

# 1 **Historical RNA expression profiles from the extinct Tasmanian tiger**

2 Emilio Mármol-Sánchez<sup>1,2\*</sup>, Bastian Fromm<sup>1,3</sup>, Nikolay Oskolkov<sup>4</sup>, Zoé Pochon<sup>2,5</sup>, Panagiotis  
3 Kalogeropoulos<sup>1</sup>, Eli Eriksson<sup>1</sup>, Inna Biryukova<sup>1</sup>, Vaishnovi Sekar<sup>1</sup>, Erik Ersmark<sup>2,6</sup>, Björn  
4 Andersson<sup>7</sup>, Love Dalén<sup>2,6,8#\*</sup> and Marc R. Friedländer<sup>1#\*</sup>

5

6 <sup>1</sup>Department of Molecular Biosciences, The Wenner-Gren Institute, Science for Life Laboratory,  
7 Stockholm University, Stockholm, Sweden. <sup>2</sup>Centre for Palaeogenetics, Stockholm, Sweden.  
8 <sup>3</sup>The Arctic University Museum of Norway, UiT, The Arctic University of Norway, Tromsø,  
9 Norway. <sup>4</sup>Department of Biology, National Bioinformatics Infrastructure Sweden, Science for  
10 Life Laboratory, Lund University, Lund, Sweden. <sup>5</sup>Department of Archaeology and Classical  
11 Studies, Stockholm University, Stockholm, Sweden. <sup>6</sup>Department of Bioinformatics and  
12 Genetics, Swedish Museum of Natural History, Stockholm, Sweden. <sup>7</sup>Department of Cell and  
13 Molecular Biology, Karolinska Institute, Stockholm, Sweden. <sup>8</sup>Department of Zoology,  
14 Stockholm University, Stockholm, Sweden. #LD and MRF contributed equally.

15

## 16 **\*Corresponding Authors:**

17 Marc R. Friedländer: [marc.friedlander@scilifelab.se](mailto:marc.friedlander@scilifelab.se)

18 Love Dalén: [love.dalen@zoologi.su.se](mailto:love.dalen@zoologi.su.se)

19 Emilio Mármol-Sánchez: [emilio.marmol.sanchez@gmail.com](mailto:emilio.marmol.sanchez@gmail.com)

20

21 **Running Title:** Palaeotranscriptomics in extinct species

22

23

24

25

## 26 **Abstract**

27 Palaeogenomics continues to yield valuable insights into the evolution, population dynamics, and  
28 ecology of our ancestors and other extinct species. However, DNA sequencing cannot reveal  
29 tissue-specific gene expression, cellular identity, or gene regulation, only attainable at the  
30 transcriptional level. Pioneering studies have shown that useful RNA can be extracted from  
31 ancient specimens preserved in permafrost and historical skins from extant canids, but no attempts  
32 have been made so far on extinct species.

33 We extract, sequence and analyze historical RNA from muscle and skin tissue of a ~130-year-old  
34 Tasmanian tiger (*Thylacinus cynocephalus*) preserved in desiccation at room temperature in a  
35 museum collection. The transcriptional profiles closely resemble those of extant species,  
36 revealing specific anatomical features such as slow muscle fibers or blood infiltration.  
37 Metatranscriptomic analysis, RNA damage, tissue-specific RNA profiles, and expression hotspots  
38 genome-wide further confirm the thylacine origin of the sequences. RNA sequences are used to  
39 improve protein-coding and noncoding annotations, evidencing missing exonic loci and the  
40 location of ribosomal RNA genes, while increasing the number of annotated thylacine  
41 microRNAs from 62 to 325. We discover a thylacine-specific microRNA isoform that could not  
42 have been confirmed without RNA evidence. Finally, we detect traces of RNA viruses, suggesting  
43 the possibility of profiling viral evolution.

44 Our results represent the first successful attempt to obtain transcriptional profiles from an extinct  
45 animal species, providing thought-to-be-lost information on gene expression dynamics. These  
46 findings hold promising implications for the study of RNA molecules across the vast collections  
47 of Natural History museums and from well-preserved permafrost remains.

48

49

50

51

## 52 **Introduction**

53 Over the past decade, high-throughput sequencing techniques have propelled the analysis of  
54 ancient DNA (aDNA) molecules, enabling the study of genomes from extinct or extant species  
55 that lived up to around two million years ago (van der Valk et al. 2021; Kjær et al. 2022), or in  
56 more recent times (Feigin et al. 2022). Studies on aDNA, alongside ancient proteins to lesser  
57 extent, have facilitated the exploration of evolutionary processes by simply examining the  
58 snapshot of time that palaeogenomes and palaeoproteomes can provide. This has allowed the  
59 reconstruction of genomes and ancestral lineages from multiple extinct species from the  
60 Pleistocene era, including Neanderthals (Green et al. 2010), woolly mammoths (Palkopoulou et  
61 al. 2015), and woolly rhinoceros (Lord et al. 2020), as well as from others that disappeared more  
62 recently, such as the quagga (Vilstrup et al. 2013; Jónsson et al. 2014) and the Tasmanian tiger  
63 (Feigin et al. 2022, 2017).

64 Aside from the well-established field of palaeogenomics, and the emerging field of  
65 palaeoproteomics (Hendy et al. 2018), the analysis of ancient RNA, a key molecule between both  
66 DNA and proteins across the central dogma of life, remains elusive. Unlike DNA, RNA provides  
67 researchers with additional layers of information so far unexplored in extinct species, such as cell  
68 and tissue identity, gene regulatory mechanisms, and evidence of the expression of coding and  
69 noncoding loci. Indeed, studies focused on ancient and/or historical RNA molecules have not  
70 experienced similar advancements witnessed in ancient DNA and proteins (Smith and Gilbert  
71 2018). While a few early and controversial studies indicated the potential presence of RNA  
72 sequences in ancient plant seeds (Rollo 1985; Venanzi and Rollo 1990; Rollo et al. 1991) and ice  
73 cores (Castello et al. 1999; Zhang et al. 2006), subsequent research reported the recovery of partial  
74 and complete genomes from RNA viruses preserved in seeds (Guy 2013; Smith et al. 2014),  
75 faeces (Ng et al. 2014) and formalin-fixed tissues (Xiao et al. 2013; Worobey et al. 2016; Gryseels  
76 et al. 2020; Patrono et al. 2022), as well as partial transcriptomes from plants (Fordyce et al.  
77 2013). However, it was not until 2017 that the first example of metazoan ancient RNAs were  
78 detected using qPCR-based methods in mummified cold-preserved remains of a human dating

79 back over 5,000 years (Keller et al. 2017), and later reproduced through sequencing techniques in  
80 humans from medieval times (Shaw et al. 2019). More recently, two additional studies have  
81 employed sequencing techniques to recover ancient RNA profiles, including messenger RNA  
82 (mRNA) and microRNA (miRNA) molecules from preserved tissues of a Late Pleistocene canid  
83 and historical wolf skins (Smith et al. 2019; Fromm et al. 2021). The presence of endogenous  
84 RNA sequences in extremely well-preserved yet ancient specimens and in historical wolf skins  
85 demonstrated that, under favorable conditions, RNA molecules could be preserved to the extent  
86 of still representing their abundance in the once-living cells of origin. Albeit promising results  
87 from both early and recent studies on sequencing and analyzing RNA in metazoan specimens of  
88 considerable age, there is currently no example of applying a palaeotranscriptomics approach to  
89 extinct metazoan species.

90 To address this gap in the emerging field of palaeotranscriptomics, we have focused on the  
91 renowned and recently extinct Tasmanian tiger (*T. cynocephalus*), also referred to as the  
92 thylacine. Thylacines were the largest carnivorous marsupials across the Holocene (Mitchell et  
93 al. 2014), and represented the only surviving species of the Thylacinidae family to survive into  
94 the modern era. They belonged to the order Dasyuromorphia, and were closely related to the  
95 extant families Dasyuridae (including Tasmanian devils, quolls, phascogales and dunnarts, among  
96 others) and Myrmecobiidae (numbats) (Feigin et al. 2017; Miller et al. 2009). These apex  
97 marsupials were once widespread all across the Australian region but eventually became restricted  
98 to an isolated population on the island of Tasmania approximately 3000 years ago (Paddle 2000).  
99 Wild thylacines persisted in Tasmania until the early twentieth century, when European colonizers  
100 classified them as an agricultural pest, and aggressively targeted their remaining populations,  
101 leading to their complete extinction. The last known thylacine died in captivity in 1936 at the  
102 Beaumaris Zoo in Hobart, Tasmania. Thylacines are particularly important because they  
103 exemplify both a recent human-driven extinction event and an evident case of convergent  
104 evolution (Newton et al. 2021; Rovinsky et al. 2021).

105 Despite diverging from placental carnivorous mammals approximately 160 million years ago  
106 (Bininda-Emonds et al. 2007), thylacines exhibited striking phenotypic similarities with extant  
107 species like those belonging to the Canidae family. This illustrates how species with distinct  
108 evolutionary relationships can undergo common selective pressures, resulting in shared  
109 adaptations (Losos 2011). Previous studies have utilized mitochondrial DNA to determine the  
110 evolutionary position of thylacines among marsupial mammals (Miller et al. 2009), and have  
111 explored their demographic history and genetic diversity by sequencing and assembling their  
112 nuclear genome (Feigin et al. 2022, 2017). This invaluable information has provided researchers  
113 with unprecedented insights into the biology of this species.

114 In this study, we present the thylacine as a proof-of-concept for expanding the field of  
115 palaeotranscriptomics into the analysis of historical RNA remains in extinct species for the first  
116 time.

117

118

## 119 **Results**

### 120 **Recovery, sequencing, and genome-wide mapping of RNA fragments**

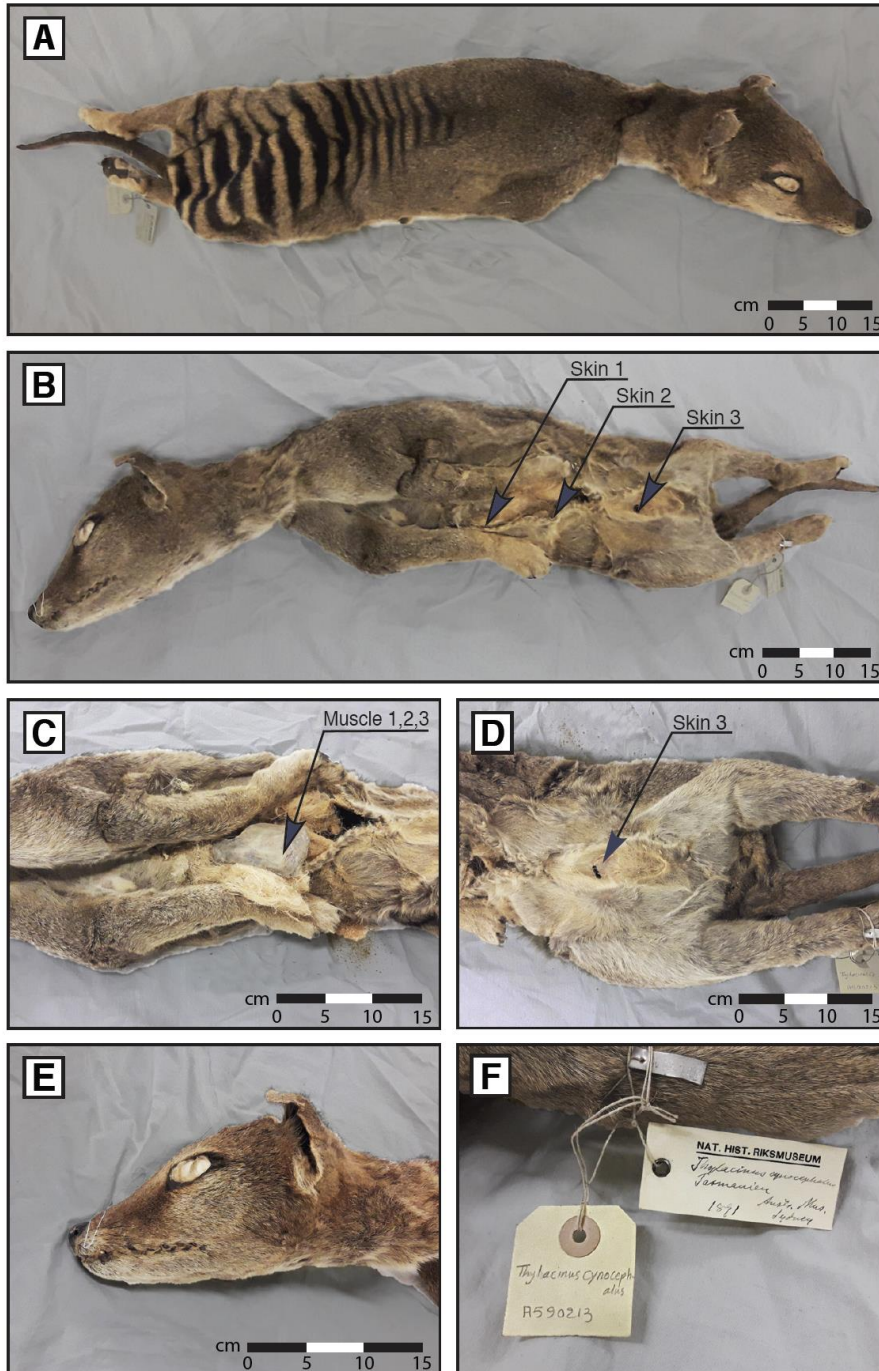
121 We first set out to investigate if it is possible to recover useful RNA molecules from a desiccated  
122 thylacine specimen stored at room temperature without specific preservative conditions. Three  
123 independent samples of both skeletal muscle and skin tissues were obtained by biopsy from a  
124 thylacine specimen available at the Stockholm Natural History Museum (NRM-MA590213, **Fig.**  
125 **1**).

126

127

128

129



130

131 **Figure 1: Thylacine specimen NRM-MA590213. (A)** Dorsal view. **(B)** Ventral view and skin  
 132 sampling areas. **(C)** Ventral view and skeletal muscle sampling area from the inner side of the left  
 133 scapula. **(D)** Inguinal region. **(E)** Head view. **(F)** Museum identification.

134

135

136 Considering the age and preservation status of the specimen, we anticipated a significant  
137 fragmentation in the historical RNA sequences possibly present within the tissue matrix.  
138 Therefore, we used an RNA extraction protocol targeting small RNA molecules, specifically  
139 designed for microRNA sequencing, on each of the six tissue samples obtained (see Methods).  
140 The samples were ground in liquid nitrogen and incubated in a digestion buffer to homogenize  
141 keratinous hard fibrous tissues (Gilbert et al. 2007; Sinding et al. 2015), while minimizing the  
142 incubation time to maximize the RNA extraction yield. From approximately 80 mg of tissue per  
143 sample, we obtained variable but substantial amounts of total RNA (**Supplemental Table 1**).  
144 Subsequently, the extracted and purified RNA fragments were prepared for high-throughput  
145 sequencing using a cDNA library protocol tailored for short RNA transcripts. The library size  
146 distribution indicated a successful extraction and library preparation, with an overall length of  
147 150 base pairs (bp) (**Supplemental Fig. 1**). We sequenced the cDNA libraries using an Illumina  
148 NextSeq 500 instrument, generating between 81.9 and 223.6 million raw sequencing reads per  
149 sample (**Supplemental Table 2**). A computational workflow for processing the RNA data is  
150 illustrated in **Fig. 2**, and will be described in subsequent sections. Initially, we trimmed the reads  
151 to remove artificial sequencing adapter sequences. Approximately 96% and 94.5% of reads had  
152 successful adapter detection and were trimmed (**Supplemental Table 2**), while the remaining  
153 ~5% were kept untrimmed, potentially originating from long RNA transcripts beyond the small  
154 RNA sequencing window employed. Trimmed sequences shorter than 18 nucleotides (nt) were  
155 discarded as they were deemed too short for reliable mapping to reference genomes. This is based  
156 on the smallest reported size for a functional microRNA transcript (~20 nt), while allowing up to  
157 2 nt loss in their 3' overhangs (Bartel 2018), and avoiding shorter RNA fragments that could  
158 probably lead to an excess of unwanted spurious mapping. This step eliminated around 30% of  
159 the trimmed reads, indicating the presence of highly degraded RNA molecules (**Supplemental**  
160 **Table 2**). PCR duplicates originating from the same RNA molecules were identified and  
161 deduplicated based on identical sequences and unique molecular identifiers (UMIs). This  
162 procedure reduced the number of sequences to 2.6-12 million per sample, indicating a PCR  
163 duplication rate of approximately 21.3 and 11 for skeletal muscle and skin tissues, respectively

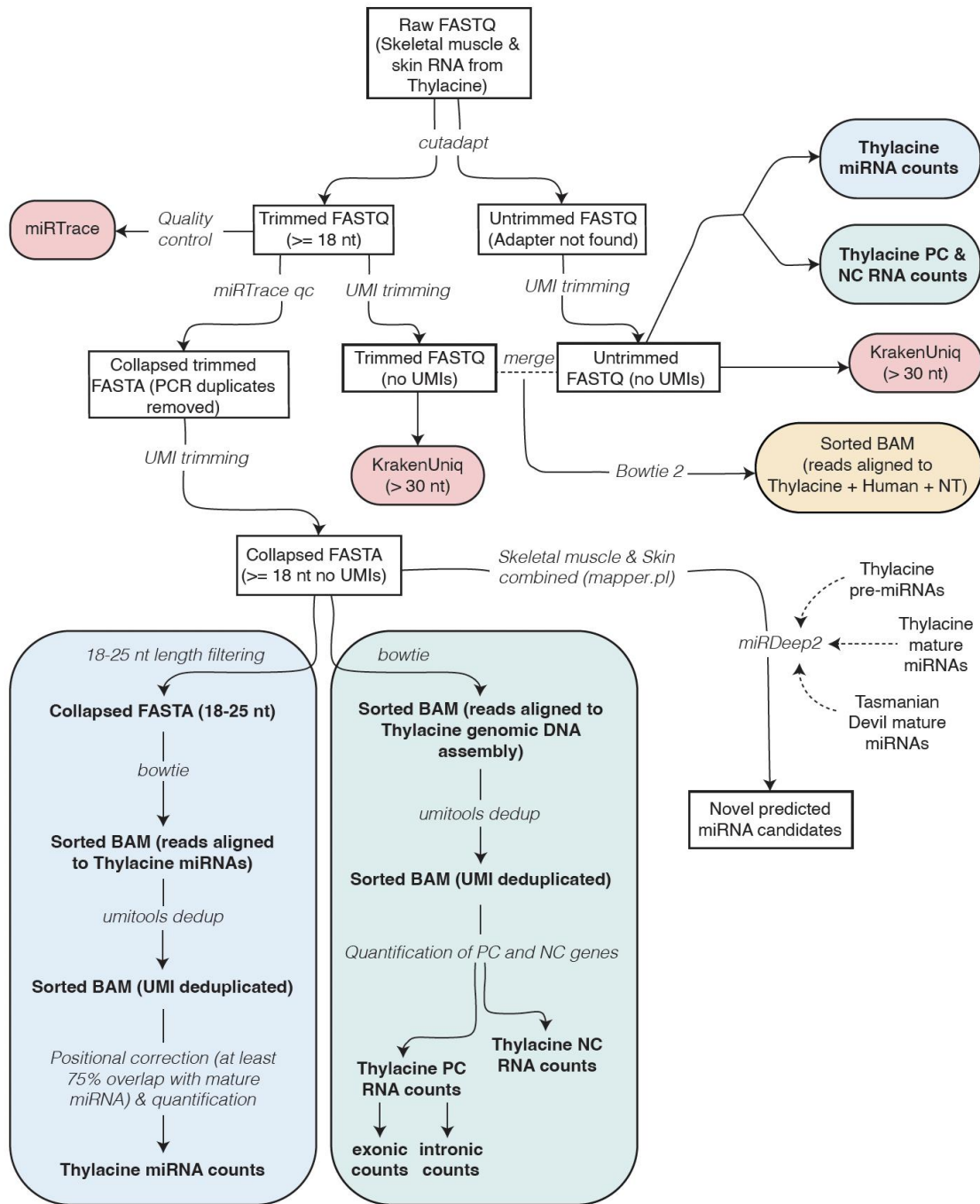
164 (**Supplemental Table 2**). Untrimmed reads exhibited a lower PCR duplication rate of ~2.5 and  
165 1.8 for skeletal muscle and skin (**Supplemental Table 3**), suggesting increased sequence  
166 variability compared to that of trimmed reads. The PCR-deduplicated and trimmed reads were  
167 then mapped to the thylacine nucleic and mitochondrial genomes, resulting in an overall  
168 successful mapping rate of 62.5% in skeletal muscle and 63.3% in skin tissues (**Supplemental**  
169 **Table 2**).

170 The read length distribution of mapped trimmed reads showed a predominance of short sequences  
171 below 30 nt, indicative of time-dependent fragmentation of RNA transcripts from thylacine  
172 origin. This pattern was more pronounced in skeletal muscle compared to skin samples  
173 (**Supplemental Fig. 2A-B**). A small increase in reads ranging from 28-35 nt and long reads of 42  
174 nt was also observed (**Supplemental Fig. 2A-B**). After UMI-based deduplication, we obtained  
175 ~1.5 and 2.8 million mapped trimmed reads for skeletal muscle and skin tissues, with an average  
176 UMI deduplication rate of 3.7 and 6.2, respectively (**Supplemental Table 2**). Most reads were  
177 short, although sample 2 from skeletal muscle and sample 1 from skin exhibited more abundant  
178 short reads, indicative of a greater degradation (**Supplemental Fig. 2C-D**). In contrast,  
179 approximately 4.5% and 7% of the PCR-deduplicated untrimmed reads were successfully mapped  
180 to the thylacine nuclear and mitochondrial genomes in skeletal muscle and skin tissues,  
181 respectively, with an average UMI deduplication rate of 3.25 and 9.1 (**Supplemental Table 3**).  
182 This represents ~11-fold decrease in successfully mapped long untrimmed reads compared to  
183 short trimmed reads with a similar UMI deduplication rate. Subsequent analyses will primarily  
184 focus on short trimmed RNA reads, unless stated otherwise. A detailed analysis of the sequenced  
185 RNAs using the miRTrace tool on trimmed and untrimmed reads is available in **Supplemental**  
186 **Files 1 and 2**, respectively. In summary, we produced millions of stringently quality-controlled  
187 sequences from thylacine tissue biopsies.

188

189

190



191

192 **Figure 2: Pre-processing, mapping, metatranscriptomics, and annotation pipeline of skin**

193 **and skeletal muscle RNA sequences from the NRM-MA590213 thylacine specimen. NT: full**

194 **NCBI non-redundant reference nucleotide database (Pochon et al. 2022); PC: protein-coding; NC:**

195 **noncoding.**

196

## 197 **Historical RNA sequences show a characteristic damage pattern**

198 The observed damage profiles of RNA reads mapping to the thylacine genome were indicative of  
199 the historical nature of the tissues, displaying increased deamination and other nucleotide  
200 substitutions (**Supplemental Fig. 3**). This pattern was more prominent towards the end of the  
201 reads, whether they were short (18-25 nt), medium-sized (26-30 nt), or long (>30 nt), in  
202 accordance with previous evidence on RNA (Smith et al. 2014, 2019; Fromm et al. 2021) and  
203 DNA damage patterns (Dabney et al. 2013). Adenosine deamination to inosine (A>I, read as A>G  
204 by the sequencer) was generally less frequent than cytidine to uridine deamination (C>U, read as  
205 C>T by the sequencer), except at the 3' end of short RNA reads (**Supplemental Fig. 3**). However,  
206 distinguishing genuine time-dependent A>I deamination from technical misincorporation is  
207 challenging (Binladen et al. 2006; Gilbert et al. 2003). Other types of misincorporations were also  
208 prevalent and roughly followed deamination events in RNA reads (**Supplemental Fig. 3**), which  
209 warrants caution about the reliability of some of the observed damage. A similar pattern emerged  
210 when examining the sequence damage of untrimmed RNA reads (**Supplemental Fig. 4**), although  
211 samples 2 and 3 from skeletal muscle showed increased damage around the 18<sup>th</sup> nucleotide, likely  
212 due to the presence of shorter untrimmed reads with undetermined ('N') nucleotides. This  
213 suggests a more pronounced and widespread occurrence of time-dependent damage in RNA  
214 compared to DNA sequences, in agreement with the increased degradation susceptibility of RNA  
215 molecules relative to DNA. Furthermore, considering that our thylacine specimen has been stored  
216 at room temperature for over a century, DNA preservation might be compromised, and it is highly  
217 likely that the same holds true for RNA (Binladen et al. 2006).

218

## 219 **Metatranscriptomic analyses reveal that thylacine-like RNAs are predominant**

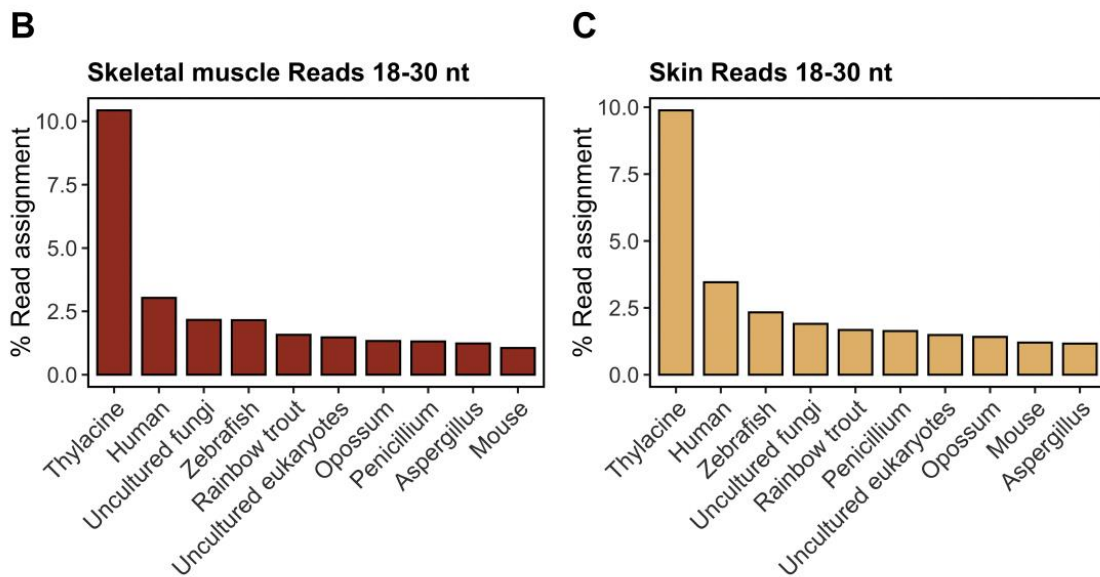
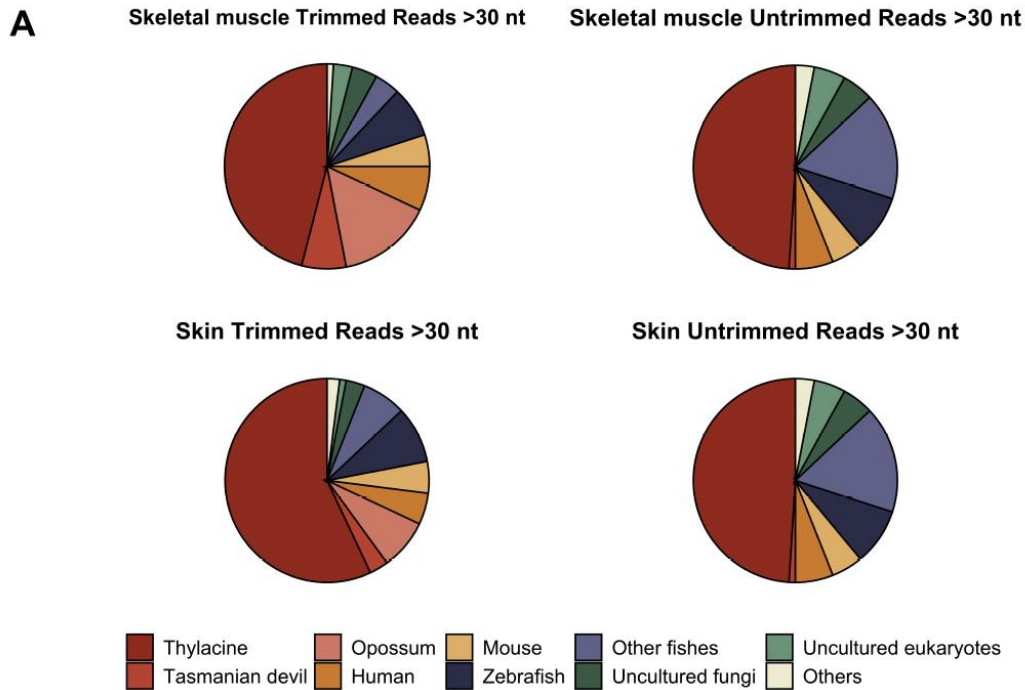
220 We aimed to determine the origins of the RNA sequences derived from thylacine tissues and  
221 distinguish endogenous reads from additional sources of contamination like other metazoans,  
222 microbes or any other microscopic life. To accomplish this, we employed a metatranscriptomic  
223 analysis based on the classification of long RNA reads (>30 nt) using KrakenUniq software

224 (Breitwieser et al. 2018) and the complete collection of sequenced organisms in the NCBI non-  
225 redundant reference nucleotide (NT) database. This analysis was performed separately on  
226 trimmed reads following sequencing adapter identification and on untrimmed reads where  
227 sequencing adapters were not detected and where hence left intact (see Methods).

228 In the case of skeletal muscle, approximately 46% of the reliably assigned sequences were  
229 attributed to thylacine contigs present in the NCBI NT database (**Fig. 3A; Supplemental File 3**).  
230 Additionally, 15% of the reads were assigned to opossum (*M. domestica*), and 7% to Tasmanian  
231 devil (*S. harrisi*), accounting for a total of 68% of reads assigned to these three marsupial species  
232 (**Fig. 3A; Supplemental File 3**). Given the phylogenetic proximity of opossums, Tasmanian  
233 devils, and thylacines as members of the metatheria clade (marsupials), it is plausible that all these  
234 sequences indeed originated from the thylacine. Moreover, the probability of cross-contamination  
235 from marsupial species other than the thylacine in our data is relatively low. The remaining  
236 assigned reads corresponded to human (7%), mouse (5%), undetermined fungi (4%), other  
237 eukaryotes (1%), as well as to zebrafish (8%) and other related fish species (4%). These findings  
238 were consistent with those obtained for untrimmed skeletal muscle reads (**Fig. 3A; Supplemental**  
239 **File 4**). The assignment obtained for skin tissue was similar to that of the muscle, although the  
240 fraction of sequences that could be reliably traced to the thylacine genome was higher for  
241 untrimmed reads (**Fig. 3A; Supplemental File 5-6**). In both skin and muscle samples, the  
242 proportion of untrimmed reads attributed to non-thylacine marsupial species was significantly  
243 smaller than that of shorter trimmed reads (**Fig. 3A**). This supports our previous hypothesis that  
244 read assignments to non-thylacine marsupial species likely represent misattributed short thylacine  
245 sequences. Consequently, reads assigned to other marsupial species might originate from  
246 alignments to highly conserved genomic loci across marsupials absent in the thylacine contigs  
247 used as reference. This outcome is expected since the thylacine sequences included in the NT  
248 database used by KrakenUniq represent only a limited portion of the entire thylacine genome.

249

250



251

252 **Figure 3: Metatranscriptomic analyses of thylacine skin and skeletal muscle samples. (A)**

253 Proportion of trimmed and untrimmed RNA reads (>30 nt) assigned using the KrakenUniq

254 pipeline from thylacine skeletal muscle and skin tissues. Only species with >1000 *k*-mers and

255 >200 species-specific reads are shown. Percentage of RNA reads aligned to the ten most abundant

256 species detected, focusing specifically on trimmed and untrimmed sequences uniquely mapped

257 (MAPQ  $\geq 1$ ) of size  $\geq 18$  and  $\leq 30$  nt in (B) skeletal muscle and (C) skin tissues. Reads with

258 multiple-species ambiguous mapping were discarded for calculating the percentage of read  
259 assignment. Global alignment was performed using Bowtie 2 with flags *--end-to-end* and *--very-*  
260 *sensitive*.

261

262 We then proceeded to analyze the shorter reads (18-30 nt) that could not be classified with  
263 KrakenUniq. Instead, we mapped them directly against a custom reference NT database using a  
264 dedicated short-read mapper (see Methods). This additional analysis is crucial since we anticipate  
265 that true historical RNAs would be fragmented, and therefore shorter in length. Consistent with  
266 our findings from the longer reads (>30 nt), the thylacine was the species with the highest  
267 proportion of uniquely traceable short reads, surpassing by approximately 3-fold the percentage  
268 of short reads unambiguously assigned to the most prevalent contaminant source detected,  
269 humans, in both skeletal muscle and skin tissues (**Fig. 3B-C**). This provides further support for  
270 the robustness of our analyses and confirms that the thylacine is indeed the primary source of our  
271 RNA sequencing data.

272 There were very few RNA reads assigned to prokaryotes or viral species, and their abundance  
273 profiles differed between muscle and skin tissues. The most abundant prokaryotic taxa were  
274 uncultured bacteria, followed by *Escherichia*, *Acinetobacter*, *Myroides*, *Streptomyces*,  
275 *Rheinheimera*, and *Corynebacterium*. These bacteria are typically associated with environmental  
276 contamination (**Supplemental Table 4**). Regarding viruses, the assigned RNA reads mostly  
277 belonged to fungi-specific viruses such as *Penicillium aurantiogriseum partitivirus 1* or *Primate*  
278 *T-lymphotropic virus 1*, potentially derived from human contamination (**Supplemental Table 5**).  
279 A small number of sequences from RNA viruses of unknown origin (picorna-like) were present,  
280 indicating that such viruses can be detected. However, these RNA viruses were shallowly  
281 supported by reads mapping to only few distinct positions in their genomes. Therefore, additional  
282 confirmation is required.

283 We also investigated an unexpected enrichment of reads assigned to fish-related species (**Fig.**  
284 **3A**), primarily to zebrafish (*D. rerio*). We discovered that most of the sequences successfully

285 aligned to the zebrafish genome (~80%) mapped to specific regions of its Chromosome 4, which  
286 contains numerous repetitive transfer RNAs (tRNAs, accounting for 50% of mapped reads) and  
287 ribosomal RNAs (rRNAs, the remaining 30%) (Howe et al. 2013). Due to the high sequence  
288 conservation of these loci across species, and the extensive presence of multiple fish species in  
289 the NT database used for *k*-mers classification, we believe that these sequences were falsely  
290 assigned to the zebrafish genome. Supporting this observation, we did not observe any enrichment  
291 of fish-related assigned reads in the skin compared to muscle tissue, as would be expected if there  
292 was a physical contamination from the environment, which is more likely to be present on the  
293 skin than deep within the muscle tissue. These findings highlight the need for caution regarding  
294 the reliability of non-endogenous contaminating sources detected by our metatranscriptomic  
295 pipeline.

296 One potential explanation for some of the reads assigned to non-thylacine genomes might be a  
297 technical bias in our NCBI NT database. Since humans, mice, and zebrafish are model organisms,  
298 their reference contigs may be over-represented. Consequently, a significant fraction of reads, not  
299 necessarily of thylacine origin, might be falsely assigned to these species by default. Supporting  
300 this hypothesis, when we realigned the reads classified as zebrafish, human, and mouse  
301 contamination to the thylacine, Tasmanian devil, and opossum assemblies, a considerable fraction  
302 of them were successfully mapped (average mapping rate of 85.74% to thylacine, 85.04% to  
303 opossum, and 70.71% to Tasmanian devil). Therefore, it would be premature to dismiss the  
304 possibility of an endogenous thylacine origin for at least some of the reads assigned to other non-  
305 marsupial species.

306 In summary, we consistently observe a predominance of sequences originating from the thylacine  
307 genome using different complementary computational approaches, thus supporting the  
308 authenticity of the historical RNA sequences.

309

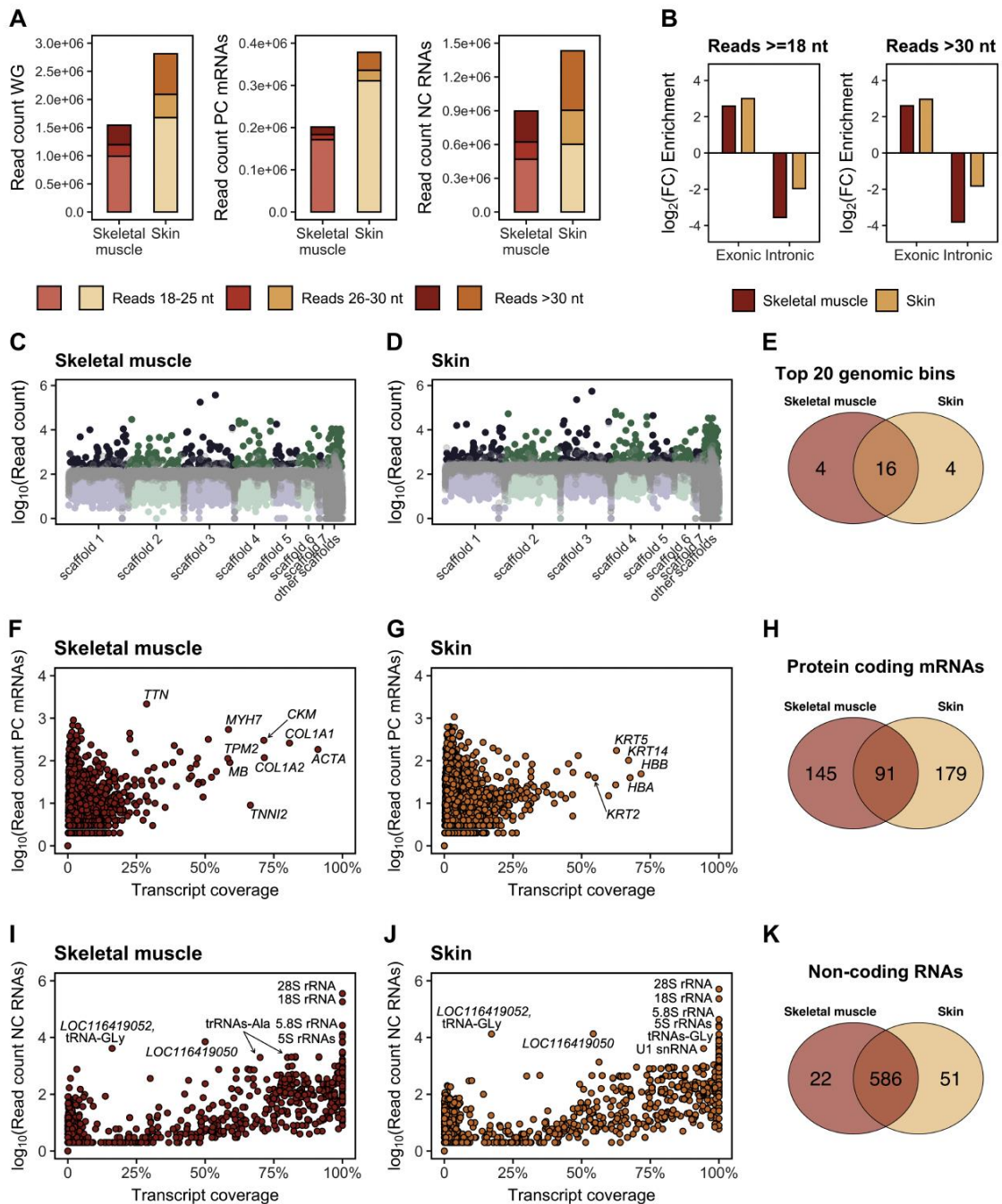
310 **Historical RNAs map to numerous classes of protein-coding and noncoding thylacine genes**

311 To understand the molecular nature of the sequenced RNAs, we traced them back to the nuclear  
312 and mitochondrial genomes of the thylacine (see Methods). Most reads mapped to noncoding  
313 RNA genes, while a smaller proportion mapped to protein-coding genes (**Fig. 4; Supplemental**  
314 **Table 6**), in agreement with the abundance of noncoding RNAs in the eukaryotic cell  
315 transcriptome. These noncoding RNAs primarily originated from highly expressed rRNA and  
316 tRNA loci. The majority of reads mapped to protein-coding genes were short, indicating age-  
317 related fragmentation (90% were under 30 nt), while fewer reads mapping to noncoding RNA  
318 genes were shorter than 30 nt (69% in muscle and 62% in skin).

319 The incorporation of reads above 30 nt is widely accepted as the lower threshold for aDNA data  
320 to avoid spurious alignments (de Filippo et al. 2018). Previous research (Smith et al. 2019) and  
321 our results demonstrate the usefulness of including ultrashort reads (18-30 nt) in obtaining  
322 comprehensive RNA abundance profiles. When mapping RNA reads to noncoding genes, the  
323 majority were assigned to rRNA genes (60% and 68% for skeletal muscle and skin tissues,  
324 respectively), followed by tRNA genes, with approximately half the hits. Other noncoding RNA  
325 genes, such as long noncoding RNAs (lncRNAs), small nucleolar RNAs (snoRNAs), and small  
326 nuclear RNAs (snRNAs) only accounted for 3% of the RNA reads in skeletal muscle and 7% in  
327 skin (**Supplemental Table 6**). This pattern was consistent across all samples except for skeletal  
328 muscle sample 3 and skin sample 1, which had a relatively low abundance of rRNAs  
329 (**Supplemental Table 6**). The length of mapped RNA reads varied among the three most  
330 abundant loci types (rRNAs, tRNAs and protein-coding mRNAs, **Supplemental Fig. 5**). Protein-  
331 coding mRNAs predominantly gathered ultrashort reads <25 nt in both skeletal muscle and skin  
332 tissues (**Supplemental Fig. 5A-B**), while rRNA loci exhibited a higher proportion of longer reads  
333 (**Supplemental Fig. 5C-D**). Reads ranging from 28-35 nt mapped to tRNA loci were particularly  
334 prevalent, especially in skeletal muscle 3 (**Supplemental Fig. 5E-F**). The latter is interesting  
335 since tRNA-derived small RNA fragments are known to have sizes compatible with such  
336 enrichment (Liu et al. 2021). Additionally, untrimmed RNA reads mapped to the thylacine  
337 assembly displayed a similar pattern to that of trimmed reads, although the overall number of

338 successfully mapped reads was significantly lower (~22-fold and 27-fold less untrimmed reads  
 339 compared to trimmed reads in skeletal muscle and skin, respectively, **Supplemental Table 7**).  
 340 Similar to trimmed reads, the dominant loci in untrimmed reads were rRNAs, accounting for an  
 341 average of ~74% of the mapped reads, followed by protein-coding mRNAs and tRNAs at a  
 342 considerable distance (**Supplemental Table 7**).

343



344

345 **Figure 4: Distribution of RNA sequences over protein-coding and noncoding genes.** (A) Read  
346 length distribution of RNA sequences mapped to the thylacine whole-genome (WG) assembly,  
347 annotated protein-coding (PC) genes (N = 19,356) and noncoding (NC) RNA genes (N = 3613).  
348 (B) Exonic enrichment and intronic depletion of RNA reads mapped to exonic and intronic  
349 regions of PC genes (at least 10% coverage) in skeletal muscle (N = 236) and skin (N = 270).  
350 Thylacine historical DNA reads mapped to the same PC genes were used as reference for  
351 comparison. (C, D) Number of RNA reads mapped to each consecutive 250 kbp window genome-  
352 wide in skeletal muscle and skin, respectively. Thylacine DNA reads (SRR5055304) mapped to  
353 each consecutive 250 kbp window genome-wide are in grey. I Venn diagram showing the top 20  
354 genomic windows (250 kbp) with the highest number of RNA reads mapped in skeletal muscle  
355 and skin. (F) Number of RNA reads mapped and coverage of each annotated PC gene (N =  
356 19,356) in skeletal muscle. (G) Number of RNA reads mapped and coverage of each annotated  
357 PC gene (N = 19,356) in skin. (H) Venn diagram showing PC genes quantified (at least 10%  
358 coverage) in skeletal muscle and skin. (I) Number of RNA reads mapped and coverage of each  
359 annotated NC RNA gene in skeletal muscle (N = 3613). (J) Number of RNA reads mapped and  
360 coverage of each annotated NC RNA gene in skin (N = 3613). (K) Venn diagram showing NC  
361 RNA genes quantified (at least 10% coverage, N = 608 for skeletal muscle and N = 637 for skin)  
362 in skeletal muscle and skin.

363

364 We further investigated an unexpected enrichment in long trimmed reads of 42 nt mapped to the  
365 thylacine assembly (~5% over all mapped trimmed reads, **Supplemental Fig. 5; Supplemental**  
366 **Table 8**) by cross-species comparison with the most probable source of exogenous modern  
367 contamination, humans. Approximately 55% of them mapped to thylacine rRNAs (**Supplemental**  
368 **Table 8**), which might explain, at least partially, the increased number of long reads >30 nt  
369 assigned to noncoding loci, as previously shown in **Fig. 4A**. We then compared the damage  
370 profiles of 42 nt reads that indistinctly mapped to the thylacine or the human assemblies (~57.5%,  
371 suspicious of human origin), with those that mapped to the thylacine but failed to align to the

372 human assembly (**Supplemental Table 8**). Indeed, we observed an increased C>U deamination  
373 profile for the 42 nt reads that mapped to the thylacine assembly but failed for the human, while  
374 those reads that mapped to both thylacine and human assemblies were less damaged  
375 (**Supplemental Fig. 6**). This supports our suspicion that a relevant proportion of these long  
376 enriched 42 nt sequences (up to around 60%) might have an exogenous origin (either human-  
377 derived or from another modern unknown contamination source).

378 In summary, we find that numerous RNA fragments can be reliably traced to thylacine protein-  
379 coding and noncoding genes.

380

### 381 **RNA sequences are enriched in exonic regions and span exon-exon junction**

382 We investigated whether RNA reads mapped to protein-coding genes were concentrated on  
383 exonic regions compared to intronic regions. This is important since potential DNA  
384 contamination traces could perturb our analyses when working with sparse amounts of highly  
385 fragmented historical RNA material. Sequencing mature mRNA molecules after splicing should  
386 show an enrichment of reads mapped to exonic regions and few reads mapped to intronic regions,  
387 as opposed to the even distribution expected for DNA sequencing data. The analysis of RNA  
388 reads mapped to protein-coding loci with reliable breadth of coverage (>10%) revealed 93% and  
389 77% of reads mapping to exonic regions in skeletal muscle and skin tissue, respectively. This  
390 finding agrees well with previous palaeotranscriptomic analyses in a Pleistocene canid (Smith et  
391 al. 2019), and indicates a significant enrichment in exonic reads compared to what was observed  
392 with thylacine DNA sequences (15% and 10% exonic reads in skeletal muscle and skin, **Fig. 4**;  
393 **Supplemental Table 9**). Conversely, RNA reads mapped to intronic regions showed a significant  
394 depletion in skeletal muscle (12-fold) and skin tissues (4-fold) compared to DNA reads mapped  
395 to introns (**Fig. 4B**).

396 We also aimed to investigate the presence of RNA reads spanning exon-exon junctions, indicative  
397 of mature intron-less mRNAs rather than nascent transcripts or DNA contamination.

398 Approximately 1% of the RNA reads mapped to exonic regions across all annotated protein-  
399 coding loci spanned exon-exon junctions (**Supplemental Table 10**). In contrast, only ~0.25% of  
400 the RNA reads mapped to intronic regions spanned exon-intron junctions, confirming that we  
401 primarily detected mature cytoplasmatic transcripts (**Supplemental Table 10**).

402 In summary, we provide evidence of the sequencing of RNA fragments from mature cytoplasmic  
403 mRNAs that are enriched in exonic sequences and span exon-exon junctions.

404

#### 405 **RNA sequences map unevenly and show evidence of unannotated loci**

406 Contrary to aDNA sequences, we expect that ancient/historical RNA fragments would map  
407 unevenly across the thylacine genome, representing variations in gene expression from distinct  
408 loci. We found that the breadth of coverage for the thylacine genome is 0.17% in skeletal muscle  
409 and 0.32% in skin. Moreover, the average depth genome-wide was ~0.008× for skeletal muscle  
410 and ~0.015× for skin (**Supplemental Table 11**). When analyzing thylacine DNA data (Feigin et  
411 al. 2017) at equivalent sequencing depth, the observed breadth of coverage genome-wide was  
412 78.35%, with an average depth of 3.45× (**Supplemental Table 11**).

413 To investigate the variation in transcriptional depth throughout the genome, we aggregated reads  
414 mapping to consecutive non-overlapping genomic windows of 250 kilobases (kbp). This analysis  
415 revealed several genome-wide expression hotspots (**Fig. 4C-D; Supplemental Fig. 7**), including  
416 two prominent ones in thylacine Scaffold 3 that consistently displayed higher read counts  
417 (**Supplemental Table 12**). Among the top twenty expression hotspots (**Supplemental Table 12**),  
418 16 (80%) were shared between skeletal muscle and skin tissues (**Fig. 4E**), indicating a common  
419 pattern of expression hotspots across the genome in both tissues. While some of these hotspots  
420 could be attributed to RNA reads mapping to the 340 annotated tRNAs or the 4 annotated 5S  
421 rRNAs in the thylacine assembly, the majority did not align with any annotated loci. The thylacine  
422 gene annotation currently includes only four 5S rRNA genes (**Supplemental Table 13**).  
423 Therefore, we speculated that the missing 18S, 28S, and 5.8S rRNA genes in the thylacine genome

424 might be located within the top two highlighted expression hotspots (**Fig. 4C-D**), potentially  
425 explaining the abundance of RNA reads mapped to these loci. Indeed, approximately 70% of the  
426 unexplained reads assigned to the expression hotspots in thylacine Scaffold 3 aligned with the  
427 Tasmanian devil 18S rRNA gene or the four human reference rRNA genes. Consequently, we  
428 determined the probable location of the missing 18S, 28S, and 5.8S rRNA genes in the thylacine  
429 genome (**Supplemental Table 13; Supplemental Fig. 8**).

430 We further investigated the causes and origin of RNA reads mapping to intergenic regions, which  
431 are not expected to produce transcriptional products. Approximately 29% of mapped reads in  
432 skeletal muscle and 38% in skin were assigned to intergenic regions (**Supplemental Table 6**). To  
433 assess the potential influence of historical thylacine DNA contamination in these mappings, we  
434 compared their distribution across genomic windows with thylacine DNA data (Feigin et al.  
435 2017). Reads mapping to intergenic regions revealed highly expressed unannotated loci,  
436 indicating the need for further annotation efforts in the thylacine assembly (**Supplemental Fig.**  
437 **9**). Reads mapping to genomic windows with breadth and depth of coverage similar to thylacine  
438 DNA reads, within a 2-fold difference (see Methods), were considered potential DNA  
439 contamination candidates. DNA reads are expected to be evenly distributed across the genome,  
440 resulting in an increased but roughly constant genome-wide breadth of coverage. In contrast, RNA  
441 reads should be concentrated in coding regions, leading to genomic hotspots with high read counts  
442 but lower genome-wide breadth of coverage compared to DNA data at an equivalent sequencing  
443 depth. After accounting for sequencing depth and read length distribution biases, around 3.2% of  
444 intergenic mapped reads in skeletal muscle and 2.1% in skin exhibited a distribution resembling  
445 that of thylacine DNA data. Thus, a limited presence of endogenous thylacine DNA  
446 contamination in our data should not be discarded. The remaining intergenic mapping reads may  
447 be attributed, at least partially, to spurious mapping caused by repetitive low-complexity genomic  
448 regions and/or unintended cross-species DNA/RNA contamination of unknown origin.

449

450 **Thylacine RNA expression profiles reflect cellular and tissue functions**

451 An advantage of RNA sequence data is that it can yield information about gene expression  
452 patterns in ancient/historical tissues, a feature that DNA alone cannot provide. We found that ~0.2  
453 million skeletal muscle RNA reads mapped to protein-coding genes, while the number was even  
454 higher for the skin sample, at ~0.4 million reads. Among the protein-coding genes quantified in  
455 muscle with reliable coverage (>10%, see Methods), several were characteristic of skeletal muscle  
456 tissue metabolism and structure (**Fig. 4F**; **Supplemental Table 14**). The most abundantly  
457 detected protein-coding transcript in the skeletal muscle was titin (*TTN*), which also showed low  
458 coverage (28.70%), and this was reproduced in both trimmed and untrimmed reads  
459 (**Supplemental Table 14**). These two observations can be explained by *TTN* being the gene with  
460 the longest known coding sequence, producing a giant protein largely abundant in striated muscle  
461 and mainly responsible for preventing overstretching of the sarcomere (Labeit and Kolmerer  
462 1995; Lee et al. 2007). To discard the presence of a confounding effect of long transcripts  
463 gathering a higher number of mapped reads leading to higher breadth of coverage, we investigated  
464 whether a positive correlation between transcript length and transcript coverage was present. We  
465 found a near-neutral correlation between these variables (**Supplemental Fig. 10**), indicating that  
466 transcript abundance is not solely influenced by sequence length. This observation dispels  
467 concerns about length-induced biases in transcript abundance and underscores the reliability of  
468 our findings.

469 Other highly expressed and well-covered transcripts in skeletal muscle included *LOC100913894*,  
470 which corresponds to a subunit of actin alpha (*ACTA*), *LOC100925998*, corresponding to myosin  
471 heavy chain beta (*MYH7*), myosin heavy chain 2 (*MYH2*), tropomyosin beta (*TPM2*), and  
472 troponin I1 and I2 (*TNNI1* and *TNNI2*), all of which form integral part of the sarcomere unit of  
473 striated muscle cells (Ahmed et al. 2022). Untrimmed reads also showed similar patterns, but with  
474 reduced resolution (**Supplemental Table 14**). *MYH7*, *TPM2*, and *TNNI1* transcripts represented  
475 the most abundant isoforms in their gene families, indicating the prevalence of type I slow muscle  
476 fibers (Bottinelli and Reggiani 2000). This is in agreement with the tentative functionality of the  
477 muscle fibers sequenced, which were probably sampled from the slow-twitched subscapularis

478 muscle according to the bone topology from where the muscle tissue was obtained. Creatine  
479 kinase M-type (*CKM*) and myoglobin (*MB*), important for muscle cell metabolism (Johnson et al.  
480 1989; Ordway and Garry 2004), were also highly abundant. Collagen alpha 1 and 2 transcripts  
481 (*COL1A1*, *COL1A2*), which contribute to the structure of tissues as type I collagen constituents  
482 (Henriksen and Karsdal 2019), were abundant as well.

483 In the skin, prominently expressed and reliably covered (>10%) loci included keratin-derived  
484 transcripts such as keratin type 1 cytoskeletal 14 (*KRT14*), and keratin type 2 cytoskeletal 2 and  
485 5 (*KRT2* and *KRT5*) consistent with the epithelial nature of the skin tissue (**Fig. 4G**;  
486 **Supplemental Table 15**). Actin gamma 1 (*ACTG1*), a cytoplasmic actin isoform expressed  
487 ubiquitously except in muscle tissues, also showed high coverage (**Supplemental Table 15**).  
488 Besides, hemoglobin alpha and beta subunits were detected, although predominantly in skin  
489 samples 2 and 3, suggesting possible blood infiltration in those samples. Comparing the reliably  
490 captured protein-coding genes between skeletal muscle (N = 236) and skin (N = 270) revealed a  
491 shared subset of only 91 genes (22%, **Fig. 4H**; **Supplemental Table 16**).

492 While tissue-specific protein-coding genes were prevalent, noncoding RNA genes displayed a  
493 different pattern. The most abundant noncoding RNA transcripts were rRNA genes in both  
494 skeletal muscle and skin tissues, followed by tRNAs and snRNAs (**Fig. 4I-J**; **Supplemental**  
495 **Tables 17-18**). Two highly expressed lncRNAs, *LOC116419050* and *LOC116419052*, likely  
496 contained reads mapping to overlapping tRNA loci. Among the reliably detected noncoding RNA  
497 genes, 586 (89%) were shared between the tissues. (**Fig. 4K**; **Supplemental Table 19**).

498 Regarding RNA reads mapping to the thylacine mitochondrial genome, the most abundant  
499 transcripts were the 16S and 12S mitochondrial rRNAs. These transcripts exhibited high breadth  
500 of coverage comparable to their nuclear rRNA counterparts in both muscle and skin tissues using  
501 trimmed and untrimmed reads (**Supplemental Tables 20-21**).

502

503 **Expanding the microRNA complement of the thylacine genome**

504 MicroRNAs (miRNAs) are short regulatory molecules involved in post-transcriptional gene  
505 expression repression (Bartel 2018). They play crucial roles in various biological processes,  
506 including development and cell identity. Our library preparation method allowed us to directly  
507 detect full miRNA transcripts, and we used our data to enhance the annotation of thylacine  
508 miRNAs. We applied three complementary annotation approaches to obtain a comprehensive  
509 complement of thylacine miRNA loci (see Methods), combining the output of MirMachine (Umu  
510 et al. 2022), which annotates miRNAs by homology search, and MirMiner (Wheeler et al. 2009),  
511 which fits high-throughput sequencing data to a model of miRNA biogenesis, along with  
512 homologous sequence search using the opossum (*M. domestica*) and Tasmanian devil (*S. harrisi*)  
513 genome assemblies. In this way, we expanded the thylacine miRNA repertoire from 62 to 325  
514 annotated miRNA genes (**Supplemental Table 22**). Among the previously described hairpins  
515 (Feigin et al. 2022), all were detected except for Mir-340 and Mir-497 (**Supplemental Table 23**).  
516 The annotated precursor miRNA sequences in the thylacine genome can be found in  
517 **Supplemental File 7**.

518 In summary, we increased the number of annotated thylacine miRNA genes by five-fold.

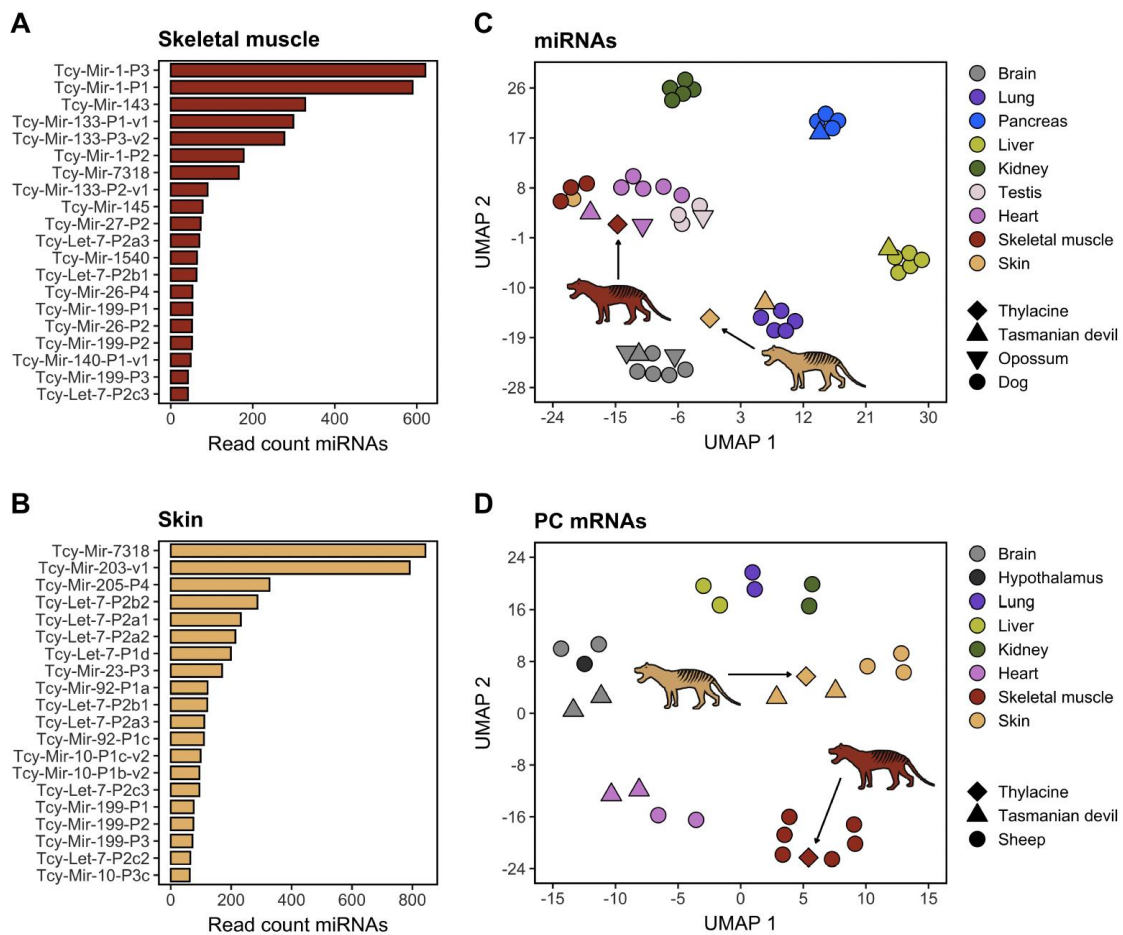
519

### 520 **Thylacine RNA profiles show high tissue-specificity**

521 To obtain a comprehensive profile of the miRNA complement present in our RNA sequencing  
522 data, we focused on reads mapped to the newly annotated thylacine miRNA loci in this study (N  
523 = 325, **Supplemental Table 24**). In skeletal muscle, we detected and quantified a total of 120  
524 distinct miRNAs, while in skin tissues, 143 distinct miRNAs were identified (**Supplemental**  
525 **Table 25**). Unlike other noncoding RNA genes such as rRNAs and tRNAs, miRNA abundance  
526 profiles varied significantly between tissues, with only 35% of the top 20 most abundant miRNAs  
527 being shared. In skeletal muscle, the MIR-1, MIR-133 and MIR-143 families were the most  
528 abundant (**Fig. 5A**), and this was observed for both trimmed and, at much more reduced  
529 resolution, untrimmed reads (**Supplemental Table 25**). From these, MIR-1 and MIR-133 are  
530 miRNAs at times referred to as myomirs, given their abundance and characteristic muscle-specific

531 functions (Chen et al. 2006; Li et al. 2018; Safa et al. 2020). In contrast, the skin miRNA profile  
 532 highlighted the abundance of MIR-7318, MIR-203, MIR-205 and LET-7 families (**Fig. 5B**;  
 533 **Supplemental Table 25**). MIR-203 and MIR-205 are also known to be abundant in the skin  
 534 (Fromm et al. 2022), contributing to epithelial growth and keratinization (Yi and Fuchs 2009;  
 535 Viticchiè et al. 2012; Wang et al. 2013; Jiang et al. 2020). Consistent with miRNA biogenesis  
 536 dynamics (Bartel 2018), we also observed dominant transcripts from either the 5' or 3' arms of  
 537 the precursor molecules (**Supplemental Fig. 11**; **Supplemental Table 26**). This shows how even  
 538 subtle details of miRNA biogenesis are reproduced in our historical RNA expression profiles, and  
 539 further supports the authenticity of the detected thylacine miRNA molecules.

540



541

542 **Figure 5: Divergent RNA profiles in thylacine skin and skeletal muscle samples.** Number of

543 RNA sequences mapped to the 20 most abundant thylacine miRNA genes profiled in (A) skeletal

544 muscle and **(B)** skin tissue. **(C)** UMAP embedding depicting diverse tissue samples clustering  
545 belonging to dog (circular shape), Tasmanian devil (triangular shape), and opossum (inverted  
546 triangular shape) miRNA expression profiles (N = 119) available at MirGeneDB2.1 (Fromm et  
547 al. 2022), as well as miRNA profiles of thylacine skeletal muscle and skin tissues (diamond  
548 shape). **(D)** UMAP embedding depicting diverse tissue samples clustering belonging to sheep  
549 (circular shape) and Tasmanian devil (triangular shape) protein-coding (PC) mRNA expression  
550 profiles (N = 261 mRNAs), as well as PC mRNA expression profiles of thylacine skeletal muscle  
551 and skin tissues (diamond shape).

552

553 In addition, we aimed to determine whether the observed divergent miRNA abundance patterns  
554 in thylacine skeletal muscle and skin tissue resemble their modern mammalian tissue  
555 counterparts. To do so, we used a miRNA-focused tissue expression atlas from the Tasmanian  
556 devil (*S. harrisi*), opossum (*M. domestica*), and dog (*C. familiaris*) species available at  
557 MirGeneDB 2.1 (Fromm et al. 2022). It is well established that tissue-specific miRNA expression  
558 patterns are conserved across evolution, therefore we hypothesized that our historical miRNA  
559 profiles should cluster according to tissues rather than to their species identity. Only shared  
560 miRNAs among the four species included and captured in the thylacine by small RNA sequencing  
561 were considered (N = 119, **Supplemental Table 27**).

562 The UMAP dimensionality reduction built using the miRNA tissue expression atlas demonstrated  
563 a close relationship in tissue identity (**Fig. 5C**). Tasmanian devil, opossum, and dog miRNA  
564 profiles largely clustered according to their tissue of origin. The thylacine skeletal muscle grouped  
565 closely to the skeletal muscle and heart tissues from dog, as well as to heart muscle samples from  
566 Tasmanian devil and opossum. Hence, this is indicative of a conserved miRNA expression profile  
567 across species for muscle-related tissues that was preserved and recovered after RNA sequencing  
568 of thylacine skeletal muscle. The thylacine skin, however, did not reproduce the same pattern  
569 (**Fig. 5C**), probably due to the limited collection of skin samples available, and the lack of  
570 homogeneous miRNA profiles among the few reference skin samples considered. A similar

571 pattern was reproduced when projecting each thylacine sample from skeletal muscle and skin  
572 tissue independently (**Supplemental Fig. 12A**).

573 We also attempted to reproduce the tissue clustering analysis using protein-coding mRNA  
574 expression profiles with sheep (*O. aries*) and Tasmanian devil (*S. harrisi*) tissue atlases,  
575 including reliably profiled protein-coding mRNA genes (breadth of coverage >10%) and with  
576 shared homologous loci among sheep, Tasmanian devil, and the thylacine (N = 261,  
577 **Supplemental Table 27**). In this case, both thylacine skeletal muscle and skin samples clustered  
578 concordantly with the corresponding sheep and Tasmanian devil muscles and skins, revealing a  
579 conserved tissue-specific abundance of protein-coding mRNA transcripts across species, at least  
580 for the protein-coding loci considered in our analyses. When attempting to reproduce the tissue  
581 clustering with each six independent samples from the thylacine, the overall tissue identity  
582 observed with merged samples was preserved (**Supplemental Fig. 12B**), as previously seen for  
583 the embedding using reads mapped to miRNA loci.

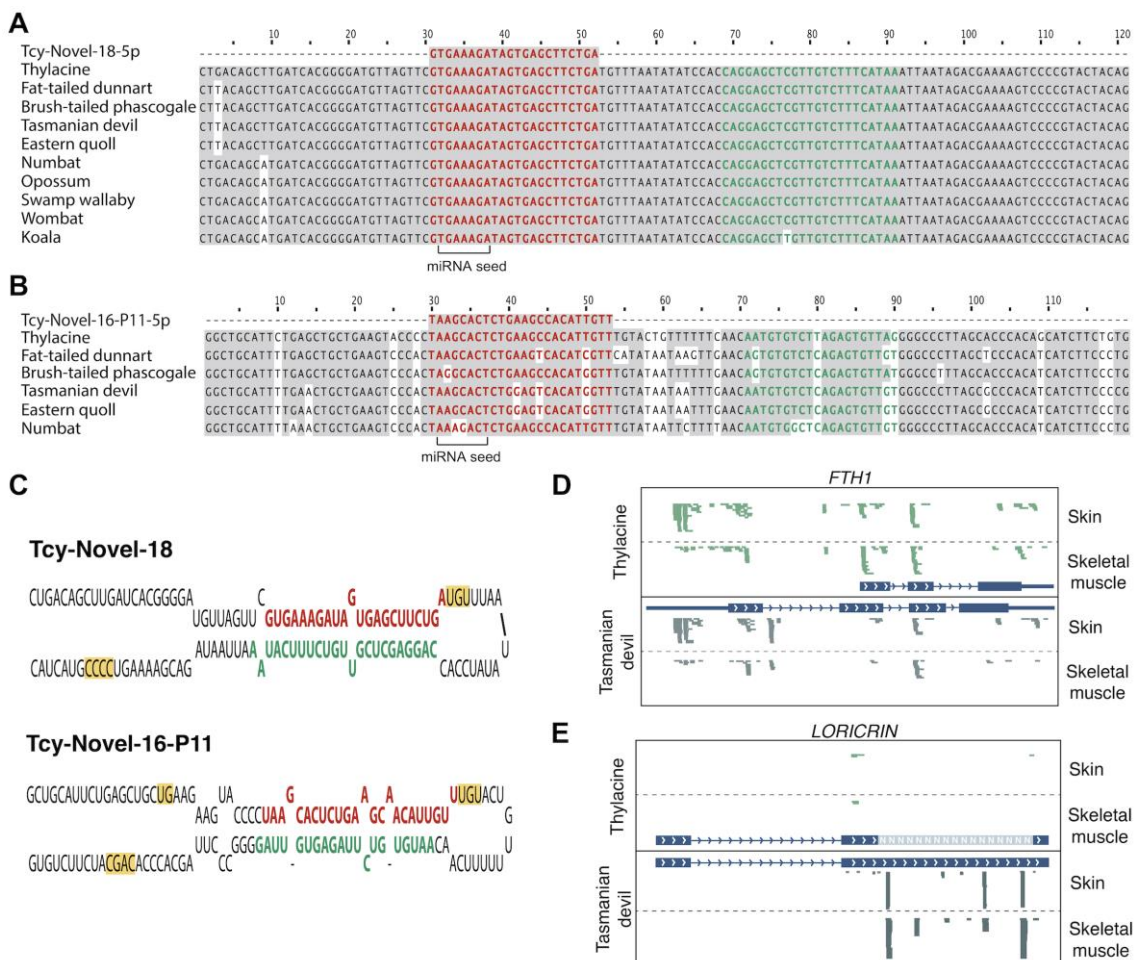
584 In summary, the miRNA and mRNA profiles in historical samples from the extinct thylacine  
585 resemble extant modern animal counterparts, supporting the authenticity of the sequences and  
586 demonstrating that tissue-specific expression profiles can be preserved in dried museum  
587 specimens.

588

### 589 **Discovery of putative novel microRNAs in the thylacine genome**

590 Species-specific miRNAs can only be confirmed through direct sequencing of their RNA  
591 molecules. Therefore, RNA sequencing data offer a unique opportunity to identify species-  
592 specific miRNA genes from extinct organisms like the thylacine. We used the miRNA discovery  
593 algorithm miRDeep2 (Friedländer et al. 2012) to predict novel miRNAs from thylacine skeletal  
594 muscle and skin. After careful curation of all miRNA candidates supported by RNA reads in both  
595 tissues, two promising loci were selected (**Supplemental Tables 28-29**).

596 One of the selected candidates, named Tcy-Novel-18, was found to be highly conserved in all  
 597 analyzed marsupial genomes (**Fig. 6A**), and located in an intron of the E3 ubiquitin ligase ring  
 598 finger 144A (*RNF144A*) gene. This putative novel miRNA did not match any known noncoding  
 599 RNA loci in commonly used RNA databases (see Methods). However, mature miRNA transcripts  
 600 from both stems of the precursor molecule were detected as expressed in Tasmanian devil  
 601 (**Supplemental Fig. 13**) and opossum tissues (**Supplemental Fig. 14**), supporting the reliability  
 602 of this locus as a true novel miRNA conserved and expressed in related extant marsupials.  
 603



604  
 605 **Fig. 6: Novel thylacine miRNAs and improved gene annotations guided by historical RNA**  
 606 **sequences.** Multiple sequence alignment for (A) Tcy-Novel-18 and (B) Tcy-Novel-16-P11  
 607 selected novel miRNA candidates predicted using miRDeep2 software. The 5p arm showing  
 608 transcriptional evidence in the thylacine RNA sequencing data from skeletal muscle and skin

609 tissue is highlighted in bold red. The 3p arm is highlighted in bold green. Nucleotides that are  
610 shared with respect to the thylacine species for each novel miRNA candidate are shown in grey.  
611 (C) Predicted secondary structure folding of the pri-miRNA precursor sequences (+/- 30 nt from  
612 the pre-miRNA) for Tcy-Novel-18 and Tcy-Novel-16-P11 novel miRNA candidates. The 5p and  
613 3p mature arms of the miRNA hairpins are highlighted in bold red and green, respectively.  
614 Processing motifs characteristic of true miRNA loci are highlighted in yellow. Examples of  
615 missing exonic annotations in the thylacine assembly are shown for transcriptional profiles (in  
616 green) obtained for (D) *FTH1* and (E) *LORICRIN* genes, using the thylacine genome assembly as  
617 a reference. RNA sequencing data of skeletal muscle and skin thylacine tissues were aligned to  
618 the Tasmanian devil assembly (in grey) for comparison.

619

620 The other novel miRNA candidate, named Tcy-Novel-16-P11, partially matched the Sha-Novel-  
621 16-P11 miRNA annotated in the Tasmanian devil complement according to MirGeneDB 2.1  
622 (Fromm et al. 2022). The thylacine sequence had three mismatches compared to its Tasmanian  
623 devil homologous miRNA, and this pattern was replicated in the Easter quoll (**Fig. 6B**). Other  
624 dasyuromorphids showed one or two mismatches at different positions, indicating variations in  
625 the seed region. There was no evidence from public sequence data that this putative miRNA is  
626 expressed in any extant species, suggesting that it may be a thylacine-specific miRNA isoform  
627 from the NOVEL-16 family. However, due to the conserved "seed" region with Sha-Novel-16-  
628 P11 and other equivalent dasyuromorphid homologous loci, it is expected to have a common  
629 repertoire of targeted mRNAs and a shared gene regulatory function.

630 Both novel miRNAs were supported by historical RNA reads with damage patterns consistent  
631 with their antiquity (**Supplemental File 8**), as well as by RNA hairpin-like secondary structures  
632 (**Fig. 6C**) resembling those of bona fide miRNAs in extant species. MiRNA-specific processing  
633 motifs (Fang and Bartel 2015) were also found in both candidates. Additional information on the  
634 two putative novel miRNAs and other candidate loci reported by the miRDeep2 software can be  
635 found in **Supplemental Table 30**.

636 In summary, we identified two novel miRNA candidates from transcriptional evidence of  
637 thylacine skeletal muscle and skin tissues.

638

### 639 **Historical RNA fragments guide improved thylacine gene annotations**

640 We performed comparative mapping of RNA sequences to the reference thylacine and Tasmanian  
641 devil assemblies (**Supplemental Tables 31-36**) to investigate potential imperfections in the  
642 thylacine genome annotations. Our analyses revealed two instances of missing thylacine  
643 annotations: The ferritin heavy chain 1 gene (*FTH1*) exhibited differences in the number of exonic  
644 regions between the thylacine and Tasmanian devil. While the current thylacine annotation had  
645 three exons, the corresponding homologous locus in the Tasmanian devil assembly had four (**Fig.**  
646 **6D**). However, transcriptional evidence from the thylacine historical RNA profiles of skeletal  
647 muscle and skin tissues supported the existence of a missing first leading exonic region in the  
648 thylacine genome for the *FTH1* gene, as annotated in the Tasmanian devil assembly.

649 Another example is the *LORICRIN* locus (**Fig. 6E**). Using the Tasmanian devil genome as a  
650 reference, we detected high expression of this gene, with a majority of the mapped reads  
651 originating from the skin tissue. However, when mapping to the thylacine assembly, the  
652 expression profile of the corresponding gene was barely detected (**Supplemental Table 32**).  
653 Upon closer inspection, we found that most of the coding sequence of the *LORICRIN* gene was  
654 missing from the thylacine assembly, represented by unknown 'N' nucleotides. This hindered the  
655 mapping of reads to this genomic region. The high abundance of RNA reads mapping to the  
656 *LORICRIN* gene in the thylacine skin tissue agrees with its relevance as the major protein  
657 component of the cornified envelope in terminal epidermal cells of mammals (Yoneda et al.  
658 1992).

659 In summary, we were able to detect, and partly correct, missing annotations in the thylacine  
660 genome assembly using historical transcriptional evidence from RNA sequencing.

661

662

## 663 **Discussion**

664 In the current study, we present the first successful transcriptomic sequencing evidence from an  
665 extinct metazoan species, the Tasmanian tiger (*T. cynocephalus*), also known as the thylacine.  
666 The recovery of RNA expression profiles no longer existing in living cells expands the possibility  
667 of delving into the biology of extinct animals. Previous studies have reported the sequencing of  
668 RNA molecules from an extremely well-preserved permafrozen canid that lived in the late  
669 Pleistocene, with an estimated age of approximately 14,300 years, as well as from historical wolf  
670 skins (Smith et al. 2019; Fromm et al. 2021). The thylacine specimen used in this study has been  
671 stored at room temperature in embalmed desiccation for more than a century. Desiccation and/or  
672 mummification have proven to be a favorable environment for the preservation of  
673 oligonucleotides (Hekkala et al. 2011; Schuenemann et al. 2017; Rossi et al. 2021; Pedersen et al.  
674 2022; Richards et al. 2022).

675 Furthermore, the presence of different ranges of tissues still available in the thylacine specimen  
676 analyzed entailed a good opportunity to explore one of the main characteristics of RNA profiles  
677 that DNA cannot provide, tissue-specific gene expression signatures. We found protein-coding  
678 transcripts highly representative of the sampled tissues, including titin and actin for skeletal  
679 muscle and keratin for skin. In addition, we found transcriptional profiles that indicate the  
680 presence of slow-acting muscle fibers in the skeletal muscle samples, and a putative source of  
681 blood infiltration in two of the three skins analyzed, showing the resolution of information that  
682 can potentially be extracted from such data.

683 We verified the authenticity of our data in several ways:

684 Metatranscriptomics analyses determined that thylacine was the main source of unambiguously  
685 assigned reads for both long and short reads, with the main contaminant source likely coming  
686 from human manipulation of the specimen. Other contamination sources are doubtful, as they  
687 might have arisen from read misassignments due to reference database biases and sequences

688 mapping to highly conserved loci across species. Despite the high percentage of successful read  
689 mapping to the thylacine genome and the damage patterns found in the RNA sequences  
690 supporting their antiquity, accurately estimating the true endogenous RNA content is challenging.  
691 Unambiguously assigned reads detected by our metatranscriptomics pipeline represent only a  
692 small proportion of sequences mapping to species-specific loci, while the remaining reads with  
693 ambiguous mapping, possibly of endogenous origin, are discarded. Targeted mapping to the  
694 thylacine assembly alone provided a high percentage of successful alignment for short trimmed  
695 reads in both tissues (~63%). However, this value was more reduced for longer untrimmed reads,  
696 possibly due to contamination from highly conserved genomic regions across species. The  
697 unexpected abundance of long (42 nt) trimmed reads mapped to the human assembly and their  
698 reduced damage profile supports this interpretation. Moreover, the reduced number of long  
699 untrimmed reads successfully mapped to the thylacine assembly, along with their damage profiles  
700 resembling those of shorter reads, indicates time-dependent nucleotide modifications in our data,  
701 further reinforcing the antiquity and the probable thylacine origin of most sequences, whether  
702 short or long.

703 Secondly, we observed distinctive abundance profiles in the skeletal muscle and skin samples,  
704 with rRNAs and tRNAs being the most abundant noncoding transcripts, as expected for eukaryote  
705 transcriptomes (Westermann et al. 2012), and protein-coding and miRNA genes showing relevant  
706 hallmarks of tissue-specificity conserved in extant species. This characteristic pattern would not  
707 have been found had the RNA reads come from undisclosed exogenous contamination or from  
708 thylacine DNA sequenced together with our RNA extracts. Nevertheless, our estimates of putative  
709 unwanted endogenous DNA contamination revealed a reduced percentage below 5% of intergenic  
710 mapped reads that might belong to sequenced thylacine DNA. We also observed marsupial-  
711 specific miRNAs. While possible, it is unlikely that marsupial-derived contamination from  
712 species other than the thylacine would have occurred, given the known history of the specimen  
713 and the procedures employed during sampling, RNA extraction, and sequencing.

714 Thirdly, the sequenced RNA molecules showed characteristic patterns of nucleotide substitutions,  
715 similar to damage patterns observed in previous ancient RNA studies (Smith et al. 2019; Fromm  
716 et al. 2021).

717 Despite these compelling results, we observed sample-dependent variability in the recovered  
718 RNA profiles. One of the three samples per tissue generally provided most of the tissue-specific  
719 resolution in coding loci, while the other two had more limited profiles. Nevertheless, the overall  
720 relative tissue-specific RNA abundance was preserved, as evidenced by the tissue clustering at  
721 the miRNA and protein-coding mRNA level using both merged samples and individual samples  
722 from skeletal muscle and skin thylacine tissues. The observed differences among samples may  
723 have originated from inherent differences in RNA preservation due to the sampling strategy, as  
724 well as technical biases introduced during RNA extraction, library preparation, or sequencing.  
725 However, for the majority of the results presented in this study, a merged sample composite for  
726 each tissue was considered, in order to maximize the resolution and statistical power of our  
727 analyses. Since no specific methodologies were applied to remove highly abundant and repetitive  
728 elements of the transcriptome, the depth of coverage for RNA transcripts other than highly  
729 expressed loci, such as rRNAs or tRNAs, is expected to be low. Future developments could  
730 benefit from applying rRNA and/or tRNA fragment depletion protocols to increase the breadth  
731 and depth of coverage of low-abundant target transcripts with a similar sequencing effort.

732 Aside from unraveling patterns of gene activity, another application of RNA sequencing is to  
733 detect molecules not present as DNA copies, such as RNA viruses. We found a limited number  
734 of RNA molecules that could be assigned to viral genomes. Their presence in such old remains  
735 suggests the potential to profile the RNA virome from specimens of extant and extinct species  
736 stored in museum dry collections. Evidence of the reconstruction of historical viral genomes has  
737 been previously reported (Smith et al. 2014; Zhang et al. 2006), and tracing the origins and  
738 evolution of relevant RNA virus families could provide knowledge to prepare mitigation  
739 measures for future pandemics (Keusch et al. 2022).

740 Ancient/historical RNA sequencing opens an unprecedented opportunity to obtain expression  
741 evidence of still unknown loci that are virtually impossible to annotate from DNA information  
742 alone. We used several complementary approaches to annotate thylacine miRNA loci from our  
743 RNA sequence data, increasing the total number of annotated thylacine miRNAs from 62 to 325,  
744 thus bringing it on par with other extant mammalian species (Fromm et al. 2022). In addition, we  
745 predicted two novel miRNA candidates. Given that one of the novel miRNA candidates contains  
746 several nucleotide substitutions relative to homologous loci in other closely related marsupial  
747 species (see Results above), its annotation would have been difficult without the expression  
748 support of the historical RNA sequences employed in this study. The NOVEL-16 family, as  
749 annotated in the Tasmanian devil according to MirGeneDB 2.1 (Fromm et al. 2022), and in the  
750 annotation produced in this study for the thylacine genome, harbors multiple copies with subtle  
751 sequence variations. This denotes an active undergoing evolutionary differentiation for this  
752 miRNA family, giving rise to dasyuromorphid-specific isoforms and novel miRNA-mRNA  
753 interactions. The other novel miRNA candidate was found to be highly conserved in all analyzed  
754 marsupial genomes, and to be broadly expressed in Tasmanian devils and opossums. While it is  
755 not clear why this miRNA has eluded previous annotation efforts in extant marsupial species, it  
756 clearly shows how analyses of historical RNA sequences from extinct ancestors and sister families  
757 can help improve genome annotations and our understanding of the gene regulatory network  
758 repertoire evolution in present-day extant species.

759 The unique characteristics of ancient/historical RNA profiles provide new opportunities to gain  
760 deeper knowledge of the genomic architecture and gene expression regulation of extinct species  
761 such as the thylacine, a still unexplored area that might benefit recent efforts in the field of de-  
762 extinction (Shapiro 2017; Seddon and King 2019). Moreover, exploring other preserved thylacine  
763 specimens available across museum collections could greatly improve the transcriptional  
764 resolution and tissue diversity reached in the present study. In general, the palaeotranscriptomics  
765 field has been neglected and underexplored, with only a few recent examples of reliable data  
766 supporting the preservation of ancient transcriptomic profiles from extant species over time

767 (Smith et al. 2014, 2019; Fromm et al. 2021). Despite the limitations in the recovery of RNA  
768 extracts from well-preserved remains, the present study adds additional proof of RNA molecules  
769 still present and recoverable at amounts sufficient to be representative of true historical  
770 transcriptomic profiles in dry museum collections. Because the timescale of RNA preservation  
771 seems to range several thousand years into the past (Smith et al. 2019; Fromm et al. 2021), we  
772 believe that a vast yet unexplored compendium of preserved tissues awaits further analysis in  
773 search of long-forgotten transcriptomes. Hence, we advocate for applying transcriptome-based  
774 approaches to recover RNA molecules from preserved specimens in dry, and possibly wet,  
775 museum collections, fostering a new era of integrative palaeo-studies covering genomics,  
776 proteomics, and transcriptomics.

777

778

## 779 **Methods**

### 780 **Sample collection, RNA extraction and sequencing**

781 Skeletal muscle and skin tissue samples were collected from an ~130-year-old embalmed  
782 desiccated thylacine specimen (NRM-MA590213), preserved at room temperature and available  
783 at the Stockholm Natural History Museum (Naturhistoriska Riksmuseet, NRM). This adult  
784 specimen was captured on the island of Tasmania and arrived at the NRM collection in 1891 as a  
785 donation from the Australian Museum of Sydney. The exact date of death and sex is unknown.  
786 Tissue samples were obtained in triplicate and stored at -20°C until use. Skeletal muscle tissue  
787 was collected from the inner surface of the left scapula bone, while skin tissue was collected from  
788 three different sections of the ventrolateral skin flaps and the inguinal region, as shown in **Fig. 1**.

789 All laboratory work, including tissue subsampling and homogenization, RNA extraction and  
790 library preparation, were performed in dedicated aDNA facilities within the Centre for  
791 Palaeogenetics (CPG) at Stockholm University, following strict standard guidelines for working  
792 with ancient/historical biomolecules (Knapp et al. 2012).

793 Approximately 80 mg of tissue per sample was sectioned into small pieces with a scalpel and  
794 pulverized in liquid nitrogen using a mortar and pestle. The resulting tissue powder was then  
795 added to 900  $\mu$ l of digestion buffer (Gilbert et al. 2007; Sinding et al. 2015). The resulting lysis  
796 mixture was then incubated for 30 minutes at 37°C. Optionally, for tissue samples that showed  
797 almost no digestion after incubation, a further homogenization process was implemented by  
798 mechanical lysis with 2 ml PowerBead Pro Tubes (Thermo Fisher Scientific) loaded with 2.38  
799 mm metallic beads in a Tissuelyser LT equipment (Qiagen). Subsequently, the total RNA fraction  
800 was isolated using the mirVana™ miRNA isolation kit (Thermo Fisher Scientific) according to  
801 the manufacturer's specifications except for the following: i) substitute the initial Lysis/Binding  
802 buffer with the previously described incubated homogenized tissue mixture, and ii) perform the  
803 final elution in 25  $\mu$ l ultrapure nuclease-free H<sub>2</sub>O and repeat the elution flow through the filter  
804 cartridge twice. The total RNA concentration from each eluted extract was determined in triplicate  
805 using both Qubit™ microRNA and Qubit™ RNA HS (High Sensitivity) Assay kits in a dedicated  
806 Qubit™ 2.0 fluorometer equipment (Thermo Fisher Scientific). Sequencing libraries were  
807 prepared using the NEXTflex™ Small RNA-seq Kit v3 protocol (Bioo Scientific) and allowing  
808 23× PCR amplification cycles with no size selection. A positive control sample was included by  
809 using a 21 nt microRNA-like sequence not matching any known microRNA in miRBase database  
810 (Kozomara et al. 2019) and provided within the NEXTflex™ Small RNA-seq Kit. The resulting  
811 library concentration was then determined with a Qubit™ dsDNA BR (Broad Range) Assay kit  
812 (Thermo Fisher Scientific), and cDNA fragment size distribution and integrity were assessed with  
813 the Agilent High Sensitivity DNA kit assay in a Bioanalyzer 2100 system (Agilent Technologies).  
814 Single-end sequencing was performed independently for each tissue on a NextSeq 500 sequencing  
815 system (Illumina) using the Illumina NextSeq high output sequencing reagent kit (75 cycles).

816

### 817 **Sequence pre-processing and quality control**

818 Raw sequenced reads were processed to remove sequencing adapters using cutadapt 3.2 software  
819 (Martin 2011), with a minimum read length after adapter trimming of 18 nt per read and a

820 maximum error rate of 10% in adapter sequence detection. Trimmed reads were collapsed and  
821 quality control-filtered using the miRTrace *qc* function (Kang et al. 2018) to remove low-quality  
822 reads and identical PCR duplicates. Unique molecular identifiers (UMIs) were then removed from  
823 the trimmed collapsed sequences using the *trimfq* function (*-b 4 -e 4*) from the Seqtk tool  
824 (<https://github.com/lh3/seqtk>) and retained as sequence ID tags for further deduplication  
825 procedures.

826

### 827 **Metatranscriptomics**

828 We performed a taxonomic analysis of the skeletal muscle and skin sequences to estimate the  
829 amount of endogenous RNA and identify potential additional RNA contamination. The analysis  
830 was conducted before PCR deduplication collapsing and after adapter trimming and UMI  
831 removal. Additionally, the same analyses were performed on RNA reads where sequencing  
832 adapters were not identified and were kept untrimmed. The taxonomic classification of the reads  
833 was carried out using the KrakenUniq software (Breitwieser et al. 2018) with the full NCBI non-  
834 redundant reference nucleotide (NT) database. The classification was based on *k*-mer mapping,  
835 commonly used with the standard BLASTN algorithm (Altschul et al. 1990). The resulting  
836 alignments were filtered based on two specific criteria: i) "species" level was selected as the  
837 taxonomic level, and ii) only "species" with >1000 *k*-mers and >200 species-specific reads  
838 (taxReads) were considered.

839 Since reads shorter than 31 nt could not be classified using the KrakenUniq methodology (which  
840 uses a default *k*-mer length of 31), we performed additional alignment using Bowtie 2 v.2.4.2  
841 (Langmead and Salzberg 2012). The reference database used for alignment included the thylacine  
842 genome assembly (Feigin et al. 2022), the hg19 human reference genome, and the full NCBI NT  
843 database as built in December 2020 (Pochon et al. 2022). The use of the hg19 human assembly,  
844 instead of more modern versions like hg38, is not expected to significantly impact the overall  
845 results obtained, as human genomic/transcriptomic content is considered to be of contaminant  
846 origin and not the main focus of our analyses. We used the global alignment Bowtie 2 mode with

847 flags *--end-to-end* and *--very-sensitive*, and selected ultrashort reads (18-30 nt) that mapped  
848 uniquely to one of the references in the merged database (MAPQ  $\geq 1$ ). The sequence IDs were  
849 matched with their corresponding taxIDs using the *seqid2taxid.map* mapping file constructed  
850 using Kraken2 software (Wood et al. 2019) from the NCBI NT taxonomy information. The  
851 number of reads assigned to each taxID was quantified, and organisms detected in skeletal muscle  
852 and skin tissue samples were ranked based on abundance. The filtered taxonomic assignment  
853 proportions were visualized using KronaTools (Ondov et al. 2011). Data visualization and  
854 subsequent analyses were conducted using *ggplot2* graphics (Wickham 2016) within R software  
855 (R Core Team 2022).

856

### 857 **Mapping and quantification**

858 Transcriptome sequence models and exon-exon maps for annotated genes in the thylacine  
859 chromosome-based nuclear genome assembly (Feigin et al. 2022) were generated using the  
860 *gffread* v0.12.6 tool (*-Z -W --force-exons --gene2exon --t-adopt --tlf*). Exonic, intronic, and intron-  
861 intron maps were deduced from genome-wide and transcriptome-wide annotations. Only the  
862 longest transcript isoforms per gene were considered for further analyses. PCR deduplicated and  
863 adapter-trimmed/untrimmed sequences without UMIs from skeletal muscle and skin tissues were  
864 mapped against a composite of the thylacine nuclear and mitochondrial genomes (NC\_011944.1)  
865 (Miller et al. 2009), as well as to the nuclear transcriptome assembly using the Bowtie aligner tool  
866 v1.3.0 (Langmead et al. 2009). The alignment allowed up to one mismatch within a seed equal to  
867 the minimum read length (18 nt) and reported a maximum of one valid alignment with high  
868 sensitivity (*-n 1 -l 18 -k 1 -y --best*). To account for sequence replication, repetitive and  
869 degenerated UMIs were removed using the *dedup* function from UMI-Tools v1.1.2 software  
870 (Smith et al. 2017) with the directional methodology for UMI clustering. UMI identifiers were  
871 obtained from trimmed reads (8-mer UMIs) and untrimmed reads (4-mer UMIs) after adapter  
872 removal. Additionally, the last 4 nucleotides of untrimmed reads were removed to prevent  
873 alignment biases.

874 RNA transcript abundance on a per gene basis was quantified based on UMI-deduplicated  
875 alignments, and the breadth of coverage for each gene was determined using the coverage function  
876 of BEDTools v2.30.0 (Quinlan and Hall 2010). Genes with a breadth of coverage of at least 10%  
877 of the entire coding sequence were considered to have reliable expression evidence. A similar  
878 procedure was applied to map RNA reads from thylacine skeletal muscle and skin tissues to the  
879 Tasmanian devil nuclear genomic DNA + mitochondrial DNA (NC\_018788.1) (Miller et al.  
880 2011), and the transcriptome assembly obtained from NCBI annotation (mSarHar1.11).

881 An investigation was conducted to explore unexpected fish-related contamination identified by  
882 KrakenUniq and Bowtie 2-based metatranscriptomic analyses. Trimmed deduplicated RNA reads  
883 from thylacine skeletal muscle and skin tissues were mapped to the zebrafish reference genome  
884 assembly (GRCz11) using the Bowtie aligner v1.3.0 (*-v 0 -k 1 --best -y*) (Langmead et al. 2009).  
885 Additionally, trimmed RNA reads assigned to zebrafish, human, or mouse genomes by  
886 KrakenUniq were realigned (*-n 1 -l 18 -k 1 -y --best*) to the thylacine (Feigin et al. 2022),  
887 Tasmanian devil (mSarHar1.11), and opossum (MonDom5) assemblies to reevaluate their origin.  
888 Additional analyses describing the identification of expression hotspots genome-wide, rRNA  
889 annotation, DNA contamination and exonic/intronic RNA enrichment are detailed in  
890 **Supplemental Methods.**

891

## 892 **RNA damage**

893 Because no UDG treatment (Briggs et al. 2010) was implemented during RNA extraction and  
894 library preparation to correct for cytosine deamination events, damage patterns in RNA sequences  
895 mapped to the thylacine assembly after PCR and UMI deduplication were assessed using the  
896 *platypus* function from PMDtools software (Skoglund et al. 2014). Damage analyses were  
897 performed for trimmed RNA reads of different length ranges and for untrimmed reads.  
898 Additionally, a comparative analysis of sequence damage patterns was conducted for the observed  
899 excess of long trimmed reads of 42 nt in length. Deamination profiles from reads of 42 nt mapped

900 to both the thylacine and human assemblies (hg19) assemblies were compared with those from  
901 reads of 42 nt that mapped only to the thylacine assembly. Reads mapping to the human genome  
902 were considered as potential modern contamination.

903

#### 904 **miRNA annotation**

905 We used the thylacine genome assembly (Feigin et al. 2022) together with MirGeneDB 2.1  
906 reference miRNA annotation (Fromm et al. 2022) as inputs of MirMachine software (Umu et al.  
907 2022) to predict miRNA loci based on homology searches across the thylacine genome. Besides,  
908 annotated miRNA hairpins from opossum and Tasmanian devil according to MirGeneDB 2.1  
909 (Fromm et al. 2022) were mapped to the thylacine assembly to identify homologous marsupial-  
910 derived miRNA loci. PCR- and UMI-deduplicated RNA reads from thylacine skeletal muscle and  
911 skin tissues were then used with the MirMiner algorithm (Wheeler et al. 2009) for the annotation  
912 of additional novel miRNA genes. Criteria for annotating novel miRNA genes included i) the  
913 expression of at least two 18-25 nt long reads from each arm of the putative miRNA hairpin  
914 precursor, ii) consistent 5'-end homogeneity of supporting RNA reads, iii) at least 16 nt  
915 complementarity between the predicted arms of the precursor miRNA, and iv) a loop sequence  
916 length between 8-40 nt. A summary of the miRNA annotation pipeline is provided in  
917 **Supplemental Fig. 15**. The predicted set of pre-miRNA hairpins was elongated by 30 nt on both  
918 sides to build primary miRNA hairpin annotations for mapping and quantification purposes.

919

#### 920 **Mapping to miRNA loci and quantification**

921 Trimmed RNA reads in the range of 18-25 nt, as well as untrimmed reads, were mapped to the  
922 set of predicted pri-miRNA hairpin precursors in the thylacine genome using the Bowtie aligner  
923 tool v1.3.0 (*-n 1 -l 18 -k 1 -y --best*) (Langmead et al. 2009). A similar procedure was carried out  
924 using Tasmanian devil pri-miRNA hairpins annotated according to MirGeneDB 2.1 (Fromm et  
925 al. 2022). Reads mapping outside mature miRNAs within the stem of pri-miRNA hairpins and

926 those with offset nucleotides covering more than 25% of their sequence length with respect to the  
927 mature miRNAs were considered unreliable mapping events and discarded. UMI deduplication  
928 was performed, and the mature 5p and 3p arms were quantified separately. The ratio between 5p  
929 and 3p abundance was then computed for each miRNA precursor with successfully mapped reads  
930 supporting their expression.

931

### 932 **Tissue clustering**

933 We aimed to determine whether the expression profiles identified in thylacine skeletal muscle and  
934 skin tissue resembled any expected miRNA and/or protein-coding mRNA expression profiles in  
935 modern tissues from extant related species. For miRNAs, samples from the miRNA abundance  
936 tissue atlas of Tasmanian devil (*S. harrisi*), opossum (*M. domestica*), and domestic dog (*C.*  
937 *familiaris*) (Koenig et al. 2016; Penso-Dolfin et al. 2016) available in MirGeneDB 2.1 (Fromm et  
938 al. 2022) were selected. Thylacine miRNA abundance profiles for all annotated miRNAs (N =  
939 325) were transformed to counts per million (CPM) estimates based on the total number of  
940 genome-wide mapped reads after UMI deduplication. Only miRNAs with shared homologous  
941 loci among dogs, Tasmanian devils, opossums, and thylacines were retained (N = 119).

942 For protein-coding mRNAs, we used a gene expression atlas in sheep (*O. aries*) (Jiang et al. 2014)  
943 available at the EMBL-EBI Expression Atlas database (Papatheodorou et al. 2020), along with  
944 Tasmanian devil RNA-seq data from different tissues (Stammnitz et al. 2023). The selected  
945 protein-coding genes were those reliably captured in our thylacine historical RNA sequencing  
946 data in both skeletal muscle and skin tissues (breadth of coverage >10%), with a shared  
947 homologous locus among sheep, Tasmanian devils, and thylacines (N = 261). Expression profiles  
948 for protein-coding mRNAs were transformed to the log<sub>2</sub> scale. To perform a dimensionality  
949 reduction of protein-coding and miRNA profiles independently, the uniform manifold  
950 approximation and projection (UMAP) algorithm (McInnes et al. 2018) was implemented using  
951 the UMAP R package (<https://github.com/tkonopka/umap>) with the following specifications:  
952 *n\_neighbors=5*, *metric="pearson"*, *spread=10*, *random\_state=30*. Thylacine tissues were

953 excluded from the initial embedding estimation. Finally, the thylacine expression profiles of  
954 protein-coding mRNAs and miRNAs were projected onto their corresponding previously learned  
955 UMAP embeddings.

956

### 957 **Novel miRNA prediction**

958 We further implemented the miRDeep2 software (Friedländer et al. 2012) in an attempt to capture  
959 additional putative miRNA loci in the thylacine genome. The *mapper.pl* function (*-d -c -m*) was  
960 used to jointly align trimmed PCR- and UMI-deduplicated thylacine skeletal muscle and skin  
961 RNA sequences against the thylacine assembly (Feigin et al. 2022). The *miRDeep2.pl* function  
962 was then run, including thylacine pre-miRNA and mature miRNA sequences, as well as  
963 Tasmanian devil mature miRNA sequences for cross-species identification. The mature miRNA  
964 sequences of the selected novel miRNA candidates were mapped to miRBase v22.1 (Kozomara  
965 et al. 2019), MirGeneDB 2.1 (Fromm et al. 2022), and Rfam 14 (Kalvari et al. 2021) to identify  
966 any overlaps with already annotated loci. Secondary structure folding was predicted using the  
967 mFold software (Zuker 2003), and the presence of conserved homologous loci in additional  
968 marsupial genome assemblies was assessed, including the Tasmanian devil (*S. harrisii*,  
969 mSarHar1.11), numbat (*M. fasciatus*), Eastern quoll (*D. viverrinus*, DasViv\_v1.0), fat-tailed  
970 dunnart (*S. crassicaudata*), brush-tailed phascogale (*P. tapoatafa*), wombat (*V. ursinus*), opossum  
971 (*M. domestica*, MonDom5), swamp wallaby (*W. bicolor*), and koala (*P. cinereus*,  
972 phaCin\_unsw\_v4.1). Brush-tailed phascogale, numbat and swamp wallaby genome assemblies  
973 were obtained from the publicly available DNA Zoo Consortium (<https://www.dnazoo.org/>)  
974 database generated through 3D de novo assembly (Dudchenko et al. 2017). Fat-tailed dunnart  
975 genome assembly was generated using an unpublished draft from the University of Melbourne  
976 and available through the DNA Zoo Consortium database. The remaining assemblies correspond  
977 to the last version available in the NCBI database. A schematic representation of the evolutionary  
978 relationship among marsupial species is shown in **Supplemental Fig. 16**. Multiple sequence  
979 alignment was performed using the Clustal Omega software (Sievers et al. 2011). Small RNA-

980 seq data from Tasmanian devil and opossum species were mapped to their corresponding genome  
981 assemblies (mSarHar1.11 and MonDom5, respectively) to detect transcriptional evidence of the  
982 selected novel miRNA candidate homologous loci. Alignment was performed using Bowtie  
983 aligner v1.3.0 (Langmead et al. 2009), allowing a maximum of two mismatches (*-v 2 -k 1 --best*  
984 *-y*). Small RNA-seq datasets are available at the NCBI Gene Expression Omnibus (GEO) under  
985 accession numbers GSE18352 (Murchison et al. 2010) and GSE40499 (Meunier et al. 2013).

986

987

## 988 **Data access**

989 The sequence data generated in this study have been submitted to the NCBI BioProject database  
990 (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number PRJNA900297.

991

## 992 **Competing interest statement**

993 The authors declare that they have no competing interests.

994

## 995 **Acknowledgments**

996 We thank Dr. Daniela C Kalthoff at the Swedish Museum of Natural History (Naturhistoriska  
997 riksmuseet, NRM) for facilitating the sampling of the thylacine specimen. The authors  
998 acknowledge the support from the National Genomics Infrastructure funded by the Swedish  
999 Research Council and the Uppsala Multidisciplinary Center for Advanced Computational Science  
1000 for assistance with massive parallel sequencing and to the UPPMAX and HPC2N computational  
1001 infrastructure. We also thank Sarah Werning and the phylopic initiative  
1002 (<https://beta.phylopic.org/>, CC BY 3.0) for facilitating the marsupial shapes used in this study.  
1003 Love Dalén acknowledges support from the Swedish Research Council (Grant nr 2021-00625).

1004 Nikolay Oskolkov is financially supported by the Knut and Alice Wallenberg Foundation as part  
1005 of the National Bioinformatics Infrastructure Sweden (NBIS) at SciLifeLab. Emilio Mármol-  
1006 Sánchez and Marc Riemer Friedländer acknowledge funding from the Strategic Research Area  
1007 (SFO) program of the Swedish Research Council (VR) through Stockholm University. Genome  
1008 assemblies and sequencing data for fat-tailed dunnart, brush-tailed phascogale, numbat, and  
1009 swamp wallaby were obtained from the DNA Zoo Consortium (<https://www.dnazoo.org/>).

1010

### 1011 **Authors' contribution**

1012 MRF, LD, BF, and EM-S conceived this study. EM-S, IB, and EE designed the experimental  
1013 protocols. MRF, LD, and BF secured funding. EM-S, NO, and ZP performed metatranscriptomic  
1014 analyses. EM-S, BF, ZP, PK, EEK, BA, and NO analyzed the data, with input from MRF and LD.  
1015 PK, and VS contributed to visualization. EM-S and MRF wrote the manuscript with contributions  
1016 from all authors.

1017

1018

### 1019 **References**

- 1020 Ahmed RE, Tokuyama T, Anzai T, Chanthra N, Uosaki H. 2022. Sarcomere maturation: function  
1021 acquisition, molecular mechanism, and interplay with other organelles. *Philos Trans R*  
1022 *Soc B* **377**: 20210325.
- 1023 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool.  
1024 *J Mol Biol* **215**: 403–410.
- 1025 Bartel DP. 2018. Metazoan MicroRNAs. *Cell* **173**: 20–51.
- 1026 Bininda-Emonds ORP, Cardillo M, Jones KE, MacPhee RDE, Beck RMD, Grenyer R, Price SA,  
1027 Vos RA, Gittleman JL, Purvis A. 2007. The delayed rise of present-day mammals. *Nature*  
1028 **446**: 507–512.

1029 Binladen J, Wiuf C, Gilbert MTP, Bunce M, Barnett R, Larson G, Greenwood AD, Haile J, Ho  
1030 SYW, Hansen AJ, et al. 2006. Assessing the fidelity of ancient DNA sequences amplified  
1031 from nuclear genes. *Genetics* **172**: 733–741.

1032 Bottinelli R, Reggiani C. 2000. Human skeletal muscle fibres: molecular and functional diversity.  
1033 *Prog Biophys Mol Biol* **73**: 195–262.

1034 Breitwieser FP, Baker DN, Salzberg SL. 2018. KrakenUniq: Confident and fast metagenomics  
1035 classification using unique k-mer counts. *Genome Biol* **19**: 1–10.

1036 Briggs AW, Stenzel U, Meyer M, Krause J, Kircher M, Pääbo S. 2010. Removal of deaminated  
1037 cytosines and detection of in vivo methylation in ancient DNA. *Nucleic Acids Res* **38**:  
1038 e87.

1039 Castello JD, Rogers SO, Starmer WT, Catranis CM, Ma L, Bachand GD, Zhao Y, Smith JE. 1999.  
1040 Detection of tomato mosaic tobamovirus RNA in ancient glacial ice. *Polar Biol* **22**: 207–  
1041 212.

1042 Chen JF, Mandel EM, Thomson JM, Wu Q, Callis TE, Hammond SM, Conlon FL, Wang DZ.  
1043 2006. The role of microRNA-1 and microRNA-133 in skeletal muscle proliferation and  
1044 differentiation. *Nat Genet* **38**: 228–233.

1045 Dabney J, Meyer M, Pääbo S. 2013. Ancient DNA damage. *Cold Spring Harb Perspect Biol* **5**.  
1046 a012567.

1047 de Filippo C, Meyer M, Prüfer K. 2018. Quantifying and reducing spurious alignments for the  
1048 analysis of ultra-short ancient DNA sequences. *BMC Biol* **16**: 121.

1049 Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, Shamim MS, Machol I,  
1050 Lander ES, Aiden AP, et al. 2017. De novo assembly of the *Aedes aegypti* genome using  
1051 Hi-C yields chromosome-length scaffolds. *Science* **356**: 92–95.

1052 Fang W, Bartel DP. 2015. The menu of features that define primary microRNAs and enable *de*  
1053 *novo* design of microRNA genes. *Mol Cell* **60**: 131–145.

1054 Feigin C, Frankenberg S, Pask A. 2022. A chromosome-scale hybrid genome assembly of the  
1055 extinct Tasmanian tiger (*Thylacinus cynocephalus*). *Genome Biol Evol* **14**: evac048.

1056 Feigin CY, Newton AH, Doronina L, Schmitz J, Hipsley CA, Mitchell KJ, Gower G, Llamas B,  
1057 Soubrier J, Heider TN, et al. 2017. Genome of the Tasmanian tiger provides insights into  
1058 the evolution and demography of an extinct marsupial carnivore. *Nat Ecol Evol* **2**: 182–  
1059 192.

1060 Fordyce SL, Ávila-Arcos MC, Rasmussen M, Cappellini E, Romero-Navarro JA, Wales N,  
1061 Alquezar-Planas DE, Penfield S, Brown TA, Vielle-Calzada JP, et al. 2013. Deep  
1062 sequencing of RNA from ancient maize kernels. *PLoS One* **8**: 50961.

1063 Friedländer MR, MacKowiak SD, Li N, Chen W, Rajewsky N. 2012. miRDeep2 accurately  
1064 identifies known and hundreds of novel microRNA genes in seven animal clades. *Nucleic  
1065 Acids Res* **40**: 37–52.

1066 Fromm B, Høye E, Domanska D, Zhong X, Aparicio-Puerta E, Ovchinnikov V, Umu SU, Chabot  
1067 PJ, Kang W, Aslanzadeh M, et al. 2022. MirGeneDB 2.1: toward a complete sampling of  
1068 all major animal phyla. *Nucleic Acids Res* **50**: D204–D210.

1069 Fromm B, Tarbier M, Smith O, Mármol-Sánchez E, Dalén L, Gilbert MTP, Friedländer MR.  
1070 2021. Ancient microRNA profiles of 14,300-yr-old canid samples confirm taxonomic  
1071 origin and provide glimpses into tissue-specific gene regulation from the Pleistocene.  
1072 *RNA* **27**: 324–334.

1073 Gilbert MTP, Hansen AJ, Willerslev E, Rudbeck L, Barnes I, Lynnerup N, Cooper A. 2003.  
1074 Characterization of genetic miscoding lesions caused by postmortem damage. *Am J Hum  
1075 Genet* **72**: 48–61.

1076 Gilbert MTP, Tomsho LP, Rendulic S, Packard M, Drautz DI, Sher A, Tikhonov A, Dalén L,  
1077 Kuznetsova T, Kosintsev P, et al. 2007. Whole-genome shotgun sequencing of  
1078 mitochondria from ancient hair shafts. *Science* **317**: 1927–1930.

1079 Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W,  
1080 Fritz MHY, et al. 2010. A draft sequence of the Neandertal genome. *Science* **328**: 710–  
1081 722.

1082 Gryseels S, Watts TD, Mpolesha JMK, Larsen BB, Lemey P, Muyembe-Tamfum JJ, Teuwen DE,  
1083 Worobey M. 2020. A near full-length HIV-1 genome from 1966 recovered from  
1084 formalin-fixed paraffin-embedded tissue. *Proc Natl Acad Sci U S A* **117**: 12222-12229.

1085 Guy PL. 2013. Ancient RNA? RT-PCR of 50-year-old RNA identifies peach latent mosaic viroid.  
1086 *Arch Virol* **158**: 691–694.

1087 Hekkala E, Shirley MH, Amato G, Austin JD, Charter S, Thorbjarnarson J, Vliet KA, Houck ML,  
1088 Desalle R, Blum MJ. 2011. An ancient icon reveals new mysteries: mummy DNA  
1089 resurrects a cryptic species within the Nile crocodile. *Mol Ecol* **20**: 4199–4215.

1090 Hendy J, Welker F, Demarchi B, Speller C, Warinner C, Collins MJ. 2018. A guide to ancient  
1091 protein studies. *Nat Ecol Evol* **2**: 791–799.

1092 Henriksen K, Karsdal MA. 2019. Type I collagen. *Biochemistry of Collagens, Laminins and*  
1093 *Elastin: Structure, Function and Biomarkers* 1–12. Academic Press.

1094 Howe K, Clark MD, Torroja CF, Torrance J, Berthelot C, Muffato M, Collins JE, Humphray S,  
1095 McLaren K, Matthews L, et al. 2013. The zebrafish reference genome sequence and its  
1096 relationship to the human genome. *Nature* **496**: 498–503.

1097 Jiang D, Guo B, Lin F, Lin S, Tao K. 2020. miR-205 inhibits the development of hypertrophic  
1098 scars by targeting THBS1. *Aging* **12**: 22046–22058.

1099 Jiang Y, Xie M, Chen W, Talbot R, Maddox JF, Faraut T, Wu C, Muzny DM, Li Y, Zhang W, et  
1100 al. 2014. The sheep genome illuminates biology of the rumen and lipid metabolism.  
1101 *Science* **344**: 1168–1173.

1102 Johnson JE, Wold BJ, Hauschka SD. 1989. Muscle creatine kinase sequence elements regulating  
1103 skeletal and cardiac muscle expression in transgenic mice. *Mol Cell Biol* **9**: 3393–3399.

1104 Jónsson H, Schubert M, Seguin-Orlando A, Ginolhac A, Petersen L, Fumagalli M, Albrechtsen  
1105 A, Petersen B, Korneliussen TS, Vilstrup JT, et al. 2014. Speciation with gene flow in  
1106 equids despite extensive chromosomal plasticity. *Proc Natl Acad Sci U S A* **111**: 18655–  
1107 18660.

1108 Kalvari I, Nawrocki EP, Ontiveros-Palacios N, Argasinska J, Lamkiewicz K, Marz M, Griffiths-  
1109 Jones S, Toffano-Nioche C, Gautheret D, Weinberg Z, et al. 2021. Rfam 14: expanded

1110 coverage of metagenomic, viral and microRNA families. *Nucleic Acids Res* **49**: D192–  
1111 D200.

1112 Kang W, Eldfjell Y, Fromm B, Estivill X, Biryukova I, Friedländer MR. 2018. MiRTrace reveals  
1113 the organismal origins of microRNA sequencing data. *Genome Biol* **19**: 1–15.

1114 Keller A, Kreis S, Leidinger P, Maixner F, Ludwig N, Backes C, Galata V, Guerriero G,  
1115 Fehlmann T, Franke A, et al. 2017. miRNAs in ancient tissue specimens of the Tyrolean  
1116 iceman. *Mol Biol Evol* **34**: 793–801.

1117 Keusch GT, Amuasi JH, Anderson DE, Daszak P, Eckerle I, Field H, Koopmans M, Lam SK,  
1118 Neves CG das, Peiris M, et al. 2022. Pandemic origins and a One Health approach to  
1119 preparedness and prevention: Solutions based on SARS-CoV-2 and other RNA viruses.  
1120 *Proc Natl Acad Sci U S A* **119**: e2202871119.

1121 Kjær KH, Winther Pedersen M, De Sanctis B, De Cahsan B, Korneliussen TS, Michelsen CS,  
1122 Sand KK, Jelavić S, Ruter AH, Schmidt AMA, et al. 2022. A 2-million-year-old  
1123 ecosystem in Greenland uncovered by environmental DNA. *Nature* **612**: 283–291.

1124 Knapp M, Clarke AC, Horsburgh KA, Matisoo-Smith EA. 2012. Setting the stage – Building and  
1125 working in an ancient DNA laboratory. *Ann Anat* **194**: 3–6.

1126 Koenig EM, Fisher C, Bernard H, Wolenski FS, Gerrein J, Carsillo M, Gallacher M, Tse A, Peters  
1127 R, Smith A, et al. 2016. The beagle dog microRNA tissue atlas: identifying translatable  
1128 biomarkers of organ toxicity. *BMC Genomics* **17**: 649.

1129 Kozomara A, Birgaoanu M, Griffiths-Jones S. 2019. miRBase: from microRNA sequences to  
1130 function. *Nucleic Acids Res* **47**: D155–D162.

1131 Labeit S, Kolmerer B. 1995. Titins: Giant proteins in charge of muscle ultrastructure and  
1132 elasticity. *Science* **270**: 293–296.

1133 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**:  
1134 357–359.

1135 Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment  
1136 of short DNA sequences to the human genome. *Genome Biol* **10**: R25.

1137 Lee EH, Hsin J, Mayans O, Schulten K. 2007. Secondary and tertiary structure elasticity of titin  
1138 Z1Z2 and a titin chain model. *Biophys J* **93**: 1719–1735.

1139 Li N, Zhou H, Tang Q. 2018. miR-133: A suppressor of cardiac remodeling? *Front Pharmacol*  
1140 **9**: 903.

1141 Liu B, Cao J, Wang X, Guo C, Liu Y, Wang T. 2021. Deciphering the tRNA-derived small RNAs:  
1142 origin, development, and future. *Cell Death Dis* **13**: 24.

1143 Lord E, Dussex N, Kierczak M, Díez-del-Molino D, Ryder OA, Stanton DWG, Gilbert MTP,  
1144 Sánchez-Barreiro F, Zhang G, Sinding MHS, et al. 2020. Pre-extinction demographic  
1145 stability and genomic signatures of adaptation in the woolly rhinoceros. *Curr Biol* **30**:  
1146 3871-3879.e7.

1147 Losos JB. 2011. Convergence, adaptation, and constraint. *Evolution* **65**: 1827–1840.

1148 Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
1149 *EMBnet J* **17**: 10–12.

1150 McInnes L, Healy J, Saul N, Großberger L. 2018. UMAP: Uniform Manifold Approximation and  
1151 Projection. *J Open Source Softw* **3**: 861.

1152 Meunier J, Lemoine F, Soumillon M, Liechti A, Weier M, Guschanski K, Hu H, Khaitovich P,  
1153 Kaessmann H. 2013. Birth and expression evolution of mammalian microRNA genes.  
1154 *Genome Res* **23**: 34–45.

1155 Miller W, Drautz DI, Janecka JE, Lesk AM, Ratan A, Tomsho LP, Packard M, Zhang Y,  
1156 McClellan LR, Qi J, et al. 2009. The mitochondrial genome sequence of the Tasmanian  
1157 tiger (*Thylacinus cynocephalus*). *Genome Res* **19**: 213–220.

1158 Miller W, Hayes VM, Ratan A, Petersen DC, Wittekindt NE, Miller J, Walenz B, Knight J, Qi J,  
1159 Zhao F, et al. 2011. Genetic diversity and population structure of the endangered  
1160 marsupial *Sarcophilus harrisii* (Tasmanian devil). *Proc Natl Acad Sci U S A* **108**: 12348–  
1161 12353.

1162 Mitchell KJ, Pratt RC, Watson LN, Gibb GC, Llamas B, Kasper M, Edson J, Hopwood B, Male  
1163 D, Armstrong KN, et al. 2014. Molecular phylogeny, biogeography, and habitat  
1164 preference evolution of marsupials. *Mol Biol Evol* **31**: 2322–2330.

1165 Murchison EP, Tovar C, Hsu A, Bender HS, Kheradpour P, Rebbeck CA, Obendorf D, Conlan  
1166 C, Bahlo M, Blizzard CA, et al. 2010. The Tasmanian devil transcriptome reveals  
1167 Schwann cell origins of a clonally transmissible cancer. *Science* **327**: 84–87.

1168 Newton AH, Weisbecker V, Pask AJ, Hipsley CA. 2021. Ontogenetic origins of cranial  
1169 convergence between the extinct marsupial thylacine and placental gray wolf. *Commun*  
1170 *Biol* **4**: 51.

1171 Ng TFF, Chen LF, Zhou Y, Shapiro B, Stiller M, Heintzman PD, Varsani A, Kondov NO, Wong  
1172 W, Deng X, et al. 2014. Preservation of viral genomes in 700-y-old caribou feces from a  
1173 subarctic ice patch. *Proc Natl Acad Sci U S A* **111**: 16842–16847.

1174 Ondov BD, Bergman NH, Phillippy AM. 2011. Interactive metagenomic visualization in a Web  
1175 browser. *BMC Bioinformatics* **12**: 385.

1176 Ordway GA, Garry DJ. 2004. Myoglobin: an essential hemoprotein in striated muscle. *J Exp Biol*  
1177 **207**: 3441–3446.

1178 Paddle R. 2000. The last Tasmanian tiger: the history and extinction of the thylacine. *J Mammol*  
1179 **83**: 634–636.

1180 Palkopoulou E, Mallick S, Skoglund P, Enk J, Rohland N, Li H, Omrak A, Vartanyan S, Poinar  
1181 H, Götherström A, et al. 2015. Complete genomes reveal signatures of demographic and  
1182 genetic declines in the woolly mammoth. *Curr Biol* **25**: 1395–1400.

1183 Papatheodorou I, Moreno P, Manning J, Fuentes AMP, George N, Fexova S, Fonseca NA,  
1184 Füllgrabe A, Green M, Huang N, et al. 2020. Expression Atlas update: from tissues to  
1185 single cells. *Nucleic Acids Res* **48**: D77–D83.

1186 Patrono L v., Vrancken B, Budt M, Dux A, Lequime S, Boral S, Gilbert MTP, Gogarten JF,  
1187 Hoffmann L, Horst D, et al. 2022. Archival influenza virus genomes from Europe reveal  
1188 genomic variability during the 1918 pandemic. *Nat Commun* **13**: 2314.

1189 Pedersen MW, Antunes C, de Cahsan B, Moreno-Mayar JV, Sikora M, Vinner L, Mann D,  
1190 Klimov PB, Black S, Michieli CT, et al. 2022. Ancient human genomes and  
1191 environmental DNA from the cement attaching 2,000-year-old head lice nits. *Mol Biol*  
1192 *Evol* **39**: msab351.

1193 Penso-Dolfin L, Swofford R, Johnson J, Alföldi J, Lindblad-Toh K, Swarbreck D, Moxon S, di  
1194 Palma F. 2016. An improved microRNA annotation of the canine genome. *PLoS One* **11**:  
1195 e0153453.

1196 Pochon Z, Bergfeldt N, Kırdök E, Vicente M, Naidoo T, Valk T van der, Altınışik NE,  
1197 Krzewińska M, Dalen L, Götherström A, et al. 2022. aMeta: an accurate and memory-  
1198 efficient ancient Metagenomic profiling workflow. *bioRxiv* 2022.10.03.510579.

1199 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic  
1200 features. *Bioinformatics* **26**: 841–842.

1201 R Core Team. 2022. R: A language and environment for statistical computing. [https://www.R-](https://www.R-project.org)  
1202 [project.org](https://www.R-project.org)

1203 Richards SM, Li L, Breen J, Hovhannisyan N, Estrada O, Gasparyan B, Gilliam M, Smith A,  
1204 Cooper A, Zhang H. 2022. Recovery of chloroplast genomes from medieval millet grains  
1205 excavated from the Areni-1 cave in southern Armenia. *Sci Rep* **12**: 15164.

1206 Rollo F. 1985. Characterisation by molecular hybridization of RNA fragments isolated from  
1207 ancient (1400 B.C.) seeds. *Theor Appl Genet* **71**: 330–333.

1208 Rollo F, Venanzi FM, Amici A. 1991. Nucleic acids in mummified plant seeds: biochemistry and  
1209 molecular genetics of pre-Columbian maize. *Genet Res* **58**: 193–201.

1210 Rossi C, Ruß-Popa G, Mattiangeli V, McDaid F, Hare AJ, Davoudi H, Laleh H, Lorzadeh Z,  
1211 Khazaeli R, Fathi H, et al. 2021. Exceptional ancient DNA preservation and fibre remains  
1212 of a Sasanian saltmine sheep mummy in Chehrābād, Iran. *Biol Lett* **17**: 20210222.

1213 Rovinsky DS, Evans AR, Adams JW. 2021. Functional ecological convergence between the  
1214 thylacine and small prey-focused canids. *BMC Ecol Evol* **21**: 58.

1215 Safa A, Bahroudi Z, Shoorei H, Majidpoor J, Abak A, Taheri M, Ghafouri-Fard S. 2020. miR-1:  
1216 A comprehensive review of its role in normal development and diverse disorders. *Biomed*  
1217 *Pharmacot* **132**: 110903.

1218 Schuenemann VJ, Peltzer A, Welte B, van Pelt WP, Molak M, Wang CC, Furtwängler A, Urban  
1219 C, Reiter E, Nieselt K, et al. 2017. Ancient Egyptian mummy genomes suggest an  
1220 increase of Sub-Saharan African ancestry in post-Roman periods. *Nat Commun* **8**: 1–11.

- 1221 Seddon PJ, King M. 2019. Creating proxies of extinct species: the bioethics of de-extinction.  
1222 *Emerg Top Life Sci* **3**: 731–735.
- 1223 Shapiro B. 2017. Pathways to de-extinction: how close can we get to resurrection of an extinct  
1224 species? *Funct Ecol* **31**: 996–1002.
- 1225 Shaw B, Burrell CL, Green D, Navarro-Martinez A, Scott D, Daroszewska A, van 'T Hof R,  
1226 Smith L, Hargrave F, Mistry S, et al. 2019. Molecular insights into an ancient form of  
1227 Paget's disease of bone. *Proc Natl Acad Sci U S A* **116**: 10463–10472.
- 1228 Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H, Remmert  
1229 M, Söding J, et al. 2011. Fast, scalable generation of high-quality protein multiple  
1230 sequence alignments using Clustal Omega. *Mol Syst Biol* **7**: 539.
- 1231 Sinding MHS, Arneborg J, Nyegaard G, Gilbert MTP. 2015. Ancient DNA unravels the truth  
1232 behind the controversial GUS Greenlandic Norse fur samples: the bison was a horse, and  
1233 the muskox and bears were goats. *J Archaeol Sci* **53**: 297–303.
- 1234 Skoglund P, Northoff BH, Shunkov M v., Derevianko AP, Pääbo S, Krause J, Jakobsson M. 2014.  
1235 Separating endogenous ancient DNA from modern day contamination in a Siberian  
1236 Neandertal. *Proc Natl Acad Sci U S A* **111**: 2229–2234.
- 1237 Smith O, Clapham A, Rose P, Liu Y, Wang J, Allaby RG. 2014. A complete ancient RNA  
1238 genome: identification, reconstruction and evolutionary history of archaeological Barley  
1239 Stripe Mosaic Virus. *Sci Rep* **4**: 1–6.
- 1240 Smith O, Dunshea G, Sinding MHS, Fedorov S, Germonpre M, Bocherens H, Gilbert MTP. 2019.  
1241 Ancient RNA from Late Pleistocene permafrost and historical canids shows tissue-  
1242 specific transcriptome survival. *PLoS Biol* **17**: e3000166.
- 1243 Smith O, Gilbert MTP. 2018. Ancient RNA. *Paleogenomics. Population Genomics.* 53–74.  
1244 Springer, Cham.
- 1245 Smith T, Heger A, Sudbery I. 2017. UMI-tools: Modelling sequencing errors in Unique Molecular  
1246 Identifiers to improve quantification accuracy. *Genome Res* **27**: gr.209601.116.

1247 Stammnitz MR, Gori K, Kwon YM, Harry E, Martin FJ, Billis K, Cheng Y, Baez-Ortega A, Chow  
1248 W, Comte S, et al. 2023. The evolution of two transmissible cancers in Tasmanian devils.  
1249 *Science* **380**: 283–293.

1250 Umu SU, Trondsen H, Paynter VM, Buschmann T, Rounge TB, Peterson KJ, Fromm B. 2022.  
1251 Accurate microRNA annotation of animal genomes using trained covariance models of  
1252 curated microRNA complements in MirMachine. *bioRxiv* doi:  
1253 10.1101/2022.11.23.517654

1254 van der Valk T, Pečnerová P, Díez-del-Molino D, Bergström A, Oppenheimer J, Hartmann S,  
1255 Xenikoudakis G, Thomas JA, Dehasque M, Sağlıcan E, et al. 2021. Million-year-old  
1256 DNA sheds light on the genomic history of mammoths. *Nature* **591**: 265–269.

1257 Venanzi FM, Rollo F. 1990. Mummy RNA lasts longer. *Nature* **343**: 25–26.

1258 Vilstrup JT, Seguin-Orlando A, Stiller M, Ginolhac A, Raghavan M, Nielsen SCA, Weinstock J,  
1259 Froese D, Vasiliev SK, Ovodov ND, et al. 2013. Mitochondrial phylogenomics of modern  
1260 and ancient equids. *PLoS One* **8**: e55950.

1261 Viticchiè G, Lena AM, Cianfarani F, Odorisio T, Annicchiarico-Petruzzelli M, Melino G, Candi  
1262 E. 2012. MicroRNA-203 contributes to skin re-epithelialization. *Cell Death Dis* **3**: e435.

1263 Wang D, Zhang Z, O’Loughlin E, Wang L, Fan X, Lai EC, Yi R. 2013. MicroRNA-205 controls  
1264 neonatal expansion of skin stem cells by modulating the PI3K pathway. *Nat Cell Biol* **15**:  
1265 1153–1163.

1266 Westermann AJ, Gorski SA, Vogel J. 2012. Dual RNA-seq of pathogen and host. *Nature Reviews*  
1267 *Microbiology* **10**: 618–630.

1268 Wheeler BM, Heimberg AM, Moy VN, Sperling EA, Holstein TW, Heber S, Peterson KJ. 2009.  
1269 The deep evolution of metazoan microRNAs. *Evol Dev* **11**: 50–68.

1270 Wickham H. 2016. ggplot2: Elegant graphics for Data Analysis. Springer-Verlag New York.  
1271 <https://ggplot2.tidyverse.org>

1272 Wood DE, Lu J, Langmead B. 2019. Improved metagenomic analysis with Kraken 2. *Genome*  
1273 *Biol* **20**: 1–13.

- 1274 Worobey M, Watts TD, McKay RA, Suchard MA, Granade T, Teuwen DE, Koblin BA, Heneine  
1275 W, Lemey P, Jaffe HW. 2016. 1970s and “Patient 0” HIV-1 genomes illuminate early  
1276 HIV/AIDS history in North America. *Nature* **539**: 98–101.
- 1277 Xiao YL, Kash JC, Beres SB, Sheng ZM, Musser JM, Taubenberger JK. 2013. High-throughput  
1278 RNA sequencing of a formalin-fixed, paraffin-embedded autopsy lung tissue sample  
1279 from the 1918 influenza pandemic. *J Pathol* **229**: 535–545.
- 1280 Yi R, Fuchs E. 2009. MicroRNA-mediated control in the skin. *Cell Death Diff* **17**: 229–235.
- 1281 Yoneda K, McBride OW, Korge BP, Kim IG, Steinert PM. 1992. The cornified cell envelope:  
1282 Loricrin and transglutaminases. *J Dermatol* **19**: 761–764.
- 1283 Zhang G, Shoham D, Gilichinsky D, Davydov S, Castello JD, Rogers SO. 2006. Evidence of  
1284 influenza a virus RNA in siberian lake ice. *J Virol* **80**: 12229–12235.
- 1285 Zuker M. 2003. Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic  
1286 Acids Res* **31**: 3406-3415.