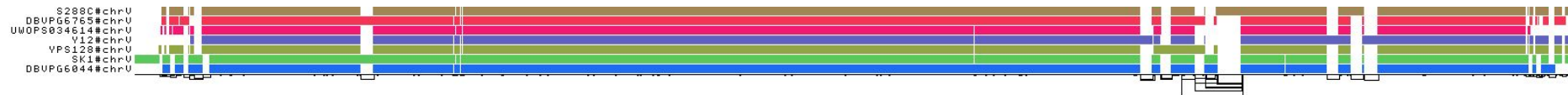
Cladogram of the *S.c.* clade



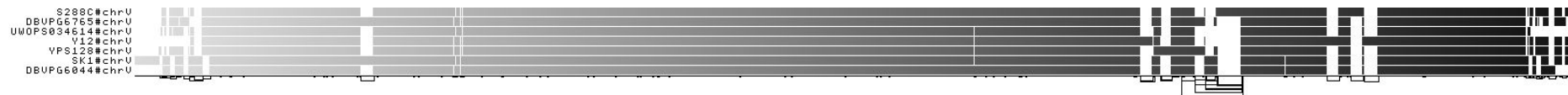
pafplot of
initial alignments
from wfmash



a bit of the 2D layout
you have to zoom in in the
web browser



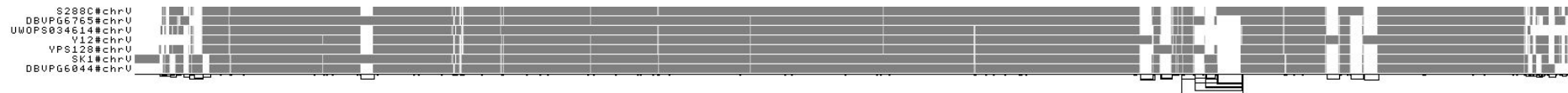
path view



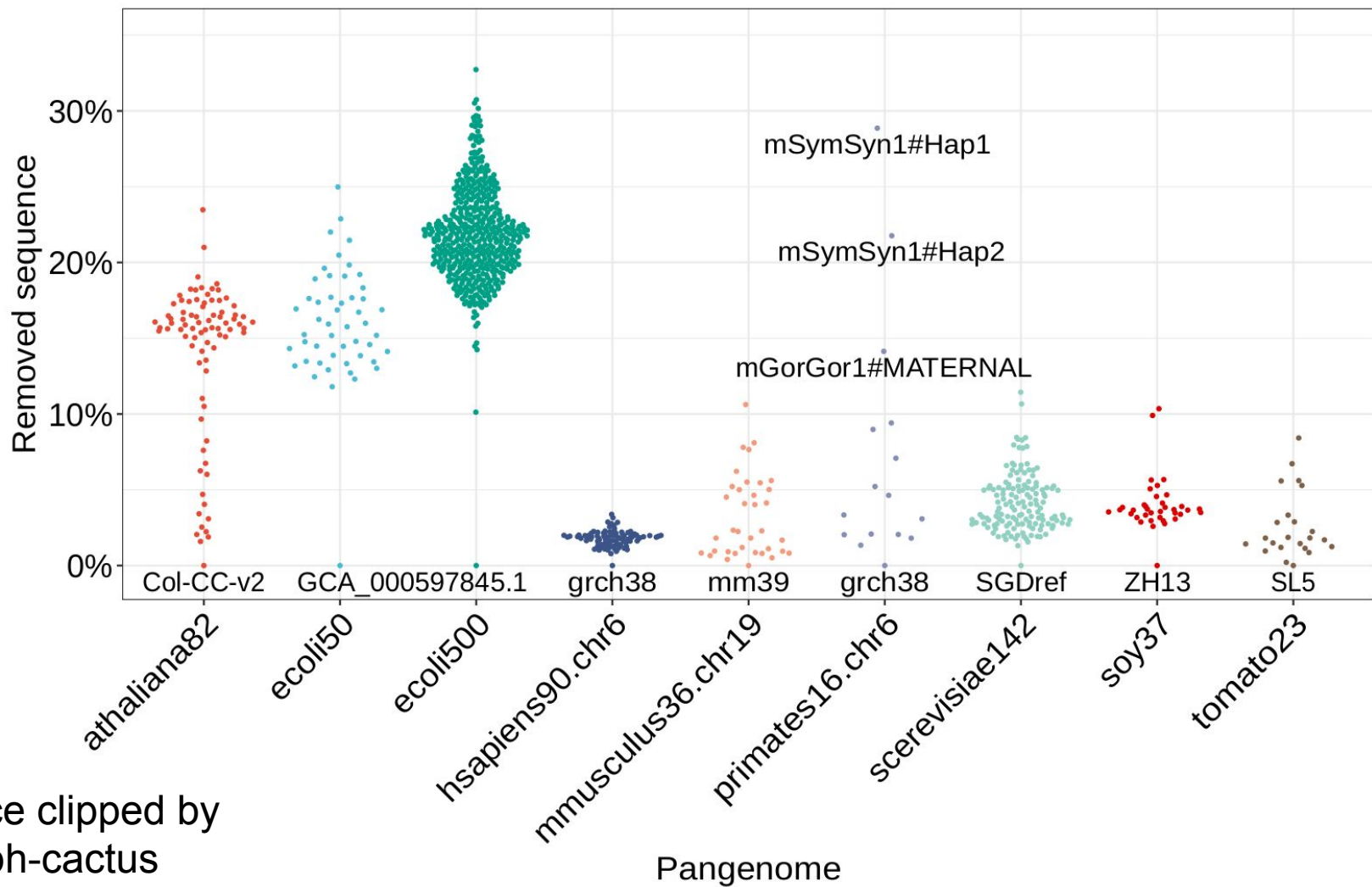
position view



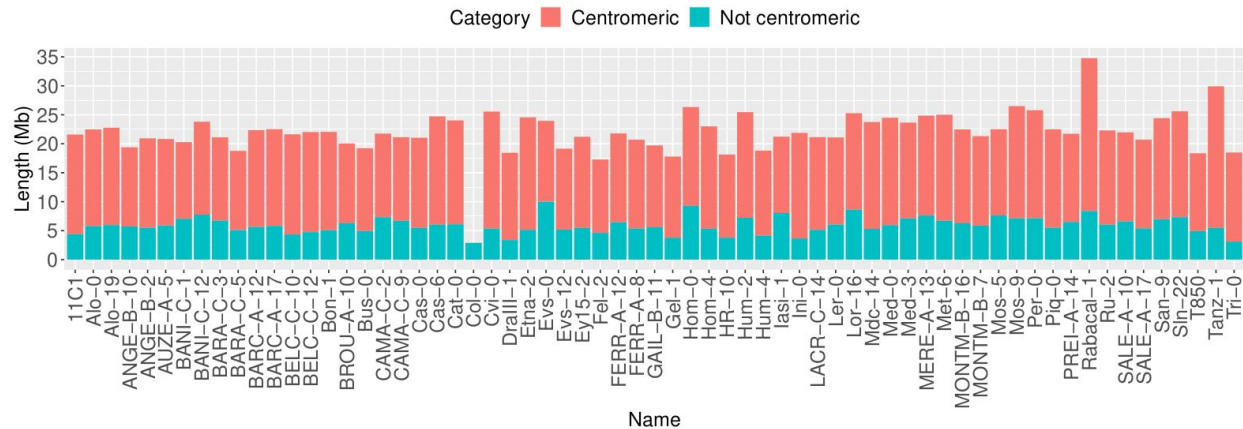
inversion view



copy number view



annotation of
clipped sequences
in minigraph-cactus
for *A. thaliana*
pangenome



(b)

