

The PanGenome Graph Builder

Andrea Guarracino^{1*}, Simon Heumos^{2*}, Flavia Villani³, Emilio Rudbeck⁴, Kaisa Thorell⁴, Lorenzo Tattini⁵, Christian Kubica⁶, Sebastian Vorbrugg⁶, Christian Fischer⁷, Njagi Mwaniki⁸, Sven Nahnsen², David Ashbrook⁷, Robert Williams⁷, Hao Chen⁷, Vincenza Colonna³, Pjotr Prins⁷, Erik Garrison⁷

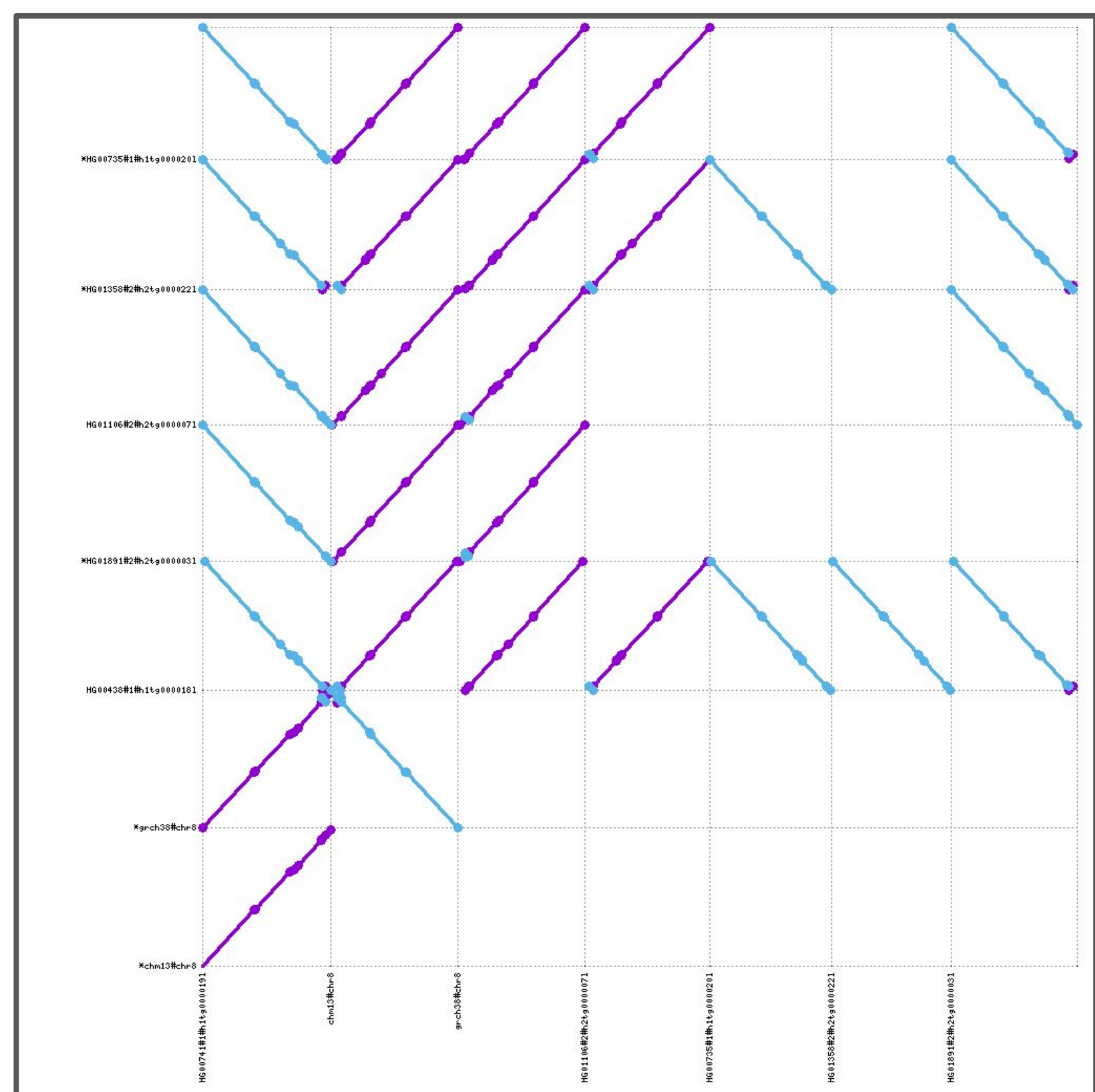
¹University of Tor Vergata, Biology, Rome, Italy, ²University of Tübingen, QBIC, Tübingen, Germany, ³CNR, IGB, Naples, Italy, ⁴University of Gothenburg, Infectious Diseases, Gothenburg, Sweden, ⁵Université Côte d'Azur, CNRS, Nice, France, ⁶MPI, Developmental Biology, Tübingen, Germany, ⁷UTHC, Genomics, Memphis, TN, USA, ⁸KEMRI-Wellcome Trust, Training, Kilifi, Kenya

*Contributed equally.

Pangenome graphs¹ contain the full genomic information of a species. For the unbiased evaluation of a species' genomic variation space, all versus all comparisons are essential. We implement the PanGenome Graph Builder (PGGB), scaling efficiently to large collections of multi-gigabase genomes. The method does not require a reference. It consists of three phases. First, we generate all versus all alignments of the input sequences. Then, the sequences and alignments are used to induce a variation graph. Finally, the graph is normalized by sorting it with an unsupervised machine learning method and applying partial order alignment to blocks in the sorted order. We apply PGGB to sequence data of different species. The resulting graphs provide excellent targets for the mapping of short and long reads, and are a basis for comparative genomic applications.

To provide a visual explanation of PGGB, we apply it to human chromosome 8. As input, we use 2 reference assemblies (GRCh38 and CHM13) and 6 *de novo* assemblies from the Human Pangenome Reference Consortium' year 1 assemblies which span all or most of the chromosome.

PGGB has three distinct phases which require 20 minutes in total to obtain this graph: (A) all-to-all alignment with [wfmash](#), (B) graph induction with [seqwish](#), and (C) normalization with [smoothxg](#), which produces the resulting graph shown. In (D) we display features that are visible in the structure of the graph³. The whole run takes around 20 minutes on a HPC compute node with an AMD EPYC 7402P 24-Core Processor and 128GB of RAM. In practice, we run with full human genomes by partitioning the input into chromosome-specific jobs, allowing turnaround of a full human pangenome from 90 haplotypes on a modest compute cluster in around a day.

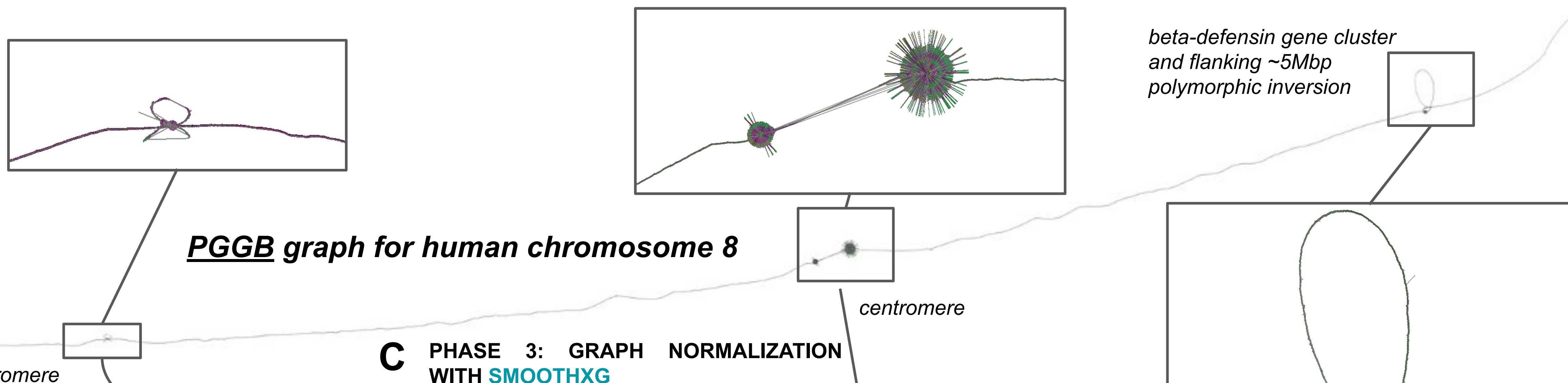
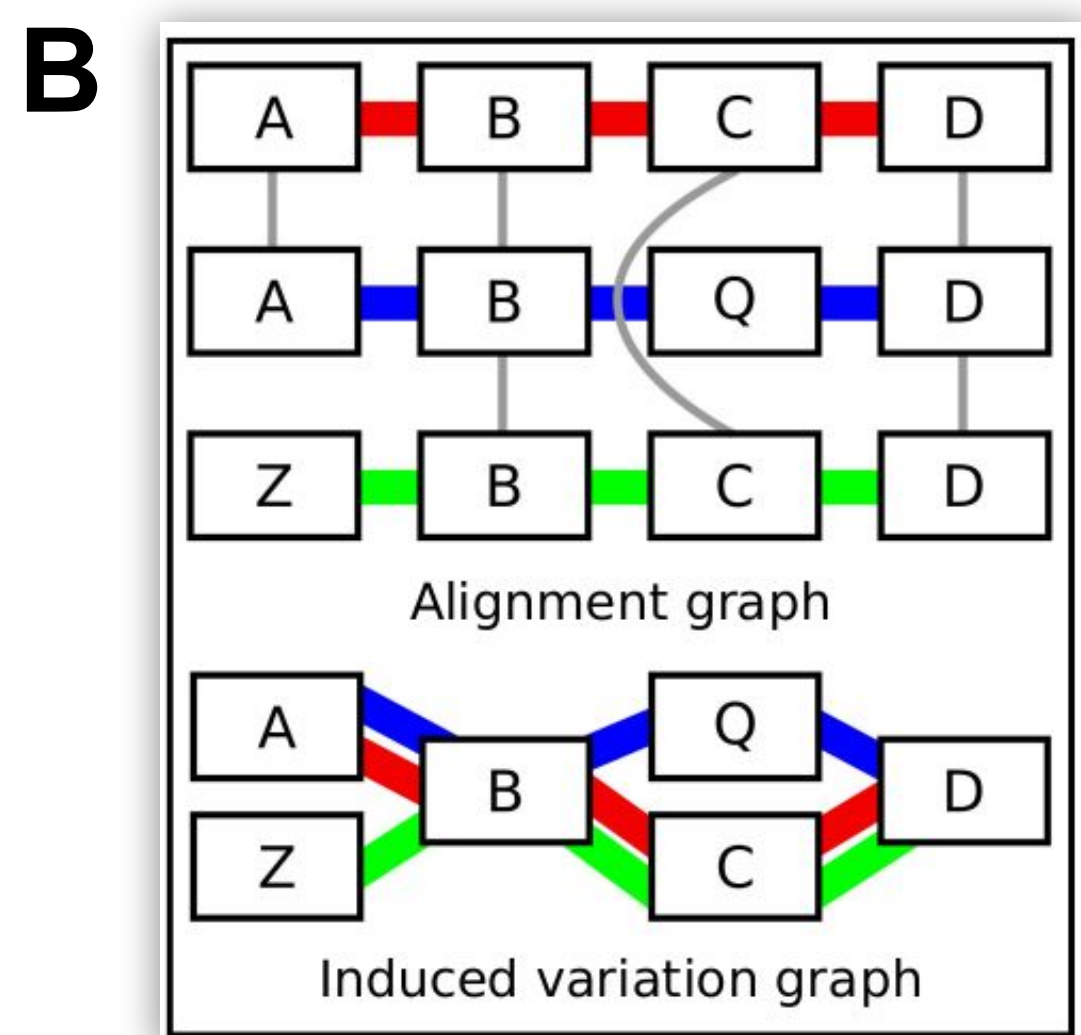


A PHASE 1: ALL-VS-ALL ALIGNMENTS WITH [WFMASH](#)

High-performance alignment of whole chromosomes is enabled by first applying [MashMap2](#) with a 100kb segment length and 98% identity filter. A hierarchical implementation of the wavefront algorithm² allows us to obtain base-level global alignments for all mappings. Here, a dotplot shows the alignment relationships from which the above graph is built.

PHASE 2: GRAPH INDUCTION WITH [SEQWISH](#)

To build a graph, we render it as an alignment graph in which input sequences are nodes and special "alignment edges" connect sequences that have been aligned. To build a variation graph, we collapse the components of the graph connected by alignments and retrace the original paths through the graph.



C PHASE 3: GRAPH NORMALIZATION WITH [SMOOTHXG](#)

A partial order alignment is run for each part of the graph, resulting in a normalized graph with base-level resolution of all variant classes.

a) binned graph visualization

```
HG00438#1#h1t#0000181
HG01358#2#h2t#0000221
HG00735#1#h1t#0000201
HG01106#2#h2t#0000071
HG01091#2#h2t#0000031
srch38#chr8
chr13#chr8
HG00741#1#h1t#0000191
```

b) colored by path position (light = start, dark = end)

```
HG00438#1#h1t#0000181
HG01358#2#h2t#0000221
HG00735#1#h1t#0000201
HG01106#2#h2t#0000071
HG01091#2#h2t#0000031
srch38#chr8
chr13#chr8
HG00741#1#h1t#0000191
```

c) colored by depth (black = 1x, green >= 1, light green = highX)

```
HG00438#1#h1t#0000181
HG01358#2#h2t#0000221
HG00735#1#h1t#0000201
HG01106#2#h2t#0000071
HG01091#2#h2t#0000031
srch38#chr8
chr13#chr8
HG00741#1#h1t#0000191
```

d) coloring by orientation (black = forward, red = reverse)

```
HG00438#1#h1t#0000181
HG01358#2#h2t#0000221
HG00735#1#h1t#0000201
HG01106#2#h2t#0000071
HG01091#2#h2t#0000031
srch38#chr8
chr13#chr8
HG00741#1#h1t#0000191
```

Graph visualizations obtained by applying [odgi viz](#). In all the visualizations

- The graph nodes' are arranged from left to right forming the pangenome sequence.
- The colored bars represent the linearized renderings of the embedded paths versus this pangenome sequence in a binary matrix
- The black lines under the paths, so called links, represent the topology of the graph.

Here (c), the copy number variation in the beta-defensin cluster is apparent as a light green stripe near the beginning of the inversion.

d) The inversion can be seen in flips of orientation between the chromosome paths embedded in the graph and the graph's orientation.

References

1. Eizenga et al. (2020). Pangenome Graphs. *Annual Reviews of Genomics and Human Genetics*, 21, 1.
2. Marco-Sola et al. (2020). Fast gap-affine pairwise alignment using the wavefront algorithm. *Bioinformatics*.
3. Logsdon et al. (2021). The structure, function and evolution of a complete chromosome 8. *Nature*. 593. 101-107.

Acknowledgements

We thank Vincenza Colonna for organizing the Crusco Summer Hackathon and the Forentum Ritrovato museum for hosting it.

We thank the deNBI cloud for providing computational resources.

S.H. acknowledges funding from the Central Innovation Programme (ZIM) for SMEs of the Federal Ministry for Economic Affairs and Energy of Germany.