## Building and understanding pangenomes

#### Erik Garrison

University of Tennessee (UTHSC), Memphis

@Workshop on Genomics, Český Krumlov May 20, 2023







#### Genomic





Reference model



Genomic

Δ: new genome; R: reference genome. Figure from <u>Eizenga et al., 2020</u>.







Genome (FASTA)

\*
HHIRLIHHINIGHGITIIGHNIGEGIFIIGHRIGEIFAIGGAIGHGIDGHGFITIIICHITIITINDIBIIGHDIBIIGHDIBIIBGBEIHIBIITIGI
grandice66, 00000001
Grandichiferandichigera • MIDHEIIIIGHRIITGHRIITGHRIFTIIFHRHITGHRIHIDHIIFFIIDAIDIFHGDIEACCIIFHG-JHICTEBFIIIECCIIICOBBITF99 @smmpl-deds.080000002 MAGAAIXTATATTSTICTTETTSTATTAGCAGGGCTSATAGGGTTGAGAGTGGAGACTGATASTGATCTTGAGTAACTGTTGGCAATACTTCCCCTTTA a 200 chr0.fa



# chell 13995 chell 13996 chell 13997 chell</t

Variation (VCF)

# **Reference** model



#### Genomic

CAAATAGACTTCCCCATAACACAAAGCCATCCTGAAAAGTTTTGTTCATTTTAGAAGAAAAAATTTTAAAACCTGAGCAC AAAAGTTCCATAATGAGTAATCAAATTTTATTTCCAATTAAATATGTTATCACCCAGTAGTATGCCATACCAGTTTCTAA TTCATACTTCAATTATCTTCTAATTTAAATTAATGACTATAATTGCTGTTATAAAAACAACAGCTCTATAGCCTGCTATTC AGACCAGTAAATAAGAGTTTTAAGGGCTTGTGATAGCAAATGAAGTTTCTTATTGGATTTTAAGAAAAATTTTTATAAAAA TATGTGAGGTTATTCAATAGAATCACATTTAATTTGCCAAGCATTTTGCAGAATGCCTAGGACTATGTAAGAAGTATTAA ATTTGCAAGCCCTTTGAATAGTTGTAATTTAAAGATAAAAATTGGTTTAATACCAGACAAAGATAGAAGCACAAGTTAGG TTATTAGAGAATTTAGCCAGTGTATCAGTTTGTATCGTAAGTCATTGGCAAGAACAACGTGTACTTTTCTGTCACCTCCC AACTAGCTATGTTTTGAGCAGTAGGAATATTTAATACCCCTTCCTCCCATTTTTCCTTTGTGTTGTCCAAATTCTGACAA CTCTACTGCCAGATAGCTCAGGGCAAAAATGATAAAGTTCAAGTTAAGAAGGCTCTGCAGTGTTCTCAGTTCTCCTCTGG TGAAAGAGGAGAAAGGTTGTGTTTAATTATGAATCTGGGATTTCCAAAACTTTACCCATGCCCTGCCCTGTCCCCTCATTA GCATGAAGCTGTTATTTAAATAGTTCAGCAATAACGACTTTAGTAGCCTCCCTAGGTTAAAAAGATTGAAATTAAATGTG TTTATCTATTGTTCTACTATTCAGTTACCTGATTATAAAATCAAAGATTATTTCATGAAACTCAGTACCCCTTCAGGGAA AAATGCTTACCCAACTTCTATTCAAAATATTTGCGCCCAGTAGTTCTGATATGACCCAAGCAGAGTTCACACATTATTAAT CTACTCCTTTCAGTCTTCTAGATGTGTTTCCTCCAAAATCTACCAGATTCTCAAATAATTTCAGGAACTTTCTCCAGAAC AGAAACAAGGTTGTTACTGATACCAACTTTGTCTCCAAACATGGGGAAGATTATCATTGGAAAGATCTATTGATGACCTA TAATACATAGTTGGAACTGTTTATCCACAGAAGTATTCCCCCAAGAATCAACCACAGAGCCAAGATGGAGCTTATGTCATT GTTATGCATACTTCTTTTACGGCTTGTGAGGGCAGGTCATACTATTCTGATTTTACAACTGAGACCCCAAGGAACCTGAGT GACTTCTAGGCTCCATTATGTCAAAAAAAACTCAAATGTGAGGCTTTGCCTACACTGAGAAACAGTAGTTCAAGAAACGG TGCCCTGGTTCTGTTAAAATAATCTGAGAGTTATGTGGTAAGTAGTTGAGAGTGAATAGGGTAGCTTTGAGAGGTGACAG CGTGCTGGCAGTCCTCACAGCCCTCGCTGGCTCCAGGCGCCTCCTCGCCTGGGCTCCCACTTGGCGGCACTTGAGGAG CCCTTCAGCCCACCACTGCACTGCGGAGCCCCTTTCTGGGCTGGCCAAGGCCGGAGCCGGCTCCCTCAGCTTGCAGGGA GCTTGGCGGGCCCCGCACTCGGAGCAGCCGGCCAGCCCTTCCAGCCCCAGGCAATGAGAGGCTTAGCACCCGGGCCAGCA GCTGCGGAGGGTGTACTCCGTCCCCAGCAGTGCCAGCTCACAGGCGCTGCGCTCAATTTCTCACCGGGCCTTAGCTGCC TTCGCGCGGGGGGGTGCTCGGGACCTGCAGCCCGCCATGCCTGAGCTCCCATCGCGCCCCGTGCGCCCCGAGC CTCCCCGATGAGCACCACCCCCTGCTCCACGGCGCCCAGTCCCATCGACCACCCAAGAGCTGAGGAGTGCGGGGGGCGCACGG CGCGGGACTGGCAGGCAGCTCCACCTGCAGCTCTCGTGCGGGATTCACTGGGGGAAGCCAGCTGGGCTCCTGAGTCTGGT GGGGACGTGGAGAACCTTTATGTCTAGCTCAGGGATTGTAAATACACCAATCGGCACTCCGTATCTAGCTCAAGGTTTGT AAACACAACAACAACACCCTGTGTCTAGCTTAGCGTTGTGTAGCGCACCAAGCCACACTCTGTATCTAGCTACTCTGG GGGGCTTTGGAGAAACCTTTGTGTCCACACTCTGTAGCCAGCTAATCTGGTGGGGACATGGAGAACCTTTGGGTGTAGCTC AGATAAGAGCATAAAAGCAGGCTGCCTGAGCCAGCAGTGGCAACCCGCTTGGGTTCCCTTCCACACTGTGGAAGGTTTG GAAGGTCTGCAGCTTCACTCCTGAAGCCAGGAGACCACGAGGCCCACCAGGAGGAACCAACTCCAGAAGCGCCGCCT ACTCCGAACACCCGAACATCAGAAGGAACAAACTCCAGATGCGCCACATTAAGAGCTGTAACACTCACCGCGAGGGTC CCTGGCTTCATTCTTGAAGTCAGTGAGACCAAGAACCCACCAATTTTGGACACAGTTTGACAATAAATTTACACTCAAAT ATCTCTAAGGAATCAAACTTACAGATTAATAATTAGTAATCAGGTCACGTAAAGTAAATTATAAAAGAGCATTGATACCA AGATTGGCAGAAAGTTTTTTGTGTGACAAAACCAAGTTTTGGCTAAGATACACACTGCTGATGGGAGTCTAAATTGCTGT TCTAGTGAAACTCTTGAACGTGTGCGTCCAAAGACATTTATAAACATGATCTTAGTAGTATTGCTTTTAGTAGCAAATTC TGGAAACATCCCAAATGTCTATCAATAGTGGAATTGATTTGAAAGGGGTGTGGAATGGTAATATAATGGAATAGCCTACA GCTGTTTAGATAAAGGAACTCCAATTAAACATACCAAAGATACATTTCAAAAAACAAGACGTTGAAAGGAAAAAAGTC ATCAAAACAATACAACAATCTACCACATTTTTATAAATTCTCAAAATATGCAATATTAAACATGCATTATTTAGGGAG GCATTCAATGTAGCAATGCATTTTTAAGAGGCTGGGATGATAAATGTAAAATTCAGAACAGGTATTATCTCTGGGAACAG GAAGAGGAGGATGCAGTGTTGGGAAGAAATACATATAAGTACAGCAGTAGAGGCAGACTTTTTTTCCTTTTCCTTTTCC CTTTATTTTCCTAGCTTTCTTTTCTTTAGCTATGGTATTTCTTTAGCTATGGTATGGTATGGTATGGTATGTACTTTCCATA GGTACATGTGCAGATTTGTTATATAGGTAAACTTGTGTCATGTGGTTTTGCTGCACAAATTATTTTCTCACCCAGGTATT TTTCTGTTTCTGCGTTAGCTTGCCAAGGATAATTGCCTCCAGCTCCATATTCCTGCAAAAGACATAATTTTGTTCCTTTA TATGGTTGCATAGTATTCCATGGTGTATATGTACCGCATTATCTTTAGCCAGTCTATCATTGATGAGCACTTAGGTTGAT TCCATGTCTTCGCTCTTAACATTTTTAAACAGTCTCTGAGTAGAATAGGGTAGGCTGGTGTAAGGAATTACTGTTTTTAA ATTTCTGGGAAGATTTGCAAGAATCTGTGGCAGTTGAGAGTAGGTTCACTTTCGCTTTATTGTGTAATTTATTGTATTTT TCATTTAATTGTACTTTGTAAACTAAATATTTATTGTATATTTTACTTCATTTTTAATTGCCATATGCAGCTTTAATT GGTCCCTCCGGCTTCCACCCCGCCCCTGCGCTCACCTGCCCGCGCGCCCCTCCCGGGGACCCGGGGCCCATGGACAC ATACACCCAGCCCTGCTGTCCCGCGCGCCAGCTCACCAGCCCTACCCAAGGGACATCATTCACGCCTGGGCGCCTCCGCC GGGCTCCGGGAGCCCAAGGTCGCGGCTGGGCCAGCGCTGAGCGTCAGAGGACGAGAGCAGGGGCCTCCCCGGTCGCCCCA

Genome (FASTA)

#### We cannot update a linear reference sequence



##sour	ce=mutat	rix popu	lation g	enome si	imulator					
##seed:	=1373972	756								
##refe	rence=ch	rQ.fa								
##phas:	ing=true									
##comma	andline=	mutatrix	-S samp	ole -p 2	-n 100 c	hrQ.fa				
##filte	er="AC >	Θ"								
##INFO:	= <id=typ< td=""><td>E, Number</td><td>=A, Type=</td><td>String,D</td><td>escripti</td><td>on="Type</td><td>of each allele (snp, ins, del, mnp, complex)"&gt;</td><td></td><td></td><td></td></id=typ<>	E, Number	=A, Type=	String,D	escripti	on="Type	of each allele (snp, ins, del, mnp, complex)">			
##INFO:	= <id=na,< td=""><td>Number=1</td><td>,Type=In</td><td>teger, De</td><td>escriptio</td><td>n="Number</td><td>of alternate alleles"&gt;</td><td></td><td></td><td></td></id=na,<>	Number=1	,Type=In	teger, De	escriptio	n="Number	of alternate alleles">			
##INFO:	= <id=len< td=""><td>,Number=</td><td>A, Type=1</td><td>Integer, D</td><td>escripti</td><td>on="Lengt</td><td>th of each alternate allele"&gt;</td><td></td><td></td><td></td></id=len<>	,Number=	A, Type=1	Integer, D	escripti	on="Lengt	th of each alternate allele">			
##INFO:	= <id=mic< td=""><td>ROSAT, Nu</td><td>mber=0,1</td><td>Type=Flag</td><td>,Descrip</td><td>tion="Gen</td><td>erated at a sequence repeat loci"&gt;</td><td></td><td></td><td></td></id=mic<>	ROSAT, Nu	mber=0,1	Type=Flag	,Descrip	tion="Gen	erated at a sequence repeat loci">			
##FORM	AT= <id=g< td=""><td>T, Number</td><td>=1, Type=</td><td>String,D</td><td>escripti</td><td>on="Genot</td><td>ype"&gt;</td><td></td><td></td><td></td></id=g<>	T, Number	=1, Type=	String,D	escripti	on="Genot	ype">			
##INFO:	= <id=ac,< td=""><td>Number=A</td><td>,Type=In</td><td>teger, De</td><td>escriptio</td><td>n="Total</td><td>number of alternate alleles in called genotype</td><td>s"&gt;</td><td></td><td></td></id=ac,<>	Number=A	,Type=In	teger, De	escriptio	n="Total	number of alternate alleles in called genotype	s">		
##INFO:	= <id=af,< td=""><td>Number=A</td><td>, Type=FI</td><td>loat, Desc</td><td>ription=</td><td>"Estimate</td><td>ed allele frequency in the range (0,1]"&gt;</td><td></td><td></td><td></td></id=af,<>	Number=A	, Type=FI	loat, Desc	ription=	"Estimate	ed allele frequency in the range (0,1]">			
##INFO:	= <id=ns,< td=""><td>Number=1</td><td>,Type=In</td><td>teger, De</td><td>escriptio</td><td>n="Number</td><td>of samples with data"&gt;</td><td></td><td></td><td></td></id=ns,<>	Number=1	,Type=In	teger, De	escriptio	n="Number	of samples with data">			
##INFO:	= <id=an,< td=""><td>Number=1</td><td>, Type=In</td><td>teger, De</td><td>escriptio</td><td>n="Total</td><td>number of alleles in called genotypes"&gt;</td><td></td><td></td><td></td></id=an,<>	Number=1	, Type=In	teger, De	escriptio	n="Total	number of alleles in called genotypes">			
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO FORMAT sample001 sample002			
chrQ	1252		C	A	99		AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	010	1 0
chrQ	3646		Т	TC	99		AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=ins	GT	010	011
chrQ	6283	2	C	т	99		AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	010	011
chrQ	7412		C	T	99		AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	010	110
chrQ	7935		Т	C	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	010	110
chrQ	8131		Т	C	99		AC=2; AF=0.5; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	011	011
chrQ	8682		AA	TG	99		AC=1;AF=0.25;AN=4;LEN=2;NA=1;NS=2;TYPE=mnp	GT	010	011
chrQ	10926		T	C	99		AC=1; AF=0, 25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	010	011
chr0	11921		G	GTT	99		AC=1; AF=0, 25; AN=4; LEN=2; NA=1; NS=2; TYPE=ins	GT	011	010
chrQ	12955		Т	G	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	00	10
chrQ	13808		т	TG	99	2	AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=ins	GT	10	00
chrQ	15271		A	G	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	01	010
chrQ	15407	1	A	C	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	10	010
chrQ	16486		C	G	99		AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	10	010
chrQ	16563		Т	A	99		AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	011	010
chrQ	16748		GTT	G	99		AC=1;AF=0.25;AN=4;LEN=2;NA=2;NS=2;TYPE=del	GT	0/0	0/1
chrQ	17697		G	C	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	010	011
chrQ	19568		A	G	99		AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	011	010
chrQ	20750	2	G	A	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	0 0	10
chrQ	21532		Т	C	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	10	010
chrQ	22291		C	Т	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	01	010
chrQ	23193		G	A	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	00	01
chrQ	23954		CTAA	TTAA	99		AC=1; AF=0.25; AN=4; LEN=4; NA=2; NS=2; TYPE=mnp	GT	0/0	0/1
chrQ	24467		C	Т	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	010	0 1
chrQ	26100		G	A	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	01	00
chrQ	29654		Т	A	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	1 0	0 0
chrQ	30062	2	Т	C	99	1	AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	01	00
chrQ	31790		A	G	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	01	010
chrQ	32792		Т	C	99		AC=3;AF=0.75;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	1 1	0 1
chrQ	33376		CC	C	99		AC=2;AF=0.5;AN=4;LEN=1;NA=1;NS=2;TYPE=del	GT	1 0	0 1
chrQ	33403	*	Т	C	99		AC=2;AF=0.5;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	01	011
chrQ	33802		A	G	99		AC=2;AF=0.5;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	01	1 0
chrQ	34450		C	Т	99		AC=4;AF=1;AN=4;LEN=1;NA=1;NS=2;TYPE=snp GT	1/1	111	
chrQ	34716		G	A	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	010	1 0
chrQ	35484	2	G	A	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	00	10
chrQ	36547		G	A	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	00	01
chrQ	38015	1	Т	A	99	1	AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	00	10
chrQ	38281	1	т	C	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	01	00
chrQ	40467		A	G	99		AC=1;AF=0.25;AN=4;LEN=1;NA=1;NS=2;TYPE=snp	GT	10	010
chrQ	40581	1	A	G	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	00	1 0
chrQ	40601		A	Т	99		AC=1; AF=0.25; AN=4; LEN=1; NA=1; NS=2; TYPE=snp	GT	010	011

Variation (VCF)

AC=1:AF=0.25:AN=d:I FN=1:NA=1:NS=2:TYPE=snn

GT BII BIB

A QQ

chr0 43268









#### 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<sup>†</sup> history.

(<sup>†</sup>for genomics)

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









Group B Streptococcus assemblies from 2002

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  $\alpha$  determines whether the pangenome is open ( $\alpha \le 1$ ) 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

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



#### Pangenome research timeline

2000-2010s: counting genes

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

2020s: complete assemblies (T2T pangenomes)



#### <u>*Wait!*</u> You can align sequences to graphs?

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



#### 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).

#### $GCSA2 \rightarrow$

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

This seeds alignment to the graph.





reads aligned to a variation graph

and long reads too!

a pacbio read vs. a yeast graph:



#### nature biotechnology

Letter | Published: 20 August 2018

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

Erik Garrison <sup>™</sup>, 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 <sup>™</sup>

Nature Biotechnology 36, 875–879 (2018) Download Citation ±

#### vg resolves reference bias at known indels in HG002



Size of deletion (negative) or insertion (positive)

50x 2x150bp Illumina sequencing of HG002





Home About Articles Submission Guidelines

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)





12880A (Iron Age Britain)



#### 15577A (Anglo-Saxon Britain)



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

#### Rui Martiniano

#### **PLOS COMPUTATIONAL BIOLOGY**

🔓 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 🔤

Α





#### Β

Pangenome variation graph (VG)



#### Retrieved motif occurrences and haplotype frequencies

Sequence	Log-odds score	<i>P</i> -value	<i>q-</i> value	Reference	Haplotype frequency
GGGCCAGCAGGGGGGGCGCTG	28.22	7.51e <sup>-12</sup>	3.86e <sup>-6</sup>	non ref.	32
TGGCCAGCAGGGGGGGCGCTG	26.16	3.12e <sup>-10</sup>	7.72e <sup>-6</sup>	ref.	5063
GGGCCAGCAGGGAGCGCTG	19.43	1.71e <sup>-7</sup>	1.73e <sup>-4</sup>	non ref.	1



Method Open Access Published: 12 February 2020

## Genotyping structural variants in pangenome graphs using the vg toolkit

<u>Glenn Hickey</u>, <u>David Heller</u>, <u>Jean Monlong</u>, <u>Jonas A. Sibbesen</u>, <u>Jouni Sirén</u>, <u>Jordan Eizenga</u>, <u>Eric T. Dawson</u>, <u>Erik Garrison</u>, <u>Adam M. Novak</u> & <u>Benedict Paten</u>

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







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

B RESEARCH ARTICLE GENOMICS

## Pangenomics enables genotyping of known structural variants in 5202 diverse genomes



#### vg giraffe: approach



Haplotype restricted sequence graph gapped alignment


### vg giraffe is accurate enough



#### vg giraffe is <u>very fast</u>

#### (A) 1KG/GRCh37 NovaSeq 6000 Runtime

VG-MAP paired			
VG-MAP single			
Bowtie2 paired			
Bowtie2 single			
Giraffe full single			
Giraffe full paired			
BWA-MEM paired			
BWA -MEM single			
Giraffe sampled paired			
Giraffe sampled single			
Giraffe primary paired			
fast Giraffe sampled paired			
Minimap2 paired			
Giraffe primary single			
fast Giraffe sampled single			
Minimap2 single			
HISAT2* paired			
HISAT2* single			
0 10 20	30	40 5	
Runtime (hours)			

#### (C) 1KG/GRCh37 NovaSeq 6000 Memory

GraphAligner	Out of memory
Giraffe full paired	
Giraffe full single	
Giraffe sampled single	
fast Giraffe sampled single	
Giraffe sampled paired	
fast Giraffe sampled paired	
Giraffe primary paired	
Giraffe primary single	
VG-MAP paired	
VG-MAP single	
- Minimap2 single	
Minimap2 paired	
BWA-MEM paired	
HISAT2* paired	
HISAT2* single	
BWA-MEM single	
Bowtie2 paired	
Bowtie2 single	
) 20 4	0 60 80 100 Memory (GB)

### vg giraffe improves variant calling



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



# Building the human pangenome

#### 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

PANGENOME

Check for updates



#### 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 7 pangenome (reference) graphs...

# Minigraph

Heng Li (author of bwa, samtools, minimap2, miniprot, and many other tools).

Idea: build the graph from the reference, progressively.

Only SVs > 50bp.

Uses a model similar to POA, but capable of dealing with structural variation, and using the minimizer chaining concept from minimap2.



### Minigraph

Heng Li (author of bwa, samtools, minimap2, miniprot, and many other tools).

Idea: build the graph from the reference, progressively.

Only SVs > 50bp.

Uses a model similar to POA, but capable of dealing with structural variation, and using the minimizer chaining concept from minimap2.



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

### Minigraph-Cactus

Add base-level variation onto the minigraph.

Uses tooling from Cactus to drive local alignment across pieces of the graph.

Makes a variation graph, but the graph has the structure of the minigraph model with some differences.

Clipping is used to make the graph easy to align to with vg giraffe.





"Hard-mode" pangenomes.

All-to-all = quadratic. ^That's an *exascale* matrix of chr6 in all great apes.

#### 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  $\rightarrow$ 



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



Wen-Wei Liao

#### The graphs accurately characterize structural variants



Variant length (bp)

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

Wen-Wei Liao

And they show very similar pictures of the pangenome

Major differences are in repetitive sequences, which tend to collapse in the PGGB graph and expand in the MC graph.

> We also look at pangenome growth via permutations, a nod to old-school pangenomics (g).



# Understanding the human pangenome

#### Highly polymorphic medically-relevant loci



View of locus that defines the Rh blood group system in the *minigraph-cactus* graph.

View of HLA/MHC in the PGGB graph

Shuangjia Lu

#### PGGB lets us look at all parts of the pangenome

(learned 2D visualization of PGGB HPRC chromosomes)





#### grch38#chr8

#### chr8's beta-defensin locus

5 Mbp polymorphic inversion flanked







The p-arm of chr8 features a region with one of the highest rates of diversity in the HPRC pangenome. (Here we see 100kb bins.)

Regionally, it's higher than chr6's MHC, although the MHC has higher peaks of diversity in HLA genes.

An implication is that this 4.7Mbp polymorphic inversion is driving diversifying processes.



https://github.com/chfi/gfaestus

DEFB\* genes are beta-defensins. The number of copies of this loop varies ~2-4 in different haplotypes.

FAM90A  $\rightarrow$  a gene of unknown function!



#### the MHC in chr6



#### C4A/B in pggb graph



https://github.com/chfi/gfaestus

Christian Fischer (UTHSC)

#### C4A/B in pggb graph



https://github.com/chfi/gfaestus

Christian Fischer (UTHSC)

# We learn that genome evolution is often <u>nonlinear</u>





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



The human pangenome illuminates the unknown



(Submit manus

# We finished the human genome

HOME > SCIENCE > VOL. 376, NO. 6588 > THE COMPLETE SEQUENCE OF A HUMAN GENOME

SPECIAL ISSUE RESEARCH ARTICLE | HUMAN GENOMICS

#### The complete sequence of a human genome



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

Filling 8% of the reference which was incomplete

> Of particular note, all of the acrocentric p-arms were assembled for the first time!



## But there are new mysteries...



2x

9q12

# But there are new mysteries...



https://github.com/lh3/minimap2/blob/67dd906a80988dddac c8c551623fdc75b0c12dd2/misc/paftools.js#L2605-L2719

# HPRC "misjoin" identification



Alignment graph of the misjoins. Every node is a contig and every edge represents the number of mapping between nodes. Alignment graph obtained with <u>pafnet</u> and visualized with <u>gephi</u>. Color code: **chr13**, **chr14**, **chr15**, **chr21**, **chr22**.

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


## <u>Chromosome</u> <u>communities</u> in the HPRC

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





## 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.





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



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

+ HG002 contigs >= 300kbps which map to acrocentrics

PanGenome Graph Builder (PGGB)



https://github.com/pangenome/pggb



# Untangling recombination in the acrocentric short arm(s)

*Untangling* extracts pairwise alignments from variation graphs.



*Untangling* extracts pairwise alignments from variation graphs.



*Untangling* extracts pairwise alignments from variation graphs.



Identify cut points in the graph

*Untangling* extracts pairwise alignments from variation graphs.



*Untangling* extracts pairwise alignments from variation graphs.



## Grounding the untangle against CHM13

We first pick a consistent set of query contig segments.

For each, we find its best mapping against each acrocentric chromosome.

We then merge this (grounding) with a "multiple" untangling result where each query contig segment picks the best *N* target (reference) segments.





Untangling the pangenome graph

chm13#chr14 chm13#chr15 chm13#chr21 chm13#chr22 Estimated identity 0.925 0.950 0.975 1.000 We look into the graph from the perspective of

the perspective of chromosome 13. Full information from pangenome plus reference annotations.



### chromosome 14









Linkage disequilibrium (LD) describes correlation between pairs of SNPs.

Higher LD (correlation) means lower rates of recombination.

We express LD in terms of its "decay" with distance. That's what we see in this plot.

LD decays faster in pseudo-homologous regions than elsewhere in the p-arms or in the q-arms. (chr15 lacks data)

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

Silvia Buonaiuto, Vincenza Colonna

### chromosome 13/14/21 Pseudo-homologous regions

The SST1 array is a GC-rich satellite (DNA repeat) derived from Alu elements.

It lies at the center of pseudo-homologous regions (PHRs) on 3 acrocentric chromosomes.

This is where we see recombination occurring in Robertsonian translocations.





## 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 turther 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 13al4g and 14a21g Robertsonian translocations

involving the common and nonrandom association of chromosomes 13 and 14, and 14 and 21 is discussed.

Chroo et al., 1988.



### distal region of NOR



**Figure 2.** Sequencing of DJ regions from individual human acrocentric chromosomes. (*A*) Schematic representations of DJ contigs. Hybrid and chromosomal identities are shown on the *left*, together with GenBank accession numbers. (*B*) The average percentage identity of 100-kb blocks, among all seven DJ contigs is shown *below* the WAV17 (HSA21) contig. (*C*) Alignment of DJ contigs demonstrating that they can be clustered into three groups, based on the sequence of their distal ends. The percentage identity of the most distal sequences within group 1 and 2 members is shown schematically on the *right*. (*D*) A schematic representation of indel distribution on the left arm inverted repeat. Hybrid and chromosomal identities are shown on the *left* (see Supplemental Fig. S7B for sequence alignments at break points).

#### https://doi.org/10.1101/gad.331892.119



#### 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).  $(Bar - 0.2 \ \mu m)$ 



TRANSACTIONS OF THE 70<sup>TH</sup> 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 🛛 A 🖂 • Theresa Naluai-Cecchini, BA

**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)



#### Article

# Recombination between heterologous human acrocentric chromosomes

https://doi.org/10.1038/s41586-023-05976-y	Andrea Guarracino <sup>1,2</sup> , Silvia Buonaiuto <sup>3</sup> , Leonardo Gomes de Lima <sup>4</sup> , Tamara Potapova <sup>4</sup> , Arang Rhie <sup>5</sup> , Sergey Koren <sup>5</sup> , Boris Rubinstein <sup>4</sup> , Christian Fischer <sup>1</sup> , Human Pangenome Reference Consortium <sup>*</sup> , Jennifer L. Gerton <sup>4</sup> , Adam M. Phillippy <sup>5</sup> , Vincenza Colonna <sup>1,3</sup> & Erik Garrison <sup>1</sup> <sup>⊠</sup>
Received: 15 August 2022	
Accepted: 17 March 2023	
Published online: 10 May 2023	
Open access	The short arms of the human acrocentric chromosomes 13, 14, 15, 21 and 22 (SAACs)
Check for updates	share large homologous regions, including ribosomal DNA repeats and extended segmental duplications <sup>1,2</sup> . Although the resolution of these regions in the first complete assembly of a human genome—the Telomere-to-Telomere Consortium's CHMI3 assembly (T2T-CHMI3)—provided a model of their homology <sup>3</sup> , 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 Consortium <sup>4</sup> (HPRC), we find that contigs from all of the SAACs form a community. A variation graph <sup>5</sup> 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 <sup>6,7</sup> . The pseudo-homologous regions in clude sequences that have previously been shown to lie at the breakpoint of Robertsonian translocations <sup>8</sup> , 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

## 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.





-T. -1 Ι 12 Ξ .... 1,2 1,1 ..... 201 P 171 I.A. 11:10 1 1.1 1 H. -H -..... -1 

### PanGenome Graph Builder





wfmash (biWFA)

### PanGenome Graph Builder





seqwish (unbiased graph builder)
#### PanGenome Graph Builder





smoothxg (graph normalization)







high-order *bidirectional* WFA (BiWFλ)



## smoothxg organizes & normalizes the graph

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



2d layout

identitv

(a)

(b) chm13#chr6: grch38#chr6. HG00438#2#J. HG00438#1#J.

HG01071#2#J. HG01071#1#J.

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

It uses the embedded genomes as references.

HG01952#2#J HG01952#1#J (c) chm13#chr6:... grch38#chr6. HG00438#2#J. position HG00438#1#J HG01071#2#J HG01071#1#J. HG01952#2#J HG01952#1#J. (f) (d) chm13#chr6: chm13#chr6: grch38#chr6. HG00438#2#J. HG00438#1#J. HG01071#2#J. orientation HG01071#1#J HG01952#2#J HG01952#1#J. (e) chm13#chr6: grch38#chr6 HG00438#2#J (g) HG00438#1#J HG01071#2#J HG01071#1#J. HG01952#2#J HG01952#1#J



copy number variation

odgi helps us understand the pangenome

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



#### Putting it all together!





## Test material today

- 1. A few genes from HLA-D (MHC class II) in humans getting started
  - a. https://github.com/pangenome/hprc-workshop/tree/evomics2023
- 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



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





these are all quite boring here...

## Conclusions

- Variation graphs solve a key conceptual problem in bioinformatics: joint representation of many genomes and their alignment.
- We can apply them to improve basic bioinformatic analyses.
- Building pangenomes benefits from an unbiased approach where many reference genomes are treated equivalently.
- Applying these methods to the human pangenome reveals patterns of recombination between heterologous acrocentric chromosomes.

#### 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)





https://doi.org/10.1111/j.1469-1809.1996.tb01180.x

Recombination between 15p and 22p.

95

### Zero-Mode Waveguide imaging of single polymerases



Time

#### https://doi.org/10.1016/S0076-6879(10)72001-2

## PacBio HiFi produces incredibly accurate long reads



https://www.pacb.com/wp-content/uploads/How-to-get-HiFi-reads v2.png

### Nanopore sequencing







## Preliminaries: Sequencing and assembling genomes

## So you want to look at genomes?

We observe genomes by sequencing fragments of DNA and assembling them computationally.



## Biology: "Sorry, it's not that simple"

Genomes evolve by copying. They are full of repeats.

Repeats make it hard to assemble genomes.

If our sequence reads are shorter than repeat sequences, the best we can do is the string graph, where repeats collapse.

<u>Before 2019</u>, our sequence reads were shorter than repeats.



Myers' string graph https://doi.org/10.1093/bioinformatics/bti1114

**Preliminaries**: Sequencing and assembling genomes completely!!

## Verkko: automated diploid assembly

Building on top of 3rd gen sequencing, we can automatically generate haplotype resolved, highly accurate (QV>50 or 1/100,000 base error).

This depends on two new bioinformatic approaches.

- long-kmer de Bruijn graphs
- sequence to graph alignment (GraphAligner) \*a pangenomic method



# Verkko assembly graph for HG002

This diploid, haplotype resolved assembly was a first. The easy sequencing (only two types) will facilitate many more in the near future!

