






# Graph construction method impacts variation representation and analyses in a bovine super-pangenome

## Journal Article

**Author(s):**

[Leonard, Alexander](#) ; [Crysnanto, Danang](#) ; [Mapel, Xena Marie](#) ; [Bhati, Meenu](#) ; [Pausch, Hubert](#) 

**Publication date:**

2023-05-22

**Permanent link:**

<https://doi.org/10.3929/ethz-b-000613247>

**Rights / license:**

[Creative Commons Attribution 4.0 International](#)

**Originally published in:**

Genome Biology 24(1), <https://doi.org/10.1186/s13059-023-02969-y>

RESEARCH

Open Access



# Graph construction method impacts variation representation and analyses in a bovine super-pangenome

Alexander S. Leonard<sup>1\*</sup>, Danang Crysnanto<sup>1</sup>, Xena M. Mapel<sup>1</sup>, Meenu Bhati<sup>1</sup> and Hubert Pausch<sup>1\*</sup> 

\*Correspondence:  
alexander.leonard@usys.ethz.ch;  
hubert.pausch@usys.ethz.ch

<sup>1</sup> Animal Genomics, ETH  
Zurich, Universitaetstrasse 2,  
8092 Zurich, Switzerland

## Abstract

**Background:** Several models and algorithms have been proposed to build pangenomes from multiple input assemblies, but their impact on variant representation, and consequently downstream analyses, is largely unknown.

**Results:** We create multi-species super-pangenomes using pggg, cactus, and mini-graph with the *Bos taurus taurus* reference sequence and eleven haplotype-resolved assemblies from taurine and indicine cattle, bison, yak, and gaur. We recover 221 k nonredundant structural variations (SVs) from the pangenomes, of which 135 k (61%) are common to all three. SVs derived from assembly-based calling show high agreement with the consensus calls from the pangenomes (96%), but validate only a small proportion of variations private to each graph. Pggg and cactus, which also incorporate base-level variation, have approximately 95% exact matches with assembly-derived small variant calls, which significantly improves the edit rate when realigning assemblies compared to minigraph. We use the three pangenomes to investigate 9566 variable number tandem repeats (VNTRs), finding 63% have identical predicted repeat counts in the three graphs, while minigraph can over or underestimate the count given its approximate coordinate system. We examine a highly variable VNTR locus and show that repeat unit copy number impacts the expression of proximal genes and non-coding RNA.

**Conclusions:** Our findings indicate good consensus between the three pangenome methods but also show their individual strengths and weaknesses that need to be considered when analysing different types of variants from multiple input assemblies.

**Keywords:** Bovinae, Graph pangenome, Long sequencing reads, Genome assembly, Structural variation, Domestic animals, VNTR profiling

## Background

Pangenomes store and represent sequences from multiple individuals and enable unbiased variation-aware sequence variant analyses [1]. Graph pangenomes represent alleles that differ between the input assemblies as nodes (representing sequence) connected by



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

edges [2]. Several methods have been proposed to construct graph pangenomes from genome-scale data. For instance, minigraph applies approximate mapping to construct structural variant-based pangenomes from multiple input assemblies [3]. The reference backbone has an impact on these pangenomes because it propagates bias when large segments are missing in the backbone assembly [4]. Recently, cactus [5] and pggg (pangenome graph builder) [6] have been proposed to construct pangenomes from multiple input assemblies using reference-free base-level alignment. The pggg and cactus pangenomes contain all types of differences found between the assemblies, ranging from single nucleotide to large structural differences with nested variation [2]. Other pangenome structures that are based on k-mers [7], maximal exact matches [8], or other algorithmic compression approaches [9] are efficient for querying sequences but lack a genomic coordinate system necessary for more general analyses.

Advancements in long-read sequencing and algorithms enable automated assembly of reference-quality genomes also for species with gigabase-sized genomes and so request for pangenomes that seamlessly accommodate and represent an increasing number of genomic resources [10, 11]. For instance, the Telomere-to-Telomere (T2T) consortium recently reported on a first complete human genome assembly [12] but routine near-T2T assembly [13] is becoming increasingly possible in humans [14] and other vertebrate species [15]. The Human Pangenome Reference Consortium (HPRC) coordinates global sequencing and assembly efforts with the goal to build a community-accepted variation-aware pangenome that integrates multiple input assemblies and represents global human genomic diversity [11]. These pangenomes have revealed homology in acrocentric chromosomes [16], trait variation resulting from variable number tandem repeats [17], and complex structural variations [18]. Pangenomes can leverage existing datasets like short reads to accurately genotype structural variants [19]. An initial draft human reference pangenome constructed from 47 phased diploid genome assemblies has recently been proposed by the HPRC [20], offering a stable, long-term replacement to the linear reference genome GRCh38 [21]. Pangenomes have also been adopted in many non-human species, revealing new clues to missing heritability in tomato [22], and evolutionary and adaptation insights into potato [23] and sheep [24], amongst others.

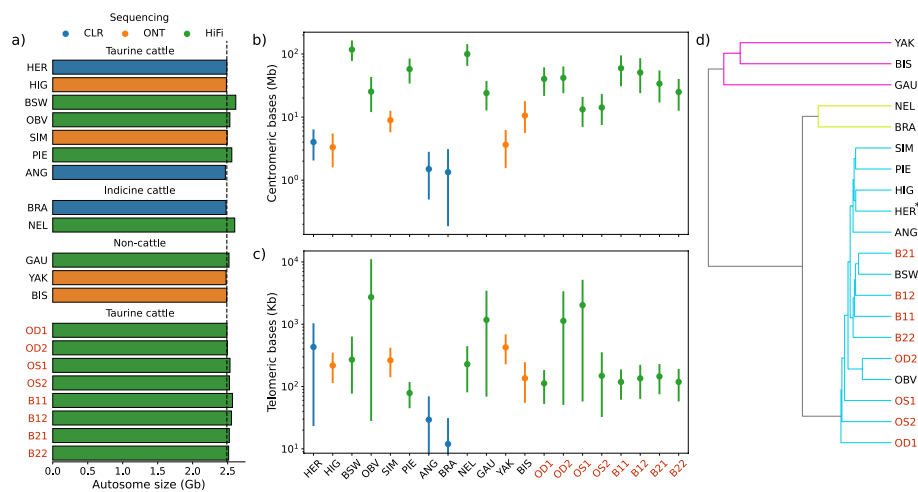
Likewise, a Bovine Pangenome Consortium (BPC) coordinates assembly efforts for the global cattle genomics community (<https://bovinepangenome.github.io/>). High nucleotide diversity and the separation of millions of global cattle into several hundred distinct breeds with unique genetic features, as well as frequent hybridization with their undomesticated relatives, make cattle an appealing species to improve assembly techniques [25] and investigate pangenome construction [4]. Bovine pangenomes constructed from multiple reference-quality assemblies revealed that the linear *Bos taurus taurus* reference sequence lacks millions of bases that are accessible in assemblies from other individuals [4, 10, 26]. However, the impacts of different construction methods on pangenome profiles, variant representation, and downstream analyses are more uncertain particularly when they include assemblies from multiple species.

Here, we apply cactus, minigraph, and pggg to build super-pangenomes with the *Bos taurus taurus* reference sequence and eleven haplotype-resolved assemblies from multiple species of the *Bos* genus including taurine and indicine cattle, bison, yak, and gaur. We assess the properties of the resulting pangenomes, recover large and small

variants, and investigate how the pangenomes represent different types of DNA variation. We then profile variable number tandem repeats (VNTR) to investigate how pangenomes integrate DNA variants that are challenging to resolve from linear alignments. Finally, we exploit phylogenetic relationships between the input assemblies to identify a highly polymorphic VNTR locus which mediates the expression of proximal genes and non-coding RNA.

## Results

We built pangenomes from autosomal sequences of domestic cattle and three of their wild relatives using minigraph, cactus, and pggp to assess how different graph pangenomes represent the same underlying input sequences. Nineteen haplotype-resolved assemblies from eight breeds of taurine (*Bos taurus taurus*) and indicine (*Bos taurus indicus*) cattle, yak (*Bos grunniens*), bison (*Bison bison bison*), and gaur (*Bos gaurus*) as well as the current Hereford-based *Bos taurus taurus* reference genome sequence were considered (Fig. 1a), of which the reference genome and eleven haplotype-resolved assemblies were used for pangenome construction and eight intermediate-quality assemblies were held out for analysis. These genomes were assembled using different sequencing and algorithmic approaches, but this has limited effect on the pangenomes [10]. However, the larger amount of centromeric/telomeric sequence in the HiFi-based assemblies (Fig. 1b, c) does require additional consideration.



**Fig. 1** Input assemblies considered for the pangenome analyses. **a** Autosomal length of the *Bos taurus taurus* reference sequence (HER, ARS-UCD1.2) and 19 haplotype resolved assemblies considered during pangenome construction and downstream analyses. Three pangenomes were created with ARS-UCD1.2 [HER] and eleven input assemblies (Brown Swiss [BSW], Piedmontese [PIE], Highland [HIG], Angus [ANG], Original Braunvieh [OBV], Simmental [SIM], Brahman [BRA], Nellore [NEL], gaur [GAU], bison [BIS], yak [YAK]), whereas eight additional Original Braunvieh or Brown Swiss assemblies indicated with red text (OD1, OD2, OS1, OS2, B11, B12, B21, B22) were only considered for downstream analyses. The colour of the bars indicates the primary sequencing technology used to construct the assemblies. The black dashed line indicates the length of ARS-UCD1.2. **b** Centromeric and **c** telomeric completeness is generally higher in the HiFi- and ONT- than CLR-based input assemblies. The marker is the sum over the autosomes and error bars indicate the 95% confidence interval. **d** A tree constructed with mash from the assemblies reveals the expected separation between taurine, indicine and non-cattle. The asterisk indicates the backbone genome used by minigraph

Minigraph used ARS-UCD1.2 as a backbone, whereas pggg and cactus are reference-free methods, although cactus is guided by an approximate phylogenetic tree (black subset of Fig. 1d).

### Pangenome construction and sequence content

The three pangenomes spanned between 427 k and 198 M nodes and contained between 2.6 and 3.0 gigabases (Table 1). Minigraph primarily incorporates larger (> 50 bp) DNA variation, and so builds a smaller graph compared to pggg and cactus which include all sizes of variation, which is also reflected in the number of path steps in the graph and final file sizes. The non-reference nodes contain nearly five times as much sequence in cactus (552 Mb) and pggg (523 Mb) than minigraph (109 Mb). Pangenome construction took 16 and 253 times more CPU hours and required 8 and 7 times more memory for cactus and pggg than minigraph.

Pggg and cactus also contain more repetitive sequence than minigraph, 46.4%, 44.3%, and 42.2% respectively, although this is largely due to including centromeric sequence. The ARS-UCD1.2 reference genome contains little centromeric sequence, preventing minigraph from anchoring centromeric sequence present in the HiFi-based (Brown Swiss, Nellore, Original Braunvieh, Piedmontese, and gaur) and to a lesser extent in ONT-based (bison and Simmental) assemblies into its pangenomes. Although pggg and cactus do include centromeric sequence into their pangenomes, they span a similar amount of bases as the sum of the input assemblies, suggesting either they are biologically distinct or cannot be effectively collapsed into a graph representation in a way similar to other repeat elements (e.g. SINE, LINE, LTR, etc., Additional file 1: Fig. S1).

### Consensus of genomic variation

To assess the commonality of genomic variation across pggg, cactus, and minigraph, as well as the suitability of their representation for downstream analyses, we decomposed and normalized the graph pangenomes into VCF files with respect to the ARS-UCD1.2 reference. Decomposing default minigraph pangenomes failed to output many expected structural variant (SV) alleles, particularly deletions in multi-node bubbles (Additional

**Table 1** Profiles of three bovine pangenomes constructed from autosomal sequences of twelve assemblies. Path steps are the total number of steps needed to trace all 12 assemblies through the graph

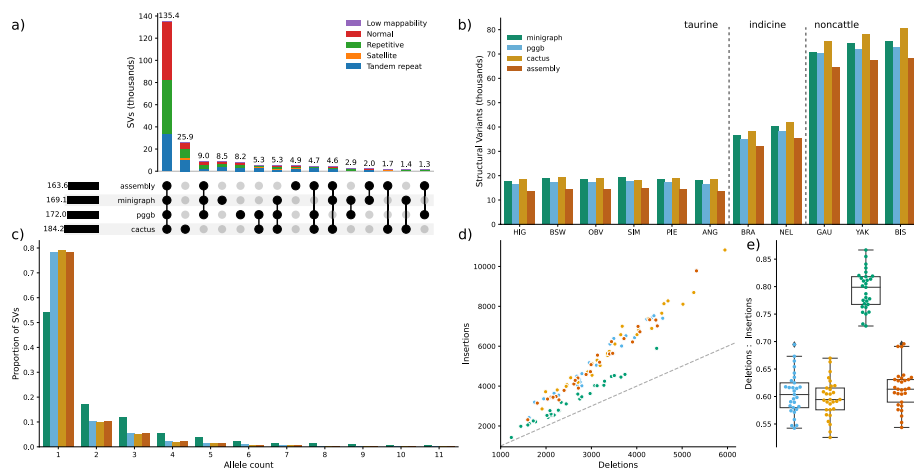
Parameter	Unit	minigraph	cactus	pggg
Nodes	N	427,012	198,431,246	179,575,371
Edges	N	606,926	272,102,708	245,150,846
Node length	bp	2,598,811,581	3,041,026,095	3,012,039,323
Path steps	N	3,358,976	1,621,936,527	1,442,793,659
Repetitive sequence	bp	1,107,501,421	1,361,489,638	1,415,552,890
Centromeric sequence	bp	2,939,789	291,982,193	255,091,362
CPU time	h	14 <sup>a</sup>	226 <sup>b</sup>	3,559
Max memory	GiB	7	54	46
GFA file size	GB	2.6	26.1	23.7

<sup>a</sup> Minigraph required an additional 18 CPU hours to add P-lines

<sup>b</sup> Cactus uses soft masked assemblies as input, which required an additional 1571 CPU hours

file 1: Fig. S2). By adding path information (P-lines in GFA) through minigraph-based realignment, we recovered an additional 21,770 SVs (13.5% increase). Realignment rarely produces paths incongruent with graph topology, where an assembly does not trace through nodes originating from that same assembly, affecting approximately 0.03%, 0.03%, and 0.05% of taurine, indicine, and non-cattle nodes (Additional file 1: Fig. S3). Cactus had the most SV genotypes marked as CONFLICT (2.2%), indicating that a haploid assembly had multiple possible ambiguous alleles, with minigraph and pggp significantly lower, 0.5% and 0.4% respectively. Such conflicts occurred more often in SV than small variation alleles, and were more frequent in divergent assemblies for cactus (Additional file 1: Fig. S4). VCF decomposition took 142 and 50 times longer for cactus and pggp compared to minigraph respectively, using approximately 26 and 22 times more memory (Additional file 2: Table S1), although pggp and cactus also contain substantial small variation to process.

Merging SVs, allowing for variation in breakpoints up to the minimum of 50% of the SV length and 1 Kb, identified 221 k nonredundant SVs  $\geq 50$  bp from the three pangenomes of which 135 k (61.2%) were common (Fig. 2a). Cactus contained 184 k SVs of which 26 k (14.1%) were private, i.e. not present in minigraph or pggp. We identified fewer SVs with pggp (172 k) and minigraph (163 k), of which 8 k (4.7%) and 9 k (5.0%) were respectively private. We classified the ARS-UCD1.2 reference genome into five regions: centromeric satellites (0.13%), tandem repeats (6.21%), repetitive elements (35.78%), low mappability (0.44%), and “normal” which encapsulated all remaining sequence (57.44%). The overlap between the three pangenomes was highest in normal (76.8%, 53.6 k) and repetitive regions (70.3%, 50.7 k), while tandem repeat regions (52.5%, 35.8 k) had lower overlap. The most challenging regions for alignment, low mappability (23.0%, 566) and especially centromeric satellites (2.1%, 81), had substantially lower overlap between the pangenomes. This also reaffirms that tandem repeats are disproportionately large contributors



**Fig. 2** Consensus of structural variation between three pangenomes. **a** Overlap of SVs between three pangenomes and assembly-based SV discovery. Overall bar height represents the total number of SVs. Variants in different classifications of genomic regions are indicated by colour of the stacked bar. **b** Number of SVs identified in the input assemblies through assembly-based mapping and the three pangenomes. **c** Allele count of the SVs (colours described in panel **b**). **d**, **e** Number of and ratio between deletions and insertions recovered from the pangenomes and assembly-based mapping (colours described in panel **b**)

to structural variation [17]. Applying stricter requirements to merge SVs across the pangenomes, requiring variant breakpoints to be within a 5 bp window, had limited effect, and even requiring exact basepair resolution had 92 k overlapping SVs, suggesting most SVs are precisely represented in the pangenomes (Additional file 1: Fig. S5).

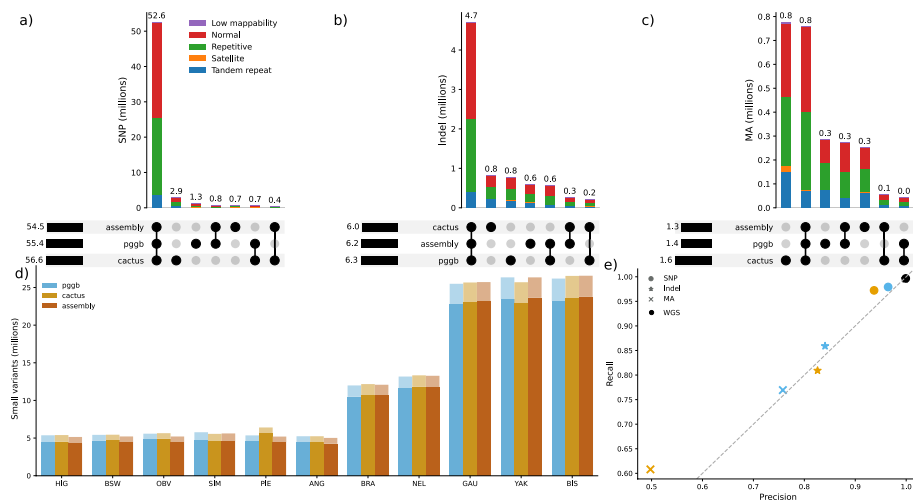
A benchmark dataset to validate SVs is not available for Bovinae, and so we assessed the SV representation accuracy in the three pangenomes comparing against SVs called directly from reference-alignment with the same set of assemblies. We identified between 13.6 k and 68.3 k SVs in the haplotype-resolved assemblies, with substantially more SVs in gaur, yak, and bison than the indicine and taurine haplotypes (Fig. 2b). More than three quarter (78.22%) of the SV alleles were observed in only one haplotype assembly. The eleven haplotype assemblies contained 163 k nonredundant SVs of which 151 k (92.30%), 150 k (91.90%) and 146 k (89.45%) were also recovered with minigraph, pggg and cactus, respectively. This approach validated the vast majority ( $N=135$  k, 96.26%) of the 141 k SVs that were found in all three pangenomes. However, it validated only 6.12%, 14.06% and 19.24% of the SVs that were respectively private to cactus, pggg and minigraph. Overall, the SV F-score for each minigraph, pggg, and cactus against the assembly truth was 0.908, 0.896, and 0.842 respectively, with normal and repetitive regions outperforming low mappability and satellite regions (Additional file 1: Fig. S6). Similar alleles are collapsed more frequently in minigraph than pggg and cactus, and so the allele frequency spectrum differs between the three tools with minigraph containing substantially less singleton SVs (Fig. 2c). While all tools recovered more insertions than deletions, minigraph contained proportionally less insertions than pggg and cactus (Fig. 2d, e).

We used optical mapping from two Nellore samples [27] unrelated to the NEL haplotype and from the BRA and ANG haplotypes [28], as well as ONT long reads from the OBV, BSW, PIE, NEL, and GAU haplotypes [10] to further validate SV consensus in the pangenomes. Since optical maps call larger SVs (typically minimum SV size of 1 Kb) than sequencing reads, and the ONT datasets had different read lengths, quality, and coverage, the overlap with the pangenome SVs was not expected to be complete. Both sets of orthogonal data substantiate that minigraph and pggg best represent SVs with nearly equal support (1355 and 1283 pangenome SVs overlapping optical map SVs, and ONT SV F-score of 0.799 and 0.802, respectively), while cactus had comparatively lower support (1045 optical map SV overlaps and ONT F-score of 0.743) (Additional file 3: Table S2).

We further assessed the commonality of small variation (SNPs and indels < 50 bp) in pggg and cactus, as minigraph primarily only represents SVs. We recovered 68.27 M nonredundant small variants from the two pangenomes. Pggg and cactus contained 63.02 M (55.37 M SNPs/7.65 M indels) and 64.27 M (56.61 M SNPs/7.66 M indels) small variants of which 59.02 M (53.28 M SNPs/5.74 M indels) were common. The overlap between pggg and cactus was substantially larger for SNPs than indels and multiallelic variants (Fig. 3a–c).

We assessed small variation representation accuracy in the two pangenomes again using assembly-derived calls as an approximate truth set. The eleven input assemblies contained 62.04 M (54.55 M SNPs / 7.49 M indels) nonredundant small variations. Each assembly had between 5.12 and 27.05 M small variations, again finding





**Fig. 3** Overlap of SNPs and indels between pangenome- and assembly-based discovery. **a–c** Overlap of small variations (SNPs and indels smaller than 50 bp) between pggp, cactus, and assembly-based discovery presented for **a** biallelic SNPs, **b** biallelic indels, and **c** multiallelic SNPs/indels. Variants in different classifications of genomic regions are indicated by colour of the stacked bar. **d** Number of small variations recovered from each haplotype from the pangenome and the assembly-based mapping. The faded and dark areas of the bars represent indels and SNPs, respectively. **e** Precision and recall of pggp and cactus for SNPs, indels, and multiallelic variants assuming the assembly-based calls are truth. The black WGS point represents gold-standard accuracy for 30 × sequencing coverage

substantially fewer in the taurine assemblies than their indicine and non-cattle counterparts (Fig. 3d). Pggp and cactus contained 96.3% and 94.8% of the assembly-based small variants, and 94.8% and 91.5% of the pangenome variation respectively were found in the assembly calls (Fig. 3e). Pggp had a higher overall F-score than cactus, 0.96 and 0.93, respectively, suggesting it encapsulates small variation more accurately. As observed for SVs, accuracies were highest in normal and repetitive regions, although pggp and cactus had F-scores of 0.93 and 0.89, respectively, in tandem repeat regions, indicating a strong ability to resolve repeat motif variation which may only differ by a single or a few bases (Additional file 1: Fig. S6).

We also simulated assemblies for chromosome 29 by introducing known variation per sample to the ARS-UCD1.2 reference genome, removing uncertainty arising from variation that was not called from the real assemblies when aligned to ARS-UCD1.2. Pangenomes constructed from these simulated assemblies had minor improvements to precision and recall for minigraph, pggp, and cactus, as well as for all types of variants (mean F1 improvement for SNP: 0.014, Indel: 0.027, multiallelic: 0.023, SV: 0.039, Additional file 1: Fig. S7) compared to pangenomes constructed from the real assemblies. The slight improvement suggests that hard-to-call regions or centromeric sequence (not included in simulated assemblies) hurt pangenome accuracy, but that there is still some loss of accuracy in pangenome construction or downstream conversion to VCF even with simulated assemblies.

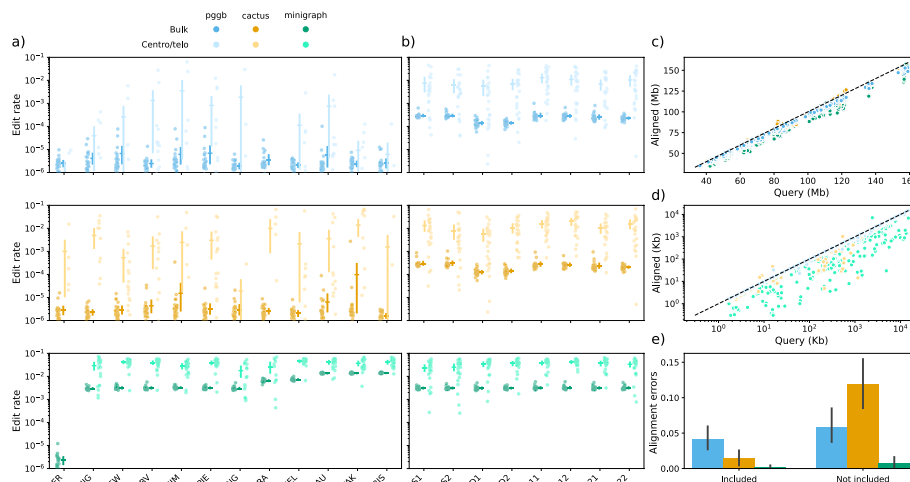


### Alignment of assemblies to pangenomes

We realigned all twelve input assemblies (including ARS-UCD1.2) against the pangenomes to calculate edit distances and quantify sequence content of the pangenomes. Since centromeric, and to a lesser extent telomeric, sequence is challenging to align, we analysed these “centro-/telomeric” regions separately to the rest of the “bulk” sequence (Additional file 4: Table S3). Minigraph had a substantially higher edit rate (0.494%) than pgggb (0.010%) and cactus (0.010%) in bulk regions and also in centro-/telomeric regions, 3.43%, 0.396%, and 0.719%, respectively (Fig. 4a). Even though cactus includes more centromeric sequence in the graph (Additional file 1: Fig. S1) and had overall comparable edit rates, its centro-/telomeric edit rates were nearly double that of pgggb. Minigraph also had a significantly lower edit rate for ARS-UCD1.2, as all reference sequence is included by definition of minigraph’s algorithm. Aligning to only the linear reference ARS-UCD1.2 resulted in higher edit rates (bulk: 0.527%, centro/telo: 3.56%) and lower query coverage compared to pangenomic alignments.

Edit rates were higher for assemblies more diverged from the reference backbone in minigraph pangenomes, indicating they do not incorporate highly diverged segments containing multiple small variations but were not large enough to form bubbles. Contrastingly, pgggb and cactus pangenomes did not have any divergence bias for edit rate, suggesting they may be more suited to include more divergent assemblies into super-pangenomes [29].

We also aligned eight additional assemblies held out from the pangenomes ( $2 \times$  haplotypes from each parent of the OBV and  $2 \times$  haplotypes for two unrelated Brown Swiss cattle (Fig. 1)), and found highly similar edit rates for minigraph, but significantly higher



**Fig. 4** Realignment of included and held-out assemblies to the pangenomes. **a** Edit rate of the twelve assemblies used in pangenome construction. Each faded dot represents one autosome of each assembly in either the bulk or centro-/telomeric ranges. Mean and 95% confidence intervals are also indicated for each category of sequence per assembly. **b** Similar plot to **a** but for eight additional assemblies held-out from pangenome construction. **c** Query coverage for the bulk sequence for each autosome of all assembly. Values above the dashed line indicate there were more aligned bases than query bases, suggesting multiple ambiguous alignments that were equally scored. **d** Similar to **c**, but for the centro-/telomeric sequence. **e** Bar plots of for the number of alignment errors reported by GraphAligner when using finite “tangle effort”. Bar heights reflect the mean across 348 (included) and 232 (not included) autosomes across the assemblies, while error bars represent 95% confidence intervals

edit rates for pggb and cactus (Fig. 4b). This is expected and reflects true genomic variation in the held-out assemblies not present in the pggb/cactus pangenomes, while minigraph never captured small variation anyway and had higher edit rates to begin with. The OD1/2 (dam of OBV sample) assemblies had a lower edit rate in pggb and cactus, as expected since the OBV haplotype is the maternal haplotype.

In addition to substantially lower edit rates, pggb and cactus also covered more of the bulk query sequence (98.4% and 98.6%, respectively) compared to minigraph (93.8%). The differences were more pronounced in the centro-/telomeric query sequence alignments (98.0%, 93.9%, and 65.8%, respectively), as expected given the relative lack of centromeric sequence in minigraph pangenomes (Fig. 4c). Some chromosomes (e.g. 6, 14, and 16) had multiple ambiguous but equally scored alignments in cactus pangenomes, leading to >100% query coverage. More query sequences can be aligned with more sensitive alignment parameters, but this requires >300 GB of memory per chromosome in cactus graphs (Additional file 5: Table S4). Even with relaxed alignment parameters, aligning to cactus and pggb pangenomes took 2.2 and 1.4 times more CPU time and 37 and 7 times more memory than to minigraph pangenomes (Additional file 6: Table S5), and was especially pronounced for aligning assemblies not included in the pangenomes. Failed alignment of 500 Kb segments was more common in pggb and cactus than minigraph, generally resulting from exceeding GraphAligner's "tangle effort" in highly complex regions (Fig. 4e). This was especially apparent when aligning not included assemblies to the cactus pangenomes.

### Pangenome resolution of VNTR

Variable number tandem repeats (VNTR) account for substantial gene expression and complex trait variation, but due to their repetitive nature and high mutation rates, these loci are difficult to resolve from linear alignments [30, 31]. A catalogue of bovine VNTR had not been established and so we identified 9568 tandem repeats (TRs) in non-masked regions of the ARS-UCD1.2 reference sequence (Additional file 7: Table S6) and investigated their prevalence and variability in the other assemblies through the three pangenomes.

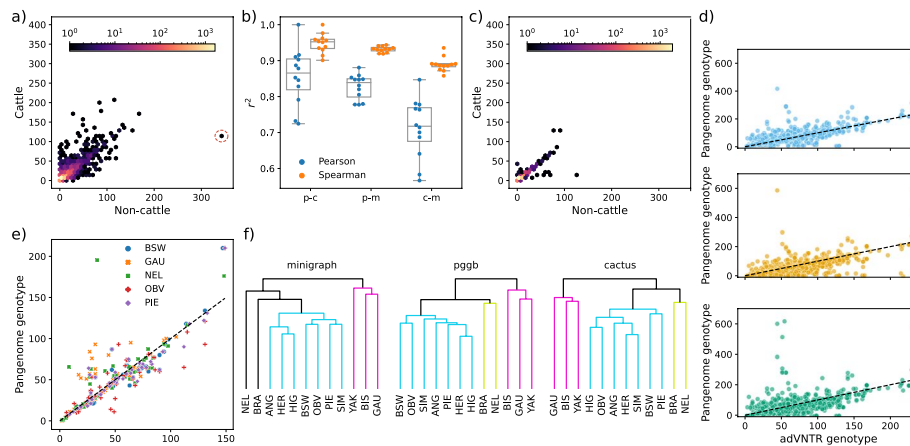
Pggb and cactus contain all assembly sequences, and so can arbitrarily convert between pangenome position and assembly coordinate with tools like odgi. As such, we can liftover the TR positions directly into each assembly through the pangenomes. On the other hand, minigraph contains information on assembly coordinates at graph bubbles, and so we can only estimate TR positions that effectively overlap with SVs. The

**Table 2** All TRs were examined from pggb and cactus, while TRs overlapping SVs were examined for minigraph. Genotypable TRs are TRs where all 12 assemblies had a known path through the local graph structure. VNTRs are genotypable TRs that had at least one sample with a different number of TR counts. CPU and Memory indicate the compute resources needed to identify the paths of all assemblies through all examined TRs

Pangenome	TRs examined	Genotypable TRs	VNTRs	CPU (h)	Memory (Gb)
Pggb	9566	8939 (93.5%)	7854 (87.9%)	18.17	7.5
Cactus	9566	8910 (93.1%)	7859 (88.2%)	24.17	7.7
Minigraph	5731	5504 (96.5%)	5457 (99.1%)	0.19	0.5

former approach is substantially more complete and accurate, but incurs a much larger compute cost (Table 2). Approximately 95% of TRs had all twelve assemblies associated with a pangenome path, while the remaining TRs suffered from reduced mappability (Table 2). All pangenomes had instances where coordinates were erroneously translated, but this was most apparent in cactus with several huge outliers (Additional file 8: Table S7). Nearly all TRs examined by minigraph were variable between the assemblies, while approximately 12% of TRs examined by pggg and cactus had zero variation between samples (Table 2).

Tandem repeat counts were similar between the three pangenomes, with 5293 commonly genotyped VNTRs. There were 3332 (63%) VNTRs with identical counts across all twelve assemblies in all three pangenomes and 4084 (77%) identical in at least two pangenomes (Fig. 5a). The remaining VNTR counts still broadly agreed with minor variability (Additional file 1: Fig. S8). The average squared Spearman correlation was 0.92 between the pangenomes, with several outliers in cactus and minigraph skewing the squared Pearson correlation, which was otherwise around an average of 0.8 (Fig. 5b, all significant). TRs present in only pggg and cactus had low variability between cattle and non-cattle (Fig. 5c), suggesting minigraph efficiently captures nearly all VNTRs of interest for further investigation. We also genotyped 465 TRs with advNTR [31] using HiFi reads from the gaur, Nellore, Piedmontese, Brown Swiss, and Original Braunvieh samples as an approximate truth set, finding good concordance (Fig. 5d, Additional file 7: Table S6), with a median difference of counts of 7, 8, and 11 for pggg, cactus, and minigraph to advNTR. Minigraph occasionally over- or underestimated TR counts if the overlapping SV was significantly larger or smaller, as that determines the translated



**Fig. 5** VNTR concordance in three pangenomes. **a** TRs with identical counts in at least two pangenomes, with the median count for the cattle and non-cattle groups. A particular VNTR with substantially more repeats in non-cattle compared to cattle that we investigated further is circled in red. **b** Pearson and Spearman squared correlation coefficients across the TR counts for pggg-cactus (p-c), pggg-minigraph (p-m), and cactus-minigraph (c-m). Each point is one assembly, with box plots over the 12 assemblies. **c** Similar to **a**, except TRs with identical counts in pggg and cactus that were not present in minigraph. **d** advNTR-derived genotypes for five HiFi samples in the three pangenomes. The black dashed line indicates the expected count using advNTR as a ground truth. **e** VNTRs where all three pangenomes agreed to a different count than advNTR, suggesting advNTR may sometimes over-/underestimate assembly-based counts. **f** Trees derived from the TR counts across different input assemblies, with colours representing clusters of taurine cattle, indicine cattle, and non-cattle

coordinates. There were also instances where all three pangenomes agreed on counts different to adVNTR, indicating even HiFi reads may not be powerful enough to genotype all VNTRs (Fig. 5e). Clustering based on TR counts produced trees that broadly grouped the taurine and indicine cattle as well as the non-cattle (Fig. 5f), although the minigraph and cactus trees had slight inconsistencies with the topology of the mash-derived tree (Fig. 1d).

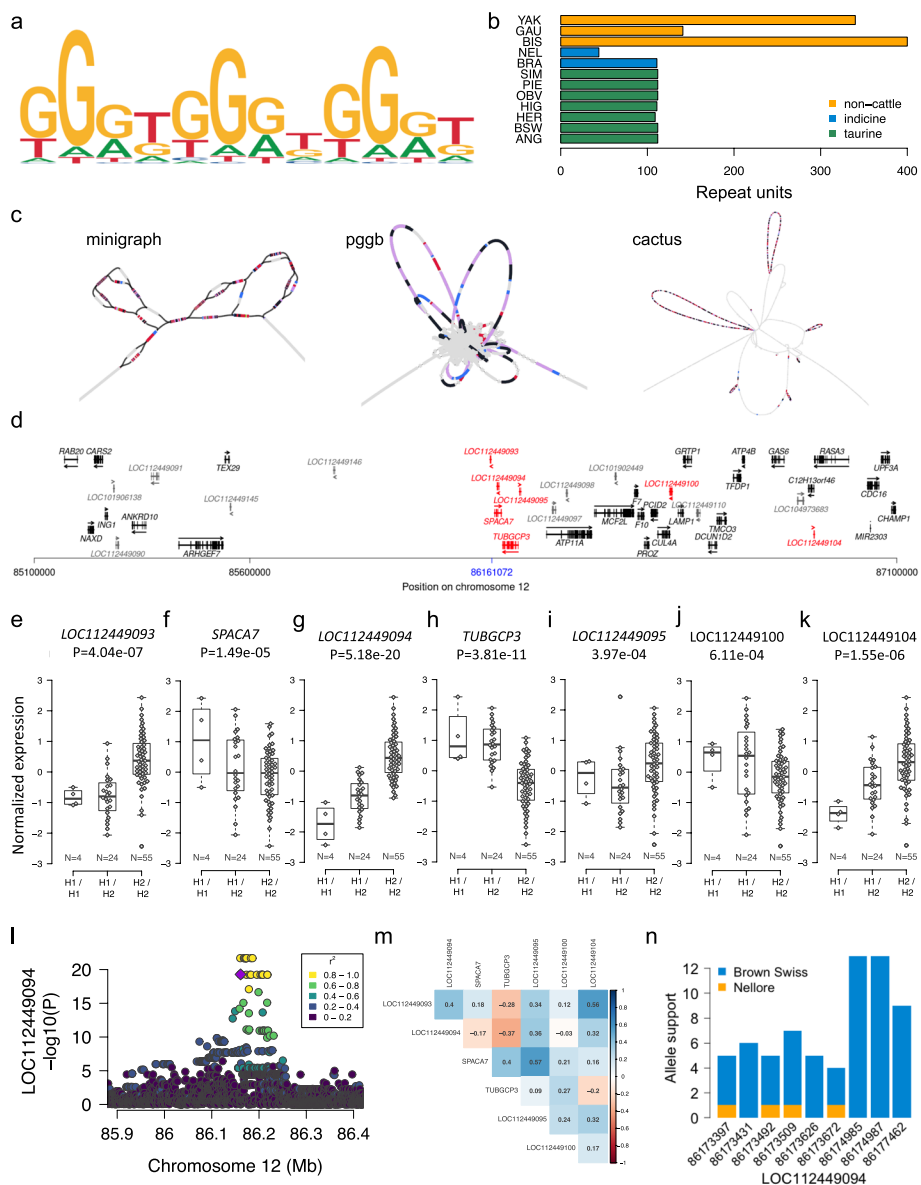
### An eVNTR mediates expression of neighbouring genes and non-coding RNA

The pangenomes identified a highly variable VNTR locus on chromosome 12 (86,161,072–86,162,351 bp), containing substantially more copies of a degenerate 12 bp motif in non-cattle and fewer in Nellore than in the taurine assemblies (Fig. 6a, b), prompting a detailed investigation. Although the genotyping was similar, bandage plots revealed significantly different graph structures of the VNTR across the three pangenomes (Fig. 6c). Minigraph added additional nodes to incorporate the (primarily non-cattle) additional tandem repeats, while pggp, and cactus to a lesser extent, incorporated the additional repeats by looping through the nodes containing the tandem repeat sequence multiple times (Additional file 1: Fig. S9). This different representation was observed generally in VNTRs with high repeat counts (Additional file 1: Fig. S10).

We observed two VNTR “haplotypes” (referred to as hap1 and hap2) in long-read alignments of additional Original Braunvieh cattle demonstrating within-breed variability. While the alignments indicated several insertions and deletions in both haplotypes with respect to ARS-UCD1.2, hap2 is 24 bp longer than hap1 due to two additional copies of the repeat motif (Additional file 1: Fig. S11). There are 63 genes annotated within the *cis*-regulatory range ( $\pm 1$  Mb) of the VNTR, of which 44 (Fig. 6d) were expressed at  $>0.2$  TPM in testis tissues of 83 mature Brown Swiss and Original Braunvieh bulls.

A degenerate repeat motif and an overall length  $>1300$  bp precluded the short sequencing read-based genotyping of the VNTR in the eQTL cohort with adVNTR. Instead, we genotyped the VNTR through two SNPs (Chr12:86,160,984 and Chr12:86,160,971) and one indel (Chr12:86,161,000) that tagged the two VNTR haplotypes. This enabled us to investigate putative *cis*- and *trans*-regulatory impacts of the VNTR on the expression of genes and long non-coding RNAs (lncRNA). Two genes (*TUBGCP3*, *SPACA7*) and five non-coding RNAs (*LOC112449093*, *LOC112449094*, *LOC112449095*, *LOC112449100*, *LOC112449104*) within  $\pm 1$  Mb of the VNTR were differentially expressed ( $P < 1.13 \times 10^{-3}$ , Bonferroni-corrected significance threshold) between the diplotypes indicating a putative *cis*-regulatory role (Fig. 6e–k). We did not detect *trans*-regulatory effects for the haplotype (Additional file 1: Fig. S12).

We mapped eQTL within the *cis*-regulatory range ( $\pm 1$  Mb of the transcription start site) of the two significant genes and five significant non-coding RNAs to investigate if variants in *cis* other than the VNTR haplotype are associated with transcript abundance. The VNTR haplotype was amongst the top variants at the *LOC112449094* eQTL, but eleven variants in strong linkage disequilibrium ( $r^2 = 0.93$ ) were more significant ( $P = 2.00 \times 10^{-22}$  vs.  $P = 5.18 \times 10^{-20}$ ) (Fig. 6l). The eQTL peak was absent when the VNTR haplotype was fitted as a covariate. We made similar observations for the *TUBGCP3* eQTL (Additional file 1: Fig. S13). Variants other than the VNTR haplotype were more strongly associated with the expression of *SPACA7*, *LOC112449093*,



**Fig. 6** A polymorphic eVNTR discovered from the pangenomes. **a** VNTR repeat motif and **b** number of repeat units in the twelve input assemblies. **c** Bandage plots of a VNTR upstream *SPACA7*. BLAST hits for the top 5 most common VNTR motifs are coloured per motif. **d** Genes (black) and (long) non-coding RNAs (dark grey) nearby the VNTR. Blue colour indicates the position of the VNTR. Arrows indicate the orientation of genes and lncRNA. Red colour indicates two genes and five lncRNAs whose expression is associated with the VNTR. **e–k** Expression (quantified in transcripts per million (TPM)) of the associated genes and ncRNAs in testis tissues of 83 Brown Swiss bulls that are either homozygous for hap1 (H1/H1), homozygous for hap2 (H2/H2) or heterozygous (H1/H2). The number of animals per diplotype is below the boxplots. **l** Cis-expression QTL mapping for *LOC112449094*. Different colours indicate the pairwise linkage disequilibrium ( $r^2$ ) between the VNTR haplotype (violet) and all other variants. **m** Correlation between the abundance of mRNA and lncRNA for the associated genes and lncRNA. **n** Allelic imbalance of nine heterozygous exonic SNPs in *LOC112449094* in testis tissue from a Nellore x Brown Swiss crossbred bull. Orange and blue colours represent paternal (Nellore) and maternal alleles (Brown Swiss)

*LOC112449095*, *LOC112449100*, and *LOC112449104* (Additional file 1: Fig. S13). These findings suggest that the VNTR haplotype primarily mediates the expression of *LOC112449094* and *TUBGCP3*.

Hap2 containing more copies of the 12 bp motif increases abundance ( $\beta_{\text{normalized}} = 1.26 \pm 0.10$ ,  $P = 5.18 \times 10^{-20}$ ) of *LOC112449094* which is lowly-expressed ( $0.86 \pm 0.31$  TPM) in 83 bull testis transcriptomes. Conversely, hap2 is associated ( $\beta_{\text{normalized}} = -1.08 \pm 0.14$ ,  $P = 3.81 \times 10^{-11}$ ) with lower *TUBGCP3* mRNA. Negative correlation ( $-0.37$ , Fig. 6m) between *LOC112449094* and *TUBGCP3* abundance suggests that *LOC112449094* could be a cis-acting lncRNA that represses expression of *TUBGCP3*. The VNTR haplotype is in strong LD with two SNPs (Chr12:86,173,431 ( $r^2 = 0.87$ ) and Chr12:86,173,509 ( $r^2 = 0.93$ )) in *LOC112449094* exons. While hap2 primarily segregates with the alternate alleles of these two SNPs, hap1 segregates with the respective reference alleles. Alternate allele support in RNA sequencing reads overlapping Chr12:86,173,431 and Chr12:86,173,509, respectively, was 71% (96 out of 135 alleles,  $P_{\text{binom}} = 1.01 \times 10^{-6}$ ) and 69% (94 out of 137 alleles,  $P_{\text{binom}} = 1.57 \times 10^{-5}$ ) in 22 animals that are heterozygous both at the VNTR haplotype and the SNPs confirming allelic imbalance due to lower *LOC112449094* expression in hap1.

We confirmed a putative cis-regulatory role of the VNTR in an individual with indicine ancestry. Expression of *LOC112449094* was relatively low (TPM = 0.66) in the testis tissue sampled from a Nellore (*Bos taurus indicus*) x Brown Swiss (*Bos taurus taurus*) crossbred bull (NxB), although it inherited hap2 from its Brown Swiss dam (Additional file 1: Fig. S11). The paternal (Nellore) haplotype is diverged from hap2 and hap1 and contains less VNTR repeat units (Fig. 6b). Alignment of parental-binned HiFi reads against ARS-UCD1.2 readily distinguished paternal from maternal alleles at nine heterozygous SNPs overlapping *LOC112449094* exons. Inspection of these variants in the RNA sequencing alignments of NxB F1 showed allelic imbalance due to a disproportionately low number or even complete absence of paternal (=Nellore) alleles (Fig. 6n) suggesting that the low VNTR repeat count of the Nellore haplotype represses expression of *LOC112449094*.

## Discussion

We constructed pangenomes using three different algorithmic approaches with minigraph, pggp, and cactus, finding each has different strengths and weaknesses. Pggp and cactus include all variations, down to SNPs, while minigraph primarily captures SVs. Consequently, pggp and cactus pangenomes had roughly 450 times more nodes and edges, 10 times the file size, and construction required 253 and 16 times as much CPU time as the minigraph pangenomes respectively, despite only including about 5 times as much non-reference sequence. High-quality assemblies are becoming increasingly easier and cheaper to produce, and so expanding pangenomes with additional assemblies will likely be common. Minigraph can iteratively add assemblies in near-linear time and memory to an existing pangenome. Cactus requires storage of modestly sized intermediate files (10 s of GB) and the complexity of adding a new assembly depends on the existing phylogenetic guide tree [5]. Pggp on the other hand, performs all-versus-all alignment, and so adding new assemblies would require building a new pangenome from scratch. As such, pggp may be suited to building reference pangenomes or periodic updates, but cactus and minigraph are more appropriate for projects with ongoing assembly efforts. Building the pangenomes per-chromosome mitigates the quadratic increase in runtime for all-versus-all alignments but neglects inter-chromosomal



connections. Building whole genome pangenomes may be necessary when investigating complex chromosomal rearrangements (e.g. Robertsonian translocations) or centromeric homology [16], but such events were not present in the assemblies assessed in this work (Additional file 1: Fig. S14).

All three pangenomes contained more SVs than the assembly-based “truth”, with private SVs in minigraph (8.5 k), pggg (8.1 k), and cactus (25.9 k). Recent work in human pangenomes suggest many of these putative errors are actually SVs missed by linear reference-based calling [20], although the substantially larger number in cactus likely indicates many true errors. However, the overall SV agreement was good, with minigraph, pggg, and cactus respectively having F-scores of 0.91, 0.90, and 0.84. Minigraph and pggg furthermore had the best overlap with orthogonal optical mapping data, supporting many of their SVs are true variation. We found that approximately 0.5% and 3.9% of alleles in pggg and cactus were removed during preliminary filtering when only keeping alleles seen in genotypes, suggesting variants may either exist in complex regions that cannot be genotyped confidently into VCF or strongly tangled regions that produce spurious variant calls. However, we found a higher overlap on small variations compared to structural variations, indicating that the cactus (F1: 0.93) and pggg (F1: 0.96) pangenomes can integrate small variations more accurately. These results demonstrate that SVs and small variation encoded in pangenome graphs, particularly cactus, would benefit from careful filtering before qualifying for downstream use in population genomic analyses.

Pggg and cactus are reference-free, while minigraph requires an initial backbone to determine variation from the input assemblies. Minigraph is thus susceptible to reference-bias, especially in the case of an incomplete reference (currently almost every reference genome except the T2T-CHM13 [12]). The ARS-UCD1.2 autosomes are interrupted by 257 gaps and contain almost no centromeric/telomeric sequence, which prevents integration of centromeric/telomeric sequence from the more complete HiFi assemblies into the minigraph pangenome. Pggg and cactus are largely able to integrate these sequences, with 255 and 292 Mb of centromeric sequence in the non-reference nodes of the pangenomes, respectively. However, cactus does use progressive alignments on the provided guide tree, and so does appear to have a slight reference-bias when projecting back into a reference coordinate system (like VCF). Assemblies more distant to the ARS-UCD1.2 branch tended to have higher genotype conflict rates, which is not observed in pggg.

Pangenomes with base-level variation like pggg and cactus can realign most of the input assemblies with minimal edit rate, but centromeric/telomeric regions are still poorly realigned. Equally, the base-level variation in these pangenomes can lead to strongly tangled graph topology, which can generate poor realignments and substantially increase computational requirements. For super-pangenomes incorporating multiple (sub-) species, containing a greater amount of variation, downstream software like Bandage, GraphAligner, and vg quickly hit CPU and memory bottlenecks working with pggg or cactus, while minigraph pangenomes are usable even without HPC infrastructure. Although beyond its design intentions, minigraph can be run with a lower “minimum variant length” parameter, increasing its sensitivity to smaller variation down to approximately 10 bp without substantially impacting compute requirements for pangenome



construction or analyses (Additional file 9: Table S8). However, this unsupported use-case cannot regularly capture SNPs and so pggb and cactus are vastly better for representing small variation accurately.

We show that the three pangenomes are generally concordant at highly polymorphic VNTR loci, although minigraph can over- or underestimate TR counts if the graph bubbles are significantly larger or smaller than the TR region, as minigraph primarily translates assembly coordinates at these bubbles. All three pangenomes had several poorly genotyped VNTRs, due to incorrect coordinate extraction in these repetitive regions, again indicating unresolved challenges with complex alignments. Furthermore, minigraph primarily captures repeat count variation rather than motif variation in TRs, the latter of which may be of greater importance in identifying eVNTRs [32]. However, the small variation present in pggb and cactus hinders visual inspection and pangenome annotation with Bandage, while TRs are simple to identify in SV-oriented minigraph pangenomes. We also show that even highly accurate long reads can fail to successfully genotype VNTRs, supporting that assembly-based approaches are the most powerful resource for detecting large or complex variation. In the case of a highly variable eVNTR upstream *SPACA7*, adVNTR predicted 116 copies for gaur, whereas both pggb and minigraph pangenomes identified 141 copies. Cactus erroneously aligned several assemblies in this VNTR region, and so only predicted 53 copies for gaur but correctly predicted the taurine cattle count. Although the eVNTR analysis partially relied on conventional linear alignment approaches, recent advances are enabling purely pangenomic approaches that can simplify and improve analyses [33].

## Conclusions

Bovine pangenomes have been utilized before to make non-reference sequences amenable to association testing and reveal trait-associated structural variants [4, 10, 26]. Our findings show good agreement between the SVs discovered from three widely used pangenome methods and so reinforce the utility of any pangenome approach to investigate variants that are difficult to resolve with a single linear reference. Pggb, and to a lesser extent cactus, pangenomes losslessly represent the input assemblies and are ideal for generating a pangenome reference containing all variation. However, minigraph has much greater utility, allowing simple expansion with additional assemblies and rapid downstream analyses with modest computational requirements. With the recent establishment of the cross-species Bovine Pangenome Consortium (<https://bovinepangenome.github.io/>) and similar efforts in plants [22, 23, 29], there is a need to emphasise pangenomes that can handle significantly higher levels of variation compared to the draft human pangenome reference [20].

## Methods

### Pangenome construction

Pangenomes were constructed per-chromosome for the 29 bovine autosomes, using the Hereford-based *Bos taurus taurus* ARS-UCD1.2 reference genome [34] and eleven other haplotype-resolved assemblies: eight domestic cattle (*Bos taurus taurus* and *Bos taurus indicus*) and three wild relatives of domestic cattle (*Bos grunniens*, *Bos gaurus*, *Bison*

**Table 3** Input assemblies for the pangenomes. The bovine reference (Hereford, ARS-UCD1.2) is a primary assembly, a collapsed haploid representation of the diploid genome. All other assemblies are haplotype-resolved, a haploid genome of either the maternal or paternal haplotype

Breed/species	Acronym	Primary sequencing strategy	Autosome length (Gb)	Reference to publication	Reference to assembly
Hereford/ <i>Bos taurus taurus</i>	UCD	CLR	2.489	[34]	[36]
Angus/ <i>Bos taurus taurus</i>	ANG	CLR	2.468	[25]	[37]
Brahman/ <i>Bos taurus indicus</i>	BRA	CLR	2.478	[25]	[38]
Highland/ <i>Bos taurus taurus</i>	HIG	ONT	2.483	[39]	[40]
Yak/ <i>Bos grunniens</i>	YAK	ONT	2.478	[39]	[41]
Simmental/ <i>Bos taurus taurus</i>	SIM	ONT	2.494	[42]	[43]
Bison/ <i>Bison bison bison</i>	BIS	ONT	2.488	[44]	[45]
Brown Swiss/ <i>Bos taurus taurus</i>	BSW	HiFi	2.617	[10]	[46]
Nellore/ <i>Bos taurus indicus</i>	NEL	HiFi	2.602	[10]	[46]
Piedmontese/ <i>Bos taurus taurus</i>	PIE	HiFi	2.560	[10]	[46]
Gaur/ <i>Bos gaurus</i>	GAU	HiFi	2.520	[10]	[46]
Original Braunvieh/ <i>Bos taurus taurus</i>	OBV	HiFi	2.532	[10]	[46]

CLR Pacific Biosciences Continuous Long Reads, HiFi Pacific Biosciences High Fidelity Reads, ONT Oxford Nanopore Technologies Reads

*bison*) (Table 3). We used these assemblies to construct pangenomes with minigraph (v0.20) [3], the PanGenome Graph Builder (pvgb, v0.5.2) [6, 35], and the Cactus Progressive pangenome pipeline (cactus v2.3.0) [5].

The minigraph pangenome was constructed using base-level alignment (‘-c’), 2% divergence level, and default parameters otherwise. The ARS-UCD1.2 reference genome was the backbone of the graph, and the other assemblies were added in order of their mash distance [47] to the reference: Highland, Brown Swiss, Original Braunvieh, Simmental, Piedmontese, Angus, Brahman, Nellore, gaur, yak, and bison.

Pvgb was run with the recommended parameters (segment length of 100,000 bp, identity mapping taken from the rounded lowest mash similarity of 98%, and target haplotype paths of 12) with all assemblies combined into a single fasta file.

Since our assemblies contain multiple species with significant divergence, we used cactus with progressive alignment using the guide tree from the mash-based distance. Assemblies were first soft-masked with the rush job mode of RepeatMasker (<http://www.repeatmasker.org>) (version 4.1.4) and rmbblast (version 2.13.0), using a database of repetitive DNA elements from Repbase (release 20181026). The cactus hierarchical alignment (HAL) was converted into GFA using hal2vg (<https://github.com/ComparativeGenomicsToolkit/hal2vg>) and vg convert [48].

### Analysis of pangenome sequence content

We extracted all fasta sequence in the graphs with odgi flatten (v0.8) [35], and then repeat masked the output as described above.

### Variation discovery from pangenomes

Pangenomes were decomposed into VCF files using vg deconstruct with --all-snarls --path-traversals --ploidy 1. We added path information (P-lines) to minigraph’s gfa

by manually curating the output of `minigraph -call` on each sample which retraces the assembly's path (available in the Github repository associated with this paper). Variants were then normalized using `vcfwave` [49], skipping variants larger than 1 Mb, and then split into small (< 50 bp) and structural variation (> 50 bp) using `bcftools norm` and `bcftools view`. Alleles not observed in any sample were dropped using `--trim-alt-alleles` during splitting, and similarly applied to per-sample VCF statistics. Multiple nucleotide polymorphisms (MNPs) were split into multiple SNPs using `bcftools norm --atomize`.

#### Variation discovery from assemblies

All non-reference assemblies were aligned to ARS-UCD1.2 using `minimap2` (v2.24 [50]) with `"-cx <preset>"`. The divergence-based preset was chosen as `asm5`, `asm10`, and `asm20` for the taurine, indicine, and non-cattle, respectively. Variants were called from the alignments using `paftools.js` and normalized and subsetted using the same approach as pangenome variation.

#### Variation discovery from ONT long reads

ONT long reads from the OBV, BSW, PIE, NEL, and GAU haplotypes were obtained from [10] to further validate SV. All samples were aligned to ARS-UCD1.2 using `minimap2` (v2.24 [50]) with `"-cx map-ont"`. Structural variants were called with `Sniffles2` [51] and then normalized and subsetted using the same approach as pangenome variation.

#### Variation discovery from optical mapping

The BRA and ANG samples had haplotype-specific optical mapping raw data available [28], which we aligned to ARS-UCD1.2 using the Bionano Solve 3.7 script `fa2cmap_multi_color.pl` with the DLE1 enzyme. We aligned the `bnx` data using the `RefAligner` script with two iterations and converted the resulting `smap` file into VCF using the script `smap_to_vcf_v2.py`. The two unrelated Nellore samples were already available in a filtered VCF [27].

#### Genomic region classification

Genomic regions were classified using centromeric satellites and repetitive elements identified with `RepeatMasker`, tandem repeats identified by `TRF` ([52], v4.10.0), and low mappability determined by `GenMap` (v1.3.0) [53] using `k`-mers of size 75 and a mappability threshold of 1 on the ARS-UCD1.2 reference genome. These regions were converted into BED format, and merged using `bedtools`, giving classification priority to `satellites > tandem repeats > repetitive elements > low mappability`. Uncovered regions were classified as "normal" using `bedtools complement`. Variants were then annotated by their genomic region classification using `bcftools annotate`.

#### Variation consensus

SVs from the three pangenomes and assemblies were merged using `Jasmine` (v.1.1.5) [54] with `--normalize_type --allow_intrasample max_dist=1000 max_dist_linear=0.5`. A more lenient SV merging was done using parameters `max_dist=10,000 max_dist_linear=1.0`. A more strict merging was done with `max_dist=5 min_seq_id=0.5`,

and a near-perfect merging was done with `max_dist=1 max_dist_linear=0 min_seq_id=0.85`.

Jasmine was also used for intersection SVs with optical mapping data, using the lenient parameters described above. The optical mapping data was subsetted per breed/species and filtered to keep only deletions and insertions between 50 bp and 1 Mb.

Small variation between individual haplotypes was intersected using `bcftools isec`, requiring exact matches of REF and ALT alleles with the `-c none` flag. The normalized cohort-level pangenome and assembly VCF files were intersected based on exact matches of small variant coordinates and ALT alleles.

### Simulated assemblies

Assemblies were simulated by using the per-sample VCF for chromosome 29 containing variant calls from the respective sample, and using `bcftools consensus` to impose the variation onto the ARS-UCD1.2 reference genome, creating pseudo-assemblies with known variation and no unknown sequence relative to ARS-UCD1.2.

### Calculation of edit distance

The twelve assemblies included in the pangenomes and eight haplotype assemblies from four taurine cattle which are not part of the pangenomes were realigned to the pangenomes using `GraphAligner` (v1.0.16, [55]) with parameters `-x dbg -C 100,000 -max-trace-count 5 -seeds-minimizer-ignore-frequent 0.001 -precise-clipping 0.9`. Centro-/telomeric sequence intervals, as classified by `RepeatMasker`, were merged within 250 Kb and 1 Kb by `bedtools` [56] respectively, and then split into centro-/telomeric sequence and bulk sequence fasta files. Both files were split every 500 Kb during realignment for computational constraints. Edit rate and query coverage were then calculated from the resulting graph alignment files.

The eight additional haplotype assemblies were constructed with the dual assembly approach in `hifiasm` (v0.16.1, [57]) using between 17- and 24-fold HiFi read coverage from two purebred Brown Swiss and two purebred Original Braunvieh cattle. The two Original Braunvieh cattle are the sire and dam (labelled as OS1, OS2, OD1, OD2) of an OxO F1 which was the source for the OBV haplotype [10] included in the pangenomes whereas the two Brown Swiss cattle are not directly related to the BSW haplotype included in the pangenomes (labelled as B11, B12, B21, B22).

### VNTR analysis

We identified TRs in the ARS-UCD1.2 reference sequence using `TRF` ([52], v4.10.0) with parameters `"2 7 7 80 10 50 500"`, allowing a minimum and maximum motif length of 10 and 100 bp respectively and a minimum and maximum tandem repeat length of 50 bp and 10 Kb. We excluded TRs which overlapped with repeat elements identified in the reference sequence using `bedtools intersect -v -f 0.5`. For `minigraph`, we then used `bedtools intersect` again for the TRs and the `minigraph` SV bed files generated by `gfatools bubble`. We then extracted the approximate coordinates using `minigraph -call` for each assembly. For `pggb` and `cactus` we used the `position` command of `odgi` to translate the TR coordinates from ARS-UCD1.2 to the paths of other assemblies in the graphs.

We calculated the number of TRs by extracting per-assembly sequence from the VNTR coordinates predicted from the pangenomes. We then used the python regex package to count how often the TR repeat motif occurs in the query sequence, allowing 25% error rate to account for substitutions, insertions, and/or deletions between repeat units.

We used adVNTR (v1.4.1, [31]) to genotype the 465 VNTRs with at least 50 repeat units in one of our HiFi samples. We built a database for these TRs using adVNTR addmodel, and then genotyped them using adVNTR genotype with the flags "--naive --pachio ---haploid" on the triobinned HiFi samples. These bam files were generated using sequences binned by Canu (v2.2, [58]) using parental reads, followed by alignment to ARS-UCD1.2 with minimap2 (as described in the "Variation discovery from assemblies" section except using the map-hifi preset).

The VNTR annotated pangenomes were generated using Bandage [59], using its built-in BLAST feature. The blastDB is created from the graph sequence, while we provide common TR motifs for given VNTRs as queries. Since motifs are "short" compared to typical blast queries, we use "-task blastn-short -word\_size <TR length> -evaluate 100" and filter minimum alignment sizes below <TR length> to increase sensitivity to specific motifs.

#### Establishing a testis eQTL cohort

Testis tissue of 83 mature bulls was collected at a commercial slaughterhouse and subjected to DNA and RNA purification as described earlier [60]. DNA samples were sequenced on an Illumina NovaSeq6000 instrument using 150 bp paired-end sequencing libraries. Following quality control, between 74,371,404 and 304,199,764 filtered read pairs per sample (mean:  $268,798,344 \pm 80,134,660$ ) were aligned to the ARS-UCD1.2 reference sequence and processed as described in Kadri et al. [60]. Single nucleotide and short insertion and deletion polymorphisms were discovered and genotyped using DeepVariant (version 1.3.0, [61]). Beagle4.1 [62] was applied to impute sporadically missing genotypes and infer haplotypes. A genomic relationship matrix was constructed and subjected to a principal components analysis using plink (v.1.9, [63]).

Total RNA sequencing libraries ( $2 \times 150$  bp) were prepared using the Illumina TruSeq Stranded Total RNA sequencing kit and sequenced on an Illumina NovaSeq6000. Following quality control, between 70,249,355 and 177,053,828 filtered read pairs per sample (mean:  $130,092,907 \pm 18,684,899$ ) were aligned to the ARS-UCD1.2 reference sequence and the Refseq gene annotation (release 106) using the splice-aware read alignment tool STAR (version 2.7.9a) [64] with options --twopassMode Basic and --wasp-OutputMode SAMtag to enable robust and unbiased allele-specific expression detection [65]. Per-individual VCF files required for WASP filtering were prepared from the cohort-level VCF files (see above).

Testis tissue from the NxB F1 was collected after regular slaughter. Total RNA was prepared and sequenced as described above. Following quality control, 136,820,136 filtered read pairs were aligned against the bovine reference sequence using STAR (see above).

### Gene expression quantification

The expression level of genes, ncRNA and lncRNA was quantified (in transcripts per million, TPM) using kallisto (version 0.46.1, [66]) and aggregated to the gene level using the R package tximport [67]. We retained genes, ncRNA and lncRNA that were expressed at an average value >0.2 TPM across the 83 transcriptomes. The raw TPM values were normalized using quantile normalization and rank-based inverse normal transformation [67]. We inferred hidden confounding variables from the normalized gene expression levels using Probabilistic Estimation of Expression Residuals (PEER) [68].

### eQTL mapping

The two VNTR haplotypes (hap1 and hap2) segregating in the Brown Swiss and Original Braunvieh breeds were derived from binned long-read alignments with minimap2 (see above). The resulting diplotypes were reconstructed for all bulls of the eQTL cohort based on the Beagle-phased genotypes (see above) for two SNPs (Chr12:86,160,984 and Chr12:86,160,971) and one indel (Chr12:86,161,000) that tagged the VNTR haplotypes. A linear model was fitted using the `lm()`-function in R to associate normalized gene expression with the VNTR haplotype (coded as 0, 1, and 2 for hap1/hap1, hap1/hap2, and hap2/hap2, respectively) while considering five PEER factors and the top three principal components of the genomic relationship matrix as covariates to account for technical bias and population stratification. An eQTL analysis within the *cis*-regulatory range of genes of interest was conducted using the linear model described above. The VNTR haplotype was fitted as an additional covariate in the conditional analyses. Bonferroni-correction was applied to determine significance thresholds.

### Allelic imbalance analysis

Polymorphic sites overlapping *LOC112449094* exons were extracted from the cohort-level genomic VCF file. Reference and alternate allele support at heterozygous genotypes was subsequently extracted from the WASP-filtered RNA sequencing read alignments using `bcftools mpileup`. An exact binomial test as implemented in the `R binom.test()`-function was applied to test for allelic imbalance.

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-023-02969-y>.

**Additional file 1: Fig. S1.** Pangenome repeat content. **Fig. S2.** Decomposing minigraph pangenomes. **Fig. S3.** Impact of adding P-lines to minigraph pangenomes. **Fig. S4.** Conflicting variation in three pangenomes. **Fig. S5.** Overlap of SVs between the three pangenomes and assemblies. **Fig. S6.** Variation representation accuracy in three pangenomes. **Fig. S7.** Pangenome performance on simulated data. **Fig. S8.** Tandem repeats that disagree between the three pangenomes. **Fig. S9.** VNTR representation in three pangenomes. **Fig. S10.** Representation of highly repetitive VNTR loci in three pangenomes. **Fig. S11.** Detection of two VNTR haplotypes in Original Braunvieh and Brown Swiss cattle. **Fig. S12.** eQTL mapping between the VNTR haplotype and normalized expression of 21,853 genes. **Fig. S13.** *cis*-expression QTL mapping for two genes and four lncRNA that were significantly associated with the VNTR haplotype. **Fig. S14.** Community detection from all-versus-all alignment of the 29 autosomes.

**Additional file 2: Table S1.** Compute resources for pangenome deconstruction. Cumulative CPU hours and maximum memory needed for vg deconstruct and the number of variants recovered for the three pangenomes.

**Additional file 3: Table S2.** Assessment of structural variant representation with orthogonal data. Orthogonal optical mapping (OM) and ONT data from relevant breeds/species to assess SV overlap with three pangenomes. Optical



mapping generates fewer SV calls, and so only the total number of overlapping SVs are shown. ONT reads generated a comparable number of SVs to the pangenomes, and we can directly calculate an F-score. The Nellore samples are unrelated to the NEL haplotype used in this work, while the remaining data are from the same individual used to generate the assembly.

**Additional file 4: Table S3.** Sequence content of input assemblies. CSV containing the number of centromeric sequence bases trimmed from the start of the chromosomes, bulk sequence bases, and telomeric sequence trimmed from the end for each chromosome and each assembly.

**Additional file 5: Table S4.** Impact of parameter settings on graph alignment. GraphAligner alignments using relaxed (-x dbg -C 100000 --max-trace-count 5 --seeds-minimizer-ignore-frequent 0.001 --precise-clipping 0.9) or strict (-x vg) alignment parameters. The edit rate and query coverage are taken from a single 500 Kb window in chromosome 12 of the Brown Swiss assembly (start at 60867381), chosen for its complexity and subsequent poor relaxed alignment. CPU time is given in seconds and Memory is peak RAM usage in GB. Strict alignment quickly requires more resources for pggg and becomes prohibitive for cactus for whole chromosome alignment.

**Additional file 6: Table S5.** Compute resources for alignment of additional assemblies. Compute resources for aligning the 12 assemblies included in pangenome construction and the 8 held out for analysis. CPU hours are averaged per assembly for both the bulk and centro-/telomeric regions, and memory is the peak RAM usage across all assemblies.

**Additional file 7: Table S6.** Overview of investigated tandem repeats. CSV containing the examined tandem repeats in pggg, cactus, and minigraph. Tandem repeats further examined by adVNTR are also indicated.

**Additional file 8: Table S7.** Several obvious misalignments found in different pangenomes which affects TR count genotyping. Part of the size difference may be true if there is variation in the number of tandem repeats, but the obvious majority of the large mis-translated regions are due to repetitive regions or spurious graph cycles.  $\Delta$  bp is the size difference for the lifted-over sample TR coordinates from the original reference TR coordinates.

**Additional file 9: Table S8.** Minigraph pangenomes constructed with different "minimum variant length" parameter (L) values. CPU hours required for pangenome construction increased dramatically for L smaller than 10 bp. Warnings refer to the number of "impossible insert" warnings issued during pangenome construction, relating to unsuitable graph topology. Bubbles and Nodes respectively refer to the number of top-level bubbles and nodes present across the autosomes. VNTR overlaps is the number of VNTRs (in total 9,568) that overlap with a graph bubble.

**Additional file 10.** Review History.

#### Acknowledgements

We thank Dr. Anna Bratus-Neuenschwander and Dr. Catharine Aquino from the ETH Zurich technology platform FGCZ (<https://fgcz.ch>) for DNA fragment analysis and DNA and RNA sequencing.

#### Peer review information

Anahita Bishop and Andrew Cosgrove were the primary editors of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

#### Review history

The review history is available as Additional file 10.

#### Authors' contributions

A.S.L., D.C., and H.P. conceived the study. A.S.L. constructed genome assemblies, constructed/decomposed the pangenomes, and performed the pangenome, variant, and VNTR analyses. D.C. contributed to pangenome construction, pangenome, variant, and VNTR analyses. X.M.M. established the eQTL cohort and aligned the RNA and DNA sequencing data. A.S.L. called variants for the eQTL cohort. H.P. performed the eQTL analyses. M.B. contributed to the design of the eQTL analyses. A.S.L. and H.P. wrote the manuscript with input from D.C. All authors read and approved the final manuscript.

#### Funding

Open access funding provided by Swiss Federal Institute of Technology Zurich. This work was financially supported by the Swiss National Sciences Foundation (SNSF), an ETH Research Grant, and Swissgenetics. The funders had no role in study design, data collection and analysis, interpretation of the data, decision to publish, or preparation of the manuscript.

#### Availability of data and materials

Scripts and pipelines used in this work are available online at a Github repository ([https://github.com/AnimalGenomicsETH/superpangenome\\_construction](https://github.com/AnimalGenomicsETH/superpangenome_construction)) under MIT license, which has been archived at zenodo [69]. HiFi reads of two BSW and two OBV cattle are available in the ENA database at the study accession PRJEB42335 [70] under sample accessions SAMEA111341789, SAMEA111341790, SAMEA111341791, and SAMEA111341792. DNA and RNA sequencing data of the eQTL cohort are available in the ENA database at the study accessions PRJEB28191 [71] and PRJEB46995 [72]. Assemblies used in pangenome construction are available at <https://zenodo.org/record/5906579> (NEL, BSW, OBV, PIE, GAU), [https://www.ncbi.nlm.nih.gov/assembly/GCF\\_002263795.1/](https://www.ncbi.nlm.nih.gov/assembly/GCF_002263795.1/) (ARS-UCD1.2), [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_003369685.2](https://www.ncbi.nlm.nih.gov/assembly/GCA_003369685.2) (ANG), [https://www.ncbi.nlm.nih.gov/assembly/GCF\\_003369695.1](https://www.ncbi.nlm.nih.gov/assembly/GCF_003369695.1) (BRA), [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_009493645.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_009493645.1) (yak), [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_009493655.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_009493655.1) (HIG), [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_018282465.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_018282465.1) (SIM), [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_018282365.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_018282365.1) (BIS). The per-chromosome pangenomes for minigraph, pggg, and cactus, as well as the genomic regions classification BED file are available online at zenodo at <https://doi.org/10.5281/zenodo.7737904> [73].



## Declarations

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Competing interests

The authors declare that they have no competing interests.

Received: 17 September 2022 Accepted: 10 May 2023

Published online: 22 May 2023

## References

- Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, et al. Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: Implications for the microbial “pan-genome.” *Proc Natl Acad Sci U S A*. 2005;102:13950–5.
- Eizenga JM, Novak AM, Sibbesen JA, Heumos S, Ghaffaari A, Hickey G, et al. Pangenome Graphs. *Annu Rev Genomics Hum Genet*. 2020;21:139–62.
- Li H, Feng X, Chu C. The design and construction of reference pangenome graphs. *Genome Biol*. 2020;21:1–19.
- Crysnanto D, Leonard AS, Fang ZH, Pausch H. Novel functional sequences uncovered through a bovine multi-assembly graph. *Proc Natl Acad Sci USA*. 2021;118:1–29.
- Armstrong J, Hickey G, Diekhans M, Fiddes IT, Novak AM, Deran A, et al. Progressive Cactus is a multiple-genome aligner for the thousand-genome era. *Nature*. 2020;587:246–51.
- Garrison E, Guarracino A, Heumos S, Villani F, Bao Z, Tattini L, et al. Building pangenome graphs. *bioRxiv*. 2023. <https://doi.org/10.1101/2023.04.05.535718>.
- Holley G, Melsted P. Bifrost: highly parallel construction and indexing of colored and compacted de Bruijn graphs. *Genome Biol*. 2020;21:1–20.
- Rossi M, Oliva M, Langmead B, Gagie T, Boucher C. MONI: a pangenomic index for finding maximal exact matches. *J Comput Biol*. 2022;29:169–87.
- Qiu Y, Kingsford C. Constructing small genome graphs via string compression. *Bioinformatics*. 2021;37(Suppl\_1):i205–13.
- Leonard AS, Crysnanto D, Fang Z-H, Heaton MP, Vander Ley BL, Herrera C, et al. Structural variant-based pangenome construction has low sensitivity to variability of haplotype-resolved bovine assemblies. *Nat Commun*. 2022;13:3012.
- Wang T, Antonacci-Fulton L, Howe K, Lawson HA, Lucas JK, Phillippy AM, et al. The Human Pangenome Project: a global resource to map genomic diversity. *Nature*. 2022;604:437–46.
- Nurk S, Koren S, Rhie A, Rautiainen M, Bizikadze AV, Mikheenko A, et al. The complete sequence of a human genome. *Science*. 2022;376:44–53.
- Rautiainen M, Nurk S, Walenz BP, Logsdon GA, Porubsky D, Rhie A, et al. Telomere-to-telomere assembly of diploid chromosomes with Verkko. *Nat Biotechnol*. 2023. <https://doi.org/10.1038/s41587-023-01662-6>.
- Jarvis ED, Formenti G, Rhie A, Guarracino A, Yang C, Wood J, et al. Semi-automated assembly of high-quality diploid human reference genomes. *Nature*. 2022;611:519–31.
- Rhie A, McCarthy SA, Fedrigo O, Damas J, Formenti G, Koren S, et al. Towards complete and error-free genome assemblies of all vertebrate species. *Nature*. 2021;592:737–46.
- Guarracino A, Buonaiuto S, Lima LG de, Potapova T, Rhie A, Koren S, et al. Recombination between heterologous human acrocentric chromosomes. *bioRxiv*. 2023; <https://doi.org/10.1101/2022.08.15.504037>.
- Lu TY, Munson KM, Lewis AP, Zhu Q, Tallon LJ, Devine SE, et al. Profiling variable-number tandem repeat variation across populations using repeat-pangenome graphs. *Nat Commun*. 2021;12:1–12.
- Porubsky D, Vollger MR, Harvey WT, Rozanski AN, Ebert P, Hickey G, et al. Gaps and complex structurally variant loci in phased genome assemblies. *bioRxiv*. 2022; <https://doi.org/10.1101/2022.07.06.498874>.
- Ebler J, Ebert P, Clarke WE, Rausch T, Audano PA, Houwaart T, et al. Pangenome-based genome inference allows efficient and accurate genotyping across a wide spectrum of variant classes. *Nat Genet*. 2022;54:518–25.
- Liao W-W, Asri M, Ebler J, Doerr D, Haukness M, Hickey G, et al. A draft human pangenome reference. *bioRxiv*. 2022; <https://doi.org/10.1101/2022.07.09.499321>.
- Schneider VA, Graves-Lindsay T, Howe K, Bouk N, Chen HC, Kitts PA, et al. Evaluation of GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly. *Genome Res*. 2017;27:849–64.
- Zhou Y, Zhang Z, Bao Z, Li H, Lyu Y, Zan Y, et al. Graph pangenome captures missing heritability and empowers tomato breeding. *Nature*. 2022;606:527–34.
- Tang D, Jia Y, Zhang J, Li H, Cheng L, Wang P, et al. Genome evolution and diversity of wild and cultivated potatoes. *Nature*. 2022;606:535–41.
- Li R, Gong M, Zhang X, Wang F, Liu Z, Zhang L, et al. A sheep pangenome reveals the spectrum of structural variations and their effects on tail phenotypes. *Genome Res*. 2023;33:463–77.
- Koren S, Rhie A, Walenz BP, Dilthey AT, Bickhart DM, Kingan SB, et al. De novo assembly of haplotype-resolved genomes with trio binning. *Nat Biotechnol*. 2018;36:1174–82.
- Talenti A, Powell J, Hemmink JD, Cook EAJ, Wragg D, Jayaraman S, et al. A cattle graph genome incorporating global breed diversity. *Nat Commun*. 2022;13:910.
- Talenti A, Powell J, Wragg D, Chepkwony M, Fisch A, Ferreira BR, et al. Optical mapping compendium of structural variants across global cattle breeds. *Scientific Data*. 2022;9:1.

28. Low WY, Tearle R, Liu R, Koren S, Rhie A, Bickhart DM, et al. Haplotype-resolved genomes provide insights into structural variation and gene content in Angus and Brahman cattle. *Nat Commun.* 2020;11:2071.
29. Khan AW, Garg V, Roorkiwal M, Golicz AA, Edwards D, Varshney RK. Super-Pangenome by integrating the wild side of a species for accelerated crop improvement. *Trends Plant Sci.* 2020;25:148–58.
30. Chaisson MJP, Sanders AD, Zhao X, Malhotra A, Porubsky D, Rausch T, et al. Multi-platform discovery of haplotype-resolved structural variation in human genomes. *Nat Commun.* 2019;10:1784.
31. Bakhtiari M, Park J, Ding YC, Shleizer-Burko S, Neuhausen SL, Halldórsson BV, et al. Variable number tandem repeats mediate the expression of proximal genes. *Nat Commun.* 2021;12:2075.
32. Lu T-Y, Smaruj PN, Fudenberg G, Mancuso N, Chaisson MJP. The motif composition of variable number tandem repeats impacts gene expression. *Genome Res.* 2023. <https://doi.org/10.1101/gr.276768.122>.
33. Sibbesen JA, Eizenga JM, Novak AM, Sirén J, Chang X, Garrison E, et al. Haplotype-aware pantranscriptome analyses using spliced pangenome graphs. *Nat Methods.* 2023;20:239–47.
34. Rosen BD, Bickhart DM, Schnabel RD, Koren S, Elisk CG, Tseng E, et al. De novo assembly of the cattle reference genome with single-molecule sequencing. *GigaScience.* 2020;9(3):giaa021.
35. Guarracino A, Heumos S, Nahnsen S, Prins P, Garrison E. ODGI: understanding pangenome graphs. *Bioinformatics.* 2022;38:3319–26.
36. Hereford assembly ARS-UCD1.2. NCBI. [https://www.ncbi.nlm.nih.gov/assembly/GCF\\_002263795.1/](https://www.ncbi.nlm.nih.gov/assembly/GCF_002263795.1/). Accessed 5 May 2023.
37. Angus assembly UOA\_Angus\_1. NCBI. [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_003369685.2](https://www.ncbi.nlm.nih.gov/assembly/GCA_003369685.2). Accessed 5 May 2023.
38. Brahman assembly UOA\_Brahman\_1. NCBI [https://www.ncbi.nlm.nih.gov/assembly/GCF\\_003369695.1](https://www.ncbi.nlm.nih.gov/assembly/GCF_003369695.1). Accessed 5 May 2023.
39. Rice ES, Koren S, Rhie A, Heaton MP, Kalbfleisch TS, Hardy T, et al. Continuous chromosome-scale haplotypes assembled from a single interspecies F1 hybrid of yak and cattle. *GigaScience.* 2020;9:giaa029.
40. Highland assembly ARS\_UNL\_Btau-highland\_paternal\_1.0\_alt. NCBI. [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_009493655.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_009493655.1). Accessed 5 May 2023.
41. Yak assembly ARS\_UNL\_BGru\_maternal\_1.0\_p. NCBI. [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_009493645.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_009493645.1). Accessed 5 May 2023.
42. Heaton MP, Smith TPL, Bickhart DM, Vander Ley BL, Kuehn LA, Oppenheimer J, et al. A reference genome assembly of simmental cattle, *Bos taurus taurus*. *J Hered.* 2021;112:184–91.
43. Simmental assembly ARS\_Simm1.0. NCBI. [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_018282465.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_018282465.1). Accessed 5 May 2023.
44. Oppenheimer J, Rosen BD, Heaton MP, Vander Ley BL, Shafer WR, Schuetze FT, et al. A reference genome assembly of American Bison, *Bison bison bison*. *J Hered.* 2021;112:174–83.
45. Bison assembly ARS-UCSC\_bison1.0. NCBI. [https://www.ncbi.nlm.nih.gov/assembly/GCA\\_018282365.1](https://www.ncbi.nlm.nih.gov/assembly/GCA_018282365.1). Accessed 5 May 2023.
46. Leonard A. Bovine pangenome assemblies, Zenodo. 2022. <https://zenodo.org/record/5906579>.
47. Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: Fast genome and metagenome distance estimation using MinHash. *Genome Biol.* 2016;17:132.
48. Garrison E, Sirén J, Novak AM, Hickey G, Eizenga JM, Dawson ET, et al. Variation graph toolkit improves read mapping by representing genetic variation in the reference. *Nat Biotechnol.* 2018;36:875–81.
49. Garrison E, Kronenberg ZN, Dawson ET, Pedersen BS, Prins P. A spectrum of free software tools for processing the VCF variant call format: vcfliib, bio-vcf, cyvcf2, hts-nim and slivar. *PLoS Comput Biol.* 2022;18:e1009123.
50. Li H. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics.* 2018;34:3094–100.
51. Smolka M, Paulin LF, Grochowski CM, Mahmoud M, Behera S, Gandhi M, et al. Comprehensive Structural Variant Detection: From Mosaic to Population-Level. *bioRxiv.* 2022; <https://doi.org/10.1101/2022.04.04.487055>.
52. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.* 1999;27:573–80.
53. Pockrandt C, Alzamel M, Iliopoulos CS, Reinert K. GenMap: ultra-fast computation of genome mappability. *Bioinformatics.* 2020;36:3687–92.
54. Kirsche M, Prabhu G, Sherman R, Ni B, Battle A, Aganezov S, et al. Jasmine and Iris: population-scale structural variant comparison and analysis. *Nat Methods.* 2023;20:408–17.
55. Rautiainen M, Marschall T. GraphAligner: rapid and versatile sequence-to-graph alignment. *Genome Biol.* 2020;21:253.
56. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics.* 2010;26:841–2.
57. Cheng H, Concepcion GT, Feng X, Zhang H, Li H. Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat Methods.* 2021;18:170–5.
58. Nurk S, Walenz BP, Rhie A, Vollger MR, Logsdon GA, Grothe R, et al. HiCanu: accurate assembly of segmental duplications, satellites, and allelic variants from high-fidelity long reads. *Genome Res.* 2020;30:1291–305.
59. Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics.* 2015;31:3350–2.
60. Kadri NK, Mapel XM, Pausch H. The intronic branch point sequence is under strong evolutionary constraint in the bovine and human genome. *Commun Biol.* 2021;4:1–13.
61. Poplin R, Chang PC, Alexander D, Schwartz S, Colthurst T, Ku A, et al. A universal snp and small-indel variant caller using deep neural networks. *Nat Biotechnol.* 2018;36:983.
62. Browning BL, Browning SR. Genotype Imputation with Millions of Reference Samples. *Am J Hum Genet.* 2016;98:116–26.
63. Chang CC, Chow CC, Tellier LCAM, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: Rising to the challenge of larger and richer datasets. *GigaScience.* 2015;4:1–16.
64. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29:15–21.

65. Van De Geijn B, Mcvicker G, Gilad Y, Pritchard JK. WASP: Allele-specific software for robust molecular quantitative trait locus discovery. *Nat Methods*. 2015;12:1061–3.
66. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol*. 2016;34:525–7.
67. Soneson C, Love MI, Robinson MD. Differential analyses for RNA-seq: Transcript-level estimates improve gene-level inferences. *F1000Research*. 2016;4:1–23.
68. Stegle O, Parts L, Piipari M, Winn J, Durbin R. Using probabilistic estimation of expression residuals (PEER) to obtain increased power and interpretability of gene expression analyses. *Nat Protoc*. 2012;7:500–7.
69. Leonard A, Crysanto D. AnimalGenomicsETH/superpangenome\_construction: v1.0. Zenodo. 2023. <https://doi.org/10.5281/zenodo.7891567>.
70. Leonard AS, Crysanto D, Mapel XM, Bhati M, Pausch H. Graph construction method impacts variation representation and analyses in a bovine super-pangenome. Datasets. European Nucleotide Archive. 2023. <https://www.ebi.ac.uk/ena/browser/view/PRJEB42335>.
71. Leonard AS, Crysanto D, Mapel XM, Bhati M, Pausch H. Graph construction method impacts variation representation and analyses in a bovine super-pangenome. Datasets. European Nucleotide Archive. 2023. <https://www.ebi.ac.uk/ena/browser/view/PRJEB28191>.
72. Leonard AS, Crysanto D, Mapel XM, Bhati M, Pausch H. Graph construction method impacts variation representation and analyses in a bovine super-pangenome. Datasets. European Nucleotide Archive. 2023. <https://www.ebi.ac.uk/ena/browser/view/PRJEB46995>.
73. Leonard AS. Graph construction method impacts variation representation and analyses in a bovine super-pangenome. Zenodo. 2023. <https://doi.org/10.5281/zenodo.7737904>.

### Publisher's Note

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

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

