

1 Genome Reference Assembly for Bottlenecked Southern Australian Koalas

2 Blanchard, A.M.,^{1*} Emes, R.D.,¹ Greenwood, A.D.,³ Holmes, N.,² Loose, M.W.,² McEwen, G.K.,³
3 Meers, J.,⁴ Speight, N.⁵ and Tarlinton, R.E.¹

4 ¹ University of Nottingham, School of Veterinary Medicine and Science, Sutton Bonington,
5 Leicestershire, LE12 5RD, UK

6 ² University of Nottingham, School of Life Sciences, Nottingham NG7 2RD, UK

7 ³ Leibniz Institute for Zoo and Wildlife Research, 10315 Berlin, Germany

8 ⁴ The University of Queensland, School of Veterinary Science, Australia

9 ⁵ The University of Adelaide, School of Animal and Veterinary Sciences, Adelaide, Australia

10

11 *Corresponding author: adam.blanchard@nottingham.ac.uk

12 **Abstract**

13 Koala populations show marked differences in inbreeding levels and in the presence or absence
14 of the endogenous Koala Retrovirus (KoRV). These genetic differences among populations may
15 lead to severe disease impacts threatening koala population viability. In addition, the recent
16 colonization of the koala genome by KoRV provides a unique opportunity to study the process of
17 retroviral adaptation to vertebrate genomes and the impact this has on speciation, genome
18 structure and function. The genome build described here is from an animal from the bottlenecked
19 “Southern” population free of endogenous and exogenous KoRV. It provides a more contiguous
20 genome build than the previous koala reference derived from an animal from a more outbred
21 “Northern” population and is the first koala genome from a KoRV polymerase free animal.

22

23

1 **Significance**

2 This high-quality genome build provides a base line comparator for studies of koala genetics and
3 retroviral integration. It is from a genetically distinct population than the current koala reference
4 genome and does not contain intact endogenous Koala retrovirus.

5 **Introduction**

6 Koalas are an iconic marsupial species classed as vulnerable on the IUCN red list. The species
7 suffers a number of threats including habitat loss and disease with climate change driven fire
8 events further decimating numbers in recent years (Charalambous & Narayan 2020). The disease
9 threats to the population are complicated by stark differences in the disease patterns in different
10 populations driven by underlying genetic differences (Sarker et al. 2020; Tarlinton et al. 2021).
11 Wild koalas are confined to the Eastern Seaboard of Australia, there are 5 major genetic groups
12 (Lott et al. 2022) but for the purposes of population management two major genetic splits are
13 recognised: “Northern” (New South Wales and Queensland) and “Southern” (Victoria and South
14 Australia), the border between the states of New South Wales and Victoria forming a hard cut off
15 between the two populations (Neaves et al. 2016; Quigley et al. 2021).

16 Koalas in the southern states were essentially extinct by 1920 due to hunting pressure and were
17 restocked across their southern range from a very small number of animals (possibly as few as
18 18) sourced from offshore island refugia (Martin et al. 1999). As a result of this animals in the
19 southern population have a markedly reduced genetic diversity compared with animals in the
20 northern population (Neaves et al. 2016; Ruiz-Rodriguez et al. 2016; Johnson et al. 2018;
21 Tarlinton et al. 2021). Our own work has demonstrated that many genes are homozygous in the
22 southern animals (Tarlinton et al. 2021). Animals in the southern populations suffer from a
23 number of diseases, such as oxalate nephrosis and testicular aplasia that are not routinely seen in

1 northern populations (Fabijan et al. 2020; Tarlinton et al. 2021) and are thought to have an
2 underlying genetic basis (Cristescu et al. 2009; Speight et al. 2020).

3 The other major difference both disease and genetics wise between Northern and
4 Southern animals is the presence of a functional recently endogenized retrovirus (Koala
5 Retrovirus or KoRV) in all Northern koalas but not in the Southern (Quigley et al. 2021; Blyton
6 et al. 2022). Both Southern and Northern animals may have exogenous infectious KoRV but the
7 rate of KoRV associated neoplasia is substantially lower in Southern koalas (Sarker et al. 2020;
8 Joyce et al. 2021; Quigley et al. 2021). While the definitive link is less clear than for neoplasia
9 (McEwen et al. 2021) KoRV is also thought to cause underlying immunosuppression
10 predisposing to chlamydia disease, which is also seen at a lower rate in Southern populations
11 (Polkinghorne et al. 2013; Sarker et al. 2020). Endogenous retroviruses are present in all
12 vertebrate genomes studied to date and the entrance of these transposable elements into genomes
13 is thought to be a major introduction of genetic diversity, potentially triggering speciation.
14 However, most examples in genomes are ancient and are essentially represented by inactive
15 viruses (Zheng et al. 2022). They are thought to be the remnants of past infectious viral
16 integrations that have managed to enter germ line cells and become fixed in a species. KoRV is
17 part of a very small group of recently endogenized viruses, integrated sometime between 200 and
18 49,000 years ago (Ishida et al. 2015), and is unique in that parts of the species range do not yet
19 have endogenous polymerase gene containing KoRV at all (Quigley et al. 2021).

20 To complicate matters further both Northern and Southern koalas have evidence of
21 historical KoRV infection as defective recombinant sequences between KoRV and another older
22 endogenous retroelement (Phascolarctos endogenous retrovirus or PhER), known as recKoRVs
23 (Löber et al. 2018; Tarlinton et al. 2022). It is not entirely clear how endogenous and exogenous
24 KoRV and recKoRV interact and whether they enhance or inhibit each other's replication and

1 disease occurrence, but the scenario provides a unique opportunity to study the impacts of the
 2 entrance of a new class of retroelements into a mammalian genome in real time rather than by
 3 phylogenetic inference of this fundamental genomic process (Tarlinton et al. 2022).

4 There are two other published koala genomes (Johnson et al. 2018) derived from northern
 5 animals “Bilbo” and “Pacific Chocolate” (Johnson et al. 2018), alongside several additional
 6 transcriptome resources (Hobbs et al. 2014; Abts et al. 2015; Tarlinton et al. 2022). The most
 7 complete existing genome for “Bilbo” is assembled at a contig level (into 1,907 contigs with an
 8 N50 of 11.6 Mb). Here we present a genome build of a Southern Australian animal “Wilpena” for
 9 use in comparative genomics of koala populations and studies of retroviral integration. This
 10 genome is more contiguous than the current reference sequence (1,265 contigs, N50 = 48.8 Mb)
 11 and from an animal known to be free of both endogenous and exogenous replication competent
 12 KoRV (Tarlinton et al. 2022).

14 Results and Discussion

15 Using 58 GB of ONT data (consisting of 2,572,260 reads with a mean read length of 24kb and
 16 mean Q score of 13.7) and 1,289 million (2x150bp >Q30) Illumina reads were assembled into a
 17 draft genome using Flye. This resulted in an N50 of 48,782,874 bases and a length of
 18 3,233,824,327 bp. A first pass polish using Medaka and final polish with Polca using the Illumina
 19 data resulting in a final high-quality genome assembly with an N50 of 48,800,306 bases, 1,265
 20 contigs and a total genome size of 3,234,982,288 bp (Table 1).

22 Table 1: Summary of the genome assembly

Genome	Wilpena This Study	Bilbo GCA_002099425	Pacific Chocolate GCA_900166895
Assembly Size	3,234,982,288 bp	3,192,581,492 bp	3,358,707,742 bp
Number of Contigs	1265	1907	796,464

Contigs \geq 5,000 bp	1222	1804	16,989
Contigs \geq 50,000 bp	651	662	8361
Contigs N50	48,800,306 bp	11,587,828 bp	880973 bp
Contigs N75	22,144,309 bp	6,857,650 bp	321,283
Contigs L50	17	85	1100
Contigs L75	41	173	2591
Largest Contig	232,027,266 bp	40,558,015 bp	5,231,295 bp
GC Content (%)	39.09	39.05	39.03
BUSCO Completeness (%)	92.9	94.0	90.0
Genes	27,669	32,109	33,654

1 The contigs were assessed for putative contamination using Conterminator (Steinegger *et al.*,
2 2020). From 1,265 contigs, 1,247 were assigned as koala and 18 were flagged as containing
3 potential contamination. Of those 18, assignments were for North American Opossum (n=2),
4 Common Brushtail (n=2), Grey Short-Tailed Opossum (n=2), Common Wombat (n=7) and Koala
5 Retrovirus (n=2). However, the same eight contigs were flagged multiple times with close
6 marsupial relatives and so are unlikely to be true contamination. There were two contigs assigned
7 as koala retrovirus, these are non-functional partial recKoRV sequences (partial KoRV env and
8 LTR) as reported previously in this animal (Tarlinton *et al.* 2022) and are not full length
9 endogenous or exogenous KoRV. The genome was soft masked using RED (Supplementary
10 Table 1) and genes predicted used braker2 along with publicly available Koala RNASeq data
11 from multiple biological sites, predicting 52,384 putative genes. Functional annotation using
12 EggNOG mapper identified 27,669 genes with transcriptional support (Supplementary Table 2).

13

14 **Conclusion**

15 A highly contiguous reference genome, from a distinct southern population, is invaluable to
16 understanding the challenges faced in conservation genetics for future breeding programmes of
17 Koalas. Not only will this enable more comprehensive comparative genomics to take place, it
18 will also allow researchers to fully understand non-functional KoRV integration sites and
19 whether they appear in similar regions of the genome to the northern population.

1

2 **Materials and Methods**

3 *Sample Collection*

4 DNA was derived from liver tissue from a 3 year old female south Australian Koala, housed in a
5 collection in the UK. The animal was originally derived from the Mt Lofty Ranges and Kangaroo
6 Island populations in South Australia. Sample collection and nanopore sequencing from this
7 animal was described in (Tarlinton et al. 2022). Ethics approval for the use of post mortem
8 material was granted by the University of Nottingham School of Veterinary Medicine and
9 Science Committee for Animal Care and Research Ethics

10

11 *Sample Preparation*

12 DNA was extracted from frozen liver tissue using the QIAGEN Genomic-tip 100/G Kit and the
13 QIAGEN Genomic Buffer Set (QIAGEN; 10243 and 19060). Frozen tissue was ground under
14 liquid and 100 mg of frozen powder was added to 9.2 ml of buffer G2 containing 5 µl of RNase
15 A (100 mg/ml) (QIAGEN; 19101) and the suspension was incubated at room temperature for 10
16 min. Proteinase K (100 µl) (QIAGEN; 19131) was added and the suspension was incubated at 50
17 °C for 1.5 h. The Genomic-tip protocol was then followed, according to the QIAGEN Genomic
18 DNA Handbook 06/2015.

19

20 *Genome Sequencing*

21 Genomic DNA was needle sheared 30 times with a 26G needle (BD; 300300) and then treated
22 with the Short Read Eliminator (SRE) Kit (Circulomics; SS-100-101-01) to remove fragments <
23 10 kb and progressively deplete fragments shorter than 25 kb. The processed DNA was used to
24 generate a sequencing library using the Genomic DNA by Ligation PromethION Kit (Oxford

1 Nanopore Technologies; SQK-LSK109). Library quantification was performed using the Qubit
2 fluorometer and the Qubit dsDNA HS Assay Kit (ThermoFisher; Q32854) and 600 ng of library
3 was run over one PromethION flow cell (Oxford Nanopore Technologies; FLO-PRO002) on a
4 PromethION Beta device. The same DNA preparation was subjected to Illumina Novaseq 6000
5 paired end 150bp read sequencing (with automated plant and whole genome library preparation)
6 by Novogene, Cambridge UK.

7

8 *Read Processing*

9 Illumina reads (both RNA and DNA) were trimmed to remove adaptors and reads with an overall
10 quality of <Q30 using FastP v0.23.1 (Chen et al. 2018). The raw Nanopore data was base called
11 using Guppy v6.1.7+21b93d1a5 and the super-accurate mode
12 (<https://community.nanoporetech.com/downloads>). Nanopore adaptors were removed using
13 Porechop v0.2.4 (Wick et al. 2017) and reads shorter than 1000bp and with a quality of <Q10
14 were removed with NanoFilt v2.6.0 (De Coster et al. 2018).

15

16 *Assembly*

17 The Nanopore reads were assembled using Flye v2.9.1 with the --nano-hq and --keep-haplotypes
18 flags. The Flye draft assembly was first polished with Medaka v1.6.1
19 (<https://github.com/nanoporetech/medaka>) and then with the Illumina reads using POLCA (from
20 MaSuRCA v4.0.9) (Zimin & Salzberg 2020). The resulting polished assembly was then gap filled
21 using Samba (from MaSuRCA v4.0.9) (Zimin & Salzberg 2022) before being assessed for
22 completeness using BUSCO v5.4.2 (Manni et al. 2021) in genome mode with the Mammalian
23 lineage database.

24

1 *Contamination assessment*

2 Each contig was assigned a taxonomic ID using blastn (Altschul *et al.*, 1990), this was parsed
3 into Conterminator (Steinegger *et al.*, 2020) to identify potential regions of contamination in the
4 genome.

6 *Annotation*

7 The final version of the genome was parsed though REpeat Detector (RED) v 1.16 (Girgis 2015)
8 to soft mask regions of repetitive elements. An index was created using HISAT2 v2.2.1 (Kim *et*
9 *al.* 2019) , and RNASeq data from accession: PRJNA230900 (Hobbs *et al.* 2017) was aligned
10 producing sam files. SAMTools v1.15 (Danecek *et al.* 2021) was used to convert sam to bam
11 before being used in Braker2 v2.1.6 (Brůna *et al.* 2021) for genome annotation. Functional
12 annotation was completed using EggNOG mapper v2.1 (Cantalapiedra *et al.* 2021).

14 **Acknowledgements**

15 Sample access was facilitated by Longleat Safari park, Funding was provided by the University
16 of Nottingham and A.D.G. and G.K.M. were supported by grant GR 3924/15-1 from the
17 Deutsche Forschungsgemeinschaft (DFG)

19 **Author Contributions**

20 RT oversaw project management, wrote funding proposals, collected post mortem samples and
21 wrote parts of the manuscript. NH performed the DNA extraction and nanopore sequencing. ML
22 and AB performed bioinformatics analysis. AB wrote parts of the manuscript and performed data
23 deposition in public repositories. AG and GM project managed the illumina sequencing. NS

1 performed initial screening testing for KoRV on the animal. NS and JM provided critical review
2 of the manuscript. All authors read and reviewed the manuscript.

3

4 **Data Availability**

5 All raw sequence data is available on the NCBI SRA under the accession SAMN30742200. The
6 final genome build (*Phascolarctos cinereus* K01) is available under the NCBI accession number
7 JAOEJA000000000.

8

9 **References**

10 Abts KC, Ivy JA, DeWoody JA. 2015. Immunomics of the koala (*Phascolarctos cinereus*).
11 *Immunogenetics*. 67:305–321. doi: 10.1007/s00251-015-0833-6.

12

13 Blyton MDJ, Young PR, Moore BD, Chappell KJ. 2022. Geographic patterns of koala retrovirus
14 genetic diversity, endogenization, and subtype distributions. *Proc. Natl. Acad. Sci.*
15 119:e2122680119. doi: 10.1073/pnas.2122680119.

16 Cantalapiedra CP, Hernández-Plaza A, Letunic I, Bork P, Huerta-Cepas J. 2021. eggNOG-
17 mapper v2: Functional Annotation, Orthology Assignments, and Domain Prediction at the
18 Metagenomic Scale Tamura, K, editor. *Mol. Biol. Evol.* 38:5825–5829. doi:
19 10.1093/molbev/msab293.

20

21 Charalambous R, Narayan E. 2020. A 29-year retrospective analysis of koala rescues in New
22 South Wales, Australia Yue, B-S, editor. *PLOS ONE*. 15:e0239182. doi:
23 10.1371/journal.pone.0239182.

24

- 1 Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor.
2 Bioinformatics. 34:i884–i890. doi: 10.1093/bioinformatics/bty560.
3
- 4 Cristescu RH et al. 2009. Inbreeding and testicular abnormalities in a bottlenecked population of
5 koalas (*Phascolarctos cinereus*). Wildl. Res. 36:299–308.
6
- 7 De Coster W, D’Hert S, Schultz DT, Cruts M, Van Broeckhoven C. 2018. NanoPack: visualizing
8 and processing long-read sequencing data Berger, B, editor. Bioinformatics. 34:2666–2669. doi:
9 10.1093/bioinformatics/bty149.
10
- 11 Fabijan J et al. 2020. Pathological Findings in Koala Retrovirus-positive Koalas (*Phascolarctos*
12 *cinereus*) from Northern and Southern Australia. J. Comp. Pathol. 176:50–66. doi:
13 10.1016/j.jcpa.2020.02.003.
14
- 15 Hobbs M et al. 2014. A transcriptome resource for the koala (*Phascolarctos cinereus*): insights
16 into koala retrovirus transcription and sequence diversity. BMC Genomics. 15:786. doi:
17 10.1186/1471-2164-15-786.
18
- 19 Ishida Y, Zhao K, Greenwood AD, Roca AL. 2015. Proliferation of Endogenous Retroviruses in
20 the Early Stages of a Host Germ Line Invasion. Mol. Biol. Evol. 32:109–120. doi:
21 10.1093/molbev/msu275.
22
- 23 Johnson RN et al. 2018. Adaptation and conservation insights from the koala genome. Nat.
24 Genet. 50:1102–1111. doi: 10.1038/s41588-018-0153-5.

- 1
- 2 Joyce BA, Blyton MDJ, Johnston SD, Young PR, Chappell KJ. 2021. Koala retrovirus genetic
3 diversity and transmission dynamics within captive koala populations. *Proc. Natl. Acad. Sci.*
4 118:e2024021118. doi: 10.1073/pnas.2024021118.
- 5
- 6 Löber U et al. 2018. Degradation and remobilization of endogenous retroviruses by
7 recombination during the earliest stages of a germ-line invasion. *Proc. Natl. Acad. Sci.*
8 115:8609–8614. doi: 10.1073/pnas.1807598115.
- 9
- 10 Lott MJ et al. 2022. Future-proofing the koala: Synergising genomic and environmental data for
11 effective species management Waits, L, editor. *Mol. Ecol.* 31:3035–3055. doi:
12 10.1111/mec.16446.
- 13
- 14 Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM. 2021. BUSCO Update: Novel and
15 Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage for Scoring of
16 Eukaryotic, Prokaryotic, and Viral Genomes Kelley, J, editor. *Mol. Biol. Evol.* 38:4647–4654.
17 doi: 10.1093/molbev/msab199.
- 18
- 19 Martin R, Handasyde KA, Lee AK. 1999. *The koala: natural history, conservation and*
20 *management*. 2nd ed. UNSW Press: Sydney, Australia.
- 21
- 22 McEwen GK et al. 2021. Retroviral integrations contribute to elevated host cancer rates during
23 germline invasion. *Nat. Commun.* 12:1316. doi: 10.1038/s41467-021-21612-7.
- 24

- 1 Neaves LE et al. 2016. Phylogeography of the Koala, (*Phascolarctos cinereus*), and Harmonising
2 Data to Inform Conservation Banks, SC, editor. PLOS ONE. 11:e0162207. doi:
3 10.1371/journal.pone.0162207.
- 4
- 5 Polkinghorne A, Hanger J, Timms P. 2013. Recent advances in understanding the biology,
6 epidemiology and control of chlamydial infections in koalas. *Vet. Microbiol.* 165:214–223. doi:
7 10.1016/j.vetmic.2013.02.026.
- 8
- 9 Quigley BL, Wedrowicz F, Hogan F, Timms P. 2021. Phylogenetic and geographical analysis of
10 a retrovirus during the early stages of endogenous adaptation and exogenous spread in a new
11 host. *Mol. Ecol.* 30:2626–2640. doi: 10.1111/mec.15735.
- 12
- 13 Ruiz-Rodriguez CT et al. 2016. Koalas (*Phascolarctos cinereus*) From Queensland Are
14 Genetically Distinct From 2 Populations in Victoria. *J. Hered.* 107:573–580. doi:
15 10.1093/jhered/esw049.
- 16 Sarker N et al. 2020. Koala retrovirus viral load and disease burden in distinct northern and
17 southern koala populations. *Sci. Rep.* 10:263. doi: 10.1038/s41598-019-56546-0.
- 18
- 19 Speight N, Bacci B, Stent A, Whiteley P. 2020. Histological survey for oxalate nephrosis in
20 Victorian koalas (*Phascolarctos cinereus*). *Aust. Vet. J.* 98:467–470. doi: 10.1111/avj.12986.
- 21
- 22 Tarlinton RE et al. 2022. Differential and defective transcription of koala retrovirus indicates the
23 complexity of host and virus evolution. *J. Gen. Virol.* 103. doi: 10.1099/jgv.0.001749.
- 24

- 1 Tarlinton RE et al. 2021. Transcriptomic and genomic variants between koala populations reveals
2 underlying genetic components to disorders in a bottlenecked population. *Conserv. Genet.*
3 22:329–340. doi: 10.1007/s10592-021-01340-7.
- 4
- 5 Wick RR, Judd LM, Gorrie CL, Holt KE. 2017. Completing bacterial genome assemblies with
6 multiplex MinION sequencing. *Microb. Genomics.* 3. doi: 10.1099/mgen.0.000132.
- 7
- 8 Zheng J, Wei Y, Han G-Z. 2022. The diversity and evolution of retroviruses: Perspectives from
9 viral “fossils”. *Viol. Sin.* 37:11–18. doi: 10.1016/j.virs.2022.01.019.
- 10
- 11 Zimin AV, Salzberg SL. 2020. The genome polishing tool POLCA makes fast and accurate
12 corrections in genome assemblies Ouzounis, CA, editor. *PLOS Comput. Biol.* 16:e1007981. doi:
13 10.1371/journal.pcbi.1007981.
- 14
- 15 Zimin AV, Salzberg SL. 2022. The SAMBA tool uses long reads to improve the contiguity of
16 genome assemblies Shao, M, editor. *PLOS Comput. Biol.* 18:e1009860. doi:
17 10.1371/journal.pcbi.1009860.
- 18