

Pangenome graphs and their applications in biodiversity genomics

Received: 23 May 2024

Accepted: 8 November 2024

Published online: 8 January 2025

 Check for updates

Simona Secomandi ^{1,8}, Guido Roberto Gallo ^{2,8}, Riccardo Rossi ³,
Carlos Rodríguez Fernandes ^{4,5}, Erich D. Jarvis ^{1,6},
Andrea Bonisoli-Alquati ⁷, Luca Gianfranceschi ² & Giulio Formenti ⁶ ✉

Complete datasets of genetic variants are key to biodiversity genomic studies. Long-read sequencing technologies allow the routine assembly of highly contiguous, haplotype-resolved reference genomes. However, even when complete, reference genomes from a single individual may bias downstream analyses and fail to adequately represent genetic diversity within a population or species. Pangenome graphs assembled from aligned collections of high-quality genomes can overcome representation bias by integrating sequence information from multiple genomes from the same population, species or genus into a single reference. Here, we review the available tools and data structures to build, visualize and manipulate pangenome graphs while providing practical examples and discussing their applications in biodiversity and conservation genomics across the tree of life.

The decrease in DNA sequencing costs is boosting biodiversity genomics^{1,2}. Large-scale international initiatives aim to generate highly contiguous, haplotype-resolved reference genomes for all species^{3–7} and to characterize biodiversity, clarify its evolution and help its conservation. Reference genomes (Box 1) are the backbones to annotate genomic variation, helping its preservation in the current biodiversity crisis². Haplotype-resolved reference genomes better capture genetic variation across individuals, enabling more effective mapping of phenomes to genomes and illuminating both the adaptive uniqueness of taxa⁸ and their position in the tree of life⁹. However, even with accurate reference genomes, sequences mapped against them can be misplaced or fail to align because of divergent or missing regions in the reference (Fig. 1a). These ‘blind spots’^{10–12} introduce a reference bias^{13,14}, potentially misrepresenting genetic variation. A solution to reference bias lies in pangenomics, the systematic capturing of genetic variation within dedicated composite assemblies, called pangenomes^{15,16}. Historically, a pangenome can be a collection of unaligned sequences^{17,18}, or a graph derived from raw reads¹⁹ or assembled sequences¹⁵ (Fig. 1b and Box 2). A pangenome graph is a graphical model storing genomic

data from different individuals, together with their relationships and variability, in a single data structure (Fig. 1c). Data embedded in the graph can be accessed to perform bioinformatic tasks, such as read mapping, by graph indexing^{15,20}. A pangenome graph accommodates multiple alternative alleles and provides a more comprehensive representation of genomic variation within a species^{21,22}, depending on the degree to which the haplotype set reflects the overall diversity of the species. At higher taxonomic ranks (for example, genus or family), super-pangenomes expand our ability to capture genetic diversity^{23–27}. While cataloging genetic variability within a species is critical², examples of pangenome graphs in non-model species are still limited^{28,29}. We expect this to change in the coming years, as long-read sequencing and high-quality genomes become more accessible^{2,5}. Pangenome graphs improve downstream analyses, providing researchers with the flexibility to analyze genetic variation at multiple levels in a single data structure. In agriculture, pangenomes for major crop plants proved effective in identifying resistance genes and beneficial alleles for crop improvement, while, in humans, they show potential to improve clinical research (Box 2). Pangenome graphs enable the accurate detection of

¹Laboratory of Neurogenetics of Language, the Rockefeller University, New York, NY, USA. ²Department of Biosciences, University of Milan, Milan, Italy. ³Department of Biotechnology and Biosciences, University of Milano–Bicocca, Milan, Italy. ⁴Centre for Ecology, Evolution and Environmental Changes (CE3C) and CHANGE, Global Change and Sustainability Institute, Departamento de Biologia Animal, Faculdade de Ciências, Universidade de Lisboa, Lisboa, Portugal. ⁵Faculdade de Psicologia, Universidade de Lisboa, Lisboa, Portugal. ⁶The Vertebrate Genome Laboratory, New York, NY, USA. ⁷Department of Biological Sciences, California State Polytechnic University, Pomona, Pomona, CA, USA. ⁸These authors contributed equally: Simona Secomandi, Guido Roberto Gallo. ✉e-mail: gformenti@rockefeller.edu

BOX 1**Glossary**

ATAC-seq: a technique used to assess chromatin accessibility by using Tn5 transposase to insert sequencing adaptors into open chromatin regions, which are then sequenced to map these accessible sites.

Bidirected graph: graph representing both DNA strands.

Bubble: region of the graph where a set of paths diverge from a common node and reconverge in the following common node. The paths walking through the bubble represent divergent haplotypes and their sequence variation.

ChIP-seq: a method used to analyze protein–DNA interactions by combining chromatin immunoprecipitation with next-generation sequencing to identify the binding sites of DNA-associated proteins.

Chromatin conformation data (Hi-C): Sequencing-based molecular technique used to detect regions in the genome where physical interactions are frequent. It measures the contact frequency between all pairs of loci, offering insights into the genome's three-dimensional organization.

Collinearity: preservation of the linear order of genes or genetic markers along chromosomes across different species.

Edges: node connectors indicating their concatenation and ultimately defining the DNA sequence of each haplotype (path). In a bidirected graph, the direction of the edges indicates the strandedness.

Haplotype false duplication: error that occurs when a single genomic region is represented twice as two distinct regions in the same assembly. This typically happens when a heterozygous region in the individual contains two highly divergent haplotypes, causing the assembler to mistakenly treat them as non-homologous regions.

Hybridization: breeding between individuals from genetically different lineages.

Introgression: gene flow between hybridizing populations or species through the backcrossing of hybrids with one or both of the parent populations.

k-mer: substrings of nucleotides of length k .

Microchromosomes: small-sized chromosomes typically found in the genomes of various animals, including birds and reptiles. They were often found to be gene rich and GC rich.

Nodes: the basic unit of a pangenome graph. They represent DNA sequences included in the graph. In a syntenic region, nodes can be traversed by multiple paths. In a bubble of variation, multiple nodes represent divergent DNA sequences, and only a subset of divergent haplotypes traverse each node.

Paths: representation of the haplotypes included in the graph. A path is defined by a concatenation of DNA sequences (nodes) connected by edges.

Phasing: process of determining which genetic variants are inherited together on the same chromosome from each parent.

Reference genome: a contiguous and accurate genome assembly, representative of a species.

Seed-and-extend aligner: an aligner that begins with small, exact alignment segments, known as seeds, and then attempts to extend or merge these segments to identify larger, highly similar regions.

Snarl: hierarchical generalization of a bubble. A snarl is a subgraph with a start and an end node, and paths traversing the snarl can have complex interconnections, representing variation.

Synteny: conservation (not necessarily in the same order) of blocks of genes or entire chromosomal regions across different species.

structural variants (SVs)^{30,31}, large (>50-bp) genomic rearrangements that underlie phenotype and fitness variation^{30,31} and disproportionately contribute to local adaptation³² and even speciation³³. Pangenome graphs will increase the accuracy and power of resequencing projects investigating population dynamics and insights into the genetic bases of phenotypic traits by streamlining variant phasing (Box 1), haplotype reconstruction and genotyping through imputation. They will aid in understanding genome evolution across species and benefit studies linking genetic variation and gene expression to phenotypes. Here, we illustrate the available data structures and tools for building, visualizing and manipulating pangenome graphs and provide practical advice for their use in downstream analyses. We also highlight their potential future contribution to conservation and biodiversity genomics.

Pangenomes as variation graphs

Over the years, the term pangenome acquired different meanings (Box 2). Here, we focus on graph-based pangenomes constructed from whole-genome alignments of assembled sequences, that is, variation graphs¹⁵. Variation graphs more comprehensively represent eukaryotic genomes by storing complete genomic sequences and their

variation^{21,34–36}. Variation graphs compress redundant sequences into bidirected networks where each node represents a sequence and edges connect nodes into complete sequences¹⁵ (Fig. 1c and the Box 1). Linear genomes or phased haplotypes are stored as explicit paths through the graph¹⁵, and sequence variation is represented by subgraphs called bubbles or snarls, in which variants are defined by alternative paths connected by shared start and end nodes (Box 1)³⁷. Here, we refer to variation graphs from whole-genome alignments simply as ‘pangenome graphs’. They were adopted by the Human Pangenome Reference Consortium (HPRC)³⁸, a global initiative that generated the first human pangenome^{36,39} (Box 2). Pangenome graphs were also assembled for the chicken (*Gallus gallus*)²¹ and, among non-model organisms, for the barn swallow (*Hirundo rustica*; Box 3)²⁸ and the house finch (*Haemorrhous mexicanus*)²⁹. Super-pangenome graphs are being generated for economically important species, such as tomato (*Solanum lycopersicum*)²⁵, grape (*Vitis* sp.; Box 3)²⁴ and cattle (*Bos taurus*)²⁶.

Maximizing capture of genetic diversity via sampling and sequencing

The sampling strategy is critical for successful biodiversity pangenomic studies. It should maximize genomic and biogeographic diversity

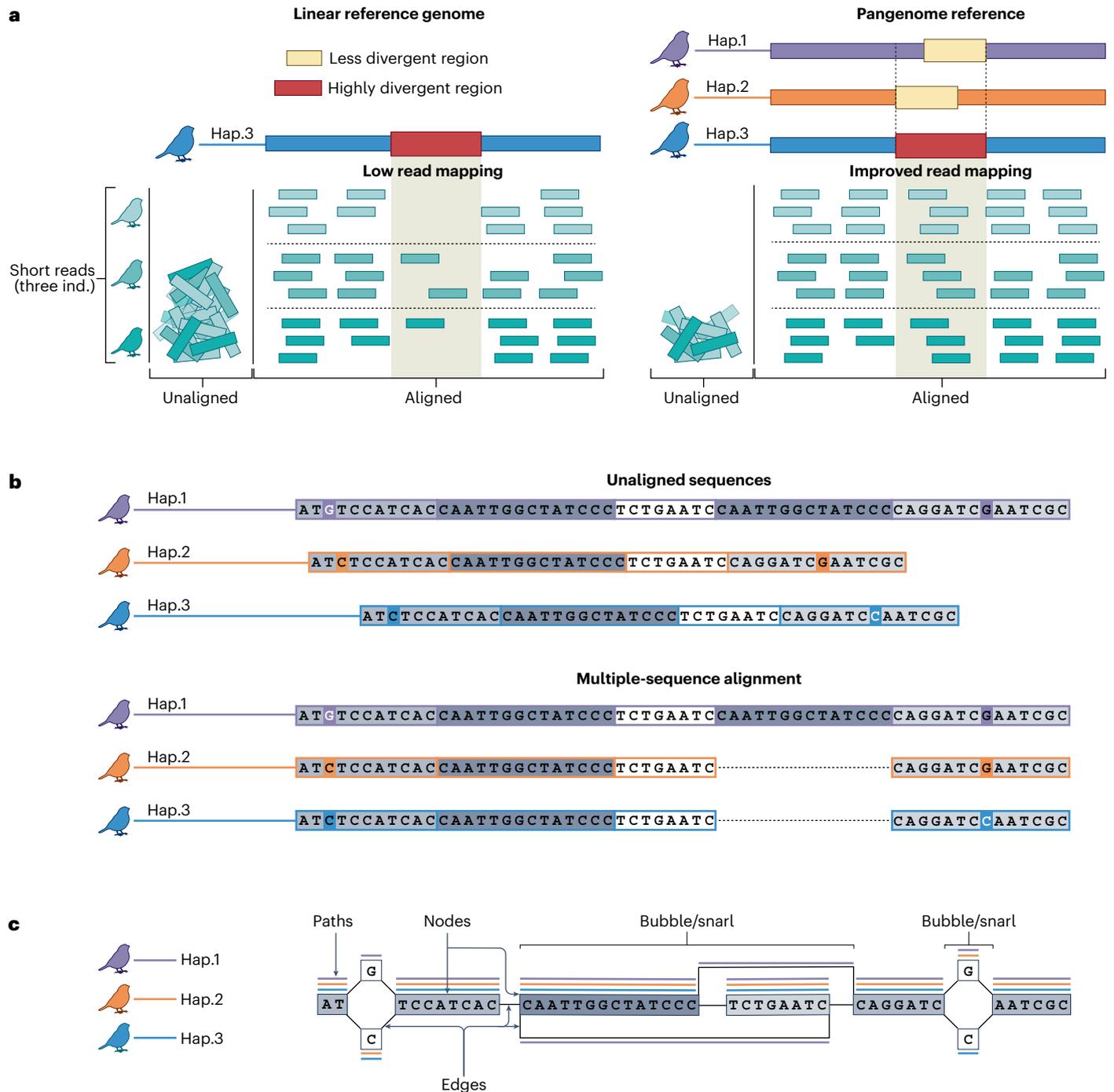


Fig. 1 | Principles of pangenome graphs. a, An example of reference bias. Short reads from three different individuals (ind.; shown here is a bird as an example) aligned to a linear reference genome do not map well to a missing region or a region with divergence (shown in red). By contrast, a pangenome reference based on multiple genomes (Hap.1, Hap.2, Hap.3; only one haplotype for each individual is represented for simplicity) improves the coverage of such regions, as less diverging regions are sampled, thereby facilitating variant calling and subsequent downstream analyses. **b**, Pangenome graphs can be constructed from unaligned raw reads (top) or from the alignment of multiple assembled sequences (bottom). **c**, A pangenome graph is a bidirected graph (Box 1) with

nodes representing DNA sequences (semi-transparent squares with base pairs) connected by bidirectional edges, which define the relationships between adjacent nodes and encode the strandedness of the sequences. The genomes walk as paths through the nodes (fine colored lines above each node), defining their base composition. Multiple genomes can share the same node sequence or take different paths across bubbles or snarls, which are subgraphs indicating the presence of variation in that region. The first bubble represents an insertion in Hap.1 or a deletion in Hap.2 and Hap.3. The second bubble represents a SNP in which Hap.1 and Hap.2 both have a G and Hap.3 has a C.

within natural populations, ideally sampling through the entire geographic range while balancing sex representation^{18,24,25,36,38,40–43} (Fig. 2a). If a panel of variants, both SNPs and SVs, is available, estimates of heterozygosity, relatedness and inbreeding offer insights into the sample

size required to achieve a comprehensive representation. Ordination and clustering analyses can help select representative individuals for inclusion in a pangenome²⁴. The ideal sample size can also be retrospectively verified by a pangenome number analysis, in which a number

BOX 2

History of the pangenome concept from bacteria to the human pangenome

The origin of the pangenome concept traces back to 2005 (ref. 122), when the bacterial genomes of *Streptococcus agalactiae* were first described as collections of core genes shared among strains, dispensable ('accessory') genes shared between some strains and strain-specific (unique) genes. In *S. agalactiae*, the core genome included 80% of the genes, with the remaining 20% categorized as accessory¹²². Work has since focused on the structure and dynamics of bacterial pangenomes^{123–127}, proving fruitful for taxonomic identification¹²⁸, to study host–pathogen interactions^{124,129} and gene families essential to pathogenicity¹²⁸ and antibiotic resistance¹³⁰. These studies resulted in biomedical applications, with new promising candidate drug targets¹³⁰ and reverse vaccinology^{131,132}. The concept was rapidly adopted by plant and animal researchers, resulting in multiple eukaryotic pangenome studies⁴⁵. To accommodate the complexity of eukaryotic genomes, with large portions (>50%¹³³) of noncoding and yet functional sequences, the term evolved to represent a complete set of sequences found in all individuals of a population or a species⁴⁵. In plants, the pangenome concept was first applied to the analysis of transposable elements, responsible for a large amount of variation in both genic and intergenic sequences¹³⁴. Plant genomes are particularly dynamic, as they undergo frequent polyploidization and diploidization events^{135,136}. Moreover, intraspecific genetic variability is often large¹³⁷. For these reasons, the concept of the pangenome has been rapidly extended from initial efforts in crops, such as rice¹³⁸ and tomato¹³⁹, to many other plant species^{140–142} for association mapping analyses^{143,144}, breeding¹⁴⁵, crop improvement¹⁴⁶ and evolutionary analyses¹⁴⁷. Recently, a pangenome comprising 69 *Arabidopsis thaliana* chromosome-level assemblies from a global range distribution was published¹⁸. Animal pangenomes have hitherto mostly focused on model species of economical value, both vertebrate and invertebrate. Among mammals, pangenome assembly and analysis projects have been conducted in the domestic

pig (*Sus scrofa*)^{22,148,149}, cattle (*B. taurus*)^{150,151} and domestic sheep (*Ovis aries*)¹⁵². Among invertebrates, studies on silkworm (*Bombyx mori*)⁴⁴, Mediterranean mussel (*Mytilus galloprovincialis*)^{171,153} and longwing butterfly (three *Heliconius* species)¹⁵⁴ pangenomes were recently published. The first step toward building a human pangenome was taken with the assembly of an Asian and an African genome and their integration within the NCBI reference human genome¹⁵⁵. This work resulted in the identification of approximately 5 Mb of novel sequences. From these early analyses, the authors estimated that a comprehensive human pangenome should contain 19–40 Mb currently missing from the state-of-art reference assembly (NCBI build 36.3)¹⁵⁵. Since then, a Danish pangenome¹⁵⁶ has been assembled and the GenomeDenmark project was initiated¹⁵⁷, both using family trios and revealing novel single-nucleotide variants and SVs. More recently, the assembly of an African pangenome¹⁵⁸ included approximately 10% more DNA (296 Mb) than the reference assembly (GRCh38), and a pangenome built from hundreds of Han Chinese individuals¹⁵⁹ identified 29.5 Mb of novel genomic sequences, including at least 188 novel protein-coding genes. The generation of a comprehensive catalog of human genetic variation has advanced through the analysis of 338 high-quality assemblies of genetically divergent populations¹⁴², demonstrating that, for a given genome sequenced to 40-fold coverage, over 400,000 previously unmapped reads could now be aligned. The authors also tested this new resource for more efficient mapping of previously discarded RNA-seq reads¹⁶⁰. Recently, the HPRC released the first draft human pangenome graph using three graph construction methods and 47 phased, diploid assemblies from genetically diverse individuals³⁶. Soon after, another human pangenome graph was built using a similar pipeline from 116 high-quality assemblies representing 36 Chinese minority ethnic groups¹⁶¹. In parallel, Seven Bridges Genomics generated population-specific genome graphs for an African pangenome¹⁶².

of representative genomes are progressively added until a pangenome plateau is reached, that is, when the full set of genes is captured and adding more individuals does not recover novel genes⁴⁴. The inclusion of more individuals augments the existing references by increasing the representation of accessory or dispensable sequences, which are shared only by a subset of individuals and often have functional and adaptive roles⁴⁵. Pangenome graphs facilitate pinpointing the functional and adaptive roles of accessory genomic regions and how they vary among geographic populations and subspecies. Accessory genomes may provide hotspots for population differentiation and speciation through divergent selection processes or via hybridization^{46,47} (Box 1). Investigating how much accessory regions are affected by introgression (Box 1) in populations and species undergoing hybridization provides insights into the dynamics of gene flow and speciation processes⁴⁸.

Priority should be given to the highest-quality samples, maximizing long-read sequencing throughput (Fig. 2a) and allowing accurate and contiguous genome assemblies^{5,49}. The quality of the input genomes minimizes noise propagation in the pangenome graph. Haplotype phasing with parental sequence data or chromatin conformation data (Hi-C; Box 1) is crucial to prevent haplotype false duplication (Box 1) and related errors^{4,5,50,51}. Sequencing coverages

of ~30-fold PacBio high-fidelity long reads (HiFi)^{11,52} and ~60-fold Oxford Nanopore Technologies (ONT)⁵³ duplex reads, in combination with ~30-fold Hi-C per haplotype and manual curation, generate reference genomes that meet the current quality standards^{5,7,54}. Genome sequencing is a rapidly evolving field, and generating complete, haplotype-resolved and near error-free genomes (telomere to telomere, T2T) is now feasible^{36,49,55} by complementing HiFi with ultra-long ONT reads. Incorporating high-quality genomes aids in the discovery of rare SVs⁵⁶, particularly in admixed populations and those with large effective population sizes^{57,58}. Moreover, it improves the representation of hard-to-sequence-and-assemble regions like centromeres, variable number tandem repeats³⁶ and other complex repeats. Examining base-level polymorphism in variable number tandem repeats may clarify their role in shaping gene expression and complex traits^{10,26,59}. Highly repetitive regions might also underlie the regulation of complex behavioral phenotypes, such as migratory behavior⁶⁰.

Overall, a pangenome graph benefits from the inclusion of all T2T-level or high-quality genomes, whereas pangenomes derived from sub-T2T genomes will limit the study of genetic diversity due to the incompleteness of challenging regions. However, acquiring multiple high-quality samples from non-model species may be difficult,

BOX 3**Case studies****The barn swallow pangenome**

The barn swallow (*H. rustica*) is a small migratory songbird, with six subspecies that differ in body size, extent and type of secondary sexual traits, and migratory behavior¹⁶³. Latitudinal clines exist, with partial overlap in some traits¹⁶⁴ and variable levels of hybridization between subspecies¹⁶⁵. To further characterize genetic variation within this species, a preliminary pangenome variation graph for the Eurasian subspecies was generated from 12 haplotypes combined with MC²⁸. The resulting pangenome increased the reference genome by approximately 500Mb (1.6Gb versus 1.1Gb). The pangenome graph enabled to infer core and accessory genes and improve read mapping and variant calling, further highlighting the potential of pangenome graphs for population genomics²⁸.

**The grape super-pangenome**

The grapevine (*Vitis vinifera*) is one of the most important fruit crops globally, with great economic and cultural value. Recently, the assembly of nine genomes from North American wild *Vitis* species aimed to comprehensively characterize genus diversity²⁴. The genomes were scaffolded at the chromosome level and fully phased. A super-pangenome was built using PGGB⁶³, enabling access to intraspecific and interspecific genetic variants and augmenting the genome size threefold (1.7Gb versus 0.5Mb). The decomposition of variants embedded in the graph captured valuable genetic variants, including those associated with flower sex phenotypes and disease resistance, thereby shedding light on species-specific adaptations. The super-pangenome also supported a pangenome-wide association study and identified variants near gene loci that are associated with chloride exclusion, potentially influencing plant salt tolerance. These findings highlight the potential of pangenomes in studying genetic variation and uncovering the genetic basis of functional traits as well as helping to address the challenges posed by climate change.



especially for rare and threatened species, and sequencing multiple individuals with different long-read technologies may currently be cost-prohibitive. We suggest that a pangenome graph should include at least one high-quality assembly as a backbone for graph construction⁵⁴, providing a robust coordinate system for downstream analyses.

Building pangenome graphs to represent complex and accessory genomic regions

Pangenome graph construction starts with the alignment of input genomes to identify sequence similarities. Alignment can be reference based^{20,61} or involve all-versus-all comparisons^{62,63}, and it can be either at

the base level^{20,62} or at a higher level (for example, including only variant sites)⁶¹. Alignment information is embedded in the graph, which can be manipulated for downstream analyses. Two main pipelines were developed by the HPRC for graph construction: Minigraph-Cactus (MC)²⁰ and the PanGenome Graph Builder (PGGB)⁶² (Fig. 2b and Supplementary Table 1). MC implements minigraph⁶¹, a sequence-to-graph aligner, as a graph constructor. In MC, a user-selected reference genome is used as the initial backbone, which is progressively augmented with structural variation from the other genomes. The resulting graph is SV only (>50 bp; Fig. 2b), and all assemblies are aligned back to the graph with a minimap2-like⁶⁴ algorithm that generates base-level alignments for each reference chromosome. MC implements a modified version of the reference-free aligner Progressive Cactus⁶⁵ to combine the alignments into base-level pangenome graphs that contain variants of all sizes (Fig. 2b). Chromosomal graphs are then combined and post-processed to reduce path complexity by collapsing redundant sequences²⁰. In addition to the chosen reference, one may specify additional assemblies with coordinates that can serve as a reference in downstream analyses²⁰. In contrast to MC, PGGB⁶² avoids using an initial reference and rather employs all-to-all genome alignments with wfasmh⁶⁶, a software for homology mapping that generates base-wise pairwise alignments. PGGB uses seqwish⁶³ as a sequence-to-graph aligner, which starts from the all-versus-all alignments to generate a complete pangenome graph, representing all variant types and sizes (Fig. 2b). The graph is then post-processed with a smoothing and normalization step^{36,62}. PGGB can be run in parallel on each chromosome community, that is, each set of sequences corresponding to each reference chromosome, to reduce computation time^{36,62}. In PGGB graphs, every genome included in the graph can be used as reference for downstream analyses⁶². However, like in MC, only one reference can be used at a time as a coordinate system for downstream analyses such as variant calling. New computational methods and file formats other than the linear binary alignment map (BAM) and variant call format (VCF) need to be developed to overcome this limitation and represent all the information embedded in the graph. Pangenome graphs can be combined with transcript annotations using the variation graph toolkit (vg)⁶⁷, a software for variation graph construction, handling and analysis, into splice-aware graphs, with paths through nodes (exons) and edges (splice junctions) representing the structure of mRNA transcripts. It is also possible to build pantranscriptomes by projecting a set of haplotype-specific transcripts onto a set of known haplotypes⁶⁷.

The size of a pangenome graph depends on the genome size of the respective species but is bound to be larger, as it incorporates accessory sequences from other individuals, and it is also influenced by the number and diversity of the individuals contributing to the pangenome as well as by the construction pipeline (Supplementary Table 1). The size of the MC graph is relatively close to the genome size of the species (-3.2 Gb versus 3.1 Gb for humans³⁶, 1.2 Gb versus 1.1 Gb for chickens²¹,

-1.6 Gb versus 1.1 Gb for barn swallows²⁸; Supplementary Table 1 and Box 3). By contrast, PGGB graph size can considerably exceed that of the genome size and MC graphs (for example, 8.4 Gb for humans³⁶). The larger size of PGGB graphs is explained by their capability to capture highly divergent satellite, centromeric, and heterochromatic regions, excluded in MC graphs due to alignment issues^{21,36} (Supplementary Table 1). The largest increase relative to true genome size was observed in grapevine²⁴, likely because of the inclusion of different species (Box 3). PGGB also has a tendency to collapse complex regions, such as copy number polymorphic loci, into a single copy, generating loops in the graph that increase its complexity^{36,62} (Fig. 2b). Given their greater size and complexity, PGGB graphs require more computational resources than MC graphs. To construct a graph based on ten human haplotypes, MC currently takes ~16 h, 154 GB of RAM and 7 GB of disk space, while PGGB takes ~117 h, 71 GB of RAM and 7.6 GB of disk space⁶⁸. A Nextflow implementation of the PGGB pipeline has recently been released to improve cluster scalability⁶⁹. PGGB has also been experimentally shown to potentially lead to an overestimation of sequence variability. For instance, the size of the PGGB chicken pangenome was larger than expected based on the estimated variation in diverse groups of chickens²¹. A preliminary estimation of the expected species genetic diversity should be computed to detect overestimation in graph construction.

Given all these differences, a careful selection of the pangenome graph construction pipeline is of utmost importance. On one hand, PGGB graphs are based on reference-free alignments and are more complete than MC graphs. On the other hand, their complexity increases the computational resources needed for graph construction and some downstream analyses, such as variant calling after read mapping, are currently computationally infeasible^{21,36}. MC graphs are easier to handle, but they omit challenging regions such as centromeres, which are hotspots of structural variation¹⁸ and may play a crucial role in adaptive evolution and speciation^{70,71}. Moreover, MC works on single-chromosome graphs during graph construction and therefore does not allow the representation of interchromosomal rearrangements³⁶, precluding, for example, the investigation of acrocentric chromosome evolution⁷². We suggest choosing the pipeline based on the desired analyses and the available computational resources. Both graphs can be used for graph decomposition and the identification of variants between the genomes in the graph; however, while MC should be used as a reference for resequencing projects, PGGB is particularly useful when focusing on regions of interest. Overall, pangenome graphs face conceptual and computational challenges and currently require substantially more resources for their construction, storage and analysis than linear genomes. While these limitations are being tackled (for example, through ‘implicit’ pangenome construction of only regions of interest⁷³), researchers need to ensure access to sufficient computing resources, such as clusters or cloud computing infrastructure.

Fig. 2 | Pangenome graph construction, manipulation and visualization.

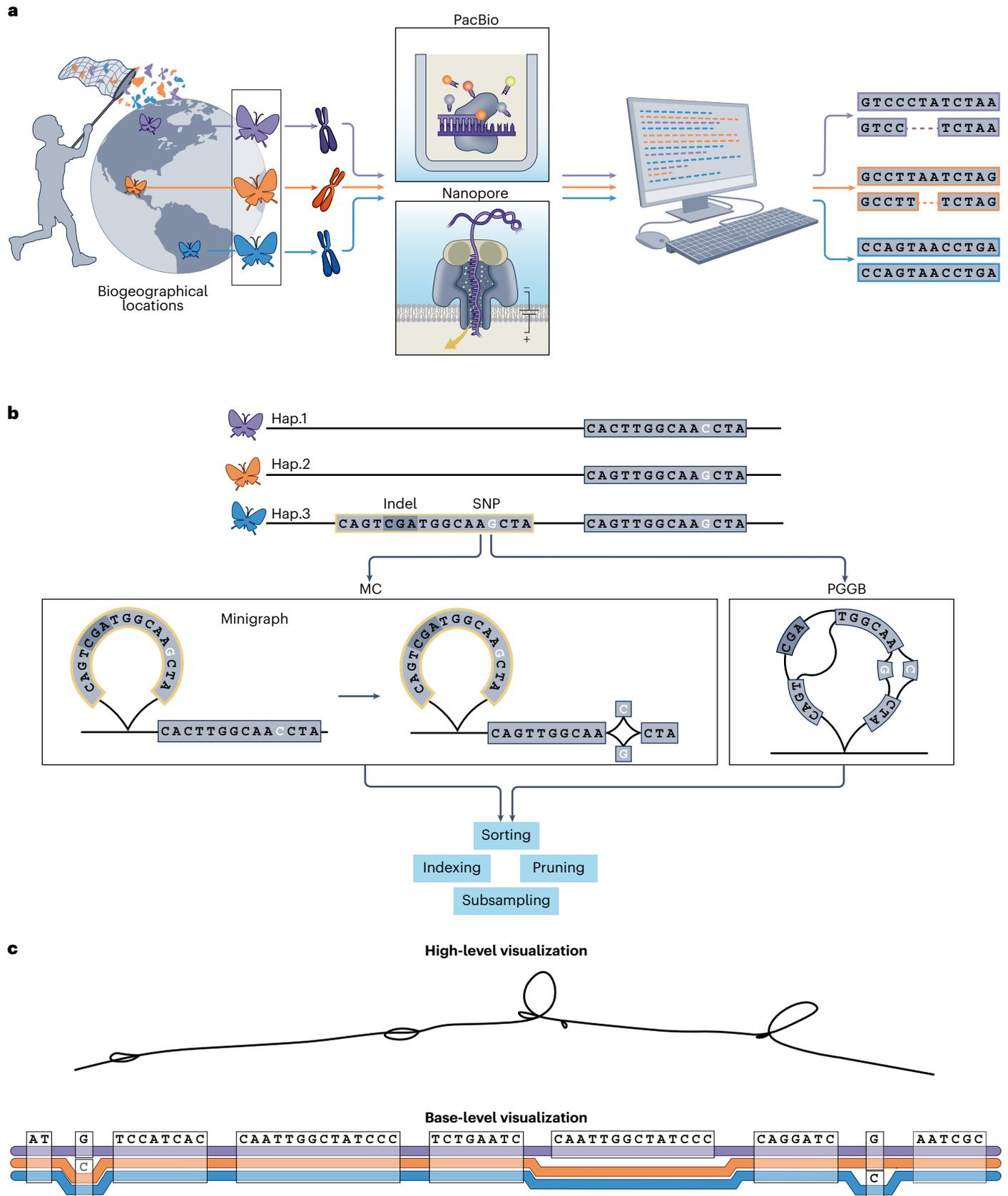
a, General representation of sample selection, sequencing and assembly of haplotype-resolved genomes (illustrated here are butterflies as an example). Samples should be collected in different geographical locations to increase variability. Next, high-quality DNA is extracted, sequenced with long-read technologies such as PacBio and Oxford Nanopore and assembled to obtain high-quality genomes. **b**, Schematic illustration of the two main methods for pangenome graph construction, reference based (MC) and reference free (PGGB). MC starts with minigraph, which generates an SV-only graph, that is, it starts with Hap.1 and adds the SVs from Hap.3 (>50 bp, highlighted in yellow). The graph construction is therefore influenced by the order of the genomes aligned. Next, Cactus adds base-level information of all sequences to the graph to also represent any SNPs. PGGB does not depend on the order of the aligned genomes and starts from all-versus-all alignments, generating a graph with loops representing complex regions of the genome, such as the centromeres. After

graph construction, typically performed operations include sorting, indexing, pruning and subsampling to correct the order of the nodes, create an index to make the pangenome elements accessible to other software, prune complex and unreliable regions and focus on regions of interest, respectively. **c**, Pangenome graphs can be visualized in many ways to help interpretation. It is possible to visualize the overall structure of the graph in 2D (high-level visualization), focusing on the relationships between nodes and edges rather than the base composition of the paths. Alternatively, it is possible to visualize the paths walking through the nodes together with their base composition in 1D (base-level visualization). The latter is a tube-like representation of the graph shown in Fig. 1c. Each colored line is a different genome (only one haplotype per individual is shown for simplicity). Nodes report the DNA sequence (semi-transparent squares), and variation is represented as divergence in the paths walking through the nodes.

Improving the accessibility of biological information in pangenome graphs

Pangenome graphs enclose extensive and complex biological information, including genomic relationships and diversity among individuals.

Their intricate and complex structure generates large data volumes that are challenging to navigate and interpret. Graph manipulation toolkits have been developed to improve graph accessibility to software used in downstream analyses, thereby facilitating the extraction of



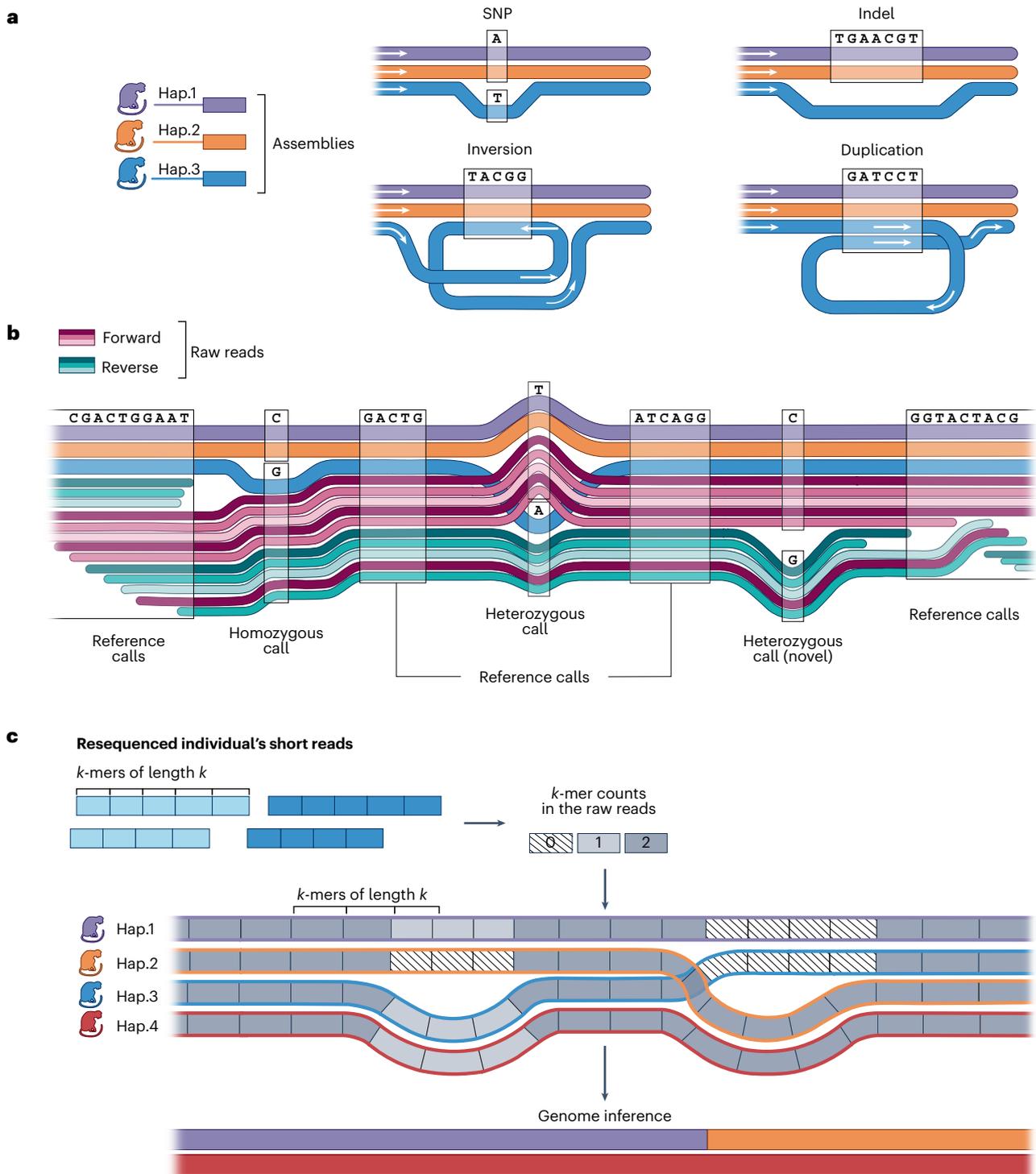


Fig. 3 | Downstream analyses. **a**, Identification of variants through graph decomposition. Sites of variation (bubbles or snarls) are identified by decomposing the pangenome graph based on its topology, that is, following genome paths across nodes and edges (illustrated here with monkeys as an example). Four bubbles are shown in a tube-like manner, as in Fig. 2c (base-level visualization). Walking through the nodes allows for the identification of a SNP, an indel, an inversion, and a duplication in Hap.3. Variants are called from the bubbles based on path divergences in relation to a selected reference genome. **b**, Use of pangenome graphs for read mapping and variant calling. Mapped reads (thin lines) after graph augmentation are visualized alongside the pangenome graph reference (thick lines). This allows us to call variants already present in the graph (for example, homozygous call in which the mapped reads share a G with Hap.3) or novel variants (for example, second heterozygous call, in which the read has a C with respect to all the haplotypes in the pangenome

graph). Variants are called by analyzing divergences between read paths and the reference genome. **c**, SV detection and sample genotyping using PanGenie, which avoids read mapping. Raw reads from a resequenced individual (blue) are divided into k -mers and assigned to nodes within the pangenome graph (thick lines with different-colored outlines). k -mers across the haplotypes are displayed in different tonalities of gray, according to their multiplicity in the raw reads. Genotypes are determined by comparing k -mer counts between paths at bubbles, which represent known haplotypes. In the first bubble, the absence of k -mers matching Hap.2 and the presence of one-copy k -mers matching Hap.1, Hap.3 and Hap.4 result in a Hap.1/Hap.4 (heterozygous) genotype. In the second bubble, two-copy k -mers matching Hap.2 and Hap.4 are observed, inferring a Hap.2/Hap.4 genotype. In the first bubble, Hap.4 is selected over Hap.3 based on neighboring bubbles, where Hap.3 lacks k -mer support.

biological information from pangenome graphs. Graph manipulation includes tasks such as sorting, indexing, pruning and subsampling (Fig. 2b). Sorting optimizes the order of graph nodes to reveal the underlying latent and sparse graph structure and minimize errors in analysis and interpretation⁷⁴. Path indexing provides faster access to specific regions of the graph, allowing software to quickly locate genes, variants and other features of interest without scanning the entire pangenome and reducing the time required to retrieve relevant information for analyses such as read mapping and variant calling⁷⁵. To further speed up computation, it is possible to simplify graph topology through pruning of complex or unreliable regions or by subsampling user-defined coordinates^{74,76}. Subsampling is particularly helpful to disentangle the complexity of a particular region of interest or when computational resources are not available to query the entire graph. After graph construction and each manipulation step, it is helpful to perform diagnostic statistics, such as on graph size, number of nodes and base content, to get a sense of the structure of the pangenome and how each step affected the graph^{74,76}. However, being a relatively new approach, pangenome graphs lack universally accepted quality metrics due to their intrinsic complexity, the multiple construction methods, and the absence of standardized benchmarking datasets. Nonetheless, as pangenome graphs become more widely adopted, efforts will develop a unified set of metrics applicable to all pangenome graphs.

Two main software packages exist to manipulate pangenome graphs: *vg*⁷⁶ and the pangenome analysis toolkit ODGI⁷⁴. *vg* relies on the *vg* format⁷⁶ and was the first tool to be scaled up to gigabase-scale graphs. ODGI operates on a node-centric object (.og) and was optimized for pangenome graphs with hundreds of haplotype-resolved genomes⁷⁴. ODGI's tools work on a graph-independent universal coordinate system that remains constant among different graphs built from the same sequences⁷⁴. This system enables coordinate translation, facilitating the liftover of coordinate-based features between genomes and graphs⁷⁴, that is, the accurate mapping of annotated features (such as genes, regulatory elements or other functional elements) from one genome assembly to another⁷⁴. The Comparative Annotation Toolkit (CAT)⁷⁷ can also annotate the haplotypes in a pangenome graph by projecting the reference gene annotation to each of the genomes, which can ease within-species annotation efforts³⁶. Feature annotations can also be injected in the graph^{74,76} and used to interpret the functional importance of paths, nodes, and edges.

Visualizing genome diversity among individuals

Graph visualization allows the inspection of homology relationships and variation between the genomes, providing insights on the latent biological data⁷⁴. For instance, it can disentangle variation at complex loci following haplotype paths along variation bubbles²¹. Visualization can occur at different scales, from the overall structure down to the base level (Fig. 2c). 2D visualizations highlight graph structure and identify complex loci, while 1D visualizations help understanding the graph topology and the relationships between genomes, potentially providing a more immediate understanding of the complexity of a region with respect to inspecting a list of variants. Various tools exist for pangenome graph visualization. *Bandage*⁷⁸ and *GfaViz*⁷⁹, originally created to visualize assembly graphs, can generate 2D graph layouts, allowing the interactive inspection of nodes and edges, with variation that appear as bubbles in the layout. *vg viz*⁷⁶ can visualize nodes, edges, paths and the base variation among sequences. *SequenceTubeMap*⁸⁰ renders these elements in a 1D 'tube map' model where paths representing genomes navigate through the sequence nodes of the graph, oriented from left to right. Read alignments and feature annotations injected in the graph can also be visualized. To scale to gigabase pangenomes, such as the HPRC graph³⁶, *MoMI-G*⁸¹ combines the base-level visualization of *SequenceTubeMap*⁸⁰ with *Circos*⁸² plot chromosome-level connections to efficiently browse SVs between genomes and aligned reads. ODGI⁷⁴ can render a raster image of the graph topology in either

2D or 1D. *Waragraph*⁸³, an interactive implementation of ODGI, is currently being developed to be able to inspect both 1D and 2D visualizations. When dealing with a large graph, rendering the entire graph at once can become impractical, and we recommend visualizing chromosome graphs or subsampled regions of interest.

Downstream analyses and their applications

Characterizing small variants and complex SVs through graph decomposition

Variant sites (SNPs, insertion–deletions (indels) and SVs; Fig. 3a) in a pangenome graph can be extracted through graph decomposition^{37,29}, the process of breaking down a pangenome graph into smaller, more manageable subgraphs or components (snarls or bubbles)³⁷. Graph decomposition can be performed with *vg snarl*³⁷ and *gfatools bubble*⁸⁴ (Supplementary Table 2). *vg deconstruct*³⁶, implemented in the MC²⁰ and PGGB pipelines⁶³, can process the output of *vg snarls* or compute snarls automatically, generating a VCF with variants called from the references chosen during graph construction with MC²⁰ or from any genome with PGGB⁶³. When working with large graphs, it is recommended to compute snarls separately before variant calling³⁶. The characterization of complex SVs, which were not previously accessible using linear reference genomes and short reads, can shed light on their role in evolution⁸⁵ and in shaping phenotypic variation, often with fitness consequences^{86–89}, as SVs can affect fitness by altering gene expression and shaping the chromosome recombination landscape⁹⁰. A complete representation of SVs can also help analyze synteny and collinearity (Box 1) within genomes. In turn, this may provide insights into chromosome evolution by encompassing the full complexity of sex chromosomes and microchromosomes (Box 1), which are typically enriched in SVs and challenging to resolve due to high repeat and GC content⁹¹. In addition, human pangenome graphs have also allowed the identification of recombination events between heterologous acrocentric chromosomes, especially at the breakpoints of Robertsonian translocations⁷². These translocations are the most common chromosomal rearrangement in humans, and a comprehensive pangenome graph has greatly enhanced the identification of the sequences and mechanisms involved⁷².

Population genomics and alignment of transcriptomics data

A pangenome graph can be used as a reference in resequencing projects to reduce mapping bias (Fig. 3b and Supplementary Table 2). Short reads map with greater confidence when more genomic sequences are represented and known variation is embedded in the reference¹⁵. However, read mapping to a pangenome is more challenging than mapping to a single reference genome, as the search space for alignment increases due to the large number of potential paths in the graph⁹². Since canonical algorithms cannot be applied directly to pangenome graphs¹⁵, new tools have been developed for sequence-to-graph alignment. Within the *vg* toolkit, the general-purpose read mapper *vg map*⁷⁶ is suitable for large and complex variation graphs, albeit slower than popular linear genome aligners with comparable accuracy⁹². *vg giraffe*⁹², currently being extended to support long reads, uses a graph Burrows–Wheeler transform⁹³, an indexing strategy that supports efficient querying and retrieval of sequences and variants from the pangenome graph, to identify the paths that represent the two observed haplotypes in an individual's reference sequences and to restrict the alignment space to these regions only, avoiding biologically unlikely allele combinations. This leads to a dramatic increase in mapping speed⁹². Long reads can also be aligned with *GraphAligner*⁹⁴, a seed-and-extend sequence-to-graph aligner (Box 1). Improved short-read mappability will benefit resequencing projects, particularly in ancient DNA studies, which face challenges of contamination, degradation, small amounts of endogenous DNA, shorter reads and, therefore, lower mappability¹⁴. Ancient DNA mapped against a variation graph has already been proven to mitigate reference biases by improving the allelic balance

in polymorphic sites¹⁴. Ancient DNA reads with non-reference alleles map in higher proportions to a graph containing alternate alleles, with respect to a linear reference genome¹⁴.

vg also allows splice-aware RNA sequencing (RNA-seq) mapping to splice-aware graphs, generating an alignment that can then be used to quantify haplotype-specific transcript expression^{67,95}. Pantranscriptomics has the potential to efficiently quantify haplotype-specific differential gene expression by exploiting the population variation that is embedded in the pantranscriptome reference⁶⁷. We anticipate that pantranscriptomic sequencing projects combining RNA-seq data with pangenome graph references will clarify the effects of gene flow⁹⁶, detecting adaptive genetic variation^{97,98}. Chromatin accessibility analyses, such as chromatin immunoprecipitation followed by high-throughput sequencing (ChIP-seq; Box 1) or assays for transposase-accessible chromatin with sequencing (ATAC-seq; Box 1), also benefit from a pangenomic approach³⁶. Their combination with RNA-seq data provides a multiomic approach that might facilitate the interpretation of regulatory events critical to a wide variety of biological processes and phenotypes^{36,99}. These approaches will enable future panepigenomic studies in non-model organisms, overcoming current limitations in handling large, multiomics datasets^{36,99}.

Detecting and genotyping variants in resequencing studies

Pangenome graphs can increase the accuracy of variant calling and genotyping in resequencing studies thanks to improved read mappability^{21,26,28,36} (Box 3, Fig. 3b and Supplementary Table 2). vg can be used to extract snarls from the graph and compute coverage and mapping quality of aligned reads to accurately identify known variants¹⁰⁰. Larger known SVs (deletions, insertions and inversions) can be genotyped by computing read coverage for each node and edge¹⁰⁰. Specifically, variable sites identified by graph decomposition are assigned the two most supported paths, representing haplotypes¹⁰⁰. For de novo small variant calling, the graph is first 'augmented' with variants identified through read alignment, after which read support is computed¹⁰⁰. A potentially more accurate approach consists of projecting the graph alignment from short-read mapping back into a linear BAM file (that is, surjection) before using traditional variant callers (for example, DeepVariant¹⁰¹, FreeBayes¹⁰² and GATK with elPrep^{103,104}) to generate a VCF referenced to the genome chosen for surjection. This approach can also be used with long reads mapped by GraphAligner⁹⁴. Given the complexity of PGGB graphs, read mapping and variant calling have so far been mostly performed on MC graphs^{21,36}. vg giraffe mapping followed by surjection and variant calling with DeepVariant is currently the state-of-the-art approach, and surjection was found as the main computational bottleneck^{21,36}. For the chicken pangenome, the number of mapped and surjected reads per CPU second dramatically decreased when mapping against a PGGB graph compared to an MC graph (1.6 reads versus 500 reads), while the memory usage increased (>250 GB versus 24 GB), making mapping and surjection with a PGGB graph currently computationally infeasible²¹. An alternative and faster approach for known variant genotyping that does not require read mapping is implemented in PanGenie¹⁰⁵. This algorithm combines long-range haplotype information embedded in the graph and *k*-mer (Box 1) counts from short-read data to jointly genotype SNPs, indels and SVs in the uncharacterized sample (Fig. 3c). Haplotypes present in the graph can support genotype inference based on neighboring bubbles in case a given bubble is poorly covered by short-read *k*-mers¹⁰⁵.

In population studies, pangenome-based variant calling increases accuracy and reduces per-sample data requirements, potentially expanding the size of assessable cohorts³⁶. Determining accurate and comprehensive variant sets increases the resolution of the analysis of demographic history, linkage disequilibrium and genome-wide selection scans. This is particularly beneficial in species with large effective population sizes, in which linkage disequilibrium is low^{28,106,107}. By improving SV genotyping, pangenome graphs can also help to

integrate SVs into genome-wide association studies, especially as long reads gradually replace short reads in resequencing projects^{24,36,105}. Performing genome-wide association studies on pangenome-based SNP and SV panels can therefore enhance our understanding of the genetic basis of complex polygenic traits and shed light on the role of natural selection and gene–environment interactions and correlations.

Conclusions and future prospects

The last few years have seen a burgeoning of both small- and large-scale projects generating high-quality reference genomes for biodiversity studies¹⁰⁸, including the Vertebrate Genomes Project^{5,7}, the Darwin Tree of Life¹⁰⁹ and the European Reference Genome Atlas¹¹⁰. Most of these projects contribute to the Earth BioGenome Project⁶, an ambitious proposal launched in 2020 to collectively sequence all named eukaryotic species within the next 10 years. While pangenome graphs are currently available only for a handful of species, recent advances in genome and pangenome assembly potentially extend this approach to most eukaryotic species. This is a desirable goal to reduce representation bias in all analyses of biodiversity, its evolution and conservation. Collecting, sequencing and assembling pangenomes from more than a few individuals could be impractical in many species due to costs and sample availability. In those cases, a pangenome from a single to few individual(s) would still increase representation and reduce reference bias, especially for highly heterozygous populations in which a single individual may carry a high amount of allelic diversity. Pangenome graphs can benefit a broad range of applications for biodiversity, from population genomics, phylogenomic, hybridization and speciation studies to conservation genomics, and will likely become the standard reference system for such research in the future. Many new directions are being investigated. For instance, panmitogenomes, that is, pangenomes constructed from thousands of mitochondrial genomes, have been shown to improve haplotyping of individuals¹¹¹ and are being considered for species identification from heterogeneous samples. Another promising new direction for the field is in super-pangenome graphs, which expand the survey of variation to taxonomic ranks above species, opening new possibilities to study the molecular and evolutionary mechanisms underlying species divergence, selection and recombination processes as well as adaptation to rapid climate changes²⁷. Dense sampling and sequencing of species within a clade have proven essential for deciphering phylogenetic and phylogeographic relationships as well as facilitating investigation of gene loss and selection events¹¹². To this end, a pangenomic approach revealed complex phylogenetic relationships among bacterial strains, allowing genetic analyses of infectious diseases to identify virulence and antimicrobial resistance genes with greater accuracy¹¹³. Despite the complexity of eukaryotic genomes, we envision that rapid improvements in the efficiency and scalability of pangenomic tools will soon allow such phylogenomic applications to be extended to eukaryotic species. In particular, owing to the ability of super-pangenome graphs to incorporate all types of genomic variation, they have the potential to elucidate complex evolutionary histories and phylogeographic relationships of large, panmictic and highly recombinant wild populations^{28,114} as well as to improve phylogenetic reconstructions of events, such as incomplete lineage sorting. Super-pangenomes can also assist in studying biodiversity in complex ecosystems where hybridization occurs^{48,115}. Hybrid zones are a prime opportunity for pinpointing the genes responsible for phenotypic traits, as genes introgress in parallel with those traits¹¹⁶. Inclusion of both hybridizing species in a pangenome graph will mitigate biases that arise from the use of the reference genome of either species¹¹⁷. Pangenomes could also help shed light on the origin of islands of divergence, highly differentiated genomic regions that might be related to reproductive isolation and, thus, to speciation processes^{118,119}. Even within the same species, a comprehensive pangenome graph that includes assemblies for all subspecies can maximize the identification of SVs that are unique to a

subspecies (Box 3). We predict that species-level pangenomes will also replace linear genomes in phylogenomic comparative genomic studies, enabled by the future development of tools for aligning pangenomes of different species. Currently, the construction of a pangenome graph lacks evolutionary information across individuals, as phylogenetic divergence is not considered in pairwise alignments, and this limitation should be taken into account when performing phylogenetic analyses.

Pangenome graphs may also effectively guide conservation strategies aimed at maximizing the preservation of genetic variation² by capturing a fuller spectrum of genetic diversity. Of particular interest is structural and functional genomic variation involved in adaptations and responses to environmental pressures. This will improve selection criteria for reintroducing and translocating individuals among populations of threatened and endangered species. Improved representation of structural elements, such as SVs, centromeres and telomeres, and copy number variations as well as SNPs, along with noncoding regulatory elements can provide comprehensive conservation-relevant information regarding inbreeding, outbreeding, deleterious mutations, introgression and local adaptation². Pangenomes can also help to identify different genomic regions in cryptic species to then develop multilocus probes that distinguish cryptic taxa and simplify conservation management¹²⁰. Moreover, we envision that pangenomics could help reconstruct the genomic blueprint of extinct biodiversity by improving the mappability of ancient DNA against a pangenome of a closely related species. A more comprehensive comparison between the extinct species and its living relative will help identify the genetic variation underlying lost traits and ecosystem functions, which are essential information for any de-extinction and restoration efforts^{14,121}. In conclusion, as methods to assemble, visualize, annotate and analyze pangenome graphs continue to improve, we recommend researchers in biodiversity genomics to embrace this new paradigm.

References

- Giani, A. M., Gallo, G. R., Gianfranceschi, L. & Formenti, G. Long walk to genomics: history and current approaches to genome sequencing and assembly. *Comput. Struct. Biotechnol. J.* **18**, 9–19 (2020).
- Theissinger, K. et al. How genomics can help biodiversity conservation. *Trends Genet.* **39**, 545–559 (2023).
- Lewin, H. A. et al. Earth BioGenome Project: sequencing life for the future of life. *Proc. Natl Acad. Sci. USA* **115**, 4325–4333 (2018).
- Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat. Methods* **18**, 170–175 (2021).
- Rhie, A. et al. Towards complete and error-free genome assemblies of all vertebrate species. *Nature* **592**, 737–746 (2021).
- Lewin, H. A. et al. The Earth BioGenome Project 2020: starting the clock. *Proc. Natl Acad. Sci. USA* **119**, e2115635118 (2022).
- Larivière, D. et al. Scalable, accessible and reproducible reference genome assembly and evaluation in Galaxy. *Nat. Biotechnol.* **42**, 367–370 (2024).
- Ghildiyal, K. et al. Genomic insights into the conservation of wild and domestic animal diversity: a review. *Gene* **886**, 147719 (2023).
- Blaxter, M. et al. Why sequence all eukaryotes? *Proc. Natl Acad. Sci. USA* **119**, e2115636118 (2022).
- Chaisson, M. J. P. et al. Multi-platform discovery of haplotype-resolved structural variation in human genomes. *Nat. Commun.* **10**, 1784 (2019).
- Wenger, A. M. et al. Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nat. Biotechnol.* **37**, 1155–1162 (2019).
- Zhao, X. et al. Expectations and blind spots for structural variation detection from long-read assemblies and short-read genome sequencing technologies. *Am. J. Hum. Genet.* **108**, 919–928 (2021).
- Günther, T. & Nettelblad, C. The presence and impact of reference bias on population genomic studies of prehistoric human populations. *PLoS Genet.* **15**, e1008302 (2019).
- Martiniano, R., Garrison, E., Jones, E. R., Manica, A. & Durbin, R. Removing reference bias and improving indel calling in ancient DNA data analysis by mapping to a sequence variation graph. *Genome Biol.* **21**, 250 (2020).
- Eizenga, J. M. et al. Pangenome graphs. *Annu. Rev. Genomics Hum. Genet.* **21**, 139–162 (2020).
- Taylor, D. J. et al. Beyond the Human Genome Project: the age of complete human genome sequences and pangenome references. *Annu. Rev. Genomics Hum. Genet.* **25**, 77–104 (2024).
- Gerdol, M. et al. Massive gene presence–absence variation shapes an open pan-genome in the Mediterranean mussel. *Genome Biol.* **21**, 275 (2020).
- Lian, Q. et al. A pan-genome of 69 *Arabidopsis thaliana* accessions reveals a conserved genome structure throughout the global species range. *Nat. Genet.* **56**, 982–991 (2024).
- Chin, C.-S. et al. A diploid assembly-based benchmark for variants in the major histocompatibility complex. *Nat. Commun.* **11**, 4794 (2020).
- Hickey, G. et al. Pangenome graph construction from genome alignments with Minigraph-Cactus. *Nat. Biotechnol.* **42**, 663–673 (2024).
- Rice, E. S. et al. A pangenome graph reference of 30 chicken genomes allows genotyping of large and complex structural variants. *BMC Biol.* **21**, 267 (2023).
- Jiang, Y.-F. et al. Pangenome obtained by long-read sequencing of 11 genomes reveal hidden functional structural variants in pigs. *iScience* **26**, 106119 (2023).
- Khan, A. W. et al. Super-pangenome by integrating the wild side of a species for accelerated crop improvement. *Trends Plant Sci.* **25**, 148–158 (2020).
- Cochetel, N. et al. A super-pangenome of the North American wild grape species. *Genome Biol.* **24**, 290 (2023).
- Li, N. et al. Super-pangenome analyses highlight genomic diversity and structural variation across wild and cultivated tomato species. *Nat. Genet.* **55**, 852–860 (2023).
- Leonard, A. S., Crysanto, D., Mapel, X. M., Bhati, M. & Pausch, H. Graph construction method impacts variation representation and analyses in a bovine super-pangenome. *Genome Biol.* **24**, 124 (2023).
- Shi, T. et al. The super-pangenome of *Populus* unveils genomic facets for its adaptation and diversification in widespread forest trees. *Mol. Plant* **17**, 725–746 (2024).
- Secomandi, S. et al. A chromosome-level reference genome and pangenome for barn swallow population genomics. *Cell Rep.* **42**, 111992 (2023).
- Fang, B. & Edwards, S. V. Fitness consequences of structural variation inferred from a House Finch pangenome. *Proc. Natl Acad. Sci. USA* **121**, e2409943121 (2024).
- Pang, A. W. et al. Towards a comprehensive structural variation map of an individual human genome. *Genome Biol.* **11**, R52 (2010).
- Chakraborty, M., Emerson, J. J., Macdonald, S. J. & Long, A. D. Structural variants exhibit widespread allelic heterogeneity and shape variation in complex traits. *Nat. Commun.* **10**, 4872 (2019).
- Tigano, A. et al. Chromosome-level assembly of the Atlantic silverside genome reveals extreme levels of sequence diversity and structural genetic variation. *Genome Biol. Evol.* **13**, evab098 (2021).
- Todesco, M. et al. Massive haplotypes underlie ecotypic differentiation in sunflowers. *Nature* **584**, 602–607 (2020).
- Garg, S., Balboa, R. & Kuja, J. Chromosome-scale haplotype-resolved pangenomics. *Trends Genet.* **38**, 1103–1107 (2022).

35. Wang, S., Qian, Y.-Q., Zhao, R.-P., Chen, L.-L. & Song, J.-M. Graph-based pan-genomes: increased opportunities in plant genomics. *J. Exp. Bot.* **74**, 24–39 (2023).
36. Liao, W.-W. et al. A draft human pangenome reference. *Nature* **617**, 312–324 (2023).
37. Paten, B. et al. Superbubbles, ultrabubbles, and cacti. *J. Comput. Biol.* **25**, 649–663 (2018).
38. Wang, T. et al. The Human Pangenome Project: a global resource to map genomic diversity. *Nature* **604**, 437–446 (2022).
39. Miga, K. H. & Wang, T. The need for a human pangenome reference sequence. *Annu. Rev. Genomics Hum. Genet.* **22**, 81–102 (2021).
40. Read, B. A. et al. Pan genome of the phytoplankton *Emiliania huxleyi* underpins its global distribution. *Nature* **499**, 209–213 (2013).
41. Liu, Y. et al. Pan-genome of wild and cultivated soybeans. *Cell* **182**, 162–176 (2020).
42. Talenti, A. et al. A cattle graph genome incorporating global breed diversity. *Nat. Commun.* **13**, 910 (2022).
43. Bozan, I. et al. Pangenome analyses reveal impact of transposable elements and ploidy on the evolution of potato species. *Proc. Natl Acad. Sci. USA* **120**, e2211117120 (2023).
44. Tong, X. et al. High-resolution silkworm pan-genome provides genetic insights into artificial selection and ecological adaptation. *Nat. Commun.* **13**, 5619 (2022).
45. Golicz, A. A., Bayer, P. E., Bhalla, P. L., Batley, J. & Edwards, D. Pangenomics comes of age: from bacteria to plant and animal applications. *Trends Genet.* **36**, 132–145 (2020).
46. Eberlein, C. et al. Hybridization is a recurrent evolutionary stimulus in wild yeast speciation. *Nat. Commun.* **10**, 923 (2019).
47. Mavárez, J. et al. Speciation by hybridization in *Heliconius* butterflies. *Nature* **441**, 868–871 (2021).
48. Hübner, S. et al. Sunflower pan-genome analysis shows that hybridization altered gene content and disease resistance. *Nat. Plants* **5**, 54–62 (2019).
49. Nurk, S. et al. The complete sequence of a human genome. *Science* **376**, 44–53 (2022).
50. Koren, S. et al. De novo assembly of haplotype-resolved genomes with trio binning. *Nat. Biotechnol.* **36**, 1174–1182 (2018).
51. Cheng, H. et al. Haplotype-resolved assembly of diploid genomes without parental data. *Nat. Biotechnol.* **40**, 1332–1335 (2022).
52. Korlach, J. et al. Real-time DNA sequencing from single polymerase molecules. *Methods Enzymol.* **472**, 431–455 (2010).
53. Mikheyev, A. S. & Tin, M. M. Y. A first look at the Oxford Nanopore MinION sequencer. *Mol. Ecol. Resour.* **14**, 1097–1102 (2014).
54. Leonard, A. S. et al. Structural variant-based pangenome construction has low sensitivity to variability of haplotype-resolved bovine assemblies. *Nat. Commun.* **13**, 3012 (2022).
55. Olagunju, T. A. et al. Telomere-to-telomere assemblies of cattle and sheep Y-chromosomes uncover divergent structure and gene content. *Nat. Commun.* **15**, 8277 (2024).
56. Beyter, D. et al. Long-read sequencing of 3,622 Icelanders provides insight into the role of structural variants in human diseases and other traits. *Nat. Genet.* **53**, 779–786 (2021).
57. Kim, J. et al. The genome landscape of indigenous African cattle. *Genome Biol.* **18**, 34 (2017).
58. Kim, K. et al. The mosaic genome of indigenous African cattle as a unique genetic resource for African pastoralism. *Nat. Genet.* **52**, 1099–1110 (2020).
59. Bakhtiari, M. et al. Variable number tandem repeats mediate the expression of proximal genes. *Nat. Commun.* **12**, 2075 (2021).
60. Caballero-López, V., Lundberg, M., Sokolovskis, K. & Bensch, S. Transposable elements mark a repeat-rich region associated with migratory phenotypes of willow warblers (*Phylloscopus trochilus*). *Mol. Ecol.* **31**, 1128–1141 (2022).
61. Li, H., Feng, X. & Chu, C. The design and construction of reference pangenome graphs with minigraph. *Genome Biol.* **21**, 265 (2020).
62. Garrison, E. et al. Building pangenome graphs. *Nat. Methods* **21**, 2008–2012 (2024).
63. Garrison, E. & Guarracino, A. Unbiased pangenome graphs. *Bioinformatics* **39**, btac743 (2023).
64. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
65. Armstrong, J. et al. Progressive Cactus is a multiple-genome aligner for the thousand-genome era. *Nature* **587**, 246–251 (2020).
66. Marco-Sola, S. et al. Optimal gap-affine alignment in O(s) space. *Bioinformatics* **39**, btad074 (2023).
67. Sibbesen, J. A. et al. Haplotype-aware pantranscriptome analyses using spliced pangenome graphs. *Nat. Methods* **20**, 239–247 (2023).
68. Andreatta, F., Lechat, P., Dufresne, Y. & Chikhi, R. Comparing methods for constructing and representing human pangenome graphs. *Genome Biol.* **24**, 274 (2023).
69. Heumos, S. et al. Cluster-efficient pangenome graph construction with nf-core/pangenome. *Bioinformatics* **40**, btae609 (2024).
70. Nergadze, S. G. et al. Birth, evolution, and transmission of satellite-free mammalian centromeric domains. *Genome Res.* **28**, 789–799 (2018).
71. Ávila Robledillo, L. et al. Extraordinary sequence diversity and promiscuity of centromeric satellites in the legume tribe *Fabeae*. *Mol. Biol. Evol.* **37**, 2341–2356 (2020).
72. Guarracino, A. et al. Recombination between heterologous human acrocentric chromosomes. *Nature* **617**, 335–343 (2023).
73. Garrison, E., Guarracino, A. & Kille, B. impg: implicit pangenome graph. *GitHub* github.com/pangenome/imp (2024).
74. Guarracino, A., Heumos, S., Nahnsen, S., Prins, P. & Garrison, E. ODGI: understanding pangenome graphs. *Bioinformatics* **38**, 3319–3326 (2022).
75. Sirén, J. Indexing variation graphs. In *2017 Proc. Meeting on Algorithm Engineering and Experiments (ALENEX)* (eds Fekete, S. & Ramachandran, V.) 13–27 (Society for Industrial and Applied Mathematics, 2017).
76. Garrison, E. et al. Variation graph toolkit improves read mapping by representing genetic variation in the reference. *Nat. Biotechnol.* **36**, 875–879 (2018).
77. Fiddes, I. T. et al. Comparative Annotation Toolkit (CAT)—simultaneous clade and personal genome annotation. *Genome Res.* **28**, 1029–1038 (2018).
78. Wick, R. R., Schultz, M. B., Zobel, J. & Holt, K. E. Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics* **31**, 3350–3352 (2015).
79. Gonnella, G., Niehus, N. & Kurtz, S. GfaViz: flexible and interactive visualization of GFA sequence graphs. *Bioinformatics* **35**, 2853–2855 (2019).
80. Beyer, W. et al. Sequence tube maps: making graph genomes intuitive to commuters. *Bioinformatics* **35**, 5318–5320 (2019).
81. Yokoyama, T. T., Sakamoto, Y., Seki, M., Suzuki, Y. & Kasahara, M. MoMI-G: modular multi-scale integrated genome graph browser. *BMC Bioinformatics* **20**, 548 (2019).
82. Krzywinski, M. et al. Circos: an information aesthetic for comparative genomics. *Genome Res.* **19**, 1639–1645 (2009).
83. Fischer, C. & Garrison, E. Waragraph. *GitHub* github.com/chfi/waragraph (2023).
84. Li, H. gfatools: tools for manipulating sequence graphs in the GFA and rGFA formats. *GitHub* github.com/lh3/gfatools (2024).
85. Sedlazeck, F. J., Lee, H., Darby, C. A. & Schatz, M. C. Piercing the dark matter: bioinformatics of long-range sequencing and mapping. *Nat. Rev. Genet.* **19**, 329–346 (2018).
86. Noor, M. A., Grams, K. L., Bertucci, L. A. & Reiland, J. Chromosomal inversions and the reproductive isolation of species. *Proc. Natl Acad. Sci. USA* **98**, 12084–12088 (2001).

87. Küpper, C. et al. A supergene determines highly divergent male reproductive morphs in the ruff. *Nat. Genet.* **48**, 79–83 (2016).
88. Weissensteiner, M. H. et al. Discovery and population genomics of structural variation in a songbird genus. *Nat. Commun.* **11**, 3403 (2020).
89. Wold, J. et al. Expanding the conservation genomics toolbox: incorporating structural variants to enhance genomic studies for species of conservation concern. *Mol. Ecol.* **30**, 5949–5965 (2021).
90. Wellenreuther, M. & Bernatchez, L. Eco-evolutionary genomics of chromosomal inversions. *Trends Ecol. Evol.* **33**, 427–440 (2018).
91. Peona, V. et al. The hidden structural variability in avian genomes. Preprint at *bioRxiv* <https://doi.org/10.1101/2021.12.31.473444> (2022).
92. Sirén, J. et al. Pangenomics enables genotyping of known structural variants in 5202 diverse genomes. *Science* **374**, abg8871 (2021).
93. Sirén, J., Garrison, E., Novak, A. M., Paten, B. & Durbin, R. Haplotype-aware graph indexes. *Bioinformatics* **36**, 400–407 (2020).
94. Rautiainen, M. & Marschall, T. GraphAligner: rapid and versatile sequence-to-graph alignment. *Genome Biol.* **21**, 253 (2020).
95. Wu, Z. et al. Human pangenome analysis of sequences missing from the reference genome reveals their widespread evolutionary, phenotypic, and functional roles. *Nucleic Acids Res.* **52**, 2212–2230 (2024).
96. Tamagawa, K., Yoshida, K., Ohnishi, S. & Takahashi, Y. Population transcriptomics reveals the effect of gene flow on the evolution of range limits. *Sci. Rep.* **12**, 1318 (2022).
97. DeBiasse, M. B., Kawji, Y. & Kelly, M. W. Phenotypic and transcriptomic responses to salinity stress across genetically and geographically divergent *Tigriopus californicus* populations. *Mol. Ecol.* **27**, 1621–1632 (2018).
98. Liu, L., Wang, Z., Su, Y. & Wang, T. Population transcriptomic sequencing reveals allopatric divergence and local adaptation in *Pseudotaxus chienii* (Taxaceae). *BMC Genomics* **22**, 388 (2021).
99. Ma, S. & Zhang, Y. Profiling chromatin regulatory landscape: insights into the development of ChIP-seq and ATAC-seq. *Mol. Biomed.* **1**, 9 (2020).
100. Hickey, G. et al. Genotyping structural variants in pangenome graphs using the vg toolkit. *Genome Biol.* **21**, 35 (2020).
101. Poplin, R. et al. A universal SNP and small-indel variant caller using deep neural networks. *Nat. Biotechnol.* **36**, 983–987 (2018).
102. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. Preprint at arxiv.org/abs/1207.3907 (2012).
103. Van der Auwera, G. A. & O'Connor, B. D. *Genomics in the Cloud: Using Docker, GATK, and WDL in Terra* (O'Reilly Media, 2020).
104. Herzeel, C. et al. Multithreaded variant calling in eRep 5. *PLoS ONE* **16**, e0244471 (2021).
105. Ebler, J. et al. Pangenome-based genome inference allows efficient and accurate genotyping across a wide spectrum of variant classes. *Nat. Genet.* **54**, 518–525 (2022).
106. Li, M.-H. & Merilä, J. Population differences in levels of linkage disequilibrium in the wild. *Mol. Ecol.* **20**, 2916–2928 (2011).
107. Charmantier, A., Garant, D. & Kruuk, L. E. B. *Quantitative Genetics in the Wild* (Oxford University Press, 2014).
108. Formenti, G. et al. The era of reference genomes in conservation genomics. *Trends Ecol. Evol.* **37**, 197–202 (2022).
109. Darwin Tree of Life Project Consortium. Sequence locally, think globally: the Darwin Tree of Life Project. *Proc. Natl Acad. Sci. USA* **119**, e2115642118 (2022).
110. Mazzoni, C. J., Ciofi, C. & Waterhouse, R. M. Biodiversity: an atlas of European reference genomes. *Nature* **619**, 252 (2023).
111. Rubin, J. D., Vogel, N. A., Gopalakrishnan, S., Sackett, P. W. & Renaud, G. HaploCart: human mtDNA haplogroup classification using a pangenomic reference graph. *PLoS Comput. Biol.* **19**, e1011148 (2023).
112. Feng, S. et al. Dense sampling of bird diversity increases power of comparative genomics. *Nature* **587**, 252–257 (2020).
113. Yang, Z. et al. Pangenome graphs in infectious disease: a comprehensive genetic variation analysis of *Neisseria meningitidis* leveraging Oxford Nanopore long reads. *Front. Genet.* **14**, 1225248 (2023).
114. Lombardo, G. et al. The mitogenome relationships and phylogeography of barn swallows (*Hirundo rustica*). *Mol. Biol. Evol.* **39**, msac113 (2022).
115. Payseur, B. A. & Rieseberg, L. H. A genomic perspective on hybridization and speciation. *Mol. Ecol.* **25**, 2337–2360 (2016).
116. Hewitt, G. M. Hybrid zones—natural laboratories for evolutionary studies. *Trends Ecol. Evol.* **3**, 158–167 (1988).
117. Sebastianelli, M. et al. A genomic basis of vocal rhythm in birds. *Nat. Commun.* **15**, 3095 (2024).
118. Irwin, D. E. et al. A comparison of genomic islands of differentiation across three young avian species pairs. *Mol. Ecol.* **27**, 4839–4855 (2018).
119. Hejase, H. A. et al. Genomic islands of differentiation in a rapid avian radiation have been driven by recent selective sweeps. *Proc. Natl Acad. Sci. USA* **117**, 30554–30565 (2020).
120. van der Sprong, J. et al. A novel target-enriched multilocus assay for sponges (Porifera): Red Sea Haplosclerida (Demospongiae) as a test case. *Mol. Ecol. Resour.* **24**, e13891 (2024).
121. Novak, B. J. De-extinction. *Genes* **9**, 548 (2018).
122. Tettelin, H. et al. Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the microbial ‘pan-genome’. *Proc. Natl Acad. Sci. USA* **102**, 13950–13955 (2005).
123. Obert, C. et al. Identification of a candidate *Streptococcus pneumoniae* core genome and regions of diversity correlated with invasive pneumococcal disease. *Infect. Immun.* **74**, 4766–4777 (2006).
124. Rasko, D. A. et al. The pangenome structure of *Escherichia coli*: comparative genomic analysis of *E. coli* commensal and pathogenic isolates. *J. Bacteriol.* **190**, 6881–6893 (2008).
125. Donati, C. et al. Structure and dynamics of the pan-genome of *Streptococcus pneumoniae* and closely related species. *Genome Biol.* **11**, R107 (2010).
126. Pinto, M. et al. Insights into the population structure and pan-genome of *Haemophilus influenzae*. *Infect. Genet. Evol.* **67**, 126–135 (2019).
127. Aggarwal, S. K. et al. Pangenomics in microbial and crop research: progress, applications, and perspectives. *Genes* **13**, 598 (2022).
128. Rouli, L., Merhej, V., Fournier, P.-E. & Raoult, D. The bacterial pangenome as a new tool for analysing pathogenic bacteria. *New Microbes New Infect.* **7**, 72–85 (2015).
129. Kiu, R., Caim, S., Alexander, S., Pachori, P. & Hall, L. J. Probing genomic aspects of the multi-host pathogen *Clostridium perfringens* reveals significant pangenome diversity, and a diverse array of virulence factors. *Front. Microbiol.* **8**, 2485 (2017).
130. Poulsen, B. E. et al. Defining the core essential genome of *Pseudomonas aeruginosa*. *Proc. Natl Acad. Sci. USA* **116**, 10072–10080 (2019).
131. Hisham, Y. & Ashhab, Y. Identification of cross-protective potential antigens against pathogenic *Brucella* spp. through combining pan-genome analysis with reverse vaccinology. *J. Immunol. Res.* **2018**, 1474517 (2018).
132. Naz, K. et al. PanRV: pangenome-reverse vaccinology approach for identifications of potential vaccine candidates in microbial pangenome. *BMC Bioinformatics* **20**, 123 (2019).
133. Francis, W. R. & Wörheide, G. Similar ratios of introns to intergenic sequence across animal genomes. *Genome Biol. Evol.* **9**, 1582–1598 (2017).
134. Morgante, M., De Paoli, E. & Radovic, S. Transposable elements and the plant pan-genomes. *Curr. Opin. Plant Biol.* **10**, 149–155 (2007).

135. Adams, K. L. & Wendel, J. F. Polyploidy and genome evolution in plants. *Curr. Opin. Plant Biol.* **8**, 135–141 (2005).
136. Cheng, F. et al. Gene retention, fractionation and subgenome differences in polyploid plants. *Nat. Plants* **4**, 258–268 (2018).
137. Shi, J., Tian, Z., Lai, J. & Huang, X. Plant pan-genomics and its applications. *Mol. Plant* **16**, 168–186 (2023).
138. Sun, C. et al. RPAN: rice pan-genome browser for ~3000 rice genomes. *Nucleic Acids Res.* **45**, 597–605 (2017).
139. Gao, L. et al. The tomato pan-genome uncovers new genes and a rare allele regulating fruit flavor. *Nat. Genet.* **51**, 1044–1051 (2019).
140. Jayakodi, M. et al. The barley pan-genome reveals the hidden legacy of mutation breeding. *Nature* **588**, 284–289 (2020).
141. Song, J.-M. et al. Eight high-quality genomes reveal pan-genome architecture and ecotype differentiation of *Brassica napus*. *Nat. Plants* **6**, 34–45 (2020).
142. Li, J. et al. Cotton pan-genome retrieves the lost sequences and genes during domestication and selection. *Genome Biol.* **22**, 119 (2021).
143. Tao, Y., Zhao, X., Mace, E., Henry, R. & Jordan, D. Exploring and exploiting pan-genomics for crop improvement. *Mol. Plant* **12**, 156–169 (2019).
144. Qin, P. et al. Pan-genome analysis of 33 genetically diverse rice accessions reveals hidden genomic variations. *Cell* **184**, 3542–3558 (2021).
145. Zhou, Y. et al. Graph pangenome captures missing heritability and empowers tomato breeding. *Nature* **606**, 527–534 (2022).
146. Zanini, S. F. et al. Pangenomics in crop improvement—from coding structural variations to finding regulatory variants with pangenome graphs. *Plant Genome* **15**, e20177 (2022).
147. Qiao, Q. et al. Evolutionary history and pan-genome dynamics of strawberry (*Fragaria* spp.). *Proc. Natl Acad. Sci. USA* **118**, e2105431118 (2021).
148. Tian, X. et al. Building a sequence map of the pig pan-genome from multiple de novo assemblies and Hi-C data. *Sci. China Life Sci.* **63**, 750–763 (2020).
149. Li, Z. et al. The pig pangenome provides insights into the roles of coding structural variations in genetic diversity and adaptation. *Genome Res.* **33**, 1833–1847 (2023).
150. Jang, J. et al. Chromosome-level genome assembly of Korean native cattle and pangenome graph of 14 *Bos taurus* assemblies. *Sci. Data* **10**, 560 (2023).
151. Dai, X. et al. A Chinese indicine pangenome reveals a wealth of novel structural variants introgressed from other *Bos* species. *Genome Res.* **33**, 1284–1298 (2023).
152. Li, R. et al. A sheep pangenome reveals the spectrum of structural variations and their effects on tail phenotypes. *Genome Res.* **33**, 463–477 (2023).
153. Saco, A. et al. Gene presence/absence variation in *Mytilus galloprovincialis* and its implications in gene expression and adaptation. *iScience* **26**, 107827 (2023).
154. Ruggieri, A. A. et al. Erratum: a butterfly pan-genome reveals that a large amount of structural variation underlies the evolution of chromatin accessibility. *Genome Res.* **32**, 2145 (2022).
155. Li, R. et al. Building the sequence map of the human pan-genome. *Nat. Biotechnol.* **28**, 57–63 (2010).
156. Besenbacher, S. et al. Novel variation and de novo mutation rates in population-wide de novo assembled Danish trios. *Nat. Commun.* **6**, 5969 (2015).
157. Maretty, L. et al. Sequencing and de novo assembly of 150 genomes from Denmark as a population reference. *Nature* **548**, 87–91 (2017).
158. Sherman, R. M. et al. Assembly of a pan-genome from deep sequencing of 910 humans of African descent. *Nat. Genet.* **51**, 30–35 (2019).
159. Duan, Z. et al. HUPAN: a pan-genome analysis pipeline for human genomes. *Genome Biol.* **20**, 149 (2019).
160. Wong, K. H. Y. et al. Towards a reference genome that captures global genetic diversity. *Nat. Commun.* **11**, 5482 (2020).
161. Gao, Y. et al. A pangenome reference of 36 Chinese populations. *Nature* **619**, 112–121 (2023).
162. Tetikol, H. S. et al. Pan-African genome demonstrates how population-specific genome graphs improve high-throughput sequencing data analysis. *Nat. Commun.* **13**, 4384 (2022).
163. Spina, F. The EURING swallow project: a largescale approach to the study and conservation of a long-distance migrant. In *Migrating Birds Know No Boundaries. Proc. International Seminar 1997* (eds Leshem, J. et al.) 151–162 (1998).
164. Møller, A. P. Sexual selection in the barn swallow (*Hirundo rustica*). IV. Patterns of fluctuating asymmetry and selection against asymmetry. *Evolution* **48**, 658–670 (1994).
165. Scordato, E. S. C. et al. Genomic variation across two barn swallow hybrid zones reveals traits associated with divergence in sympatry and allopatry. *Mol. Ecol.* **26**, 5676–5691 (2017).

Acknowledgements

We are grateful to the HPRC community for useful discussions over the years that helped shape the Review. We particularly thank E. Garrison for his input to the Review. We thank C. Di Pietro for drawing the final figures. C.R.F. thanks the support of CE3C through an assistant researcher contract (FCiência.ID contract 366) and the Fundação para a Ciência e a Tecnologia for Portuguese National Funds attributed to CE3C within the projects UIDB/00329/2020, UIDP/00329/2020 and LA/P/0121/2020 and the FPUL for a contract of invited assistant professor.

Author contributions

G.R.G., G.F. and S.S. conceived the Review. S.S., G.R.G., G.F. and L.G. drafted the manuscript and ideated the figures, with substantial contributions from A.B.-A., C.R.F., R.R. and E.D.J. All authors reviewed and approved the final text.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41588-024-02029-6>.

Correspondence and requests for materials should be addressed to Giulio Formenti.

Peer review information *Nature Genetics* thanks Alan Hoelzel and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.

© Springer Nature America, Inc. 2025