

Pangenome graphs and their applications in biodiversity genomics

Simona Secomandi^{1,*}, Guido Roberto Gallo^{2,*}, Riccardo Rossi³, Carlos Rodríguez Fernandes^{4,5}, Erich D. Jarvis^{1,6}, Andrea Bonisoli-Alquati⁷, Luca Gianfranceschi², Giulio Formenti^{6,†}

* Co-first authors

† Corresponding author (e-mail address: gformenti@rockefeller.edu)

¹Laboratory of Neurogenetics of Language, The Rockefeller University, New York, NY, USA

²Department of Biosciences, The University of Milan, Via Giovanni Celoria 26, Milan MI, 20133 Italy

³Department of Biotechnology and Biosciences, University of Milano - Bicocca, Piazza della Scienza 2, 20126 Milan, Italy

⁴CE3C - Centre for Ecology, Evolution and Environmental Changes & CHANGE - Global Change and Sustainability Institute, Departamento de Biologia Animal, Faculdade de Ciências, Universidade de Lisboa, 1749-016 Lisboa, Portugal

⁵Faculdade de Psicologia, Universidade de Lisboa, Alameda da Universidade, 1649-013 Lisboa, Portugal

⁶The Vertebrate Genome Laboratory, 1230 York Ave, New York, NY 10065 USA

⁷Department of Biological Sciences, California State Polytechnic University, Pomona. Pomona, CA

Abstract

Complete datasets of genetic variants are key to biodiversity genomics 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.

32 **Main**

33

34 The decrease in DNA sequencing costs is boosting biodiversity genomics¹². Large-scale
35 international initiatives aim to generate highly-contiguous, haplotype-resolved reference
36 genomes for all species³⁻⁷, to characterize biodiversity, clarify its evolution, and help its
37 conservation. Reference genomes are the backbones to annotate genomic variation, helping its
38 preservation in the current biodiversity crisis². Haplotype-resolved reference genomes better
39 capture genetic variation across individuals, enabling more effective mapping of phenomes to
40 genomes, and illuminating both the adaptive uniqueness of taxa⁸ and their position in the tree
41 of life⁹. However, even with accurate reference genomes, sequences mapped against them can
42 be misplaced or fail to align because of divergent or missing regions in the reference (**Fig. 1a**).
43 These ‘blind spots’¹⁰⁻¹² introduce a reference bias^{13,14}, potentially misrepresenting genetic
44 variation. A solution to reference bias lies in pangenomics, the systematic capturing of genetic
45 variation within dedicated composite assemblies, called pangenomes^{15,16}. Historically, a
46 pangenome can be a collection of unaligned sequences^{17,18}, or a graph derived from raw reads¹⁹
47 or assembled sequences¹⁵ (**Fig. 1b; Box 1**). A pangenome graph is a graphical model storing
48 genomic data from different individuals, together with their relationships and variability, in a
49 single data structure (**Fig. 1c**). Data embedded in the graph can be accessed to perform
50 bioinformatics tasks, such as read mapping, by graph indexing^{15,20}. A pangenome graph
51 accommodates multiple alternative alleles and provides a more comprehensive representation
52 of genomic variation within a species^{21,22}, depending on the degree to which the haplotypes set
53 reflects the overall diversity of the species. At higher taxonomic ranks (e.g. genus or family),
54 super-pangenomes expand our ability to capture genetic diversity²³⁻²⁷. While cataloging
55 genetic variability within a species is critical²⁸, examples of pangenome graphs in non-model
56 species are still limited²⁹. We expect this to change in the coming years, as long-read
57 sequencing and high-quality genomes become more accessible^{2,30}. Pangenome graphs improve
58 downstream analyses, providing researchers with the flexibility to analyze genetic variation at
59 multiple levels in a single data structure. In agriculture, pangenomes for major crop plants
60 proved effective in identifying resistance genes and beneficial alleles for crop improvement,
61 while in humans they show potential to improve clinical research (**Box 1**). Pangenome graphs
62 enable the accurate detection of structural variants (SVs)^{31,32}, large (> 50 bp) genomic
63 rearrangements that underlie phenotype and fitness variation^{31,32} and disproportionately
64 contribute to local adaptation³³, and even speciation³⁴. Pangenomes graphs will increase the
65 accuracy and power of resequencing projects investigating population dynamics and insights

66 into the genetic bases of phenotypic traits, by streamlining variant phasing, haplotype
67 reconstruction, and genotyping through imputation. They will aid in understanding genome
68 evolution across species and benefit studies linking genetic variation and gene expression to
69 phenotypes. Here, we illustrate the available data structures and tools for building, visualizing,
70 and manipulating pangenome graphs, and provide practical advice for their use in downstream
71 analyses. We also highlight their potential future contribution to conservation and biodiversity
72 genomics.

73

74 **Pangenomes as variation graphs**

75 Over the years, the term pangenome acquired different meanings (**Box 1**). Here, we focus on
76 graph-based pangenomes constructed from whole-genome alignments of assembled sequences,
77 i.e. variation graphs¹⁵. Variation graphs more comprehensively represent eukaryotic genomes
78 by storing complete genomic sequences and their variation^{21,35–37}. Variation graphs compress
79 redundant sequences into bidirected networks where each node represents a sequence and
80 edges connect nodes into complete sequences¹⁵ (**Fig. 1c**). Linear genomes or phased haplotypes
81 are stored as explicit paths through the graph¹⁵, and sequence variation is represented by
82 subgraphs called bubbles, or snarls, in which variants are defined by alternative paths
83 connected by shared start and end nodes³⁸. Here, we refer to variation graphs from whole-
84 genome alignments simply as “pangenome graphs”. They were adopted by the Human
85 Pangenome Reference Consortium (HPRC)³⁹, a global initiative that generated the first human
86 pangenome^{37,40} (**Box 1**). Pangenome graphs were also assembled for the chicken (*Gallus*
87 *gallus*)²¹, and, among non-model organisms, for the barn swallow (*Hirundo rustica*; **Box 2**)²⁹.
88 Super-pangenome graphs are being generated for economically important species, such as
89 tomato (*Solanum lycopersicum*)²⁵, grape (*Vitis sp.*; **Box 2**)²⁴ and cattle (*Bos taurus*)²⁶.

90 **Maximizing capture of genetic diversity via sampling and sequencing**

91 The sampling strategy is critical for successful biodiversity pangenomic studies. It should
92 maximize genomic and biogeographic diversity within natural populations, ideally sampling
93 through the entire geographic range, while balancing sex representation^{18,24,25,37,39,41–44} (**Fig.**
94 **2a**). If a panel of variants, both SNPs and SVs, is available, estimates of heterozygosity,
95 relatedness, and inbreeding offer insights into the sample size required to achieve a
96 comprehensive representation. Ordination and clustering analyses can help select
97 representative individuals for inclusion in a pangenome²⁴. The ideal sample size can also be

98 retrospectively verified by a pangene number analysis, where a number of representative
99 genomes are progressively added until a pangene plateau is reached, i.e., when the full set of
100 genes is captured and adding more individuals does not recover novel genes⁴⁵. The inclusion
101 of more individuals augment the existing references by increasing the representation of
102 accessory, or dispensable, sequences, which are shared only by a subset of individuals and
103 often have functional and adaptive roles⁴⁶. Pangenome graphs facilitate the pinpointing of
104 functional and adaptive roles of accessory genomic regions and how they vary among
105 geographic populations and subspecies. Accessory genomes may provide hotspots for
106 population differentiation and speciation through divergent selection processes or via
107 hybridization^{47,48}. Investigating how much accessory regions are affected by introgression in
108 populations and species undergoing hybridization provides insights into the dynamics of gene
109 flow and speciation processes⁴⁹.

110 Among the collected samples, priority should be given to highest-quality ones, maximizing
111 long-read sequencing throughput (**Fig. 2a**) and allowing accurate and contiguous genome
112 assemblies^{5,50}. The quality of the input genomes minimizes noise propagation in the
113 pangenome graph. Haplotype phasing with parental sequence data or chromatin conformation
114 data (Hi-C), is crucial to prevent haplotype false duplication and related errors^{4,5,51,52}.
115 Sequencing coverages of ~30-fold PacBio high-fidelity long reads (HiFi)^{11,53} and ~60-fold
116 Oxford Nanopore Technologies (ONT)⁵⁴ Duplex reads, in combination with ~30-fold Hi-C per
117 haplotype and manual curation, generate reference genomes that meet the current quality
118 standards^{5,7,55}. Genome sequencing is a rapidly evolving field, and generating complete,
119 haplotype-resolved, and near-error-free genomes (telomere-to-telomere, T2T) is now
120 feasible^{37,50,56}, by complementing HiFi with ultra-long ONT reads. Incorporating high-quality
121 genomes aids in the discovery of rare SVs⁵⁷, particularly in admixed populations and those
122 with large effective population sizes^{58,59}. Moreover, it will improve the representation of hard-
123 to-sequence and assemble regions like centromeres, variable number tandem repeats
124 (VNTRs)³⁷ and other complex repeats. Examining base-level polymorphism in VNTRs may
125 clarify their role in shaping gene expression and complex traits^{10,26,60}. Highly repetitive regions
126 might also underlie the regulation of complex behavioral phenotypes, such as migratory
127 behavior⁶¹.

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

132 species, and sequencing multiple individuals with different long-read technologies may
133 currently be cost-prohibitive. We suggest that a pangenome graph should include at least one
134 high-quality assembly as a backbone for graph construction⁵⁵, providing a robust coordinate
135 system for downstream analyses.

136 **Building pangenome graphs to represent complex and accessory genomic** 137 **regions**

138 Pangenome graph construction starts with the alignment of the input genomes to identify
139 sequence similarities. Alignment can be reference-based^{20,62} or involve all-versus-all
140 comparisons^{63,64}, and it can be either at the base-level^{20,63} or at a higher level (e.g. including
141 only variant sites)⁶². Alignment information is embedded in the graph, which can be
142 manipulated for downstream analyses. Two main pipelines were developed by the HPRC for
143 graph construction: Minigraph-Cactus (MC)²⁰ and the PanGenome Graph Builder (PGGB)⁶³
144 (**Fig. 2b; Supplementary Table 1**). MC implements Minigraph⁶², a sequence-to-graph aligner,
145 as a graph constructor. In MC, a user-selected reference genome is used as the initial backbone,
146 which is progressively augmented with structural variation from the other genomes. The
147 resulting graph is SV-only (> 50 bp, **Fig. 2b**) and all assemblies are aligned back to the graph
148 with a minimap2-like⁶⁵ algorithm that generates base-level alignments for each reference
149 chromosome. MC implements a modified version of the reference-free aligner Progressive
150 Cactus⁶⁶ to combine the alignments into base-level pangenome graphs that contain variants of
151 all sizes (**Fig. 2b**). Chromosomal graphs are then combined and post-processed to reduce path
152 complexity by collapsing redundant sequences²⁰. In addition to the chosen reference, one may
153 specify additional assemblies whose coordinates can serve as reference in downstream
154 analyses²⁰. In contrast to MC, PGGB⁶³ avoids using an initial reference, and rather employs
155 all-to-all genome alignments with wfmash⁶⁷, a software for homology mapping that generates
156 base-wise pairwise alignments. PGGB uses seqwish⁶⁴ as a sequence-to-graph aligner, which
157 starts from the all-versus-all alignments to generate a complete pangenome graph, representing
158 all variant types and sizes (**Fig. 2b**). The graph is then post-processed with a smoothing and
159 normalization step^{37,63}. PGGB can be run in parallel on each chromosome community to reduce
160 computation time^{37,63}. In PGGB graphs, every genome included in the graph can be used as
161 reference for downstream analyses⁶³. However, like in MC, only one reference can be used at
162 a time as a coordinate system for downstream analyses such as variant calling. New
163 computational methods and file formats other than the linear binary alignment map (BAM) and
164 variant call format (VCF) need to be developed to overcome this limitation and represent all

165 the information embedded in the graph. Pangenome graphs can be combined with transcript
166 annotations using the variation graph toolkit (vg)⁶⁸, a software for variation graph construction,
167 handling and analysis, into splice-aware graphs, with paths through nodes (exons) and edges
168 (splice junctions) representing the structure of mRNA transcripts. It is also possible to build
169 pantranscriptomes by projecting a set of haplotype-specific transcripts onto a set of known
170 haplotypes⁶⁸.

171 The size of a pangenome graph depends upon the genome size of the respective species, but is
172 bound to be larger, as it incorporates accessory sequences from other individuals, and is also it
173 is influenced by the number and diversity of the individuals contributing to a pangenome, as
174 well as the construction pipeline (**Supplementary Table 1**). The size of the MC graph is
175 relatively close to the genome size of the species, (~3.2 Gbp versus 3.1 Gbp for human³⁷, 1.2
176 Gbp versus 1.1 Gbp for chicken²¹, ~1.6 Gbp versus 1.1 Gbp for barn swallow²⁹;
177 **Supplementary Table 1; Box 2**). In contrast, PGGB graph size can considerably exceed that
178 of the genome size and MC graphs (e.g. 8.4 Gbp for human³⁷). The larger size of PGGB graphs
179 is explained by their capability to capture highly divergent satellite, centromeric and
180 heterochromatic regions, excluded in MC graphs due to alignment issues^{21,37} (**Supplementary**
181 **Table 1**). The largest increase relative to true genome size was observed in grapevine²⁴, likely
182 because of the inclusion of different species (**Box 2**). PGGB also has a tendency to collapse
183 complex regions, such as copy-number polymorphic loci, into a single copy, generating loops
184 in the graph that increase its complexity^{37,63} (**Fig. 2b**). Given their greater size and complexity,
185 PGGB graphs require more computational resources than MC graphs. To construct a graph
186 based on 10 human haplotypes, MC currently takes ~16 hours, 154 Gb of RAM and 7 Gb of
187 disk space, while PGGB takes 117 hours, 71 Gb of RAM and 7.6 Gb of disk space⁶⁹. PGGB
188 has also been experimentally shown to potentially lead to an overestimation of sequence
189 variability. For instance, the size of the PGGB chicken pangenome was larger than expected
190 based on the estimated variation in diverse groups of chickens²¹. A preliminary estimation of
191 the expected species variability should be computed to detect overestimation in graph
192 construction.

193 Given all these differences, a careful selection of the pangenome graph construction pipeline
194 is of utmost importance. On one hand, PGGB graphs are based on reference-free alignments
195 and are more complete than MC graphs. On the other hand, their complexity increases the
196 computational resources needed for graph construction and some downstream analyses, such
197 as variant calling after read mapping, are currently computationally infeasible^{21,37}. MC graphs
198 are easier to handle, but they omit challenging regions such as centromeres, which are hotspots

199 of structural variation¹⁸ and may play a crucial role in adaptive evolution and speciation^{70,71}.
200 Moreover, MC works on single chromosome graphs during graph construction and therefore
201 does not allow the representation of interchromosomal rearrangements³⁷, precluding, for
202 example, the investigation of acrocentric chromosome evolution⁷². We suggest choosing the
203 pipeline based on the desired analyses and the available computational resources. Both graphs
204 can be used for graph decomposition and the identification of variants between the genomes in
205 the graph; however, while MC should be used as a reference for resequencing projects, PGGB
206 is particularly useful when focusing on regions of interest. Overall, pangenome graphs face
207 conceptual and computational challenges and currently require significantly more resources for
208 their construction, storage, and analysis than linear genomes. While these limitations are being
209 tackled (e.g. through ‘implicit’ pangenome construction of only regions of interest⁷³),
210 researchers need to ensure access to sufficient computing resources, such as clusters or cloud
211 computing infrastructures.

212 **Improving the accessibility of biological information in pangenome graphs**

213 Pangenome graphs enclose extensive and complex biological information, including genomic
214 relationships and diversity among individuals. Their intricate and complex structure generates
215 large data volumes that are challenging to navigate and interpret. Graph manipulation toolkits
216 have been developed to improve graph accessibility to software used in downstream analyses,
217 thereby facilitating the extraction of biological information from pangenome graphs. Graph
218 manipulation includes tasks such as sorting, indexing, pruning and subsampling (**Fig. 2b**).
219 Sorting optimizes the order of graph nodes to reveal the underlying latent and sparse graph
220 structure and minimize errors in analysis and interpretation⁷⁴. Path indexing provides faster
221 access to specific regions of the graph, allowing software to quickly locate genes, variants, and
222 other features of interest without scanning the entire pangenome and reducing the time required
223 to retrieve relevant information for analyses such as read mapping and variant calling⁷⁵. To
224 further speed up computation, it is possible to simplify graph topology through pruning of
225 complex or unreliable regions, or by subsampling user-defined coordinates^{74,76}. Subsampling
226 is particularly helpful to disentangle the complexity of a particular region of interest, or when
227 computational resources are not available to query the entire graph. After graph construction
228 and each manipulation step, it is helpful to perform diagnostic statistics, such as graph size,
229 number of nodes and base content to get a sense of the structure of the pangenome and how
230 each step affected the graph^{74,76}. However, being a relatively new approach, pangenome graphs
231 lack universally accepted quality metrics due to their intrinsic complexity, the multiple

232 construction methods, and the lack of standardized benchmarking datasets. Nonetheless, as
233 pangenome graphs become more widely adopted, efforts will develop a unified set of metrics
234 applicable to all pangenome graphs.

235 Two main software packages exist to manipulate pangenome graphs: vg⁷⁶ and the pangenome
236 analysis toolkit ODGI⁷⁴. Vg relies on the .vg format⁷⁶ and was the first tool to be scaled up to
237 gigabase-scale graphs. ODGI operates on a node-centric object (.og), and was optimized for
238 pangenome graphs with hundreds of haplotype-resolved genomes⁷⁴. ODGI's tools work on a
239 graph-independent universal coordinate system that remains constant among different graphs
240 built from the same sequences⁷⁴. This system enables coordinate translation, facilitating the
241 lift-over of coordinate-based features between genomes and graphs⁷⁴, i.e., the accurate mapping
242 of annotated features (such as genes, regulatory elements, or other functional elements) from
243 one genome assembly to another⁷⁴. CAT (Comparative Annotation Toolkit)⁷⁷ can also annotate
244 the haplotypes in a pangenome graph by projecting the reference gene annotation to each of
245 the genomes, which can ease within-species annotation efforts³⁷. Feature annotations can also
246 be injected in the graph^{74,76} and used to interpret the functional significance of paths, nodes,
247 and edges.

248 **Visualizing genome diversity among individuals**

249 Graph visualization allows the inspection of homology relationships and variation between the
250 genomes, providing insights on the latent biological data⁷⁴. For instance, it can disentangle
251 variation at complex loci following haplotype paths along variation bubbles²¹. Visualization
252 can occur at different scales, from the overall structure down to the base-level (**Fig. 2c**). 2D
253 visualizations highlight graph structure and identify complex loci, while 1D visualizations help
254 understanding the graph topology and the relationships between genomes, potentially
255 providing a more immediate understanding of the complexity of a region with respect to
256 inspecting a list of variants. Various tools exist for pangenome graph visualization. Bandage⁷⁸
257 and GfaViz⁷⁹, originally created to visualize assembly graphs, permit the rendering of the 2D
258 graph layout and the interactive inspection of nodes and edges with variation that appear as
259 bubbles in the layout. Vg viz⁷⁶ can visualize nodes, edges, paths, and the base variation among
260 sequences. SequenceTubeMap⁸⁰ renders these elements in a 1D "tube map" model where paths
261 representing genomes navigate through the sequence nodes of the graph, oriented from left to
262 right. Read alignments and feature annotations injected in the graph can also be visualized. To
263 scale to gigabase pangenomes, such as the HPRC graph³⁷, MoMI-G⁸¹ combines the base-level
264 visualization of SequenceTubeMap⁸⁰ with CIRCOS⁸² plots chromosome-level connections to

265 efficiently browse SVs between genomes and aligned reads. ODGI⁷⁴ can render a raster image
266 of the graph topology in either 2D or 1D. Waragraph⁸³, an interactive implementation of ODGI,
267 is currently being developed to be able to inspect both 1D and 2D visualizations. When dealing
268 with a large graph, rendering the entire graph at once can become impractical and we
269 recommend visualizing chromosome graphs or subsampled regions of interest.

270

271 **Downstream analyses and their applications**

272 **Characterizing small variants and complex SVs through graph decomposition**

273 Variant sites (SNPs, indels, and SVs; **Fig. 3a**) in a pangenome graph can be extracted through
274 graph decomposition³⁸, the process of breaking down a pangenome graph into smaller, more
275 manageable subgraphs or components (snarls or bubbles)³⁸. Graph decomposition can be
276 performed with vg snarl³⁸ and gfatools bubble⁸⁴ (**Supplementary Table 2**). Vg deconstruct³⁷,
277 implemented in MC²⁰ and PGGB pipelines⁶⁴, can process the output of vg snarls or compute
278 snarls automatically generating a VCF with variants called from the references chosen during
279 graph construction with MC²⁰ or from any genome with PGGB⁶⁴. When working with large
280 graphs, it is recommended to compute snarls separately before variant calling³⁷. The
281 characterization of complex SVs, which were not previously accessible using linear reference
282 genomes and short reads, can shed light on their role in evolution⁸⁵ and in shaping phenotypic
283 variation, often with fitness consequences⁸⁶⁻⁸⁹, as SVs can affect fitness by altering gene
284 expression and shaping the chromosome recombination landscape⁹⁰. A complete representation
285 of SVs can also help analyze synteny and collinearity within genomes. In turn, this may provide
286 insights into chromosome evolution by encompassing the full complexity of sex chromosomes
287 and microchromosomes, which typically are enriched in SVs and challenging to resolve due to
288 high-repeat and GC content⁹¹. In addition, human pangenome graphs have also allowed the
289 identification of recombination events between heterologous acrocentric chromosomes,
290 especially at the breakpoints of Robertsonian translocations⁷². These translocations are the
291 most common chromosomal rearrangement in humans, and a comprehensive pangenome graph
292 has greatly enhanced the identification of the sequences and mechanisms involved⁷².

293 **Population genomics and alignment of transcriptomics data**

294 A pangenome graph can be used as a reference in resequencing projects to reduce mapping
295 bias (**Fig. 3b**; **Supplementary Table 2**). Short-reads map with greater confidence when more

296 genomic sequences are represented and known variation is embedded in the reference¹⁵.
297 However, read mapping to a pangenome is more challenging than mapping to a single reference
298 genome, as the search space for alignment increases due to the large number of potential paths
299 in the graph⁹². Since canonical algorithms cannot be applied directly to pangenome graphs¹⁵,
300 new tools have been developed for sequence-to-graph alignment. Within the vg toolkit, the
301 general-purpose read mapper `vg map`⁷⁶ is suitable for large and complex variation graphs, albeit
302 slower than popular linear-genome aligners with comparable accuracy⁹². `Vg Giraffe`⁹²,
303 currently being extended to support long-reads, uses a graph Burrows-Wheeler transform
304 (GBWT)⁹³, an indexing strategy that supports efficient querying and retrieval of sequences and
305 variants from the pangenome graph, to identify the paths that represent the two observed
306 haplotypes in an individual's reference sequences and to restrict the alignment space to these
307 regions only, avoiding biologically unlikely allele combinations. This leads to a dramatic
308 increase in mapping speed⁹². Long reads can also be aligned with `GraphAligner`⁹⁴, a seed-and-
309 extend sequence-to-graph aligner. Improved short-read mappability will benefit resequencing
310 projects, particularly in ancient DNA (aDNA) studies, which face challenges of contamination,
311 degradation, small amounts of endogenous DNA, shorter reads, and, therefore, lower
312 mappability¹⁴. aDNA mapped against a variation graph has already been proven to mitigate
313 reference biases by improving the allelic balance in polymorphic sites¹⁴. aDNA reads with non-
314 reference alleles maps in higher proportions to a graph containing alternate alleles, with respect
315 to a linear reference genome¹⁴.

316 `Vg` also allows splice-aware RNA-seq mapping to splice-aware graphs, generating an
317 alignment that can then be used to quantify haplotype-specific transcript expression^{68,95}.
318 Pantranscriptomics has the potential to efficiently quantify haplotype-specific differential gene
319 expression by exploiting the population variation that is embedded in the pantranscriptome
320 reference⁶⁸. We anticipate that pantranscriptomic sequencing projects combining RNA-seq
321 data with pangenome graph references will clarify the effects of gene flow⁹⁶, detecting adaptive
322 genetic variation^{97,98}. Chromatin accessibility analyses, such as chromatin
323 immunoprecipitation and sequencing (ChIP-seq) or assay for transposase-accessible chromatin
324 with sequencing (ATAC-seq), also benefit from a pangenomic approach³⁷. Their combination
325 with RNA-seq data provides a multi-omic approach that might facilitate the interpretation of
326 regulatory events critical to a wide variety of biological processes and phenotypes^{37,99}. These
327 approaches will enable future panepigenomics studies in non-model organisms, overcoming
328 current limitations in handling large, multi-omics data sets^{37,99}.

329 Detecting and genotyping variants in resequencing studies

330 Pangenome graphs can increase the accuracy of variant calling and genotyping in resequencing
331 studies thanks to improved read mappability^{21,26,29,37} (**Box 2; Fig. 3b; Supplementary Table**
332 **2**). Vg can be used to extract snarls from the graph and compute coverage and mapping quality
333 of aligned reads to accurately identify known variants¹⁰⁰. Larger known SVs (deletions,
334 insertions and inversions) can be genotyped by computing read coverage for each node and
335 edge¹⁰⁰. Specifically, variable sites identified by graph decomposition are assigned the two
336 most supported paths, representing haplotypes¹⁰⁰. For *de novo* small variant calling, the graph
337 is first “augmented” with variants identified through read alignment, after which read support
338 is computed¹⁰⁰. A potentially more accurate approach consists of projecting the graph
339 alignment from short-read mapping back into a linear BAM file (i.e., surjection) before using
340 traditional variant callers (e.g. Deepvariant¹⁰¹, Freebayes¹⁰² and GATK with Elprep^{103,104}) to
341 generate a VCF referenced to the genome chosen for surjection. This approach can be also used
342 with long reads mapped by GraphAligner⁹⁴. Given the complexity of PGGB graphs, read
343 mapping and variant calling have been mostly performed on MC graphs so far^{21,37}. Vg Giraffe
344 mapping followed by surjection and variant calling with Deepvariant is currently the state-of-
345 the-art approach and surjection was found as the main computational bottleneck^{21,37}. For the
346 chicken pangenome, the number of mapped and surjected reads per CPU-second dramatically
347 decreased when mapping against a PGGB graph with respect to a MC graph (1.6 reads versus
348 500 reads), while the memory usage increased (>250 Gb versus 24 Gb), making mapping and
349 surjection with a PGGB graph computationally infeasible²¹. An alternative and faster approach
350 for known variant genotyping that does not require read mapping is implemented in
351 PanGenie¹⁰⁵. This algorithm combines long-range haplotype information embedded in the
352 graph and *k*-mer counts from short-read data to jointly genotype SNPs, indels and SVs in the
353 uncharacterized sample (**Fig. 3c**). Haplotypes present in the graph can support genotype
354 assignment based on neighboring bubbles in case a given bubble is poorly covered by short-
355 read *k*-mers¹⁰⁵.

356 In population studies, pangenome-based variant calling increases accuracy and reduces the per-
357 sample data requirements, potentially expanding the size of assessable cohorts³⁷. Determining
358 accurate and comprehensive variant sets increases the resolution of the analysis of demographic
359 history, linkage-disequilibrium (LD), and genome-wide selection scans. This is particularly
360 beneficial in species with large effective population sizes, where LD is low^{29,106,107}. By
361 improving SVs genotyping, pangenome graphs can also help to integrate SVs into GWAS,

362 especially as long-reads gradually replace short reads in resequencing projects^{24,37,105}.
363 Performing GWAS on pangenome-based SNP and SV panels can therefore enhance our
364 understanding of the genetic basis of complex polygenic traits, and shed light on the role of
365 natural selection and gene-environment interactions and correlations.

366

367 **Conclusions and future prospects**

368 The last few years have seen a burgeoning of both small- and large-scale projects generating
369 high-quality reference genomes for biodiversity studies¹⁰⁸, including the Vertebrate Genomes
370 Project^{5,7}, the Darwin Tree of Life¹⁰⁹, and the European Reference Genome Atlas¹¹⁰. Most of
371 these projects contribute to the Earth Biogenome Project⁶, an ambitious proposal launched in
372 2020 to collectively sequence all named eukaryotic species within the next ten years. While
373 pangenome graphs are currently available only for a handful of species, recent advances in
374 genome and pangenome assembly potentially extend this approach to most eukaryotic species.
375 This is a desirable goal to reduce representation bias in all analyses of biodiversity, its
376 evolution, and conservation. Collecting, sequencing and assembling pangenomes from more
377 than a few individuals could be impractical in many species due to costs and sample
378 availability. In those cases, a pangenomes from single-to-few individuals would still increase
379 representation and reduce reference bias, especially for highly heterozygous populations where
380 a single individual may carry a high amount of allelic diversity. Pangenome graphs can benefit
381 a broad range of applications for biodiversity, from population genomics, phylogenomics,
382 hybridization and speciation studies, to conservation genomics, and will likely become the
383 standard reference system for such research in the future. Many new directions are being
384 investigated. For instance, panmitogenomes, i.e., pangenomes constructed from thousands of
385 mitochondrial genomes, have been shown to improve haplotyping of individuals¹¹¹ and are
386 being considered for species identification from heterogeneous samples. Another promising
387 new direction for the field is in super-pangenome graphs, which expand the survey of variation
388 to taxonomic ranks above species, opening new possibilities to study the molecular and
389 evolutionary mechanisms underlying species divergence, selection and recombination
390 processes, as well as adaptation to rapid climate changes²⁷. Dense sampling and sequencing of
391 species within a clade have proven essential for deciphering phylogenetic and phylogeographic
392 relationships, as well as facilitating investigation of gene loss and selection events¹¹². To this
393 end, a pangenomic approach revealed complex phylogenetic relationships among bacterial

394 strains, allowing genetic analyses of infectious diseases to identify virulence and antimicrobial
395 resistance genes with greater accuracy¹¹³. Despite the complexity of eukaryotic genomes, we
396 envision that rapid improvements in the efficiency and scalability of pangenomic tools will
397 soon allow such phylogenomic applications to be extended to eukaryotic species. In particular,
398 owing to the ability of super-pangenome graphs to incorporate all types of genomic variation,
399 they have the potential to elucidate complex evolutionary histories and phylogeographic
400 relationships of large, panmictic, and highly recombinant wild populations^{29,114}, as well as to
401 improve phylogenetic reconstructions of events, such as incomplete lineage sorting. Super-
402 pangenomes can also assist in studying biodiversity in complex ecosystems where
403 hybridization occurs^{49,115}. Hybrid zones are a prime opportunity for pinpointing the genes
404 responsible for phenotypic traits, as genes introgress in parallel with those traits¹¹⁶. Inclusion
405 of both hybridizing species in a pangenome graph will mitigate biases that arise from the use
406 of the reference genome of either species¹¹⁷. Pangenomes could also help shed light on the
407 origin of islands of divergence, highly differentiated genomic regions that might be related to
408 reproductive isolation and, thus, to speciation processes^{118,119}. Even within the same species, a
409 comprehensive pangenome graph that includes assemblies for all subspecies can maximize the
410 identification of structural genetic variants that are unique to a particular subspecies (**Box 2**).
411 We predict that species-level pangenomes will also replace linear genomes in phylogenomic
412 comparative genomics studies, enabled by the future development of tools for aligning
413 pangenomes of different species. Currently, the construction of a pangenome graph lacks
414 evolutionary information across the individuals, as phylogenetic divergence is not considered
415 in pairwise alignments, and this limitation should be taken into account when performing
416 phylogenetic analyses.

417 Pangenome graphs may also effectively guide conservation strategies aimed at maximizing the
418 preservation of genetic variation², by capturing a fuller spectrum of genetic diversity. Of
419 particular interest are structural and functional genomic variation involved in adaptations and
420 responses to environmental pressures. This will improve selection criteria for reintroducing and
421 translocating individuals among populations of threatened and endangered species. Improved
422 representation of structural elements, such as SVs, centromeres and telomeres, and CNVs, as
423 well as SNPs, along with non-coding regulatory elements, can provide comprehensive
424 conservation-relevant information regarding inbreeding, outbreeding, deleterious mutations,
425 introgression, and local adaptation². Pangenomes can also help to identify different genomic
426 regions in cryptic species, to then develop multilocus probes that distinguish cryptic taxa and
427 simplify conservation management¹²⁰. Moreover, we envision that pangenomics could help

428 reconstruct the genomic blueprint of extinct biodiversity by improving the mappability of
429 aDNA against a pangenome of a closely-related species. A more comprehensive comparison
430 between the extinct species and its living relative will help identify the genetic variation
431 underlying lost traits and ecosystem functions, which are essential information for any de-
432 extinction and restoration efforts^{14,121}. In conclusion, as methods to assemble, visualize,
433 annotate, and analyze pangenomes graphs continue to improve, we recommend researchers in
434 biodiversity genomics to embrace this new paradigm.

435

436 **Author contributions**

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

440 **Acknowledgments**

441 We are grateful to the HPRC community for the useful discussions over the years that help
442 shape the manuscript. We would like to particularly thank Erik Garrison for his input to the
443 manuscript. We thank Caterina Di Pietro for drawing the final figures. C.R.F. thanks the
444 support of CE3C through an assistant researcher contract (FCiência.ID contract #366) and FCT
445 (Fundação para a Ciência e a Tecnologia) for Portuguese National Funds attributed to CE3C
446 within the projects UIDB/00329/2020, UIDP/00329/2020, and LA/P/0121/2020, and FPUL for
447 a contract of invited assistant professor.

448

449 **Competing interests**

450

451 The authors declare no competing interests.

452

453

456 Box 1: History of the pangenome concept from bacteria to the human
457 pangenome

458 The origin of the pangenome concept traces back to 2005¹²², when the bacterial genomes of
459 *Streptococcus agalactiae* were first described as collections of core genes shared among
460 strains, dispensable ('accessory') genes shared between some strains, and strain-specific
461 (unique) genes. In *S. agalactiae*, the core genome included 80% of the genes, with the
462 remaining 20% categorized as accessory¹²². Work has since focused on the structure and
463 dynamics of bacterial pangenomes^{123–127}, proving fruitful for taxonomic identification¹²⁸, to
464 study host-pathogen interactions^{124,129} and gene families essential to pathogenicity¹²⁸ and
465 antibiotic resistance¹³⁰. These studies resulted in biomedical applications, with new promising
466 candidate drug targets¹³⁰ and reverse vaccinology^{131,132}. The concept was rapidly adopted by
467 plant and animal researchers, resulting in multiple eukaryotic pangenome studies⁴⁶. To
468 accommodate the complexity of eukaryotic genomes, with large portions (>~50%¹³³) of non-
469 coding and yet functional sequences, the term evolved to represent a complete set of sequences
470 found in all individuals of a population or a species⁴⁶. In plants, the pangenome concept was
471 first applied to the analysis of transposable elements, responsible for a large amount of variation
472 in both genic and intergenic sequences¹³⁴. Plant genomes are particularly dynamic as they
473 undergo frequent polyploidization and diploidization events^{135,136}. Moreover, intraspecific
474 genetic variability is often large¹³⁷. For these reasons, the concept of the pangenome has been
475 rapidly extended from initial efforts in crops, such as rice¹³⁸ and tomato¹³⁹, to many other plant
476 species^{140–142}, for association mapping analyses^{143,144}, breeding¹⁴⁵, crop improvement¹⁴⁶ and
477 evolutionary analyses¹⁴⁷. Very recently, a pangenome comprising 69 *A. thaliana* chromosome-
478 level assemblies from a global range distribution was published¹⁸. Animal pangenomes have
479 hitherto mostly focused on model species of economical value, both vertebrates and
480 invertebrates. Among mammals, pangenome assembly and analysis projects have been
481 conducted in the domestic pig (*Sus scrofa*)^{22,148,149}, cattle (*Bos taurus*)^{150,151}, and domestic
482 sheep (*Ovis aries*)¹⁵². Among invertebrates, studies on silkworm (*Bombyx mori*)⁴⁵,
483 Mediterranean mussel (*Mytilus galloprovincialis*)^{17,153} and longwing butterflies (three
484 *Heliconius* species)¹⁵⁴ pangenomes were recently published. The first step towards building a
485 human pangenome was taken with the assembly of an Asian and an African genome and their
486 integration within the NCBI reference human genome¹⁵⁵. This work resulted in the
487 identification of approximately 5 Mbp (mega base pairs) of novel sequences. From these early
488 analyses, the authors estimated that a comprehensive human pangenome should contain 19-40
489 Mb currently missing from the state-of-art reference assembly (NCBI Build 36.3)¹⁵⁵. Since
490 then, a Danish pangenome¹⁵⁶ has been assembled and the GenomeDenmark project initiated¹⁵⁷,
491 both utilizing family trios and revealing novel single-nucleotide and structural variants. More
492 recently, the assembly of an African pangenome¹⁵⁸ included approximately 10% more DNA
493 (296 Mbp) than the reference assembly (GRCh38), and a pangenome built from hundreds of
494 Han Chinese individuals¹⁵⁹ identified 29.5 Mbp of novel genomic sequences, including at least
495 188 novel protein-coding genes. The generation of a comprehensive catalog of human genetic

496 variation has advanced through the analysis of 338 high-quality assemblies of genetically
497 divergent populations¹⁴³, demonstrating that, for a given genome sequenced to 40-fold
498 coverage, over 400,000 previously unmapped reads, could now be aligned. The authors also
499 tested this new resource for more efficient mapping of previously discarded RNA-Seq reads¹⁶⁰.
500 Recently, the HPRC released the first draft human pangenome graph using three graph
501 construction methods and 47 phased, diploid assemblies from genetically diverse individuals³⁷.
502 Soon after, another human pangenome graph was built from 116 high-quality assemblies
503 representing 36 Chinese minority ethnic groups¹⁶¹. In parallel, Seven Bridges Genomics
504 Company generated population-specific genome graphs for a panAfrican genome¹⁶².

505

506 **Box 2: Case studies**

507

508 **The barn swallow pangenome**

509

510 The barn swallow (*Hirundo rustica*) is a small migratory songbird, with six subspecies that
511 differ in body size, extent and type of secondary sexual traits, and migratory behavior¹⁶³.
512 Latitudinal clines exist, with partial overlap in some traits¹⁶⁴ and variable levels of
513 hybridization between subspecies¹⁶⁵. In a first attempt to expand the characterization of the
514 genetic variation within this species, we generated a preliminary pangenome variation graph
515 for the Eurasian subspecies by using 12 haplotypes combined with MC²⁹. The resulting
516 pangenome increased the reference genome by approximately 500 Mbp (1.6 Gbp versus 1.1
517 Gbp). We were able to use the pangenome graph to infer core and accessory genes, and tested
518 it as a reference for read mapping and variant calling, further highlighting the potential of
519 pangenome graphs for population genomics²⁹.

520

521 **The grape super-pangenome**

522 The grapevine (*Vitis vinifera*) is one of the most important fruit crops globally, with great
523 economic and cultural value. Recently, the assembly of nine genomes from North American
524 wild *Vitis* species aimed to comprehensively characterize the genus diversity²⁴. The genomes
525 were scaffolded at the chromosome-level and fully phased. A super-pangenome was built using
526 PGGB⁶⁴, enabling access to intra- and inter-specific genetic variants and augmenting the
527 genome size threefold (1.7 Gbp versus 0.5 Mbp). The decomposition of variants embedded in
528 the graph captured valuable genetic variants, including those associated with flower sex
529 phenotype and disease resistance, thereby shedding light on species-specific adaptations. The
530 super-pangenome also supported a pan-genome-wide association study (panGWAS), and
531 identified variants near gene loci that are associated with chloride exclusion, potentially
532 influencing plant salt tolerance. These findings highlight the potential of pangenomes in
533 studying genetic variation and uncovering the genetic basis of functional traits, as well as
534 helping to address the challenges posed by climate change.

535 **Figure legends**

536 **Fig. 1: Principles of pangenome graphs**

537 a) An example of reference bias. Short reads from three different individuals (shown here is a
538 bird as an example) aligned to a linear reference genome do not map well to a missing region
539 or a region with divergence (shown in red). In contrast, a pangenome reference based on
540 multiple genomes (Hap.1, Hap.2, Hap.3; only one haplotype for each individual is represented
541 for simplicity) improves the coverage of such regions as less diverging regions are sampled,
542 thereby facilitating variant calling and subsequent downstream analyses. b) Pangenome graphs
543 can be constructed from unaligned raw reads (top) or from the alignment of multiple assembled
544 sequences (bottom). c) A pangenome graph is a bidirected graph with nodes representing DNA
545 sequences (semi-transparent squares with base pairs) connected by bidirectional edges, which
546 define the relationships between adjacent nodes and encode for the strandedness of the
547 sequences. The genomes walk as paths through the nodes (fine coloured lines above each node),
548 defining their base composition. Multiple genomes can share the same node sequence or take
549 different paths across bubbles, or snarls, which are subgraphs indicating the presence of
550 variation in that region. The first bubble represents an insertion in Hap.1, or deletion in Hap.2
551 and 3. The second bubble represents a SNP where Hap.1 and Hap.2 both have a G and Hap.3
552 has a C.

553

554 **Fig. 2: Pangenome graph construction, manipulation and visualization.**

555 a) General representation of sample selection, sequencing and assembly of haplotype-resolved
556 genomes (illustrated here are butterflies as an example). Samples should be collected in
557 different geographical locations to increase variability. Then, high-quality DNA is extracted,
558 sequenced with long-read technologies such as PacBio and Oxford Nanopore, and assembled
559 to obtain high-quality genomes. b) Schematic illustration of the two main methods for
560 pangenome graph construction, reference-based (MC) and reference-free (PGGB). MC starts
561 with Minigraph, which generates a SV-only graph, i.e. it starts with Hap.1 and adds the SVs
562 from Hap.3 (> 50 bp, highlighted in yellow). The graph construction is therefore influenced by
563 the order of the genomes aligned. Then, Cactus adds base-level information of all sequences to
564 the graph to also represent any SNPs. PGGB does not depend on the order of the aligned
565 genomes and starts from all-versus-all alignments generating a graph with loops representing
566 complex regions of the genome, such as the centromeres. After graph construction, typically
567 performed operations include sorting, indexing, pruning and subsampling in order to correct
568 the order of the nodes, create an index to make the pangenome elements accessible to other
569 software, prune complex and unreliable regions, and focus on regions of interest, respectively.
570 c) Pangenome graphs can be visualized in many ways to help interpretation. Is it possible to
571 visualize the overall structure of the graph in 2D (high-level visualization) focusing on the
572 relationships between nodes and edges rather than the base composition of the paths.
573 Alternatively it is possible to visualize the paths walking through the nodes together with their

574 base composition in 1D (base-level visualization). The latter is a tube-like representation of the
575 graph shown in Fig. 1c. Each coloured line is a different genomes (only one haplotype per
576 individual is shown for simplicity). Nodes report the DNA sequence (semi-transparent squares)
577 and variation is represented as divergence in the paths walking through the nodes.

578

579 **Fig. 3: Downstream analyses.** a) Identification of variants between the genomes included in
580 the pangenome. Sites of variation (bubble or snarls) are identified through graph decomposition
581 based on the graph topology, i.e. following the paths of the genomes across the nodes and
582 edges. Four different bubbles are represented in a tube-like manner as in Fig. 2c (base-level
583 visualization). By walking through the nodes of different genomes, it is possible to identify a
584 SNP, an indel (deletion/insertion), an inversion and a duplication, present here in Hap.3 with
585 respect to Hap.1 and Hap.2. Variants are called from the bubbles by looking at the path
586 divergence of a genome with respect to a chosen reference. b) Use of pangenome graphs for
587 read mapping and variant calling. Mapped reads (thin lines) after graph augmentation can be
588 visualized together with the pangenome graph reference (thick lines). This allows to call
589 variants already present in the graph, either a homozygous call, in which the mapped reads
590 share a G with Hap. 3, or novel variants (e.g. second heterozygous call, in which the read has
591 a C with respect to all the haplotypes in the pangenome graph reference). Variants are called
592 based on the divergence between the paths of the reads and the chosen reference genome. c)
593 Identification of structural variants and sample genotyping using PanGenie, which avoids read
594 mapping. Raw reads of a resequenced individual are divided into k -mers and assigned to nodes
595 in the pangenome graph. To assign a genotype, k -mer counts are compared between paths at
596 bubbles in the graph, which usually represent known haplotypes. Genomes are represented as
597 thick lines with different coloured outlines. In the first bubble, the raw reads of the individual
598 lack k -mers that match Hap.2 in that region, but have a 1-copy k -mer matching Hap.1, Hap.3
599 and Hap.4. Hence, the genotype inferred is Hap.1/Hap.4 (note that Hap.4 and Hap.3 share the
600 same sequence in the bubble, so here it could have also been Hap.1/Hap.3). In the second
601 bubble, the individual has a 2-copy k -mer matching Hap.2 and Hap.4. Therefore the inferred
602 genotype is Hap.2/Hap.4. In the first bubble, Hap.4 has been chosen during genotype inference
603 over Hap.3 as a possible divergent haplotype based on the neighboring bubble. In fact Hap.3
604 genotype is less likely because it is not supported by any k -mer in the second bubble.

605

606 **Glossary**

607

608 Bidirected graph: graph representing both DNA strands.

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

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

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

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

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

623 k-mer: substrings of nucleotides of length k

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

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

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

634 Hybridization: breeding between individuals from genetically different lineages

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

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

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

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

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

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

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

652 -

653 **References**

- 654 1. Giani, A. M., Gallo, G. R., Gianfranceschi, L. & Formenti, G. Long walk to genomics: History
655 and current approaches to genome sequencing and assembly. *Comput. Struct. Biotechnol. J.* **18**,
656 9–19 (2020).
- 657 2. Theissinger, K. *et al.* How genomics can help biodiversity conservation. *Trends Genet.* **0**, (2023).
- 658 3. Lewin, H. A. *et al.* Earth BioGenome Project: Sequencing life for the future of life. *Proc. Natl.*
659 *Acad. Sci. U. S. A.* **115**, 4325–4333 (2018).
- 660 4. Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de novo
661 assembly using phased assembly graphs with hifiasm. *Nat. Methods* **18**, 170–175 (2021).
- 662 5. Rhie, A. *et al.* Towards complete and error-free genome assemblies of all vertebrate species.
663 *Nature* **592**, 737–746 (2021).
- 664 6. Lewin, H. A. *et al.* The Earth BioGenome Project 2020: Starting the clock. *Proc. Natl. Acad. Sci.*
665 *U. S. A.* **119**, (2022).
- 666 7. Larivière, D. *et al.* Scalable, accessible and reproducible reference genome assembly and
667 evaluation in Galaxy. *Nat. Biotechnol.* **42**, 367–370 (2024).
- 668 8. Ghildiyal, K. *et al.* Genomic insights into the conservation of wild and domestic animal
669 diversity: A review. *Gene* **886**, 147719 (2023).
- 670 9. Blaxter, M. *et al.* Why sequence all eukaryotes? *Proc. Natl. Acad. Sci. U. S. A.* **119**, (2022).
- 671 10. Chaisson, M. J. P. *et al.* Multi-platform discovery of haplotype-resolved structural variation in
672 human genomes. *Nat. Commun.* **10**, 1784 (2019).
- 673 11. Wenger, A. M. *et al.* Accurate circular consensus long-read sequencing improves variant
674 detection and assembly of a human genome. *Nat. Biotechnol.* **37**, 1155–1162 (2019).
- 675 12. Zhao, X. *et al.* Expectations and blind spots for structural variation detection from long-read
676 assemblies and short-read genome sequencing technologies. *Am. J. Hum. Genet.* **108**, 919–928
677 (2021).
- 678 13. Günther, T. & Nettelblad, C. The presence and impact of reference bias on population genomic
679 studies of prehistoric human populations. *PLoS Genet.* **15**, e1008302 (2019).

- 680 14. Martiniano, R., Garrison, E., Jones, E. R., Manica, A. & Durbin, R. Removing reference bias and
681 improving indel calling in ancient DNA data analysis by mapping to a sequence variation graph.
682 *Genome Biol.* **21**, 250 (2020).
- 683 15. Eizenga, J. M. *et al.* Pangenome Graphs. *Annu. Rev. Genomics Hum. Genet.* **21**, 139–162 (2020).
- 684 16. Taylor, D. J. *et al.* Beyond the Human Genome Project: The Age of Complete Human Genome
685 Sequences and Pangenome References. *Annu. Rev. Genomics Hum. Genet.* (2024)
686 doi:10.1146/annurev-genom-021623-081639.
- 687 17. Gerdol, M. *et al.* Massive gene presence-absence variation shapes an open pan-genome in the
688 Mediterranean mussel. *Genome Biol.* **21**, 275 (2020).
- 689 18. Lian, Q. *et al.* A pan-genome of 69 *Arabidopsis thaliana* accessions reveals a conserved genome
690 structure throughout the global species range. *Nat. Genet.* **56**, 982–991 (2024).
- 691 19. Chin, C.-S. *et al.* A diploid assembly-based benchmark for variants in the major
692 histocompatibility complex. *Nat. Commun.* **11**, 4794 (2020).
- 693 20. Hickey, G. *et al.* Pangenome graph construction from genome alignments with Minigraph-
694 Cactus. *Nat. Biotechnol.* **42**, 663–673 (2024).
- 695 21. Rice, E. S. *et al.* A pangenome graph reference of 30 chicken genomes allows genotyping of
696 large and complex structural variants. *BMC Biol.* **21**, 267 (2023).
- 697 22. Jiang, Y.-F. *et al.* Pangenome obtained by long-read sequencing of 11 genomes reveal hidden
698 functional structural variants in pigs. *iScience* **26**, 106119 (2023).
- 699 23. Khan, A. W. *et al.* Super-Pangenome by Integrating the Wild Side of a Species for Accelerated
700 Crop Improvement. *Trends Plant Sci.* **25**, 148–158 (2020).
- 701 24. Cochetel, N. *et al.* A super-pangenome of the North American wild grape species. *Genome Biol.*
702 **24**, 290 (2023).
- 703 25. Li, N. *et al.* Super-pangenome analyses highlight genomic diversity and structural variation
704 across wild and cultivated tomato species. *Nat. Genet.* **55**, 852–860 (2023).
- 705 26. Leonard, A. S., Crysanto, D., Mapel, X. M., Bhati, M. & Pausch, H. Graph construction method
706 impacts variation representation and analyses in a bovine super-pangenome. *Genome Biol.* **24**,
707 124 (2023).

- 708 27. Shi, T. *et al.* The super-pangenome of *Populus* unveils genomic facets for its adaptation and
709 diversification in widespread forest trees. *Mol. Plant* (2024) doi:10.1016/j.molp.2024.03.009.
- 710 28. Theissinger, K. *et al.* How genomics can help biodiversity conservation. *Trends Genet.* (2023)
711 doi:10.1016/j.tig.2023.01.005.
- 712 29. Secomandi, S. *et al.* A chromosome-level reference genome and pangenome for barn swallow
713 population genomics. *Cell Rep.* **42**, 111992 (2023).
- 714 30. Rhie, A. *et al.* Towards complete and error-free genome assemblies of all vertebrate species.
715 *Nature* **592**, 737–746 (2021).
- 716 31. Pang, A. W. *et al.* Towards a comprehensive structural variation map of an individual human
717 genome. *Genome Biol.* **11**, R52 (2010).
- 718 32. Chakraborty, M., Emerson, J. J., Macdonald, S. J. & Long, A. D. Structural variants exhibit
719 widespread allelic heterogeneity and shape variation in complex traits. *Nat. Commun.* **10**, 4872
720 (2019).
- 721 33. Tigano, A. *et al.* Chromosome-Level Assembly of the Atlantic Silverside Genome Reveals
722 Extreme Levels of Sequence Diversity and Structural Genetic Variation. *Genome Biol. Evol.* **13**,
723 (2021).
- 724 34. Todesco, M. *et al.* Massive haplotypes underlie ecotypic differentiation in sunflowers. *Nature*
725 **584**, 602–607 (2020).
- 726 35. Garg, S., Balboa, R. & Kuja, J. Chromosome-scale haplotype-resolved pangenomics. *Trends*
727 *Genet.* **38**, 1103–1107 (2022).
- 728 36. Wang, S., Qian, Y.-Q., Zhao, R.-P., Chen, L.-L. & Song, J.-M. Graph-based pan-genomes:
729 increased opportunities in plant genomics. *J. Exp. Bot.* **74**, 24–39 (2023).
- 730 37. Liao, W.-W. *et al.* A draft human pangenome reference. *Nature* **617**, 312–324 (2023).
- 731 38. Paten, B. *et al.* Superbubbles, Ultrabubbles, and Cacti. *J. Comput. Biol.* **25**, 649–663 (2018).
- 732 39. Wang, T. *et al.* The Human Pangenome Project: a global resource to map genomic diversity.
733 *Nature* **604**, 437–446 (2022).
- 734 40. Miga, K. H. & Wang, T. The Need for a Human Pangenome Reference Sequence. *Annu. Rev.*
735 *Genomics Hum. Genet.* **22**, 81–102 (2021).

- 736 41. Read, B. A. *et al.* Pan genome of the phytoplankton *Emiliana* underpins its global distribution.
737 *Nature* **499**, 209–213 (2013).
- 738 42. Liu, Y. *et al.* Pan-Genome of Wild and Cultivated Soybeans. *Cell* **182**, 162–176.e13 (2020).
- 739 43. Talenti, A. *et al.* A cattle graph genome incorporating global breed diversity. *Nat. Commun.* **13**,
740 910 (2022).
- 741 44. Bozan, I. *et al.* Pangenome analyses reveal impact of transposable elements and ploidy on the
742 evolution of potato species. *Proc. Natl. Acad. Sci. U. S. A.* **120**, e2211117120 (2023).
- 743 45. Tong, X. *et al.* High-resolution silkworm pan-genome provides genetic insights into artificial
744 selection and ecological adaptation. *Nat. Commun.* **13**, 5619 (2022).
- 745 46. Golicz, A. A., Bayer, P. E., Bhalla, P. L., Batley, J. & Edwards, D. Pangenomics Comes of Age:
746 From Bacteria to Plant and Animal Applications. *Trends Genet.* **36**, 132–145 (2020).
- 747 47. Eberlein, C. *et al.* Hybridization is a recurrent evolutionary stimulus in wild yeast speciation.
748 *Nat. Commun.* **10**, 923 (2019).
- 749 48. Mavárez, J. *et al.* Speciation by hybridization in *Heliconius* butterflies. *Nature* **441**, 868–871
750 (2021).
- 751 49. Hübner, S. *et al.* Sunflower pan-genome analysis shows that hybridization altered gene content
752 and disease resistance. *Nat Plants* **5**, 54–62 (2019).
- 753 50. Nurk, S. *et al.* The complete sequence of a human genome. *Science* **376**, 44–53 (2022).
- 754 51. Koren, S. *et al.* De novo assembly of haplotype-resolved genomes with trio binning. *Nat.*
755 *Biotechnol.* (2018) doi:10.1038/nbt.4277.
- 756 52. Cheng, H. *et al.* Haplotype-resolved assembly of diploid genomes without parental data. *Nat.*
757 *Biotechnol.* **40**, 1332–1335 (2022).
- 758 53. Korlach, J. *et al.* Real-time DNA sequencing from single polymerase molecules. *Methods*
759 *Enzymol.* **472**, 431–455 (2010).
- 760 54. Mikheyev, A. S. & Tin, M. M. Y. A first look at the Oxford Nanopore MinION sequencer. *Mol.*
761 *Ecol. Resour.* **14**, 1097–1102 (2014).
- 762 55. Leonard, A. S. *et al.* Structural variant-based pangenome construction has low sensitivity to
763 variability of haplotype-resolved bovine assemblies. *Nat. Commun.* **13**, 3012 (2022).

- 764 56. Smith, T. *et al.* The first complete T2T Assemblies of Cattle and Sheep Y-Chromosomes
765 uncover remarkable divergence in structure and gene content. *Res Sq* (2024)
766 doi:10.21203/rs.3.rs-4033388/v1.
- 767 57. Beyter, D. *et al.* Long-read sequencing of 3,622 Icelanders provides insight into the role of
768 structural variants in human diseases and other traits. *Nat. Genet.* **53**, 779–786 (2021).
- 769 58. Kim, J. *et al.* The genome landscape of indigenous African cattle. *Genome Biol.* **18**, 34 (2017).
- 770 59. Kim, K. *et al.* The mosaic genome of indigenous African cattle as a unique genetic resource for
771 African pastoralism. *Nat. Genet.* **52**, 1099–1110 (2020).
- 772 60. Bakhtiari, M. *et al.* Variable number tandem repeats mediate the expression of proximal genes.
773 *Nat. Commun.* **12**, 2075 (2021).
- 774 61. Caballero-López, V., Lundberg, M., Sokolovskis, K. & Bensch, S. Transposable elements mark a
775 repeat-rich region associated with migratory phenotypes of willow warblers (*Phylloscopus*
776 *trochilus*). *Mol. Ecol.* **31**, 1128–1141 (2022).
- 777 62. Li, H., Feng, X. & Chu, C. The design and construction of reference pangenome graphs with
778 minigraph. *Genome Biol.* **21**, 265 (2020).
- 779 63. Garrison, E. *et al.* Building pangenome graphs. *bioRxiv* (2023) doi:10.1101/2023.04.05.535718.
- 780 64. Garrison, E. & Guarracino, A. Unbiased pangenome graphs. *Bioinformatics* **39**, (2023).
- 781 65. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100
782 (2018).
- 783 66. Armstrong, J. *et al.* Progressive Cactus is a multiple-genome aligner for the thousand-genome
784 era. *Nature* **587**, 246–251 (2020).
- 785 67. Guarracino, A., Mwaniki, N., Marco-Sola, S. & Garrison, E. *Wfmash: Base-Accurate DNA*
786 *Sequence Alignments Using WFA and mashmap2*. (Github).
- 787 68. Sibbesen, J. A. *et al.* Haplotype-aware pantranscriptome analyses using spliced pangenome
788 graphs. *Nat. Methods* **20**, 239–247 (2023).
- 789 69. Andreade, F., Lechat, P., Dufresne, Y. & Chikhi, R. Comparing methods for constructing and
790 representing human pangenome graphs. *Genome Biol.* **24**, 274 (2023).
- 791 70. Nergadze, S. G. *et al.* Birth, evolution, and transmission of satellite-free mammalian centromeric

- 792 domains. *Genome Res.* **28**, 789–799 (2018).
- 793 71. Ávila Robledillo, L. *et al.* Extraordinary Sequence Diversity and Promiscuity of Centromeric
794 Satellites in the Legume Tribe Fabaeae. *Mol. Biol. Evol.* **37**, 2341–2356 (2020).
- 795 72. Guarracino, A. *et al.* Recombination between heterologous human acrocentric chromosomes.
796 *Nature* **617**, 335–343 (2023).
- 797 73. Garrison, E., Guarracino, A. & Kille, B. *Impg: Implicit Pangenome Graph*. (Github).
- 798 74. Guarracino, A., Heumos, S., Nahnsen, S., Prins, P. & Garrison, E. ODGI: understanding
799 pangenome graphs. *Bioinformatics* **38**, 3319–3326 (2022).
- 800 75. Sirén, J. Indexing Variation Graphs. in *2017 Proceedings of the Meeting on Algorithm*
801 *Engineering and Experiments (ALENEX)* 13–27 (Society for Industrial and Applied
802 Mathematics, 2017).
- 803 76. Garrison, E. *et al.* Variation graph toolkit improves read mapping by representing genetic
804 variation in the reference. *Nat. Biotechnol.* **36**, 875–879 (2018).
- 805 77. Fiddes, I. T. *et al.* Comparative Annotation Toolkit (CAT)-simultaneous clade and personal
806 genome annotation. *Genome Res.* **28**, 1029–1038 (2018).
- 807 78. Wick, R. R., Schultz, M. B., Zobel, J. & Holt, K. E. Bandage: interactive visualization of de
808 novo genome assemblies. *Bioinformatics* **31**, 3350–3352 (2015).
- 809 79. Gonnella, G., Niehus, N. & Kurtz, S. GfaViz: flexible and interactive visualization of GFA
810 sequence graphs. *Bioinformatics* **35**, 2853–2855 (2019).
- 811 80. Beyer, W. *et al.* Sequence tube maps: making graph genomes intuitive to commuters.
812 *Bioinformatics* **35**, 5318–5320 (2019).
- 813 81. Yokoyama, T. T., Sakamoto, Y., Seki, M., Suzuki, Y. & Kasahara, M. MoMI-G: modular multi-
814 scale integrated genome graph browser. *BMC Bioinformatics* **20**, 548 (2019).
- 815 82. Krzywinski, M. *et al.* Circos: an information aesthetic for comparative genomics. *Genome Res.*
816 **19**, 1639–1645 (2009).
- 817 83. Fischer, C. *Waragraph*. (Github).
- 818 84. Li, H. gfatools: Tools for manipulating sequence graphs in the GFA and rGFA formats.
819 <https://github.com/lh3/gfatools>.

- 820 85. Sedlazeck, F. J., Lee, H., Darby, C. A. & Schatz, M. C. Piercing the dark matter: bioinformatics
821 of long-range sequencing and mapping. *Nat. Rev. Genet.* **19**, 329–346 (2018).
- 822 86. Noor, M. A., Grams, K. L., Bertucci, L. A. & Reiland, J. Chromosomal inversions and the
823 reproductive isolation of species. *Proc. Natl. Acad. Sci. U. S. A.* **98**, 12084–12088 (2001).
- 824 87. Küpper, C. *et al.* A supergene determines highly divergent male reproductive morphs in the ruff.
825 *Nat. Genet.* **48**, 79–83 (2016).
- 826 88. Weissensteiner, M. H. *et al.* Discovery and population genomics of structural variation in a
827 songbird genus. *Nat. Commun.* **11**, 3403 (2020).
- 828 89. Wold, J. *et al.* Expanding the conservation genomics toolbox: Incorporating structural variants to
829 enhance genomic studies for species of conservation concern. *Mol. Ecol.* **30**, 5949–5965 (2021).
- 830 90. Wellenreuther, M. & Bernatchez, L. Eco-Evolutionary Genomics of Chromosomal Inversions.
831 *Trends Ecol. Evol.* **33**, 427–440 (2018).
- 832 91. Peona, V. *et al.* The hidden structural variability in avian genomes. *bioRxiv* 2021.12.31.473444
833 (2022) doi:10.1101/2021.12.31.473444.
- 834 92. Sirén, J. *et al.* Pangenomics enables genotyping of known structural variants in 5202 diverse
835 genomes. *Science* **374**, abg8871 (2021).
- 836 93. Sirén, J., Garrison, E., Novak, A. M., Paten, B. & Durbin, R. Haplotype-aware graph indexes.
837 *Bioinformatics* **36**, 400–407 (2020).
- 838 94. Rautiainen, M. & Marschall, T. GraphAligner: rapid and versatile sequence-to-graph alignment.
839 *Genome Biol.* **21**, 253 (2020).
- 840 95. Wu, Z. *et al.* Human pangenome analysis of sequences missing from the reference genome
841 reveals their widespread evolutionary, phenotypic, and functional roles. *Nucleic Acids Res.* **52**,
842 2212–2230 (2024).
- 843 96. Tamagawa, K., Yoshida, K., Ohnishi, S. & Takahashi, Y. Population transcriptomics reveals the
844 effect of gene flow on the evolution of range limits. *Sci. Rep.* **12**, 1318 (2022).
- 845 97. DeBiasse, M. B., Kawji, Y. & Kelly, M. W. Phenotypic and transcriptomic responses to salinity
846 stress across genetically and geographically divergent *Tigriopus californicus* populations. *Mol.*
847 *Ecol.* **27**, 1621–1632 (2018).

- 848 98. Liu, L., Wang, Z., Su, Y. & Wang, T. Population transcriptomic sequencing reveals allopatric
849 divergence and local adaptation in *Pseudotaxus chienii* (Taxaceae). *BMC Genomics* **22**, 388
850 (2021).
- 851 99. Ma, S. & Zhang, Y. Profiling chromatin regulatory landscape: insights into the development of
852 ChIP-seq and ATAC-seq. *Mol Biomed* **1**, 9 (2020).
- 853 100. Hickey, G. *et al.* Genotyping structural variants in pangenome graphs using the vg toolkit.
854 *Genome Biol.* **21**, 35 (2020).
- 855 101. Poplin, R. *et al.* A universal SNP and small-indel variant caller using deep neural networks. *Nat.*
856 *Biotechnol.* **36**, 983–987 (2018).
- 857 102. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. *arXiv*
858 [*q-bio.GN*] (2012).
- 859 103. Van der Auwera, G. A. & O'Connor, B. D. *Genomics in the Cloud: Using Docker, GATK, and*
860 *WDL in Terra*. ('O'Reilly Media, Inc.', 2020).
- 861 104. Herzeel, C. *et al.* Multithreaded variant calling in elPrep 5. *PLoS One* **16**, e0244471 (2021).
- 862 105. Ebler, J. *et al.* Pangenome-based genome inference allows efficient and accurate genotyping
863 across a wide spectrum of variant classes. *Nat. Genet.* **54**, 518–525 (2022).
- 864 106. Li, M.-H. & Merilä, J. Population differences in levels of linkage disequilibrium in the wild.
865 *Mol. Ecol.* **20**, 2916–2928 (2011).
- 866 107. Charmantier, A., Garant, D. & Kruuk, L. E. B. *Quantitative Genetics in the Wild*. (Oxford
867 University Press, 2014).
- 868 108. Formenti, G. *et al.* The era of reference genomes in conservation genomics. *Trends Ecol. Evol.*
869 (2022) doi:10.1016/j.tree.2021.11.008.
- 870 109. Darwin Tree of Life Project Consortium. Sequence locally, think globally: The Darwin Tree of
871 Life Project. *Proc. Natl. Acad. Sci. U. S. A.* **119**, (2022).
- 872 110. Mazzoni, C. J., Ciofi, C. & Waterhouse, R. M. Biodiversity: an atlas of European reference
873 genomes. *Nature* **619**, 252 (2023).
- 874 111. Rubin, J. D., Vogel, N. A., Gopalakrishnan, S., Sackett, P. W. & Renaud, G. HaploCart: Human
875 mtDNA haplogroup classification using a pangenomic reference graph. *PLoS Comput. Biol.* **19**,

- 876 e1011148 (2023).
- 877 112. Feng, S. *et al.* Dense sampling of bird diversity increases power of comparative genomics.
878 *Nature* **587**, 252–257 (2020).
- 879 113. Yang, Z. *et al.* Pangenome graphs in infectious disease: a comprehensive genetic variation
880 analysis of *Neisseria meningitidis* leveraging Oxford Nanopore long reads. *Front. Genet.* **14**,
881 1225248 (2023).
- 882 114. Lombardo, G. *et al.* The Mitogenome Relationships and Phylogeography of Barn Swallows
883 (*Hirundo rustica*). *Mol. Biol. Evol.* **39**, (2022).
- 884 115. Payseur, B. A. & Rieseberg, L. H. A genomic perspective on hybridization and speciation. *Mol.*
885 *Ecol.* **25**, 2337–2360 (2016).
- 886 116. Hewitt, G. M. Hybrid zones-natural laboratories for evolutionary studies. *Trends Ecol. Evol.* **3**,
887 158–167 (1988).
- 888 117. Sebastianelli, M. *et al.* A genomic basis of vocal rhythm in birds. *Nat. Commun.* **15**, 3095
889 (2024).
- 890 118. Irwin, D. E. *et al.* A comparison of genomic islands of differentiation across three young avian
891 species pairs. *Mol. Ecol.* **27**, 4839–4855 (2018).
- 892 119. Hejase, H. A. *et al.* Genomic islands of differentiation in a rapid avian radiation have been driven
893 by recent selective sweeps. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 30554–30565 (2020).
- 894 120. van der Sprong, J. *et al.* A novel target-enriched multilocus assay for sponges (Porifera): Red Sea
895 Haplosclerida (Demospongiae) as a test case. *Mol. Ecol. Resour.* **24**, e13891 (2024).
- 896 121. Novak, B. J. De-Extinction. *Genes* **9**, (2018).
- 897 122. Tettelin, H. *et al.* Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*:
898 implications for the microbial ‘pan-genome’. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 13950–13955
899 (2005).
- 900 123. Obert, C. *et al.* Identification of a Candidate *Streptococcus pneumoniae* core genome and regions
901 of diversity correlated with invasive pneumococcal disease. *Infect. Immun.* **74**, 4766–4777
902 (2006).
- 903 124. Rasko, D. A. *et al.* The pangenome structure of *Escherichia coli*: comparative genomic analysis

- 904 of *E. coli* commensal and pathogenic isolates. *J. Bacteriol.* **190**, 6881–6893 (2008).
- 905 125. Donati, C. *et al.* Structure and dynamics of the pan-genome of *Streptococcus pneumoniae* and
906 closely related species. *Genome Biol.* **11**, R107 (2010).
- 907 126. Pinto, M. *et al.* Insights into the population structure and pan-genome of *Haemophilus*
908 *influenzae*. *Infect. Genet. Evol.* **67**, 126–135 (2019).
- 909 127. Aggarwal, S. K. *et al.* Pangenomics in Microbial and Crop Research: Progress, Applications, and
910 Perspectives. *Genes* **13**, (2022).
- 911 128. Rouli, L., Merhej, V., Fournier, P.-E. & Raoult, D. The bacterial pangenome as a new tool for
912 analysing pathogenic bacteria. *New Microbes New Infect* **7**, 72–85 (2015).
- 913 129. Kiu, R., Caim, S., Alexander, S., Pachori, P. & Hall, L. J. Probing Genomic Aspects of the
914 Multi-Host Pathogen *Clostridium perfringens* Reveals Significant Pangenome Diversity, and a
915 Diverse Array of Virulence Factors. *Front. Microbiol.* **8**, 2485 (2017).
- 916 130. Poulsen, B. E. *et al.* Defining the core essential genome of *Pseudomonas aeruginosa*. *Proc. Natl.*
917 *Acad. Sci. U. S. A.* **116**, 10072–10080 (2019).
- 918 131. Hisham, Y. & Ashhab, Y. Identification of Cross-Protective Potential Antigens against
919 Pathogenic *Brucella* spp. through Combining Pan-Genome Analysis with Reverse Vaccinology.
920 *J Immunol Res* **2018**, 1474517 (2018).
- 921 132. Naz, K. *et al.* PanRV: Pangenome-reverse vaccinology approach for identifications of potential
922 vaccine candidates in microbial pangenome. *BMC Bioinformatics* **20**, 123 (2019).
- 923 133. Francis, W. R. & Wörheide, G. Similar Ratios of Introns to Intergenic Sequence across Animal
924 Genomes. *Genome Biol. Evol.* **9**, 1582–1598 (2017).
- 925 134. Morgante, M., De Paoli, E. & Radovic, S. Transposable elements and the plant pan-genomes.
926 *Curr. Opin. Plant Biol.* **10**, 149–155 (2007).
- 927 135. Adams, K. L. & Wendel, J. F. Polyploidy and genome evolution in plants. *Curr. Opin. Plant*
928 *Biol.* **8**, 135–141 (2005).
- 929 136. Cheng, F. *et al.* Gene retention, fractionation and subgenome differences in polyploid plants. *Nat*
930 *Plants* **4**, 258–268 (2018).
- 931 137. Shi, J., Tian, Z., Lai, J. & Huang, X. Plant pan-genomics and its applications. *Mol. Plant* **16**,

- 932 168–186 (2023).
- 933 138. Sun, C. *et al.* RPAN: rice pan-genome browser for ~3000 rice genomes. *Nucleic Acids Res.* **45**,
934 597–605 (2017).
- 935 139. Gao, L. *et al.* The tomato pan-genome uncovers new genes and a rare allele regulating fruit
936 flavor. *Nat. Genet.* **51**, 1044–1051 (2019).
- 937 140. Jayakodi, M. *et al.* The barley pan-genome reveals the hidden legacy of mutation breeding.
938 *Nature* **588**, 284–289 (2020).
- 939 141. Song, J.-M. *et al.* Eight high-quality genomes reveal pan-genome architecture and ecotype
940 differentiation of *Brassica napus*. *Nat Plants* **6**, 34–45 (2020).
- 941 142. Li, J. *et al.* Cotton pan-genome retrieves the lost sequences and genes during domestication and
942 selection. *Genome Biol.* **22**, 119 (2021).
- 943 143. Tao, Y., Zhao, X., Mace, E., Henry, R. & Jordan, D. Exploring and Exploiting Pan-genomics for
944 Crop Improvement. *Mol. Plant* **12**, 156–169 (2019).
- 945 144. Qin, P. *et al.* Pan-genome analysis of 33 genetically diverse rice accessions reveals hidden
946 genomic variations. *Cell* **184**, 3542–3558.e16 (2021).
- 947 145. Zhou, Y. *et al.* Graph pangenome captures missing heritability and empowers tomato breeding.
948 *Nature* **606**, 527–534 (2022).
- 949 146. Zanini, S. F. *et al.* Pangenomics in crop improvement—from coding structural variations to
950 finding regulatory variants with pangenome graphs. *Plant Genome* **15**, e20177 (2022).
- 951 147. Qiao, Q. *et al.* Evolutionary history and pan-genome dynamics of strawberry (*Fragaria* spp.).
952 *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).
- 953 148. Tian, X. *et al.* Building a sequence map of the pig pan-genome from multiple de novo assemblies
954 and Hi-C data. *Sci. China Life Sci.* **63**, 750–763 (2020).
- 955 149. Li, Z. *et al.* The pig pangenome provides insights into the roles of coding structural variations in
956 genetic diversity and adaptation. *Genome Res.* **33**, 1833–1847 (2023).
- 957 150. Jang, J. *et al.* Chromosome-level genome assembly of Korean native cattle and pangenome graph
958 of 14 *Bos taurus* assemblies. *Sci Data* **10**, 560 (2023).
- 959 151. Dai, X. *et al.* A Chinese indicine pangenome reveals a wealth of novel structural variants

- 960 introgressed from other *Bos* species. *Genome Res.* **33**, 1284–1298 (2023).
- 961 152. Li, R. *et al.* A sheep pangenome reveals the spectrum of structural variations and their effects on
962 tail phenotypes. *Genome Res.* **33**, 463–477 (2023).
- 963 153. Saco, A. *et al.* Gene presence/absence variation in *Mytilus galloprovincialis* and its implications
964 in gene expression and adaptation. *iScience* **26**, 107827 (2023).
- 965 154. Ruggieri, A. A. *et al.* Erratum: A butterfly pan-genome reveals that a large amount of structural
966 variation underlies the evolution of chromatin accessibility. *Genome Res.* **32**, 2145 (2022).
- 967 155. Li, R. *et al.* Building the sequence map of the human pan-genome. *Nat. Biotechnol.* **28**, 57–63
968 (2010).
- 969 156. Besenbacher, S. *et al.* Novel variation and de novo mutation rates in population-wide de novo
970 assembled Danish trios. *Nat. Commun.* **6**, 5969 (2015).
- 971 157. Maretty, L. *et al.* Sequencing and de novo assembly of 150 genomes from Denmark as a
972 population reference. *Nature* **548**, 87–91 (2017).
- 973 158. Sherman, R. M. *et al.* Assembly of a pan-genome from deep sequencing of 910 humans of
974 African descent. *Nat. Genet.* **51**, 30–35 (2019).
- 975 159. Duan, Z. *et al.* HUPAN: a pan-genome analysis pipeline for human genomes. *Genome Biol.* **20**,
976 149 (2019).
- 977 160. Wong, K. H. Y. *et al.* Towards a reference genome that captures global genetic diversity. *Nat.*
978 *Commun.* **11**, 5482 (2020).
- 979 161. Gao, Y. *et al.* A pangenome reference of 36 Chinese populations. *Nature* **619**, 112–121 (2023).
- 980 162. Tetikol, H. S. *et al.* Pan-African genome demonstrates how population-specific genome graphs
981 improve high-throughput sequencing data analysis. *Nat. Commun.* **13**, 4384 (2022).
- 982 163. Spina, F. The EURING swallow project: a large-scale approach to the study and conservation of
983 a long-distance migrant. *Migrating birds know no boundaries Proceedings of the*.
- 984 164. M Ller, A. P. SEXUAL SELECTION IN THE BARN SWALLOW (*HIRUNDO RUSTICA*). IV.
985 PATTERNS OF FLUCTUATING ASYMMETRY AND SELECTION AGAINST
986 ASYMMETRY. *Evolution* **48**, 658–670 (1994).
- 987 165. Scordato, E. S. C. *et al.* Genomic variation across two barn swallow hybrid zones reveals traits

988 associated with divergence in sympatry and allopatry. *Mol. Ecol.* **26**, 5676–5691 (2017).

989