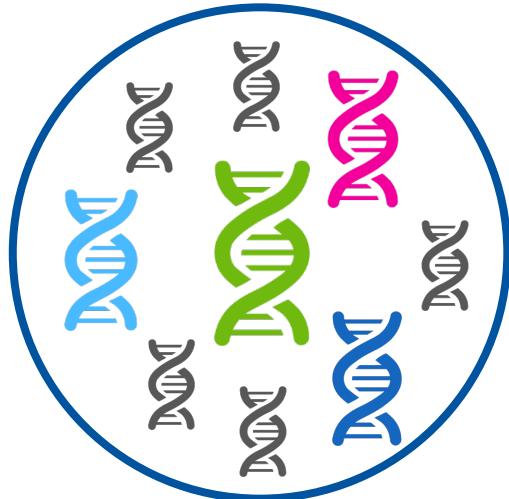


Pangenomics

Sandra Smit

12 Oct 2023, BioSB Algorithms for Genomics



Program today

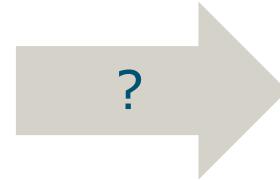
- 9.00 Lecture
- 11.30 Break
- 11.00 Computer exercises (rosalind or project)
- 12.30 Lunch
- 13.30 Guest lecture by Erwin Datema from Keygene
- 14.15 Lecture about PanTools
- 15.00 Break
- 15.30 Computer exercises (rosalind or project or PanTools)
- 17.30 Closing

Learning outcomes

- Understand the need for pangenomes
- Know what a pangenome is
- Be aware of various applications
- Get a handle on this fast-developing field
 - Tools
 - Concepts
 - Data structures
 - Algorithms

Comparative genomics

Genotype x Environment



Phenotype

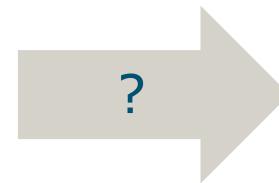


Comparative genomics

Genotype x Environment



Phenotype

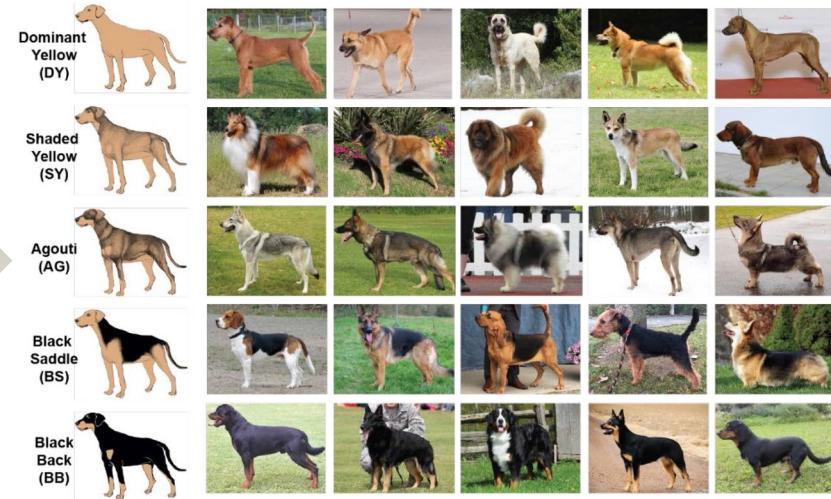


Comparative genomics

Genotype x Environment

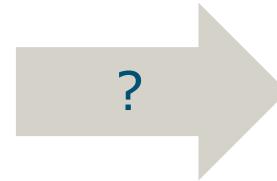


Phenotype



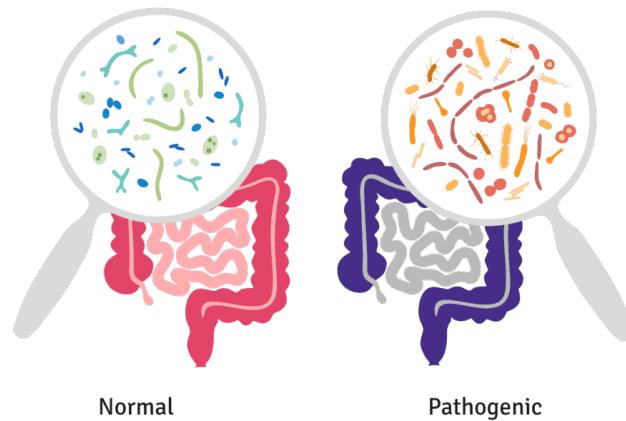
Comparative genomics

Genotype x Environment



Phenotype

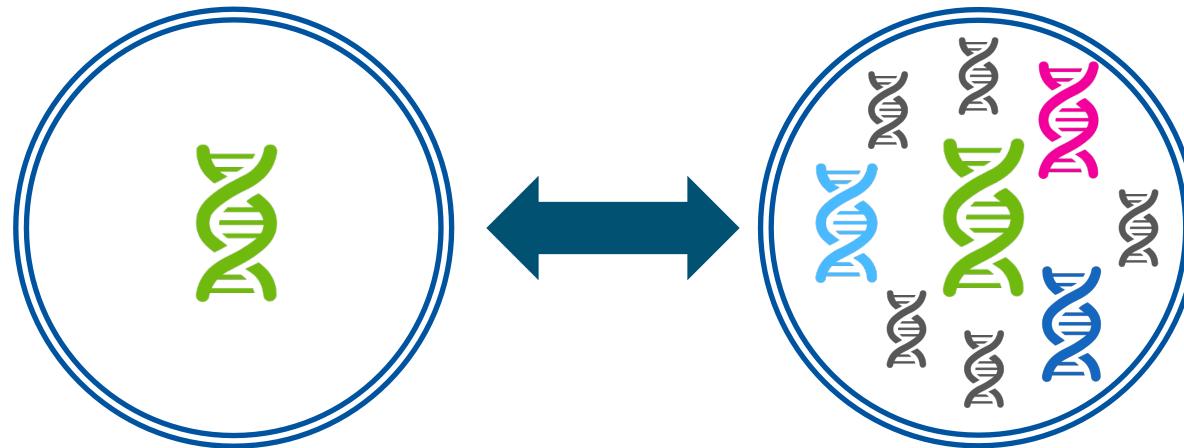
Gut Microbiome



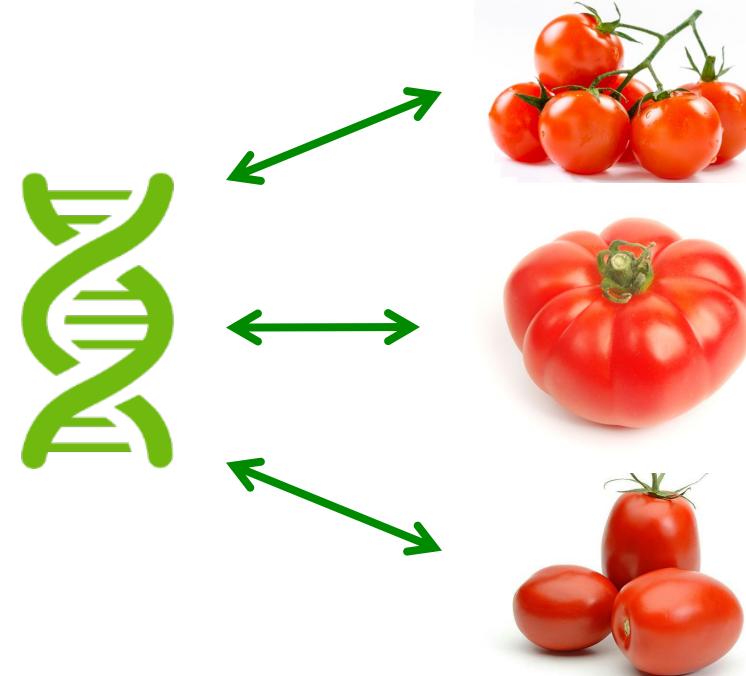
Normal

Pathogenic

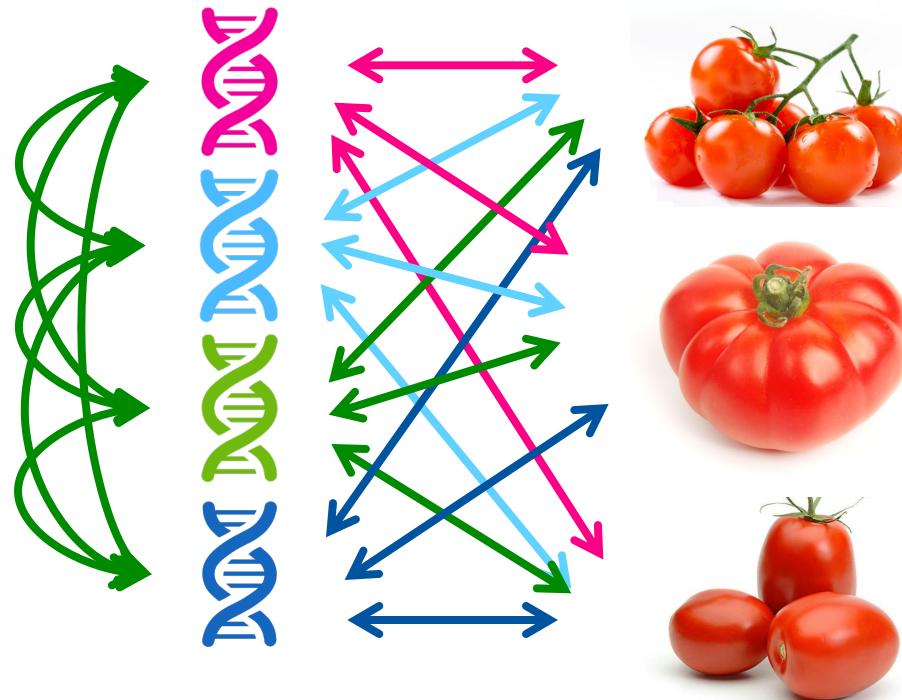
A paradigm shift...



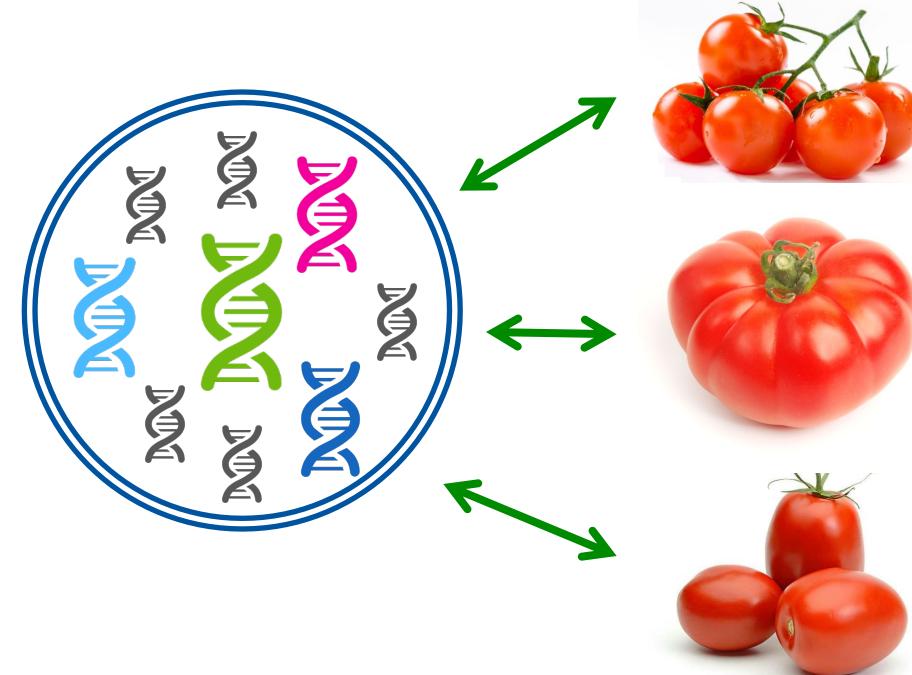
The reference genome



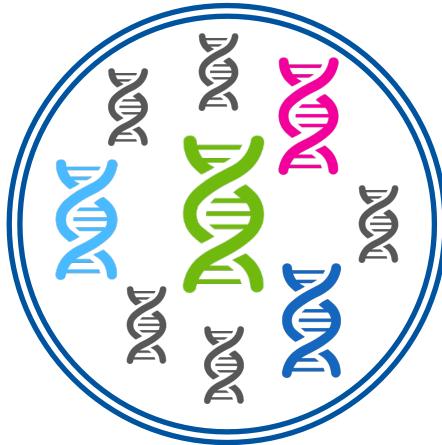
Variation and diversity



Pangenome



Pangenome



Single representation for
many genomes

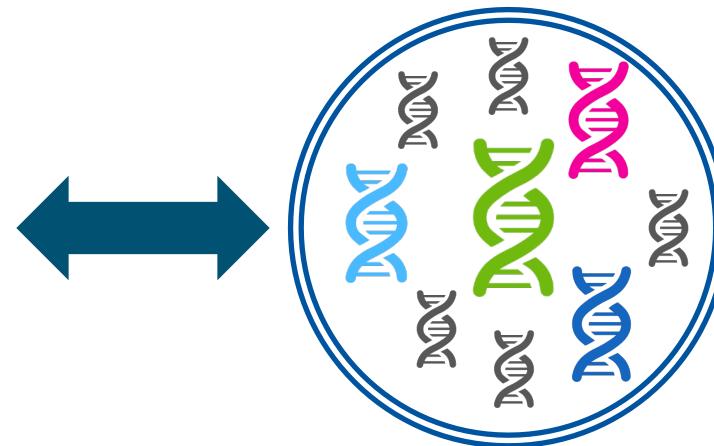
- Compress similarities
- Retain (all) variations

Why pangenomes?

Reference genome



Pangenome



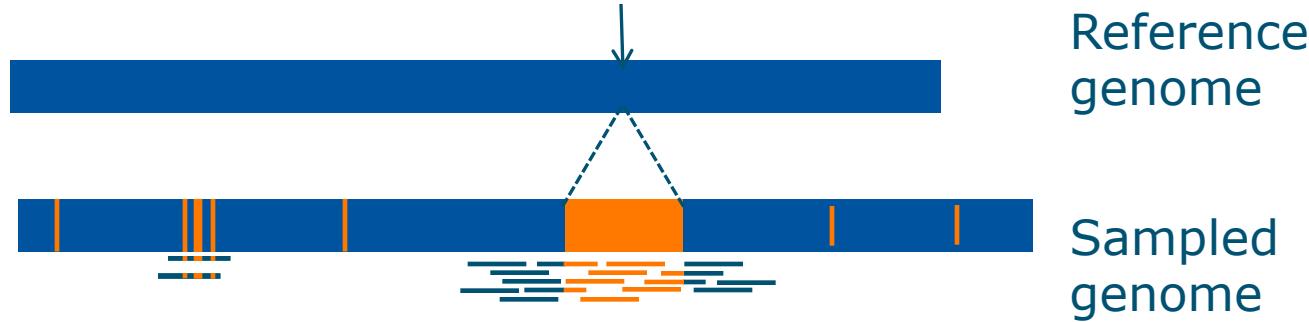
Reference genome benefits

Reference genome

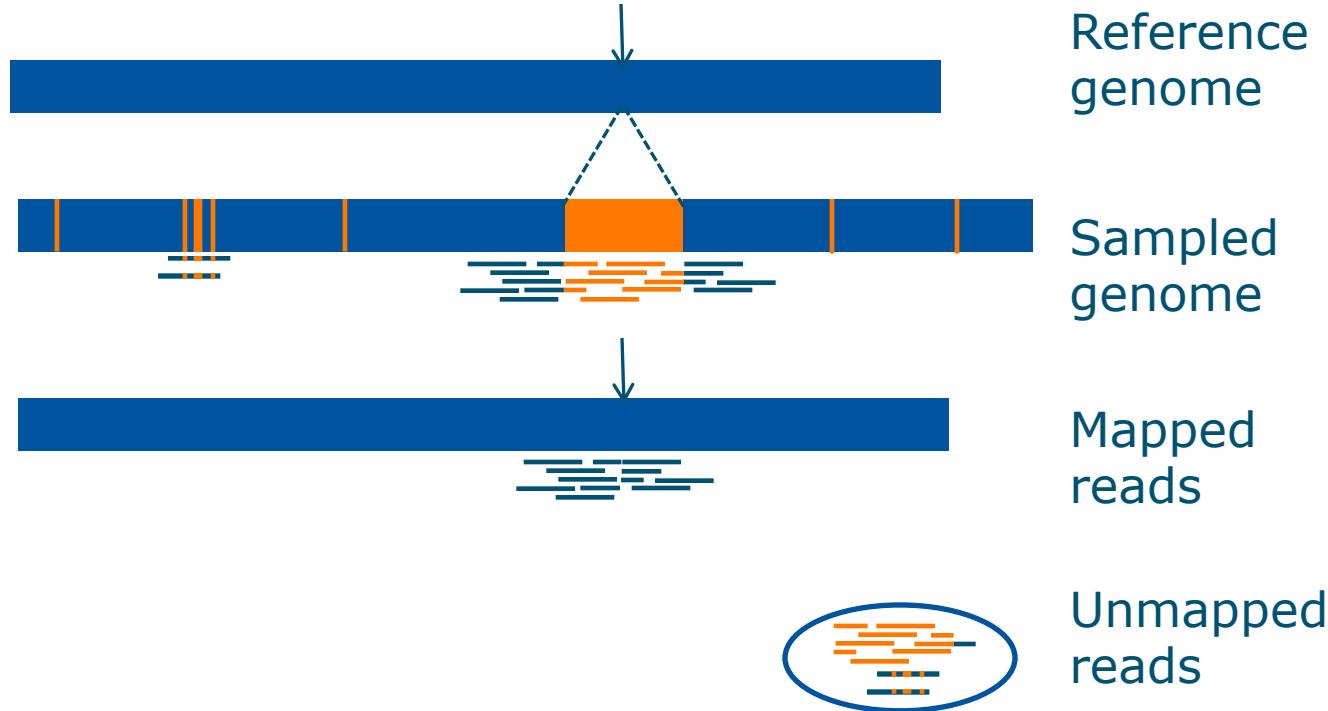


- Linear haploid representation
- Fast searching/mapping
- Indexing
- Sequence/string algorithms

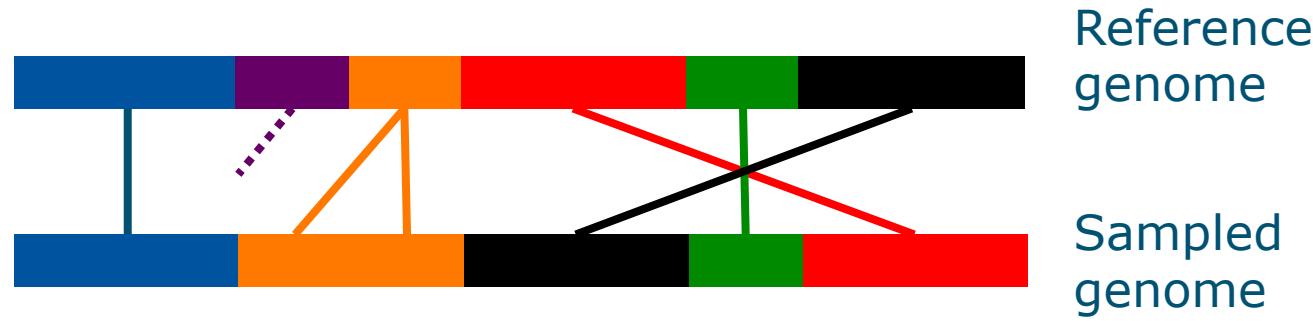
Reference bias



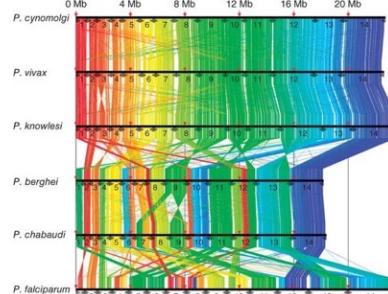
Reference bias



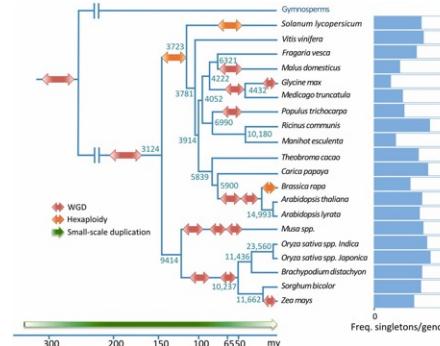
There is more than sequence variation



There is more than sequence variation



structural variations



ancient duplications

CLASS I - Retroelements (RNA intermediate)

Endogenous retrovirus (HIV, HERV)
► GAG PR RT RH IN ENV ►

Ty3/gypsy - BEL transposons

► GAG PR RT RH IN ►

Ty1/copia retrotransposons

► GAG PR IN RT RH ►

DIRS1-like retrotransposons

► GAG RT RH MT TR ►

Non-LTR retrotransposons (L1, L2)

GAG APE RT RT REL

Penelope-like retrotransposons

RT EN

Non-autonomous retrotransposons (SINES)

(Alu, B2, B4)

CLASS II - DNA transposons

DDE transposons (piggyBac, Mariner, Transib)

◀ TNP ▶

Helitrons

RPA Z PR REP-HEL IN

Polintons / Mavericks

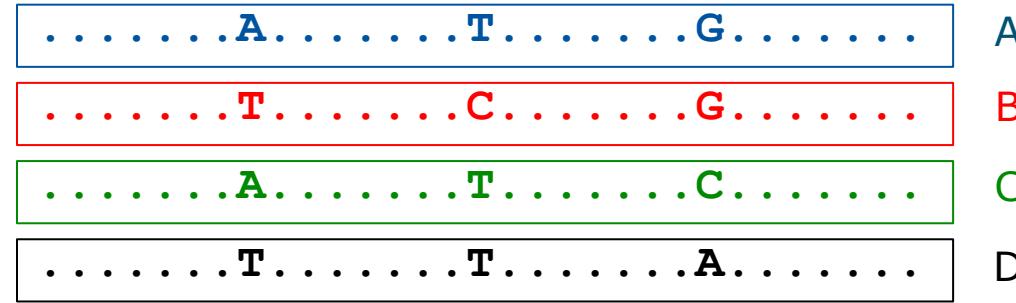
◀ IN Z CC PR B-POL ▶

Non-autonomous DNA transposons (MITEs)

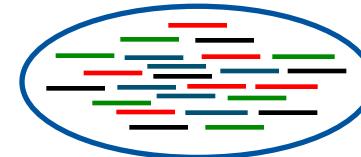
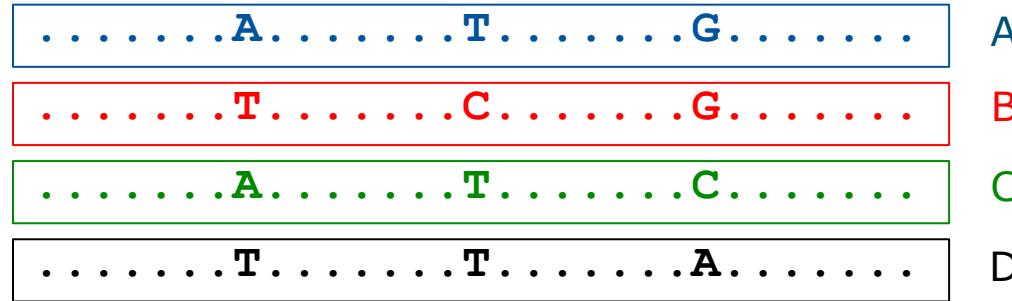
(MADE1)

repeats

Polyplody and heterozygosity



Polyplody and heterozygosity



Tetraploid sample

Polyplody and heterozygosity



.....**A**.....**T**.....**G**..... A
.....**T**.....**C**.....**G**..... B
.....**A**.....**T**.....**C**..... C
.....**T**.....**T**.....**A**..... D



Reference

.....**T**.....**T**.....**G**.....

SNPs

A **C** **C**
A

Why pangenomes?



Why pangenomes?

Reference genome



Pangenome



Catch **intragenomic** and
intergenomic variation

Why pangenomes?

The Pangenome: Are Single Reference Genomes Dead?

Researchers are abandoning the concept of a list of genes sequenced from a single individual, instead aiming for a way to describe all the genetic variation within a species.

Dec 1, 2016

CATHERINE OFFORD



22

<https://www.the-scientist.com/features/the-pangenome-are-single-reference-genomes-dead-32458>

Why pangenomes?

The Pangenome: Are Single Reference Genomes Dead?

Researchers are abandoning the concept of a list of genes sequenced from a single individual, instead aiming for a way to describe all the genetic variation within a species.

Dec 1, 2016

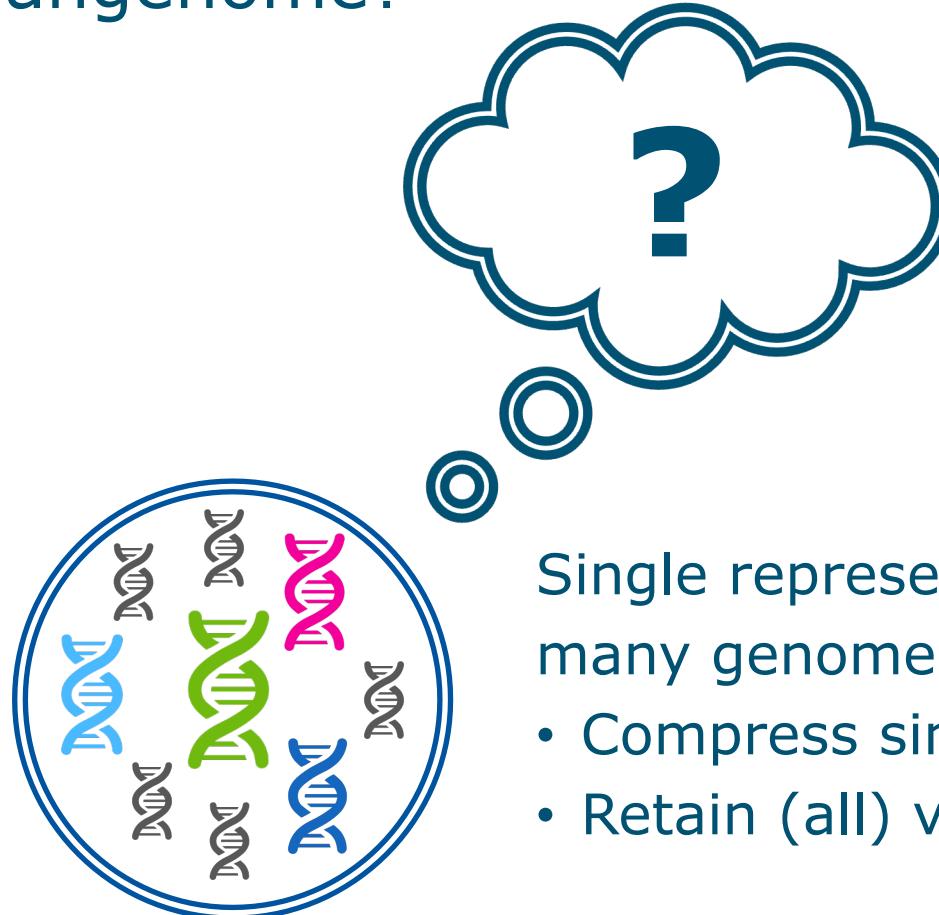
CATHERINE OFFORD



No!

- High-quality reference genomes are crucial input for pangenomes
- But only a single reference is not enough
- THE pangenome does not exist; it is a moving target
- Putting every genome available in a pangenome might not be useful

What is a pangenome?



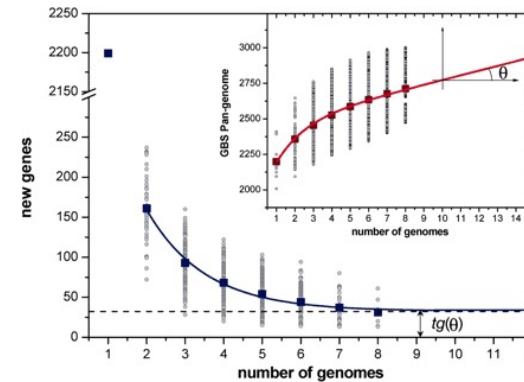
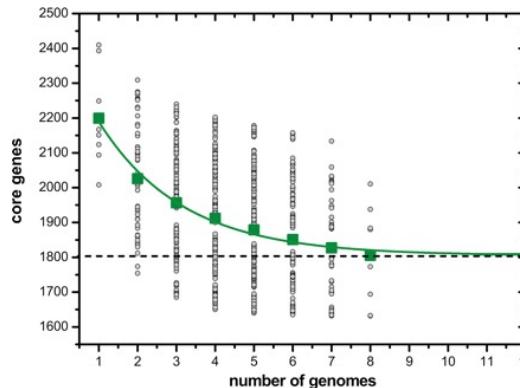
Single representation for
many genomes

- Compress similarities
- Retain (all) variations

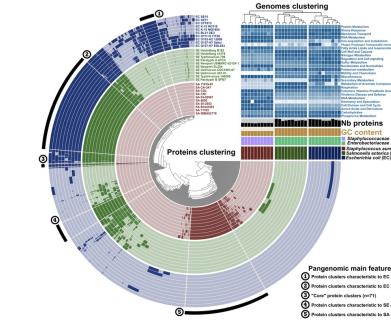
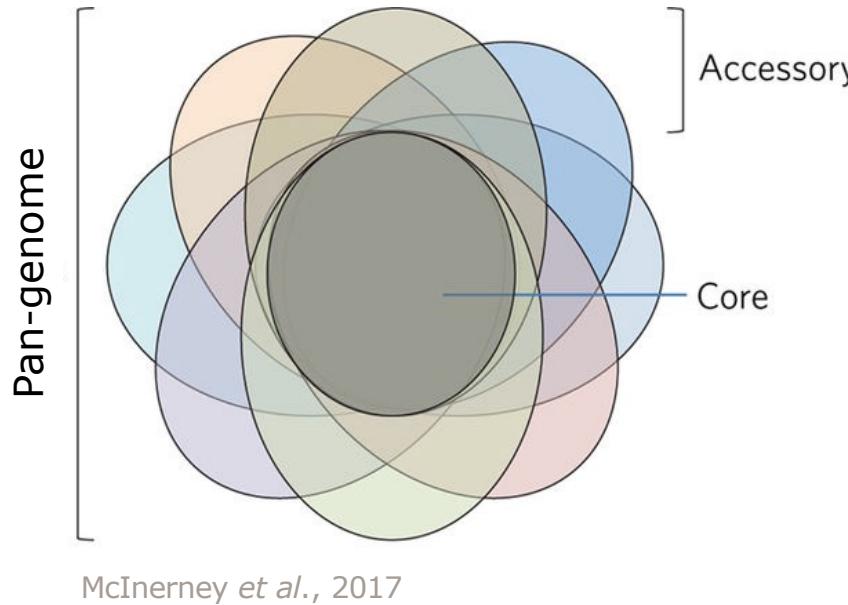
The first pangenome (Tettelin, 2005)

Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: Implications for the microbial “pan-genome”

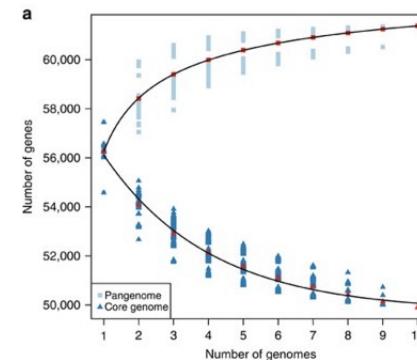
How many genomes are necessary to fully describe a bacterial species?



Gene-based pangenome

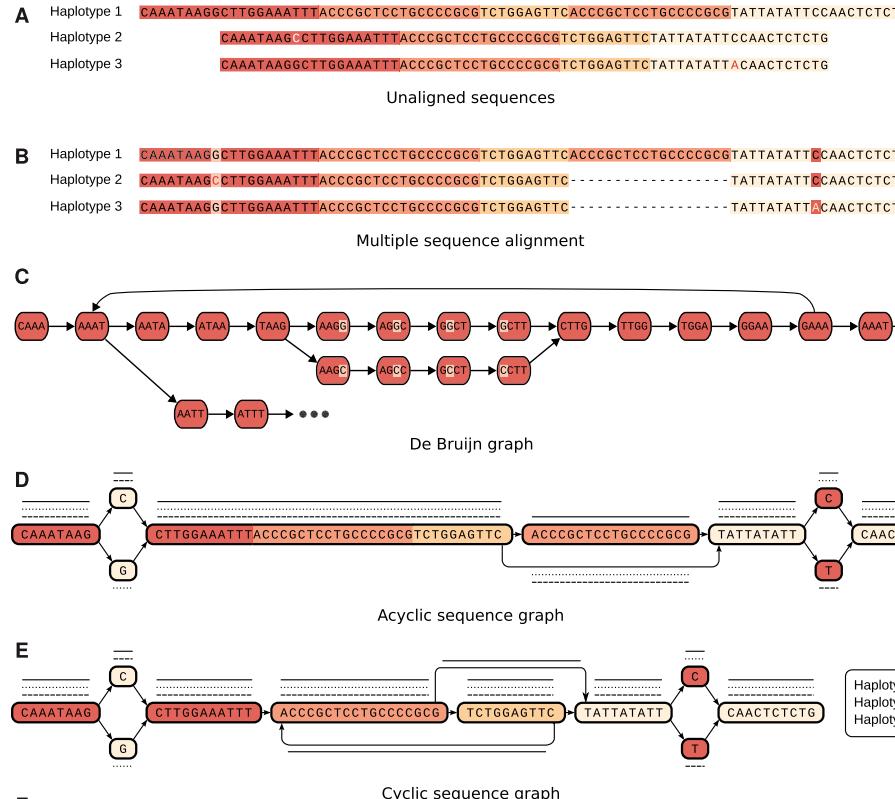


Eren et al., 2015



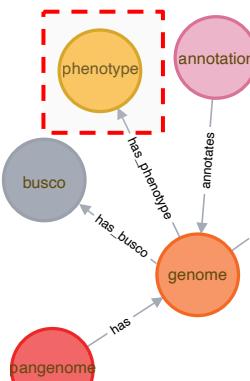
Golicz et al., 2017

Sequence-level pangenome

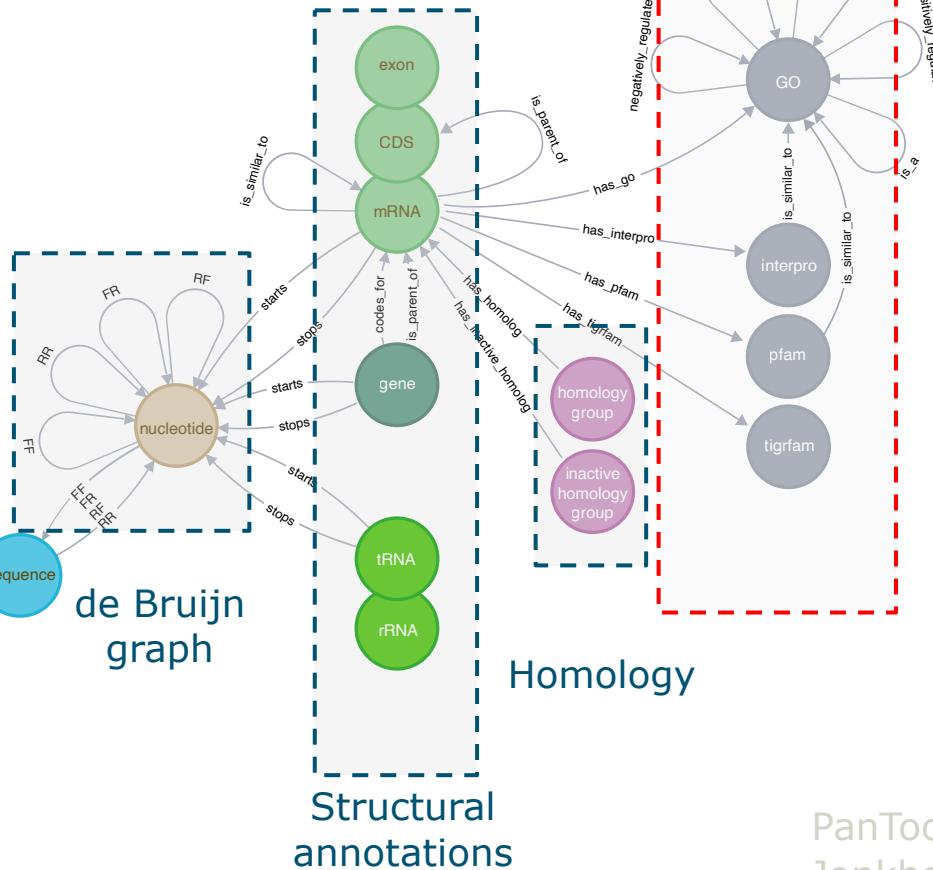


Sequences and annotations

Phenotypes



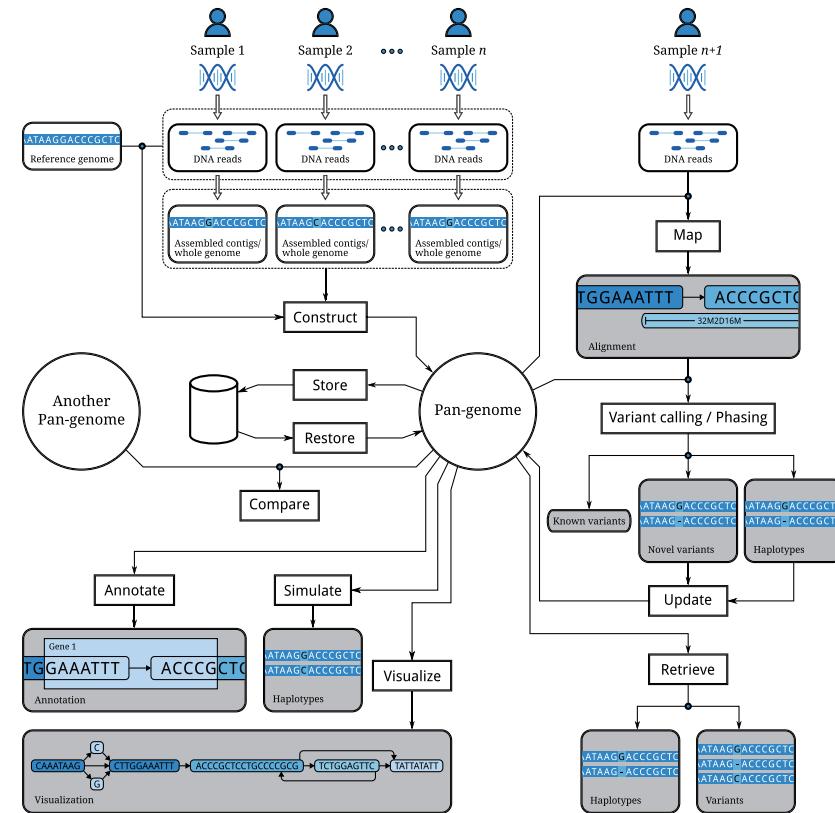
de Bruijn graph



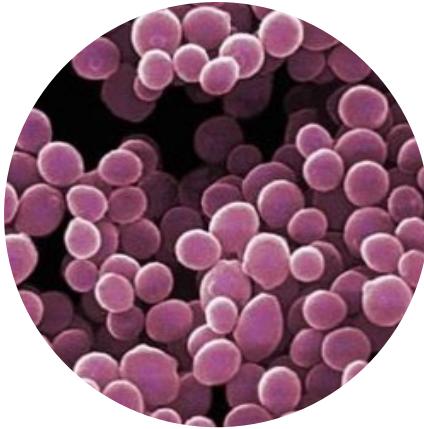
What is a pangenome?

■ Variation-aware reference

- Construct
- Annotate
- Map
- Variant calling
- Visualize
- ...



Pangenome applications



Health



Agriculture

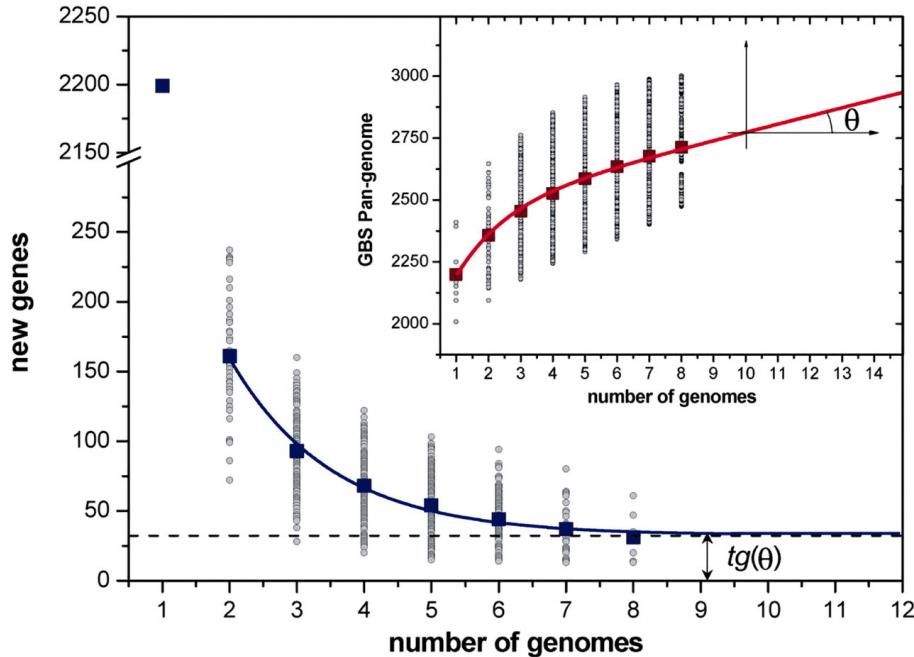


Biotechnology

The first microbial pangenome (2005)

- *Streptococcus agalactiae*
- Tettelin et al. 2005
- Gene-level pangenome

Unfortunately, the sequence of a single genome does not reflect how genetic variability drives pathogenesis within a bacterial species and also limits genome-wide screens for vaccine candidates or for antimicrobial targets.



Microbial pangenomes

PanRV: Pangenome-reverse vaccinology approach for identifications of potential vaccine candidates in microbial pangenome

[Kanwal Naz](#), [Anam Naz](#), [Shifa Tariq Ashraf](#), [Muhammad Rizwan](#), [Jamil Ahmad](#), [Jan Baumbach](#) & [Amjad Ali](#)✉

[BMC Bioinformatics](#) 20, Article number: 123 (2019) |

panX: pan-genome analysis and exploration ⓘ

[Wei Ding](#), [Franz Baumdicker](#), [Richard A Neher](#)✉

Nucleic Acids Research, Volume 46, Issue 1, 9 January 2018, Page e5,

<https://doi.org/10.1093/nar/gkx977>

Published: 25 October 2017 Article history ▾

Article | Published: 24 November 2021

Aspergillus fumigatus pan-genome analysis identifies genetic variants associated with human infection

[Amelia E. Barber](#), [Tongta Sae-Ong](#), [Kang Kang](#), [Bastian Seelbinder](#), [Jun Li](#), [Grit Walther](#), [Gianni Panagiotou](#)✉ & [Oliver Kurzai](#)✉

[Nature Microbiology](#) 6, 1526–1536 (2021) | [Cite this article](#)

Human applications

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

Erik Garrison¹, Jouni Sirén¹, Adam M Novak²✉, Glenn Hickey², Jordan M Eizenga², Eric T Dawson^{1,3,4}, William Jones¹, Shilpa Garg⁵, Charles Markello², Michael F Lin⁶, Benedict Paten² & Richard Durbin^{1,4}✉

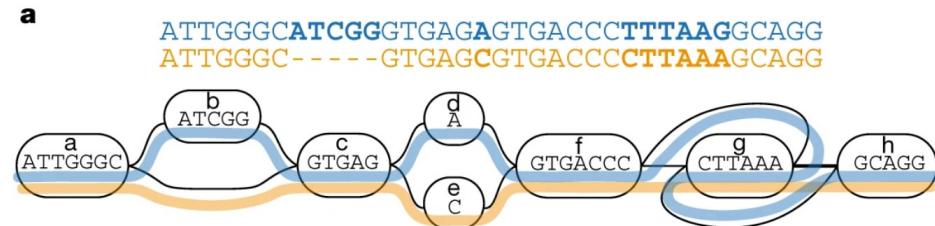
Assembly of a pan-genome from deep sequencing of 910 humans of African descent

Rachel M. Sherman ✉, Juliet Forman, [...] Steven L. Salzberg ✉

Nature Genetics **51**, 30–35 (2019) | Download Citation ↓

The first human pangenome (2023)

- 47 phased, diploid assemblies from a cohort of genetically diverse individuals



Article | [Open Access](#) | Published: 10 May 2023

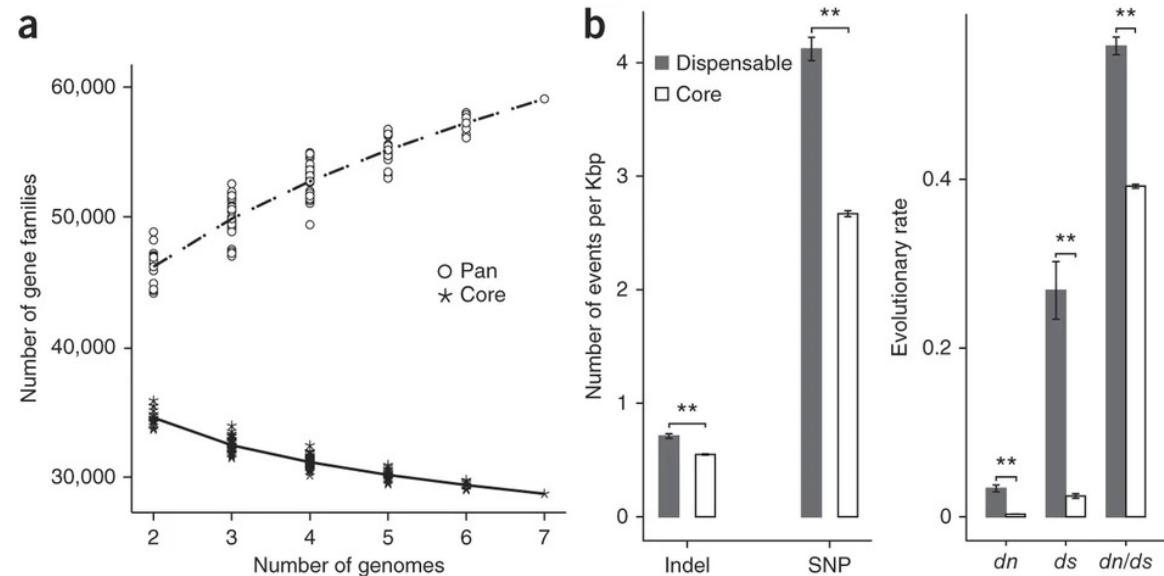
A draft human pangenome reference

[Wen-Wei Liao](#), [Mobin Asri](#), [Jana Ebler](#), [Daniel Doerr](#), [Marina Haukness](#), [Glenn Hickey](#), [Shuangjia Lu](#), [Julian K. Lucas](#), [Jean Monlong](#), [Haley J. Abel](#), [Silvia Buonaiuto](#), [Xian H. Chang](#), [Haoyu Cheng](#), [Justin Chu](#), [Vincenza Colonna](#), [Jordan M. Eizenga](#), [Xiaowen Feng](#), [Christian Fischer](#), [Robert S. Fulton](#), [Shilpa Garg](#), [Cristian Groza](#), [Andrea Guerracino](#), [William T. Harvey](#), [Simon Heumos](#), ... [Benedict Paten](#)

+ Show authors

The first plant pangenome (2014)

- Soybean
- Li et al. (2014)
- *Glycine soja*, the wild relative of cultivated soybean *Glycine max*
- 7 genomes



Since then...

Year published	Species	#accessions	Core genome
2014	<i>Brassica rapa</i>	3	87%
	Soybean	7	49%
	Rice	3	92%
	Maize	503	39%
2015	Rice	1483	?
2016	<i>Brassica oleracea</i>	10	81%
	Poplar	7	<90%
2017	<i>Brachypodium distachyon</i>	54	55%
	<i>Medicago truncatula</i>	15	33%
	Wheat	19	64%
2018	<i>Brassica napus</i>	53	62%
	Pepper	383	56%
	Rice	67	62%
	Rice	3010	54-62%
2019	Sesame	5	58%

Pangenomics – plants

Article | [Published: 13 May 2019](#)

The tomato pan-genome uncovers new genes and a rare allele regulating fruit flavor

[Lei Gao](#), [Itay Gonda](#), [Honghe Sun](#), [Qiyue Ma](#), [Kan Bao](#), [Denise M. Tieman](#), [Elizabeth A. Burzynski-Chang](#), [Tara L. Fish](#), [Kaitlin A. Stromberg](#), [Gavin L. Sacks](#), [Theodore W. Thanhhauser](#), [Majid R. Foolad](#), [Maria Jose Diez](#), [Jose Blanca](#), [Joaquin Canizares](#), [Yimin Xu](#), [Esther van der Knaap](#), [Sanwen Huang](#), [Harry J. Klee](#), [James J. Giovannoni](#)  & [Zhangjun Fei](#) 

[Nature Genetics](#) 51, 1044–1051 (2019) | [Cite this article](#)

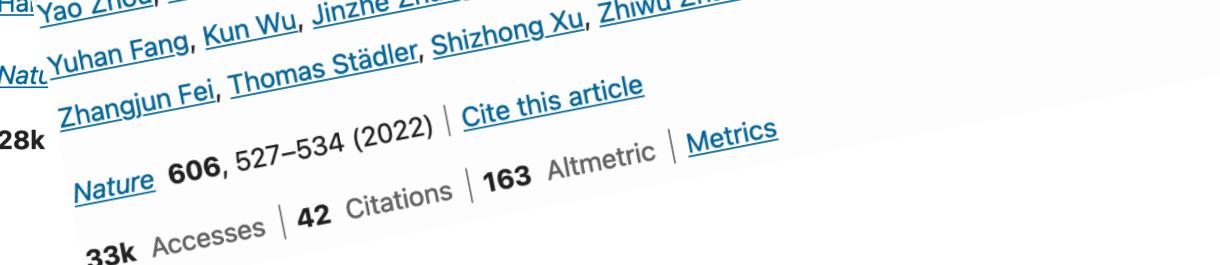
28k Accesses | **308** Citations | **555** Altmetric | [Metrics](#)

Pangenomics – plants

Article | Published: 13 May 2019

The tomato pan-genome allele  | Published: 08 June 2022

Article | Open Access

Hai Yao Zhou, Zhiyang Zhang, Zhigui Bao, Hongbo Li, Yaqing Lyu, Yanjun Zan, Yaoyao Wu, Lin Cheng, Nat Yuhan Fang, Kun Wu, Jinzhe Zhang, Hongjun Lyu, Tao Lin, Qiang Gao, Surya Saha, Lukas Mueller, Zhangjun Fei, Thomas Städler, Shizhong Xu, Zhiwu Zhang, Doug Speed & Sanwen Huang 

28k      

Nature 606, 527–534 (2022) | Cite this article

33k Accesses | 42 Citations | 163 Altmetric | Metrics

Pangenomics – plants

Article | Published: 13 May 2019

Article | Open Access | ~~Open Access~~ | Published: 06 April 2023 | Published: 08 June 2022
The tomato pan-genome captures missing heritability and breeding diversity and structural variation across wild and cultivated tomato species

Ning Li, Qiang He, Juan Wang, Baike Wang, Jiantao Zhao, Shaoyong Huang, Tao Yang, Yaping Tang, Shengbao Yang, Patiguli Aisimutuola, Ruiqiang Xu, Jiahui Hu, Chunping Jia, Kai Ma, Zhiqiang Li, Fangling Jiang, Jie Gao, Haiyan Lan, Yongfeng Zhou, Xinyan Zhang, Sanwen Huang, Zhangjun Fei, Huan Wang✉, Hongbo Li✉ & Qinghui Yu✉

Nature Genetics 55, 852–860 (2023) | Cite this article
12k Accesses | 27 Altmetric | Metrics

Animal pangenomes

[Science China Life Sciences](#)

pp 1–14 | [Cite as](#)

Building a sequence map of the pig pan-genome from multiple *de novo* assemblies and Hi-C data

Authors

The Chicken Pan-Genome Reveals Gene Content Variation and a Promoter Region Deletion in *IGF2BP1* Affecting Body Size ♂

Kejun Wang, Haifei Hu, Yadong Tian, Jingyi Li, Armin Scheben, Chaoqi Zhang, Weiwei Li

Junfeng Wu, Lan Yang, Xuewei Fan ... [Show more](#)

[Author Notes](#)

Molecular Biology and Evolution, Volume 38, Issue 11, Nov
<https://doi.org/10.1093/molbev/msab231>

Published: 30 July 2021

Authors and affiliations

Wei Fu, Yan Li, Xihong Wang, Ming Li, Duo Du, Qianzi Tang, Yudong Cai, Yiming Long, Yue Zhao,

]

Genome Res. 2022 Aug; 32(8): 1585–1601.

doi: [10.1101/gr.276550.122](https://doi.org/10.1101/gr.276550.122)

PMCID: PMC9435747

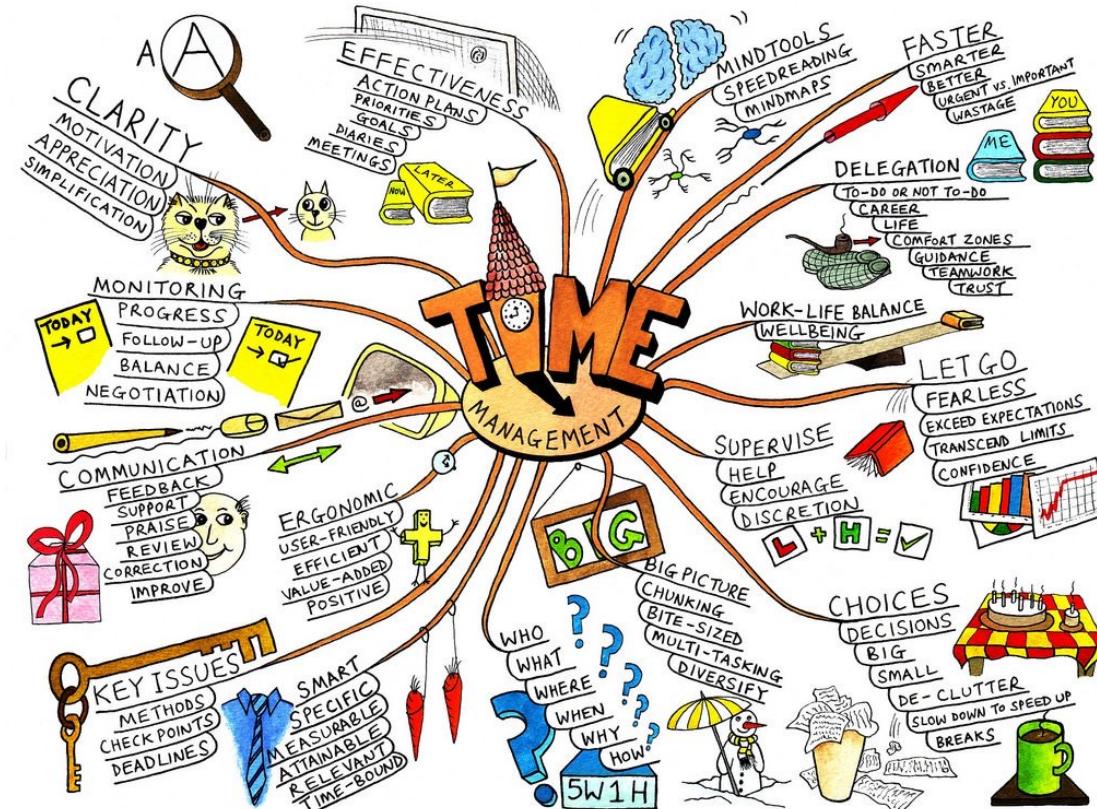
PMID: [35977842](https://pubmed.ncbi.nlm.nih.gov/35977842/)

Assembly of a pangenome for global cattle reveals missing sequences and novel structural variations, providing new insights into their diversity and evolutionary history

Yang Zhou,^{#1,6} Lv Yang,^{#1,6} Xiaotao Han,¹ Jiazheng Han,¹ Yan Hu,¹ Fan Li,¹ Han Xia,¹ Lingwei Peng,¹ Clarissa Boschiero,² Benjamin D. Rosen,² Derek M. Bickhart,³ Shujun Zhang,¹ Aizhen Guo,⁴ Curtis P. Van Tassell,² Timothy P.L. Smith,⁵ Liguo Yang,¹ and George E. Liu²

► [Author information](#) ► [Article notes](#) ► [Copyright and License information](#) ► [PMC Disclaimer](#)

Assignment: Mindmap around pangenomics



Pangenomics

A short history

- Concepts
- Algorithms
- Data structures
- Tools



Highly recommended review

Pangenome Graphs.

Eizenga JM, Novak AM, Sibbesen JA, Heumos S, Ghaffaari A, Hickey G, Chang X, Seaman JD,

Rounthwaite R, Ebler J, Rautiainen M, Garg S, Paten B, Marschall T, Sirén J, Garrison E.

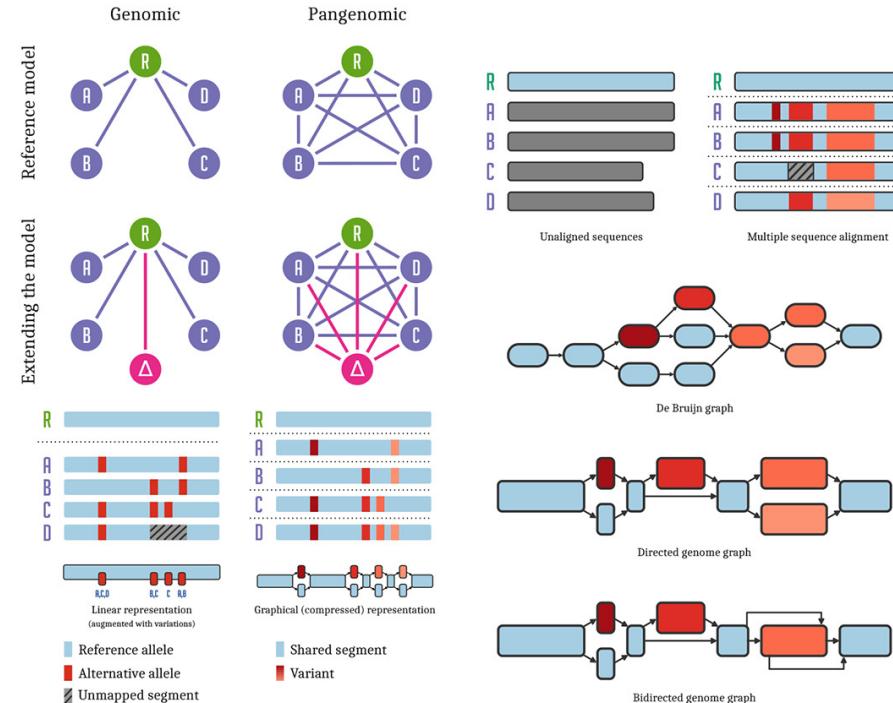
Annu Rev Genomics Hum Genet. 2020 Aug 31;21:139-162. doi: 10.1146/annurev-genom-120219-

080406. Epub 2020 May 26.

PMID: 32453966

Free PMC article.

Review.



Highly recommended review

[Home](#) > [Natural Computing](#) > Article

Computational graph pangenomics: a tutorial on data structures and their applications

Open access | Published: 04 March 2022 | **21**, 81–108 (2022)

[Download PDF](#) 

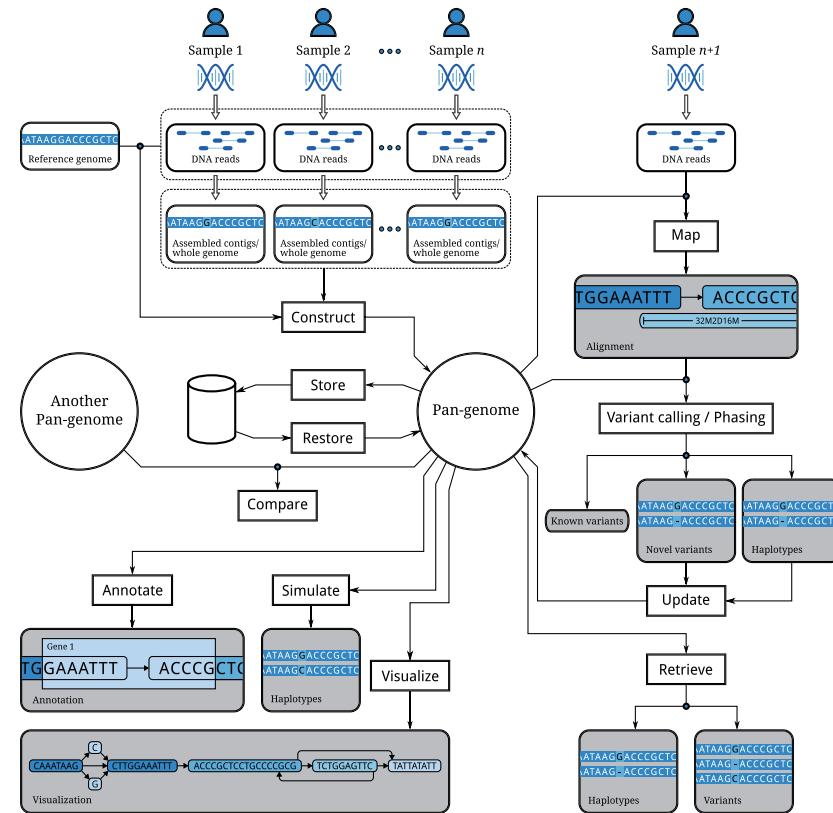
 You have full access to this [open access](#) article

Jasmijn A. Baaijens, [Paola Bonizzoni](#) , Christina Boucher, [Gianluca Della Vedova](#), [Yuri Pirola](#), Raffaella Rizzi & [Jouni Sirén](#)

What is a pangenome?

■ Variation-aware reference

- Construct
- Annotate
- Map
- Variant calling
- Visualize
- ...



Legend

○ Abstract datatype
□ Operation

→ Operation input
→ Operation output

Set
● Groups multiple inputs/outputs

Pangenome related data formats
□ Other data formats
Database
↓ Data processing

Pangenome construction

- Linear pangenome sequence
 - Iterative alignment
- Colored compacted De Bruijn Graph (cDBG)
 - Typically built from whole genomes
- Variation graph (Directed Acyclic Graph)
 - Typically built from reference + vcfs
 - Nowadays built from multiple assemblies

Iterative approach: concept

- Start with reference genome
- Map reads of another genome to reference
- Assemble the unmapped reads
- Extend the reference genome with the assembled sequences
- Continue with the next genome

Iterative approach: example

- Iterative mapping and assembly approach
- Nine morphologically diverse *B. oleracea* varieties and a wild relative—*Brassica macrocarpa*.

[Open access](#) | Published: 11 November 2016

The pangenome of an agronomically important crop plant *Brassica oleracea*

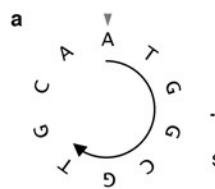
[Agnieszka A. Golicz](#), [Philipp E. Bayer](#), [Guy C. Barker](#), [Patrick P. Edger](#), [HyeRan Kim](#), [Paula A. Martinez](#),
[Chon Kit Kenneth Chan](#), [Anita Severn-Ellis](#), [W. Richard McCombie](#), [Isobel A. P. Parkin](#), [Andrew H. Paterson](#), [J. Chris Pires](#), [Andrew G. Sharpe](#), [Haibao Tang](#), [Graham R. Teakle](#), [Christopher D. Town](#),
[Jacqueline Batley](#) & [David Edwards](#) 

[Nature Communications](#) 7, Article number: 13390 (2016) | [Cite this article](#)

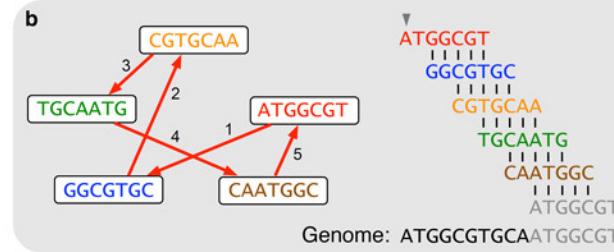
Iterative approach: pros and cons

- Little storage required
- Can be processed with conventional reference-based tools
- Easy to annotate
- Great for gene content and PAV
- Reference bias
- Processing order matters
- Filtering of novel seqs is crucial
- Does not keep allelic variants
- Ignores intragenomic variation

De Bruijn graph (known from assembly)

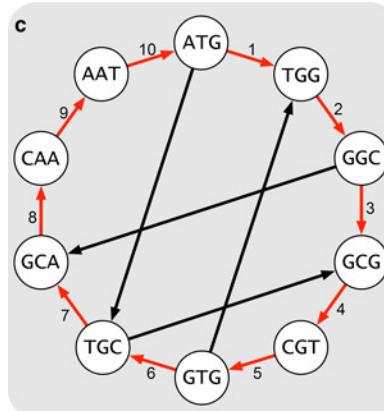


Short-read sequencing



Vertices are k -mers
Edges are pairwise alignments

Vertices are $(k-1)$ -mers
Edges are k -mers

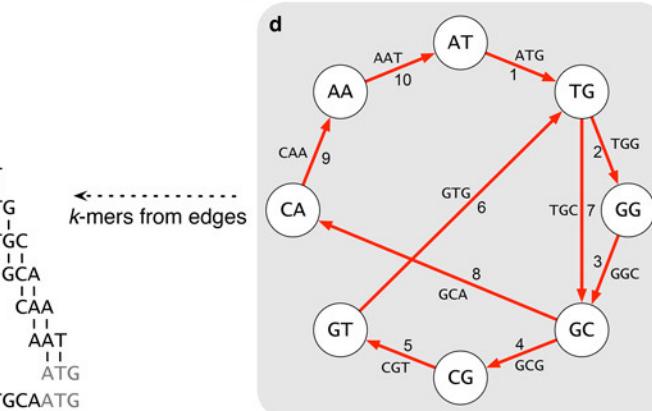


Hamiltonian cycle
Visit each vertex once
(harder to solve)

k -mers from vertices

k -mers from edges

Genome: ATGGCGTGCATGGCGT

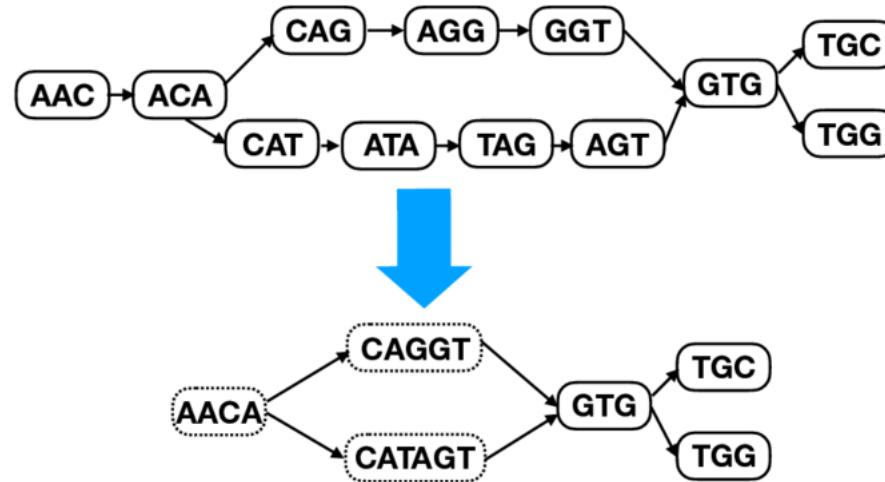


Eulerian cycle
Visit each edge once
(easier to solve)

Compeau et al. (2017)

Compacted De Bruijn graph

de Bruijn graph for the two sequences AACAGGTGC and AACATAGTGG



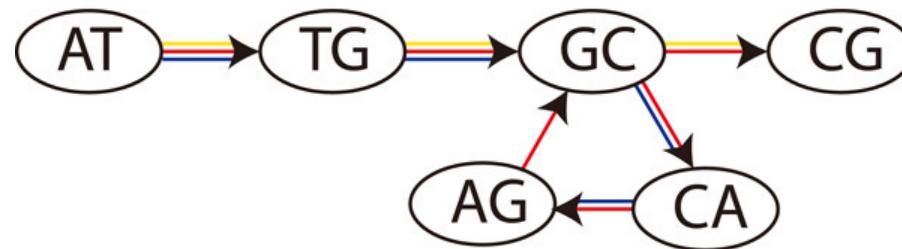
corresponding compacted de Bruijn graph

Colored (compacted) De Bruijn Graph

Reads



de Bruijn Graph



Analysis

ATGGCG ATGCAGCG ATGCAG

De Bruijn Graph: tools

> [Bioinformatics](#). 2016 Jun 15;32(12):i201-i208. doi: 10.1093/bioinformatics/btw279.

Compacting de Bruijn graphs from sequencing data quickly and in low memory

bcalm 2

Rayan Chikhi ¹, Antoine Limasset ², Paul Medvedev ³

Affiliations + expand

PMID: 27307618 PMCID: [PMC4908363](#) DOI: [10.1093/bioinformatics/btw279](#)

> [Bioinformatics](#). 2016 Sep 1;32(17):i487-i493. doi: 10.1093/bioinformatics/btw455.

PanTools: representation, storage and exploration of pan-genomic data

Siavash Sheikhzadeh ¹, M Eric Schranz ², Mehmet Akdel ¹, Dick de Ridder ¹, Sandra Smit ¹

Affiliations + expand

PMID: 27587666 DOI: [10.1093/bioinformatics/btw455](#)

Bifrost: highly parallel construction and indexing of colored and compacted de Bruijn graphs

[Guillaume Holley](#)  & [Páll Melsted](#)

Genome Biology 21, Article number: 249 (2020) | [Cite this article](#)

Method | [Open access](#) | Published: 08 September 2022

Scalable, ultra-fast, and low-memory construction of compacted de Bruijn graphs with Cuttlefish 2

[Jamshed Khan](#), [Marek Kokot](#) , [Sebastian Deorowicz](#) & [Rob Patro](#) 

Genome Biology 23, Article number: 190 (2022) | [Cite this article](#)

Extremely fast construction and querying of compacted and colored de Bruijn graphs with GGCAT

[Andrea Cracco](#)¹ and [Alexandru I. Tomescu](#)²

 Author Affiliations

- Corresponding authors: alexandru.tomescu@helsinki.fi, andrea.cracco@univr.it

De Bruijn Graph: pros and cons

- Strong compression
- Many applications
- Not alignment-based
- Lose sequence information

Variation graphs



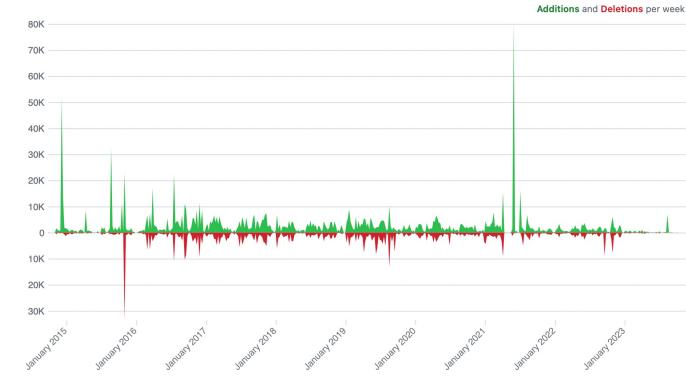
Produces a directed acyclic graph (DAG), ordered along the reference genome, with bubbles at the sites of variation

Variation graphs provide a succinct encoding of the sequences of many genomes. A variation graph (in particular as implemented in vg) is composed of:

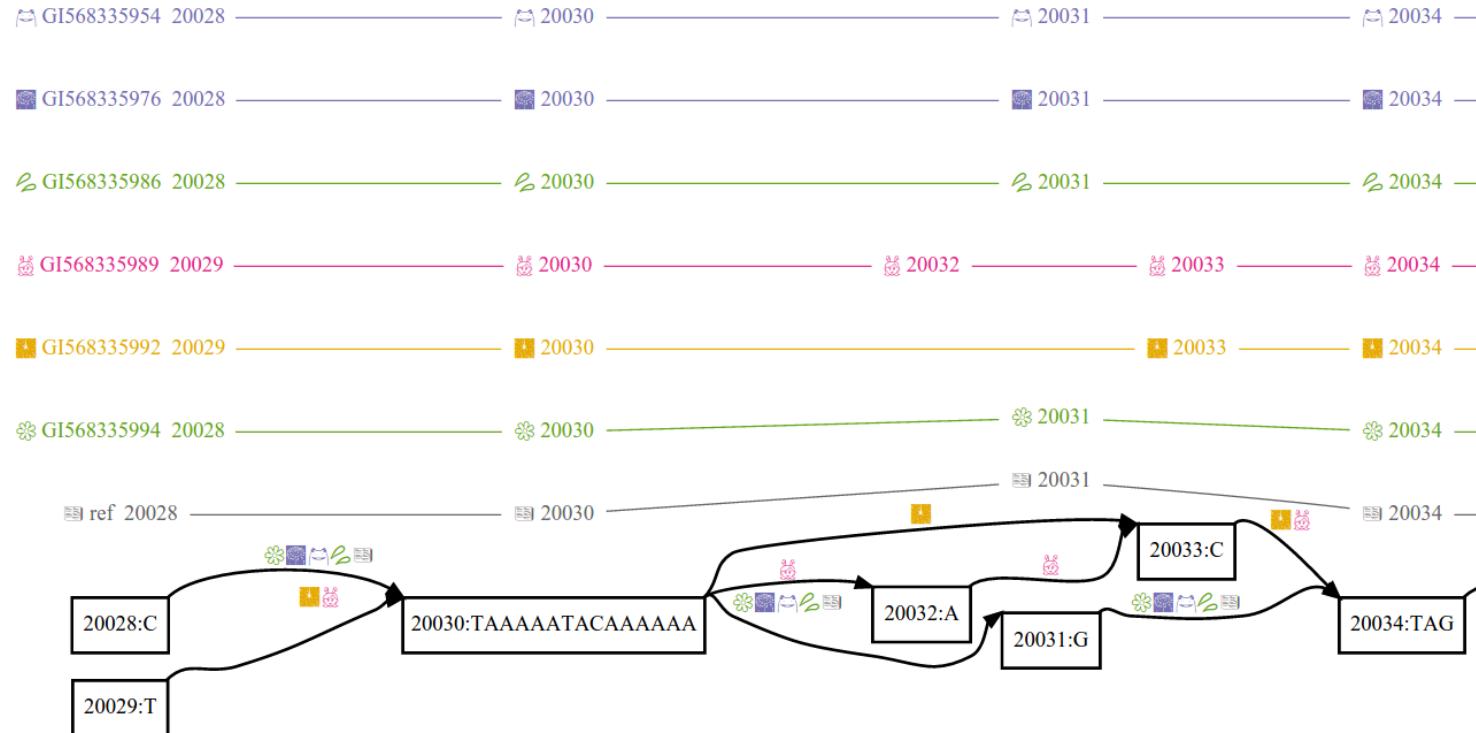
- *nodes*, which are labeled by sequences and ids
- *edges*, which connect two nodes via either of their respective ends
- *paths*, describe genomes, sequence alignments, and annotations (such as gene models and transcripts) as walks through nodes connected by edges

This model is similar to sequence graphs that have been used in assembly and multiple sequence alignment.

Code frequency over the history of vgteam/vg



Variation graphs



Variation graphs: algorithm description

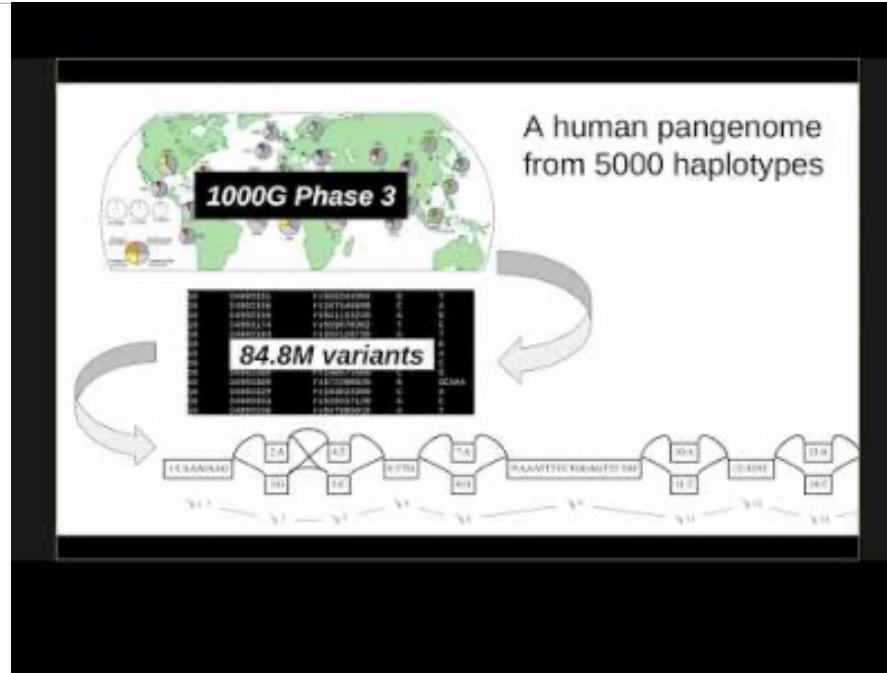
Published: 01 October 2018

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) | [Cite this article](#)

Building and understanding pangenome variation graphs - Erik Garrison



https://www.youtube.com/watch?v=S8oSL4_Bqqw

Minigraph

- Motivation: there are no satisfactory solutions to the construction of reference pangenome graphs.

Method | [Open access](#) | Published: 16 October 2020

The design and construction of reference pangenome graphs with minigraph

[Heng Li](#) , [Xiaowen Feng](#) & [Chong Chu](#)

Genome Biology **21**, Article number: 265 (2020) | [Cite this article](#)

- Our implementation, minigraph (<https://github.com/lh3/minigraph>), can construct a pangenome graph from twenty human assemblies in 3 h.

PanGenome Graph Builder (PGGB)

Building pangenome graphs

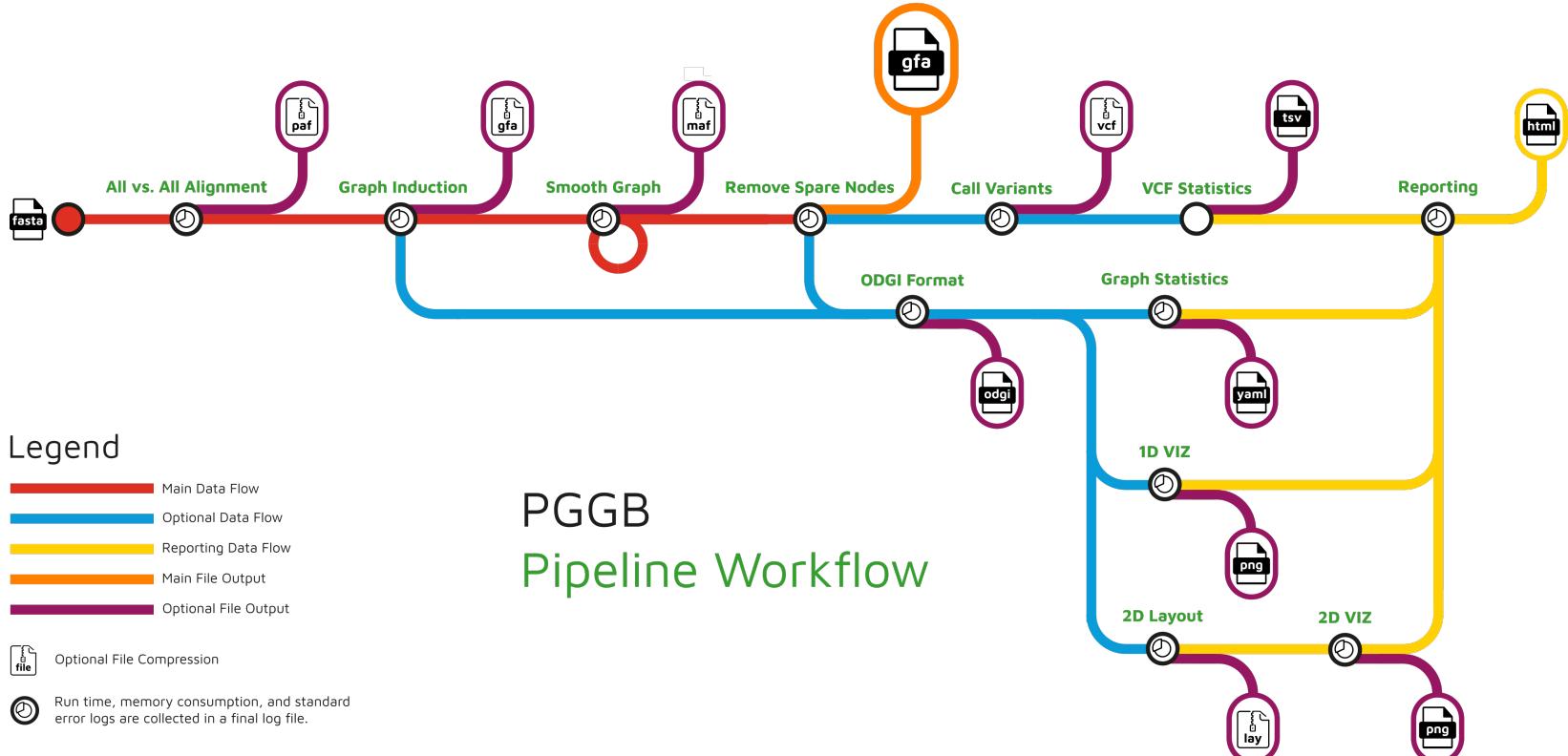
BioRxiv, April 2023

 Erik Garrison,  Andrea Guerracino,  Simon Heumos,  Flavia Villani,  Zhigui Bao,  Lorenzo Tattini,
 Jörg Hagmann,  Sebastian Vorbrugg,  Santiago Marco-Sola,  Christian Kubica,  David G. Ashbrook,
 Kaisa Thorell,  Rachel L. Rusholme-Pilcher,  Gianni Liti, Emilio Rudbeck,  Sven Nahnsen,
 Zuyu Yang,  Mwaniki N. Moses,  Franklin L. Nobrega,  Yi Wu,  Hao Chen,  Joep de Ligt,
 Peter H. Sudmant,  Nicole Soranzo,  Vincenza Colonna,  Robert W. Williams,  Pjotr Prins

doi: <https://doi.org/10.1101/2023.04.05.535718>

- Motivation: existing methods for constructing them are biased due to reference-guided approaches
- PGGB uses all-to-all whole-genome alignments and learned graph embeddings to build and iteratively refine a model in which we can identify variation, measure conservation, detect recombination events, and infer phylogenetic relationships

PanGenome Graph Builder (PGGB)



<https://pggb.readthedocs.io>

Minigraph-cactus

Article | Published: 10 May 2023

Pangenome graph construction from genome alignments with Minigraph-Cactus

Glenn Hickey , Jean Monlong, Jana Ebler, Adam M. Novak, Jordan M. Eizenga, Yan Gao, Human

Pangenome Reference Consortium, Tobias Marschall, Heng Li & Benedict Paten 

[Nature Biotechnology](#) (2023) | [Cite this article](#)

Here we present the Minigraph-Cactus pangenome pipeline, which creates pangomes directly from whole-genome alignments, and demonstrate its ability to scale to 90 human haplotypes from the Human Pangenome Reference Consortium.

Research Briefing | Published: 22 May 2023

Combining reference genomes into a pangenome graph improves accuracy and reduces bias

[Nature Biotechnology](#) (2023) | [Cite this article](#)

1583 Accesses | 19 Altmetric | [Metrics](#)

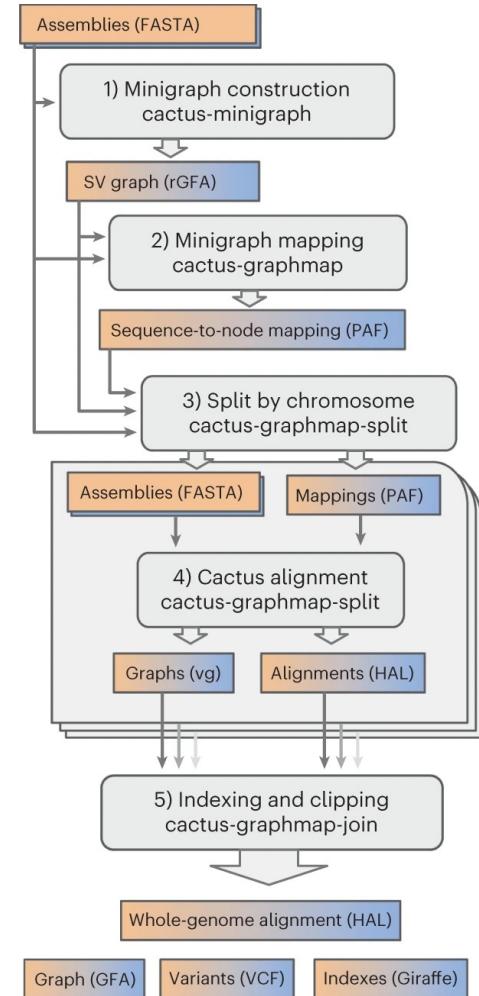
Minigraph-Cactus, a method to efficiently combine multiple reference genome assemblies into a pangenome reference graph, can be used to improve accuracy of read mapping and variant calling compared with a single reference in downstream applications.

Minigraph-cactus

This method proceeds by iteratively adding structural variants from each sample, then using these variants as anchors for computing a base-level realignment (Fig. 1).

The resulting graph consistently represents variation at different scales, from single-nucleotide polymorphisms to complex, nested intrachromosomal rearrangements.

Crucially, existing pangenomics software can be used with this method, as well as for read mapping and variant calling.

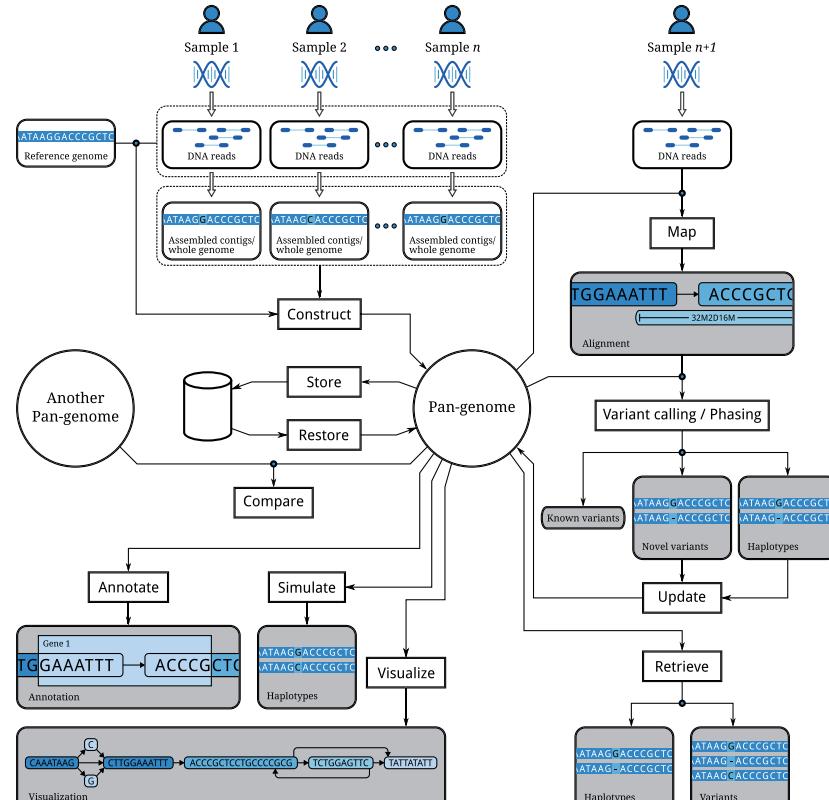


Not discussed

- Annotation
- Indexing

What is a pangenome?

- Variation-aware reference
 - Construct
 - Annotate
 - Map
 - Variant calling
 - Visualize
 - ...



Read mapping: tools

Table 1 Overview of some state-of-the-art sequence-to-graph aligners, categorized according to the type of graph they handle.

acyclic variation graph	general variation graph	assembly de Bruijn graph	de Bruijn graph as a reference
GenomeMapper Graph Genome Aligner HISAT2 PaSGAL V-MAP Vargas	vg GraphAligner Minigraph Giraffe	BGREAT BrownieAligner GraphAligner SPAligner	deBGA PuffAligner Beller & Ohlebusch Nexus (this paper)

From: Depuydt et al. (2023).

<https://assets.researchsquare.com/files/rs-2583159/v1/61d0b7c9b76925e42ee15cdd.pdf?c=1677080263>

Read mapping: tools

[Genome Biol.](#) 2020; 21: 253.

Published online 2020 Sep 24. doi: [10.1186/s13059-020-02157-2](https://doi.org/10.1186/s13059-020-02157-2)

GraphAligner: rapid and versatile sequence-to-graph alignment

[Mikko Rautiainen](#)^{✉^{1,2,3}} and [Tobias Marschall](#)^{✉⁴}

Mapping short reads with Giraffe

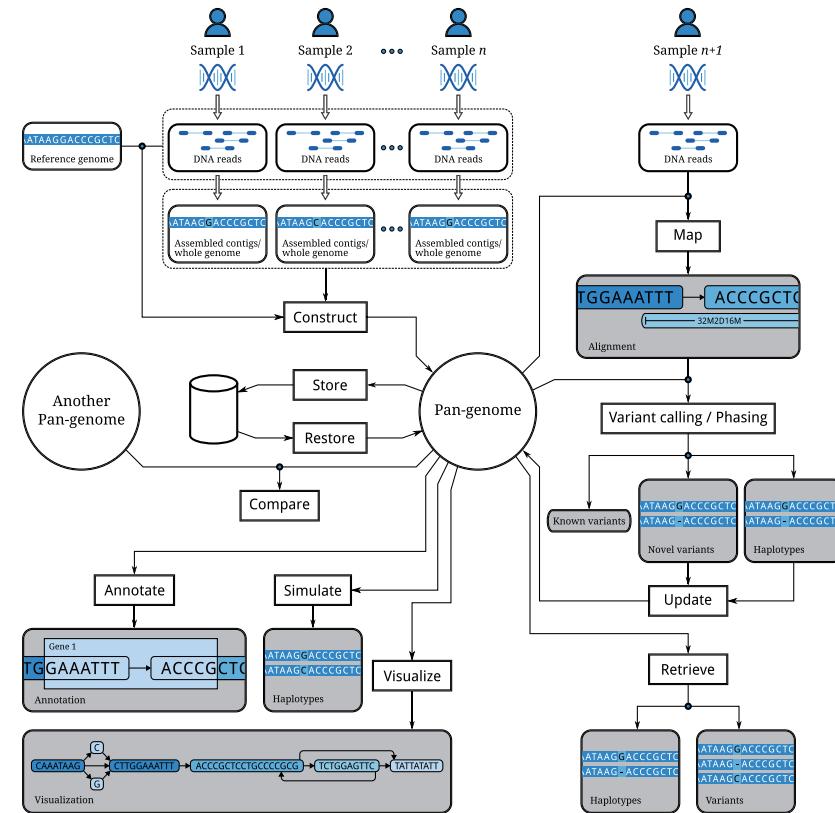
Jouni Siren edited this page on May 28 · 10 revisions

This tutorial will explain how to use `vg giraffe` to map short reads to a pangenome graph. See also [Giraffe best practices](#).

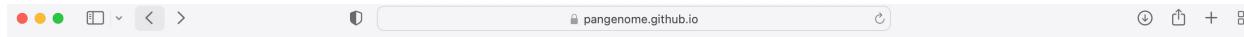
What is a pangenome?

■ Variation-aware reference

- Construct
- Annotate
- Map
- Variant calling
- Visualize
- ...



<https://pangenome.github.io>

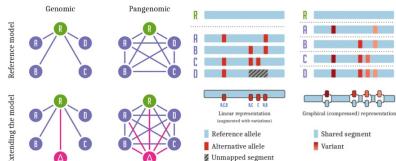


Practical Graphical Pangenomics

tools and workflows based on genome variation graphs

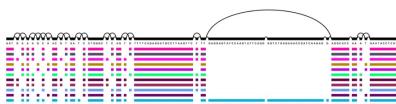
Pangenomic methods

Standard approaches to genome inference and analysis relate sequences to a single linear reference genome. This is efficient but has a fundamental problem: Differences from this reference are hard to observe and describe in a coherent way. Variation and sequence are separated.



Pangenomic methods allow us to relate all genomes or sequences in our analysis directly to each other. Sequence and variation are combined into a coherent data structure. This practice is still new, and research into ways to design, implement, and apply this model is ongoing. However, there is a growing consensus around best practices. Many methods work on an augmented sequence graph model and use a handful of common data formats for input and output.

The *variation graph* data model describes the all-to-all alignment of many sequences (genomes or genes for instance) as walks through a graph whose nodes are labeled with DNA sequences:



vg

The variation graph toolkit **vg** provides computational methods for creating and manipulating of genome variation graphs. It's pangenome representation of a set of genomes overcomes reference bias and improves read mapping. This is highlighted in the [Nature Biotechnology publication](#). Users can receive support on [vg's BioStars page](#).

PanGenome Graph Evaluator (pgge)

This pangenome graph evaluation pipeline measures the reconstruction accuracy of a pangenome graph (in the variation graph model). Its goal is to give guidance in finding the best pangenome graph construction tool for a given input data and task.

xg

The succinct graph index **xg** presents a static index of nodes, edges and paths of a variation graph. **xg** can be used to annotate graph nodes with their reference path relative positions. It was a key component of early development in **vg**, and was used to scale sparse read mapping to large genomes. It implements the [libhandleread API](#).

GWBT

[GWBT](#) - GWAS GWV BT is a command-line tool for GWAS analysis.

PanGenome Graph Builder (pggb)

This pangenome graph construction pipeline renders a collection of sequences into a pangenome graph (in the variation graph model). Its goal is to build a graph that is locally directed and acyclic while preserving large-scale variation. Maintaining local linearity is important for the interpretation, visualization, and reuse of pangenome variation graphs. A Nextflow version of the pipeline is also available [nf-core/pangenome](#).

Pangenome Graph Variation Format (PGVF)

PGVF is a hard fork of the GFAv1 format that allows the description of graph-to-graph alignments. It represents a collection of aligned graphs as a network of walks through an underlying merged sequence graph. While pangenome graphs let us represent differences between genomes, we have no mechanism to represent differences between pangenome graphs, or to combine multiple pangenome graphs into one structure without losing information. This motivates the development of a new biological data format.

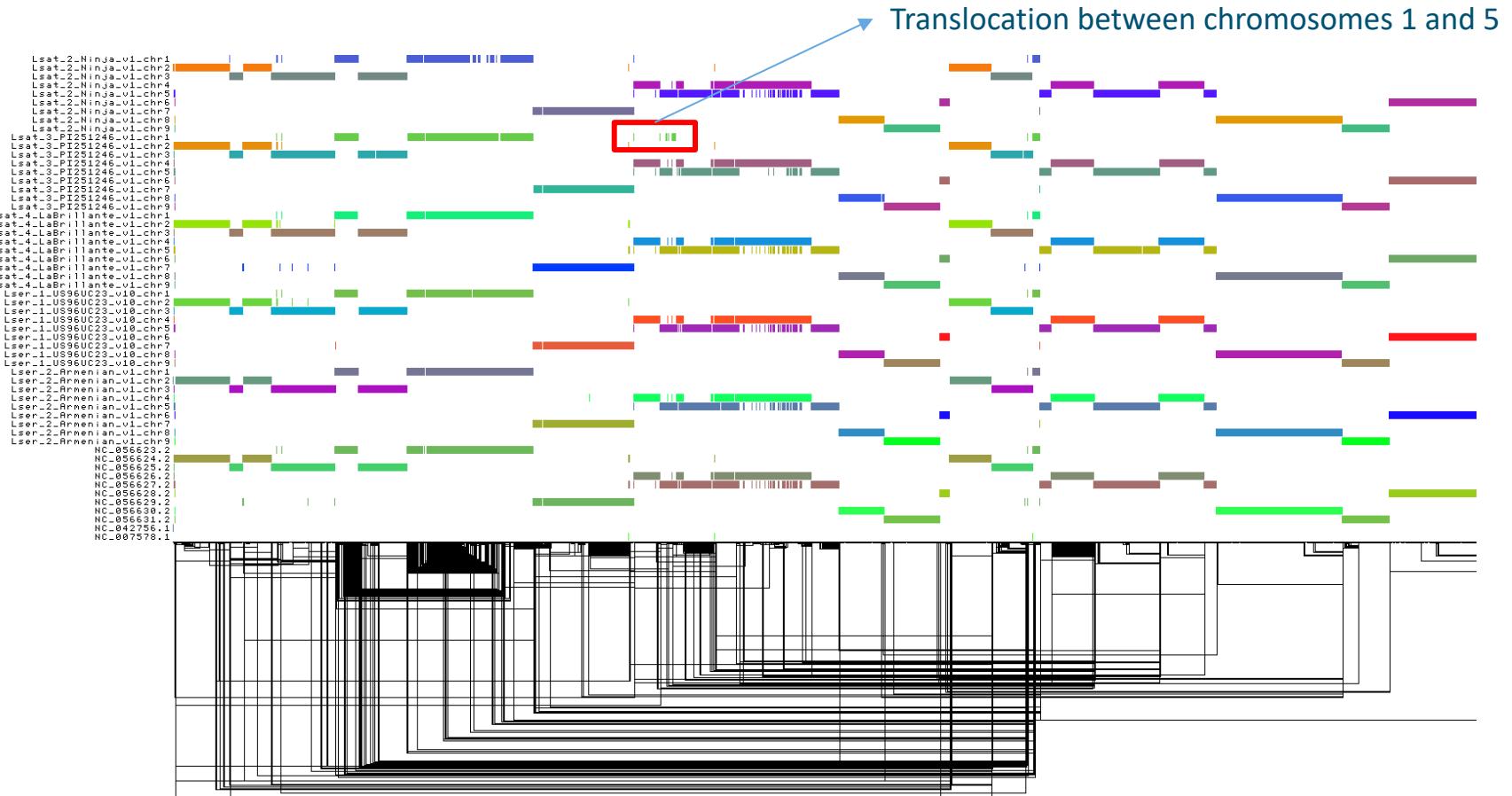
odgi

odgi, the Optimized Dynamic (genome) Graph Interface, links a thrifty dynamic in-memory variation graph data model to a set of algorithms designed for scalable sorting, pruning, transformation, and visualization of very large genome graphs. **odgi** includes [python bindings](#) that can be used to [directly interface with its data model](#). The [odgi manual](#) provides detailed information about its features and subcommands, including examples.

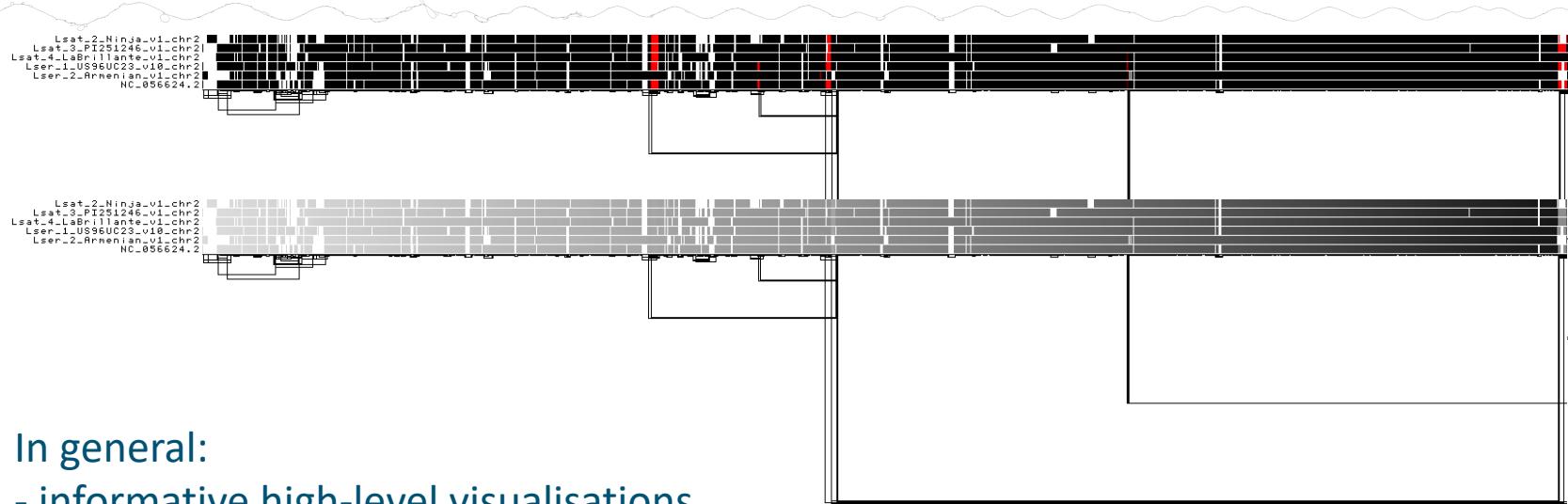
spodgi

[spodgi](#) - SPace-Oriented Dynamic Graph Interface

PGGB pangenome



PGGB pangenome chr2



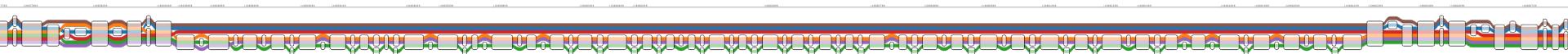
In general:

- informative high-level visualisations
- takes a long time to compute (± 5 days for all chromosomes)
- alignment extremely dependent on parameters

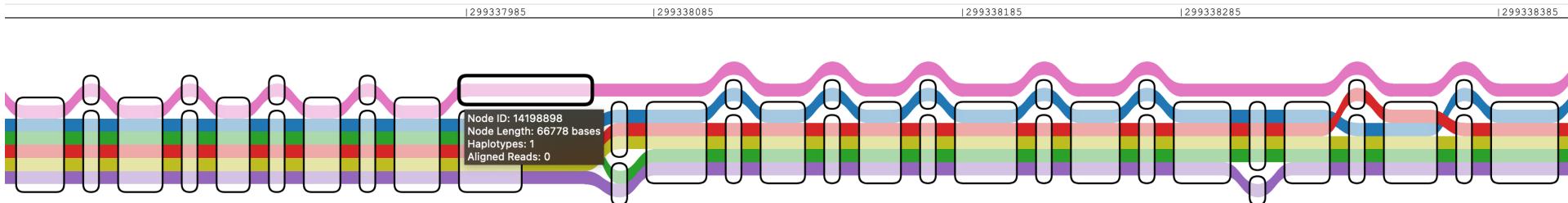
MC pangenome

- “Minigraph-Cactus” pipeline (using “sequenceTubeMap” visualisation)

- 4kb transposon insertion (*LsKN1*) clearly visible
(insertion present in Salinas, Ninja, LaBrillante)



- g2349.t1 is non-reference because of ±67kb part in Salinas which doesn't have it
(transcript present in Ninja, PI251246, LaBrillante, US96UC23, Armenian 999)



Interactive visualization: PanX



Interactive visualization: PanVA



Conclusions on pangenomics

- Fast-developing field
 - There are still many challenges
- Graph representations are dominant
- Dominant tools: vg, pggb, minigraph-cactus
- Downstream applications are crucial to do comparative genomics
- Visualization is an important aspect of pangenome analysis

Program today

- 9.00 Lecture
- 11.30 Break
- 11.00 Computer exercises (rosalind or project)
- 12.30 Lunch
- 13.30 Guest lecture by Erwin Datema from Keygene
- 14.15 Lecture about PanTools
- 15.00 Break
- 15.30 Computer exercises (rosalind or project or PanTools)
- 17.30 Closing

Exercises

- Rosalind
 - Graphs (pangenomics): Complete questions 35 and higher
 - Supporting slides on graph algorithms available
- Project
 - Minimap & miniasm
- PanTools (afternoon)
 - Tutorial on <https://pantools.readthedocs.io>