1    **Unlocking inaccessible historical genomes preserved in formalin**

2

3    Erin E. Hahn[1], Marina R. Alexander[1], Alicia Grealy[1], Jiri Stiller[2], Donald M. Gardiner[2] and

4    Clare E. Holleley[1*]

5

6    Institutional Addresses:

7    [1] National Research Collections Australia, Commonwealth Scientific Industrial Research

8    Organisation, Canberra, ACT 2601, Australia

9    [2] Agriculture and Food, Commonwealth Scientific Industrial Research Organisation, St Lucia,

10    Queensland 4067, Australia

11

12    Email addresses:

13    EEH – erin.hahn@csiro.au; MRA – marina.alexander@csiro.au; AG –

14    Alicia.grealy@csiro.au; JS – jiri.stiller@csiro.au; DMG – Donald.gardiner@csiro.au; CEH –

15    clare.holleley@csiro.au*

16    * Corresponding author

17 **Abstract**

18 **Background**

19 Museum specimens represent an unparalleled record of historical genomic data. However, the

20 wide-spread practice of formalin preservation has thus far impeded genomic analysis of a large

21 proportion of specimens. Limited DNA sequencing from formalin-preserved specimens has

22 yielded low genomic coverage with unpredictable success. We set out to refine sample

23 processing methods and to identify specimen characteristics predictive of sequencing success.

24 With a set of taxonomically diverse specimens collected between 1936 and 2015 and ranging

25 in preservation quality, we compared the efficacy of several end-to-end whole genome

26 sequencing workflows alongside a k-mer-based trimming-free read alignment approach to

27 maximize mapping of endogenous sequence.

28 **Results**

29 We recovered complete mitochondrial genomes and up to 3X nuclear genome coverage from

30 formalin-fixed tissues. Hot alkaline lysis coupled with phenol-chloroform extraction out-

31 performed proteinase K digestion in recovering DNA, while library preparation method had

32 little impact on sequencing success. The strongest predictor of DNA yield was overall

33 specimen condition, which additively interacts with preservation conditions to accelerate DNA

34 degradation.

35 **Conclusions**

36 We demonstrate a significant advance in capability beyond limited recovery of a small number

37 of loci via PCR or target-capture sequencing. To facilitate strategic selection of suitable

38 specimens for genomic sequencing, we present a decision-making framework that utilizes

39 independent and non-destructive assessment criteria. Sequencing of formalin-fixed specimens

40 will contribute to a greater understanding of temporal trends in genetic adaptation, including

41 those associated with a changing climate. Our work enhances the value of museum collections

42    worldwide by unlocking genomes of specimens that have been disregarded as a valid molecular

43    resource.

44

45    **Keywords**: DNA, formaldehyde, formalin-fixed, genome, hot alkali, museum, museomics,

46    preservation media

## Background

48 Natural history collections are a window into the recent past, offering a view of historical

49 biodiversity that is unparalleled in its detail. Collected over the last 250 years, voucher

50 specimens document a period of time over which humans have had a devastating impact on the

51 natural world (1). The comprehensive metadata associated with each specimen (collection date,

52 location, sex, weight, age, etc.), phenotypic data (e.g., color, size, gut contents) and genomic

53 data can be used to monitor ecosystem health and study the mechanisms driving adaptation,

54 evolution, speciation and extinction (2,3). The value of collections as sources of historical

55 genetic material has been recognized for the past 30 years, with numerous pathways emerging

56 to retrieve high-quality DNA from challenging archival vertebrate tissues such as skins (4),

57 feathers (5,6), eggshells (7,8) and toe pads (9).

58 DNA degradation associated with preservation method and aging has limited most genetic

59 studies of museum specimens to interrogation of relatively few loci via PCR amplification,

60 often targeting the high copy mitochondrial genome. For phylogenetic studies where a survey

61 of many-fold more loci improves understanding of species' evolutionary history (10–12),

62 genome-wide analyses are increasingly becoming common place. With demand for historical

63 genome-wide data on the rise, newly-developed target-capture approaches now facilitate

64 broader genomic survey from degraded museum specimens (13–15). In some cases, recovery

65 and assembly of whole historical genomes has been achieved (16,17), including, extinct from

66 species (e.g., the Tasmanian tiger (18)). While technological advances are enabling recovery

67 of genomic data from many museum specimens, genomic study of those preserved with 10%

68 formalin (3.4% w/v formaldehyde) has thus far been very limited.

69 Formalin-fixation, followed by storage in ethanol, is a common curatorial method used to

70 preserve soft tissue structure. Of the 1.9 million records of preserved chordates within the open-

71    access Atlas of Living Australia (ALA) specimen database (19), 33% are classified as "spirit-

72    preserved" (preserved in ethanol with or without prior formalin-fixation). A search for

73    "formalin" preparation within the ALA's chordate records indicates at least 4% of specimens

74    (N = 77,301) have been formalin-fixed. This is likely a severe underestimate because formalin-

75    fixation is not consistently recorded by all collections. Notably, for fish, reptiles and

76    amphibians, formalin-fixation has historically been the primary method used to preserve tissues

77    long-term while mammals and birds are commonly dry-preserved. Most collections now

78    archive frozen fresh tissue specifically as a genomic resource. However, prior to the 1980s,

79    spirit-preservation was the only method used to preserve soft tissue. Thus, spirit collections

80    offer the only opportunity to obtain genetic data from a large proportion of older specimens,

81    holotypes and some of the world's most biodiverse vertebrate taxonomic groups.

82    Genomic study of formalin-preserved museum specimens has lagged behind because DNA

83    extracted from such tissues is typically low-yield and highly fragmented. PCR amplification of

84    formalin-degraded DNA templates is generally restricted to few, short genomic loci, which

85    provide limited phylogenetic resolution (20). Formalin fixation presents further challenges by

86    inducing numerous molecular lesions, such as strand breaks, base misincorporation, and both

87    intra- and intermolecular cross-links (21–23). Formaldehyde damage to DNA templates can

88    result in sequencing artefacts that are difficult to differentiate from true genetic variants

89    (22,23). Because PCR amplification of damaged DNA is particularly prone to sequencing

90    artefacts, it is preferable to perform deep next-generation sequencing of amplicons (20) or to

91    avoid amplicon approaches altogether through whole genome sequencing (WGS) of degraded

92    templates (24). Coupled with library preparation methods optimized for low-input and

93    damaged DNA templates (25,26), high-throughput sequencing can generate enough coverage

94    to call genomic variants with high confidence (27). Thus, WGS and reduced representations of

95    genomes could provide a way to overcome the challenges associated with formalin damage

96    and accurately reconstruct historical genetic variation from formalin-preserved tissues.

97    Promisingly, WGS of formalin-fixed paraffin-embedded (FFPE) archival tissues has become

98    routine in clinical and medical contexts (28). However, museum specimens are often older,

99    exposed to higher concentrations of formaldehyde, incubated in the fixative for longer (29) and

100    in most cases have not been preserved in ideal conditions. Common museum practices, such

101    as failure to rinse specimens prior to permanent storage in ethanol, result in prolonged

102    formaldehyde exposure (30). Indeed, many specimens can be in contact with formaldehyde (or

103    its derivatives, such as formic acid) for the entirety of their tenure in a collection. Prolonged

104    formaldehyde exposure, especially under acidic conditions, is thought to result in more extreme

105    DNA degradation (20,31). The damage resulting from the preservation process compounds

106    with DNA damage due to natural decomposition, which can be extensive and often precedes

107    any obvious visual indicators of decomposition (32). Unfortunately, the time between death

108    and preservation (post-mortem interval) is highly variable and rarely recorded. In light of these

109    additional challenges, WGS methods used with FFPE tissues are relevant but not directly

110    transferable to formalin-fixed museum tissues.

111    Of the few genetic studies of formalin-fixed museum specimens, most have targeted nuclear

112    (33–37) and high copy mitochondrial (20,38,39) loci via PCR amplification due to the

113    difficulty and unpredictability of nuclear DNA extraction. There are few examples of broader-

114    scale genomic sequencing of formalin-fixed museum specimens and none have recovered

115    whole vertebrate genomes. Hot alkaline extraction followed by WGS of a single 30-year-old

116    formalin-preserved *Anolis* lizard yielded sufficient coverage to reconstruct the entire

117    mitochondrial genome (40). Using the same method, whole genomes were recovered for the

118    bioluminescent bacterial symbionts contained within light organs of formalin-preserved

119    cardinalfish (41). Using a proteinase K digestion method, sufficient gDNA was recovered for

6

120    capture and sequencing of ultra-conserved elements from formalin-preserved snakes (42).

121    Hybridization capture baits have also been used to recover the mitochondrial genome from a

122    120-year-old formalin-preserved Crimean green lizard (43). Highlighting the difficulty of

123    recovering gDNA from formalin-preserved specimens, numerous studies have reported failure

124    to extract and amplify gDNA from formalin-preserved museum tissues (20,44,45). In this

125    context, it is unfortunate yet wise to be hesitant to conduct destructive sampling of formalin-

126    preserved specimens for the purposes of costly WGS.

127    Recent reports of successful, albeit limited, genomic sequencing in formalin-preserved

128    specimens indicate WGS of higher quality specimens is possible. However, without a

129    framework to guide specimen selection, genomic work on formalin-preserved museum tissues

130    will remain infeasible. It is likely impossible to fully know the numerous and interdependent

131    factors driving sequencing success (e.g., age of the specimen (46,47), method of preservation

132    (48), post-mortem interval (32) and heat and light exposure during storage). However,

133    identification of metrics with which to pre-screen specimens for sequencing suitability will

134    improve yield of genomic data while reducing unnecessary destruction of specimens. With

135    screening criteria in hand, museum curators will be less reluctant to grant destructive sampling

136    (49) and researchers will be more inclined to include historical specimens in their analyses.

137    To facilitate informed-selection of formalin-preserved museum specimens for WGS, we set

138    out to further refine appropriate extraction and library preparation methods and to identify

139    specimen characteristics predictive of DNA extraction and sequencing success. First, we

140    investigated the relationship between residual formaldehyde concentration and pH in

141    preservation media through a survey of specimens in the Australian National Wildlife

142    Collection (ANWC; Crace, Australia). Next, in a phased approach, we compared DNA yield

143    achieved with three extraction methods - (1) hot alkaline lysis digestion followed by phenol-

144    chloroform extraction, (2) proteinase K digestion followed by phenol-chloroform extraction

145  and (3) proteinase K digestion followed by silica spin column purification. We then applied

146  the best-performing DNA extraction method to terrestrial vertebrate specimens representing

147  the broad range of tissue quality observed in museum specimens and tested performance of two

148  library preparation methods – (1) single-stranded method v2.0 (ss2) (25) and (2) BEST double-

149  stranded method (dsBEST) (26). Placing our results into context with a comprehensive and

150  unbiased survey of collection-wide spirit preservation conditions, we present a decision-

151  making framework to accelerate and facilitate genomic research using formalin-preserved

152  specimens.

153

## Results

### Preservation media condition survey

156  Within 149 ANWC specimen jars surveyed (23 amphibian, 40 mammal, 40 reptile, and 46

157  avian), preservation media pH ranged from 4.8–8.4 with 70 (47%), 61 (41%) and 18 (12%)

158  having neutral (6.5–7.5), low (< 6.5) and high (> 7.5) pH, respectively. Residual formaldehyde

159  concentration ([F]) ranged from 0–40,000 mg/L. High [F] (> 1000 mg/L) was detected in 61%

160  of low pH jars, 6% of neutral pH jars and 0% of high pH jars. We assumed specimens in jars

161  yielding [F] = 0 (n = 82) were preserved with ethanol and without exposure to formaldehyde.

162  Consistent with the practice of fixing specimens with unbuffered formalin combined with the

163  gradual degeneration of formaldehyde to formic acid, the pH of the formalin-preserved samples

164  (range 4.8–7.1; mean = 6.2) was significantly lower than for the ethanol-preserved samples

165  (range 6.1–8.4; mean = 7.1) (T-test; p < 0.0001; Supplementary Figure 1A). The recorded

166  collection date of the specimens ranged from 1936–2015. The time since collection (age) of

167  the ethanol-preserved specimens (mean = 40.1 years) was not significantly different than the

168  formalin-preserved specimens (mean = 36.1 years) (Supplementary Figure 1B). Among the

8

169    formalin-preserved samples, [F] and pH were negatively correlated (R = -0.6, p < 0.001; Figure

170    1). Age was not significantly correlated with either [F] or pH. Of the 12 specimens selected for

171    sequencing, collection date ranged from 1962–2006 and pH ranged from 4.9–8.2. Three

172    sequenced specimens were ethanol-preserved and nine sequenced specimens were formalin-

173    preserved with [F] ranging from 325–20,000 mg/L (Table 1).

174

175    **Table 1. Specimen metadata and independently assessed preservation quality metrics**

176    **for samples selected for sequencing**

177    Twelve specimens (three ethanol-preserved and nine formalin-preserved) from the ANWC

178    spirit vault were selected for DNA extraction and sequencing. Unique ANWC specimen IDs,

179    species names, common name, recorded year of collection, residual formaldehyde

180    concentration in the preservation media (mg/L), pH and tissue sampled for extraction are

181    given.

| Preservation | Specimen ID | Species name | Common name | Collection year | [F] (mg/L) | pH | Tissue sampled |
|---|---|---|---|---|---|---|---|
|  | ANWC B30438 | *Phalacrocorax carbo* | Great black cormorant | 1977 | 0 | 8.2 | Skin |
| Ethanol | ANWC B00001 | *Aquila audax* | Wedge-tailed eagle | 1973 | 0 | 7.68 | Liver |
|  | ANWC M15492 | *Phascolarctos cinereus* | Koala | 1971 | 0 | 7 | Muscle |
|  | ANWC A02522 | *Rhinella marina* | Cane toad | 2002 | 2050 | 6.41 | Liver |
|  | ANWC M11465 | *Macropus eugenii* | Tammar wallaby | 1989 | 8000 | 5.26 | Liver |
|  | ANWC R03280 | *Crocodylus porosus* | Saltwater crocodile | 1973 | 4000 | 6.31 | Liver |
|  | ANWC B47838 | *Melopsittacus undulatus* | Budgerigar | 1996 | 5000 | 6.3 | Liver |
| Formalin | ANWC R06312 | *Pogona minima* | Dwarf dragon | 1986 | 1800 | 7.04 | Liver |
|  | ANWC R01545 | *Pogona vitticeps* | Central dragon | 1971 | 325 | 6.24 | Liver |
|  | ANWC B40690 | *Taeniopygia guttata* | Zebra finch | 1986 | 20000 | 4.86 | Muscle |
|  | ANWC B34691 | *Falco cenchroides* | Australian kestrel | 2006 | 2000 | 5.45 | Liver |
|  | ANWC M03973 | *Ornithorhynchus anatinus* | Platypus | 1962 | 10000 | 5.79 | Muscle |

182

**DNA quantification**

We compared DNA yield from the hot alkaline lysis (HA), proteinase K plus phenol-chloroform (proK-PC) and proteinase K plus column (proK-col) extraction methods for the *Rhinella marina*, *Macropus eugenii* and *Crocodylus porosus* specimens and observed no significant differences between extraction methods (one-way ANOVA; Supplemental Figure 3A). However, the HA method produced more DNA from the two poor quality specimens (*M. eugenii* and *C. porosus*) compared to either of the proteinase K methods (Table 2). Thus, we predicted the HA method would perform better on specimens ranging broadly in preservation quality and we used this method to extract the remaining nine specimens. HA extraction yielded DNA detectable by high sensitivity Qubit for all twelve specimens. Two ethanol-preserved specimens (*Aquila audax* and *Phascolarctos cinereus*) and two formalin- preserved specimens (*R. marina* and *Melopsittacus undulatus*) yielded > 1,000 ng total DNA from 50 mg of tissue (Table 2). Three specimens, *Phalacrocorax carbo*, *Taeniopygia guttata* and *Ornithorhynchus anatinus*, yielded particularly low (< 100 ng) total DNA from 50 mg of tissue (Table 2). We observed no significant difference in DNA yield between ethanol and formalin-preserved specimens (T-test; Supplemental Figure 3B). However, mean DNA yield from ethanol-preserved specimens was more than double that from formalin-preserved specimens. Mean DNA yield from formalin-preserved specimens in preservation media with low pH (< 6) was not significantly different from those in media with neutral to high pH (> 6) (Supplemental Figure 3C). DNA yield was significantly higher from formalin-preserved liver tissue compared to non-liver tissue (T-test; $p < 0.05$; Supplemental Figure 3D). Both [F] and age showed a negative but non-significant correlation with DNA yield from formalin-preserved specimens (Supplemental Figures 3E and 3F).

206  **Table 2. Sequencing and alignment statistics**

207  For all specimens, DNA yield is given for the individual extractions of 50 mg of tissue. For

208  the remaining metrics, the values shown were calculated having combined both the ss2 and

209  dsBEST libraries. The number of raw reads is given as a sum of all single reads (R1 and R2)

210  from the paired-end sequencing run. Reads aligned indicates the percent of raw reads aligned

211  to reference genome after removal of PCR and optical duplicates. The mean aligned insert

212  length is the mean length (in bp) of the aligned portion of the read. $C_{nuc}$ is the coverage of the

213  nuclear genome. $C_{mt}$ is the proportion of mitochondrial genome with greater than 30X

214  coverage. $C_{pot}$ is the estimated potential genomic coverage if the full library had been

215  sequenced, calculated from the estimated library complexity. MRM is the number of reads

216  aligned to the mitochondrial genome per one million raw reads.

| Preservation | Species | Extraction method | DNA yield from 50 mg (ng) | Raw reads (million) | Reads aligned (%) | Mean aligned insert length (bp) | $C_{nuc}$ | $C_{mt}$ | $C_{pot}$ | MRM |
|---|---|---|---|---|---|---|---|---|---|---|
| Formalin | *Rhinella marina* | HA | 1,860 | 434 | 21 | 65 | 2.2 | 0.78 | 6.2 | 52 |
| | | proK-col | 666 | 77 | 40 | 81 | 1 | 0 | 6.2 | 14 |
| | | proK-PC | 2,550 | 321 | 15 | 74 | 1.2 | 0.42 | 11.4 | 29 |
| | *Macropus eugenii* | HA | 271 | 306 | 8 | 56 | 0.5 | 0.59 | 2.7 | 50 |
| | | proK-col | 4 | 17 | 1 | 67 | 0 | 0 | 0.1 | 3 |
| | | proK-PC | 33 | 801 | < 1 | 65 | 0 | 0 | 0.1 | 2 |
| | *Crocodylus porosus* | HA | 130 | 23 | < 1 | 67 | 0 | 0 | 0 | 11 |
| | | proK-col | None detected | 160 | < 1 | 70 | 0 | 0 | 0.1 | 12 |
| | | proK-PC | 79 | 294 | < 1 | 62 | 0 | 0 | 0 | 2 |
| | *Melopsittacus undulatus* | HA | 2,400 | 318 | 20 | 60 | 3.1 | 0.94 | 23.6 | 201 |
| | *Pogona minima* | HA | 521 | 367 | 7 | 58 | 0.8 | 0.51 | 7.5 | 29 |
| | *Pogona vitticeps* | HA | 672 | 432 | 15 | 59 | 2.1 | 0.85 | 7.9 | 52 |
| | *Taeniopygia guttata* | HA | 15 | 62 | < 1 | 66 | 0 | 0 | 0 | 1 |
| | *Falco cenchroides* | HA | 690 | 303 | 5 | 56 | 0.7 | 0.12 | 2.1 | 14 |
| | *Ornithorhynchus anatinus* | HA | 22 | 520 | < 1 | 70 | 0 | 0.13 | 0.8 | 20 |
| Ethanol | *Phalacrocorax carbo* | HA | 57 | 292 | < 1 | 69 | 0.10 | 0.90 | 0.60 | 50.00 |
| | *Aquila audax* | HA | 1,932 | 282 | 67 | 76 | 11.3 | 0.98 | 323 | 2515 |
| | *Phascolarctos cinereus* | HA | 1,254 | 423 | 60 | 76 | 5.4 | 0.94 | 93 | 2606 |

11

**Pre-alignment library quality assessment**

Prior to alignment, we used FastQC to assess the quality of paired-end reads from ss2 and dsBEST libraries. All libraries contained a high proportion of adapter content and low read quality score beginning at roughly 50 bp, consistent with highly fragmented input DNA. Focusing on the first 75 bp of the raw reads, mean sequence quality was slightly but significantly higher for read 2 (mean Phred score = 34.3) than for read 1 (mean Phred score = 33.7) across all libraries (paired T-test; $p < 0.001$). Likewise, the mean sequence quality was significantly higher in ss2 libraries compared to the corresponding dsBEST libraries for both read 1 (mean of the differences = 2.1; paired T-test; $p < 0.001$) and read 2 (mean of the differences = 0.79; paired T-test; $p < 0.01$). Mean sequence quality was not significantly different between reads derived from ethanol and formalin-preserved tissues, even when excluding libraries prepared from less than 200 ng of input DNA (paired T-test). We found evidence of cross-contamination in several libraries prepared from low DNA yield extractions. Compared to negative controls, both *O. anatinus* libraries and all but two *C. porosus* libraries showed a higher number of reads classified as genus *Mus* by Kraken2 (Supplementary Table 2). The *O. anatinus* libraries also contained a high percentage of reads classified as *Homo sapiens* (9.7% and 25%). The *O. anatinus* and *C. porosus* tissues were among those that yielded the least DNA. The *O. anatinus* HA extraction yielded just 22 ng. The *C. porosus* HA and proK-PC extractions yielded 130 and 79 ng, respectively, while the proK-col extraction yielded no detectable DNA. The only other specimens to yield less than 500 ng were the *P. carbo*, *T. guttata* and *M. eugenii*.

**Relative alignment quality from three extraction methods**

We used three indicators of alignment quality to compare the relative success of the three extraction methods on the *R. marina*, *M. eugenii* and *C. porosus* specimens: percent of raw reads aligned to the genome (% alignment), the number of reads aligned to the mitochondrial

12

242  genome per million raw reads (MRM) and the mean aligned insert length. Among these three

243  specimens, we observed no significant differences between library preparation methods in any

244  of the three alignment quality indicators (paired T-tests). Therefore, we took the mean of the

245  two library preparations to compare extraction methods across each alignment quality

246  indicator. Again, we observed no significant difference between the three extraction methods

247  applied to the *R. marina*, *M. eugenii* and *C. porosus* specimens in any of the three alignment

248  quality indicators (one-way ANOVA). All six *C. porosus* libraries yielded < 1% alignment

249  (Figure 2A and Table 2), indicating failure of all extraction and library preparation methods on

250  this specimen. Excluding the *C. porosus* libraries, we observed significant differences in MRM

251  between the extraction methods (one-way ANOVA; $p < 0.05$) with the HA method producing

252  significantly more MRM than both the proK-col and proK-PC methods (Tukey tests; $p < 0.05$).

253  We observed no significant difference in MRM between the proK-col and proK-PC methods

254  (Tukey tests) nor in % alignment or mean insert length between the three extraction methods

255  (one-way ANOVA).

**Effect of specimen quality on sequencing success**

257  The percentage of aligned reads removed by optical and PCR de-duplication varied between

258  8.8% and 99.5% across all libraries. Among the HA alignments, de-duplication reduced

259  significantly more mapped reads from dsBEST libraries than from ss2 libraries (paired T-test;

260  $p < 0.01$). Combining the ss2 and dsBEST libraries for each HA extraction, de-duplication

261  removed more than double the percentage of reads (69.8% versus 32.8%) from poor quality

262  specimens (those yielding < 1% reads aligned) compared to better quality specimens (those

263  yielding > 1% reads aligned). However, this difference was not significant (T-test). De-

264  duplication removed significantly more reads from the formalin-preserved specimens (mean =

265  54.6%) than from the ethanol-preserved specimens (mean = 16.7%) (T-test; $p < 0.01$).

266  Following de-duplication, the mean percent of mapped reads remaining was 44% and 59% for

267    the dsBEST and ss2 HA libraries, respectively. Across all specimens extracted using the HA

268    method, we observed no significant differences between library preparation methods in any of

269    the three alignment quality indicators (paired T-tests). Therefore, we conducted further

270    comparison of the effect of specimen quality on alignment success taking the mean of each

271    alignment quality indicator from the two HA library preps.

272    HA extraction of one of three ethanol-preserved specimens (*P. carbo*) and three of nine

273    formalin-preserved specimens (*C. porosus*, *T. guttata* and *O. anatinus*) produced < 1% aligned

274    reads (Table 2), indicating equal rates of very poor sequencing success with ethanol- and

275    formalin-preserved tissues. Excluding the specimens with < 1% aligned reads, the ethanol-

276    preserved specimens produced a significantly higher percentage of aligned reads (T-test; $p <$

277    $0.01$). Two of the three ethanol-preserved specimens (*A. audax* and *P. cinereus*) produced >

278    60% aligned reads while the remaining six formalin-preserved specimens (*R. marina*, *M.*

279    *eugenii*, *M. undulatus*, *Pogona minima*, *Pogona vitticeps* and *Falco cenchroides*) produced

280    between 5% and 21% aligned reads (Table 2). Excluding the specimens with < 1% aligned

281    reads, the mean insert length was significantly longer for the ethanol-preserved specimens

282    (mean = 76 bp) compared to the formalin-preserved specimens (mean = 59 bp) (T-test; $p <$

283    $0.0001$). MRM was also significantly higher for the ethanol-preserved specimens (mean =

284    2,560) compared to the formalin-preserved specimens (mean = 43) (T-test: $p < 0.01$).

285    The percentage of reads aligned increased with preservation media pH (R = 0.44; Figure 3A),

286    decreased with preservation media [F] (R = -0.53; Figure 3B) and decreased with specimen

287    age (R = -0.46; Figure 3C), although these correlations were not statistically significant. The

288    percentage of aligned reads was significantly higher in specimens sampled with liver than those

289    sampled with muscle and skin (T-test; $p < 0.05$; Figure 3D). Of the specimens yielding poor

290    sequencing success (< 1% reads aligned), all but *C. porosus* were sampled with either muscle

14

291    or skin as liver was not present. The only specimen sampled with a tissue other than liver to

292    yield a percent of reads aligned > 1% was the ethanol-preserved *P. cinereus*.

293    **Genome sequencing coverage**

294    Nuclear genome coverage ($C_{nuc}$) of the deduplicated alignments was < 1X for the majority of

295    libraries. Since raw read yield was highly variable, $C_{nuc}$ is not an appropriate measure with

296    which to compare the extraction or library preparation methods. However, it is noteworthy that

297    we achieved $C_{nuc}$ > 1X for two of the ethanol-preserved specimens and three of formalin-

298    preserved specimens. Combining all libraries for a given specimen, we achieved a total of 5.4X

299    and 11.3X $C_{nuc}$ for the ethanol-preserved *P. cinereus* and *A. audax* specimens, respectively

300    (Table 2). Likewise, we achieved a total of 2.1X, 3.1X and 4.4X $C_{nuc}$ for the formalin-preserved

301    *P. vitticeps*, *M. undulatus* and *R. marina* specimens, respectively (Table 2). To estimate the

302    potential for improving $C_{nuc}$ through re-sequencing of the prepared libraries, we calculated

303    potential genomic coverage ($C_{pot}$) (Table 2). Combining all libraries for a given specimen, $C_{pot}$

304    exceeded 20X for the *R. marina* and *M. undulatus* and exceeded 75X for the *P. cinereus* and

305    *A. audax*. Focussing on the mitochondrial genome, the proportion of sites with 30X or higher

306    coverage ($C_{mt}$) was nearly complete (> 0.9) for all three ethanol-preserved specimens (Table

307    2). $C_{mt}$ for the formalin-preserved *M. undulatus* (0.94) was comparable to that of the ethanol-

308    preserved specimens. $C_{mt}$ was moderate to high (> 0.5) for five of the formalin-preserved

309    specimens (Table 2). Only the *C. porosus*, *T. guttata*, *F. cenchroides* and *O. anatinus* yielded

310    very poor $C_{mt}$ (< 0.15).

311    **Read length periodicity**

312    From the aligned insert lengths estimated with Picard, we plotted the frequency of reads

313    between 50 and 100 bp (Figure 4). This plot revealed a pattern of read length periodicity in

314    several specimens, notably those that resulted in higher mapping success. We observed

315    prominent periodicity of approximately 10.1 bp in the *R. marina* specimen extracted with the

15

316  proK-PC method. While less pronounced, we observed read length periodicity of

317  approximately 10.8 bp in the HA extractions of *R. marina*, *P. vitticeps*, *P. minima*, *F.*

318  *cenchroides*, *A. audax* and *P. cinereus*. The pattern of periodicity was observed in both the

319  dsBEST and ss2 libraries, however, it was slightly more pronounced in the dsBEST libraries.

## Discussion

321  In this study, we present evidence challenging the common perception that formalin-preserved

322  museum specimens are devoid of accessible DNA. Processed with a tailored molecular and

323  bioinformatic workflow, formalin-preserved specimens had an overall sequencing success rate

324  equivalent to ethanol-preserved specimens, albeit with recovery of a lower percentage of

325  sequence reads mapping to the reference genome. Contrary to popular belief, we found

326  genome-wide nuclear data is retrievable from some formalin-preserved museum specimens,

327  even with a moderate investment of sequencing effort (with 30% of formalin-preserved

328  specimens, we achieved > 2X nuclear genome coverage from 300-500 million raw reads). We

329  also show reconstruction of large sections of the mitochondrial genome is possible even in poor

330  quality specimens where limited nuclear data were recovered (with 55% of formalin-preserved

331  specimens, we achieved > 30X coverage of more than 50% of the mitochondrial genome).

332  Investigating specimens covering a range of preservation quality, we also developed a decision-

333  making framework to improve sequencing success rate and prioritize suitable specimens. Our

334  findings support a considered and targeted sequencing approach that transforms thousands of

335  spirit collection specimens into a new molecular resource. Improved access to genomic data

336  held in these specimens has the potential to inform research into the mechanisms driving

337  adaptation, evolution, speciation and extinction.

16

338 **Hot alkaline lysis effectively recovers gDNA from formalin-preserved archival tissues**

339 **suitable for next generation sequencing.**

340 Originally developed for DNA extraction from FFPE sections, the HA method relies on high

341 heat (120°C) under alkaline conditions (pH = 13) to break strong inter- and intramolecular

342 cross links and utilizes organic extraction to maximize capture of fragmented gDNA from

343 formalin-preserved tissues (50–52). This method has been applied to museum specimens to

344 successfully recover sections of the mitochondrial genome in trout (53) and full mitochondrial

345 genomes from lizards (40) and bacterial symbionts (41). Here we show the HA yields gDNA

346 in adequate quantities for WGS from higher-quality formalin-preserved museum specimens.

347 Coupled with library preparation methods designed to efficiently convert degraded DNA, we

348 produced complex sequencing libraries with the potential to recover full vertebrate genomes

349 when mapped using a strategy optimized to maximize recovery of endogenous sequence. Our

350 results indicate that the HA method is appropriate for DNA extraction from a broad range of

351 taxa preserved under various conditions, making it well-suited for application in both museum

352 and pathological settings.

353 In a small-scale comparison to proK digestion with either phenol-chloroform extraction or

354 column purification, the HA method performed superiorly for poor quality formalin-preserved

355 specimens. We experienced equal success rates with the HA method in formalin and ethanol-

356 preserved tissues. It is not standard practice to apply the HA method to ethanol-preserved

357 specimens, which do not suffer from cross-linking, but we implemented it in this study to serve

358 as a comparison to formalin-fixed tissues. Thus, while the HA method is likely unnecessarily

359 harsh for recovery of DNA from tissues not crosslinked with formaldehyde, we propose this

360 extraction method is suitable across a wide range of tissue qualities and preservation conditions

361 observed in museum spirit collections. And, given that we achieved relatively high yield from

362 the ethanol-preserved tissues, we propose that the HA method is appropriate in cases where

17

363   contact with formalin cannot be determined. We caution; however, the HA method's success

364   may be limited to DNA-rich tissues such as liver. Our HA extractions of formalin-preserved

365   muscle and ethanol-preserved skin tissue failed to yield adequate gDNA for sequencing, while

366   our HA extraction of ethanol-preserved muscle tissue was less successful than our extraction

367   of ethanol-preserved liver tissue. HA extraction has been previously observed to perform

368   poorly compared to cetyltrimethylam-monium bromide (CTAB) protocols on formalin-

369   preserved mammalian heart tissue (54). We also note that, preservation conditions being equal,

370   DNA yield may differ between taxonomic groups due to factors such blood cell nucleation.

371   Due to low sample size, we were not able to test if the lack of nucleated red blood cells in

372   mammal tissues impacted DNA yield.

373   **aDNA library preparation methods effectively capture DNA extracted from formalin-**

374   **preserved archival tissues**

375   DNA degradation in museum specimens is a significant challenge to genome sequencing. To

376   improve our conversion of degraded DNA from formalin-preserved tissues into high quality

377   library molecules, we utilized two library preparation methods developed specifically for

378   degraded aDNA templates. We tested the ss2 (46) and dsBEST (47) methods on DNA extracted

379   from both ethanol and formalin-preserved archival tissues. Sequence quality was significantly

380   higher for libraries prepared using the ss2 method compared to the dsBEST protocol. However,

381   this quality difference did not result in significantly lower rates of read alignment or reduced

382   mapped insert length for the dsBEST libraries. While we did not see differences in

383   contamination rates between the two methods, an advantage of the dsBEST method is its

384   reliance on fewer tube transfers and additions of solution, thus reducing opportunities to lose

385   DNA and introduce contaminants. The ss2 and dsBEST methods performed similarly on all

386   twelve of our archival templates, indicating both are well-suited to prepare libraries from DNA

387   extracted from ethanol and formalin-preserved tissues. Alternative library preparation methods

18

388  developed specifically for degraded DNA may prove equally effective. To maximize

389  conversion of fragmented archival DNA template, we advise using a library preparation

390  method designed to capture small fragments whilst minimising contamination risk. Overall, we

391  observed samples with very low DNA yield (< 200 ng from 50 mg of tissue) did not produce

392  libraries with high rates of mapping success. Thus, as a cost-saving measure, we advise

393  quantifying DNA templates prior to library preparation and focussing sequencing effort on

394  higher yielding samples.

395  **High alignment rates of fragmented DNA are achieved through exhaustive match**

396  **searching**

397  Removal of adapter sequence and low-quality bases via read-trimming is a standard pre-

398  processing procedure conducted on raw sequencing reads prior to mapping. In the context of

399  libraries prepared from highly degraded templates, filtering and trimming can reduce the

400  dataset substantially. For example, pre-processing of the library prepared from a formalin-

401  preserved *Anolis* lizard reduced the dataset to 13.5% of the raw data (40). Although filtering

402  and trimming are effective at removing PCR duplicates and erroneous bases introduced through

403  library preparation and sequencing, quality control parameters should be optimized to avoid

404  removing informative endogenous sequence, particularly with data derived from highly

405  fragmented low-input templates. Compared to DNA extractions from fresh tissue, our

406  extractions from formalin-preserved specimens were highly fragmented as is typical of aDNA

407  sources (55). We opted to trial a computationally efficient approach that eliminates loss of

408  endogenous sequence during pre-processing. The kalign function from the open source kit4b

409  toolkit performs alignments of raw reads by searching for the maximum length match within

410  the read to the reference sequence regardless of the match's position within the read. For each

411  raw read, kalign performs a rapid complete exhaustive match search across the indexed

412  reference genome. The match search is performed recursively through seed expansions

413    generated along the read length. The longest match to endogenous sequence is retrieved while

414    satisfying the minimum length threshold of the match. Using this approach, we aligned up to

415    21% and 67% of raw reads from formalin and ethanol-preserved tissues, respectively. These

416    alignment rates are consistent with the degree of degradation in the DNA we extracted from

417    spirit-preserved museum specimens being intermediate between that of fresh and truly ancient

418    tissues. A previous application of the ss2 method yielded a maximum of 11.3% mappable reads

419    from libraries prepared from aDNA tissue sources (25). The same study yielded 60% and 68%

420    mappable reads from libraries prepared from horse and pig liver stored in buffered formalin for

421    5 and 11 years, respectively (25). In comparison, our modest alignment rates may be the result

422    of tissues of intermediate age and using a different metric of calculating the percent of mapped

423    reads.

424    **Sequencing success is strongly influenced by specimen integrity prior to fixation**

425    To explore the effects of formalin-fixation on sequencing success, we selected three specimens

426    preserved with ethanol only and nine specimens preserved with formalin. We found no

427    significant difference in DNA yield between the ethanol and formalin-preserved specimens and

428    the differences we observed in DNA fragment lengths were minimal. Furthermore, we

429    observed equal rates of very poor sequencing success within ethanol and formalin-preserved

430    specimens, indicating preservation method is not a strict determinant of sequencing success.

431    Older, poor-quality ethanol-preserved specimens have previously been shown to be as

432    problematic for genomic analyses as formalin-preserved specimens (42,56). This is not to say

433    preservation method does not impact sequencing success. Two of our ethanol-preserved

434    specimens (*P. cinereus* and *A. audax*) had much higher mapping rates (60% and 67% reads

435    aligned, respectively) than even our most successful formalin-preserved specimens (*R. marina,*

436    produced 21% reads aligned with the HA method). Our findings indicate WGS of formalin-

437    preserved museum specimens is possible using HA extraction paired with a library preparation

438     optimized for conversion of degraded DNA. However, as with all potential DNA sources, the

439     overall integrity of the tissue will ultimately determine sequencing success.

440     The specimens with poor sequencing success (< 1% reads aligned) were largely older, their

441     preservation media had lower pH and higher [F] and they were sampled with a tissue other than

442     liver. On the contrary, the specimens with better sequencing success were preserved more

443     recently, their preservation media had neutral pH and lower [F] and the tissue sampled was

444     liver. We calculated the correlation between specimen quality measures ([F], pH, age and tissue

445     type) and both DNA yield and mapping success. Tissue type was the only quality measure

446     significantly associated with lower DNA yield, with liver yielding significantly more DNA

447     than either muscle or skin. Our higher success with liver is consistent with findings of a

448     previous study comparing sequencing success from liver, muscle and tail-tip in a formalin-

449     preserved *Anolis* lizard (40). However, in that study, the tissues were extracted using different

450     methods and thus it could not be determined if success was driven by tissue type or extraction

451     method.

452     Post-mortem DNA degradation occurs more rapidly in liver relative to other bodily tissues

453     including skeletal muscle, heart and brain (57,58). In the museum curatorial setting, specimens

454     undergo varying degrees of post-mortem decay prior to fixation. As is the case for most

455     museum specimens, the length of the post-mortem interval (PMI) was not recorded for the

456     specimens used in this study. Given expected rapid decay of the viscera, we used the visual

457     appearance of the gut contents as a reasonable proxy for the length of the PMI. The four

458     specimens used in this study that lacked liver tissue were visibly more degraded than those

459     with intact liver tissue (Supplementary Figure 2). In the case of the *P. cinereus*, *P. carbo* and

460     *O. anatinus*, the complete absence of viscera indicated the internal organs were likely well-

461     degraded and discarded prior to fixation. For specimens preserved after a long PMI, DNA

462     integrity throughout the carcass would be lower than in specimens preserved after a short PMI.

463   Therefore, we conclude that the higher yield from specimens sampled with liver is a reflection

464   of overall specimen quality and DNA damage occurring post-mortem but prior to fixation.

**Re-thinking formalin damage**

466   Formalin-preserved museum specimens have long been considered intractable sources of

467   gDNA. Encouragingly, we found specimen contact with formaldehyde does not prohibit DNA

468   sequencing if tissue decomposition occurring prior to fixation is minimized. With appropriate

469   sample vetting (Figure 5), HA extraction and DNA library preparation optimized for degraded

470   DNA, historical genomic data may be extracted from many formalin-preserved specimens.

471   These data will not be of similar quality to those recovered from fresh or ethanol-preserved

472   tissues. However, higher sequencing volume and borrowing of analytical methods from the

473   field of aDNA may facilitate reconstruction of historical genomes from formalin-preserved

474   tissues. We found evidence that DNA damage in formalin-preserved specimens shares

475   characteristics with that of aDNA. In addition to capturing shorter fragments with low mapping

476   rates, we observed a pattern of read length periodicity of approximately 10 bp. This is

477   consistent with observations in aDNA specimens (59) and is an interval that coincides with the

478   length of a turn of the DNA helix. Pederson et al (2014) attributed the 10 bp read periodicity

479   in specimens greater than 4,000 years old to protection of the DNA by nucleosomes

480   preferentially positioned at 10 bp intervals. We observed a striking periodicity pattern

481   averaging 10.8 bp in HA extracted samples and 10.1bp in the proK-PC samples. The shorter

482   periodicity in the proK treated samples may be due to reduced protection of the ends of DNA

483   fragments by digestion of the nucleosomes during extraction. We did not observe a signal of

484   nucleosome occupancy in read depth or in enrichment of fragments of nucleosome length (147

485   bp) as did Pederson et al., perhaps because we sequenced shorter fragments to comparatively

486   low depth. However, the appearance of 10 bp periodicity suggests it may be possible to infer

487   nucleosome occupancy from patterns of DNA degradation observed in formalin-preserved

488   specimens if higher coverage is achieved.

**489   Managing expectations**

490   We have shown WGS of formalin-preserved museum specimens is feasible and success can be

491   improved through specimen quality vetting. We stress; however, measures of specimen quality

492   are imperfect and the key parameters may vary between and within museum collections.

493   Modern collection institutions aim to limit light exposure and temperature variation within their

494   spirit vaults. With older specimens, the likelihood they have been exposed to undocumented

495   DNA-degrading conditions increases. We found the age of the specimen was not strongly

496   predictive of sequencing success, however, we did not sample specimens collected prior to the

497   1960s. This warrants further investigation into the extent to which intact DNA can be extracted

498   from much older formalin-preserved specimens.

499   While preservation media pH and [F] were not predictive of sequencing success in our

500   specimens, we note these measures do not always accurately reflect preservation condition.

501   Most institutions periodically top up the specimen jars in their spirit vaults to replace ethanol

502   lost through evaporation. In some cases, the preservation media is replaced entirely. Thus,

503   media pH and [F] values at the time of sampling for sequencing may not reflect preservation

504   and long-term storage conditions. With additional sampling of older and more varied

505   specimens, it may be possible to establish clear correlates of sequencing success associated

506   with pH and [F].

507   Both researchers and museums would benefit from an improved set of guidelines for strategic

508   decision making based on independent quality metrics rather than qualitative *ad hoc*

509   assessments. This will empower researchers to most effectively deploy their sequencing

510   budgets and support museums in deciding when to grant requests for destructive sampling. A

511    cost-benefit analysis should be conducted prior to genomic sequencing of museum specimens.

512    From the perspective of the museum, destructive sampling should be avoided if the specimen

513    is unlikely to yield sufficient DNA to achieve a project's aims. From the perspective of the

514    researcher, sequencing of high-quality specimens should be prioritized to generate high-quality

515    data. To assist in making these assessments, we provide a decision-making tree (Figure 5) for

516    use by both curators and researchers to determine which specimens are likely to be appropriate

517    for genomic analyses.

518    Ultimately, museum curators decide if the potential benefit of sequencing outweighs the

519    damage to the specimen through destructive sampling. Once sampling and DNA extraction has

520    been completed, the decision to proceed with library preparation and sequencing can be made

521    on the basis of DNA yield. We found specimens with high DNA yield (> 1,500 ng/50 mg

522    tissue) produced a high percentage (> 20%) of mappable reads while specimens with low DNA

523    yield (< 200 ng/50 mg tissue) produced virtually no mappable reads. While specimens yielding

524    between 200–1,500 ng of DNA per 50 mg tissue produced relatively low genomic coverage,

525    they did produce high coverage of the mitochondrial genome. Thus, reconstruction of historical

526    mitochondrial haplotypes may be possible from specimens yielding low quantities of DNA.

527    When nuclear data is required, high-volume sequencing should be reserved for high-quality

528    specimens. Generally speaking, most research projects aim to sequence a small number of

529    museum specimens with which to provide a base-line for comparison to contemporary

530    specimens. In light of the limited availability of historical specimens in collections, it is often

531    reasonable and feasible to allocate a relatively large budget to conduct deep sequencing of a

532    small number of specimens.

24

## Conclusions

534  Our results demonstrate formalin-fixation is not a complete barrier to WGS in museum

535  specimens. While success is not a guarantee, the use of HA lysis for DNA extraction

536  followed by an appropriate sequencing library preparation optimized for degraded DNA can

537  produce libraries of sufficient complexity for genomic analyses. When selecting specimens

538  for sequencing, our results indicate those with poor gut integrity are least likely to yield

539  sufficient DNA for sequencing.

## Methods

**Preservation media condition survey**

542  We conducted an unbiased survey of the ANWC spirit vault to measure variation in

543  preservation characteristics that can be sampled without disturbing the specimen. We randomly

544  selected 149 specimen jars spanning a range of taxonomic groups and ages, and removed a 25

545  mL aliquot of preservation media. We measured pH using an Orion$^{TM}$ Versa Star Pro$^{TM}$

546  benchtop pH meter (*Thermo Scientific*) and [F] using MQuant® test strips (*Merck*). Where [F]

547  was at the upper detection limit of the test strips, we diluted the aliquot 1:10 with ultrapure

548  water and remeasured, extrapolating the neat concentration of the media by multiplying the

549  measurement by the dilution factor.

**Specimen selection**

551  To select specimens for genomic sequencing, we first identified those with a publicly available

552  whole-genome reference for the specimen species or closely related species. Of these

553  specimens, we selected 12 representing a range of taxonomic groups, preservation conditions

554  and ages and sampled 50 mg of tissue. We sampled liver tissue when it was available. Muscle

555  was sampled from an ethanol-preserved *P. cinereus* specimen and from formalin-preserved *T.*

556  *guttata* and *O. anatinus* specimens. Skin was sampled from an ethanol-preserved *P. carbo*. All

557  specimens sampled with liver were preserved as whole animals whereas substantial portions of

558  the body were absent from those specimens sampled with muscle or skin (Supplementary

559  Figure 2). From the nine formalin-preserved specimens, we selected three with which to test

560  the relative success of three DNA extraction methods. To represent "good" quality formalin-

561  preserved specimens, we selected a cane toad (*R. marina*) preserved in 2002. Visually, this

562  specimen appeared minimally degraded and measurements of the storage media indicated low

563  [F] and a neutral pH. To represent "poor" quality formalin-preserved specimens, we selected a

564  tammar wallaby (*M. eugenii*) preserved in 1989 and a saltwater crocodile (*C. porosus*)

565  preserved in 1973. Visually, these two "poor" specimens were reasonably well-preserved,

566  however, measurements of the storage media indicated substantial [F] in both specimen jars

567  and mildly acidic pH in that of the wallaby.

**Tissue preparation**

569  Prior to DNA extraction, we liquid nitrogen pulverized all dissected tissue into a fine powder

570  using a cryoPREP® (*Covaris*) dry pulverizer (three impacts to a TT05 tissueTUBE™ on

571  intensity setting three; 10 sec in liquid nitrogen between impacts). We then stored the

572  pulverized tissue powder in 70% ethanol at -80°C until further processing. We re-hydrated the

573  pulverized tissue by stepping it into 50% ethanol, 30% ethanol then TE buffer with rocking for

574  10 min intervals. For the nine formalin-fixed tissues, we quenched excess formaldehyde by

575  rocking for 2 hrs in 1 mL GTE buffer (100 mM glycine, 10 mM Tris-HCL, pH 8.0, 1 mM

576  EDTA), followed by a further wash in fresh GTE for 2 hrs and a final fresh GTE wash overnight

577  at room temperature. We removed the GTE buffer and washed with rocking in sterile water for

578  10 min.

**Proteinase K DNA extraction**

580  We conducted two variations on a standard proteinase K (proK) digestion. For each specimen,

581  we digested two 50 mg (wet weight) aliquots of tissue overnight at 55°C with 30 μL of 20

582    mg/mL proteinase K in 970 µL lysis buffer (10 mM NaCl, 20 mM Tris-HCl, pH 8.0, 1 mM

583    EDTA, 1% SDS). We isolated DNA from the proK lysates with either (A) three extractions of

584    phenol-chloroform followed by ethanol precipitation (proK-PC), resuspending the DNA in 30

585    µL TE, or (B) a QIAquick PCR purification column (*Qiagen*) (proK-col), following the

586    manufacturer's instructions and eluting the DNA in 30 µL TE. Alongside the museum tissues,

587    we processed tissue-free controls. We quantified extracted dsDNA using a Qubit fluorometer

588    and high sensitivity (HS) DNA kit (*Invitrogen*).

589    **Hot alkaline lysis DNA extraction**

590    For the hot alkaline lysis (HA) extractions, we heated 50 mg (wet weight) tissue aliquots to

591    120°C for 25 min in 500 µL of alkali buffer (0.1 M NaOH with 1% SDS, pH 13) according to

592    methods described in (52). We purified DNA from the lysate with three phenol-chloroform

593    extractions followed by ethanol precipitation, resuspending the DNA in 30 µL TE. Alongside

594    the museum tissues, we processed tissue-free controls. We quantified extracted dsDNA using

595    a Qubit fluorometer and HS DNA kit.

596    **Library preparation methods**

597    To avoid cross-contamination, we prepared all sequencing libraries in the Ecogenomics and

598    Bioinformatics Laboratory trace facility at the Australian National University following

599    standard anti-contamination procedures.  We prepared libraries from all DNA extracts and

600    tissue-free controls using two methods developed for high efficiency conversion of fragmented

601    aDNA; the single-stranded method v2.0 (ss2) (25) and the BEST double-stranded method

602    (dsBEST) (26). Concurrently, we prepared DNA-free control libraries. For sequencing of Read

603    1 in both library preparation methods, we used an adapter with the sequence 5`–

604    AGATCGGAAGAGCACACGTCTGAACTCCAGTCAC–3`. For sequencing of Read 2, we

605    used adapters with the sequences 5`–GGAAGAGCGTCGTGTAGGGAAAGAGTGT–3` and

606    5`–AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT–3`  for  the  ss2  and  dsBEST

27

607    methods, respectively. We removed excess adapter and primer dimer by isolating fragments

608    between 160 bp and 400 bp from the resulting libraries using the PippinHT size-selection

609    system (*Sage Science*). We further purified the libraries with a MinElute PCR purification kit

610    (*Qiagen)* and quantitated the library concentrations using the LabChip GXII (*PerkinElmer*)

611    capillary electrophoresis system. We then pooled the libraries in approximately equimolar

612    concentrations and measured the concentration of the final pooled library using a Qubit

613    fluorometer and HS DNA kit. The Australian Genome Research Facility sequenced the pooled

614    library on a 150 bp paired-end S4 flow cell on the Illumina NovaSeq 6000 platform.

615    **Quality control of raw reads**

616    We computed quality control metrics for the raw reads using FastQC v.0.11.8 (60). Our adapter

617    content analysis included both default Illumina adapters and our custom library adapters. To

618    rapidly detect library contamination by non-target species' DNA, we classified the taxonomic

619    origin of reads using Kraken2 v.2.0.9b (61). We estimated the number of unique fragments

620    present in the raw sequence libraries with the EstimateLibraryComplexity function of PICARD

621    v.2.9.2 (62).

622    **Alignment**

623    We aligned reads to reference nuclear and mitochondrial genomes obtained from the DNA Zoo

624    Consortium (63,64) and GenBank (65) (Supplementary Table 1). Species-specific reference

625    genomes were not available for three of the specimens. For *A. audax*, *F. cenchroides* and *P.*

626    *minima,* we used the reference genomes of species in the same genera- *A. chrysaetos*, *F.*

627    *perigrinus* and *P. vitticeps*, respectively (Supplementary Table 1). We hard-masked the eleven

628    genomes with RepeatMasker v.4.1.0 (66) including our ss2 and dsBEST library adapters in the

629    repeat database and applying the -qq option allowing 10% less sensitivity while decreasing

630    processing time. We aligned raw reads with the kalign function of the ngskit4b tool suite

631    v.200218 (67) with options -c25 (--minchimeric=<int>; minimum chimeric length as a

632   percentage of probe length) -l25 (--minacceptreadlen=<int>; after any end trimming only

633   accept read for further processing if read is at least this length) -d50 (--pairminlen=<int>;

634   accept paired end alignments with observed insert sizes of at least this) -U4 (--pemode=<int>;

635   paired end processing mode: 4 - paired end no orphan recovery treating orphan ends as SE).

636   We removed PCR and optical duplicates from the alignments using the MarkDuplicates

637   function of PICARD enabling REMOVE_DUPLICATES=TRUE. For each de-duplicated

638   alignment, we computed a histogram of aligned insert lengths and calculated the mean aligned

639   insert length using the CollectInsertSizeMetrics function of PICARD.

640   **Genome coverage analyses**

641   We estimated nuclear genome coverage ($C_{nuc}$) as the number of unique aligned reads multiplied

642   by the mean insert length divided by unmasked genome size. To estimate how much genomic

643   coverage could be achieved by increasing sequencing depth, we calculated the sequenced

644   proportion of the prepared library as the number of read pairs examined divided by the

645   estimated library size. We estimated the number of possible reads represented in the prepared

646   library by dividing the number of actual reads aligned by the sequenced proportion of the

647   library. We then roughly estimated the potential genomic coverage represented in the full

648   prepared    library    ($C_{pot}$)    as:    $(\# \ possible \ reads \ \times \ mean \ insert \ length \ (bp)) \div$

649   $genome \ size \ (bp)$. To calculate the proportion of mitochondrial genome sites with 30X or

650   greater coverage ($C_{mt}$), we executed the Samtools *depth* function (68) on SAM files for the

651   mitochondrial contigs for each species combined across all libraries.

652   **Statistical analyses**

653   We performed statistical analyses in the R environment, v.4.0.2 (69) and produced figures

654   using the packages *ggplot2* (70) and *ggpubr* (71). To test if the residuals of data were normally

655   distributed, we ran Shapiro-Wilk tests with the function *shapiro.test*. We conducted T-tests

656   with the function *t.test*, analyses of variance (ANOVA) with the function *aov* and computed

29

657    confidence intervals using Tukey's Honest Significant Difference method (Tukey test) with the

658    function *TukeyHSD* in the base package *stats*. We computed Pearson correlation coefficients

659    with associated p-values with the *ggpubr* function *stat_cor*.

660 **Figures**

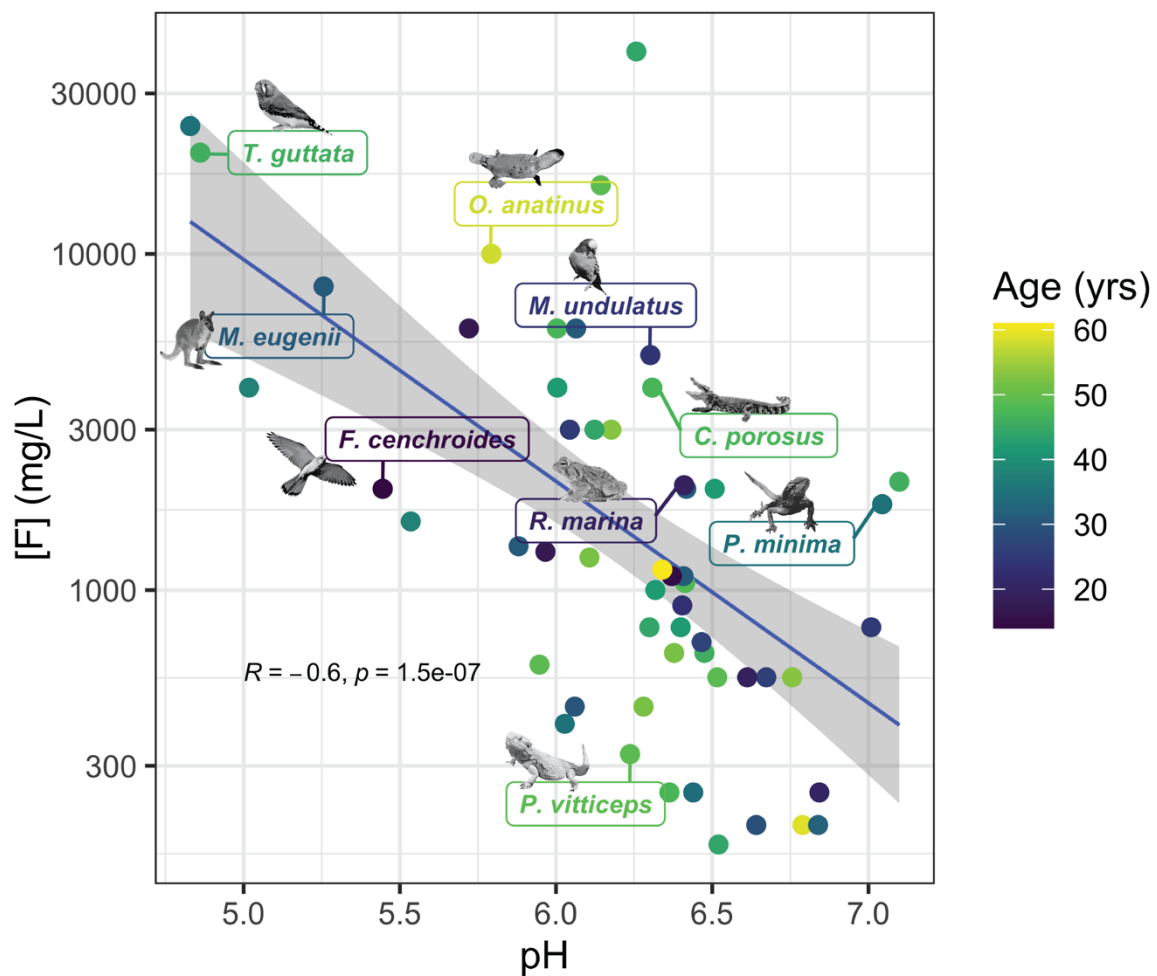661 **Figure 1. Preservation media survey results of formalin-fixed specimens in the**

662 **Australian National Wildlife Collection**

663 Residual formaldehyde concentration [F] (mg/L) is shown on a log-scale in relation to pH.

664 Individual specimens (N = 65) are colored by the time since their collection (age) and the

665 specimens selected for sequencing are indicated by species name. A linear model was used to

666 fit a regression line and standard error is shown in grey; R = Pearson's correlation coefficient.



667

668

669 **Figure 2. Effectiveness of extraction and library preparation methods for *R. marina*, *C.**

670 *porosus* and *M. eugenii* specimens.**

671 (A) Alignment to the whole genome expressed as the percentage of reads aligning (B)

672 Alignment to the mitochondrial genome expressed as the number of reads aligned per million

673 raw reads (MRM). dsBEST = BEST double-stranded method (26); ss2 = single-stranded

674 method v2.0 (25); HA = hot alkaline lysis; proK-col = proteinase K digestion followed by

675 column purification; proK-PC = proteinase K digestion followed by phenol-chloroform

676 extraction.



677

678    **Figure 3. Alignment results for hot alkali extracted samples**

679    The correlation between the percentage of reads aligned to the whole genome (combining

680    both library preparations of the hot alkali extracted specimens) and (A) preservation media

681    pH, (B) preservation media formaldehyde concentration (g/L), (C) number of years in the

682    collection and (D) tissue sampled. In A-C, all specimens are shown colored by their fixation

683    type and R = Pearson's correlation coefficient for the formalin-fixed specimens. In D, only

684    the formalin-preserved specimens are plotted and individual specimens are shown with black

685    dots, * = p < 0.05.



686

687 **Figure 4. Libraries with read periodicity**

688 The frequency of insert lengths, in bp, estimated from the mapped dsBEST libraries is shown

689 for six preserved specimens. Read periodicity in the *R. marina* libraries from the proteinase K

690 with phenol-chloroform (proK-PC) extractions averages 10.1 bp while periodicity in libraries

691 from the hot alkali extractions of six specimens averages 10.8 bp.



692

693 **Figure 5. Decision-making tree for a priori estimation of likely sequencing success in**

694 **spirit-preserved museum specimens.**

695 Green ticks indicate the specimen is well-suited the sequencing application and there is a high

696 likelihood of success. Black question marks indicate the specimen is marginal for the

697 sequencing application and there is high variation in the likelihood of success. Red crosses

698 indicate the specimen is not well-suited for the sequencing application and there is a low

699 likelihood of success.



700

## Declarations

**Ethics approval and consent to participate**

Not applicable.

**Consent for publication**

Not applicable.

**Availability of data and materials**

The sequencing data generated and analysed in this study are archived in the CSIRO Data Access Portal. Correspondence and requests for materials should be addressed to CEH (clare.holleley@csiro.au)

**Competing interests**

The authors declare they have no competing interests.

**Funding**

Funding for this study was provided by the Environomics CSIRO Future Science Platform (grants R-10011 and R-14486) awarded to CEH.

**Authors' contributions**

This study was conceived by CEH. Experiments were designed by CEH, MRA and AG and conducted by MRA and AG. Data analysis was conducted by EEH, JS, MRA and AG and advised by DMG and CEH. All authors contributed to the writing and editing of the manuscript.

**Acknowledgements**

We thank Olly Berry and Andrew Young for their leadership within the Environomics Future Science Platform. We thank the director of the Australian National Wildlife Collection, Leo Joseph, and the ANWC staff (specifically, Margaret Cawsey, Alex Drew, Tonya Haff, Dave Spratt and Chris Wilson) for their contributions of curatorial expertise, metadata management and sampling assistance. We thank Kerensa McElroy for her assistance and guidance in data

36

# References

735

736   1.   Shaffer HB, Fisher RN, Davidson C. The role of natural history collections in documenting
737        species declines. Trends Ecol Evol. 1998;13(1):27–30.

738   2.   Meineke EK, Davies TJ, Daru BH, Davis CC. Biological collections for understanding
739        biodiversity in the Anthropocene. Philos Trans R Soc Lond B Biol Sci. 2018;
740        doi:10.1098/rstb.2017.0386

741   3.   Holmes MW, Hammond TT, Wogan GOU, Walsh RE, LaBarbera K, Wommack EA,
742        Martins FM, Crawford JC, Mack KL, Bloch LM, Nachman MW. Natural history
743        collections as windows on evolutionary processes. Mol Ecol. 2016;25(4):864–81.

744   4.   Martínková N, Searle JB. Amplification success rate of DNA from museum skin
745        collections: a case study of stoats from 18 museums. Mol Ecol Notes. 2006;6(4):1014–7.

746   5.   Rawlence NJ, Wood JR, Armstrong KN, Cooper A. DNA content and distribution in
747        ancient feathers and potential to reconstruct the plumage of extinct avian taxa. Proc Biol
748        Sci. 2009;276(1672):3395–402.

749   6.   Sefc KM, Payne RB, Sorenson MD. Microsatellite amplification from museum feather
750        samples: effects of fragment size and template concentration on genotyping errors. Auk.
751        2003;120(4):982–9.

752   7.   Grealy A, Bunce M, Holleley CE. Avian mitochondrial genomes retrieved from museum
753        eggshell. Mol Ecol Resour. 2019;19(4):1052–62.

754   8.   Grealy A, Langmore NE, Joseph L, Holleley CE. Genetic barcoding of museum eggshell
755        improves data integrity of avian biological collections. Sci Rep. 2021;11(1):1605.

756   9.   Tsai WLE, Schedl ME, Maley JM, McCormack JE. More than skin and bones: comparing
757        extraction methods and alternative sources of DNA from avian museum specimens. Mol
758        Ecol Resour. 2020;20(5):1220–7.

759   10.  McElroy K, Beattie K, Symonds MRE, Joseph L. Mitogenomic and nuclear diversity in
760        the Mulga Parrot of the Australian arid zone: cryptic subspecies and tests for selection.
761        Emu - Austral Ornithology. 2018;118(1):22–35.

762   11.  Morgan CC, Creevey CJ, O'Connell MJ. Mitochondrial data are not suitable for resolving
763        placental mammal phylogeny. Mamm Genome. 2014;25(11–12):636–47.

764   12.  Hackett SJ, Kimball RT, Reddy S, Bowie RCK, Braun EL, Braun MJ, et al. A
765        phylogenomic study of birds reveals their evolutionary history. Science.
766        2008;320(5884):1763–8.

767   13.  Derkarabetian S, Benavides LR. Sequence capture phylogenomics of historical ethanol-
768        preserved museum specimens: unlocking the rest of the vault. Mol Ecol. 2019; doi:
769        10.1111/1755-0998.13072.

770   14.  Linck EB, Hanna ZR, Sellas A, Dumbacher JP. Evaluating hybridization capture with
771        RAD probes as a tool for museum genomics with historical bird specimens. Ecol Evol.
772        2017;7(13):4755–67.

773  15. Wood HM, González VL, Lloyd M, Coddington J, Scharff N. Next-generation museum
774      genomics: phylogenetic relationships among palpimanoid spiders using sequence capture
775      techniques (Araneae: Palpimanoidea). Mol Phylogenet Evol. 2018;127:907–18.

776  16. Parejo M, Wragg D, Henriques D, Charrière J-D, Estonba A. Digging into the genomic
777      past of Swiss honey bees by whole-genome sequencing museum specimens. Genome Biol
778      Evol. 2020;12(12):2535–51.

779  17. Staats M, Erkens RHJ, van de Vossenberg B, Wieringa JJ, Kraaijeveld K, Stielow B, Geml
780      J, Richardson JE, Bakker FT. Genomic treasure troves: complete genome sequencing of
781      herbarium and insect museum specimens. PLoS One. 2013;8(7):e69189.

782  18. Feigin CY, Newton AH, Doronina L, Schmitz J, Hipsley CA, Mitchell KJ, Gower G,
783      Llamas B, Soubrier J, Heider TN, Menzies BR, Cooper A, O'Neill RJ, Pask AJ. Genome
784      of the Tasmanian tiger provides insights into the evolution and demography of an extinct
785      marsupial carnivore. Nat Ecol Evol. 2018;2(1):182–92.

786  19. Atlas of Living Australia. https://www.ala.org.au/. Accessed 28 January 2021.

787  20. Appleyard SA, Maher S, Pogonoski JP, Bent SJ, Chua X-Y, McGrath A. Assessing DNA
788      for fish identifications from reference collections: the good, bad and ugly shed light on
789      formalin fixation and sequencing approaches. J Fish Biol. 2021; doi:10.1111/jfb.14687.

790  21. Srinivasan M, Sedmak D, Jewell S. Effect of fixatives and tissue processing on the content
791      and integrity of nucleic acids. Am J Pathol. 2002;161(6):1961–71.

792  22. Williams C, Ponten F, Moberg C, Soderkvist P, Uhlen M, Ponten J, Sitbon G, Lundeberg
793      J. A high frequency of sequence alterations is due to formalin fixation of archival
794      specimens. Am J Pathol. 1999;155(5):1467–71.

795  23. Do H, Dobrovic A. Sequence artifacts in DNA from formalin-fixed tissues: causes and
796      strategies for minimization. Clin Chem. 2015;61(1):64–71.

797  24. Burrell AS, Disotell TR, Bergey CM. The use of museum specimens with high-throughput
798      DNA sequencers. J Hum Evol. 2015;79:35–44.

799  25. Gansauge M-T, Gerber T, Glocke I, Korlevic P, Lippik L, Nagel S, Riehl LM, Schmidt A,
800      Meyer M. Single-stranded DNA library preparation from highly degraded DNA using T4
801      DNA ligase. Nucleic Acids Res. 2017;45(10):e79.

802  26. Carøe C, Gopalakrishnan S, Vinner L, Mak SST, Sinding MHS, Samaniego JA, Wales N,
803      Sicheritz-Pontén, Gilbert MTP. Single-tube library preparation for degraded DNA.
804      Methods in Ecology and Evolution. 2018;9(2):410–419.

805  27. Parks M, Lambert D. Impacts of low coverage depths and post-mortem DNA damage on
806      variant calling: a simulation study. BMC Genomics. 2015;16:19.

807  28. Robbe P, Popitsch N, Knight SJL, Antoniou P, Becq J, He M, et al. Clinical whole-genome
808      sequencing from routine formalin-fixed, paraffin-embedded specimens: pilot study for the
809      100,000 Genomes Project. Genet Med. 2018; doi:10.1038/gim.2017.241.

810  29. Simmons JE. Fluid preservation : a comprehensive reference. Lanham: Rowman &
811      Littlefield; 2014.

812   30.  MacLeod ID. Washing formaldehyde from fixed spirit specimens: a mechanism for the
813        preservation of Megamouth III. AICCM Bulletin. 2008;31(1):36–43.

814   31.  Koshiba M, Ogawa K, Hamazaki S, Sugiyama T, Ogawa O, Kitajima T. The effect of
815        formalin fixation on DNA and the extraction of high-molecular-weight DNA from fixed
816        and embedded tissues. Pathol Res Pract. 1993;189(1):66–72.

817   32.  Bär W, Kratzer A, Mächler M, Schmid W. Postmortem stability of DNA. Forensic Sci Int.
818        1988;39(1):59–70.

819   33.  Palmer ADN. DNA isolation and amplification from formaldehyde-fixed animal tissues
820        rich in mucopolysaccharides, pigments, and chitin. Prep Biochem Biotechnol.
821        2009;39(1):72–80.

822   34.  Bibi SS, Rehman A, Minhas RA, Janjua S. Evaluation of DNA extraction method from
823        formalin preserved skin samples of *Panthera pardus* for molecular genetic assessment.
824        The Journal of Animal & Plant Sciences. 2015;25:1196–9.

825   35.  Lutterschmidt WI, Cureton JC Ii, Gaillard AR. "Quick" DNA extraction from claw
826        clippings of fresh and formalin-fixed box turtle (*Terrapene ornata*) specimens. Herpetol
827        Rev. 2010;41(3):313–5.

828   36.  Scatena MP, Morielle-Versute E. Suitability of DNA extracted from archival specimens
829        of fruit-eating bats of the genus Artibeus (Chiroptera, Phyllostomidae) for polymerase
830        chain reaction and sequencing analysis. Genet Mol Biol. 2008;31(1):160–5.

831   37.  Joshi BD, Mishra S, Singh SK, Goyal SP. An effective method for extraction and
832        polymerase chain reaction (PCR) amplification of DNA from formalin preserved tissue
833        samples of snow leopard. Afr J Biotechnol. 2013;12(22):3399–404.

834   38.  Shedlock AM, Haygood MG, Pietsch TW, Bentzen P. Enhanced DNA extraction and PCR
835        amplification of mitochondrial genes from formalin-fixed museum specimens.
836        Biotechniques. 1997;22(3):394–400.

837   39.  Boyle EE, Zardus JD, Chase MR, Etter RJ, Rex MA. Strategies for molecular genetic
838        studies of preserved deep-sea macrofauna. Deep Sea Res Part I. 2004;51(10):1319–36.

839   40.  Hykin SM, Bi K, McGuire JA. Fixing formalin: a method to recover genomic-scale DNA
840        sequence data from formalin-fixed museum specimens using high-throughput sequencing.
841        PLoS One. 2015;10(10):e0141579.

842   41.  Gould AL, Fritts-Penniman A, Gaisiner A. Museum genomics illuminate the high
843        specificity of a bioluminescent symbiosis across a genus of reef fish. Front Ecol Evol.
844        2021; doi:10.3389/fevo.2021.630207.

845   42.  Ruane S, Austin CC. Phylogenomics using formalin-fixed and 100+ year-old intractable
846        natural history specimens. Mol Ecol Resour. 2017;17(5):1003–8.

847   43.  Kehlmaier C, Zinenko O, Fritz U. The enigmatic Crimean green lizard (*Lacerta viridis
848        magnifica*) is extinct but not valid: mitogenomics of a 120-year-old museum specimen
849        reveals historical introduction. J Zoolog Syst Evol Res. 2020;58(1)303–7.

850   44. Diaz-Viloria N, Sanchez-Velasco L, Perez-Enriquez R. Inhibition of DNA amplification
851       in marine fish larvae preserved in formalin. J Plankton Res. 2005;27(8):787–92.

852   45. Pierson TW, Kieran TJ, Clause AG, Castleberry NL. Preservation-induced morphological
853       change in salamanders and failed DNA extraction from a decades-old museum specimen:
854       implications for *Plethodon ainsworthi*. J Herpetol. 2020;54(2):137–43.

855   46. McGaughran A. Effects of sample age on data quality from targeted sequencing of
856       museum specimens: what are we capturing in time? BMC Genomics. 2020;21(1):188.

857   47. Watanabe M, Hashida S, Yamamoto H, Matsubara T, Ohtsuka T, Suzawa K, Maki Y, Soh
858       J, Asano H, Tsukuda K, Toyooka S, Miyoshi S. Estimation of age-related DNA
859       degradation from formalin-fixed and paraffin-embedded tissue according to the extraction
860       methods. Exp Ther Med. 2017;14(3):2683–8.

861   48. Zimmermann J, Hajibabaei M, Blackburn DC, Hanken J, Cantin E, Posfai J, Evans TC.
862       DNA damage in preserved specimens and tissue samples: a molecular assessment. Front
863       Zool. 2008;5:18.

864   49. Freedman J, van Dorp LB, Brace S. Destructive sampling natural science collections: an
865       overview for museum professionals and researchers. Journal of Natural Science
866       Collections. 2018;5:21–34.

867   50. Shi S-R, Datar R, Liu C, Wu L, Zhang Z, Cote RJ, Taylor CR. DNA extraction from
868       archival formalin-fixed, paraffin-embedded tissues: heat-induced retrieval in alkaline
869       solution. Histochem Cell Biol. 2004;122(3):211–8.

870   51. Shi S-R, Cote RJ, Wu L, Liu C, Datar R, Shi Y, Liu D, Lim H, Taylor CR. DNA extraction
871       from archival formalin-fixed, paraffin-embedded tissue sections based on the antigen
872       retrieval principle: heating under the influence of pH. J Histochem Cytochem.
873       2002;50(8):1005–11.

874   52. Campos PF, Gilbert TMP. DNA extraction from formalin-fixed material. Ancient DNA:
875       Methods and Protocols. 2012;840:81–5.

876   53. Splendiani A, Fioravanti T, Giovannotti M, Olivieri L, Ruggeri P, Nisi Cerioni P, Vanni
877       S, Enrichetti F, Caputo Barucchi V. Museum samples could help to reconstruct the original
878       distribution of *Salmo trutta* complex in Italy. J Fish Biol. 2017;90(6):2443–51.

879   54. Paireder S, Werner B, Bailer J, Werther W, Schmid E, Patzak B, Cichna-Markl M.
880       Comparison of protocols for DNA extraction from long-term preserved formalin fixed
881       tissues. Anal Biochem. 2013;439(2):152–60.

882   55. Prüfer K, Stenzel U, Hofreiter M, Pääbo S, Kelso J, Green RE. Computational challenges
883       in the analysis of ancient DNA. Genome Biol. 2010;11(5):R47.

884   56. McGuire JA, Cotoras DD, O'Connell B, Lawalata SZS, Wang-Claypool CY, Stubbs A,
885       Huang X, Wogan GOU, Hykin SM, Reilly SB, Bi K, Riyanto A, Arida E, Smith LL, Milne
886       H, Streicher JW, Iskandar DT. Squeezing water from a stone: high-throughput sequencing
887       from a 145-year old holotype resolves (barely) a cryptic species problem in flying lizards.
888       PeerJ. 2018;6:e4470.

889  57.  Johnson LA, Ferris JAJ. Analysis of postmortem DNA degradation by single-cell gel
890       electrophoresis. Forensic Sci Int. 2002;126(1):43–7.

891  58.  Itani M, Yamamoto Y, Doi Y, Miyaishi S. Quantitative analysis of DNA degradation in
892       the dead body. Acta Med Okayama. 2011;65(5):299–306.

893  59.  Pedersen JS, Valen E, Velazquez AMV, Parker BJ, Rasmussen M, Lindgreen S, Lilje B,
894       Tobin DJ, Kelly TK, Vang S, Andersson R, Jones PA, Hoover CA, Tikhonov A,
895       Prokhortchouk E, Rubin EM, Sandelin A, Gilbert MTP, Krogh A, Willerslev E, Orlando
896       L. Genome-wide nucleosome map and cytosine methylation levels of an ancient human
897       genome. Genome Res. 2014;24(3):454–66.

898  60.  Andrews S. FastQC: a quality control tool for high throughput sequence data.
899       https://www.bioinformatics.babraham.ac.uk/projects/fastqc/. Accessed 1 August 2019.

900  61.  Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. Genome
901       Biol. 2019;20(1):257.

902  62.  Picard. http://broadinstitute.github.io/picard/. Accessed 8 September 2020.

903  63.  Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, Shamim MS,
904       Machol I, Lander ES, Aiden AP, Aiden EL. De novo assembly of the *Aedes aegypti*
905       genome using Hi-C yields chromosome-length scaffolds. Science. 2017;356(6333):92–5.

906  64.  DNA Zoo. https://www.dnazoo.org. Accessed 7 August 2019.

907  65.  GenBank. https://www.ncbi.nlm.nih.gov/genbank/. Available from the National Center
908       for Biotechnology Information. Accessed 7 August 2019.

909  66.  Smit AFA, Hubley R, Green P. Repeat-Masker Open-4.0. http://www.repeatmasker.org.
910       Accessed 15 August 2019.

911  67.  Stephen S. kit4b. https://github.com/kit4b. Accessed 16 January 2020.

912  68.  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,
913       Durbin N, 1000 Genome Project Data Processing Subgroup. The sequence alignment/map
914       format and SAMtools. Bioinformatics. 2009;25(16):2078–9.

915  69.  R Core Team. R: A language and environment for statistical computing. https://wwwR-
916       project.org/. Accessed 19 September 2019.

917  70.  Wickham H. ggplot2. https://ggplot2.tidyverse.org/. Accessed 20 November 2020.

918  71.  Kassambara A. ggpubr. https://github.com/kassambara/ggpubr/. Accessed 20 November
919       2020.