# Pangenome Graphs

IBMI PhD Talks 2023
22 February 2023

<u>Simon Heumos</u>

# *De novo* assembly and a pangenomic model
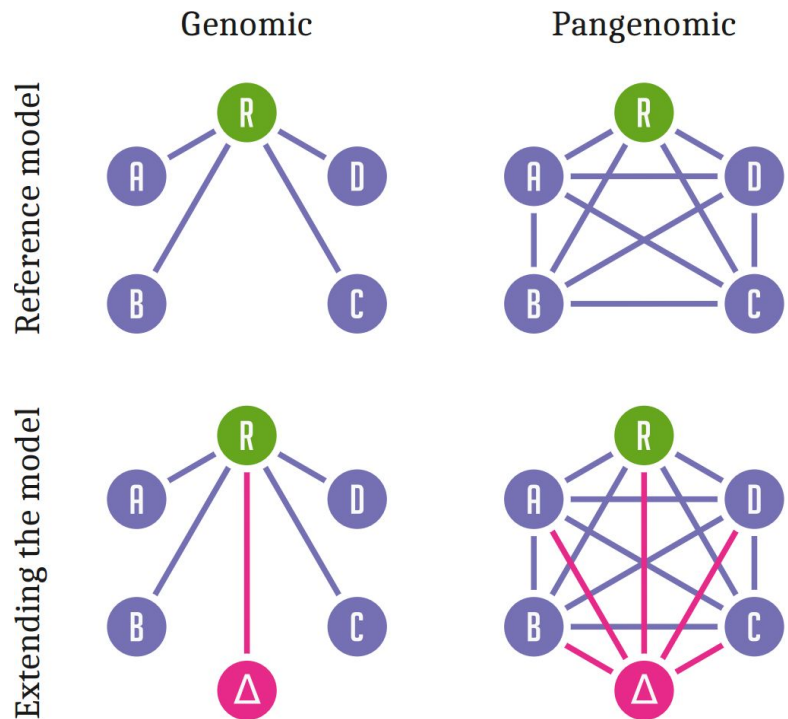
Genomic

Reference model



Extending the model

Thanks to advances in sequencing technology, new **telomere-to-telomere** genome assemblies are produced at a high rate.

Δ: new genome; R: reference genome.
Figure from Eizenga et al., 2020.
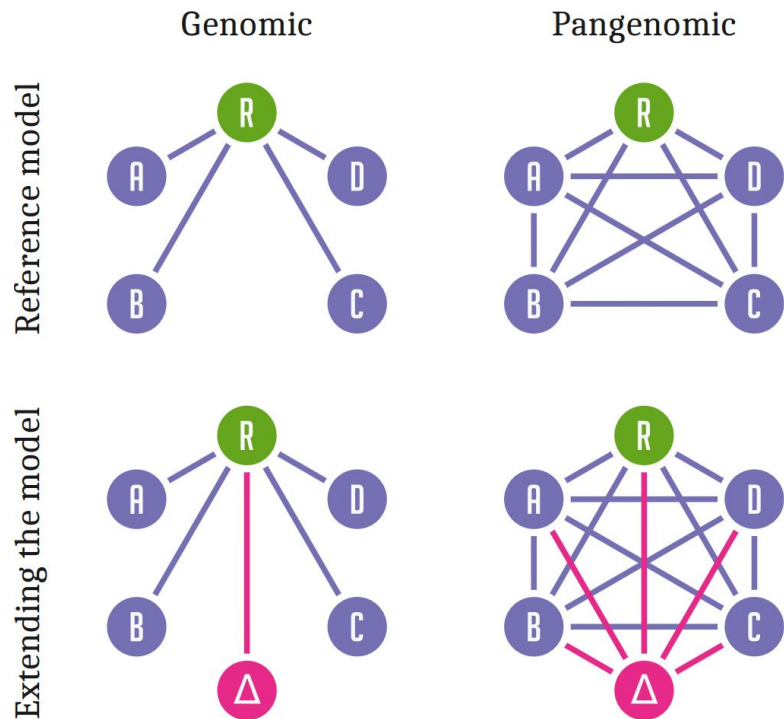
# *De novo* assembly and a pangenomic model



Genomic  Pangenomic

Reference model

Extending the model

Thanks to advances in sequencing technology, new **telomere-to-telomere** genome assemblies are produced at a high rate.

**Pangenomes** can **model** the full set of genomic elements in a given species or clade, reducing the **reference-bias**.

Δ: new genome; R: reference genome.
Figure from Eizenga et al., 2020.

# A pangenome encoded as a graph



Δ: new genome; R: reference genome.
Figure from Eizenga et al., 2020.

Figure from Eizenga et al., 2020.

# Pangenome graphs - representation

Pangenomes can take many forms, including **graph-based** data structures.

**Pangenome graphs** compress redundant sequences into a smaller data structure that is still representative of the full set.
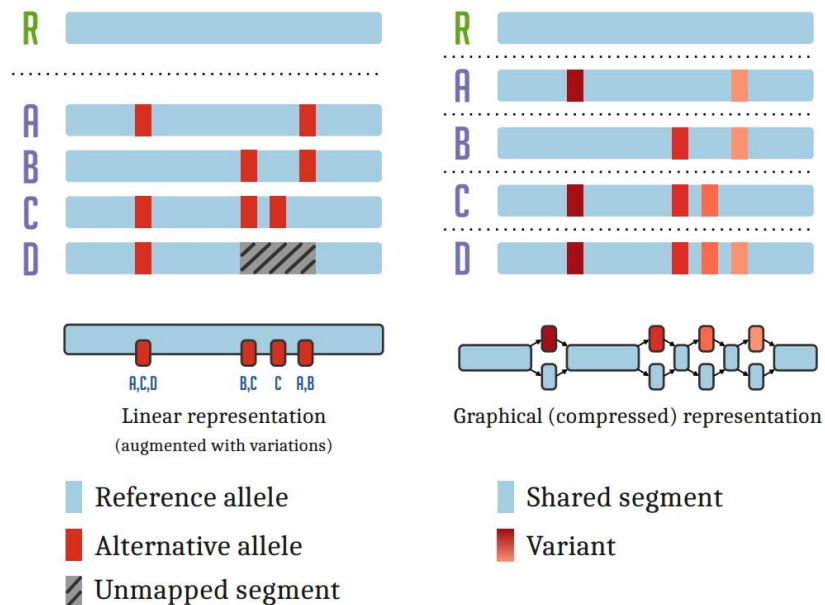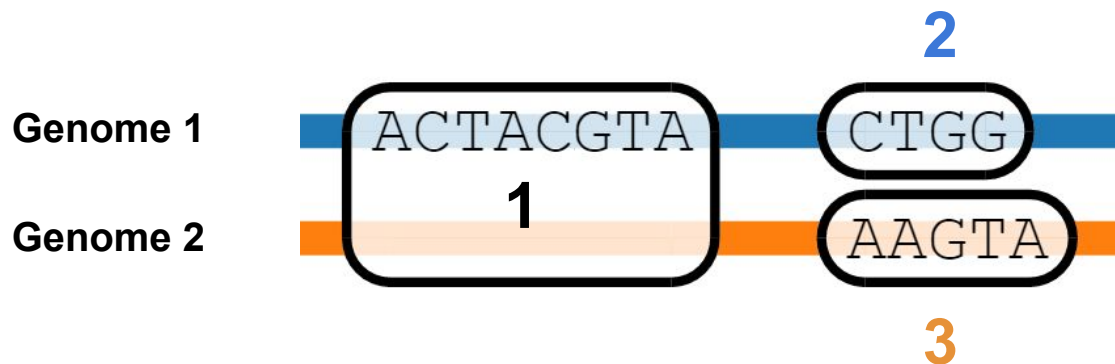


Figure from Eizenga et al., 2020.

# Variation graphs



Genome 1: **ACTACGTA**<span style="color:blue">**CTGG**</span>  Path: **1** <span style="color:blue">**2**</span>

Genome 2: **ACTACGTA**<span style="color:orange">**AAGTA**</span>  Path: **1** <span style="color:orange">**3**</span>

Linear sequences are **paths** through nodes.

**2**

**Genome 1**  ACTACGTA  CTGG

**1**

**Genome 2**  AAGTA

**3**

Graph topology is not directly shown.

The nodes represent DNA sequences.

Sketch made using SequenceTubeMap.

**Paths** can be contigs, haplotypes, reads, or whole chromosomes.

# Towards a 1D visualization



Genome 1: **ACTACGTA** **CTGG**   Path: **1 2**

Genome 2: **ACTACGTA** **AAGTA**   Path: **1 3**

Genome 1

Genome 2

2

1
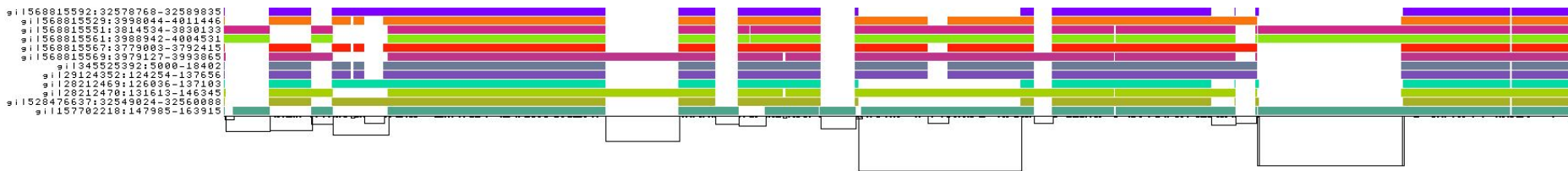
3

**ACTACGTACTGGAAGTA**

Concatenate nucleotides to a pangenome sequence.

Presence - absence matrix encodes actual genomic sequence.
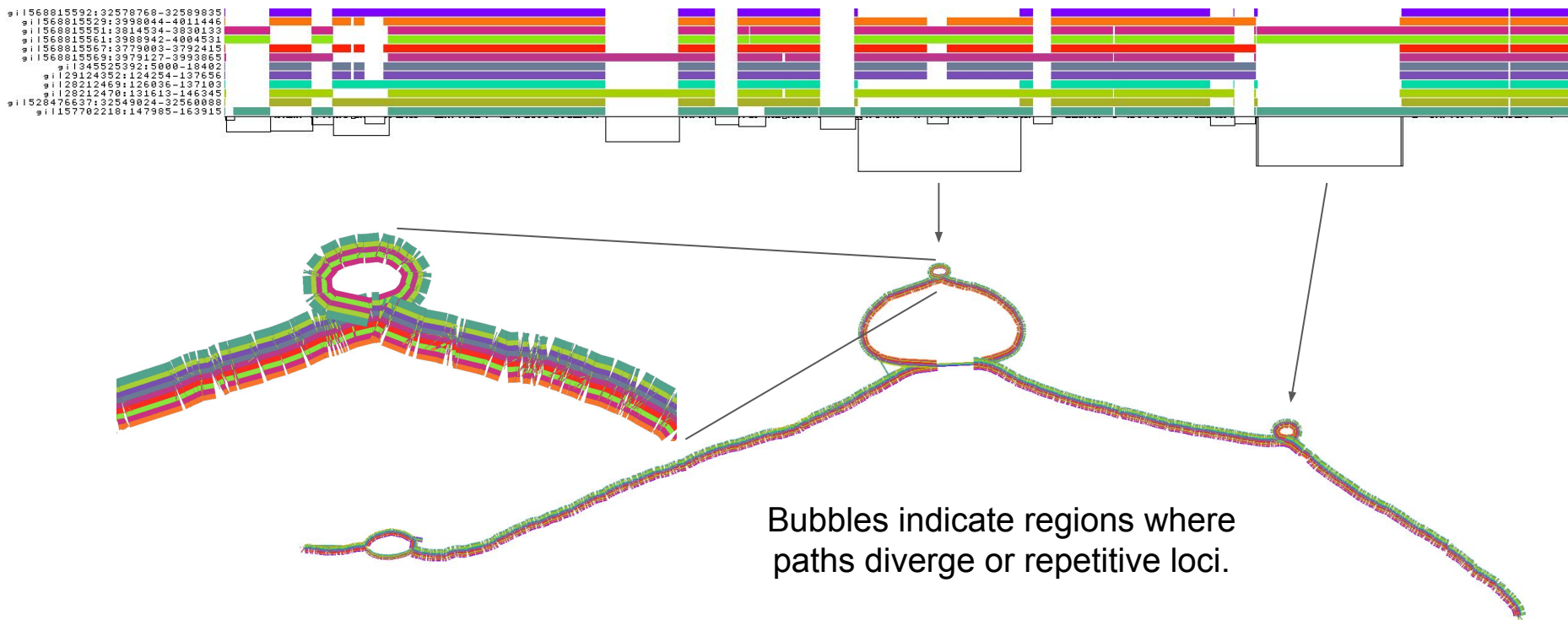
# 1D Graph visualization explained

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



- Graph nodes are arranged from left to right forming the pangenome sequence

- Colored bars are the paths versus the pangenome sequences in a binary matrix

- Path names are left

- The black lines under the paths are the links representing the graph topology

# 2D Graph visualization explained

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



Bubbles indicate regions where paths diverge or repetitive loci.

# Building Pangenome Graphs

Solving the whole genome alignment problem in 3 steps.

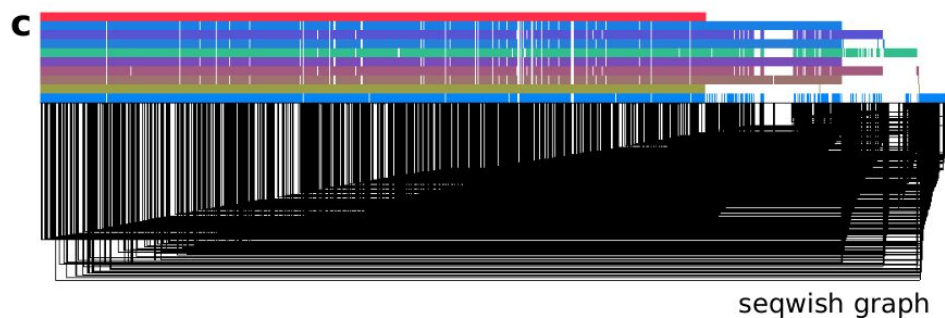a) all-to-all alignment
b) graph induction
c-f) normalization

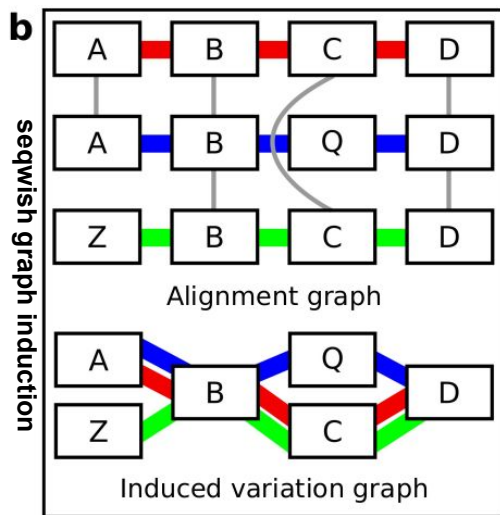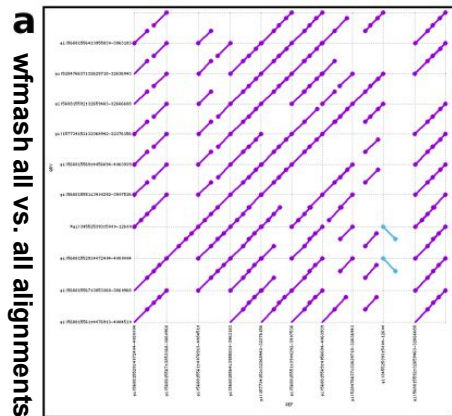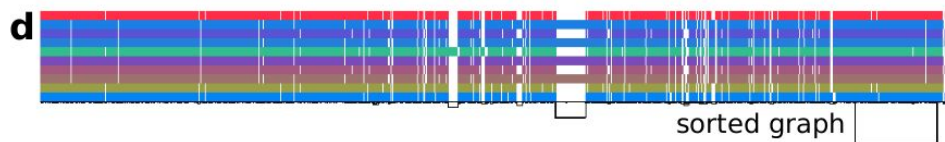implemented in the PanGenome Graph Builder ([PGGB](PGGB))



Erik Garrison    Andrea Guarracino

wfmash all vs. all alignments

seqwish graph induction

a

b

A — B — C — D
A — B — Q — D
Z — B — C — D

Alignment graph
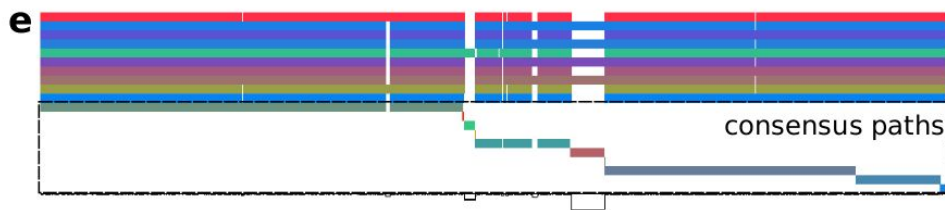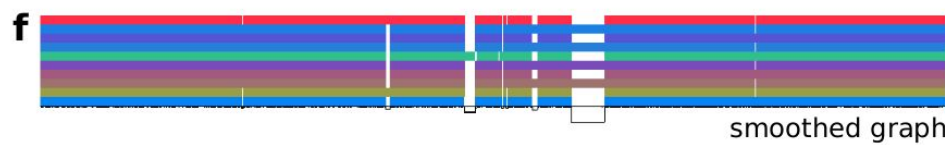
A
Z — B — Q — D
        C

Induced variation graph

c
seqwish graph

d
sorted graph

e
consensus paths

f
smoothed graph

**smoothxg graph normalization**

# Building Pangenome Graphs

Solving the whole genome alignment problem in 3 steps.

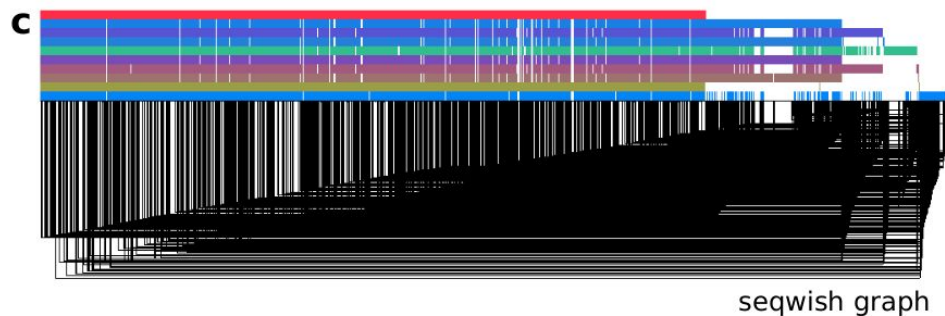a) all-to-all alignment
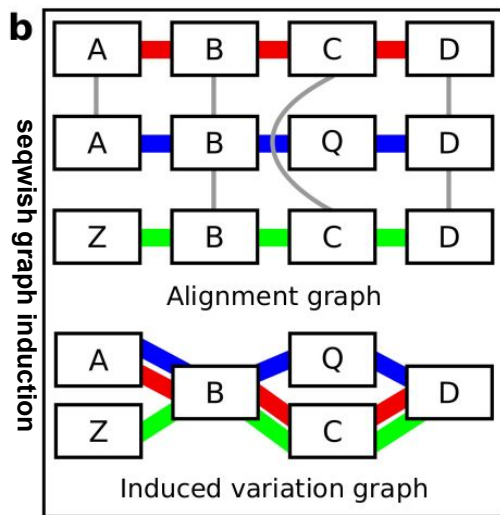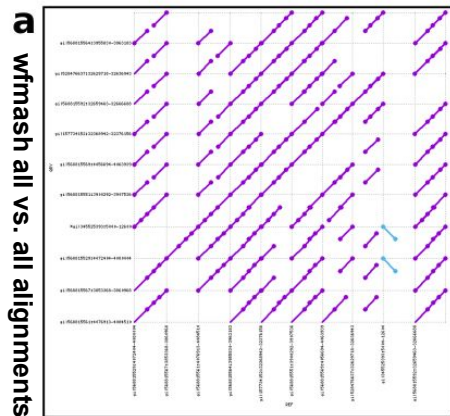b) graph induction
c-f) normalization
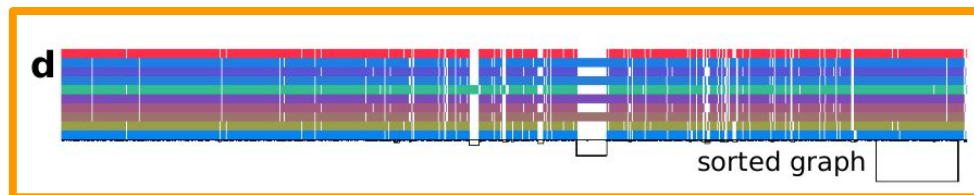
implemented in the PanGenome Graph Builder ([PGGB](#))

Erik Garrison    Andrea Guarracino
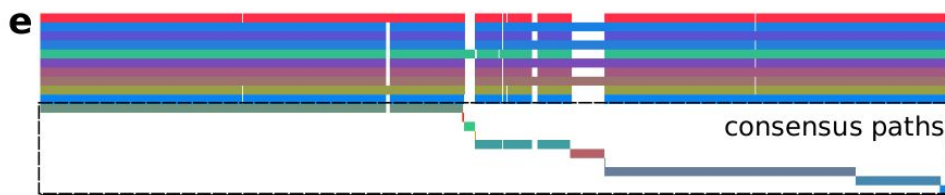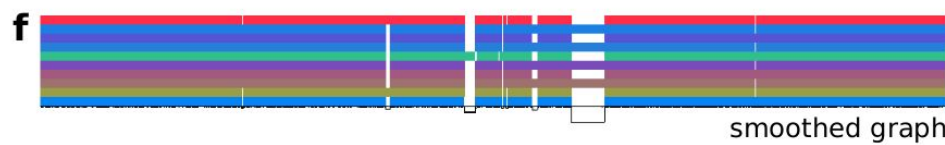


a

wfmash all vs. all alignments

b

seqwish graph induction

A — B — C — D
A — B — Q — D
Z — B — C — D

Alignment graph

A — B — Q — D
Z —    — C —

Induced variation graph

c

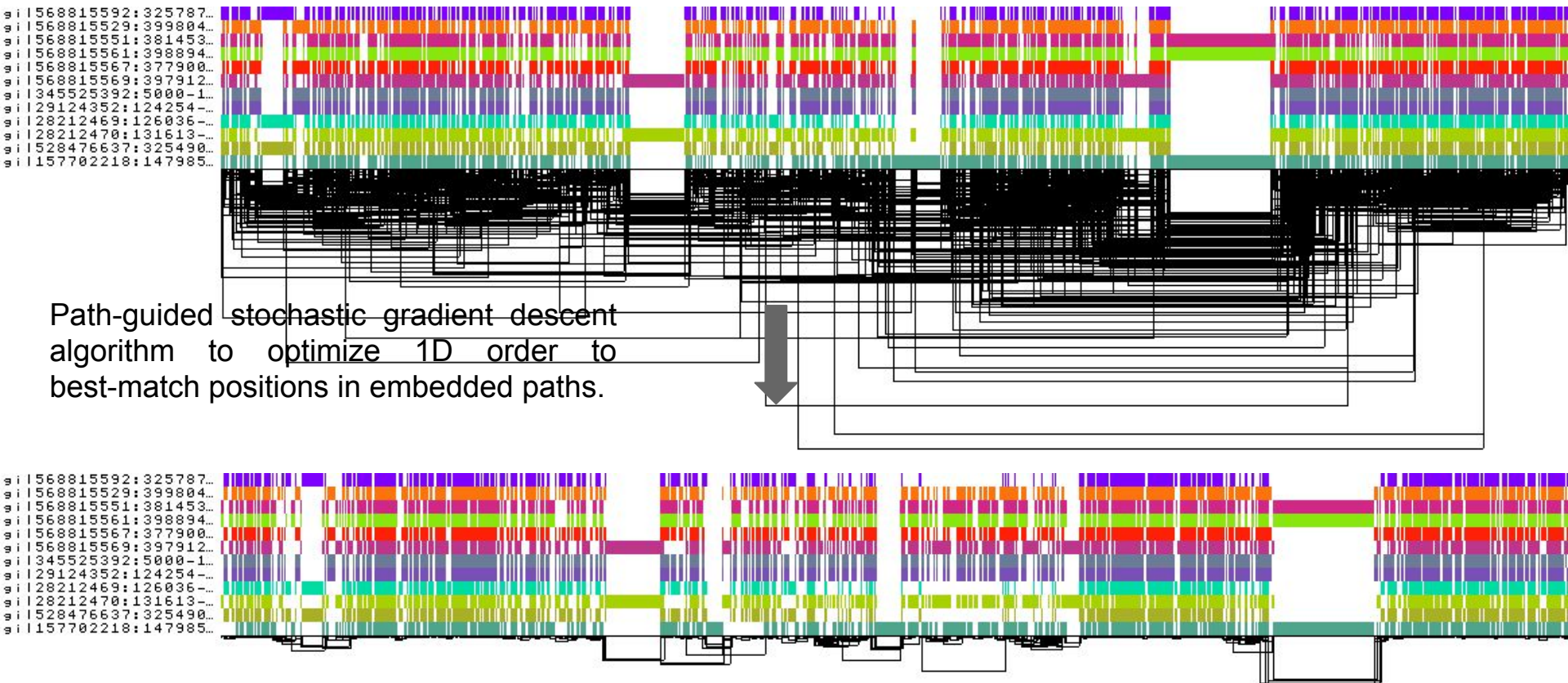seqwish graph

d

sorted graph

e

consensus paths
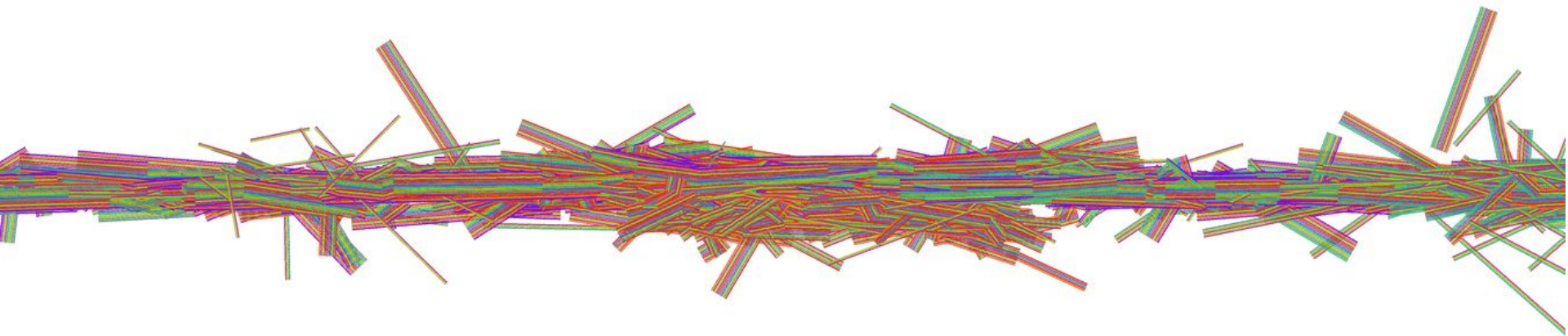
f

smoothed graph

**smoothxg graph normalization**

# Path-Guided Stochastic Gradient Descent (PG-SGD) in 1D

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



Path-guided stochastic gradient descent algorithm to optimize 1D order to best-match positions in embedded paths.
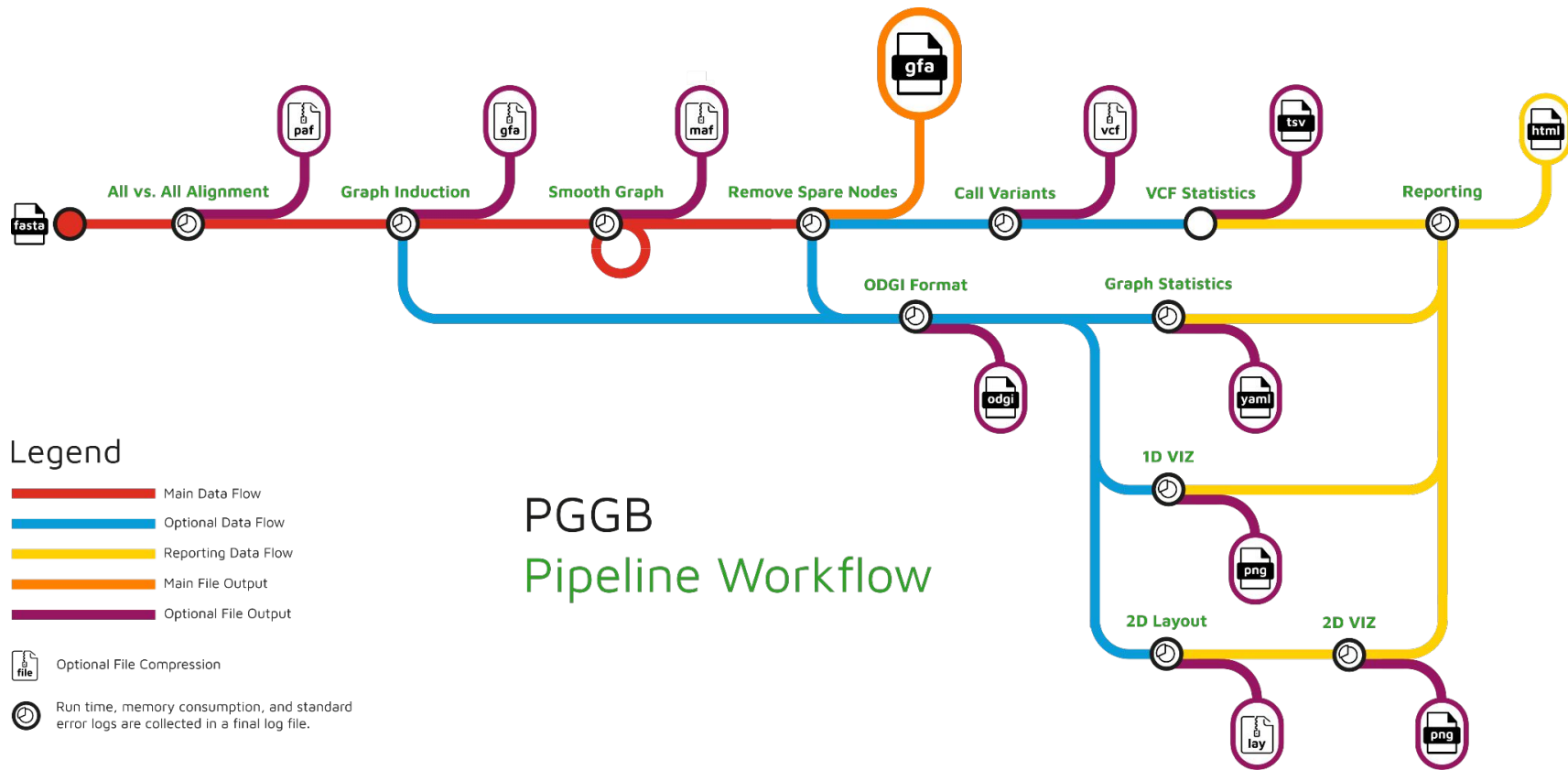
# Bonus: 2D Graph layout by PG-SGD

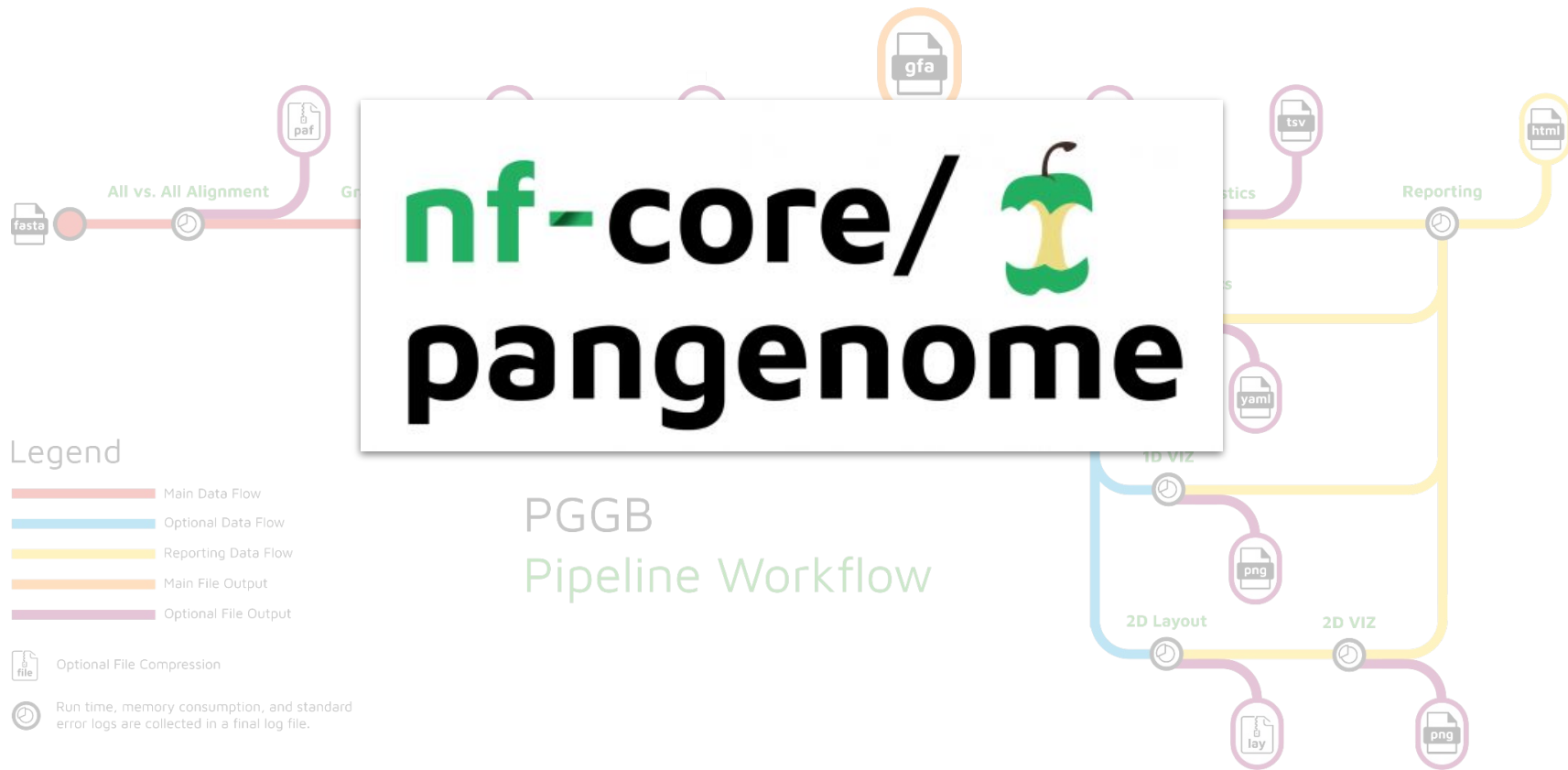Pangenome graph with 12 ALT sequences of the HLA-DRB1 gene from the GRCh38 reference genome.



Path-guided stochastic gradient descent algorithm to optimize 2D layout. Path-labeled rendering with `odgi draw`.

The layout can be plugged into gfaestus for interactive visualization.

# The PanGenome Graph Builder (PGGB) - overview

# The PanGenome Graph Builder (PGGB) - overview

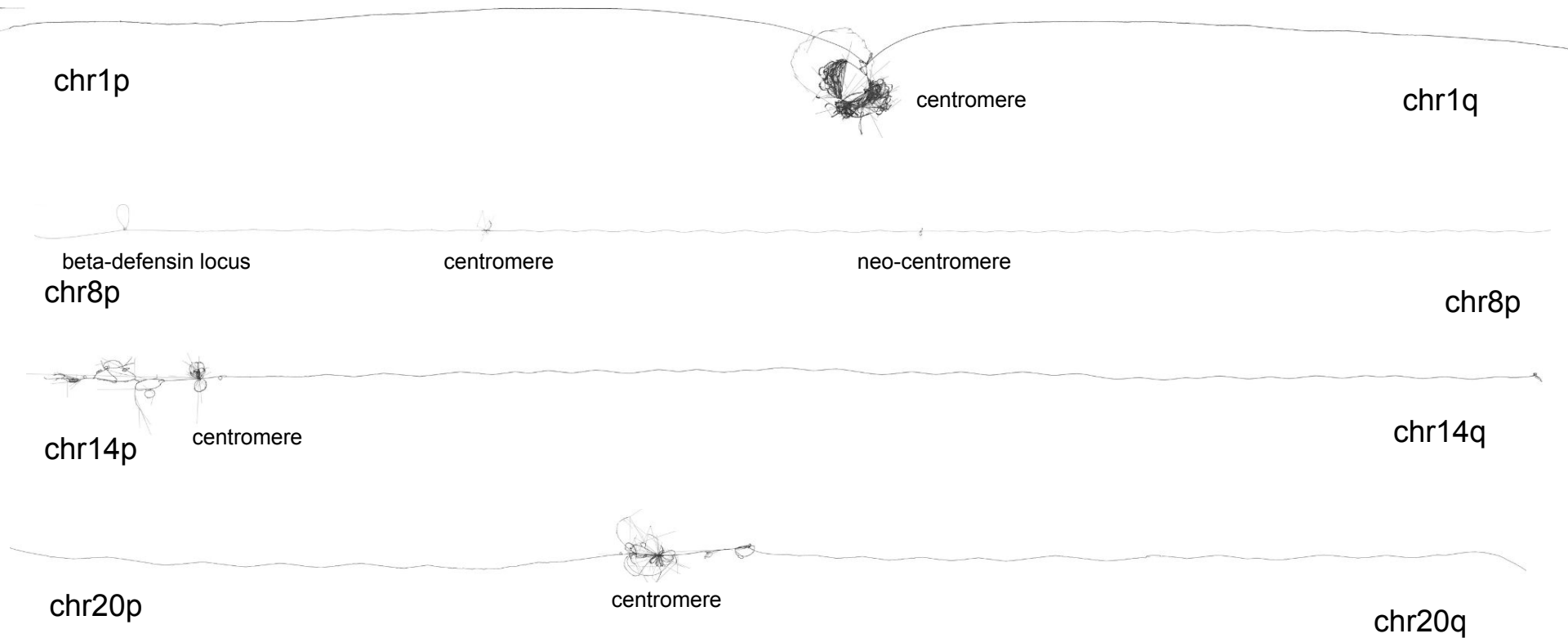# A 90 Haplotypes Human Pangenome Graph

Erik Garrison

(learned 2D visualization of PGGB HPRC chromosomes)

chr1p

centromere

chr1q

beta-defensin locus

centromere

neo-centromere

chr8p

chr8p

chr14p

centromere

chr14q

chr20p

centromere

chr20q

# Acknowledgements

Erik Garrison

Andrea Guarracino

Pjotr Prins
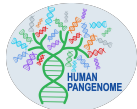
Vincenza Colonna

Flavia Villani

David G. Ashbrook

Robert W.Williams

Christian Fischer

HPRC

Sven Nahnsen

Oliver Kohlbacher

Michael Krone

Gisela Gabernet

Friederike Hanssen

Antonia Schuster

Lukas Heumos

Philipp Ehmele

Christian Kubica

Sebastian Vorbrugg

Jörg Hagmann

Jerven Bollemann

Toshiyuki T. Yokoyama

Torsten Pook

Franziska Huth

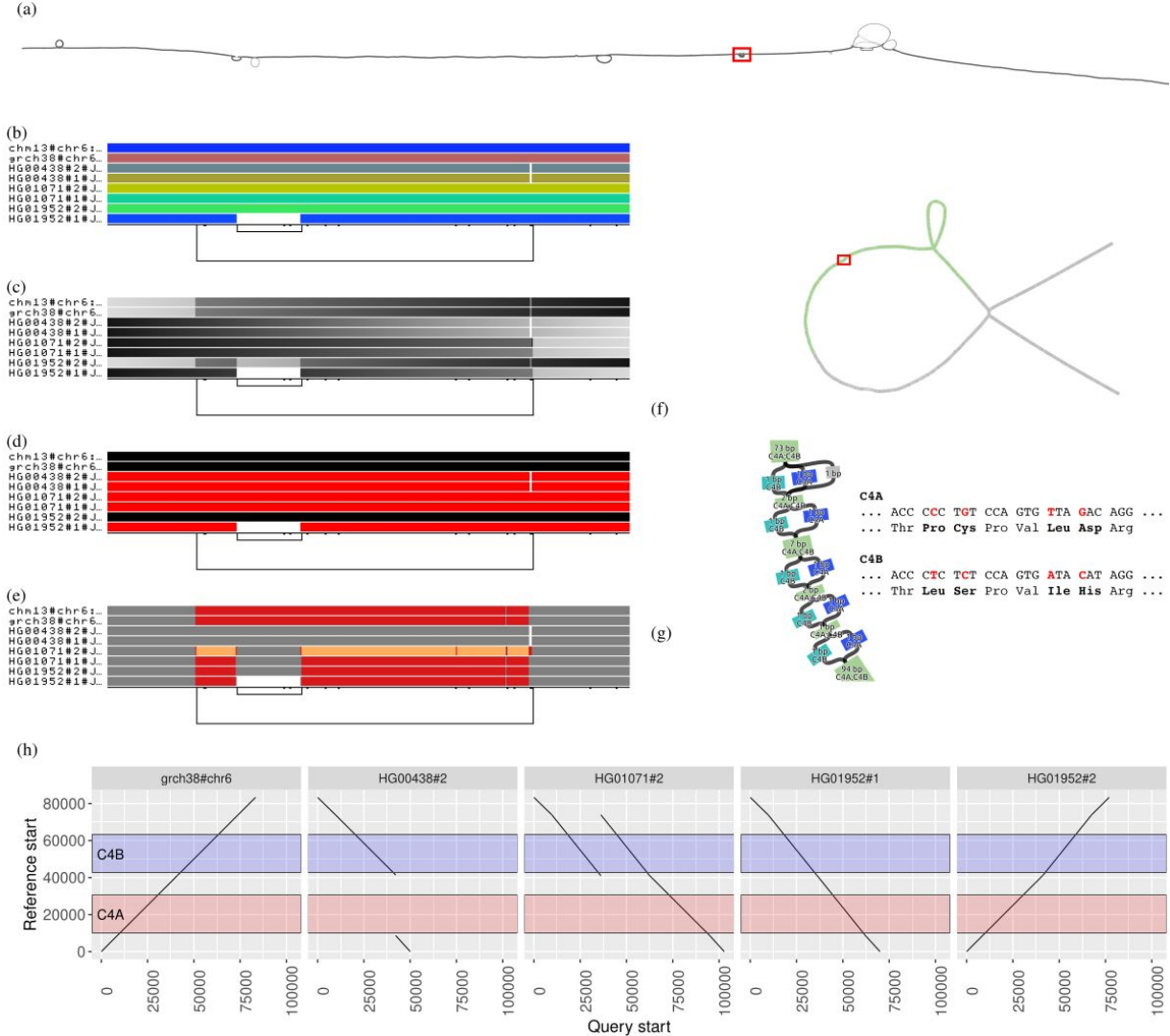Visualizing MHC and C4 pangenome graphs.

(a) `odgi extract`, `odgi layout`, and `odgi draw`

(b)-(e) `odgi viz` of 8 paths

(f) `Bandage` layout annotated with `odgi position` in green including the HERV sequence

(g) Annotated `Bandage` layout indicating single nucleotide differences in C4A and C4B

(h) `odgi untangle` output showing copy number state with respect to CHM13

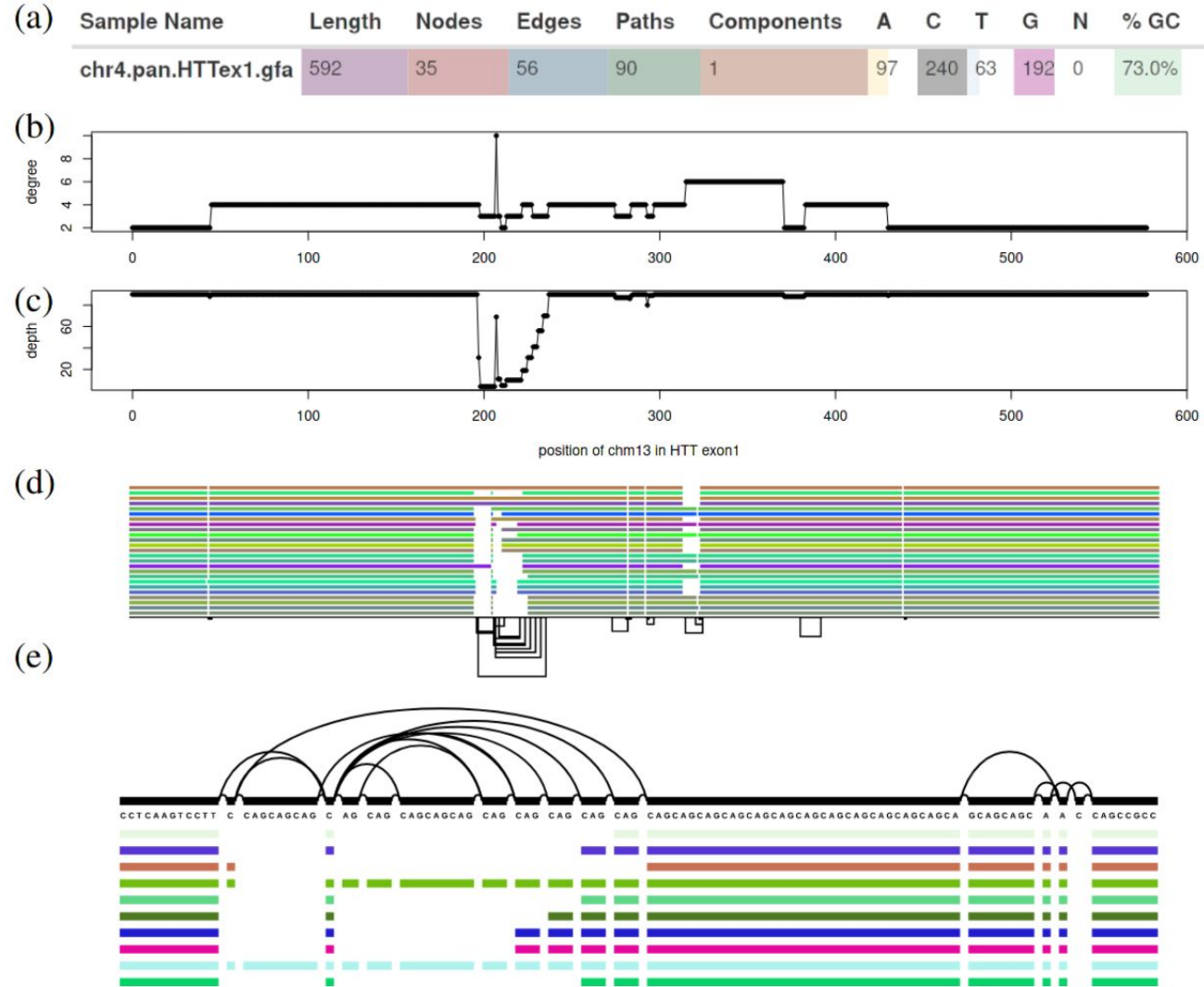Features of a 90-haplotype human pangenome graph of exon 1 of the huntingtin gene (*HTTexon1).*

(a) ODGI `MultiQC` report excerpt

(b) Node degree distribution of CHM13 in the graph

(c) Node depth distribution of CHM13 in the graph

(d) `odgi viz` – 23 largest gene alleles

(e) `vg viz` – nucleotide level zoomed in on 10 largest gene alleles



| (a) | Sample Name | Length | Nodes | Edges | Paths | Components | A | C | T | G | N | % GC |
|-----|-------------|--------|-------|-------|-------|------------|-----|-----|-----|-----|-----|------|
| | chr4.pan.HTTex1.gfa | 592 | 35 | 56 | 90 | 1 | 97 | 240 | 63 | 192 | 0 | 73.0% |

# 1D Sorting process explained



AAGTACTGGACTACGTA
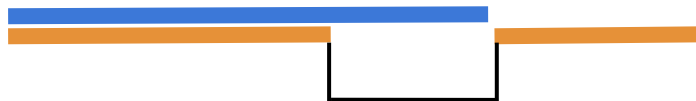
3 -> 1
2 -> 2
1 -> 3

ACTACGTACTGGAAGTA

S    1    AAGTA
S    2    CTGG
S    3    ACTACGTA
P    Genome1    3+,2+
P    Genome2    3+,1+
L    3    +    2    +
L    3    +    1    +
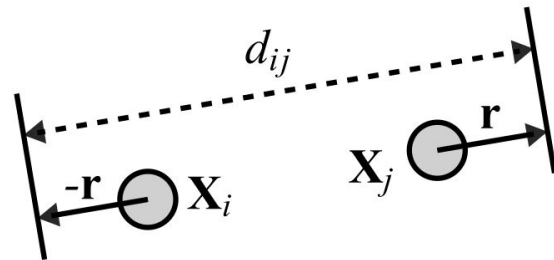

S    1    ACTACGTA
S    2    CTGG
S    3    AAGTA
P    Genome1    1+,2+
P    Genome2    1+,3+
L    1    +    2    +
L    1    +    3    +

# 1D Graph Sorting by P-SGD - The Algorithm Explained

**Objective:** Move a single pair of nodes at a time.

Optimizing the disparity between the layout distance of a node pair and the actual nucleotide distance of a path traversing these nodes.

- The first node $X_i$ of a pair is a uniform path step pick from all nodes.

# 1D Graph Sorting by P-SGD - The Algorithm Explained

**Objective:** Move a single pair of nodes at a time.

Optimizing the disparity between the layout distance of a node pair and the actual nucleotide distance of a path traversing these nodes.

- The first node $X_i$ of a pair is a uniform path step pick from all nodes.
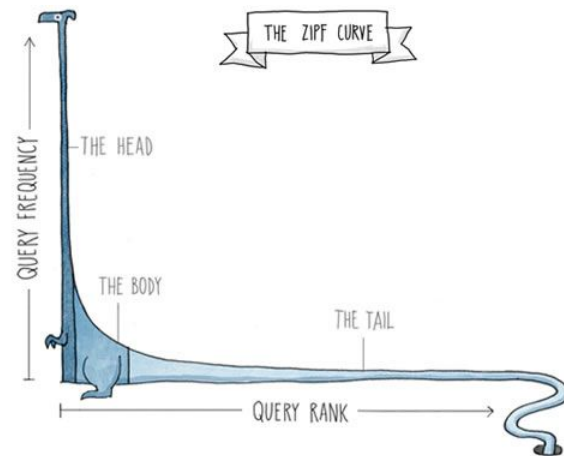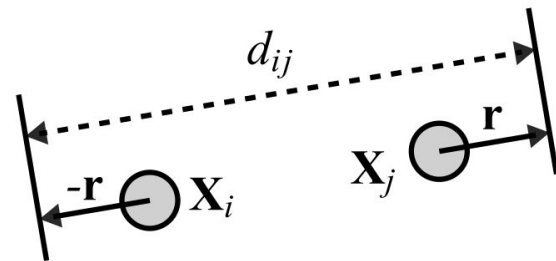- The second node $X_j$ of a pair is sampled from the same path following a Zipfian distribution.

# 1D Graph Sorting by P-SGD - The Algorithm Explained

**Objective:** Move a single pair of nodes at a time.

Optimizing the disparity between the layout distance of a node pair and the actual nucleotide distance of a path traversing these nodes.



- The first node $X_i$ of a pair is a uniform path step pick from all nodes.
- The second node $X_j$ of a pair is sampled from the same path following a Zipfian distribution.
- The path nucleotide distance of the nodes in the pair guides the actual layout distance $d_{ij}$ update of these nodes. The magnitude $r$ of the update depends on the current learning rate of the SGD.