


Draft genome of the lowland anoa (*Bubalus depressicornis*) and comparison with buffalo genome assemblies (Bovidae, Bubalina)

Stefano Porrelli ¹, Michèle Gerbault-Seureau,² Roberto Rozzi ^{3,4}, Rayan Chikhi ⁵, Manon Curaudeau ², Anne Ropiquet ¹, Alexandre Hassanin ^{2,*}

¹Department of Natural Sciences, Faculty of Science and Technology, Middlesex University, London NW4 4BT, UK,

²Institut Systématique Evolution Biodiversité (ISYEB), Sorbonne Université, MNHN, CNRS, EPHE, UA, 75005 Paris, France,

³Museum für Naturkunde, Leibniz-Institut für Evolutions- und Biodiversitätsforschung, 10115 Berlin, Germany,

⁴German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, 04103 Leipzig, Germany,

⁵Institut Pasteur, Université Paris Cité, Sequence Bioinformatics, 75015 Paris, France

*Corresponding author: Institut Systématique Evolution Biodiversité (ISYEB), Sorbonne Université, MNHN, CNRS, EPHE, UA, 57 rue Cuvier, CP 51, 75005 Paris, France. Email: alexandre.hassanin@mnhn.fr

Abstract

Genomic data for wild species of the genus *Bubalus* (Asian buffaloes) are still lacking while several whole genomes are currently available for domestic water buffaloes. To address this, we sequenced the genome of a wild endangered dwarf buffalo, the lowland anoa (*Bubalus depressicornis*), produced a draft genome assembly and made comparison to published buffalo genomes. The lowland anoa genome assembly was 2.56 Gbp long and contained 103,135 contigs, the longest contig being 337.39 kbp long. N50 and L50 values were 38.73 and 19.83 kbp, respectively, mean coverage was 44× and GC content was 41.74%. Two strategies were adopted to evaluate genome completeness: (1) determination of genomic features with de novo and homology-based predictions using annotations of chromosome-level genome assembly of the river buffalo and (2) employment of benchmarking against universal single-copy orthologs (BUSCO). Homology-based predictions identified 94.51% complete and 3.65% partial genomic features. De novo gene predictions identified 32,393 genes, representing 97.14% of the reference's annotated genes, whilst BUSCO search against the mammalian orthologs database identified 71.1% complete, 11.7% fragmented, and 17.2% missing orthologs, indicating a good level of completeness for downstream analyses. Repeat analyses indicated that the lowland anoa genome contains 42.12% of repetitive regions. The genome assembly of the lowland anoa is expected to contribute to comparative genome analyses among bovid species.

Keywords: Bovidae; *Bubalus depressicornis*; lowland anoa; genome assembly; de novo assembly

Introduction

The lowland anoa, *Bubalus depressicornis* (Smith 1827), is a wild dwarf buffalo endemic to Sulawesi and Buton Islands, where it can be found in sympatry with the mountain anoa, *Bubalus quarlesi* (Ouwens 1910). Both anoa species are currently classified as endangered with declining populations due to hunting and habitat loss (Burton et al. 2016). Because of their singular appearance, they were initially described in their own genus *Anoa* (Ouwens 1910). However, *Anoa* was not regarded as a valid genus in more recent classifications, in which both anoa species were ascribed to the genus *Bubalus*, together with the wild water buffalo—*Bubalus arnee* (Kerr 1792) and the tamaraw—*Bubalus mindorensis* (Heude 1888; Groves 1969; IUCN 2022). Molecular studies based on mitochondrial sequences have supported a sister-group relationship between *B. depressicornis* and *B. quarlesi* (Schreiber et al. 1999; Priyono et al. 2020). In addition, the mitogenome of the lowland anoa was found to be equally distant from those of the

2 types of domestic water buffalo, the river buffalo from the Indian subcontinent and Mediterranean countries and the swamp buffalo from China and Southeast Asia (Hassanin et al. 2012). Since the same phylogenetic pattern was recovered from the analyses of 2 nuclear datasets, one based on 30 autosomal genes and the other based on 2 genes of the Y chromosome, Curaudeau et al. (2021) have concluded the existence of 2 species of domestic buffaloes: *Bubalus bubalis* (Linnaeus 1758) for the river buffalo and *Bubalus kerabau* (Fitzinger 1860) for the swamp buffalo, which diverged during the Pleistocene at around 0.84 Mya. As discussed in Curaudeau et al. (2021), the 2 domestic species can easily be distinguished based on coat and horn characteristics (Castelló 2016), and they have different karyotypes: *B. bubalis* has $2n = 50$ chromosomes with a fundamental number (FN) equal to 58; whereas *B. kerabau* has $2n = 48$ chromosomes and FN = 56 (Nguyen et al. 2008).

With rapid progress and cost reduction in sequencing technologies, many whole genomes of domestic bovid species have been

sequenced. Whole-genome sequencing has allowed the identification of variants involved in domestication and genetic improvement for several livestock species such as cattle and buffaloes (Zimin *et al.* 2009; Canavez *et al.* 2012; Li *et al.* 2020; Rosen *et al.* 2020). Chromosome-level genome assemblies include those of the domestic cow, *Bos taurus* (Zimin *et al.* 2009), the domestic river buffalo, *B. bubalis* (Deng *et al.* 2016), the swamp buffalo, *B. kerabau* [reported as *Bubalus carabanensis* in Luo *et al.* (2020) but see Curaudeau *et al.* (2021) for further taxonomic information], the domestic Yak, *Bos grunniens* (Zhang *et al.* 2021) and the zebu cattle, *Bos indicus* (Canavez *et al.* 2012). Whereas a total of 8 chromosome- and scaffold-level genome assemblies are publicly available for domestic buffaloes, there are currently no genome data available for wild species of the genus *Bubalus*. To fill this gap, a biopsy of a living lowland anoa was used for next-generation sequencing, and a draft genome was assembled de novo for comparison to other buffalo genome assemblies available in international databases such as NCBI (National Center for Biotechnology Information) and BIG_GWH (Beijing Institute of Genomics Genome Warehouse database).

Materials and methods

DNA extraction, library preparation, and genome sequencing

A living male adult of lowland anoa, named Yannick, was sampled at the *Ménagerie du Jardin des Plantes* of the Muséum national d'Histoire naturelle (MNHN, Paris, France; Fig. 1). A skin biopsy was performed in 2006 by a veterinary surgeon following protocols approved by the MNHN and in line with ethical guidelines. The same biopsy was previously used to determine its karyotype ($2n = 48$; FN = 58; Nguyen *et al.* 2008). DNA was extracted using the DNeasy Blood and Tissue Kit (Qiagen, Hilden, Germany) following the manufacturer's protocol. DNA quantification was performed with a Qubit 2.0 Fluorometer with Qubit dsDNA HS Assay Kit (Thermo Fischer Scientific, Waltham, MA, USA). Library preparation and sequencing were conducted at the *Institut du Cerveau et de la Moelle épinière*. The sample was sequenced on a NextSeq 500 Illumina system generating 2×151 bp reads using the NextSeq 500 High Output Kit v2 with 300 cycles and aiming for an insert size of 350 bp.



Fig. 1. Lowland anoa (*Bubalus depressicornis*) housed at the *Ménagerie du Jardin des Plantes* (© Alexandre Hassanin—MNHN).

De novo assembly

Data quality was assessed with FastQC v.0.11.5 (<https://www.bioinformatics.babrah.am.ac.uk/projects/fastqc/>) and results were collated with MultiQC v1.12 (Ewels *et al.* 2016). Raw reads were quality-trimmed and adapter sequences and contaminants removed with Trimmomatic v.0.36 (Bolger *et al.* 2014) with the following parameters: “ILLUMINACLIP: TruSeq3 -PE.fa:2:30:10 LEADING:33 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36.” Data quality of quality-trimmed reads was reassessed with FastQC. A de novo assembly was performed with MaSuRCA v.3.3.1 (Zimin *et al.* 2013, 2017) using recommended parameters for mammalian genomes and paired-end illumina-only data, as indicated in Zimin *et al.* (2017). The mean and standard deviation for the Insert size were estimated with an “estimate-insert-size” script (<https://gist.github.com/rchikhi/7281991>). Paired-end reads were error corrected using Quorum (Marçais *et al.* 2015) and assembled into super-reads using a k-mer size of 99, as selected by the MaSuRCA assembler. The super-reads were then assembled into contigs using the CABOG assembler, part of the MaSuRCA pipeline (Zimin *et al.* 2017), followed by gap closing with the paired-end information (Zimin *et al.* 2013).

Assembly quality assessment

Genome assemblies publicly available for *Bubalus* and *Syncerus* genera were retrieved from NCBI and BIG_GWH for quality comparison and assessment. The dataset included 2 assemblies at the chromosome level for the river buffalo (*B. bubalis*) with a coverage of 100× and 572×, 4 scaffold-level draft assemblies of river buffalo with coverage ranging between 69× and 119×, one chromosome-level assembly of swamp buffalo (*B. kerabau*) with a mean coverage of 65×, and one scaffold-level draft assembly of the African buffalo (*Syncerus caffer*) with 162× coverage. The 8 retrieved assemblies were sequenced and assembled with different methods, summarized in Table 1.

The quality of the lowland anoa genome assembly was assessed with QAST-LG v.5.0.1 (Mikheenko *et al.* 2018) using the river buffalo NDDB_SH_1 genome assembly (Deng *et al.* 2016) as a reference. The default parameters for mammalian genomes were used to compare all assemblies in QAST-LG: “MODE: large, threads: 50, eukaryotic: true, minimum contig length: 3,000, minimum alignment length: 500, ambiguity: 1, threshold for extensive misassembly size: 7,000.” All analyzed assemblies were aligned to the river buffalo NDDB_SH_1 assembly and results were plotted with Circos v. 0.69.8 (Krzywinski *et al.* 2009) and Jupiter consistency plots (Chu 2018).

We adopted 2 different strategies to evaluate genome completeness. Firstly, genomic features were predicted with the homology-based method by aligning the lowland anoa genome to that of the annotated river buffalo reference genome (NDDB_SH_1 and relative annotations retrieved from NCBI). Secondly, we used a de novo gene prediction method with GlimmerHMM v3.0.4 (Majoros *et al.* 2004). Thirdly, we employed benchmarking against universal single-copy orthologs (BUSCO v5.2.2; Manni *et al.* 2021) using the mammalia_odb10 dataset (2021 February 19, number of genomes: 24, number of BUSCOs: 9,226) from OrthoDB (Kriventseva *et al.* 2019) and compared to other buffalo genome assemblies already deposited on NCBI and BIG_GWH (Table 1).

Repeats and gene annotation

Repetitive regions in the lowland anoa genome were identified, annotated, and masked with RepeatMasker v.4.1.2-p1

Table 1. Information regarding genome assemblies available for buffalo species.

Species/assembly name	Breed	Geographic location	ID	Assembly accession no	Sequencing technology	Assembly method	Coverage	Assembly level
<i>Bubalus bubalis</i> NDDB_SH_1_ (RefSeq)	Murrah	India	NDDB_SH_1	GCF_019923935.1	PacBio Sequel; 10X and BioNano Optical Map	Falcon+Scaff10X +BioNano v. 2019-02-25	572×	Chromosome
<i>Bubalus bubalis</i> Jaffrabadi_v3.0	Jaffrabadi	India	AAUIN_1	GCA_000180995.3	454; Illumina NextSeq 500	MaSuRCA v. 2.3.2b	100×	Scaffold
<i>Bubalus bubalis</i> UOA_WB_1	Mediterranean	Italy	UOA_WB_1	GCA_003121395.1	PacBio	Falcon-Unzip v. 1.8.7	69×	Chromosome
<i>Bubalus bubalis</i> Bubbub1.0	Bangladesh	Bangladesh	Bubbub1.0	GCA_004794615.1	Illumina HiSeq 2000	Soapdenovo v. 2.04	119×	Scaffold
<i>Bubalus bubalis</i> ASM299383v1	Egyptian	Egypt	EGYBUF_1.0	GCA_002993835.1	SOLiD	Velvet v. 1.1; Bowtie2 v. 2.1.0; SHRiMP v. 2.2.3	70×	Scaffold
<i>Bubalus bubalis</i> UMD_CASPUR_WB_2.0	Mediterranean	United States	UMD_CASPUR_WB_2.0	GCA_000471725.1	Illumina GAIIx; Illumina HiSeq; 454	MaSuRCA v. 1.8.3	70×	Scaffold
<i>Bubalus depressicornis</i> * MNHNYannick_LA_1	—	Indonesia	MNHNYannick_LA_1	Assembled MaSuRCA	Illumina NextSeq 500	MaSuRCA v. 3.3.1	44×	Scaffold
<i>Bubalus kerabau</i> CUSA_SWP	Fuzhong	China	CUSA_SWP	GWHAJZ00000000	PacBio 57.8	Wtdbg 1.2.8	65×	Chromosome
<i>Syncerus caffer</i> ASM640878v2	African Buffalo	South Africa	ABF221	GCA_006408785.2	Illumina HiSeq	Platanus v. 1.2.4	162×	Scaffold

* This study.

(Tarailo-Graovac and Chen 2009). Firstly, a de novo repeat library was constructed from the genome assembly with RepeatModeler v.2.0.2a. RepeatMasker was used with default parameters to produce a homolog-based repeat library and mask the genome's repetitive regions. The scripts "calcDivergenceFromAlign.pl" and "createRepeatLandscape.pl" were used to calculate the Kimura divergence values and to plot the resulting repeat landscape. The repeat landscape of *B. taurus* was retrieved from the RepeatMasker database for visual comparison.

Results and discussion

Whole-genome sequencing and data QC

Whole-genome sequencing generated 991,437,058 paired-end reads with a length of 151bp. Quality trimming removed 46,616,722 low-quality, adapter-contaminated, and PCR-duplicated reads, representing approximately 0.5% of the total reads. A total of 944,820,336 clean paired-end reads were generated, covering the lowland anoa genome with an estimated 56× depth based on a genome size of 2.56 Gbp. The estimation of insert size using in-house script returned a mean of 377 and a standard deviation of 83.

De novo assembly quality metrics

The final lowland anoa genome assembly generated here contained 103,135 contigs, the largest being 337.39 kbp long, an N50 of 38.73 kbp and an L50 of 19.83 kbp (Table 2). The total length was 2.56 Gbp with a mean coverage of 44×, and GC content was 41.74%, in agreement with other published assemblies (between 41.60% and 41.92%, Table 3). When aligned to the NDDB_SH_1 genome assembly, the fraction of the anoa genome assembly was 95.41%, a value comparable to other buffalo genome assemblies (Fig. 2), with a total alignment length of 2,515,453,843 bp. A total of 886 contigs could not be aligned to the river buffalo genome assembly, whilst 8,085 contigs were only partially aligned, resulting in a total unaligned length of 45,224,171 bp, which reflects the discrepancy between the total length of the lowland anoa genome and the total aligned length to the reference river buffalo genome assembly. Partially aligned and unaligned contigs could have resulted from structural variations between the lowland anoa and the reference river buffalo assembly, such as large INDELS (insertion/deletions), as well as repetitive regions and/or alternative haplotypes causing assembly errors. The nature of short-read technology causes difficulties in characterizing genomic regions such as telomeres, centromeres, repetitive, and highly heterochromatic regions (Johnson et al. 2005; Low et al. 2019; Weissensteiner and Suh 2019), which are notoriously difficult to assemble and could be better resolved with long-read sequencing.

The lowland anoa genome assembly has a modest N50 compared to other buffalo genome assemblies (Table 3), indicating lower levels of contiguity, which is expected due to the short-read output of Illumina sequencing technology (read

length = 151 bp). In addition, repeat analysis revealed that 42.12% of the lowland anoa genome is composed of repetitive regions. This, coupled with low-sequence coverage, sequencing and assembly errors, causes breaks in the assembly contiguity (Gnerre et al. 2011; Low et al. 2019). This is apparent even in high-quality chromosome-level genome assemblies that use multiple sequencing libraries and multiple sequencing technologies, such as the previous human genome assembly GRCh38, which contained hundreds of gaps (International Human Genome Sequencing Consortium 2004). In addition, the chromosome-level genome assemblies retrieved from NCBI (NDDB_SH_1, UOA_WB_1) were sequenced using multiple insert size libraries and sequencing technologies and were intensively verified with multiple methods such as optical mapping, Hi-C, and RH (Deng et al. 2016; Low et al. 2019).

Moreover, quality metrics of publicly available assemblies are usually limited to reporting N50 and L50 values, which represent the shortest contig length needed to cover 50% of the total assembly size, and the number of contigs whose cumulative length covers 50% of the total assembly size, respectively (Bradnam et al. 2013). Such metrics are often used to compare and evaluate performances of the ever-growing assembly and annotation methods and software (Manchanda et al. 2020). However, we hereby show that reporting N50 and L50 metrics exclusively can be misleading, as they only provide a standard measure of assembly contiguity whilst omitting information such as gene content and completeness, as well as assembly correctness. Furthermore, N50 values can be artificially raised by deliberately excluding short contigs from analyses and by the presence of undetermined nucleotides (Ns) linking the scaffolded contigs (Gurevich et al. 2013). Therefore, to assess the quality of the lowland anoa genome assembly, we generated conventional N50 and L50 metrics and also determined genome completeness in terms of gene content and genome correctness by comparing our assembly to a chromosome-level genome assembly of the river buffalo (*B. bubalis*). In addition, a swamp buffalo (*B. kerabau*, CUSA_SWP) and a more distantly related African buffalo species (*S. caffer*, ABF221) were also included in our comparison.

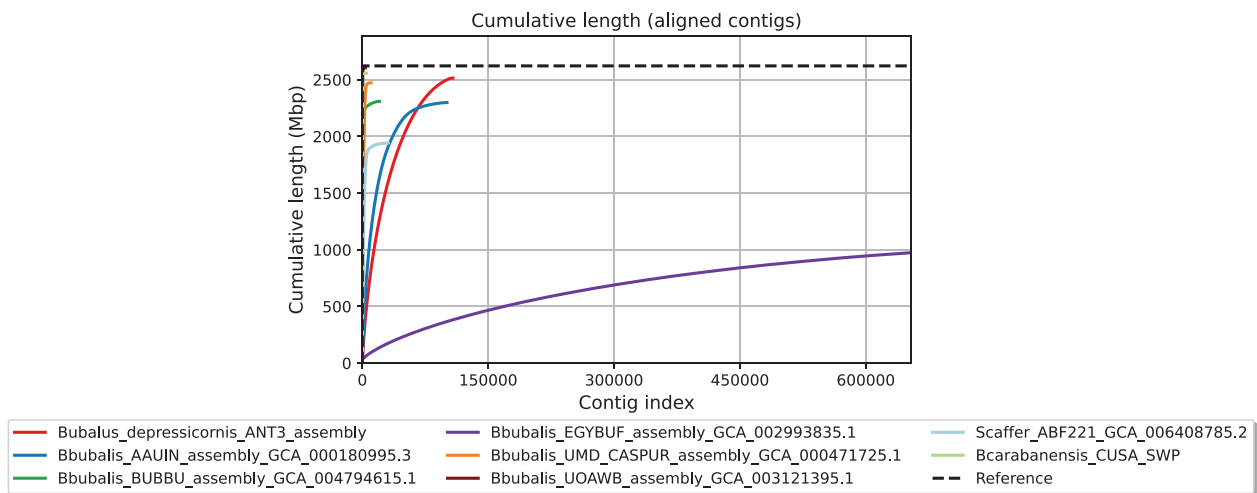
Regardless of the modest N50 value, the lowland anoa genome assembly is in good agreement with the NDDB_SH_1 assembly, with 95.91% of contigs correctly mapped to the 25 reference chromosomes of the river buffalo and fewer misassembled blocks compared to other draft assemblies (Fig. 3). The genome assembly of the Egyptian river buffalo (EGYBUF_1.0) had an abnormally high number of misassembled blocks with respect to the reference genome, followed by the genome assembly of a female Italian river buffalo (UOA_WB_1). To investigate this, misassemblies and structural variation metrics were computed in QUILT (Table 4). The Egyptian river buffalo assembly (EGYBUF_1.0) showed the highest number of mismatches and the highest number of Ns, followed by the Jaffrabadi river buffalo (AAUIN_1). The genome assembly of the African buffalo (*S. caffer*, ABF221) showed a larger number of mismatches (Table 4), but this can be explained by the higher sequence divergence between *Syncerus* and *Bubalus*, as the 2 genera have separated in the Late Miocene (Hassanin et al. 2012). Misassemblies and structural variation metrics could not explain the misassembled blocks of the UOA_WB_1 assembly observed in the Circos plot of Fig. 3. However, some of these misassembled blocks could be due to unplaced contigs. To investigate this, the UOA_WB_1 assembly was aligned to the NDDB_SH_1 reference to generate Jupiter consistency plots. When using the largest 26 contigs of the UOA_WB_1 assembly to cover 100% of the reference river buffalo genome, an

Table 2. Draft assembly statistics of the lowland anoa genome.

Contig statistics	value
Total length	2,565,510,706
Number of contigs	103,135
Largest contig	337,395
GC (%)	41.74
N50	38,737
L50	19,832

Table 3. Comparison of assembly quality metrics of the lowland anoa (*Bubalus depressicornis*) and other buffalo assemblies.

Name/assembly name (NCBI)	ID	Genome fraction %	Total aligned length	Largest alignment	Scaffolds count	N50	L50	GC%
<i>Bubalus bubalis</i> NDDDB_SH1 (RefSeq)	NDDDB_SH_1	—	—	—	26	116,997,125	9	41.75
<i>Bubalus bubalis</i> Jaffrabadi_v3.0	AAUIN_1	83.189	2,299,810,356	834,863	75,621	104,127	9,942	41.78
<i>Bubalus bubalis</i> UOA_WB_1	UOA_WB_1	98.851	2,605,694,501	34,949,624	509	117,219,835	9	41.81
<i>Bubalus bubalis</i> Bubbub1.0	Bubbub1.0	86.537	2,309,804,413	9,328,338	14,905	7,025,746	116	41.6
<i>Bubalus bubalis</i> ASM299383v1	EGYBUF_1.0	36.01	974,053,149	2,013,276	6,313	3,666,815	234	41.92
<i>Bubalus bubalis</i> UMD_CASPUR_WB_2.0	UMD_CASPUR_WB_2.0	93.634	2,473,056,510	7,952,377	5,714	1,545,294	508	41.73
<i>Bubalus depressicornis</i> MNHNYannick_LA_1	MNHNYannick_LA_1	95.415	2,515,453,834	337,395	103,135	38,737	19,832	41.74
<i>Bubalus kerabau</i> CUSA_SWP	CUSA_SWP	97.086	2,557,653,758	23,566,932	1,534	117,253,548	8	41.83
<i>Syncerus caffer</i> ASM640878v2	ABF221	73.046	1,942,672,810	4,692,267	13,167	2,448,414	351	41.72

**Fig. 2.** Cumulative length of aligned contigs of the lowland anoa (red line) against the river buffalo NDDDB_SH_1 reference genome assembly (dashed line) and compared to other buffalo genome assemblies available on NCBI.

almost perfect level of synteny was observed (Fig. 4a). Although this result was expected for genomes of the same species, it also indicates a good level of assembly quality in terms of correctness. However, when including all 509 contigs of the UOA_WB_1 assembly, several misassembled regions were observed (Fig. 4b). Three nonexclusive hypotheses can be advanced to interpret this result: possible genomic rearrangements, genome assembly errors, and repetitive regions. Whether the results of the consistency plots are due to the factors mentioned above or other factors, such as contamination, remains speculative. Nevertheless, the results of the quality metric comparison conducted here further indicate the unreliability of using exclusively N50 and L50 metrics when assessing assembly quality. Instead, contiguity metrics should be supplemented with genome completeness and correctness metrics.

Genomic features, gene prediction, and annotation

Homology and de novo gene predictions performed on the lowland anoa genome assembly were in agreement with each other and indicated a good level of genome completeness. Results were

comparable to other published genome assemblies (Tables 5 and 6), and an improvement over the Bangladeshi river buffalo (Bubbub1.0), the Egyptian river buffalo (EGYBUF_1.0), and Mediterranean river buffalo (UMD_CASPUR_WB_2.0) assemblies.

Interestingly, these 3 assemblies showed higher contiguity (N50) than the draft assembly of the lowland anoa, further indicating the unreliability of using exclusively N50 and L50 metrics when assessing genome assembly quality.

Out of the 1,921,249 genomic features annotations of the reference assembly NDDDB_SH_1, homology prediction identified 1,815,794 (94.51%) complete and 69,929 (3.63%) partial features in the lowland anoa genome assembly, which is comparable to other published assemblies (Fig. 5), indicating a good level of genome completeness. GlimmerHMM de novo predicted 1,027,469 unique genomic features (mRNA and coding sequences, CDS), which is an improvement over some of the water buffalo assemblies used for quality comparison (Table 5). Homology-based gene prediction identified 32,393 genes in the lowland anoa genome assembly, representing 97.14% of the genes annotated in NDDDB_SH_1 ($n = 33,348$). Of these, 59.11% (19,148) were complete and 40.88% (13,245) were partial, probably reflecting the level of fragmentation



Fig. 3. Circos plot of scaffolds mapped to NDDB_SH_1 reference genome assembly (*Bubalus bubalis*). Outer circle represents reference sequence with GC% heatmap (0% = white, 69% = black). Inner circles represent assembly tracks, with heatmap representing correct contigs (green) and misassembled blocks (red).

Table 4. QUAST-LG statistics of all buffalo assemblies with respect to the river buffalo NDDB_SH_1 reference.

	<i>B. depressicornis</i> MNHNYannick_ LA_1	<i>B. bubalis</i> AAUIN_1	<i>B. bubalis</i> Bubbub1.0	<i>B. bubalis</i> EGYBUF_1.0	<i>B. bubalis</i> UMD_CASPUR_ WB_2.0	<i>B. bubalis</i> UOA_WB_1	<i>B. kerabau</i> CUSA_SWP	<i>S. caffer</i> ABF221
Misassemblies	4,949	19,238	3,561	131	4,040	1,724	2,111	6,565
Relocations	1,447	13,540	2,761	85	1,434	1,051	1,199	3,397
Translocations	3,203	4,714	757	10	2,569	647	896	3,032
Inversions	299	984	43	36	37	26	16	136
Misassembled contigs	4,550	15,988	1,049	45	1,943	255	533	1,727
Misassembled contigs length	159,179,266	1,334,096,556	2,506,642,146	55,459,162	1,891,377,139	2,639,940,877	2,594,120,526	2,486,555,687
Local misassemblies	7,014	73,267	241,261	6,933	7,100	4,870	9,940	435,454
Possible TEs	164	874	886	10	544	136	158	654
Unaligned mis. contigs	287	2,378	548	2,522	63	104	381	1,324
Unaligned contigs partial	886 + 8,085	2,555 + 57,865	297 + 7,280	2,806 + 3,472	182 + 3,290	1 + 416	140 + 1110	900 + 7,314
Unaligned length	45,224,171	596,227,806	299,544,303	1,673,093,194	82,826,374	49,291,638	51,316,520	779,611,955
Genome fraction (%)	95.415	83.189	86.537	36.01	93.634	98.851	97.086	73.046
Duplication ratio	1.007	1.425	1.076	1.36	1.034	1.005	1.013	1.045
Mismatches	16,233,421	19,654,061	23,375,163	17,890,296	10,863,130	10,118,782	15,844,866	114,608,168
Indels	1,578,224	746,243	705,955	6,440,610	1,136,878	1,400,310	1,534,735	2,128,964
Indels length	12,654,316	56,163,406	24,209,936	35,356,432	24,745,254	23,411,739	33,123,824	18,236,722
Mismatches per 100 kbp	649	901	1,030	1,895	442	390	622	5,983
Indels per 100 kbp	63	34	31	682	46	54	60	111
Indels (<= 5 bp)	1,297,998	598,354	515,830	5,758,980	893,802	1,227,309	1,269,689	1,641,754
Indels (> 5 bp)	280,226	147,889	190,125	681,630	243,076	173,001	265,046	487,210
N's	493,027	850,098,824	138,209,713	328,128,682	73,946,361	373,500	22,116,406	59,283,755
N's per 100 kbp	19.22	22,942	5,040.03	11,097	2,820.18	14.06	840.50	2,131.26

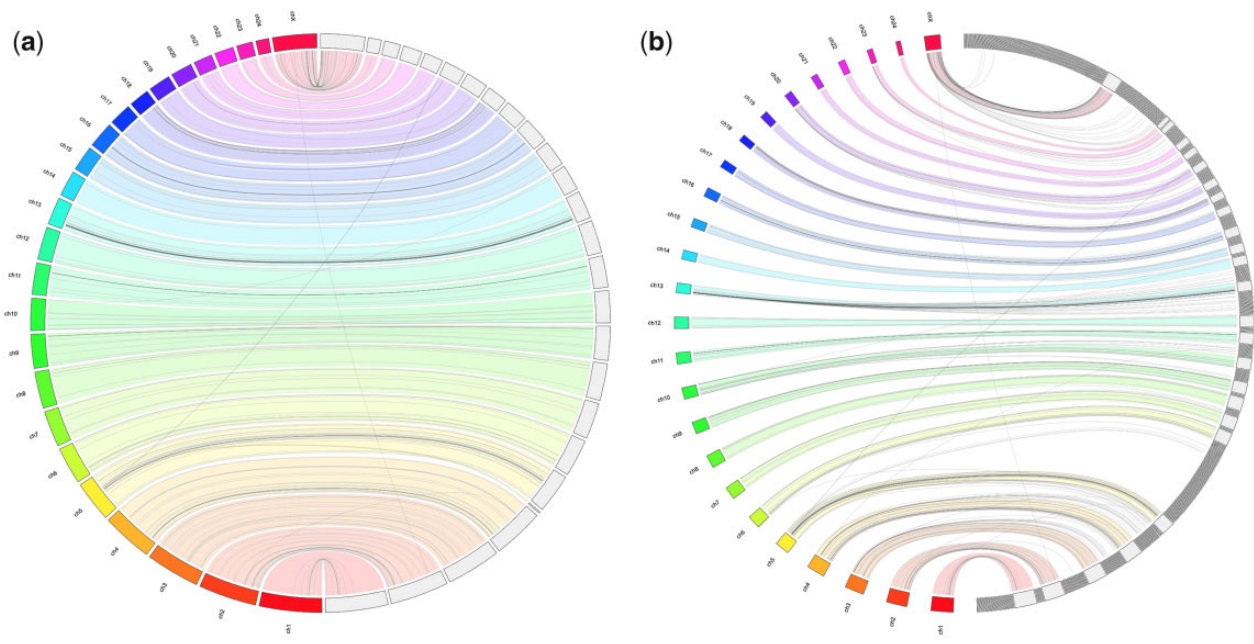


Fig. 4. Jupiter consistency plot showing alignment between the river buffalo genome assemblies UO_AWB_1 and NDDB_SH_1. The left of the plots shows the numbered NDDB_SH_1 chromosomes. The right of the plots shows (a) the 26 longest contigs of the UO_AWB_1 assembly needed to cover 100% of the reference genome and (b) all the 509 contigs of the UO_AWB_1 assembly. Colored bands represent synteny between the genomes. Lines represent genomic rearrangements, break points in the scaffolds or assembly errors. The absence of lines connecting the UO_AWB_1 blocks to the NDDB_SH_1 chromosomes indicates contigs that could not be aligned to the reference.

Table 5. Gene features (CDS and mRNA) predicted with GlimmerHMM.

Name/assembly name (NCBI)	ID	Predicted gene features (unique)	Predicted gene features (≥ 0 bp)	Predicted gene features (≥ 300 bp)	Predicted gene features (≥ 1500 bp)	Predicted gene features ($\geq 3,000$ bp)
<i>Bubalus bubalis</i> Jaffrabadi_v3.0	AAUIN_1	1,065,654	1,087,174 + 1,214 part	719,235 + 911 part	129,801 + 19 part	24,579 + 7 part
<i>Bubalus bubalis</i> UOA_WB_1	UOA_WB_1	1,055,791	1,059,972 + 21 part	762,464 + 17 part	154,594 + 0 part	29,659 + 0 part
<i>Bubalus bubalis</i> Bubbub1.0	Bubbub1.0	948,732	958,663 + 101 part	655,839 + 73 part	136,045 + 4 part	27,867 + 1 part
<i>Bubalus bubalis</i> ASM299383v1	EGYBUF_1.0	826,048	826,155 + 69 part	530,835 + 37 part	96,365 + 0 part	16,243 + 0 part
<i>Bubalus bubalis</i> UMD_CASPUR_WB_2.0	UMD_CASPUR_WB_2.0	963,177	964,473 + 138 part	669,508 + 117 part	134,780 + 5 part	26,448 + 2 part
<i>Bubalus depressicornis</i> MNHNYannick_LA_1	MNHNYannick_LA_1	1,027,469	1,023,163 + 5,278 part	702,282 + 4,582 part	131,966 + 204 part	24,994 + 37 part
<i>Bubalus kerabau</i> CUSA_SWP	CUSA_SWP	1,042,862	1,046,662 + 87 part	752,170 + 70 part	151,809 + 10 part	29,488 + 6 part
<i>Syncerus caffer</i> ASM640878v2	ABF221	1,061,091	1,064,542 + 229 part	750,719 + 171 part	150,033 + 10 part	29,460 + 1 part

Table 6. Genes predicted with homology-based prediction method.

Name/assembly name (NCBI)	ID	Genes	Partial genes	Total	% of reference's annotated genes (n = 33,348)
<i>Bubalus bubalis</i> Jaffrabadi_v3.0	AAUIN_1	10,804	20,895	31,699	95.05
<i>Bubalus bubalis</i> UOA_WB_1	UOA_WB_1	30,810	1,955	32,765	98.25
<i>Bubalus bubalis</i> Bubbub1.0	Bubbub1.0	11,039	20,983	32,022	96.02
<i>Bubalus bubalis</i> ASM299383v1	EGYBUF_1.0	1,345	23,770	25,115	75.31
<i>Bubalus bubalis</i> UMD_CASPUR_WB_2.0	UMD_CASPUR_WB_2.0	18,656	13,271	31,927	95.74
<i>Bubalus depressicornis</i> MNHNYannick_LA_1	MNHNYannick_LA_1	19,148	13,245	32,393	97.14
<i>Bubalus kerabau</i> CUSA_SWP	CUSA_SWP	28,349	3,419	31,768	95.26
<i>Syncerus caffer</i> ASM640878v2	ABF221	8,763	21,575	30,338	90.97

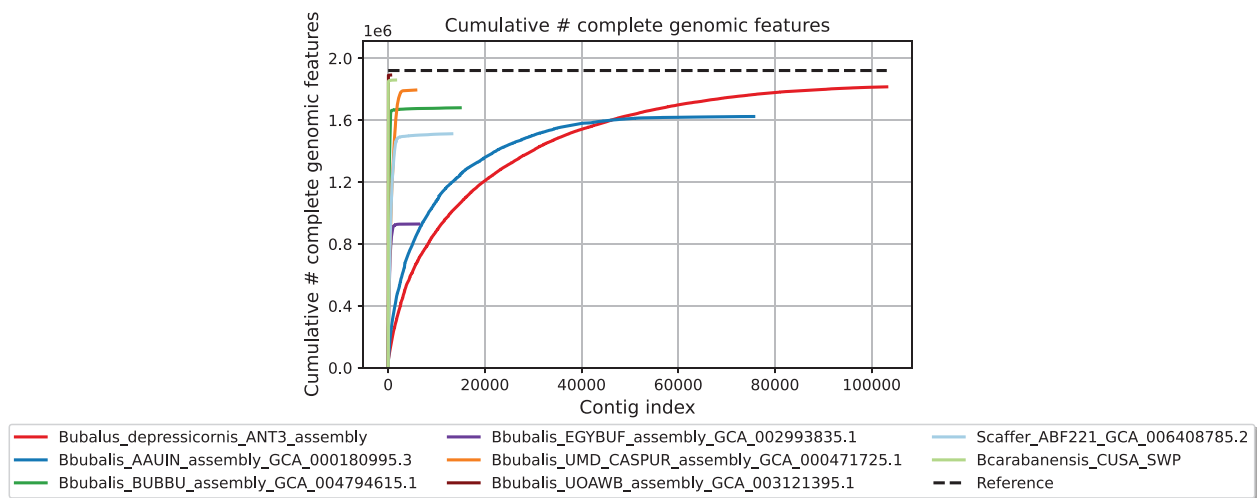


Fig. 5. Complete genomic features identified in the lowland anoa assembly and compared to other assemblies using the river buffalo (*Bubalus bubalis*) NDD_SH1 reference sequence and annotations.

BUSCO Assessment Results

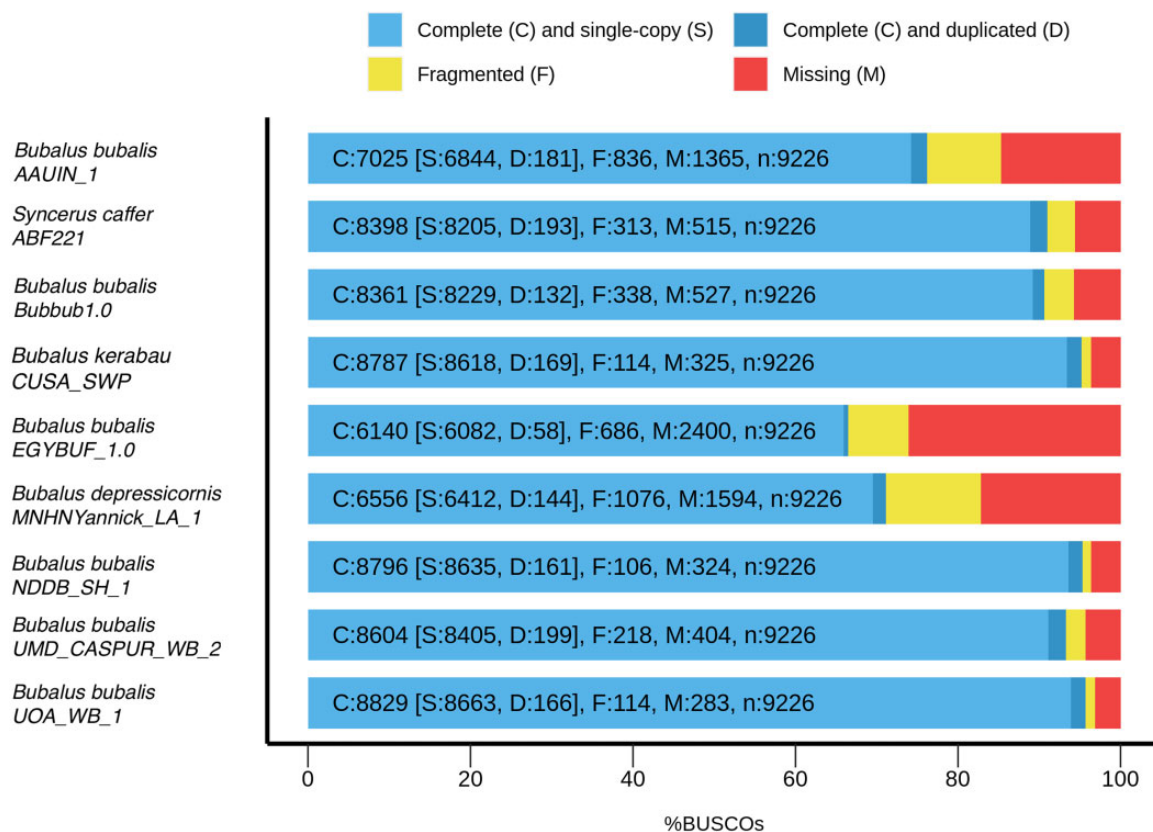


Fig. 6. BUSCO results of the genome assembly of the lowland anoa (*Bubalus depressicornis*) compared to other publicly available buffalo genome assemblies.

of the lowland anoa genome assembly. Nevertheless, the total number of genes predicted still represents an improvement over some of the compared assemblies (Table 6).

When predicting mammalian orthologs with BUSCO, the lowland anoa genome assembly contained 6,556 (71.1%) complete BUSCOs, of which 6,412 (69.5%) were single copy and 144 (1.6%) were duplicated. The number of fragmented BUSCOs was 1,076 (11.7%), whilst 1,594 (17.2%) were missing. The BUSCO results

indicate an acceptable level of genome completeness (<70%, Simão et al. 2015) for downstream analyses for the anoa genome assembly, and a slight improvement over the Egyptian river buffalo assembly (EGYBUF_1.0, Fig. 6).

Mammalian genomes contain large families of repeats (Goodier and Kazazian 2008), such as long interspersed nuclear elements (LINEs), short interspersed nuclear elements (SINEs), and long-terminal repeats (LTRs). RepeatMasker revealed that

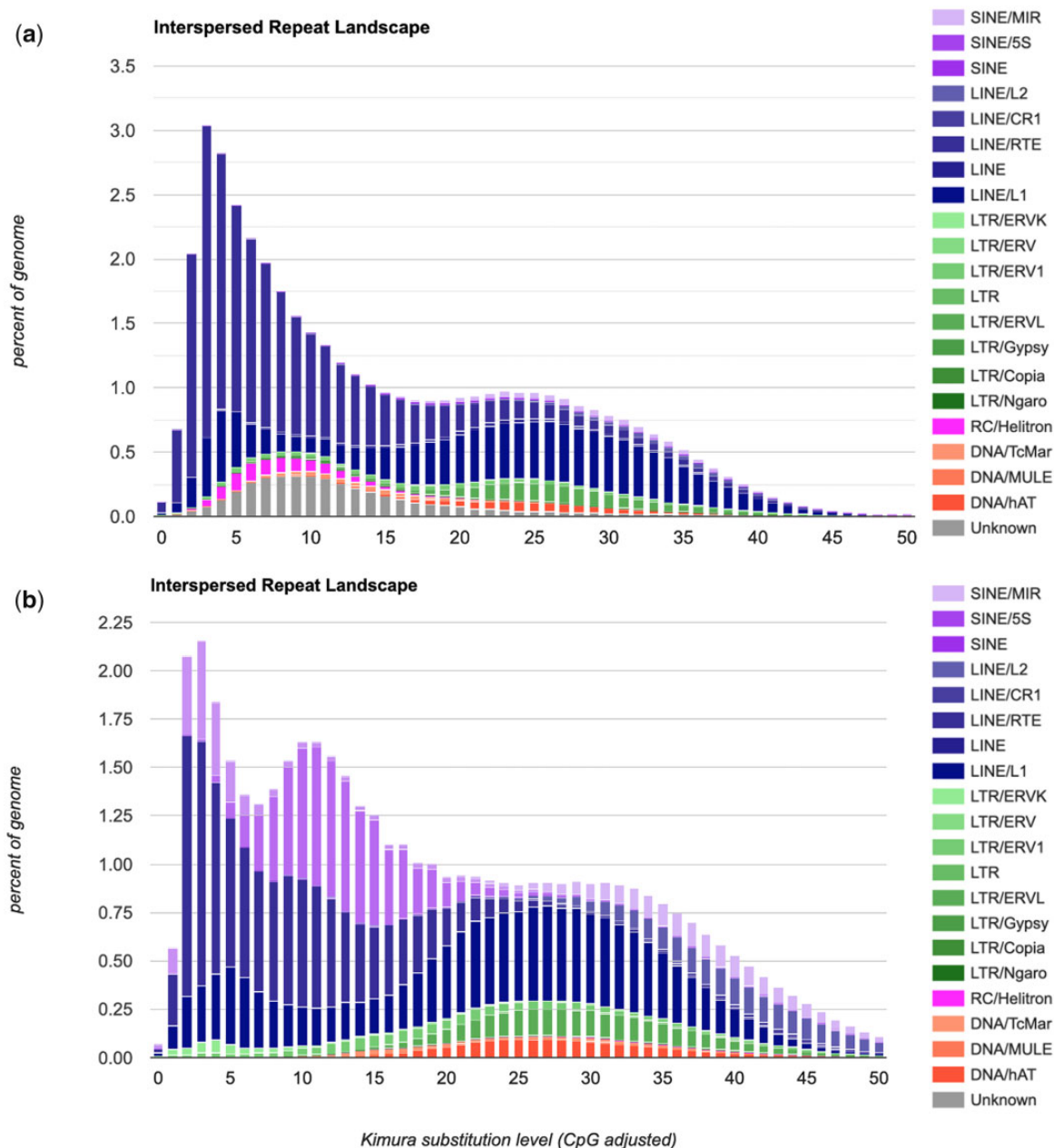
Table 7. Repeat sequence composition of the lowland anoa genome.

Family	Copy number of elements	Length occupied (bp)	% Genome
SINEs	296,064	26,945,915	1.03
LINES	2,864,468	786,815,034	30.04
LINE1	1,203,360	282,366,346	10.78
LINE2	101,415	13,911,301	0.53
RTE/Bov-B	1,461,651	481,114,012	18.37
LTR elements	362,123	81,208,077	3.10
DNA transposon	255,003	38,433,935	1.47
Small RNA	139,586	14,174,190	0.54
Satellites	269	52,169	0.00
Simple repeats	500,363	20,187,327	0.77
Low complexity	81,685	3,956,146	0.15
Unclassified	611,789	100,086,577	3.82
Total			42.12

42.12% of the lowland anoa genome is composed of repetitive regions (Table 7), which is comparable to data previously published for genome assemblies of river buffalo and other bovids (Deng et al. 2016; Low et al. 2019; Minto et al. 2019; El-Khishin et al. 2020). Results also agree with the repetitive content in the cattle genome (Fig. 7b). Both lowland anoa and cattle genomes showed 2 waves of repeat expansion in their repeat landscape (Fig. 7, a and b), suggesting a shared inheritance of such repeats. In the lowland anoa, the LINES were more abundant, representing 30.04% of the repeats, followed by LTRs representing 3.10% and SINEs representing 1.03% (Table 7).

Conclusion

To date, whole-genome sequencing has allowed the identification of variants involved in domestication and genetic improvement

**Fig. 7.** Interspersed repeat landscape of (a) the lowland anoa genome assembled in this study and (b) *Bos taurus*.

for several livestock species (Zimin *et al.* 2009; Canavez *et al.* 2012; Li *et al.* 2020; Rosen *et al.* 2020). However, the lack of wild buffalo genomes hinders further analyses addressing functional and evolutionary aspects of this group, as well as possible conservation efforts. The draft genome assembly of the lowland anoa reported here is expected to contribute to this gap in data availability, as this is the first draft genome assembly for wild Asian buffaloes. Furthermore, we showed that short-read Illumina sequencing data can still provide a cost-effective way of sequencing mammalian genomes to an adequate level of completeness for downstream comparative analyses.

Data availability

The raw data and assembly are available on NCBI under BioProject PRJNA849775. The genome assembly of the lowland anoa is available on NCBI under BioSample accession SAMN29133250. The raw data are available on the Sequence Read Archive (SRA) on NCBI under accession SRR21016826.

Acknowledgments

The authors thank the people of the *Ménagerie du Jardin des Plantes* who helped to collect the biopsy of the lowland anoa used in this study: Norin Chai, Gerard Dousseau, Christelle Hano, Abderrahmane Latreche, Claire Rejaud, Roland Simon, and Rudy Wedlarski. They would like to thank Huw Jones for the proofreading of the manuscript.

Funding

RR was supported by sDiv, Synthesis Centre of the German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, funded by the German Research Foundation (DFG–FZT 118, 202548816), and the German Research Foundation (DFG Research grant RO 5835/2-1).

Conflicts of interest

None declared.

Literature cited

Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*. 2014;30(15):2114–2120. doi:10.1093/bioinformatics/btu170.

Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, Boisvert S, Chapman JA, Chapuis G, Chikhi R, *et al.* Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate species. *Gigascience*. 2013;2(1):10–31. doi:10.1186/2047-217X-2-10.

Burton J, Wheeler P, Mustari A. *Bubalus depressicornis*. The IUCN Red List of Threatened Species™. 2016. doi:10.2305/IUCN.UK.2016-2.RLTS.T3126A46364222.

Canavez FC, Luche DD, Stothard P, Leite KRM, Sousa-Canavez JM, Plastow G, Meidanis J, Souza MA, Feijao P, Moore SS, *et al.* Genome sequence and assembly of *Bos indicus*. *J Hered*. 2012; 103(3):342–348. doi:10.1093/jhered/esr153.

Castelló JR. *Bovids of the World: Antelopes, Gazelles, Cattle, Goats, Sheep, and Relatives*. Princeton (NJ): Princeton University Press; 2016.

Chu J. Jupiter plot: a Circos-Based tool to Visualize Genome Assembly Consistency (Version 1.0). Github; 2018. [accessed 2022 Feb 22]. <https://github.com/JustinChu/JupiterPlot>.

Curaudeau M, Rozzi R, Hassanin A. The genome of the lowland anoa (*Bubalus depressicornis*) illuminates the origin of river and swamp buffalo. *Mol Phylogenet Evol*. 2021;161(March):107170. doi:10.1016/j.ympev.2021.107170.

Deng T, Pang C, Lu X, Zhu P, Duan A, Tan Z, Huang J, Li H, Chen M, Liang X. De novo transcriptome assembly of the Chinese swamp buffalo by RNA sequencing and SSR marker discovery. *PLoS One*. 2016;11(1):e0147132. doi:10.1371/journal.pone.0147132.

El-Khishin DA, Ageez A, Saad ME, Ibrahim A, Shokrof M, Hassan LR, Abouelhoda MI. Sequencing and assembly of the Egyptian buffalo genome. *PLoS One*. 2020;15(8):e0237087. doi:10.1371/journal.pone.0237087.

Ewels P, Magnusson M, Lundin S, Käller M. MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics*. 2016;32(19):3047–3048. doi:10.1093/bioinformatics/btw354.

Fitzinger LJ. Der Sunda-Büffel (*Bubalus kerabau*). In: *Wissenschaftlich-populäre Naturgeschichte der Säugethiere in ihren sämtlichen Hauptformen*, V. Kaiserlich-Königlichen Hof- und Staatsdruckerei. Wien. p. 329. 1860.

Gnerre S, MacCallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, Sharpe T, Hall G, Shea TP, Sykes S, *et al.* High-quality draft assemblies of mammalian genomes from massively parallel sequence data. *Proc Natl Acad Sci U S A*. 2011;108(4):1513–1518. doi:10.1073/pnas.1017351108.

Goodier JL, Kazazian HH. Retrotransposons revisited: the restraint and rehabilitation of parasites. *Cell*. 2008;135(1):23–35. doi:10.1016/j.cell.2008.09.022.

Groves CP. Systematics of the anoa (Mammalia, Bovidae). *Beaufortia*. 1969;17:1–12.

Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics*. 2013;29(8):1072–1075. doi:10.1093/bioinformatics/btt086.

Hassanin A, Delsuc F, Ropiquet A, Hammer C, Jansen Van Vuuren B, Matthee C, Ruiz-Garcia M, Catzeflis F, Areskouk V, Nguyen TT, *et al.* Pattern and timing of diversification of Cetartiodactyla (Mammalia, Laurasiatheria), as revealed by a comprehensive analysis of mitochondrial genomes. *C R Biol*. 2012;335(1):32–50. doi:10.1016/j.crv.2011.11.002.

Heude PM. Note sur le petit buffle sauvage de l'île de Mindoro (Philippines). *Mémoires Concern l'histoire Nat L'Empire Chinois*. 1888;2(4):50.

International Human Genome Sequencing Consortium. Finishing the euchromatic sequence of the human genome. *Nature*. 2004; 431(7011):931–945.

IUCN. Red List Threat Species. The IUCN Red List of Threatened Species; 2022. [accessed 2022 Feb 15]. <https://www.iucnredlist.org/>

Johnson JM, Edwards S, Shoemaker D, Schadt EE. Dark matter in the genome: evidence of widespread transcription detected by microarray tiling experiments. *Trends Genet*. 2005;21(2):93–102. doi:10.1016/j.tig.2004.12.009.

Kerr R. *Arnee Bos arnee*. In: Strahan A, Cadell T, Creech W, editors. *The Animal Kingdom or Zoological System of the Celebrated Sir Charles Linnaeus*. Class I. Mammalia. Edinburgh and London. 1792. p. 336.

Kriventseva EV, Kuznetsov D, Tegenfeldt F, Manni M, Dias R, Simão FA, Zdobnov EM. OrthoDB v10: sampling the diversity of animal, plant, fungal, protist, bacterial and viral genomes for

- evolutionary and functional annotations of orthologs. *Nucleic Acids Res.* 2019;47(D1):D807–D811. doi:[10.1093/nar/gky1053](https://doi.org/10.1093/nar/gky1053).
- Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA. Circos: an information aesthetic for comparative genomics. *Genome Res.* 2009;19(9):1639–1645. doi:[10.1101/gr.092759.109](https://doi.org/10.1101/gr.092759.109).
- Li X, Yang J, Shen M, Xie XL, Liu GJ, Xu YX, Lv FH, Yang H, Yang YL, Liu CB, et al. Whole-genome resequencing of wild and domestic sheep identifies genes associated with morphological and agronomic traits. *Nat Commun.* 2020;11(1):1–16. doi:[0.1038/s41467-020-16485-11](https://doi.org/10.1038/s41467-020-16485-11).
- Linnaeus. *Bubalus bubalis*. GBIF Secr; 1758. [accessed 2022 Mar 14]. <https://www.gbif.org/species/7422937>.
- Low WY, Tearle R, Bickhart DM, Rosen BD, Kingan SB, Swale T, Thibaud-Nissen F, Murphy TD, Young R, Lefevre L, et al. Chromosome-level assembly of the water buffalo genome surpasses human and goat genomes in sequence contiguity. *Nat Commun.* 2019;10(1):260–211. doi:[10.1038/s41467-018-08260-0](https://doi.org/10.1038/s41467-018-08260-0).
- Luo X, Zhou Y, Zhang B, Zhang Y, Wang X, Feng T, Li Z, Cui K, Wang Z, Luo C, et al. Understanding divergent domestication traits from the whole-genome sequencing of swamp- and river-buffalo populations. *Natl Sci Rev.* 2020;7(3):686–701. doi:[10.1093/nsr/nwaa024](https://doi.org/10.1093/nsr/nwaa024).
- Majoros WH, Pertea M, Salzberg SL. TigrScan and GlimmerHMM: two open source ab initio eukaryotic gene-finders. *Bioinformatics.* 2004;20(16):2878–2879. doi:[10.1093/bioinformatics/bth315](https://doi.org/10.1093/bioinformatics/bth315).
- Manchanda N, Portwood JL, Woodhouse MR, Seetharam AS, Lawrence-Dill CJ, Andorf CM, Hufford MB. GenomeQC: a quality assessment tool for genome assemblies and gene structure annotations. *BMC Genomics.* 2020;21(1):1–9. doi:[10.1186/s12864-020-6568-2](https://doi.org/10.1186/s12864-020-6568-2).
- Manni M, Berkeley MR, Seppely M, Simão FA, Zdobnov EM. BUSCO update: novel and streamlined workflows along with broader and deeper phylogenetic coverage for scoring of eukaryotic, prokaryotic, and viral genomes. *Mol Biol Evol.* 2021;38(10):4647–4654. doi:[10.1093/molbev/msab199](https://doi.org/10.1093/molbev/msab199).
- Marçais G, Yorke JA, Zimin A. QuorUM: an error corrector for Illumina reads. *PLoS One.* 2015;10(6):e0130821. doi:[10.1371/journal.pone.0130821](https://doi.org/10.1371/journal.pone.0130821).
- Mikheenko A, Prjibelski A, Saveliev V, Antipov D, Gurevich A. Versatile genome assembly evaluation with QUAST-LG. *Bioinformatics.* 2018;34(13):i142–i150. doi:[10.1093/bioinformatics/bty266](https://doi.org/10.1093/bioinformatics/bty266).
- Mintoo AA, Zhang H, Chen C, Moniruzzaman M, Deng T, Anam M, Emdadul Huque QM, Guang X, Wang P, Zhong Z, et al. Draft genome of the river water buffalo. *Ecol Evol.* 2019;9(6):3378–3388. doi:[10.1002/ece3.4965](https://doi.org/10.1002/ece3.4965).
- Nguyen TT, Aniskin VM, Gerbault-Seureau M, Planton H, Renard JP, Nguyen BX, Hassanin A, Volobouev VT. Phylogenetic position of the saola (*Pseudoryx nghetinhensis*) inferred from cytogenetic analysis of eleven species of Bovidae. *Cytogenet Genome Res.* 2008;122(1):41–54. doi:[10.1159/000151315](https://doi.org/10.1159/000151315).
- Ouwens PA. Contribution a la connaissance des mammifères de Célébes. *Bull Dépt Agric Indes Néerl.* 1910;38(Zool. 6):1–7.
- Priyono DS, Solihin DD, Farajallah A, Purwantara B. The first complete mitochondrial genome sequence of the endangered mountain anoa (*Bubalus quarlesi*) (Artiodactyla: Bovidae) and phylogenetic analysis. *J Asia-Pacific Biodivers.* 2020;13(2):123–133. doi:[10.1016/j.japb.2020.01.006](https://doi.org/10.1016/j.japb.2020.01.006).
- Rosen BD, Bickhart DM, Schnabel RD, Koren S, Elsik CG, Tseng E, Rowan TN, Low WY, Zimin A, Couldrey C, et al. De novo assembly of the cattle reference genome with single-molecule sequencing. *Gigascience.* 2020;9(3):1–9. doi:[10.1093/gigascience/giaa021](https://doi.org/10.1093/gigascience/giaa021).
- Schreiber A, Seibold I, Nötzold G, Wink M. Cytochrome b gene haplotypes characterize chromosomal lineages of anoa, the Sulawesi dwarf buffalo (Bovidae: *Bubalus* sp.). *J Hered.* 1999;90(1):165–176. doi:[10.1093/jhered/90.1.165](https://doi.org/10.1093/jhered/90.1.165).
- Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics.* 2015;31(19):3210–3212. doi:[10.1093/bioinformatics/btv351](https://doi.org/10.1093/bioinformatics/btv351).
- Smith CH. The seventh order of the Mammalia. The Ruminantia. In: E Griffith, CH Smith, E Pidgeon, editors. *The Animal Kingdom Arranged in Conformity with Its Organization*, by the Baron Cuvier, Member of the Institute of France, with Additional Descriptions of All the Species Hitherto Named, and of Many Not before Noticed. London: Whittaker G.B; 1827. p. 293.
- Tarailo-Graovac M, Chen N. Using RepeatMasker to identify repetitive elements in genomic sequences. *Curr Protoc Bioinforma.* 2009;25:1–14. doi:[10.1002/0471250953.bi0410s25](https://doi.org/10.1002/0471250953.bi0410s25).
- Weissensteiner MH, Suh A. Repetitive DNA: the dark matter of avian genomics. In: Kraus R, editor. *Avian Genomics in Ecology and Evolution*. Cham: Springer; 2019. pp. 93–150. https://doi.org/10.1007/978-3-030-16477-5_5.
- Zhang S, Liu W, Liu X, Du X, Zhang K, Zhang Y, Song Y, Zi Y, Qiu Q, Lenstra JA, et al. Structural variants selected during yak domestication inferred from long-read whole-genome sequencing. *Mol Biol Evol.* 2021;38(9):3676–3680. doi:[10.1093/molbev/msab134](https://doi.org/10.1093/molbev/msab134).
- Zimin AV, Delcher AL, Florea L, Kelley DR, Schatz MC, Puiu D, Hanrahan F, Pertea G, Van Tassell CP, Sonstegard TS, et al. A whole-genome assembly of the domestic cow, *Bos taurus*. *Genome Biol.* 2009;10(4):R42. doi:[10.1186/gb-2009-10-4-r42](https://doi.org/10.1186/gb-2009-10-4-r42).
- Zimin AV, Marçais G, Puiu D, Roberts M, Salzberg SL, Yorke JA. The MaSuRCA genome assembler. *Bioinformatics.* 2013;29(21):2669–2677. doi:[10.1093/bioinformatics/btt476](https://doi.org/10.1093/bioinformatics/btt476).
- Zimin AV, Puiu D, Luo MC, Zhu T, Koren S, Marçais G, Yorke JA, Dvořák J, Salzberg SL. Hybrid assembly of the large and highly repetitive genome of *Aegilops tauschii*, a progenitor of bread wheat, with the MaSuRCA mega-reads algorithm. *Genome Res.* 2017;27(5):787–792. doi:[10.1101/gr.213405.116](https://doi.org/10.1101/gr.213405.116).

Communicating editor: D.-J. de Koning