

Bovine pangenome reveals trait-associated structural variation from diverse assembly inputs

Alexander S. Leonard¹, Danang Crysanto¹, Zih-Hua Fang¹, Michael P Heaton², Brian L. Vander Ley³, Carolina Herrera⁴, Heinrich Bollwein⁴, Derek M. Bickhart⁵, Kristen L. Kuhn², Timothy PL. Smith², Benjamin D. Rosen⁶, Hubert Pausch¹

¹Animal Genomics, ETH Zurich, Universitaetstrasse 2, 8006 Zurich, Switzerland

²U.S. Meat Animal Research Center, USDA-ARS, 844 Road 313, Clay Center, NE 68933, USA

³Great Plains Veterinary Educational Center, University of Nebraska-Lincoln, Lincoln, NE, 68588, USA

⁴Clinic of Reproductive Medicine, Department for Farm Animals, University of Zurich, 8057, Zurich, Switzerland

⁵Dairy Forage Research Center, USDA-ARS, 1925 Linden Drive, Madison, WI, 53706, USA

⁶USDA-ARS, Beltsville, MD, 20705-2350, Animal Genomics and Improvement Laboratory, USDA-ARS, 10300 Baltimore Ave, Beltsville, MD 20705, USA

Abstract

Advantages of pangenomes over linear reference assemblies for genome research have recently been established. However, potential effects of sequence platform and assembly approach, or of combining assemblies created by different approaches, on pangenome construction have not been investigated. Ten haplotype-resolved assemblies of three bovine trios representing increasing levels of heterozygosity were generated that each demonstrate a substantial improvement in contiguity and accuracy over the current *Bos taurus* reference genome, with more telomere and centromere content and an average 2.5x increase in NG50 and 11x decrease in base errors. Downsampling analysis demonstrated that diploid coverage as low as 20x for HiFi or 60x for ONT was sufficient to produce two haplotype-resolved assemblies meeting the standards set by the Vertebrate Genome Project. The assemblies were integrated into structural variant-based pangenomes that demonstrated significant consensus regardless of sequence platform, assembler algorithm, or coverage. Inspecting pangenome topologies identified approximately 900 structural variants overlapping with coding sequences; this approach revealed variants affecting *QRICH2*, *PRDM9*, *HSPA1A*, *TAS2R46* and *GC* that have potential to affect phenotype.

Introduction

Cattle are a substantial component of global animal-based food production, and are raised for meat, milk, or both. Two subspecies of cattle, taurine and indicine, have emerged from at least two distinct domestication events (Loftus et al., 1994; Pitt et al., 2019)), with

artificial selection for production goals or environmental adaptation contributing to diversity within cattle, resulting in the current existence of hundreds of taurine and indicine cattle breeds. Interbreeding and introgressions with other bovids, like yak and banteng (N. Chen et al., 2018; Wu et al., 2018), further drive an increase in genetic diversity within Bovinae.

The *Bos taurus* reference genome was first drafted in 2004 and was based on whole genome shotgun sequence of a Hereford cow, supplemented with sequences of bacterial artificial chromosome clones prepared from DNA of her sire (Elsik et al., 2009). A major revision using long read sequencing of the same cow was recently reported (ARS-UCD1.2; (Rosen et al., 2020)) and remains the accepted reference for conducting genomic studies in cattle due to extensive annotation efforts and connections to historical analyses, despite more recent bovine assemblies having higher contiguity and accuracy (Heaton et al., 2021; Oppenheimer et al., 2021; Rice et al., 2020).

Sequence variability between cattle breeds at both the single nucleotide (SNP) and short insertion or deletion (indel) level has been extensively characterized through reference-guided approaches (Daetwyler et al., 2014; Hayes & Daetwyler, 2019; Kim et al., 2017). However, these studies suffer from potential reference bias resulting from the use of a single taurine cattle reference assembly for SNP and indel discovery. Larger structural variants (SVs) and variation located in repetitive or challenging regions have rarely been studied across Bovinae due to the inherent limitations of short sequencing reads and incomplete reference genomes. Moreover, no reference assembly of a single individual can reflect the immense genomic diversity present in global breeds of domestic cattle (Crysnanto et al., 2019, 2021; Crysnanto & Pausch, 2020; Talenti et al., 2021).

Pangenomes have long been proposed (Tettelin et al., 2005) as a way to better reflect variation present in a group of individuals (e.g., breed, species, clade, etc.). Pangenomes can be constructed from variants called through reference-guided approaches (Garrison et al., 2018; Hickey et al., 2020), contigs assembled from reads which failed to align to the reference (Sherman et al., 2018), or multiple whole-genome assemblies (Armstrong et al., 2020; Li et al., 2020). The latter approach may more faithfully capture challenging regions and SVs due to the complexity of calling and representing nested variation. Third-generation sequencing continues to become more cost-effective and accessible, making population-scale *de novo* assemblies more feasible. An influx of high-quality assemblies makes the need for pangenome representations more pressing, although the effects of integrating disparate assemblies into pangenomes are unknown.

The present study employed different sequencing and assembly approaches to produce genomes from three bovine trios of varying heterozygosity. The relative strengths of Pacific Biosciences high-fidelity (HiFi) (Eid et al., 2009; Wenger et al., 2019) and Oxford Nanopore Technologies (ONT) (Mikheyev & Tin, 2014) sequencing were examined through generation of haplotype-resolved, reference-quality assemblies for the same individuals sequenced with both platforms and assembled with various algorithms. The set of assemblies were then used to evaluate the effects on pangenome construction depending on assembly approach, or when constructed from a combination of approaches. The utility of a bovine pangenome is then demonstrated by analyses that reveal novel insights into the evolutionary relationships between Bovinae.

Results

The three examined bovine trios (Figure 1a-c) reflect diverse breeding strategies (within-breed, inter-subspecies, and inter-species) and increasing heterozygosity. The first F1 (OxO) was a cross between two Original Braunvieh cattle (*Bos taurus taurus*), but was still substantially less inbred compared to the cow used for the ARS-UCD1.2 reference (pedigree-based coefficient of inbreeding of 0.07 compared to 0.30 (Rosen et al., 2020) respectively). The second F1 (NxB) was a cross between Nellore (*Bos taurus indicus*) and Brown Swiss (*Bos taurus taurus*) cattle and the third F1 (GxP) was an inter-species cross between a gaur (*Bos gaurus*) bull and a Piedmontese (*Bos taurus taurus*) cow. HiFi and ONT reads were collected for each F1, and short reads were collected for all animals in the trios (Figure 1d, Supplementary Table 1). F1 long reads were separated into paternal, maternal, and unknown origin (Koren et al., 2018). The success of parent-of-origin assignment improved significantly from 81.1% to 99.9% with increasing heterozygosity for HiFi reads, but was near perfectly separable for ONT at all examined heterozygosities (Figure 1e, Supplementary Table 2).

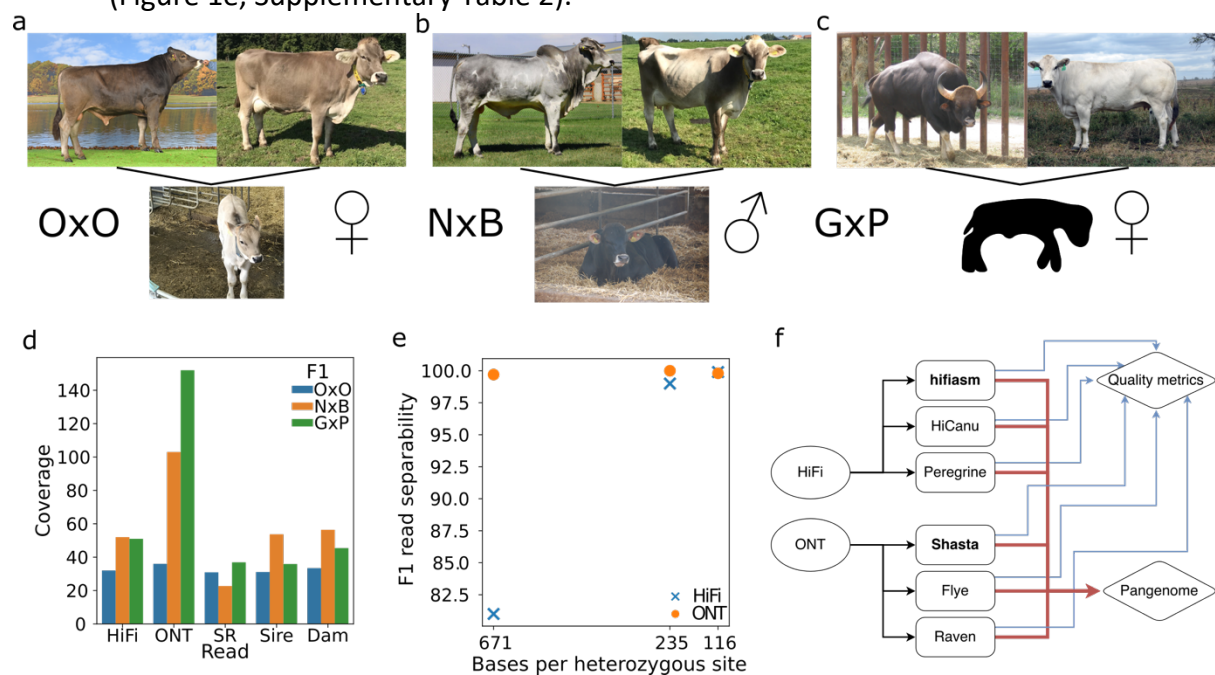


Figure 1. Overview of bovine trios. a) – c) Representative animals for the parents of the three bovine trios and the respective F1s (OxO, NxB, and GxP) examined in this study. The OxO and GxP were female, while the NxB was male. d) The three respective F1s were sequenced to 32-, 52-, and 51- fold HiFi coverage, with read N50 of 20, 21, and 14 Kb. ONT sequencing was performed to 36-, 103-, and 152-fold coverage respectively, with read N50 of 65, 45, and 49 Kb. Coverage is determined with respect to an assumed genome size of 2.7 Gb. F1 short reads were collected to 31-, 23-, and 37-fold coverage. e) Separating F1 reads into parental haplotype bins improved with increasing heterozygosity for HiFi, but F1 reads were nearly 100% separable for ONT even at low heterozygosity. f) HiFi reads were assembled with hifiasm, HiCanu, and Peregrine, while ONT reads were assembled with Shasta, Flye, and Raven. Bold type indicates the tools used to produce the assemblies that are discussed in detail. Assemblies were assessed individually for quality metrics (blue lines) as well as integrated together into pangenome analyses (red lines).

Assembling bovine genomes

Contigs were assembled for each F1 using multiple assemblers for HiFi and ONT data (Figure 1f, see methods). Contigs were scaffolded by alignment to ARS-UCD1.2 to produce the final assemblies, an approach with minimal bias given the well-curated reference and highly contiguous assemblies (Alonge et al., 2019). The haplotype-resolved assemblies for the inter-subspecies NxS and inter-species GxP are directly breed/species-specific, while both haplotypes of the OxO constitute the same breed, resulting in a total of five novel breed/species assemblies. Variation between the paternal and maternal haplotypes of the OxO, reflecting within-breed diversity, was sufficient to generate haplotype-resolved assemblies, although only the maternal haplotype was used to represent the Original Braunvieh breed in subsequent analyses except where explicitly stated. The “quality” of each assembly was assessed by several widely used metrics: contiguity (NG50) representing the size distribution of contigs, phasing (PG50) to characterize haplotype separation, correctness (QV) quantified as Phred-scaled base-error rate, and completeness (BUSCO) which approximates the percentage of near-universal, single-copy genes that were identified. These metrics provide a coarse-grained summary of assembly quality amenable to comparisons.

Table 1 The assemblies produced by hifiasm for HiFi data and Shasta for ONT data were selected for further analysis based on having the best quality metrics (Table 1) and computational tractability compared to four other tested assemblers (Supplementary Table 3). These assemblies were of reference quality with every examined metric exceeding those of the current *Bos taurus* Hereford-based reference. Improvements over the current reference are substantial, including a 3-11x reduction in autosomal gaps, 1.8-3.6x increase in NG50, and 3-22x reduction in base errors. Furthermore, they all exceed the current standards set by the Vertebrate Genome Project (VGP) (Rhie et al., 2021). The ONT-based assemblies were marginally above the QV targeted by the VGP, but other metrics for these assemblies such as the contig or phased contiguity are orders of magnitude greater than VGP thresholds.

Table 1. Quality metrics for ten haplotype-resolved assemblies. The assembly haplotype is either maternal or paternal (indicating either an “X” or “Y” paternal sex chromosome). The ARS-UCD1.2 reference is not haplotype-resolved and lacks sufficient parental data to assess phasing, hence the N/A. Size and contigs refer to the entire genome assembly, while the autosomal values only measure chromosomes 1 through 29. NG50 is the contig N50 using the ARS-UCD1.2 reference sequence as the expected length. PG50 is NG50 after splitting contigs into haplotype-phased blocks. Phasing and QV are determined through merqury using parental and F1 short reads. Scaffolded NG50 is not shown, as all assemblies are effectively end-to-end (excluding centromeres and telomeres), with values greater than 100 Mb. CLR: PacBio Continuous Long Read.

Breed or species	Haplotype	Read technology	Size (autosomal size)	Contigs (autosomal contigs)	NG50	PG50	QV	BUSCO (single copy)	Repeat
Original Braunvieh Nellore	Maternal	HiFi	3.11 (2.57)	1706 (105)	47.0	23.6	49.7	95.7 (93.9)	48.95
	Maternal	ONT	2.70 (2.48)	2622 (109)	71.6	2.8	40.7	95.1 (93.4)	43.19
	Paternal (Y)	HiFi	2.95 (2.60)	1217 (52)	94.4	79.1	46.1	93.3 (91.8)	47.81
	Paternal (Y)	ONT	2.57 (2.49)	1457 (67)	68.5	64.9	42.4	92.8 (91.3)	42.64
Brown Swiss gaur	Maternal	HiFi	3.07 (2.62)	1045 (58)	86.7	81.1	45.6	95.9 (94.2)	48.43
	Maternal	ONT	2.67 (2.48)	1268 (71)	64.0	53.0	42.5	95.3 (93.7)	42.85
	Paternal (X)	HiFi	3.02 (2.52)	1352 (75)	73.5	61.2	48.4	95.7 (94.1)	47.73

<i>Piedmontese</i>	Paternal (X)	ONT	2.64 (2.48)	532 (89)	68.1	68.1	41.2	95.1 (93.3)	42.26
	Maternal	HiFi	3.10 (2.56)	1427 (90)	52.0	47.6	48.3	95.8 (94.1)	48.43
	Maternal	ONT	2.66 (2.48)	782 (64)	82.8	82.8	40.9	95.3 (93.6)	43.06
<i>Hereford</i> (ARS-UCD1.2)	(N/A)	CLR	2.72 (2.49)	2597 (289)	25.9	N/A	35.8	95.7 (93.9)	42.96
VGP Standards					1	0.1	40	90	N/A

The HiFi- and ONT-based assemblies were generally comparable, however there were notable differences in the average assembly correctness and genome size metrics. Shasta assemblies averaged QV 41.5 after one round each of polishing with ONT reads and short reads, while hifiasm assemblies reach QV 47.6 without any polishing. The log scale of QV means that the hifiasm assemblies had a 4-fold reduction in base errors compared to the Shasta assemblies, indicating the ability of HiFi data to achieve higher quality in fewer steps. In contrast, phasing in Shasta assemblies is better compared to hifiasm, while both platforms showed improved phasing at higher heterozygosity in agreement with the relative ability to sort F1 reads by parental origin prior to assembly. The lower PG50 of the Original Braunvieh Shasta assembly reflects limited ONT coverage that resulted in the need to perform diploid polishing, rather than directly assembling trio binned reads, and is not due to an inherent limitation imposed by the level of heterozygosity.

The mean autosomal genome size of the assemblies generated by hifiasm and Shasta was 2.57 ± 0.03 Gb and 2.48 ± 0.004 Gb respectively, such that on average each hifiasm autosome was 2.8 Mb longer than ARS-UCD1.2, and 0.34 Mb shorter for Shasta. The additional length of hifiasm autosomes was primarily due to the presence of more repetitive sequences (43.5% repetitive content in hifiasm versus 41.7% in Shasta), especially centromeric repeat sequence. Previous studies have shown that the higher accuracy of HiFi reads allows assemblers to confidently assemble through more repeats despite having shorter read lengths (Chu et al., 2021), leading to extension of autosomal contigs into flanking centromere sequence. Hifiasm assemblies also contained more sequence in contigs not assigned to chromosomes (average of 300 Mb) compared to Shasta assemblies (50 Mb). These unassigned contigs were composed primarily of repetitive sequences (Supplementary Table 4) including novel centromeric sequence and long terminal repeats. Unplaced contigs were generally higher in repeat content than the scaffolds (88% versus 48% for hifiasm and 87% versus 42% for Shasta), and thus would present a challenge to scaffolding by any technology including the reference alignment approach applied here.

Bovine autosomes are acrocentric (Blazak & Eldridge, 1977), and so a complete bovine assembly is conceptually closer to “centromere-to-telomere”. Hifiasm assemblies contained substantially more centromeric sequence than Shasta assemblies, respectively averaging 2.01 and 0.14 Mb per autosome, compared to 0.08 in the ARS-UCD1.2 reference (Figure 2a). Similarly, hifiasm autosomes average 2.6 Kb of vertebrate telomeric repeats (TTAGGG) within 10 Kb of the chromosome end, compared to 0.8 Kb for the ARS-UCD1.2 reference. Telomeres were almost entirely missing in Shasta assemblies, averaging only 88 bp of telomeric repeats per autosome (Figure 2b). Chromosomes which contain at least 50 kb of centromeric repeats at the proximal end and 500 bp of telomeric repeats at the distal end

are considered to be end-to-end (but not necessarily “complete”), of which there were 5 for ARS-UCD1.2, and a mean of 13.2 and 1.2 for hifiasm and Shasta across the five breeds/species. Hifiasm and Shasta assemblies had a near equal distribution of gaps, averaging about 1.5 gaps per autosome, compared to nearly 9 in the ARS-UCD1.2 reference (Figure 2c). These observations hold in general for all HiFi- and ONT-based bovine assemblies investigated (Supplementary Figure 1). These differences are visible on the example Brown Swiss chromosome ideograms in Figure 2d. Both HiFi- and ONT-based assemblies were able to routinely assemble the 16 Kb bovine mitochondrial DNA (Anderson et al., 1982).

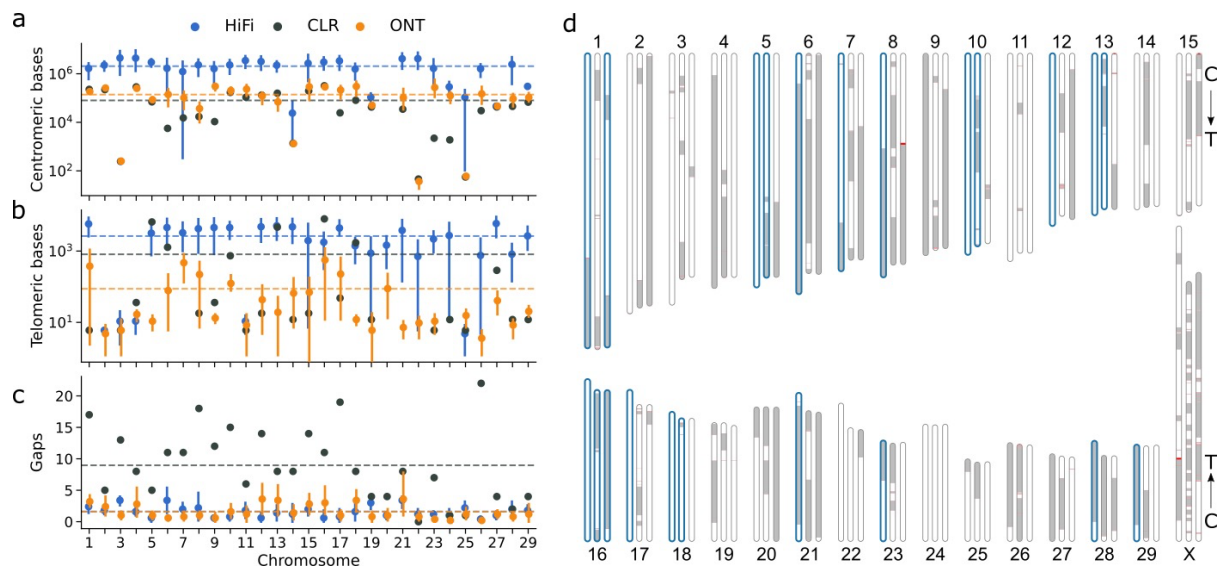


Figure 2. Centromeric and telomeric completeness of assemblies produced by hifiasm and Shasta. a) The mean number of bases identified per autosome as “Satellite” by RepeatMasker for the five hifiasm (blue) and five Shasta (orange) assemblies, where error bars indicate the 95% confidence interval. The black dots represent values from the CLR-based ARS-UCD1.2. Dashed lines indicate the autosome-wide mean for the respective color of points. Mean values of 0 (e.g., chromosome 20) are not shown due to the log scale. b) Similar to a), but the number of bases in telomeric repeats within 10 Kb of chromosome ends. c) Similar to a), but the number of scaffold gaps. d) Chromosome ideograms for ARS-UCD1.2 (center), and Brown Swiss assemblies produced by hifiasm (left) and Shasta (right). Scaffolded contigs alternate white/grey across gapped regions, which are colored red. Chromosomes which are predicted to extend from centromere to telomere are bolded in blue. Arrows indicate the centromere-to-telomere directionality of the chromosomes.

Optimal sequencing coverage depths

Impact of sequencing depth on assembly quality was assessed using the Brown Swiss haplotype of the NxS as an example. Subsets of the 52x diploid HiFi and 55x haploid (trio binned) ONT reads, respectively, were randomly sampled to mimic lower sequencing depths (Figure 3). Completeness metrics (e.g., BUSCO or k-mer content) plateau when coverage increased above 25-fold, however other metrics like contiguity or correctness continued to benefit from higher sequencing coverage with diminishing returns. The trio aware mode of hifiasm only required about 19x diploid coverage of HiFi reads to meet the VGP targets. Shasta required around 28x haploid coverage (corresponding to roughly 56x diploid coverage) to achieve the necessary QV, although 17x haploid (34x diploid) coverage fulfills the contiguity and completeness targets. The minimum VGP-satisfying coverage varies slightly for the different F1s due to different input sequencing read properties but is approximately consistent across all examined bovine trios (Supplementary Figure 2).

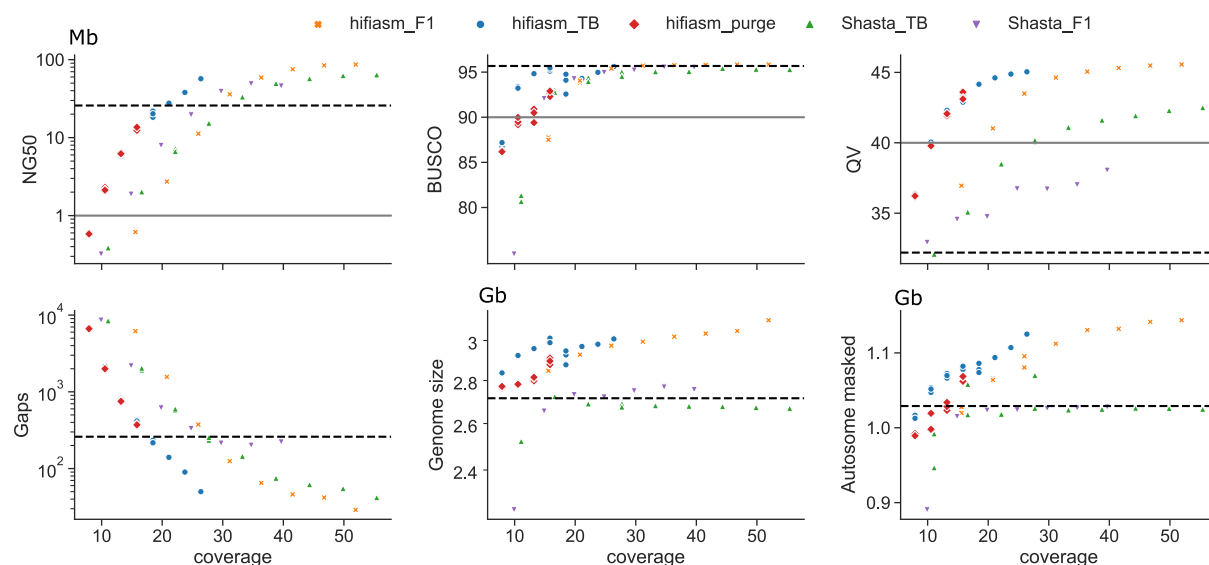


Figure 3. Assembly quality at subsampled coverages. Trio aware hifiasm (hifiasm_F1) uses diploid coverage while Shasta (Shasta_TB) uses haploid coverage. We additionally examined trio binned hifiasm (hifiasm_TB) using haploid coverage and the polish-phased Shasta approach using diploid coverage (Shasta_F1). NG50, BUSCO, QV, Autosomal gaps, and Genome size are defined in Table 1, while Autosomal masked is the number of autosomal bases within repetitive elements as identified by RepeatMasker. The black dashed line represents the relevant value for the ARS-UCD1.2 reference and the grey solid line is the VGP target where applicable. Three subsampling replicates were performed for lower coverage assemblies (<20x for HiFi and <30x for ONT) due to their higher stochasticity. For trio binned hifiasm assemblies below 15x coverage, we manually set the duplication purging parameter (hifiasm_purge) and reran on the same subsamplings.

The higher accuracy of HiFi reads allows hifiasm to exploit sequence common to both haplotypes during trio aware assembly. Reaching a comparable quality through a trio binned approach required approximately 16% higher coverage, with 11x haploid coverage (22x diploid) necessary. The higher error rate of ONT reads makes haplotype-aware correction and phasing challenging, and so diploid assembly followed by haplotype separation or diploid-aware polishing is less effective than trio binning (Figure 3, up triangles versus down triangles). Phasing is particularly poor, with the PG50 25 times smaller on average compared to trio binning approaches.

There is an increased risk of coverage gaps when sequencing coverage is reduced, even if the resulting assemblies achieve certain genome-wide standards. When 11x HiFi coverage is aligned to ARS-UCD1.2 and binned into 10 Kb windows, there are 950 regions on average of near total dropout (<1x coverage) across the autosomes. This drops by 23% at 13.5x coverage and by 30% at 16x coverage, as the effects of stochasticity are reduced. While the overall assembly quality does not fluctuate substantially between random subsamplings (Figure 3), it may overestimate the quality at specific regions. Furthermore, default parameters of assemblers are typically tuned to higher coverage and assembling at low coverage can introduce subtle issues. Hifiasm could underestimate a parameter related to duplication purging at coverages below 15x, resulting in a sharp transition to larger assemblies with more duplicated BUSCO genes (Figure 3, blue circles versus red diamonds). Manually setting the parameter to its expected value recovered similar behavior seen in higher coverage assemblies.

Constructing a bovine pangenome

Sequencing technologies and assemblers evolve rapidly, and so even recently generated bovine assemblies, including the ones reported here, have been produced under non-uniform conditions (e.g., (Crysanto et al., 2021; Talenti et al., 2021)). Given the differences we observed between HiFi- and ONT-based assemblies, especially in comparison to the CLR-based ARS-UCD1.2 reference, it is crucial to examine how pangenomes respond to different assembly inputs.

Pangenomes were constructed with minigraph (Li et al., 2020) using ARS-UCD1.2 as the initial backbone of the graph structure. Assemblies were iteratively added into the graph, and regions of synteny were ignored while sufficiently diverged subsequences (>50 bp) instead augmented the graph with new nodes (“bubbles”). Pangenomes for each autosome were constructed from all hifiasm assemblies, all Shasta assemblies, or random mixtures of the two. Graph properties of each pangenome in terms of the amount of non-reference sequence added, were generally robust to the input assemblies (Figure 4a,b). Pangenomes constructed from five hifiasm assemblies had more non-reference sequence added compared to five Shasta assemblies (82.5 Mb across 88.5k nodes versus 63.5 Mb across 90.2k nodes), in agreement with the greater completeness observed in hifiasm assemblies. Approximately 92.1% of identified SVs were common between hifiasm and Shasta pangenomes, with 3.6% and 4.3% unique to each respectively (Figure 4c). There was not a clear bias between HiFi- and ONT-based pangenomes, with only 1.8 and 0.39 Mb non-reference sequence uniquely identified in each. These bubbles were more repetitive than autosome-wide averages (53% and 45% respectively, Supplementary Table 5). Minigraph is sensitive to the order of integration, and on rare occasions constructed significantly different bubbles, particularly for palindromic sequences, resulting in larger variance on chromosomes 7 and 12 (Supplementary Figure 3).

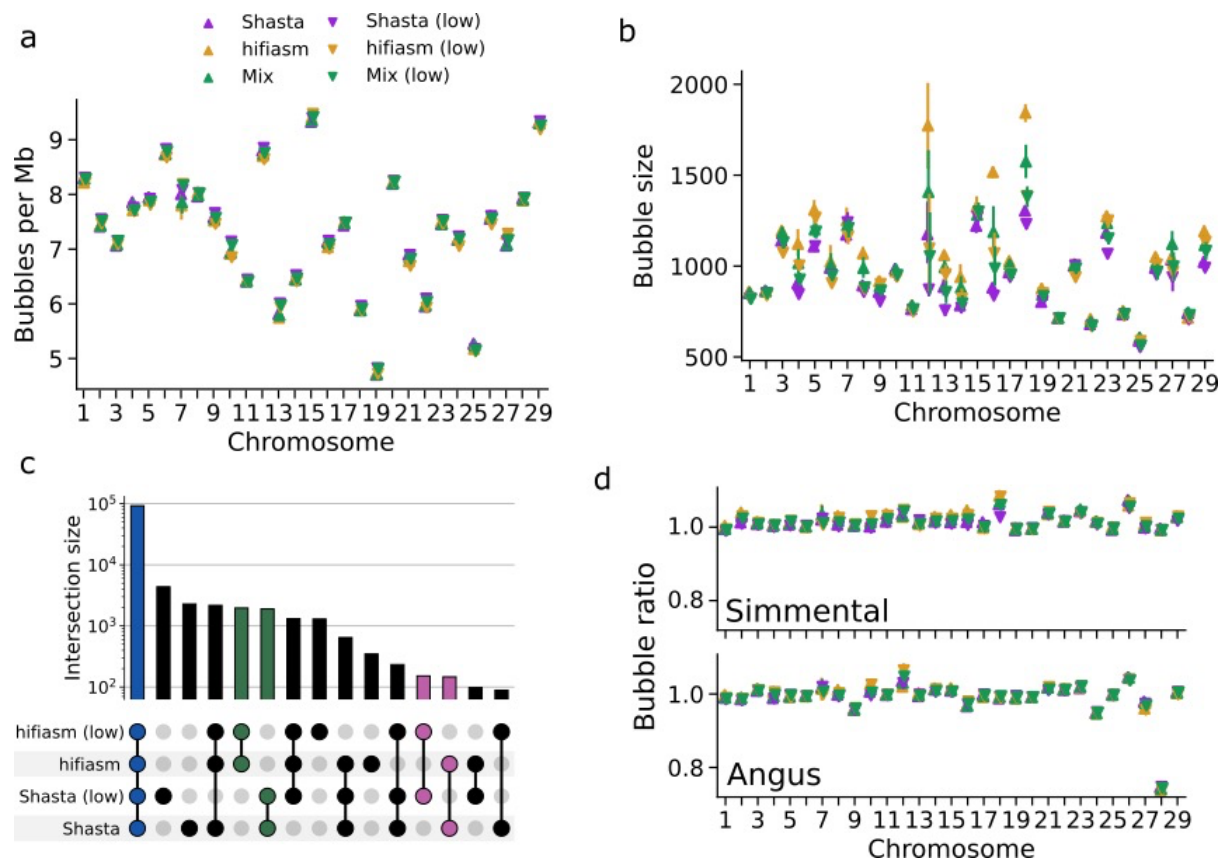


Figure 4. Pangenomes are generally robust to different input assemblies. a) The number of large (>1 Kb) bubbles is highly consistent across hifiasm, Shasta, and mixed pangenomes at both full (up triangles) and lower (down triangles) coverages. b) The mean bubble size is also consistent across different inputs, but bubbles are larger on average in hifiasm pangenomes compared to Shasta pangenomes. c) The vast majority (84.5%) of SVs identified through minigraph are present in all pangenomes (blue). SVs unique to either hifiasm or Shasta (green) only account for about 3.5% of all SVs, while SVs only identified through either full or lower coverage pangenomes are negligible (pink). d) Comparing the number of bubbles present in Simmental- or Angus-backed pangenomes to the ARS-UCD1.2-backed pangenome in a) shows consistency. Angus chromosome 28 is the only exception due to its incomplete sequence. All points reflect the mean over 20 stochastic pangenome constructions, and error bars indicate the 95% confidence interval.

Pangenomes constructed from lower-coverage assemblies remained robust. We selected hifiasm and Shasta assemblies generated with an average of 21.6x diploid and 24.9x haploid coverage respectively, which approximately satisfy the VGP standards. The hifiasm assemblies have a higher average QV compared to the Shasta assemblies (41.9 versus 39.3), but lower NG50 (2.3 Mb versus 19.5 Mb) and more autosomal gaps (1869 versus 323) (Supplementary Table 6). Although these assemblies are substantially worse than their full-coverage counterparts, the resulting pangenomes are similar (Figure 4a-c). Taking the high coverage SVs as the truth set, the low coverage hifiasm and Shasta pangenomes have an F1 score of 98.5 and 94.7 for SV discovery respectively. The low coverage Shasta pangenomes tend to identify a greater number of SVs not present in other pangenomes, and so may be false positives.

Pangenomes constructed using the existing Angus (CLR-based) (Low et al., 2020) or Simmental (ONT-based) (Heaton et al., 2021) assemblies as backbones (Figure 4d) produced similar results compared to using the Hereford-based ARS-UCD1.2. More non-reference sequence was identified in the lower-quality Angus-backed pangenome (+13%), while the more complete Simmental-backed pangenomes had less (-6%). Reference-bias propagates through minigraph's pangenomes, such as the missing sequence in Angus chromosome 28

(Lloret-Villas et al., 2021) resulting in 25% fewer bubbles compared to using ARS-UCD1.2 (Figure 4d).

Quantifying bovine structural diversity through the pangenome

Structural variation bubbles were associated with their source assembly by retracing each “haplotype walk” through the graph. The phylogenetic topology of their evolutionary relationship was then estimated by counting the number of mutually exclusive bubbles any two assemblies have to construct a condensed distance matrix. The results are consistent with expectations, with the gaur (*Bos gaurus*) largely separate, followed by the Nellore (*Bos taurus indicus*), and then the three taurine cattle (Figure 5a). All constructed pangenomes unequivocally predicted the first two branches, while there was also good agreement within the closely related taurine cattle (Figure 5b).

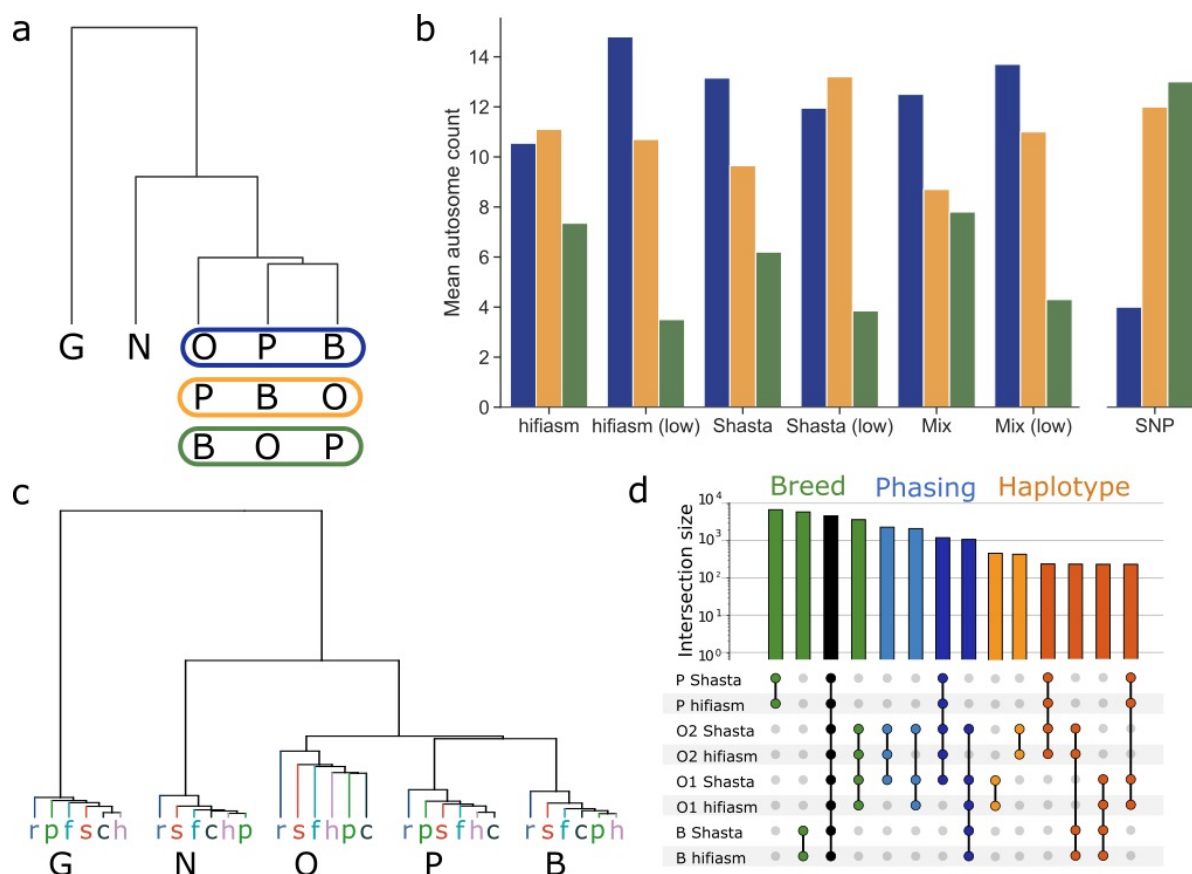


Figure 5. SV-based dendrograms. *a*) All dendrograms followed the same overall topology, with gaur (G) and Nellore (N) clearly differentiated while the taurine cattle displayed three possible arrangements, with either the Original Braunvieh (O), Piedmontese (P), or Brown Swiss (B) more distantly related. *b*) Across 20 constructed pangenomes with all hifiasm, all Shasta, or mixing hifiasm and Shasta assemblies, as well as the low coverage equivalents, there was good agreement on the average number of autosomes displaying a specific topology (color from panel a). The SNP dendrogram, based on parental short reads, generally predicted different topologies. *c*) Pangenomes including all 30 assemblies from hifiasm (h), Shasta (s), Peregrine (p), Flye (f), HiCanu (c), and Raven (r) predict the same overall topology without ONT or HiFi specific branches. *d*) An UpSet plot of a taurine cattle pangenome (Piedmontese, Brown Swiss, and the paternal [O1] and maternal [O2] haplotypes of Original Braunvieh) reveals inter-breed variation (green) as well as intra-breed variation in the Original Braunvieh haplotypes (red & orange). We can also identify phasing error candidates in the Original Braunvieh Shasta assemblies, where SVs are common to both Shasta assemblies but not both hifiasm assemblies (light and dark blue).

Several chromosomes (e.g., 2, 6, 9, etc.) are observed to only predict a single taurine topology, even across different assembly inputs and coverages, while other chromosomes (e.g., 1, 12, 16, etc.) have multiple predicted topologies with similar frequencies. This may indicate that certain chromosomes harbor greater structural variation between specific taurine breeds, which may be reflected in phenotypic differences. We compared against a conventional approach, calling small variants from the parental short reads corresponding to the five haplotypes against ARS-UCD1.2 (Figure 5b). While the overall topology and magnitude is similar to that found through the pangenome, there was no significant concordance within the taurine cattle, as might be expected given the low linkage disequilibrium between small and structural variants (Yan et al., 2021).

Genome level dendrograms were mostly consistent across different combinations of inputs, demonstrating that pangenome robustness extends beyond graph properties and into applications. The same general topology was recovered for all examined assemblers, both individually and combined, as well as the lower coverage assemblies. The dendrogram correctly places all assemblies of each breed on their own branch (Figure 5c) when 31 pangenome assemblies are included (1 reference backbone plus 5 breeds/species x 6 assemblers). Some HiFi- and all ONT-based assemblies for Original Braunvieh are only partially phased due to limited coverage, resulting in greater structural variation within these assemblies (Figure 5d). The variation between the parental haplotypes of the OxO results in the higher branching point for “O” in Figure 5c, and highlights the importance of cleanly resolved haplotypes for downstream analyses.

Pangenome topology at trait-associated SVs

Pangenome integration of haplotype-resolved assemblies representing multiple breeds/species supported investigation of the evolutionary history of a multi-allelic copy number variation (CNV) at 86.96 Mb on BTA6 encompassing an enhancer of the group-specific component (*GC*) gene. This CNV has pleiotropic effects on mastitis resistance and other dairy traits, and segregates in different breeds of cattle (Lee et al., 2021; Olsen et al., 2016; Pausch et al., 2016), although the prevalence of the CNV has not yet been determined for the breeds/species included in our bovine pangenome. The reported CNV was only observed in the Shasta assemblies of Brown Swiss and Original Braunvieh cattle (Figure 6a,b) indicating it may contribute to variation in somatic cell score also in these breeds (Fang & Pausch, 2019). The CNV formed part of a complex superbubble (Paten et al., 2018), suggesting substantial allelic diversity between the haplotype-resolved assemblies. In addition to the two (Brown Swiss) and three (Original Braunvieh) copies of the 12 Kb CNV segment, there were also two shorter insertion elements detected within the superbubble from the Nellore and gaur assemblies (Figure 6a,b).

Inspection of coverage data retrieved from HiFi- and ONT-binned long read alignments against the ARS-UCD1.2 reference suggests that the Brown Swiss and Original Braunvieh haplotypes harbor between 2 and 4 additional copies of the 12 Kb segment (Figure 6c, Supplementary Notes). Coverages derived from F1 short read alignments are consistent with these values but are unable to resolve the disagreement observed between HiFi and ONT coverage (Figure 6d). While most examined assemblers predicted duplication of the 12 Kb segment in Brown Swiss and Original Braunvieh, there was poor consensus in copy number, supporting that this region remains challenging to resolve even with long reads and

de novo assembly (Supplementary Figure 4). The same inspection did not identify the duplication in the gaur, Piedmontese and Nellore haplotypes, but did validate the two insertions of additional repeat elements exclusively in the gaur and Nellore haplotypes. We retrieved various repetitive elements in the tandemly duplicated 12 Kb segments of the haplotype-resolved Brown Swiss and Original Braunvieh assemblies (Figure 6e) broadly confirming the results of (Lee et al., 2021).

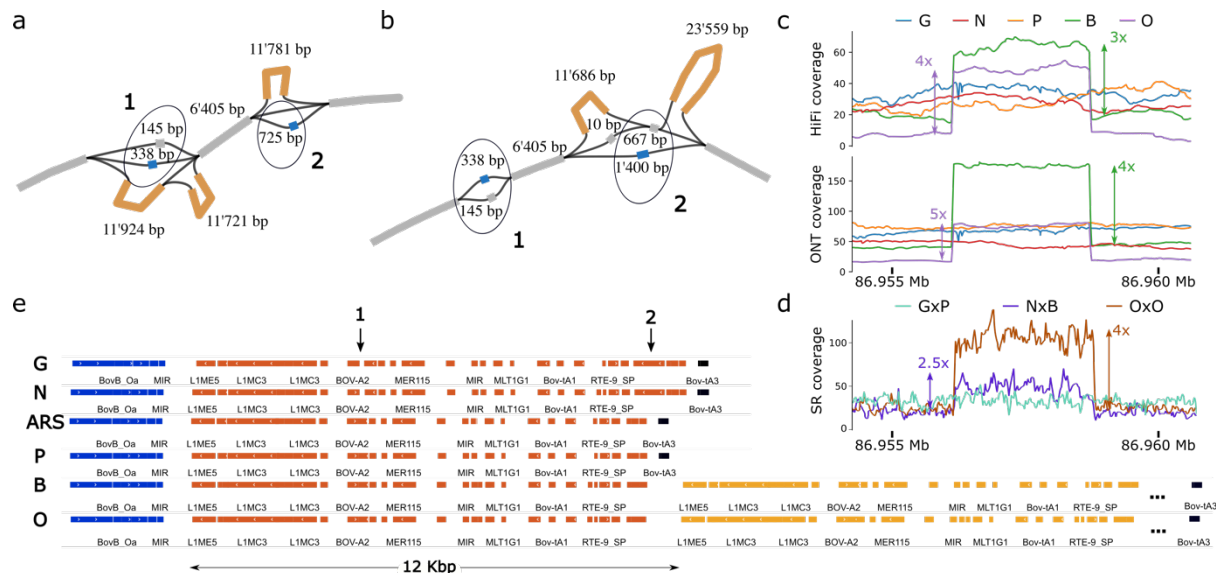


Figure 6. Topology of a tandem duplication on BTA6. a-b) Two example pangenome subgraphs of the promoter region of GC. Reference paths (including those in bubbles) are colored gray, while the tandem duplications are orange and the insertions are blue. Complex superbubbles generally have suboptimal topologies due to the lack of base-level alignment. For example, the 725 bp insertion is obvious in a), but appears as the difference between a 1,400 bp and 667 bp path in b). However, both subgraphs identify the approximately 200 bp (1) and 700 bp (2) insertions in Nellore and gaur, as well as the tandem duplication in Brown Swiss and Original Braunvieh. c) The identified CNV region shows clear coverage increase in only Brown Swiss and Original Braunvieh, across both HiFi and ONT haplotype-resolved reads. d) F1 short reads also show increased coverage for the NxP and OxO trios. The NxP coverage increase is consistent with only the Brown Swiss haplotype carrying additional copies. e) The 12 Kb repeat structure (orange) is clearly identified by RepeatMasker across all assemblies, shown here for the ARS-UCD1.2 reference and Shasta assemblies for gaur, Nellore, Piedmontese, Brown Swiss, and Original Braunvieh. The two marked gaur/Nellore insertions (1&2) are consistent with the pangenomes in a-b). One additional copy in Brown Swiss and Original Braunvieh is shown (yellow), while the tandem duplication eventually ends with a similar repeat (Bov-tA3, black) to the other assemblies.

An SV at the *ASIP* gene, encoding the Agouti-signaling protein involved in mammalian pigmentation and located on chromosome 13, was also investigated. This genomic sequence harbors alleles associated with coat color variation in many species including cattle (Girardot et al., 2006; Trigo et al., 2021). Both Nellore and Brown Swiss cattle present variability in coat color ranging from near white to almost black. The NxP was born with a light coat, which darkened as the bull aged (Supplementary Notes). Our pangenome confirmed great allelic diversity between the bovine assemblies involving insertions and deletions of repetitive elements upstream of *ASIP*. However, the previously described variants associated with coat color variation (Trigo et al., 2021) were not in the pangenomes. Short read alignments of the Nellore sire confirmed it carried the previously described SV in the heterozygous state, but the F1 inherited the other haplotype (Supplementary Notes). Thus, the darkening of the coat we observed in the NxP is not due to previously described *ASIP* alleles.

We identified 900 and 922 genes in the ARS-UCD1.2 genome annotation (Refseq release 106) whose coding regions overlap bubbles in hifiasm and Shasta based pangenomes respectively (Supplementary Figure 5). Of these, 808 and 845 are bubbles with haplotype path information for each assembly, and so are amenable to association investigations. Several of these genes are listed in the OMIA (Online Mendelian Inheritance in Animals) database as they harbor alleles causing phenotypic variation in cattle; for instance, both hifiasm and Shasta pangenomes contained a bubble encompassing a 36 bp deletion in *ACAN* resulting in the in-frame deletion of 12 amino acids in the gaur assemblies. Variants in *ACAN* are associated with stature in cattle (Cavanagh et al., 2007) and other species (Gibson & Briggs, 2016). The pangenome also recovered the insertion of an 11 Kb segment on chromosome 23 in all assemblies. This segment encompasses *HSPA1B* (Hess et al., 2018) which is not annotated in ARS-UCD1.2 (Suqueli García et al., 2017). This 11 Kb insertion is challenging to identify by inspection of short and long read alignments, mainly appearing as elevated coverage over the 2294 bp segment encompassing *HSPA1A*, indicating a similarly sized duplication, with soft-clipped bases extending on both sides (Supplementary Note).

The bubbles in the pangenome indicate further putatively trait-associated regions. For instance, the coding sequence of *QRICH2* overlapped with multiple bubbles indicating tandem duplications of a 30 bp region of the fifth exon (Figure 7a). Loss of function alleles in mammalian *QRICH2* orthologs lead to multiple morphological abnormalities of the sperm flagella (Shen et al., 2019). We find that the fifth exon of *QRICH2*, which is affected by the coding sequence expansion, is transcribed in high abundance (>30 transcripts per million) in testes of mature taurine bulls. Inspection of long read alignments confirms an expansion of the coding sequence relative to ARS-UCD1.2 that extends the high molecular weight glutenin subunit of the protein by 10 amino acids in our taurine cattle, 60 amino acids in Nellore, and 50 amino acids in gaur (Figure 7b,c). This SV is challenging to resolve with short reads (Supplementary Notes). Another example of potential trait-associated SV lies in *TAS2R46*, related to bitter taste receptors and associated with adaptation to dietary habitats (Dong et al., 2009), which overlapped with a 17 Kb deletion in gaur (Figure 7d,e). This deletion also spanned ENSBTAG000000001761. A final example is found in *PRDM9*, the only known speciation gene in mammals and known to harbor alleles associated with variation in meiotic recombination within and between Bovinae (Ahlawat et al., 2016; Sandor et al., 2012; Zhou et al., 2018), where an SV overlapped with copy number variation in the zinc finger array domain (Figure 7f,g). The Nellore and gaur assemblies contained one zinc finger less, while the paternal haplotype of Original Braunvieh carried one more relative to ARS-UCD1.2. The maternal haplotype of Original Braunvieh contained the same number of zinc fingers as ARS-UCD1.2, supporting the intra- and inter-breed/species variation observed for *PRDM9*.

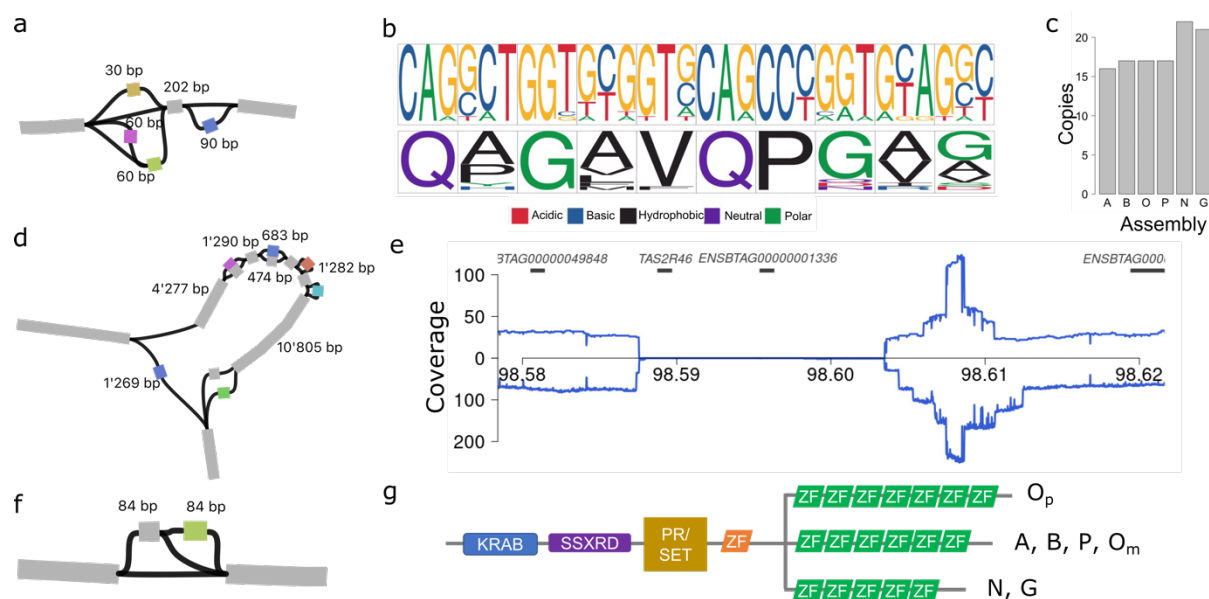


Figure 7. Identification of structural variation in coding sequences of QRIH2, TAS2R46 and PRDM9 through the HiFi-based pangenome. a) Pangenome topology in the fifth exon of bovine QRIH2 revealed tandem repeats of 30 bp sequence. b) Nucleotide (upper) and protein (lower) sequence logo plot of the repeat motif. c) While the ARS-UCD1.2 reference sequence contains 15 copies of the repeat motif, the pangenome revealed 1, 5, and 6 additional copies in the five haplotype-resolved assemblies (A – ARS-UCD1.2, B – Brown Swiss, O – Original Braunvieh, P – Piedmontese, N – Nellore, G – gaur). d) Representation of a 17 kb deletion on BTA5 encompassing TAS2R46 and ENSBTAG00000001761. e) Coverage of binned HiFi (above horizontal line) and ONT (below horizontal line) long read alignments in gaur indicate a large deletion between 98,587,384 and 98,604,401 bp. f) Pangenome topology at the eleventh exon of PRDM9 indicating paths with gain and loss of 84 bp sequence. g) Representation of the domains of PRDM9 in the haplotype-resolved assemblies including a variable number of zinc fingers (ZF) in the different assemblies, where O_m and O_p are the maternal and paternal haplotypes of the *OxO*.

Discussion

Ten haplotype-resolved assemblies for cattle and related species were constructed using recent sequencing and assembling technologies. Assemblies produced by hifiasm (HiFi) and Shasta (ONT) were substantially more contiguous and correct than the current *Bos taurus* reference sequence, which is a haplotype-merged assembly based on older CLR technology. The higher accuracy of HiFi reads was found to be generally more advantageous to quality measures of assembly and increased completeness of centromeric and telomeric regions compared to the longer but higher-error ONT reads. HiFi-based assemblers also required less compute and storage resources compared to ONT-based assemblers; producing haplotype-resolved hifiasm assemblies required approximately 600 CPU hours and 200 GB of peak memory usage, while the equivalent Shasta (plus polishing) assemblies took 2200 CPU hours and 750 GB of peak memory usage. Correct-then-assemble approaches like Canu (Koren et al., 2017) can be practical for smaller genomes (Wick & Holt, 2021), but on gigabase-sized mammalian genomes like in Bovinae we observed >20 Tb of peak temporary storage and >25k CPU hours for correcting only 30-fold ONT reads. Even recent reference-guided correction approaches like Ratatosk (Holley et al., 2021) still needed approximately 15k CPU hours to correct 55-fold ONT reads. Cutting-edge sequencing and bioinformatic improvements (Baid et al., 2021; Silvestre-Ryan & Holmes, 2021), like the ONT Guppy5 basecaller, will likely assist more efficient assembly, resulting in higher QV and reduced computational load; however, currently the ONT specific requirements might be computationally prohibitive, especially when assembling many samples.

The phased assembly graph approach of hifiasm is most efficient with lower heterozygosity samples, where HiFi reads are least sortable by parental haplotype and there is more mutual sequence to exploit, but still functions with highly heterozygous F1s. The inter-species GxP trio binned assemblies were more contiguous (+25% NG50) compared to the trio aware assemblies, but contained several large (>5 Mb) misassemblies which the latter did not (Supplementary Figure 6). Phasing in both HiFi- and ONT-based assemblies improved as heterozygosity increased and allowed more cleanly resolved haplotypes, which can be beneficial to downstream analyses (Feng & Li, 2021; Yang & Chaisson, 2021). The ability of hifiasm, and to a lesser extent Shasta with diploid-aware polishing, to assemble phased haplotypes from purebred individuals also avoids ethical and logistical concerns regarding the higher heterozygosity crosses previously targeted (Koren et al., 2018).

The minigraph pangenomes were strongly comparable whether the input assemblies were all HiFi-based, ONT-based, or a mix. The greater completeness observed in hifiasm assemblies is reflected in those pangenomes containing more non-reference sequence, but no notable HiFi or ONT specific biases were observed in the pangenomes. The quality and completeness of the pangenome backbone can have an impact, seen on the incomplete Angus chromosome 28 or the ARS-UCD1.2 chromosomes generally lacking centromeric sequence, but again we found no specific bias between CLR- or ONT-based backbones. These results indicate that optimal minigraph pangenomes would use high-quality, complete genomes as the backbone, like emerging T2T assemblies (Belser et al., 2021; Nurk et al., 2021; Rengs et al., 2021). Alternatively, reference-free approaches to pangenomes (Armstrong et al., 2020) may also circumvent these issues.

Mutual variation identified through shared paths in the pangenome provides opportunities to study the phylogeny of Bovinae beyond SNPs and indels. The ability to accurately separate and represent paths within SVs also enables pangenome-based GWAS (PWAS) (Gupta, 2021), as recently explored in crops (Della Coletta et al., 2021; Hufford et al., 2021; Song et al., 2020). We identified multiple SVs overlapping annotated coding sequences in different bovine pangenomes, demonstrating that pangenomes provide a framework to make them amenable to association mapping. Furthermore, some of these SVs (e.g., tandem duplications in *QRICH2*) are inaccessible from short or noisy long read alignments and some (e.g., an insertion of *HSPA1B*) are challenging to resolve even with long read alignments. These cases highlight the benefits of *de novo*, haplotype-resolved assemblies and pangenome integration.

Two haplotype-resolved assemblies which satisfy VGP quality standards can be produced with diploid coverage of approximately 20x for hifiasm (HiFi) or 60x for Shasta (ONT). These lower coverage assemblies are nearly as efficient (90+%) for SV discovery as the full coverage assemblies, with approximately 88% efficiency for genic SV identification. Combined, these results support that cataloguing most within-breed structural variation is a tractable goal through *de novo* long read assembly of 50 samples per breed, sufficient to represent the effective population size. Extensive existing short-read sequencing data can then be leveraged to genotype SVs present in the pangenome (Ebler et al., 2020; Hickey et al., 2020) and then imputed into tens of thousands of cattle previously genotyped with microarrays.

Methods

Ethics statement

The sampling of blood from the NxB and OxO trios was approved by the veterinary office of the Canton of Zurich (animal experimentation permit ZH 200/19). All GxP protocols were approved by the Institutional Animal Care and Use Committee (IACUC) of the University of Nebraska–Lincoln, an AAALAC International Accredited institution (IACUC Project ID Project ID 1697). Gaur semen collections were approved by the Henry Doorly Zoo IACUC in 1992.

Animals

Cows from the Original Braunvieh (O), Brown Swiss (B), and Piedmontese (P) breeds were inseminated with semen samples from Original Braunvieh, Nellore (N), and gaur (G) sires, respectively. A female (OxO) and a male (NxB) calf were delivered at term. A female fetus (GxP) cross was collected by cesarean section at 119 days of gestation. Blood samples were taken from the calves and dams by trained veterinarians. Lung tissue was taken from the fetus and DNA extracted as described (Logsdon, 2020). Semen samples were available for the bulls. High-molecular weight DNA was extracted from blood and semen samples, respectively, using Qiagen's MagAttract HMW DNA Kit as described earlier (Crysnanto et al., 2021).

Sequencing

The genomes of the F1s were sequenced with long reads using PacBio HiFi and ONT. The PromethION 1D Genomic DNA by ligation SQK-LSK110 library prep kit was used for the OxO, NxB, and GxP F1s. The libraries were respectively sequenced on PromethION (R9.4.1) flowcells (nuclease wash was applied to one cell). The OxO and NxB reads were basecalled on Guppy4, while the GxP reads were basecalled with Guppy3.

Paired-end libraries (2x150 bp) were produced using parental DNA samples and sequenced on Illumina instruments.

Genome assembly

HiFi reads were first filtered with fastp (version 0.21.0) (S. Chen et al., 2018), removing reads below 1 Kb length or QV20. Nanopore reads were pre-filtered to a minimum of QV7, with no length restrictions. Short reads were also filtered with fastp, using the “-g” parameter to trim polyG tails.

Both HiFi and ONT reads were trio binned into paternal and maternal haplotypes using (Trio)Canu (Koren et al., 2018) (version e0d6bb0), using parental short reads and an estimated genome size of 2.7g.

HiCanu (Nurk et al., 2020) was used with default settings and “genomeSize=2.7g -pacbio-hifi” to produce a draft set of contigs from the HiFi reads. Only contigs which were not suggested as bubbles (suggestBubble=no) were retained, which improved automated

coverage detection when using `purge_dups` (version 1.2.3) (Guan et al., 2020). The final set of “assembled” contigs was obtained after following the default purging pipeline.

The trio aware hifiasm (Cheng et al., 2021) (version 0.15.3-r339) assemblies were produced from HiFi reads with non-default parameters “-a 5 -n 5”. The resulting graph information was then used in the trio-mode, with parental k-mers constructed with `yak` (version 0.1-r62-dirty) specified through the “-1 {paternal} -2 {maternal}” parameters. The haplotype gfa files were then converted to fasta contigs with `gfa2fa` command of `gfatools` (version 0.5-r234).

The trio binned hifiasm assemblies were obtained with the same assembly parameters, with the addition of the “--primary” flag to produce a primary and alternate assembly. The primary gfa file was converted to fasta as described above. Only the primary assembly was retained.

Peregrine (Chin & Khalak, 2019) (version main:2aefc14+) was used with default settings. Only the primary assembly was retained.

Shasta (Shafin et al., 2020) (version 0.7) assemblies were produced using the standard configuration file “Nanopore-Sep2020.conf”. The assemblies used the haplotype-binned nanopore reads, except for the OBV assembly, which used the total read set. Due to the high coverage for the gaur and Piedmontese haplotypes, the parameter “minReadLength=30000” was used instead.

All nanopore assemblies except the OBV Shasta were polished using PEPPER (Shafin et al., 2021) (version 0.1) and the haplotype-binned nanopore reads. Each assembly was then polished with `bcftools` (Danecek et al., 2021) (version 1.12) using the `merfin` (Formenti et al., 2021) (commit version 1331fa5) filtered vcf output of `DeepVariant` (Poplin et al., 2018) (version 1.1). The `DeepVariant` input was F1 short reads where reads containing k-mers specific to the other haplotype were excluded with `meryl`. The OBV Shasta assembly was polished with PEPPER-DV (version 0.4), tagging the intermediate bam file with haplotype-binned read IDs. The subsequent short read polishing was as described above.

Raven (Vaser & Šikić, 2021) (version 1.5.0) assemblies were obtained from nanopore reads with default parameters, except for “-p 0”. Polishing with the default two rounds of `racon` (Vaser et al., 2017) resulted in slightly lower QV compared to *post hoc* polishing with PEPPER.

Flye (Kolmogorov et al., 2019) (version 2.8.3-b1725) assemblies were constructed with “--genome-size=2.7g --nano-corr” from `Ratatosk` (version 0.4) (Holley et al., 2021) error-corrected nanopore reads. The nanopore reads were corrected using a reference-guided approach, taking the haplotype-specific hifiasm assembly as the reference. Ambiguous IUPAC codes were randomly replaced with equal probability of an appropriate nucleotide. The contig set was taken from the pre-scaffolding result of Flye and was not polished due to the pre-corrected reads having sufficient accuracy.

Scaffolding

Contigs were scaffolded into chromosomes by the reference-guided approach of RagTag (Alonge et al., 2019) (v2.0.1) to ARS-UCD1.2, using the additional parameters “--mm2-params "-cx asm5" -r -m 1000000”.

Quality metrics

Completeness was assessed with BUSCO (version 5.1.2), using the metaeuk backend (commit 9dee7a7) and odb10 cetartiodactyla database (e96dfc6299c567768085ee9569b6ab15). Quality value and phasing were calculated with merqury (version 1.3) (Rhie et al., 2020), with k-mer databases constructed from short reads using meryl (<https://github.com/marbl/meryl>) (version r953). Contiguity was assessed with calN50.js (<https://github.com/lh3/calN50>). Repetitive elements were identified by RepeatMasker (<http://www.repeatmasker.org>) (version 4.1.1) and rmbblast (version 2.10.0), using a modified version of the 2018 Repbase database. For whole genome analysis, the rush job mode was used, while for pangenome region analysis the slow mode was used.

Coverage downsampling

Sequencing subsets were made through seqtk (<https://github.com/lh3/seqtk>) (version 1.3-r115-dirty) with the command “seqtk seq -f {sample}”, including the “-A” flag to drop quality scores where appropriate. In the case of repeat trials, the “-s {seed}” flag was set with a randomly generated 64-bit integer to ensure a unique subsampling of data. The reduced coverage assemblies were conducted identically to the full coverage methods unless explicitly mentioned otherwise.

Coverage depth estimation

Sequencing reads were mapped against ARS-UCD1.2 with minimap2 (version 2.19-r1059-dirty) (Li, 2018). Coverage depth was determined by megadePTH (version 1.1.0c) (Wilks et al., 2021), and averages over 10 Kb windows were estimated with pyBigWig (version 0.3.18) (Ramírez et al., 2016).

Pangenome construction and bubble extraction

Pangenomes were constructed on a per chromosome basis using minigraph (version 0.15-r426) (Li et al., 2020), with default parameters. The selected assemblies were added to the graph in a randomly shuffled order, where repeated constructions from the same input set may differ due to the ordering. Pangenomes were visualized with Bandage (version 0.8.1) (Wick et al., 2015).

Haplotype paths were called with minigraph, using the “--call --xasm” option. Bubble intersections across assemblies were created with UpSetPlot (<https://github.com/jnothman/UpSetPlot>) (version 0.6.0) and converted into a condensed distance matrix by assessing when two assemblies did not take the same path through a bubble.

Pangenome bubble and gene overlapping

All entries in the ARS-UCD1.2 annotation (RefSeq release 106) with a “CDS” description were extracted into a bed file. Bubbles from the pangenome were also extracted into a bed file using “gfatools bubble”. The two bed files were then intersected using bedtools (version 2.30.0) (Quinlan & Hall, 2010) with the intersect command and “-wo” to find potential trait-associated SVs. We then filtered the list by checking all five breed/species were successfully assigned a haplotype path through the bubble, to remove incomplete associations.

SNP

Parental short reads were mapped to the ARS-UCD1.2 reference using BWA-mem2 (Li, 2013; Md et al., 2019). Variants were called using DeepVariant (version 1.1) with the “WGS” model and merged using GLnexus (version 1.3.1) (Lin et al., 2018). Phylograms were constructed per chromosome with vcft-kit (version 0.2.9) (Cook & Andersen, 2017), using the command “vk phylo tree upgma”.

Miscellaneous analysis

Assembly, validation, and pangenome workflows have been implemented using Snakemake (Köster et al., 2021), and are available online (<https://github.com/AnimalGenomicsETH/bovine-assembly>). Nucleotide and protein sequence logos were generated using the R package ggseqlogo (Wagih, 2017).

Availability of data

HiFi reads for the OxO and NxB F1s are available at the study accession PRJEB42335 under sample accession SAMEA7759028 and SAMEA7765441.

ONT reads for the OxO and NxB are available at the study accession PRJEB42335 under sample accession SAMEA10017983 and SAMEA10017982.

Short reads for the OxO and NxB are available under accession number SAMEA9986200 and SAMEA7589752. Parental short reads are available at SAMEA9986201 & SAMEA9986199 (OxO) and at SAMEA6163185 & SAMEA9533783 (NxB).

Long and short read sequencing data for the GxP trio are available at the study accession PRJEB48481 under secondary accessions SAMEA10563833, SAMEA10563834, and SAMEA10563835.

Acknowledgements

We thank Dr. Melissa Terranova and Flavio Ferrari (AgroVet-Strickhof) for animal handling and Dr. Sandra Milena Bernal Ulloa for sampling blood. We are thankful for the technical

support provided by Dr. Anna Bratus-Neuenschwander from the ETH Zürich technology platform FGCZ (<https://fgcz.ch>) for sequencing and DNA fragment analysis. The results reported here were made possible with resources provided by the USDA shared compute cluster (Ceres) as part of the ARS SciNet initiative. We thank the USMARC Core Facility staff for outstanding technical assistance. Also thank B. Lee, J. Carlson, K. McClure, H. Clark, H. Sadd, M. Sadd, and B. Shuck for outstanding technical support. We thank the North American Piedmontese Association for their enthusiastic support and assistance.

Funding

This work was financially supported from the Federal Office for Agriculture (FOAG), Bern, Switzerland, and the Swiss National Sciences Foundation (SNSF).

Competing Interests

The authors declare no competing interests.

Bibliography

- Ahlawat, S., De, S., Sharma, P., Sharma, R., Arora, R., Kataria, R. S., Datta, T. K., & Singh, R. K. (2016). Evolutionary dynamics of meiotic recombination hotspots regulator PRDM9 in bovids. *Molecular Genetics and Genomics* 292:1, 292(1), 117–131. <https://doi.org/10.1007/S00438-016-1260-6>
- Alonge, M., Soyk, S., Ramakrishnan, S., Wang, X., Goodwin, S., Sedlazeck, F. J., Lippman, Z. B., & Schatz, M. C. (2019). RaGOO: fast and accurate reference-guided scaffolding of draft genomes. *Genome Biology* 20:1, 20(1), 1–17. <https://doi.org/10.1186/S13059-019-1829-6>
- Anderson, S., de Bruijn, M. H. L., Coulson, A. R., Eperon, I. C., Sanger, F., & Young, I. G. (1982). Complete sequence of bovine mitochondrial DNA conserved features of the mammalian mitochondrial genome. *Journal of Molecular Biology*, 156(4), 683–717. [https://doi.org/10.1016/0022-2836\(82\)90137-1](https://doi.org/10.1016/0022-2836(82)90137-1)
- Armstrong, J., Hickey, G., Diekhans, M., Fiddes, I. T., Novak, A. M., Deran, A., Fang, Q., Xie, D., Feng, S., Stiller, J., Genereux, D., Johnson, J., Marinescu, V. D., Alföldi, J., Harris, R. S., Lindblad-Toh, K., Haussler, D., Karlsson, E., Jarvis, E. D., ... Paten, B. (2020). Progressive Cactus is a multiple-genome aligner for the thousand-genome era. *Nature* 2020 587:7833, 587(7833), 246–251. <https://doi.org/10.1038/s41586-020-2871-y>
- Baid, G., Cook, D. E., Shafin, K., Yun, T., Llinares-Lopez, F., Berthet, Q., Wenger, A. M., Rowell, W. J., Nattestad, M., Yang, H., Kolesnikov, A., Topfer, A., Ammar, W., Vert, J.-P., Vaswani, A., McLean, C. Y., Chang, P.-C., & Carroll, A. (2021). DeepConsensus: Gap-Aware Sequence Transformers for Sequence Correction. *BioRxiv*, 2021.08.31.458403. <https://doi.org/10.1101/2021.08.31.458403>
- Belser, C., Baurens, F.-C., Noel, B., Martin, G., Cruaud, C., Istace, B., Yahiaoui, N., Labadie, K., Hřibová, E., Doležel, J., Lemainque, A., Wincker, P., D'Hont, A., & Aury, J.-M. (2021). Telomere-to-telomere gapless chromosomes of banana using nanopore sequencing. *Communications Biology* 2021 4:1, 4(1), 1–12. <https://doi.org/10.1038/s42003-021->

02559-3

- Blazak, W. F., & Eldridge, F. E. (1977). A Robertsonian Translocation and its Effect upon Fertility in Brown Swiss Cattle. *Journal of Dairy Science*, 60(7), 1133–1142. [https://doi.org/10.3168/jds.S0022-0302\(77\)83999-4](https://doi.org/10.3168/jds.S0022-0302(77)83999-4)
- Cavanagh, J. A. L., Tammen, I., Windsor, P. A., Bateman, J. F., Savarirayan, R., Nicholas, F. W., & Raadsma, H. W. (2007). Bulldog dwarfism in Dexter cattle is caused by mutations in ACAN. *Mammalian Genome*, 18(11), 808–814. <https://doi.org/10.1007/s00335-007-9066-9>
- Chen, N., Cai, Y., Chen, Q., Li, R., Wang, K., Huang, Y., Hu, S., Huang, S., Zhang, H., Zheng, Z., Song, W., Ma, Z., Ma, Y., Dang, R., Zhang, Z., Xu, L., Jia, Y., Liu, S., Yue, X., ... Lei, C. (2018). Whole-genome resequencing reveals world-wide ancestry and adaptive introgression events of domesticated cattle in East Asia. *Nature Communications*, 9(1), 1–13. <https://doi.org/10.1038/s41467-018-04737-0>
- Chen, S., Zhou, Y., Chen, Y., & Gu, J. (2018). fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics*, 34(17), i884–i890. <https://doi.org/10.1093/BIOINFORMATICS/BTY560>
- Cheng, H., Concepcion, G. T., Feng, X., Zhang, H., & Li, H. (2021). Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nature Methods* 2021 18:2, 18(2), 170–175. <https://doi.org/10.1038/s41592-020-01056-5>
- Chin, C. S., & Khalak, A. (2019). Human genome assembly in 100 minutes. *BioRxiv*. <https://doi.org/10.1101/705616>
- Chu, C., Borges-Monroy, R., Viswanadham, V. V., Lee, S., Li, H., Lee, E. A., & Park, P. J. (2021). Comprehensive identification of transposable element insertions using multiple sequencing technologies. *Nature Communications* 2021 12:1, 12(1), 1–12. <https://doi.org/10.1038/s41467-021-24041-8>
- Cook, D. E., & Andersen, E. C. (2017). VCF-kit: Assorted utilities for the variant call format. *Bioinformatics*, 33(10), 1581–1582. <https://doi.org/10.1093/bioinformatics/btx011>
- Crysnanto, D., Leonard, A. S., Fang, Z. H., & Pausch, H. (2021). Novel functional sequences uncovered through a bovine multiassembly graph. *Proceedings of the National Academy of Sciences of the United States of America*, 118(20), 2101056118. <https://doi.org/10.1073/pnas.2101056118>
- Crysnanto, D., & Pausch, H. (2020). Bovine breed-specific augmented reference graphs facilitate accurate sequence read mapping and unbiased variant discovery. *Genome Biology*, 21(1), 184. <https://doi.org/10.1186/s13059-020-02105-0>
- Crysnanto, D., Wurmser, C., & Pausch, H. (2019). Accurate sequence variant genotyping in cattle using variation-aware genome graphs. *Genetics Selection Evolution*, 51(1), 21. <https://doi.org/10.1186/s12711-019-0462-x>
- Daetwyler, H. D., Capitan, A., Pausch, H., Stothard, P., Van Binsbergen, R., Brøndum, R. F., Liao, X., Djari, A., Rodriguez, S. C., Grohs, C., Esquerré, D., Bouchez, O., Rossignol, M. N., Klopp, C., Rocha, D., Fritz, S., Eggen, A., Bowman, P. J., Coote, D., ... Hayes, B. J. (2014). Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and complex traits in cattle. *Nature Genetics*, 46(8), 858–865. <https://doi.org/10.1038/ng.3034>
- Danecek, P., Bonfield, J. K., Liddle, J., Marshall, J., Ohan, V., Pollard, M. O., Whitwham, A., Keane, T., McCarthy, S. A., Davies, R. M., & Li, H. (2021). Twelve years of SAMtools and BCFtools. *GigaScience*, 10(2), 1–4. <https://doi.org/10.1093/GIGASCIENCE/GIAB008>
- Della Coletta, R., Qiu, Y., Ou, S., Hufford, M. B., & Hirsch, C. N. (2021). How the pan-genome is changing crop genomics and improvement. *Genome Biology* 2021 22:1, 22(1), 1–19.

- <https://doi.org/10.1186/S13059-020-02224-8>
- Dong, D., Jones, G., & Zhang, S. (2009). Dynamic evolution of bitter taste receptor genes in vertebrates. *BMC Evolutionary Biology* 2009 9:1, 9(1), 1–9.
<https://doi.org/10.1186/1471-2148-9-12>
- Ebler, J., Clarke, W. E., Rausch, T., Audano, P. A., Houwaart, T., Korbel, J., Eichler, E. E., Zody, M. C., Dilthey, A. T., & Marschall, T. (2020). Pangenome-based genome inference. *BioRxiv*, 2020.11.11.378133. <https://doi.org/10.1101/2020.11.11.378133>
- Eid, J., Fehr, A., Gray, J., Luong, K., Lyle, J., Otto, G., Peluso, P., Rank, D., Baybayan, P., Bettman, B., Bibillo, A., Bjornson, K., Chaudhuri, B., Christians, F., Cicero, R., Clark, S., Dalal, R., DeWinter, A., Dixon, J., ... Turner, S. (2009). Real-time DNA sequencing from single polymerase molecules. *Science*, 323(5910), 133–138.
<https://doi.org/10.1126/science.1162986>
- Elsik, C. G., Tellam, R. L., Worley, K. C., Gibbs, R. A., Muzny, D. M., Weinstock, G. M., Adelson, D. L., Eichler, E. E., Elnitski, L., Guigó, R., Hamernik, D. L., Kappes, S. M., Lewin, H. A., Lynn, D. J., Nicholas, F. W., Raymond, A., Rijkels, M., Skow, L. C., Zdobnov, E. M., ... Zhao, F. Q. (2009). The genome sequence of taurine cattle: A window to ruminant biology and evolution. *Science*, 324(5926), 522–528.
<https://doi.org/10.1126/science.1169588>
- Fang, Z. H., & Pausch, H. (2019). Multi-trait meta-analyses reveal 25 quantitative trait loci for economically important traits in Brown Swiss cattle. *BMC Genomics*, 20(1).
<https://doi.org/10.1186/s12864-019-6066-6>
- Feng, X., & Li, H. (2021). Higher Rates of Processed Pseudogene Acquisition in Humans and Three Great Apes Revealed by Long-Read Assemblies. *Molecular Biology and Evolution*, 38(7), 2958–2966. <https://doi.org/10.1093/molbev/msab062>
- Formenti, G., Rhie, A., Walenz, B. P., Thibaud-Nissen, F., Shafin, K., Koren, S., Myers, E. W., Jarvis, E. D., & Phillippy, A. M. (2021). Merfin: improved variant filtering and polishing via k-mer validation. *BioRxiv*, 2021.07.16.452324.
<https://doi.org/10.1101/2021.07.16.452324>
- Garrison, E., Sirén, J., Novak, A. M., Hickey, G., Eizenga, J. M., Dawson, E. T., Jones, W., Garg, S., Markello, C., Lin, M. F., Paten, B., & Durbin, R. (2018). Variation graph toolkit improves read mapping by representing genetic variation in the reference. *Nature Biotechnology*, 36(9), 875–881. <https://doi.org/10.1038/nbt.4227>
- Gibson, B. G., & Briggs, M. D. (2016). The aggrecanopathies; an evolving phenotypic spectrum of human genetic skeletal diseases. *Orphanet Journal of Rare Diseases* 2016 11:1, 11(1), 1–8. <https://doi.org/10.1186/S13023-016-0459-2>
- Girardot, M., Guibert, S., Laforet, M.-P., Gallard, Y., Larroque, H., & Oulmouden, A. (2006). The insertion of a full-length *Bos taurus* LINE element is responsible for a transcriptional deregulation of the Normande Agouti gene. *Pigment Cell Research*, 19(4), 346–355. <https://doi.org/10.1111/J.1600-0749.2006.00312.X>
- Guan, D., Guan, D., McCarthy, S. A., Wood, J., Howe, K., Wang, Y., Durbin, R., & Durbin, R. (2020). Identifying and removing haplotypic duplication in primary genome assemblies. *Bioinformatics*, 36(9), 2896–2898. <https://doi.org/10.1093/bioinformatics/btaa025>
- Gupta, P. K. (2021). Quantitative genetics: pan-genomes, SVs, and k-mers for GWAS. In *Trends in Genetics* (Vol. 37, Issue 10, pp. 868–871). Elsevier.
<https://doi.org/10.1016/j.tig.2021.05.006>
- Hayes, B. J., & Daetwyler, H. D. (2019). 1000 Bull Genomes Project to Map Simple and Complex Genetic Traits in Cattle: Applications and Outcomes.

- <https://doi.org/10.1146/Annurev-Animal-020518-115024>, 7, 89–102.
<https://doi.org/10.1146/ANNUREV-ANIMAL-020518-115024>
- Heaton, M. P., Smith, T. P. L., Bickhart, D. M., Vander Ley, B. L., Kuehn, L. A., Oppenheimer, J., Shafer, W. R., Schuetze, F. T., Stroud, B., McClure, J. C., Barfield, J. P., Blackburn, H. D., Kalbfleisch, T. S., Davenport, K. M., Kuhn, K. L., Green, R. E., Shapiro, B., & Rosen, B. D. (2021). A Reference Genome Assembly of Simmental Cattle, *Bos taurus taurus*. *Journal of Heredity*, 112(2), 184–191. <https://doi.org/10.1093/JHERED/ESAB002>
- Hess, K., Oliverio, R., Nguyen, P., Le, D., Ellis, J., Kdeiss, B., Ord, S., Chalkia, D., & Nikolaidis, N. (2018). Concurrent action of purifying selection and gene conversion results in extreme conservation of the major stress-inducible Hsp70 genes in mammals. *Scientific Reports*, 8(1). <https://doi.org/10.1038/s41598-018-23508-x>
- Hickey, G., Heller, D., Monlong, J., Sibbesen, J. A., Sirén, J., Eizenga, J., Dawson, E. T., Garrison, E., Novak, A. M., & Paten, B. (2020). Genotyping structural variants in pangenome graphs using the vg toolkit. *Genome Biology* 2020 21:1, 21(1), 1–17. <https://doi.org/10.1186/S13059-020-1941-7>
- Holley, G., Beyter, D., Ingimundardottir, H., Møller, P. L., Kristmundsdottir, S., Eggertsson, H. P., & Halldorsson, B. V. (2021). Ratatosk: hybrid error correction of long reads enables accurate variant calling and assembly. *Genome Biology* 2021 22:1, 22(1), 1–22. <https://doi.org/10.1186/S13059-020-02244-4>
- Hufford, M. B., Seetharam, A. S., Woodhouse, M. R., Chougule, K. M., Ou, S., Liu, J., Ricci, W. A., Guo, T., Olson, A., Qiu, Y., Coletta, R. Della, Tittes, S., Hudson, A. I., Marand, A. P., Wei, S., Lu, Z., Wang, B., Tello-Ruiz, M. K., Piri, R. D., ... Dawe, R. K. (2021). De novo assembly, annotation, and comparative analysis of 26 diverse maize genomes. *Science*, 373(6555), 655–662. <https://doi.org/10.1126/science.abg5289>
- Kim, J., Hanotte, O., Mwai, O. A., Dessie, T., Bashir, S., Diallo, B., Agaba, M., Kim, K., Kwak, W., Sung, S., Seo, M., Jeong, H., Kwon, T., Taye, M., Song, K.-D., Lim, D., Cho, S., Lee, H.-J., Yoon, D., ... Kim, H. (2017). The genome landscape of indigenous African cattle. *Genome Biology* 2017 18:1, 18(1), 1–14. <https://doi.org/10.1186/S13059-017-1153-Y>
- Kolmogorov, M., Yuan, J., Lin, Y., & Pevzner, P. A. (2019). Assembly of long, error-prone reads using repeat graphs. *Nature Biotechnology* 2019 37:5, 37(5), 540–546. <https://doi.org/10.1038/s41587-019-0072-8>
- Koren, S., Rhie, A., Walenz, B. P., Dilthey, A. T., Bickhart, D. M., Kingan, S. B., Hiendleder, S., Williams, J. L., Smith, T. P. L., & Phillippy, A. M. (2018). De novo assembly of haplotype-resolved genomes with trio binning. *Nature Biotechnology*, 36(12), 1174–1182. <https://doi.org/10.1038/nbt.4277>
- Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., & Phillippy, A. M. (2017). Canu: Scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Research*, 27(5), 722–736. <https://doi.org/10.1101/gr.215087.116>
- Köster, J., Mölder, F., Jablonski, K. P., Letcher, B., Hall, M. B., Tomkins-Tinch, C. H., Sochat, V., Forster, J., Lee, S., Twardziok, S. O., Kanitz, A., Wilm, A., Holtgrewe, M., Rahmann, S., & Nahnsen, S. (2021). Sustainable data analysis with Snakemake. *F1000Research*, 10, 33. <https://doi.org/10.12688/f1000research.29032.2>
- Lee, Y.-L., Takeda, H., Moreira, G. C. M., Karim, L., Mullaart, E., Coppieters, W., consortium, T. G., Appeltant, R., Veerkamp, R. F., Groenen, M. A. M., Georges, M., Bosse, M., Druet, T., Bouwman, A. C., & Charlier, C. (2021). A 12 kb multi-allelic copy number variation encompassing a GC gene enhancer is associated with mastitis resistance in dairy cattle.

- PLOS Genetics*, 17(7), e1009331. <https://doi.org/10.1371/JOURNAL.PGEN.1009331>
- Li, H. (2013). *Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM*. <https://arxiv.org/abs/1303.3997v2>
- Li, H. (2018). Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18), 3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>
- Li, H., Feng, X., & Chu, C. (2020). The design and construction of reference pangenome graphs with minigraph. *Genome Biology*, 21(1). <https://doi.org/10.1186/s13059-020-02168-z>
- Lin, M. F., Rodeh, O., Penn, J., Bai, X., Reid, J. G., Krasheninina, O., & Salerno, W. J. (2018). GLnexus: joint variant calling for large cohort sequencing. *BioRxiv*, 343970. <https://doi.org/10.1101/343970>
- Lloret-Villas, A., Bhati, M., Kadri, N. K., Fries, R., & Pausch, H. (2021). Investigating the impact of reference assembly choice on genomic analyses in a cattle breed. *BMC Genomics* 2021 22:1, 22(1), 1–17. <https://doi.org/10.1186/S12864-021-07554-W>
- Loftus, R. T., MacHugh, D. E., Bradley, D. G., Sharp, P. M., & Cunningham, P. (1994). Evidence for two independent domestications of cattle. *Proceedings of the National Academy of Sciences*, 91(7), 2757–2761. <https://doi.org/10.1073/PNAS.91.7.2757>
- Logsdon, G. (2020). *HMW gDNA purification and ONT ultra-long-read data generation*. <https://www.protocols.io/view/hmw-gdna-purification-and-ont-ultra-long-read-data-bchhit36>
- Low, W. Y., Tearle, R., Liu, R., Koren, S., Rhie, A., Bickhart, D. M., Rosen, B. D., Kronenberg, Z. N., Kingan, S. B., Tseng, E., Thibaud-Nissen, F., Martin, F. J., Billis, K., Ghurye, J., Hastie, A. R., Lee, J., Pang, A. W. C., Heaton, M. P., Phillippy, A. M., ... Williams, J. L. (2020). Haplotype-resolved genomes provide insights into structural variation and gene content in Angus and Brahman cattle. *Nature Communications*, 11(1), 1–14. <https://doi.org/10.1038/s41467-020-15848-y>
- Md, V., Misra, S., Li, H., & Aluru, S. (2019). Efficient Architecture-Aware Acceleration of BWA-MEM for Multicore Systems. *Proceedings - 2019 IEEE 33rd International Parallel and Distributed Processing Symposium, IPDPS 2019*, 314–324. <https://arxiv.org/abs/1907.12931v1>
- Mikheyev, A. S., & Tin, M. M. Y. (2014). A first look at the Oxford Nanopore MinION sequencer. *Molecular Ecology Resources*, 14(6), 1097–1102. <https://doi.org/10.1111/1755-0998.12324>
- Nurk, S., Koren, S., Rhie, A., Rautiainen, M., Bzikadze, A. V., Mikheenko, A., Vollger, M. R., Altemose, N., Uralsky, L., Gershman, A., Aganezov, S., Hoyt, S. J., Diekhans, M., Logsdon, G. A., Alonge, M., Antonarakis, S. E., Borchers, M., Bouffard, G. G., Brooks, S. Y., ... Phillippy, A. M. (2021). The complete sequence of a human genome. *BioRxiv*, 2021.05.26.445798. <https://doi.org/10.1101/2021.05.26.445798>
- Nurk, S., Walenz, B. P., Rhie, A., Vollger, M. R., Logsdon, G. A., Grothe, R., Miga, K. H., Eichler, E. E., Phillippy, A. M., & Koren, S. (2020). HiCanu: Accurate assembly of segmental duplications, satellites, and allelic variants from high-fidelity long reads. *Genome Research*, 30(9), 1291–1305. <https://doi.org/10.1101/GR.263566.120>
- Olsen, H. G., Knutsen, T. M., Lewandowska-Sabat, A. M., Grove, H., Nome, T., Svendsen, M., Arnyasi, M., Sodeland, M., Sundsaasen, K. K., Dahl, S. R., Heringstad, B., Hansen, H. H., Olsaker, I., Kent, M. P., & Lien, S. (2016). Fine mapping of a QTL on bovine chromosome 6 using imputed full sequence data suggests a key role for the group-specific component (GC) gene in clinical mastitis and milk production. *Genetics Selection*

- Evolution*, 48(1). <https://doi.org/10.1186/s12711-016-0257-2>
- Oppenheimer, J., Rosen, B. D., Heaton, M. P., Vander Ley, B. L., Shafer, W. R., Schuetze, F. T., Stroud, B., Kuehn, L. A., McClure, J. C., Barfield, J. P., Blackburn, H. D., Kalbfleisch, T. S., Bickhart, D. M., Davenport, K. M., Kuhn, K. L., Green, R. E., Shapiro, B., & Smith, T. P. L. (2021). A Reference Genome Assembly of American Bison, *Bison bison bison*. *Journal of Heredity*, 112(2), 174–183. <https://doi.org/10.1093/JHERED/ESAB003>
- Paten, B., Eizenga, J. M., Rosen, Y. M., Novak, A. M., Garrison, E., & Hickey, G. (2018). Superbubbles, Ultrabubbles, and Cacti. *Journal of Computational Biology*, 25(7), 649. <https://doi.org/10.1089/CMB.2017.0251>
- Pausch, H., Emmerling, R., Schwarzenbacher, H., & Fries, R. (2016). A multi-trait meta-analysis with imputed sequence variants reveals twelve QTL for mammary gland morphology in Fleckvieh cattle. *Genetics Selection Evolution*, 48(1). <https://doi.org/10.1186/s12711-016-0190-4>
- Pitt, D., Sebane, N., Nicolazzi, E. L., MacHugh, D. E., Park, S. D. E., Colli, L., Martinez, R., Bruford, M. W., & Orozco-terWengel, P. (2019). Domestication of cattle: Two or three events? *Evolutionary Applications*, 12(1), 123. <https://doi.org/10.1111/EVA.12674>
- Poplin, R., Chang, P.-C., Alexander, D., Schwartz, S., Colthurst, T., Ku, A., Newburger, D., Dijamco, J., Nguyen, N., Afshar, P. T., Gross, S. S., Dorfman, L., McLean, C. Y., & DePristo, M. A. (2018). A universal SNP and small-indel variant caller using deep neural networks. *Nature Biotechnology* 2018 36:10, 36(10), 983–987. <https://doi.org/10.1038/nbt.4235>
- Quinlan, A. R., & Hall, I. M. (2010). BEDTools: A flexible suite of utilities for comparing genomic features. *Bioinformatics*, 26(6), 841–842. <https://doi.org/10.1093/bioinformatics/btq033>
- Ramírez, F., Ryan, D. P., Grüning, B., Bhardwaj, V., Kilpert, F., Richter, A. S., Heyne, S., Dündar, F., & Manke, T. (2016). deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Research*, 44(W1), W160–W165. <https://doi.org/10.1093/nar/gkw257>
- Rengs, W. M. J. van, Schmidt, M. H.-W., Effgen, S., Wang, Y., Zaidan, M. W. A. M., Huettel, B., Schouten, H. J., Usadel, B., & Underwood, C. J. (2021). A gap-free tomato genome built from complementary PacBio and Nanopore long DNA sequences reveals extensive linkage drag during breeding. *BioRxiv*, 2021.08.30.456472. <https://doi.org/10.1101/2021.08.30.456472>
- Rhie, A., McCarthy, S. A., Fedrigo, O., Damas, J., Formenti, G., Koren, S., Uliano-Silva, M., Chow, W., Fungtammasan, A., Kim, J., Lee, C., Ko, B. J., Chaisson, M., Gedman, G. L., Cantin, L. J., Thibaud-Nissen, F., Haggerty, L., Bista, I., Smith, M., ... Jarvis, E. D. (2021). Towards complete and error-free genome assemblies of all vertebrate species. *Nature* 2021 592:7856, 592(7856), 737–746. <https://doi.org/10.1038/s41586-021-03451-0>
- Rhie, A., Walenz, B. P., Koren, S., & Phillippy, A. M. (2020). Merqury: Reference-free quality, completeness, and phasing assessment for genome assemblies. *Genome Biology*, 21(1). <https://doi.org/10.1186/s13059-020-02134-9>
- Rice, E. S., Koren, S., Rhie, A., Heaton, M. P., Kalbfleisch, T. S., Hardy, T., Hackett, P. H., Bickhart, D. M., Rosen, B. D., Ley, B. Vander, Maurer, N. W., Green, R. E., Phillippy, A. M., Petersen, J. L., & Smith, T. P. L. (2020). Continuous chromosome-scale haplotypes assembled from a single interspecies F1 hybrid of yak and cattle. *GigaScience*, 9(4), 1–9. <https://doi.org/10.1093/GIGASCIENCE/GIAA029>
- Rosen, B. D., Bickhart, D. M., Schnabel, R. D., Koren, S., Elsik, C. G., Tseng, E., Rowan, T. N.,

- Low, W. Y., Zimin, A., Couldrey, C., Hall, R., Li, W., Rhie, A., Ghurye, J., McKay, S. D., Thibaud-Nissen, F., Hoffman, J., Murdoch, B. M., Snelling, W. M., ... Medrano, J. F. (2020). De novo assembly of the cattle reference genome with single-molecule sequencing. *GigaScience*, 9(3), 1–9. <https://doi.org/10.1093/gigascience/giaa021>
- Sandor, C., Li, W., Coppieters, W., Druet, T., Charlier, C., & Georges, M. (2012). Genetic Variants in REC8, RNF212, and PRDM9 Influence Male Recombination in Cattle. *PLOS Genetics*, 8(7), e1002854. <https://doi.org/10.1371/JOURNAL.PGEN.1002854>
- Shafin, K., Pesout, T., Chang, P.-C., Nattestad, M., Kolesnikov, A., Goel, S., Baid, G., Eizenga, J. M., Miga, K. H., Carnevali, P., Jain, M., Carroll, A., & Paten, B. (2021). Haplotype-aware variant calling enables high accuracy in nanopore long-reads using deep neural networks. *BioRxiv*, 2021.03.04.433952. <https://doi.org/10.1101/2021.03.04.433952>
- Shafin, K., Pesout, T., Lorig-Roach, R., Haukness, M., Olsen, H. E., Bosworth, C., Armstrong, J., Tigyi, K., Maurer, N., Koren, S., Sedlazeck, F. J., Marschall, T., Mayes, S., Costa, V., Zook, J. M., Liu, K. J., Kilburn, D., Sorensen, M., Munson, K. M., ... Paten, B. (2020). Nanopore sequencing and the Shasta toolkit enable efficient de novo assembly of eleven human genomes. *Nature Biotechnology*, 38(9), 1044–1053. <https://doi.org/10.1038/s41587-020-0503-6>
- Shen, Y., Zhang, F., Li, F., Jiang, X., Yang, Y., Li, X., Li, W., Wang, X., Cheng, J., Liu, M., Zhang, X., Yuan, G., Pei, X., Cai, K., Hu, F., Sun, J., Yan, L., Tang, L., Jiang, C., ... Xu, W. (2019). Loss-of-function mutations in QRICH2 cause male infertility with multiple morphological abnormalities of the sperm flagella. *Nature Communications* 2019 10:1, 10(1), 1–15. <https://doi.org/10.1038/s41467-018-08182-x>
- Sherman, R. M., Forman, J., Antonescu, V., Puiu, D., Daya, M., Rafaels, N., Boorgula, M. P., Chavan, S., Vergara, C., Ortega, V. E., Levin, A. M., Eng, C., Yazdanbakhsh, M., Wilson, J. G., Marrugo, J., Lange, L. A., Williams, L. K., Watson, H., Ware, L. B., ... Salzberg, S. L. (2018). Assembly of a pan-genome from deep sequencing of 910 humans of African descent. *Nature Genetics* 2018 51:1, 51(1), 30–35. <https://doi.org/10.1038/s41588-018-0273-y>
- Silvestre-Ryan, J., & Holmes, I. (2021). Pair consensus decoding improves accuracy of neural network basecallers for nanopore sequencing. *Genome Biology* 2021 22:1, 22(1), 1–6. <https://doi.org/10.1186/S13059-020-02255-1>
- Song, J.-M., Guan, Z., Hu, J., Guo, C., Yang, Z., Wang, S., Liu, D., Wang, B., Lu, S., Zhou, R., Xie, W.-Z., Cheng, Y., Zhang, Y., Liu, K., Yang, Q.-Y., Chen, L.-L., & Guo, L. (2020). Eight high-quality genomes reveal pan-genome architecture and ecotype differentiation of *Brassica napus*. *Nature Plants*, 6(1), 34. <https://doi.org/10.1038/S41477-019-0577-7>
- Suqueli García, M. F., Castellote, M. A., Feingold, S. E., & Corva, P. M. (2017). Characterization of a deletion in the Hsp70 cluster in the bovine reference genome. *Animal Genetics*, 48(4), 377–385. <https://doi.org/10.1111/age.12561>
- Talenti, A., Powell, J., Hemmink, J. D., Cook, E. A. J., Wragg, D., Jayaraman, S., Paxton, E., Ezeasor, C., Obishakin, E. T., Agusi, E. R., Tijjani, A., Marshall, K., Fisch, A., Ferreira, B., Qasim, A., Chaudhry, U. N., Wiener, P., P., T., Morrison, L. J., ... Prendergast, J. (2021). A cattle graph genome incorporating global breed diversity. *BioRxiv*, 2021.06.23.449389. <https://doi.org/10.1101/2021.06.23.449389>
- Tettelin, H., Maignani, V., Cieslewicz, M. J., Donati, C., Medini, D., Ward, N. L., Angiuoli, S. V., Crabtree, J., Jones, A. L., Durkin, A. S., DeBoy, R. T., Davidsen, T. M., Mora, M., Scarselli, M., Margarit Y Ros, I., Peterson, J. D., Hauser, C. R., Sundaram, J. P., Nelson, W. C., ... Fraser, C. M. (2005). Genome analysis of multiple pathogenic isolates of

- Streptococcus agalactiae: Implications for the microbial “pan-genome.” *Proceedings of the National Academy of Sciences of the United States of America*, 102(39), 13950–13955. <https://doi.org/10.1073/pnas.0506758102>
- Trigo, B. B., Utsunomiya, A. T. H., Fortunato, A. A. A. D., Milanesi, M., Torrecilha, R. B. P., Lamb, H., Nguyen, L., Ross, E. M., Hayes, B., Padula, R. C. M., Sussai, T. S., Zavarez, L. B., Cipriano, R. S., Caminhas, M. M. T., Lopes, F. L., Pelle, C., Leeb, T., Bannasch, D., Bickhart, D., ... Utsunomiya, Y. T. (2021). Variants at the ASIP locus contribute to coat color darkening in Nellore cattle. *Genetics Selection Evolution* 2021 53:1, 53(1), 1–12. <https://doi.org/10.1186/S12711-021-00633-2>
- Vaser, R., & Šikić, M. (2021). Time- and memory-efficient genome assembly with Raven. *Nature Computational Science*, 1(5), 332–336. <https://doi.org/10.1038/S43588-021-00073-4>
- Vaser, R., Sovic, I., Nagarajan, N., & Sikic, M. (2017). Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Research*, 27(5), gr.214270.116. <https://doi.org/10.1101/GR.214270.116>
- Wagih, O. (2017). Ggseqlogo: A versatile R package for drawing sequence logos. *Bioinformatics*, 33(22), 3645–3647. <https://doi.org/10.1093/bioinformatics/btx469>
- Wenger, A. M., Peluso, P., Rowell, W. J., Chang, P.-C., Hall, R. J., Concepcion, G. T., Ebler, J., Fungtammasan, A., Kolesnikov, A., Olson, N. D., Töpfer, A., Alonge, M., Mahmoud, M., Qian, Y., Chin, C.-S., Phillippy, A. M., Schatz, M. C., Myers, G., DePristo, M. A., ... Hunkapiller, M. W. (2019). Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nature Biotechnology* 2019 37:10, 37(10), 1155–1162. <https://doi.org/10.1038/s41587-019-0217-9>
- Wick, R. R., & Holt, K. E. (2021). Benchmarking of long-read assemblers for prokaryote whole genome sequencing. *F1000Research* 2021 8:2138, 8, 2138. <https://doi.org/10.12688/f1000research.21782.4>
- Wick, R. R., Schultz, M. B., Zobel, J., & Holt, K. E. (2015). Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics*, 31(20), 3350–3352. <https://doi.org/10.1093/BIOINFORMATICS/BTV383>
- Wilks, C., Ahmed, O., Baker, D. N., Zhang, D., Collado-Torres, L., & Langmead, B. (2021). Megadepth: efficient coverage quantification for BigWigs and BAMs. *Bioinformatics*. <https://doi.org/10.1093/BIOINFORMATICS/BTAB152>
- Wu, D. D., Ding, X. D., Wang, S., Wójcik, J. M., Zhang, Y., Tokarska, M., Li, Y., Wang, M. S., Faruque, O., Nielsen, R., Zhang, Q., & Zhang, Y. P. (2018). Pervasive introgression facilitated domestication and adaptation in the Bos species complex. *Nature Ecology and Evolution*, 2(7), 1139–1145. <https://doi.org/10.1038/s41559-018-0562-y>
- Yan, S. M., Sherman, R. M., Taylor, D. J., Nair, D. R., Bortvin, A. N., Schatz, M. C., & McCoy, R. C. (2021). Local adaptation and archaic introgression shape global diversity at human structural variant loci. *ELife*, 10. <https://doi.org/10.7554/ELIFE.67615>
- Yang, J., & Chaisson, M. (2021). *TT-Mars: Structural Variants Assessment Based on Haplotype-resolved Assemblies*. <https://doi.org/10.1101/2021.09.27.462044>
- Zhou, Y., Shen, B., Jiang, J., Padhi, A., Park, K.-E., Oswald, A., Sattler, C. G., Telugu, B. P., Chen, H., Cole, J. B., Liu, G. E., & Ma, L. (2018). Construction of PRDM9 allele-specific recombination maps in cattle using large-scale pedigree analysis and genome-wide single sperm genomics. *DNA Research: An International Journal for Rapid Publication of Reports on Genes and Genomes*, 25(2), 183. <https://doi.org/10.1093/DNARES/DSX048>

