1 **A draft genome assembly of the eastern banjo frog *Limnodynastes dumerilii***

2 ***dumerilii* (Anura: Limnodynastidae)**

3 Qiye Li[1,2], Qunfei Guo[1,3], Yang Zhou[1], Huishuang Tan[1,4], Terry Bertozzi[5,6], Yuanzhen Zhu[1,7],

4 Ji Li[2,8], Stephen Donnellan[5], Guojie Zhang[2,8,9,10*]

5

6 [1] BGI-Shenzhen, Shenzhen 518083, China

7 [2] State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology,

8 Chinese Academy of Sciences, Kunming 650223, China

9 [3] College of Life Science and Technology, Huazhong University of Science and Technology,

10 Wuhan 430074, China

11 [4] Center for Informational Biology, University of Electronic Science and Technology of China,

12 Chengdu 611731, China

13 [5] South Australian Museum, North Terrace, Adelaide 5000, Australia

14 [6] School of Biological Sciences, University of Adelaide, North Terrace, Adelaide 5005,

15 Australia

16 [7] School of Basic Medicine, Qingdao University, Qingdao 266071, China

17 [8] China National Genebank, BGI-Shenzhen, Shenzhen 518120, China

18 [9] Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences,

19 650223, Kunming, China

20 [10] Section for Ecology and Evolution, Department of Biology, University of Copenhagen, DK-

21 2100 Copenhagen, Denmark

22 [*] Correspondence: guojie.zhang@bio.ku.dk (G.Z.).

23

24 e-mail addresses for all authors:

25 Qiye Li <liqiye@genomics.cn>, Qunfei Guo <guoqunfei@genomics.cn>, Yang Zhou

26 <zhouyang@genomics.cn>, Huishuang Tan < tanhuishuang@genomics.cn>, Terry Bertozzi

27 <Terry.Bertozzi@samuseum.sa.gov.au>, Yuanzhen Zhu <zhuyuanzhen@genomics.cn>, Ji Li

28 <liji1@genomics.cn>, Stephen Donnellan <Steve.Donnellan@samuseum.sa.gov.au> and

29 Guojie Zhang <guojie.zhang@bio.ku.dk>

30

31 ORCIDs:

32 Qiye Li: 0000-0002-5993-0312; Yang Zhou: 0000-0003-1247-5049; Terry Bertozzi: 0000-0001-6665-3395;

33 Stephen Donnellan: 0000-0002-5448-3226; Guojie Zhang: 0000-0001-6860-1521

## Abstract

Amphibian genomes are usually challenging to assemble due to large genome size and high repeat content. The Limnodynastidae is a family of frogs native to Australia, Tasmania and New Guinea. As an anuran lineage that successfully diversified on the Australian continent, it represents an important lineage in the amphibian tree of life but lacks reference genomes. Here we sequenced and annotated the genome of the eastern banjo frog *Limnodynastes dumerilii dumerilii* to fill this gap. The total length of the genome assembly is 2.38 Gb with a scaffold N50 of 285.9 kb. We identified 1.21 Gb of non-redundant sequences as repetitive elements and annotated 24,548 protein-coding genes in the assembly. BUSCO assessment indicated that more than 94% of the expected vertebrate genes were present in the genome assembly and the gene set. We anticipate that this annotated genome assembly will advance the future study of anuran phylogeny and amphibian genome evolution.

## Introduction

46

47 The recent powerful advances in genome sequencing technology have allowed efficient

48 decoding of the genomes of many species [1, 2]. So far, genome sequences are available

49 publicly for more than one thousand species sampled across the animal branch of the tree of

50 life. These genomic resources have provided vastly improved perspectives on our knowledge

51 of the origin and evolutionary history of metazoans [3, 4], facilitated advances in agriculture

52 [5], enhanced approaches for conservation of endangered species [6], and uncovered the

53 genomic changes underlying the evolutionary successes of some clades such as birds [7] and

54 insects [8]. However, amphibian genomes are still challenging to assemble due to their large

55 genome sizes, high repeat content and sometimes high heterozygosity if specimens are

56 collected from wild populations [9]. This also accounts for the scarcity of reference genomes

57 for Anura (frogs and toads) — the most species-rich order of amphibians including many

58 important models for developmental biology and environmental monitoring [10]. Specifically,

59 despite the existence of more than 7,000 living species of Anura [11], only 10 species have

60 their genomes sequenced and annotated to date [12-21], which cover only 8 out of the 54 anuran

61 families. Moreover, genomes of Neobatrachia, which contains more than 95% of the anuran

62 species [11], are particularly under-represented. Only 5 of the 10 publicly available anuran

63 genomes belong to Neobatrachia [22]. This deficiency of neobatrachian genomes would

64 undoubtedly restrict the study of the genetic basis underlying the great diversification of this

65 amphibian lineage, and our understanding of the adaptive genomic changes that facilitate the

66 aquatic to terrestrial transition of vertebrates and the numerous unique reproductive modes

67 found in this clade.

68 As a candidate species proposed for genomic analysis by the Genome 10K (G10K) initiative

69 [9], we sequenced and annotated the genome of the Australian banjo frog *Limnodynastes*

70 *dumerilii* (also called the pobblebonk; NCBI:txid104065) to serve as a representative species

71 of the neobatrachian family Limnodynastidae. This burrowing frog is endemic to Australia and

72 named after its distinctive "bonk" call, which is likened to a banjo string being plucked. It

73 mainly occurs along the southeast coast of Australia, from the coast of New South Wales,

74 throughout Victoria and into the southwest corner of South Australia and Tasmania [23]. Five

75 subspecies of *L. dumerilii* are recognized, including *Limnodynastes dumerilii dumerilii*, *L.*

76 *dumerilii grayi*, *L. dumerilii fryi*, *L. dumerilii insularis* and *L. dumerilii variegata* [24]. The

77 subspecies chosen for sequencing is the eastern banjo frog *L. dumerilii dumerilii*

78 (NCBI:txid104066), as it is the most widespread among the five subspecies and forms hybrid

79  zones with a number of the other subspecies [23]. We believe that the release of genomic

80  resources from this neobatrachian frog will benefit the future studies of phylogenomics and

81  comparative genomics of anurans, and also facilitate other research related to the evolutionary

82  biology of *Limnodynastes*.

83

## Methods

85  **Sample collection, library construction and sequencing**

86  Genomic DNA was extracted from the liver of an adult female *Limnodynastes dumerilii*

87  *dumerilii* (Fig. 1) using the Gentra Puregene Tissue Kit (QIAGEN, Hilden, Germany)

88  according to manufacturer's instructions with the following exceptions: following the DNA

89  precipitation step, DNA was spooled onto a glass rod, washed twice in 70% ethanol and dried

90  before dissolving in 100 ul of the recommended elution buffer [25]. The specimen was

91  originally caught in River Torrens, Adelaide, South Australia, Australia, and is archived in the

92  South Australian Museum (registration number: SAMAR66870).

93  A total of 211 Gb of sequence was generated from four short-insert libraries (170 bp × 1, 250

94  bp × 1, 500 bp × 1, and 800 bp × 1), and 185 Gb of sequence from ten mate-paired libraries (2

95  kb × 3, 5 kb × 3, 10 kb × 2, and 20 kb × 2). All the 14 libraries were subject to paired-end

96  sequencing on the HiSeq 2000 platform following the manufacturer's instructions (Illumina,

97  San Diego, CA, USA), using PE100 or PE150 chemistry for the short-insert libraries and PE49

98  for the mate-paired libraries [26] (Table 1). Low-quality reads, adapter-contaminated reads,

99  and duplicated reads arising from polymerase chain reaction (PCR) amplification during

100 library construction were removed by SOAPnuke (v1.5.3, RRID:SCR_015025) [27] prior to

101 downstream analyses. This yielded a total of 180 Gb of clean sequence for genome assembly,

102 which represents 71 times coverage of the estimated haploid genome size of *L. d. dumerilii* in

103 terms of sequence depth, and 1,198 times in terms of physical depth (Table 1).

104

105 **Genome size estimation and genome assembly**

106 To obtain a robust estimation of the genome size of *L. d. dumerilii*, we conducted *k*-mer

107 analysis with all of the clean sequence (130 Gb) from the four short-insert libraries using a

108 range of *k* values (17, 19, 21, 23, 25, 27, 29 and 31). The *k*-mer frequencies were counted by

109 Jellyfish (v2.2.6) [28] with the -*C* setting. The genome size of *L. d. dumerilii* was estimated to

110 be around 2.54 Gb (Table 2), which was calculated as the number of effective *k*-mers (i.e. total

111 *k*-mers – erroneous *k*-mers) divided by the homozygous peak depth following Cai *et al* [29]. It

112  is noteworthy that, the presence of a distinct heterozygous peak, which displayed half of the
113  depth of the homozygous peak in the *k*-mer frequency distribution, suggests that the diploid
114  genome of this wild-caught individual has a high level of heterozygosity (Fig. 2). The rate of
115  heterozygosity was estimated to be around 1.17% by GenomeScope (v1.0.0,
116  RRID:SCR_017014) [30] (Table 2).

117  We then employed Platanus (v1.2.1, RRID:SCR_015531) [31] to assemble the genome of *L.*
118  *d. dumerilii*. Briefly, all the clean sequence from the four short-insert libraries were first
119  assembled into contigs using *platanus assemble* with parameters *-t 20 -k 29 -u 0.2 -d 0.6 -m*
120  *150*. Then paired-end reads from the four short-insert and ten mate-paired libraries were used
121  to connect contigs into scaffolds by *platanus scaffold* with parameters *-t 20 -u 0.2 -l 3* and the
122  insert size information of each library. Finally, *platanus gap_close* was employed to close
123  intra-scaffold gaps using the paired-end reads from the four short-insert libraries with default
124  settings. This Platanus assembly was further improved by Kgf (version 1.16) [9] followed by
125  GapCloser (v1.10.1, RRID:SCR_015026) [9] for gap filling with the clean reads from the four
126  short-insert libraries.

127

128  **Repetitive element annotation**
129  Both homology-based and *de novo* predictions were employed to identify repetitive elements
130  in the *L. d. dumerilii* genome assembly [32]. For homology-based prediction, known repetitive
131  elements were identified by aligning the *L. d. dumerilii* genome sequences against the Repbase-
132  derived RepeatMasker libraries using RepeatMasker (v4.1.0, RRID:SCR_012954; setting *-*
133  *nolow -norna -no_is*) [33], and against the transposable element protein database using
134  RepeatProteinMask (an application within the RepeatMasker package; setting *-noLowSimple -*
135  *pvalue 0.0001 -engine ncbi*). For *de novo* prediction, RepeatModeler (v2.0,
136  RRID:SCR_015027) [34] was first executed on the *L. d. dumerilii* assembly to build a *de novo*
137  repeat library for this species. Then RepeatMasker was employed to align the *L. d. dumerilii*
138  genome sequences against the *de novo* library for repetitive element identification. Tandem
139  repeats in the *L. d. dumerilii* genome assembly were identified by Tandem Repeats Finder
140  (v4.09) [35] with parameters *Match=2 Mismatch=7 Delta=7 PM=80 PI=10 Minscore=50*
141  *MaxPeriod=2000*.

142

143  **Protein-coding gene annotation**

144    Similar to repetitive element annotation, both homology-based and *de novo* predictions were

145    employed to build gene models for the *L. d. dumerilii* genome assembly [36]. For homology-

146    based prediction, protein sequences from diverse vertebrate species, including *Danio rerio,*

147    *Xenopus tropicalis*, *Xenopus laevis*, *Nanorana parkeri*, *Microcaecilia unicolor*, *Rhinatrema*

148    *bivittatum*, *Anolis carolinensis*, *Gallus gallus* and *Homo sapiens*, were first aligned to the *L.*

149    *d. dumerilii* genome assembly using TBLASTN (blast-2.2.26, RRID:SCR_011822) [37] with

150    parameters *-F F -e 1e-5*. Then the genomic sequences of the candidate loci together with 5

151    kb flanking sequences were extracted for exon-intron structure determination, by aligning

152    the homologous proteins to these extracted genomic sequences using GeneWise (wise-2.2.0,

153    RRID:SCR_015054) [38]. For *de novo* prediction, we randomly picked 1,000 homology-

154    derived gene models of *L. d. dumerilii* with complete open reading frames (ORFs) and

155    reciprocal aligning rates exceeding 90% against the *X. tropicalis* proteins to train

156    AUGUSTUS (v3.3.1, RRID:SCR_008417) [39]. The obtained gene parameters were then

157    used by AUGUSTUS to predict protein-coding genes on the repeat-masked *L. d. dumerilii*

158    genome assembly. Finally, gene models derived from the above two methods were

159    combined into a non-redundant gene set using a similar strategy to Xiong *et al*. (2016) [40].

160    Genes showing BLASTP (blast-2.2.26, RRID:SCR_001010; parameters *-F F -e 1e-5*) hits

161    to transposon proteins in the UniProtKB/Swiss-Prot database (v2019_11), or with more than

162    70% of their coding regions overlapping repetitive sequences, were removed from the

163    combined gene set.

164

## Results and Discussion

### Assembly and annotation of the *L. d. dumerilii* genome

167    We assembled the nuclear genome of a female eastern banjo frog *L. d. dumerilii* (Fig. 1) with

168    ~180 Gb (71X) clean Hiseq data from four short-insert libraries (170 bp × 1, 250 bp × 1, 500

169    bp × 1, and 800 bp × 1) and ten mate-paired libraries (2 kb × 3, 5 kb × 3, 10 kb × 2, and 20 kb

170    × 2) (Table 1). The final genome assembly comprised 520,896 sequences with contig and

171    scaffold N50s of 10.2 kb and 286.0 kb, respectively, and a total length of 2.38 Gb, which is

172    close to the estimated genome size of 2.54 Gb by *k*-mer analysis (Table 2 and Fig. 2). There

173    are 242 Mb of regions present as unclosed gaps (Ns), accounting for 10.2% of the assembly.

174    The GC content of the *L. d. dumerilii* assembly excluding gaps was estimated to be 41.0%. The

175    combination of homology-based and *de novo* prediction methods masked 1.21 Gb of non-

176    redundant sequences as repetitive elements, accounting for 56.4 % of the *L. d. dumerilii*

177 genome assembly excluding gaps (Table 3). We also obtained 24,548 protein-coding genes

178 in the genome assembly, of which 67% had complete ORF. Functional annotation by

179 searching the *L. d. dumerilii* proteins against public databases of UniProtKB/Swiss-Prot

180 (v2019_11, RRID:SCR_004426) [41], NCBI nr (v20191030), and KEGG (v93.0,

181 RRID:SCR_012773) [42] with BLASTP (blast-2.2.26; parameters *-F F -e 1e-5*) successfully

182 annotated almost all of the *L. d. dumerilii* gene loci (Table 4).

183

184 **Data validation and quality control**

185 Two strategies were employed to estimate the completeness of the *L. d. dumerilii* genome

186 assembly. First, all the clean reads from the short-insert libraries were aligned to the genome

187 assembly using BWA-MEM (BWA, version 0.7.16, RRID:SCR_010910) with default

188 parameters [43]. We observed that 99.6 % of reads could be mapped back to the assembled

189 genome and 85.6 % of the inputted reads were mapped in proper pairs as accessed by samtools

190 flagstat (SAMtools v1.7, RRID:SCR_002105), suggesting that most sequences of the *L. d.*

191 *dumerilii* genome were present in the current assembly. Secondly, we assessed the *L. d.*

192 *dumerilii* assembly with Benchmarking Universal Single-Copy Orthologs (BUSCO; v3.0.2,

193 RRID:SCR_015008), a software package that can quantitatively measure genome assembly

194 completeness based on evolutionarily informed expectations of gene content [44], and found

195 that up to 94.7 % of the 2,586 expected vertebrate genes were present in the *L. d. dumerilii*

196 assembly. Furthermore, 85.5% and 84.5 % of the expected genes were identified as complete

197 and single-copy genes, respectively. This BUSCO assessment further highlighted the

198 comprehensiveness of the current *L. d. dumerilii* genome assembly in terms of gene space.

199 We then evaluated the completeness of the *L. d. dumerilii* protein-coding gene set with BUSCO

200 (v3.0.2) and DOGMA (v3.0, RRID:SCR_015060) [45], a program that measures the

201 completeness of a given transcriptome or proteome based on a core set of conserved domain

202 arrangements (CDAs). BUSCO analysis showed that 97.1 % of the expected vertebrate genes

203 were present in the *L. d. dumerilii* protein-coding gene set with 88.5 % and 84.5% identified

204 as complete and single-copy genes, respectively, close to that estimated for the genome

205 assembly. Meanwhile, DOGMA analysis based on PfamScan Annotations (PfamScan v1.5;

206 Pfam v32.0, RRID:SCR_015060) [46] and the eukaryotic core set identified 95.4 % of the

207 expected CDAs in the annotated gene set. These results demonstrated the high completeness

208 of the *L. d. dumerilii* protein-coding gene set.

209

**Re-use potential**

Here, we report a draft genome assembly of the eastern banjo frog *L. d. dumerilii*. It represents the first genome assembly from the family Limnodynastidae (Anura: Neobatrachia). Although the continuity of the assembly in terms of contig and scaffold N50s is modest, probably due to the high repeat content (56%) and heterozygosity (1.17%), the completeness of this draft assembly is demonstrated to be high according to read mapping and BUSCO assessment. Thus, it is suitable for phylogenomics and comparative genomics analyses with other available anuran genomes or phylogenomic datasets. In particular, the high-quality protein-coding gene set derived from the genome assembly will be useful for deducing orthologous relationships across anuran species or reconstructing the ancestral gene content of anurans. Due to evolutionary importance of *Limnodynastes* frogs in Australia, the genomic resources released in this study will also support further research on the biogeography of speciation, evolution of male advertisement calls, hybrid zone dynamics, and conservation of *Limnodynastes* frogs.

## Availability of supporting data

The raw sequencing reads are deposited in NCBI under the BioProject accession PRJNA597531 and are also deposited in the CNGB Nucleotide Sequence Archive (CNSA) with accession number CNP0000818. Genome assembly, protein-coding gene and repeat annotations are deposited in the *GigaScience* GigaDB [47] and NCBI under accession number GCA_011038615.1.

## List of abbreviations

BUSCO: Benchmarking Universal Single-Copy Orthologs; G10K: Genome 10K; NCBI: National Center for Biotechnology Information; PCR: Polymerase Chain Reaction; ORF: Open Reading Frame; KEGG: Kyoto Encyclopedia of Genes and Genomes; DOGMA: DOmain-based General Measure for transcriptome and proteome quality Assessment; CDA: Conserved Domain Arrangement; CNGB: China National GeneBank; CNSA: CNGB Sequence Archive.

## Funding

## Competing interests

246    The authors declare that they have no competing interests.

247

## Author contributions

249    G.Z. and Q.L. conceived and supervised the study; T.B. and S.D. prepared the DNA samples;

250    Y.Z. and Q.G. performed *k*-mer analysis and genome assembly; Q.G. and J.L. conducted

251    assessment of assembly quality; H.T. performed protein-coding gene annotation; Y.Z.

252    performed repeat annotation; G.Z. and S.D. contributed reagents/materials/analysis tools; Q.L.

253    wrote the manuscript with the inputs from all authors. All authors read and approved the final

254    manuscript.

255

## References

257    1.    Goodwin S, McPherson JD and McCombie WR. Coming of age: ten years of next-
258          generation sequencing technologies. Nature reviews Genetics. 2016;17 6:333-51.
259          doi:10.1038/nrg.2016.49.

260    2.    van Dijk EL, Jaszczyszyn Y, Naquin D and Thermes C. The Third Revolution in
261          Sequencing Technology. Trends in genetics : TIG. 2018;34 9:666-81.
262          doi:10.1016/j.tig.2018.05.008.

263    3.    Sebe-Pedros A, Degnan BM and Ruiz-Trillo I. The origin of Metazoa: a unicellular
264          perspective. Nature reviews Genetics. 2017;18 8:498-512. doi:10.1038/nrg.2017.21.

265    4.    Laumer CE, Fernandez R, Lemer S, Combosch D, Kocot KM, Riesgo A, et al.
266          Revisiting metazoan phylogeny with genomic sampling of all phyla. Proceedings
267          Biological sciences / The Royal Society. 2019;286 1906:20190831.
268          doi:10.1098/rspb.2019.0831.

269    5.    Beiki H, Eveland AL and Tuggle CK. Recent advances in plant and animal genomics
270          are taking agriculture to new heights. Genome Biol. 2018;19 1:48.
271          doi:10.1186/s13059-018-1427-z.

272    6.    Supple MA and Shapiro B. Conservation of biodiversity in the genomics era. Genome
273          Biol. 2018;19 1:131. doi:10.1186/s13059-018-1520-3.

274    7.    Zhang G, Li C, Li Q, Li B, Larkin DM, Lee C, et al. Comparative genomics reveals
275          insights into avian genome evolution and adaptation. Science. 2014;346 6215:1311-
276          20. doi:10.1126/science.1251385.

277  8.   Thomas GWC, Dohmen E, Hughes DST, Murali SC, Poelchau M, Glastad K, et al.
278       Gene content evolution in the arthropods. Genome Biol. 2020;21 1:15.
279       doi:10.1186/s13059-019-1925-7.

280  9.   Koepfli KP, Paten B, Genome KCoS and O'Brien SJ. The Genome 10K Project: a
281       way forward. Annu Rev Anim Biosci. 2015;3:57-111. doi:10.1146/annurev-animal-
282       090414-014900.

283  10.  Carroll R. The rise of amphibians: 365 million years of evolution. Johns Hopkins
284       University Press. 2009.

285  11.  AmphibiaWeb. <https://amphibiaweb.org> University of California, Berkeley, CA,
286       USA. (Accessed 18 Feb 2020).

287  12.  Li J, Yu H, Wang W, Fu C, Zhang W, Han F, et al. Genomic and transcriptomic
288       insights into molecular basis of sexually dimorphic nuptial spines in *Leptobrachium*
289       *leishanense*. Nat Commun. 2019;10 1:5551. doi:10.1038/s41467-019-13531-5.

290  13.  Li Y, Ren Y, Zhang D, Jiang H, Wang Z, Li X, et al. Chromosome-level assembly of
291       the mustache toad genome using third-generation DNA sequencing and Hi-C analysis.
292       Gigascience. 2019;8 9 doi:10.1093/gigascience/giz114.

293  14.  Seidl F, Levis NA, Schell R, Pfennig DW, Pfennig KS and Ehrenreich IM. Genome
294       of *Spea multiplicata*, a Rapidly Developing, Phenotypically Plastic, and Desert-
295       Adapted Spadefoot Toad. G3 (Bethesda). 2019;9 12:3909-19.
296       doi:10.1534/g3.119.400705.

297  15.  Edwards RJ, Tuipulotu DE, Amos TG, O'Meally D, Richardson MF, Russell TL, et
298       al. Draft genome assembly of the invasive cane toad, *Rhinella marina*. Gigascience.
299       2018;7 9 doi:10.1093/gigascience/giy095.

300  16.  Rogers RL, Zhou L, Chu C, Marquez R, Corl A, Linderoth T, et al. Genomic
301       Takeover by Transposable Elements in the Strawberry Poison Frog. Mol Biol Evol.
302       2018;35 12:2913-27. doi:10.1093/molbev/msy185.

303  17.  Denton RD, Kudra RS, Malcom JW, Du Preez L and Malone JH. The African
304       Bullfrog (*Pyxicephalus adspersus*) genome unites the two ancestral ingredients for
305       making vertebrate sex chromosomes. bioRxiv. 2018:329847.

306  18.  Hammond SA, Warren RL, Vandervalk BP, Kucuk E, Khan H, Gibb EA, et al. The
307       North American bullfrog draft genome provides insight into hormonal regulation of
308       long noncoding RNA. Nat Commun. 2017;8 1:1433. doi:10.1038/s41467-017-01316-
309       7.

310  19.  Session AM, Uno Y, Kwon T, Chapman JA, Toyoda A, Takahashi S, et al. Genome
311       evolution in the allotetraploid frog *Xenopus laevis*. Nature. 2016;538 7625:336-43.
312       doi:10.1038/nature19840.

313  20.  Sun YB, Xiong ZJ, Xiang XY, Liu SP, Zhou WW, Tu XL, et al. Whole-genome
314       sequence of the Tibetan frog *Nanorana parkeri* and the comparative evolution of
315       tetrapod genomes. Proceedings of the National Academy of Sciences of the United
316       States of America. 2015;112 11:E1257-62. doi:10.1073/pnas.1501764112.

317   21.   Hellsten U, Harland RM, Gilchrist MJ, Hendrix D, Jurka J, Kapitonov V, et al. The
318         genome of the Western clawed frog *Xenopus tropicalis*. Science. 2010;328 5978:633-
319         6. doi:10.1126/science.1183670.

320   22.   The NCBI Assembly database: https://www.ncbi.nlm.nih.gov/assembly/?term=Anura;
321         access on February 18, 2020.

322   23.   Martin A. Studies in Australian amphibia III. The limnodynastes dorslis complex
323         (Anura: Leptodactylidae). Australian Journal of Zoology. 1972;20 2:165-211.

324   24.   Schauble CS, Moritz C and Slade RW. A molecular phylogeny for the frog genus
325         Limnodynastes (Anura: myobatrachidae). Mol Phylogenet Evol. 2000;16 3:379-91.
326         doi:10.1006/mpev.2000.0803.

327   25.   Bertozzi T and Donnellan S. DNA extraction protocol for the eastern banjo frog using
328         the Gentra Puregene Tissue Kit. protocols.io 2020;
329         doi:dx.doi.org/10.17504/protocols.io.bcy6ixze.

330   26.   Li Q, Guo Q, Zhou Y, Tan H, Bertozzi T, Zhu Y, et al. Construction and sequencing
331         of DNA libraries on Hiseq 2000 platform for the eastern banjo frog. protocols.io
332         2020;  doi:dx.doi.org/10.17504/protocols.io.bc22iyge.

333   27.   Chen Y, Chen Y, Shi C, Huang Z, Zhang Y, Li S, et al. SOAPnuke: a MapReduce
334         acceleration-supported software for integrated quality control and preprocessing of
335         high-throughput sequencing data. Gigascience. 2018;7 1:1-6.
336         doi:10.1093/gigascience/gix120.

337   28.   Marcais G and Kingsford C. A fast, lock-free approach for efficient parallel counting
338         of occurrences of k-mers. Bioinformatics. 2011;27 6:764-70.
339         doi:10.1093/bioinformatics/btr011.

340   29.   Cai H, Li Q, Fang X, Li J, Curtis NE, Altenburger A, et al. A draft genome assembly
341         of the solar-powered sea slug *Elysia chlorotica*. Sci Data. 2019;6:190022.
342         doi:10.1038/sdata.2019.22.

343   30.   Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, et al.
344         GenomeScope: fast reference-free genome profiling from short reads. Bioinformatics.
345         2017;33 14:2202-4. doi:10.1093/bioinformatics/btx153.

346   31.   Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, et al. Efficient
347         de novo assembly of highly heterozygous genomes from whole-genome shotgun short
348         reads. Genome Res. 2014;24 8:1384-95. doi:10.1101/gr.170720.113.

349   32.   Li Q, Guo Q, Zhou Y, Tan H, Bertozzi T, Zhu Y, et al. Repetitive element annotation
350         protocol for the eastern banjo frog. protocols.io 2020;
351         doi:dx.doi.org/10.17504/protocols.io.bc4niyve.

352   33.   Smit AF, Hubley R and Green P. Available fom http://www.repeatmasker.org. 20
353         September 2019 date last accessed.

354   34.   Smit A and Hubley R. Available fom http://www.repeatmasker.org/RepeatModeler/.
355         20 September 2019 date last accessed.

356  35.  Benson G. Tandem repeats finder: a program to analyze DNA sequences. Nucleic
357       acids research. 1999;27 2:573-80.

358  36.  Li Q, Guo Q, Zhou Y, Tan H, Bertozzi T, Zhu Y, et al. Protein-coding gene
359       annotation protocol for the eastern banjo frog. protocols.io 2020;
360       doi:dx.doi.org/10.17504/protocols.io.bc38iyrw.

361  37.  Altschul SF, Gish W, Miller W, Myers EW and Lipman DJ. Basic local alignment
362       search tool. Journal of molecular biology. 1990;215 3:403-10. doi:10.1016/S0022-
363       2836(05)80360-2.

364  38.  Birney E, Clamp M and Durbin R. GeneWise and Genomewise. Genome Res.
365       2004;14 5:988-95. doi:10.1101/gr.1865504.

366  39.  Stanke M, Diekhans M, Baertsch R and Haussler D. Using native and syntenically
367       mapped cDNA alignments to improve de novo gene finding. Bioinformatics. 2008;24
368       5:637-44. doi:10.1093/bioinformatics/btn013.

369  40.  Xiong Z, Li F, Li Q, Zhou L, Gamble T, Zheng J, et al. Draft genome of the leopard
370       gecko, *Eublepharis macularius*. Gigascience. 2016;5 1:47. doi:10.1186/s13742-016-
371       0151-4.

372  41.  UniProt Consortium T. UniProt: the universal protein knowledgebase. Nucleic acids
373       research. 2018;46 5:2699. doi:10.1093/nar/gky092.

374  42.  Kanehisa M and Goto S. KEGG: kyoto encyclopedia of genes and genomes. Nucleic
375       acids research. 2000;28 1:27-30. doi:10.1093/nar/28.1.27.

376  43.  Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-
377       MEM. arXiv preprint arXiv:13033997. 2013.

378  44.  Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO:
379       assessing genome assembly and annotation completeness with single-copy orthologs.
380       Bioinformatics. 2015;31 19:3210-2. doi:10.1093/bioinformatics/btv351.

381  45.  Dohmen E, Kremer LP, Bornberg-Bauer E and Kemena C. DOGMA: domain-based
382       transcriptome and proteome quality assessment. Bioinformatics. 2016;32 17:2577-81.
383       doi:10.1093/bioinformatics/btw231.

384  46.  Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, et al. Pfam: the
385       protein families database. Nucleic acids research. 2014;42 Database issue:D222-30.
386       doi:10.1093/nar/gkt1223.

387  47.  Li Q, Guo Q, Zhou Y, Tan H, Bertozzi T, Zhu Y, et al. Genomic data from the
388       Eastern banjo frog *Limnodynastes dumerilii dumerilii* (Anura: Limnodynastidae).
389       GigaScience Database. 2020; doi:http://dx.doi.org/10.5524/100717.
390

## Figures



**Figure 1. Photograph of an adult *Limnodynastes dumerilii dumerilii* from the Adelaide region (image from Stephen Mahony).**
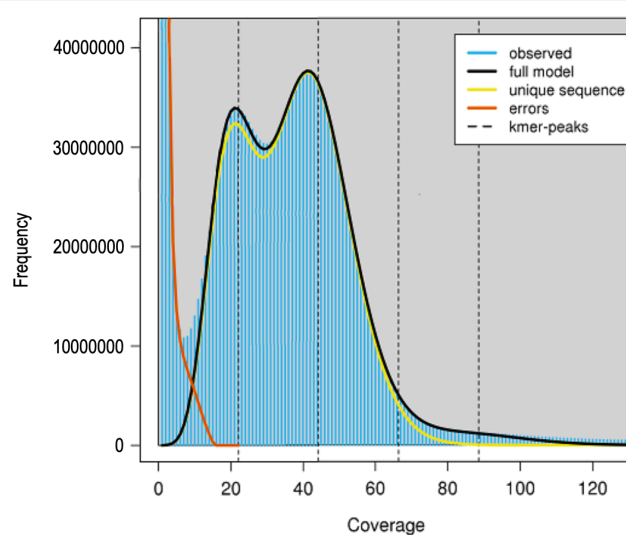


**Figure 2. A 21-mer frequency distribution of the *L. d. dumerilii* genome data.** The first peak at coverage 21X corresponds to the heterozygous peak. The second peak at coverage 42X corresponds to the homozygous peak.

## Tables

**Table 1. Statistics of DNA reads produced for the *L. d. dumerilii* genome.**

| Insert size (bp) | No. of Libraries | Read length (bp) | Raw data | | | Clean data | | |
|---|---|---|---|---|---|---|---|---|
| | | | Total bases (Gb) | Sequence depth (X) | Physical depth (X) | Total bases (Gb) | Sequence depth (X) | Physical depth (X) |
| 170 | 1 | 100 | 43.45 | 17.11 | 14.54 | 36.20 | 14.25 | 12.75 |
| 250 | 1 | 150 | 67.56 | 26.60 | 22.17 | 46.22 | 18.20 | 15.69 |
| 500 | 1 | 150 | 61.47 | 24.20 | 40.33 | 31.02 | 12.21 | 24.43 |
| 800 | 1 | 150 | 38.34 | 15.09 | 40.25 | 16.73 | 6.59 | 21.08 |
| 2,000 | 3 | 49 | 59.90 | 23.58 | 481.28 | 29.85 | 11.75 | 301.33 |
| 5,000 | 3 | 49 | 52.11 | 20.52 | 1046.72 | 11.88 | 4.68 | 299.82 |
| 10,000 | 2 | 49 | 36.81 | 14.49 | 1478.79 | 5.27 | 2.07 | 266.00 |
| 20,000 | 2 | 49 | 36.02 | 14.18 | 2894.10 | 2.54 | 1.00 | 256.41 |
| Total | 14 | | 395.66 | 155.77 | 6018.18 | 179.71 | 70.75 | 1197.50 |

Note: Coverage calculation was based on the estimated haploid genome size of 2.54 Gb according to $k$-mer analysis. Sequence coverage is the average number of times a base is read, while physical coverage is the average number of times a base is spanned by sequenced fragments.

**Table 2. Estimation of genome size and heterozygosity of *L. d. dumerilii* by *k*-mer analysis.**

| $k$ | Total number of $k$-mers | Minimum coverage (X) | Number of erroneous $k$-mers | Homozygous peak | Estimated genome size (Gb) | Estimated heterozygosity (%) |
|---|---|---|---|---|---|---|
| 17 | 112,401,363,509 | 9 | 1,418,748,938 | 45 | 2.47 | 1.10 |
| 19 | 110,136,516,133 | 8 | 2,588,664,358 | 43 | 2.50 | 1.23 |
| 21 | 107,871,808,889 | 7 | 3,023,604,282 | 42 | 2.50 | 1.24 |
| 23 | 105,607,392,491 | 7 | 3,286,834,146 | 40 | 2.56 | 1.22 |
| 25 | 103,343,108,760 | 7 | 3,501,481,190 | 39 | 2.56 | 1.19 |
| 27 | 101,078,882,097 | 7 | 3,689,197,189 | 38 | 2.56 | 1.16 |
| 29 | 98,815,880,190 | 6 | 3,839,002,752 | 37 | 2.57 | 1.14 |
| 31 | 96,552,885,503 | 6 | 3,986,778,359 | 36 | 2.57 | 1.11 |

Note: $k$-mer frequency distributions were generated by Jellyfish (v2.2.6) using 130 Gb clean sequences as input. Minimum coverage was the coverage depth value of the first trough in $k$-mer frequency distribution. $k$-mers with coverage depth less than the minimum coverage were regarded as erroneous $k$-mers. Estimated genome size was calculated as (Total number of $k$-mers – Number of erroneous $k$-mers) / Homozygous peak.

**Table 3. Statistics of repetitive sequences identified in the *L. d. dumerilii* genome.**

| Category | Total repeat length (bp) | % of assembly |
|---|---|---|
| DNA | 155,988,597 | 7.30% |
| LINE | 242,754,702 | 11.36% |
| SINE | 11,761,904 | 0.55% |
| LTR | 97,615,246 | 4.57% |
| Tandem repeats | 178,355,571 | 8.35% |
| Unknown | 704,263,255 | 32.96% |
| Combined | 1,205,873,056 | 56.43% |

Note: DNA: DNA transposon; LINE: long interspersed nuclear element; SINE: short interspersed nuclear elements; LTR: long terminal repeat.

**Table 4. Summary of protein-coding genes annotated in the *L. d. dumerilii* genome.**

| Characteristics of protein-coding genes | |
|---|---|
| Total number of protein-coding genes | 24,548 |
| Gene space (exon + intron; Mb) | 634.6 (26.7 % of assembly) |
| Mean gene size (bp) | 25,851 |
| Mean CDS length (bp) | 1,552 |
| Exon space (Mb) | 38.1 (1.6 % of assembly) |
| Mean exon number per gene | 8.6 |
| Mean exon length (bp) | 181 |
| Mean intron length (bp) | 3,217 |
| Functional annotation by searching public databases | |
| % of proteins with hits in UniProtKB/Swiss-Prot | 95.8 |
| % of proteins with hits in NCBI nr database | 99.6 |
| % of proteins with KO assigned by KEGG | 71.3 |
| % of proteins with functional annotation (combined) | 99.9 |