

Genome Analysis

Cluster-efficient pangenome graph construction with nf-core/pangenome

Simon Heumos ^{ID 1,2,3,4,*}, Michael L. Heuer ^{ID 5}, Friederike Hanssen ^{ID 1,2,3,4},
Lukas Heumos ^{ID 6,7,8}, Andrea Guarracino ^{ID 9,10}, Peter Heringer ^{ID 1,2,3,4},
Philipp Ehmele ^{ID 6}, Pjotr Prins ^{ID 9}, Erik Garrison ^{ID 9}, Sven Nahnsen ^{ID 1,2,3,4,*}

¹Quantitative Biology Center (QBiC) Tübingen, University of Tübingen, Tübingen, Germany

²Biomedical Data Science, Dept. of Computer Science, University of Tübingen, Tübingen, Germany

³M3 Research Center, University Hospital Tübingen, Tübingen, Germany

⁴Institute for Bioinformatics and Medical Informatics (IBMI), Eberhard-Karls University of Tübingen, Tübingen, Germany

⁵University of California, Berkeley, Berkeley, 94720, California, USA

⁶Institute of Computational Biology, Department of Computational Health, Helmholtz Munich, Germany

⁷Comprehensive Pneumology Center with the CPC-M bioArchive, Helmholtz Zentrum Munich, Member of the German Center for Lung Research (DZL), Munich, Germany

⁸TUM School of Life Sciences Weihenstephan, Technical University of Munich, Freising, Germany

⁹Department of Genetics, Genomics and Informatics, University of Tennessee Health Science Center, 71 S Manassas St, Memphis, 38163, Tennessee, USA

¹⁰Human Technopole, Viale Rita Levi-Montalcini 1, 20157, Milan, Italy

*To whom correspondence should be addressed.

Associate Editor: XXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: Pangenome graphs offer a comprehensive way of capturing genomic variability across multiple genomes. However, current construction methods often introduce biases, excluding complex sequences or relying on references. The PanGenome Graph Builder (PGGB) addresses these issues. To date, though, there is no state-of-the-art pipeline allowing for easy deployment, efficient and dynamic use of available resources, and scalable usage at the same time.

Results: To overcome these limitations, we present *nf-core/pangenome*, a reference-unbiased approach implemented in Nextflow following *nf-core*'s best practices. Leveraging biocontainers ensures portability and seamless deployment in HPC environments. Unlike PGGB, *nf-core/pangenome* distributes alignments across cluster nodes, enabling scalability. Demonstrating its efficiency, we constructed pangenome graphs for 1000 human chromosome 19 haplotypes and 2146 *E. coli* sequences, achieving a two to threefold speedup compared to PGGB without increasing greenhouse gas emissions.

Availability: *nf-core/pangenome* is released under the MIT open-source license, available on GitHub and Zenodo, with documentation accessible at <https://nf-co.re/pangenome/1.1.2/docs/usage>.

Supplementary: Supplementary data are available at *Bioinformatics* online.

Contact: simon.heumos@qbic.uni-tuebingen.de, sven.nahnsen@qbic.uni-tuebingen.de

1 Introduction

The availability of high-quality population-wide whole-genome assemblies (Liao *et al.*, 2023; Kang *et al.*, 2023; Weller *et al.*, 2023; Zhou *et al.*, 2022; Liu *et al.*, 2020; Leonard *et al.*, 2022) offers new opportunities to study sequence evolution and variation within and between genomic populations. A challenge is simultaneously representing and analyzing

hundreds to thousands of genomes at a gigabase scale. One solution here is a pangenome. It models a population's entire set of genomic sequences (Ballouz *et al.*, 2019). In contrast to reference-based genomic approaches, which relate sequences to a linear genome, pangenomics relates each new sequence to all the others represented in the pangenome (The Computational Pan-Genomics Consortium, 2016; Eizenga *et al.*, 2020; Sherman and Salzberg, 2020) minimizing reference-bias. Pangenomes can be described as sequence graphs which store DNA sequences in nodes with

edges connecting the nodes as they occur in the individual sequences (Hein, 1989). Genomes are encoded as paths traversing the nodes (Garrison *et al.*, 2018).

Current pangenome graph construction methods exclude complex sequences or are reference-biased (Chin *et al.*, 2023; Minkin *et al.*, 2016). One recent approach that overcomes such limitations is the PanGenome Graph Builder (PGGB) pipeline (Garrison *et al.*, 2023). PGGB iteratively refines an all-to-all whole-genome alignment graph that lets us explore sequence conservation and variation, infer phylogeny, and identify recombination events. PGGB has been extensively evaluated (Garrison *et al.*, 2023; Andreace *et al.*, 2023) and applied to build the first draft human pangenome reference (Liao *et al.*, 2023). However, PGGB is implemented in bash, which (a) makes it difficult to deploy on HPC systems, (b) does not allow for a fine granular tuning of computing resources for different steps of the pipeline (Sztuka *et al.*, 2024), and (c) limits its cluster scalability to one node. These limitations greatly hinder the broad application of large-scale pangenomes.

To compensate for that, we wrote *nf-core/pangenome*, a reference-unbiased approach to construct pangenome graphs. Mirroring PGGB, *nf-core/pangenome* is implemented in Nextflow (Di Tommaso *et al.*, 2017). In contrast to PGGB, *nf-core/pangenome* can distribute the quadratic all-to-all base-level alignments across nodes of a cluster by splitting the approximate alignments into problems of equal size. We benchmarked the time spent on base-pair level alignments and show that it is reduced linearly with an increase in alignment problem chunks (Suppl. 5.5). We showcase the workflow’s scalability by applying it to 1000 chromosome 19 human haplotypes and 2146 *E. coli* sequences, which were built in less than half the time PGGB required while not increasing the CO₂ equivalent (CO₂e) emissions.

2 Material and Methods

2.1 Pipeline overview

The pipeline’s (Fig. 1a) input is a FASTA file compressed with *bgzip* (Li *et al.*, 2009) containing the sequences to create the graph. Sequence names should follow the Pangenome Sequence Naming specification (PanSN-spec) (<https://github.com/pangenome/PanSN-spec>, last accessed October 2024). The primary output is a pangenome variation graph (Garrison *et al.*, 2018) in the Graphical Fragment Assembly (GFA) format version 1 (<http://gfa-spec.github.io/GFA-spec/GFA1.html>, last accessed October 2024).

2.1.1 Core workflow

The core workflow of *nf-core/pangenome* mirrors PGGB (Fig. 1a) with additional enhancements: (a) All concurrent processes can be run in parallel. (b) Each process can be given individual computing resources.

The pipeline begins with an all-to-all alignment of the input sequences using the whole-chromosome pairwise sequence aligner WFMASH (<https://github.com/waveygang/wfmash>, last accessed October 2024), avoiding reference, order, or orientation bias, allowing every sequence to serve as a reference. In the pangenome graph induction step SEQWISH (Garrison and Guarracino, 2022), an alignment to variation graph inducer, converts the sequence alignments into a variation graph. The graph is then simplified using SMOOTHXG (Garrison *et al.*, 2023): A 1-dimensional (1D) graph embedding (Heumos *et al.*, 2024) orders the graph’s nodes to best match the nucleotide distances of the genomic paths of the graph. Next, the graph is split into partially overlapping segments. The sequences of each segment are realigned with a local Multiple Sequence Alignment (MSA) kernel, partial order alignment (POA) (Lee *et al.*, 2002). Afterwards, the segments are laced back together into a variation graph. By default, the SMOOTHXG process is applied 3 times in order to smooth the edge effects at the boundaries of the segments. Finally, we employ

GFAFFIX (Liao *et al.*, 2023) to systematically condense redundant nodes within the graph.

Graph quality is assessed with ODGI (Guarracino *et al.*, 2022), which provides statistics and visualizations. Optionally, variants can be called against any (reference) path(s) in the graph using *vg deconstruct* (Garrison *et al.*, 2018). Results are summarized in a MultiQC (Ewels *et al.*, 2016) report. Pipeline implementation details are given in Suppl. 5.1.

If desired, the pipeline performs community detection to identify clusters of related sequences in the pangenome graph, revealing biological patterns such as conserved or divergent regions across genomes (Supplementary 5.2), with the core workflow executed for each community in parallel.

3 Results

3.1 Building a 1000 haplotypes chr19 pangenome graph

The Human Pangenome Resource Consortium (HPRC) recently built a draft human pangenome of 90 haplotypes. However, haplotype data for thousands of individuals was already generated by the 1000 Genomes Project (1KGP) (Durbin *et al.*, 2010). As a use case, we used *nf-core/pangenome* to build a graph of 1000 chromosome 19 haplotypes (Kuhnle *et al.*, 2020) in 3 days, emitting 51.07 kg CO₂e. PGGB took 7 days for the same task (56.32 kg CO₂e). In Fig. 1b the pangenome growth curve generated with PANACUS (Liao *et al.*, 2023) shows nucleotide growth as more haplotypes are added. The softcore pangenome, defined as sequences traversed by 95% of haplotypes, comprises the majority of the pangenome even with 1000 haplotypes. This stability may be due to the exclusion of complex regions like the centromere in the short-read data.

3.2 Building a 2146 sequences *E. coli* pangenome graph

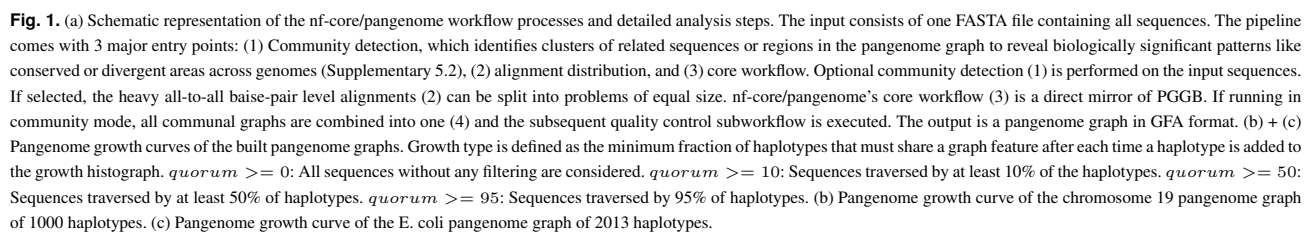
To evaluate the pipeline’s scalability, we built a pangenome graph of 2146 *E. coli* sequences. The *nf-core/pangenome* graph was completed in 10 days, emitting 175.18 kg CO₂e, while PGGB could not finish within 30 days due to cluster time restrictions. For the growth curve (Fig. 1c) we excluded 130 plasmid sequences. The softcore pangenome remains stable at ~3Mb, but its size constitutes less than 10% of the total pangenome. This substantial pangenomic growth is likely driven by horizontal gene transfer, as bacteria incorporate genes from one another at various genomic locations. Other reasons could be sequencing errors or human contamination (Breitwieser *et al.*, 2019).

4 Discussion

We implemented *nf-core/pangenome*, an easy-to-install, portable, and cluster-scalable pipeline for unbiased pangenome variation graph construction. It is the first pangenomic pipeline within the *nf-core* framework that enables the comparative analysis of gigabase-scale pangenome datasets. While tools like Minigraph (Li *et al.*, 2020) or PGR-TK (Chin *et al.*, 2023) also address pangenome analysis, *nf-core/pangenome* uniquely integrates this capability into the standardized *nf-core* framework, offering compatibility with a wide range of modular workflows and community-developed best practices.

The pipeline’s core workflow has been successfully applied to *Neisseria meningitidis* (Yang *et al.*, 2023), wild grapes (Cochetel *et al.*, 2023), humans (Guarracino *et al.*, 2023; Liao *et al.*, 2023), grapevines (Guo *et al.*, 2024), taurines (Milia *et al.*, 2024), and rats (Villani *et al.*, 2024) underpinning the community effort to focus on a best-practice workflow to create reference-unbiased and sequence complete pangenome graphs. The modular domain-specific language (DSL) 2 pipeline structure facilitates easy exchange of processes with alternative tools, expanding its functionality and integration with other (sub-)workflows.

We have shown that we are able to perform all-vs-all base pair level alignments of thousands of sequences. When executed on an HPC,



would allow us to break the whole genome multiple alignments into smaller pieces, construct a pangenome graph for each piece, and lace these together into a full graph with gfalace (<https://github.com/pangenome/gfalace>, last accessed October 2024).

Software and data availability

Code and data resources for this manuscript and its figures are available in the public repository: <https://github.com/subwaystation/pangenome-paper>.

Acknowledgments

We thank Matthias Seybold from QBIC for maintaining the Core Facility Cluster. We thank Sabrina Krakau from QBIC for giving feedback to the nf-co2footprint plugin section. We are grateful to the nf-core community

for their support during the implementation of the pipeline. From the nf-core community, we want to thank Matthias Hörtenhuber, Maxime Garcia, Susanne Jodoin, Julia Mir Petrol, Adam Talbot, and Gisela Gabernet.

Funding

S.H. acknowledges funding from the Central Innovation Programme (ZIM) for SMEs of the Federal Ministry for Economic Affairs and Energy of Germany. This work was supported by the BMBF-funded de.NBI Cloud within the German Network for Bioinformatics Infrastructure (de.NBI) (031A532B, 031A533A, 031A533B, 031A534A, 031A535A, 031A537A, 031A537B, 031A537C, 031A537D, 031A538A). A.G. acknowledges support from the Human Technopole. S.N. acknowledges support from iFIT funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy—EXC 2180—390900677 and CMFI under EXC 2124—390838134. We gratefully acknowledge support from NIH/NIDA U01DA047638 (E.G.), NIH/NIGMS R01GM123489 (E.G., P.P.), and NSF PPOSS Award #2118709 (E.G., P.P.), and the CITG (E.G.).

Competing interests

Author L.H. is employed by LaminLabs.

References

- Andreace, F. *et al.* (2023). Comparing methods for constructing and representing human pangenome graphs. *Genome Biology*, **24**(1).
- Ballouz, S. *et al.* (2019). Is it time to change the reference genome? *Genome Biology*, **20**(1), 159.
- Breitwieser, F. P. *et al.* (2019). Human contamination in bacterial genomes has created thousands of spurious proteins. *Genome Research*, **29**(6), 954–960.
- Chin, C.-S. *et al.* (2023). Multiscale analysis of pangenomes enables improved representation of genomic diversity for repetitive and clinically relevant genes. *Nature Methods*, **20**(8), 1213–1221.
- Cochetel, N. *et al.* (2023). A super-pangenome of the North American wild grape species. *Genome Biology*, **24**(1).
- Di Tommaso, P. *et al.* (2017). Nextflow enables reproducible computational workflows. *Nature Biotechnology*, **35**(4), 316–319.
- Durbin, R. M. *et al.* (2010). A map of human genome variation from population-scale sequencing. *Nature*, **467**(7319), 1061–1073.
- Eizenga, J. M. *et al.* (2020). Pangenome graphs. *Annual Review of Genomics and Human Genetics*, **21**(1), 139–162.
- Ewels, P. *et al.* (2016). MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics*, **32**(19), 3047–3048.
- Garrison, E. and Guarracino, A. (2022). Unbiased pangenome graphs. *Bioinformatics*, **39**(1).
- Garrison, E. *et al.* (2018). Variation Graph Toolkit Improves Read Mapping by Representing Genetic Variation in the Reference. *Nature Biotechnology*, **36**(9), 875–879.
- Garrison, E. *et al.* (2023). Building pangenome graphs. *bioRxiv*.
- Guarracino, A. *et al.* (2022). ODGI: understanding pangenome graphs. *Bioinformatics*, **38**(13), 3319–3326.
- Guarracino, A. *et al.* (2023). Recombination between heterologous human acrocentric chromosomes. *Nature*, **617**(7960), 335–343.
- Guo, L. *et al.* (2024). Super Pangenome of Grapevines Empowers Improvement of the Oldest Domesticated Fruit. *bioRxiv*.
- Hein, J. (1989). A new method that simultaneously aligns and reconstructs ancestral sequences for any number of homologous sequences, when the phylogeny is given. *Molecular Biology and Evolution*.
- Heumos, S. *et al.* (2024). Pangenome graph layout by path-guided stochastic gradient descent. *Bioinformatics*, **40**(7).
- Hickey, G. *et al.* (2023). Pangenome graph construction from genome alignments with minigraph-cactus. *Nature Biotechnology*, **42**(4), 663–673.
- Kang, M. *et al.* (2023). The pan-genome and local adaptation of *Arabidopsis thaliana*. *Nature Communications*, **14**(1).
- Kuhnle, A. *et al.* (2020). Efficient construction of a complete index for pan-genomics read alignment. *Journal of Computational Biology*, **27**(4), 500–513.
- Lannelongue, L. *et al.* (2021). Green algorithms: Quantifying the carbon footprint of computation. *Advanced Science*, **8**(12).
- Lee, C. *et al.* (2002). Multiple sequence alignment using partial order graphs. *Bioinformatics*, **18**(3), 452–464.
- Leonard, A. S. *et al.* (2022). Structural variant-based pangenome construction has low sensitivity to variability of haplotype-resolved bovine assemblies. *Nature Communications*, **13**(1).
- Li, H. *et al.* (2009). The sequence alignment/map format and samtools. *Bioinformatics*, **25**(16), 2078–2079. 19505943[pmid].
- Li, H. *et al.* (2020). The Design and Construction of Reference Pangenome Graphs with Minigraph. *Genome Biology*, **21**(1), 265.
- Liao, W.-W. *et al.* (2023). A draft human pangenome reference. *Nature*, **617**(7960), 312–324.
- Liu, Y. *et al.* (2020). Pan-genome of wild and cultivated soybeans. *Cell*, **182**(1), 162–176.e13.
- Milia, S. *et al.* (2024). Taurine pangenome uncovers a segmental duplication upstream of *KIT* associated with depigmentation in white-headed cattle. *bioRxiv*.
- Minkin, I. *et al.* (2016). TwoPaCo: an efficient algorithm to build the compacted de Bruijn graph from many complete genomes. *Bioinformatics*, **33**(24), 4024–4032.
- Sayers, E. W. *et al.* (2021). Database resources of the national center for biotechnology information. *Nucleic Acids Research*, **50**(D1), D20–D26.
- Sherman, R. M. and Salzberg, S. L. (2020). Pan-genomics in the human genome era. *Nature Reviews Genetics*, **21**(4), 243–254.
- Sirén, J. *et al.* (2024). Personalized pangenome references. *Nature Methods*.
- Sztuka, M. *et al.* (2024). Nextflow vs. plain bash: different approaches to the parallelization of SNP calling from the whole genome sequence data. *NAR Genomics and Bioinformatics*, **6**(2), lqae040.
- The Computational Pan-Genomics Consortium (2016). Computational pan-genomics: status, promises and challenges. *Briefings in Bioinformatics*, page bbw089.
- Traag, V. A. *et al.* (2019). From Louvain to Leiden: guaranteeing well-connected communities. *Scientific Reports*, **9**(1).
- Villani, F. *et al.* (2024). Pangenome reconstruction in rats enhances genotype-phenotype mapping and novel variant discovery. *bioRxiv*.
- Vivian, J. *et al.* (2017). Toil enables reproducible, open source, big biomedical data analyses. *Nature Biotechnology*, **35**(4), 314–316.
- Weller, C. A. *et al.* (2023). Highly complete long-read genomes reveal pangenomic variation underlying yeast phenotypic diversity. *Genome Research*, **33**(5), 729–740.
- Wratten, L. *et al.* (2021). Reproducible, scalable, and shareable analysis pipelines with bioinformatics workflow managers. *Nature Methods*, **18**(10), 1161–1168.
- Yang, Z. *et al.* (2023). Pangenome graphs in infectious disease: a comprehensive genetic variation analysis of *Neisseria meningitidis* leveraging Oxford Nanopore long reads. *Frontiers in Genetics*, **14**.
- Zhou, Y. *et al.* (2022). Assembly of a pangenome for global cattle reveals missing sequences and novel structural variations, providing new insights into their diversity and evolutionary history. *Genome Research*, **32**(8), 1585–1601.