

1 A comparative genomics multitool for scientific discovery and conservation

2 *Zoonomia Consortium**

3 *A list of authors and their affiliations appears at the end of the paper

4 **The Zoonomia Project is investigating the genomics of shared and specialized traits in**
5 **eutherian mammals. We describe a whole-genome alignment of 240 species with**
6 **unprecedented phylogenetic diversity, with over 80% of mammalian families represented,**
7 **and new genome assemblies for 131 species. We find that regions of reduced genetic**
8 **diversity are more abundant in species with high extinction risk, discern signals of**
9 **evolutionary selection at unprecedented resolution, and describe insights enabled by**
10 **individual reference genomes. By prioritizing phylogenetic diversity and making data**
11 **available quickly, without restriction, the Zoonomia Project aims to support biological**
12 **discovery, medical research, and biodiversity conservation.**

13 The genomics revolution is enabling advances not just in medical research¹, but also in basic
14 biology² and in biodiversity conservation, where genomic tools have helped apprehend poachers³
15 and protect endangered populations⁴. Yet we still have only limited ability to predict which
16 genomic variants lead to changes in organism-level phenotypes, such as increased disease risk —
17 a task complicated by the sheer size of the human genome (~3 billion nucleotides)⁵.

18 Comparative genomics can address this challenge by identifying nucleotide positions unchanged
19 across millions of years of evolution⁶, suggesting that changes negatively impact fitness and
20 focusing the search for disease-causing variants. In 2011, the 29 Mammals Project⁷ identified 12
21 base pair (bp) regions of evolutionary constraint, totalling 5.4% of the genome, by measuring
22 sequence conservation in humans plus 28 other mammals. These regions proved to be more
23 enriched for complex disease heritability than any other functional mark, including coding
24 status⁸. By expanding the number of species, and making an alignment independent of any
25 single reference genome, the Zoonomia Project was designed to detect evolutionary constraint in
26 the eutherian lineage at unprecedented resolution, while providing new genomic resources for
27 over 130 species.

28 **Designing a comparative-genomics multitool**

29 When selecting species, we sought to maximize evolutionary branch length, to include at least
30 one species from each eutherian family, and to prioritize species of medical, biological, or
31 biodiversity conservation interest. Our assemblies increase the percentage of eutherian families
32 with a representative genome from 49% to 82%, and include nine species that are the sole extant
33 member of their family and seven that are critically endangered (**Figure 1**)⁹: the Mexican howler
34 monkey (*Alouatta palliata mexicana*), Hirola (*Beatragus hunteri*), Russian saiga (*Saiga tatarica*
35 *tatarica*), social Tuco-tuco (*Ctenomys sociabilis*), indri (*Indri indri*), northern white rhino
36 (*Ceratotherium simum cottoni*) and black rhinoceros (*Diceros bicornis*).

37 We collaborated with 28 different institutions to collect samples, with nearly half (47%)
38 provided by The Frozen Zoo[®] of San Diego Zoo Global (**Supplementary Table 1**). Since 1975,
39 The Frozen Zoo[®] has stored renewable cell cultures for about 10,000 vertebrate animals
40 representing over 1,100 taxa, including more than 200 species classified as vulnerable,
41 endangered, critically endangered, or extinct¹⁰. For 36 target species we were unable to acquire a
42 DNA sample of sufficient quality, even though our requirements were modest (see below),
43 highlighting a major impediment to expanding the phylogenetic diversity of genomics.

44 We used two complementary approaches to generate genome assemblies (**Extended Data Table**
45 **1**). First, for 131 genomes, we generated shorter contiguity assemblies (“DISCOVAR
46 assemblies”) by performing a single lane of sequencing (2x250bp reads) on PCR-free libraries,
47 and assembling with *DISCOVAR de novo*¹¹. This method does not require intact cells and uses
48 less than two micrograms of medium-quality DNA (most fragments >5 kilobases (kb)), which
49 allowed us to include difficult-to-access species (**Extended Data Figure 1; Extended Data**
50 **Figure 2**) while achieving contig (contiguous sequences constructed from overlapping short
51 reads) lengths comparable to existing assemblies (median contig N50 of 46.8kb; compared to
52 47.9kb for Refseq genome assemblies).

53 For nine DISCOVAR genomes, and one pre-existing assembly, the lesser hedgehog tenrec
54 (*Echinops telfairi*), we increased contiguity 200-fold (median scaffold length increased from
55 90.5kb to 18.5Mb) through proximity-ligation, which uses chromatin interaction data to capture
56 the physical relationships among genomic regions¹². Unlike short-contiguity genomes, these
57 assemblies capture structural changes like chromosomal rearrangements¹³. The upgraded
58 assemblies increase the number of eutherian orders represented by a long-range assembly (contig
59 N50 > 20kb, scaffold N50 > 10Mb) from 12 to 18 out of 19. We are working on upgrading the
60 large treeshrew (*Tupaia tana*) assembly for the remaining order, Scandentia.

61

62 **Comparative power of 240 species**

63 The Zoonomia alignment includes 120 of our new assemblies and 121 existing assemblies
64 representing a total of 240 species (the dataset includes assemblies for two different dogs)
65 spanning ~110 million years of mammal evolution (**Supplementary Table 2**). With a total
66 evolutionary branch length of 16.6 substitutions per site, we expect just 191 positions in the
67 human genome (0.000006%) to be identical across the aligned species due to chance (false
68 positives) rather than evolutionary constraint (**Extended Data Table 2**). We applied this same
69 calculation to The Exome Aggregation Consortium (ExAC), which analyzed exomes for 60,706
70 humans¹⁴, and estimated that 88% of positions would be expected to have no variation. This
71 illustrates the potential for relatively small cross-species datasets to inform human genetic
72 studies — even for diseases driven by high penetrance coding mutations, for which ExAC is
73 optimally powered¹⁵.

74 **Biological insights from new assemblies**

75 The scope and species diversity in the Zoonomia project supports evolutionary studies in many
76 different lineages. Papers published to date, and the demonstrated utility of existing comparative
77 genomics resources^{16,17}, illustrate the benefits of making new genome assemblies and alignments
78 accessible to all researchers, without restrictions on use.

79 **Speciation.** Comparing our assembly for the endangered mantled howler monkey subspecies
80 *Alouatta palliata mexicana* with the Guatemalan black howler monkey (*Alouatta pigra*), which
81 has a neighboring range, suggests different forms of selection shape reproductive isolation¹⁸.
82 Initial divergence in allopatry was followed by positive selection on prezygotic isolating
83 mechanisms, offering empirical support for a speciation process first outlined by Dobzhansky in
84 1935¹⁹.

85 **Protection from cancer.** Using our assembly for the capybara (*Hydrochoerus hydrochaeris*), a
86 giant rodent, Herrera-Alvarez *et al.* identified positive selection on anti-cancer pathways,
87 echoing earlier reports that another large mammal, the elephant, carries extra copies (retrogenes)
88 of the tumor-suppressor gene *TP53*^{20,21}. This offers a possible resolution to Peto's paradox — the
89 observation that cancer in large mammals is rarer than expected — and could reveal new anti-
90 cancer mechanisms.

91 **Convergent evolution of venom.** Casewell *et al.* used our assembly for the Hispaniolan
92 solenodon (*Solenodon paradoxus*; **Extended Data Figure 2**), to investigate venom production, a
93 trait in just a few eutherian lineages, including shrews and solenodons²². They identified
94 paralogous copies of a kallikrein 1 serine protease (*KLK1*) that, together, encode solenodon
95 venom, and showed that the *KLK1s* were independently co-opted in solenodons and shrews in an
96 intriguing example of molecular convergence.

97 **Informing biodiversity conservation strategies.** Beichman *et al.* analyzed our giant otter
98 (*Pteronura brasiliensis*) assembly and found low diversity and elevated burden of putatively
99 deleterious genetic variants, consistent with the otter's recent population decline through
100 overhunting and habitat loss²³. Intriguingly, the giant otter had fewer putatively deleterious
101 variants than either the southern or northern sea otter, suggesting highest potential for recovery if
102 populations are protected.

103 **Rapid assessment of species infection risk during COVID-19 pandemic.** Using the Zoonomia
104 alignment, and public genomic data from hundreds of other vertebrates, Damas *et al.* compared
105 the structure of ACE2 (the SARS-CoV-2 receptor) and identified 47 mammals that have a high
106 or very high likelihood of being virus reservoirs, intermediate hosts, or good model organisms to
107 study COVID-19, and detected positive selection in the ACE2 receptor binding domain specific
108 to bats²⁴.

109 **Genetic diversity and extinction risk**

110 We next asked whether a reference genome from a single individual can help identify
111 populations with low genetic diversity to prioritize in biodiversity-conservation efforts. Diversity
112 metrics reflect demographic history^{25,26}, and heterozygosity is lower in threatened species²⁷. This

113 analysis was feasible because we used a single sequencing and assembly protocol for all
114 DISCOVAR assemblies, minimizing variation in accuracy, completeness, and contiguity due to
115 the sequencing technology and the assembly process that would otherwise confound species
116 comparisons.

117 We estimated genetic diversity for 130 of our DISCOVAR assemblies, each representing a
118 different species (**Supplementary Table 3**). Four failed during analysis. For the remaining 126,
119 we calculated two metrics: (1) the fraction of sites at which the sequenced individual is
120 heterozygous (“overall heterozygosity”); and (2) the proportion of the genome residing in an
121 extended region without any variation (“segments of homozygosity”, or SoH). SoH is designed
122 for short-contiguity assemblies, where scaffolds are potentially shorter than runs of
123 homozygosity. Overall, heterozygosity and SoH are correlated (Pearson correlation $r=-0.56$;
124 $p=1.8 \times 10^{-9}$; $N=98$). However, while overall heterozygosity is correlated with contig N50
125 (Pearson correlation $r_{\text{het}}=-0.39$; $p_{\text{het}}=4 \times 10^{-5}$; $N_{\text{het}}=105$), likely due to the difficulty of assembling
126 more heterozygous genomes²⁸, SoH is not (Pearson correlation; $r_{\text{SoH}}=0.09$ $p_{\text{SoH}}=0.38$; $N_{\text{SoH}}=98$).
127 Overall heterozygosity and SoH are highly correlated with between the lower- and high-
128 contiguity versions of the upgraded assemblies (Pearson correlation; $r_{\text{het}}=0.999$; $p_{\text{het}}=5 \times 10^{-7}$;
129 $N_{\text{het}}=7$; $r_{\text{SoH}}=0.996$; $p_{\text{SoH}}=1.4 \times 10^{-6}$ $N_{\text{SoH}}=7$).

130 Genomic diversity varies significantly among species in different International Union of
131 Conservation Nature (IUCN) conservation categories, as measured by overall heterozygosity
132 (**Figure 2A**), and SoH (**Figure 2B**). SoH increases ($p=0.0235$; $R^2=0.055$; $N=94$) with increasing
133 levels of conservation concern; heterozygosity decreases ($p=0.011$; $R^2=0.064$; $N=101$). There is
134 no significant difference between wild and captive populations in overall heterozygosity (**Figure**
135 **2C**) or SoH (**Figure 2D**).

136 Unusual diversity values can suggest particular population demographics, although data from
137 more than a single individual is needed to confirm these inferences. All seven critically
138 endangered species have SoH higher than the median for species categorized as Least Concern
139 (**Figure 2E**). The genomes with the lowest heterozygosity and highest SoH were the social tuco-
140 tuco (*Ctenomys sociabilis*; $\text{het}=0.00063$; $\text{SoH}=78.7\%$), sampled from small laboratory colony
141 with just 12 founders²⁹, and the eastern mole (*Scalopus aquaticus*; $\text{het}=0.0008$; $\text{SoH}=81.3\%$),
142 supplied by a professional mole catcher and likely from a population bottlenecked by pest
143 control measures.

144 The correlation between diversity metrics and IUCN is not explained by other species-level
145 phenotypes. For Least Concern ($N=75$) species, we assessed 21 phenotypes cataloged in the
146 Pantheria³⁰ database for correlation with heterozygosity or SoH. The most significant was
147 between SoH and litter size, a trait also shown to predict extinction risk ($p_{\text{SoH}}=0.02$)³¹, but none
148 is significant after Bonferroni correction (**Extended Data Table 3**).

149 Our inference that diversity trends lower in species at higher risk of extinction comes from a
150 small fraction (2.6%) of threatened mammals⁹. Whether a direct correlation with extinction risk,
151 or arising from association of diversity with species-level phenotypes such as litter size, it

152 suggests valuable information can be gleaned from sequencing just a single individual. Should
153 this pattern prove robust across more species, diversity metrics from a single reference genome
154 could help identify populations at-risk, even when few species-level phenotypes are documented,
155 and prioritize species for population-level follow-up.

156 **Resources for biodiversity conservation**

157 For each genome assembly, we cataloged all high-confidence variant sites (broad.io/variants) to
158 support the design of cost-effective, accurate genetic assays that are usable even when sample
159 quality is low³² and often preferable to designing expensive custom tools, relying on tools from
160 related species, or sequencing random regions³³. The reference genomes themselves support
161 development of technologies such as using gene drive to control invasive species, or “de-
162 extinction” through cloning and genetic engineering³⁴.

163 Our genomes have two notable limitations: We sequenced only a single individual, which is
164 insufficient for studying population origins, population structure and recent demographic
165 events^{35,36}, and the shorter contiguity of our assemblies prevented us from analyzing runs of
166 homozygosity (RoH)²⁶. This highlights a dilemma facing all large-scale genomics initiatives:
167 determining when the value of sequencing additional individuals exceeds the value of improving
168 the reference genome itself.

169 **Whole-genome alignment**

170 We aligned the genomes of 240 species (our assemblies and other mammalian genomes released
171 when we started the alignment) as part of a 600-way pan-amniote alignment using the Cactus
172 alignment software (**Supplementary Table 2**)³⁷. Rather than aligning to a single anchor genome,
173 Cactus infers an ancestral genome for each pair of assemblies (**Figure 3A**). Consistent with our
174 predictions, we have increased power to detect sequence constraint at individual bases relative to
175 earlier studies. We detect 3.1% of bases in the human genome to be under purifying selection in
176 the eutherian lineage (FDR < 5%) without using windowing or other means to integrate
177 contextual information across neighbouring bases. This is more than double the number from the
178 largest previous 100-vertebrate alignment (**Figure 3B**), with improvement most notable in
179 noncoding sequence (**Figure 3C**), and in increased resolution of individual features (**Figure 3D**).
180 This represents a substantial proportion, but not all, of the 5 to 8% of the human genome
181 suggested to be under purifying selection^{7,38}.

182 *Next steps*

183 Using our alignment of 240 mammalian genomes, we are pursuing four key analysis strategies.
184 (1) **Largest nuclear genome eutherian phylogeny**: build a comprehensive phylogeny and
185 timetree, including trees partitioned by functional annotations, mode of inheritance, and long-
186 term recombination rates. (2) **Detailed map of evolutionary constraint**: identify highly
187 conserved genomic regions, regions under accelerated evolution in particular lineages, and
188 changes that likely impact phenotype, leveraging functional data from ENCODE³⁹, GTEx⁴⁰ and
189 the Human Cell Atlas⁴¹. (3) **Genotype-phenotype correlations**: investigate patterns of constraint

190 in human disease-associated regions, identify patterns of convergent adaptive evolution², and
191 apply a forward genomics strategy to link functional elements to traits. (4) **Evolution of genome**
192 **structure**: map syntenic regions between genomes, identify evolutionary breakpoints, and
193 characterize the repeat landscape.

194 **Conclusion**

195 The Zoonomia Project has captured mammalian diversity at unprecedented scope, and is among
196 the first of many projects underway to sequence, catalog, and characterize whole branches of
197 Earth's eukaryotic biodiversity. Based on our experience, we propose the following principles
198 for realizing the full value of large-scale comparative genomics:

199 (1) **Prioritizing sample collection**: We must support field researchers who collect samples and
200 understand species ecology and behaviour, develop strategies for sample collection absent bulky
201 laboratory equipment or cold chains, develop technology for using non-invasive sample types,
202 and establish more repositories of renewable cell cultures¹⁰.

203 (2) **Accessible, scalable tools for computational analysis**: Few research groups have access to
204 computational resources necessary for work with massive genomic datasets. We must address
205 the shortage of skilled computational scientists, and design software and data-storage systems to
206 make powerful computational pipelines accessible to all researchers.

207 (3) **Rapid data-sharing**: Data embargoes must not be permitted to delay analyses that directly
208 benefit conservation of endangered species, human health, or progress in basic science. Genomic
209 data should be shared as quickly as possible, and without restrictions on use.

210 Numerous large-scale genome sequencing efforts are now underway, including the Earth
211 BioGenome Project⁴², Genome 10K⁴³, the Vertebrate Genomes Project, Bat 1K⁴⁴, Bird 10K, and
212 DNA Zoo. As the number of genomes grows, so will the usefulness of comparative genomics in
213 disease research and therapeutic development. Preserving, rather than merely recording, Earth's
214 biodiversity must be a priority. Through global scientific collaborations, and by making genomic
215 resources available and accessible to all research communities, we can ensure that the legacy of
216 genomics is not a digital archive of lost species.

217

218

219 **Figure 1. The Zoonomia Projects brings the fraction of eutherian families represented by at**
220 **least one assembly to 83%.**

221 Phylogenetic tree of the mammalian families in the Zoonomia Project alignment, including both
222 our new assemblies and all other high-quality mammalian genomes publicly available in
223 Genbank when we started the alignment (March 2018; Supplementary Table 2). Tree topology is
224 based on data from timetree.org⁴⁵. Existing taxonomic classifications recognize a total of 127
225 extant eutherian mammalian families⁴⁶, including 43 families not previously represented in

226 Genbank (red boxes) and 41 families with additional representative genome assemblies (pink
227 boxes). Of the remaining families, 21 had Genbank genome assemblies but no Zoonomia Project
228 assembly (grey boxes) and 22 had no representative genome assembly (white boxes).
229 Parenthetical numbers indicate the number of species with genome assemblies in a given family.
230 Image credits: Fossa: Bertal/Wikimedia [CC BY-SA]; Arctic fox: Michael
231 Haferkamp/Wikimedia [CC BY-SA]; (F) Hirola: JRProbert/Wikimedia [CC BY-SA];
232 Bumblebee bat: Sébastien J. Puechmaille [CC BY-SA]; Snowshoe hare: Denali National Park
233 and Preserve/Wikimedia [Public domain]; Aye-aye: Tom Junek/Wikimedia [CC BY-SA];
234 Geoffroy's spider monkey: Patrick Gijsbers/Wikimedia [CC BY-SA]; Southern three-banded
235 armadillo: Hedwig Storch/Wikimedia [CC BY-SA]; Giant Anteater: Graham Hughes/Wikimedia
236 [CC BY-SA]; Brown-throated Sloth: Dick Culbert from Gibsons, B.C., Canada/Wikimedia [CC
237 BY].

238 **Figure 2. Genetic diversity varies across IUCN conservation categories.**

239 **(A)** Heterozygosity declines and **(B)** SoH increases with level of concern for species
240 conservation, as assessed by IUCN conservation categories. Horizontal gray lines indicate
241 median. Comparing individuals sampled from wild and captive populations, we saw no
242 statistically significant difference (independent samples t-test) in either **(C)** overall
243 heterozygosity or **(D)** % segments of homozygosity, with similar means (horizontal gray lines)
244 between birth population types. For A-D, total n=105 species, with n for each tested category
245 indicated on the x axis. Statistical tests were two-sided. **(E)** Overall heterozygosity and SoH for
246 all genomes analyzed (including those with high allelic balance ratio; n=124 species), with
247 median SoH (17.1%, horizontal dashed line) and median overall heterozygosity (0.0026, vertical
248 dashed line) for species categorized as Least Concern (dashed lines). Values for individuals from
249 the seven critically endangered species are shown in red, with red text labels.

250 **Figure 3. The Zoonomia alignment doubles the fraction of the human genome predicted to**
251 **be under purifying selection at single base pair resolution.**

252 **(A)** Cactus alignments are reference-genome-free, enabling detection of sequence absent from
253 human (red) or other clades (purple), lineage-specific innovations (orange, green) and eutherian
254 mammal-specific sequence (blue). **(B)** We compared phyloP predictions of conserved positions
255 for a widely-used 100-vertebrate alignment (n=100 vertebrate species; grey) to the Zoonomia
256 alignment (n=240 eutherian species; red). The cumulative portion of the genome expected to be
257 covered by true vs. false positive calls is shown, starting from the highest confidence calls (solid
258 line) and proceeding to calls with lower confidence (dashed lines), with a horizontal line
259 indicating the point at which the confidence level drops below an expected FDR of 0.05 (two-
260 sided). **(C)** A higher proportion of functionally annotated bases are detected as highly conserved
261 (FDR<0.05) in the Zoonomia alignment (red) than the 100-vertebrate alignment (grey), most
262 notably in noncoding regions. **(D)** At a putative androgen receptor binding site overlapping a
263 ChIP-seq peak and a phastCons constrained element prediction, phyloP scores predict seven
264 bases are under purifying selection in the Zoonomia alignment (red; FDR=0.05; two-sided),

265 peaking in positions with the most information content in the androgen receptor JASPAR⁴⁷
266 motif, compared to one (grey) in the 100-vertebrate alignment, with scores at FDR > 0.05 in light
267 red and light grey.

268

269 **Methods**

270 *Species selection, sample shipping, and regulatory approvals.*

271 Species were selected to maximize branch length across the eutherian mammal phylogeny, and
272 to capture genomes of species from previously unrepresented eutherian families. Of 172 species
273 initially selected for inclusion, we obtained sufficiently high quality DNA samples for genome
274 sequencing for 137. DNA samples from collaborating institutions were shipped to either the
275 Broad Institute (N=69) or Uppsala University (N=68). For samples received at Broad then sent to
276 Uppsala, shipping approval was secured from the US Fish and Wildlife Service. IACUC
277 approval was not required.

278 *Sample quality control, library construction, and sequencing.*

279 DNA integrity for each sample was visualized via agarose gel (at Broad) or Agilent tape station
280 (at Uppsala). Samples passed QC if the bulk of DNA fragments were greater than 5kb. DNA
281 concentration was then determined using Invitrogen Qubit dsDNA HS assay kit. For each of the
282 samples that passed QC, 1-3µg of DNA was fragmented on the Covaris E220 Instrument using
283 the 400bp standard program (10% Duty cycle, 140 PIP, 200 cycles per burst, 55s). Fragmented
284 samples underwent SPRI double size selection (0.55X, 0.7Xf) followed by PCR-free Illumina
285 library construction following the manufacturer's instructions (Kapa #KK8232) using PCR-free
286 adapters from Illumina (#FC-121-3001). Final library fragment size distribution was determined
287 on Agilent 2100 Bioanalyzer with High Sensitivity DNA Chips. Paired-end libraries were
288 pooled, then sequenced on a single-lane of the Illumina HiSeq2500, set for Version 2 chemistry
289 and 2x250bp reads. This yielded a total of mean 375M (standard deviation = 125M) reads per
290 sample.

291 *Assembly and validation*

292 For each species, we applied DISCOVAR *de novo*¹¹ (discoverdenovo-52488;
293 <ftp://ftp.broadinstitute.org/pub/crd/DiscoverDeNovo/>) to assemble the 2x250bp read group, using
294 the following command: DiscoverDeNovo READS=[READFILE]
295 OUT_DIR=[SPECIES_ID]/[SPECIES_ID].discover_files NUM_THREADS=24
296 MAX_MEM_GB=200G

297 Coverage for each genome was automatically calculated by DISCOVAR, with a mean coverage
298 of 40.1x (s.d. +/- 14x). We assessed genome assembly, gene set, and transcriptome completeness
299 using BUSCO, which provides quantitative measures based on gene content from near-universal
300 single-copy orthologs⁴⁸. BUSCO was run with default parameters, using the mammalian gene

301 model set (mammalia_odb9, n=4104), using the following command: python ./BUSCO.py -i
302 [input fasta] -o [output_file] -l ./mammalia_odb9/ -m genome -c 1 -sp human.

303 Median contig N50 for existing RefSeq assemblies was calculated using the assembly statistics
304 for the most recent release of 118 eutherian mammals with RefSeq assembly accession numbers.
305 Assemblies were all classified as either “Reference Genome” or “Representative Genome”.
306 Assembly statistics downloaded from NCBI on April 10, 2019.

307 *Genome upgrades.* We selected genomes from each eutherian order without a preexisting long-
308 contiguity assembly based on (1) whether the underlying assembly met the minimum quality
309 threshold needed for HiRise upgrades; (2) whether a second sample of sufficient quality could be
310 obtained from that individual. All upgrades were done with Dovetail Chicago libraries and
311 assembled with HiRise 2.1, as previously described⁴⁹.

312 *Estimating heterozygosity*

313 *Selection of assemblies for heterozygosity analysis.* Heterozygosity statistics were calculated for
314 all but four of our short read assemblies (N=126) as well as 8 Dovetail-upgraded genomes. Four
315 failed because they were either too fragmented to analyze (N=3) or due to undetermined errors
316 (N=1). One assembly was excluded because it was a second individual from an already
317 represented species.

318 *Heterozygosity calls.* We applied the standard GATK pipeline with genotype quality banding to
319 identify the callable fraction of the genome^{50,51}. First, we used *samtools* to subsample paired
320 reads from the unmapped bam files⁵². After removing adapter sequences from the selected reads,
321 we used BWA-MEM to map reads to the reference genome scaffolds of >10kb, removing
322 duplicates using the PicardTools MarkDuplicates utility⁵³. We then called heterozygous sites
323 using standard GATK-Haplotypecaller specifications, and with additional gVCF banding at 0,
324 10, 20, 30, 40, 50 and 99 qualities. We used the fraction of the genome with genotype quality
325 >15 callable for subsequent analyses. For the lists of high-confidence variant sites, we include
326 only heterozygous positions after filtering at GQ>20, max DP<100, min DP>6, as described in
327 the README file at broad.io/variants.

328 *Inferring overall heterozygosity.* To avoid confounding by sex chromosomes or complex regions,
329 we excluded all scaffolds with less than 0.5 or greater than 2x of the average sample read depth,
330 then calculated global heterozygosity as the fraction of heterozygous calls over the whole
331 callable genome.

332 *Calling Segments of Homozygosity (SoH).* We estimated the proportion of the genome within
333 segments of homozygosity (SoH) using a metric designed for genomes with scaffold N50 shorter
334 than the expected maximum length of runs of homozygosity (our median scaffold N50 is 62kb).
335 We first split all scaffolds into windows with a maximum length of 50kb, with windows ranging
336 from 20kb-50kb for scaffolds <50kb. For each window, we calculated the average number of
337 heterozygous sites per bp. We discriminated windows with extremely low heterozygosity by
338 using the Python 3.5.2 pomegranate package to fit a two-component Gaussian Mixture Model to

339 the joint distribution of window heterozygosity, forcing the first component to be centered
340 around the lower tail of the distribution and allowing the second to freely capture all the
341 remaining heterozygosity variability^{54,55}. As heterozygosity cannot be negative, and normal
342 distributions near zero can cross into negative values, we used the normal cumulative distribution
343 function to correct the posterior distribution by the negative excess -- effectively fitting a
344 truncated normal to the first component. The final SoH value was calculated using the posterior
345 maximum likelihood classification between both components. We see no significant correlation
346 between contig N50 and SoH (Pearson correlation=0.055, p=0.57, N=112).

347 *Assessing the impact of % callable genome.* We assessed whether the % of the genome that was
348 callable (Supplementary Table 3) was likely to impact our analysis. The callable % was
349 correlated with heterozygosity (r=-0.80, p<2.2e-16, N=130), and weakly with SoH (r=0.18,
350 p=0.06, N=112). There is no significant difference in callable % among IUCN categories
351 (p_{anova}=0.98; N=122) or between captive and wild populations (p_{t-test}=0.81; N=120).

352 *Analyzing patterns of diversity.* We excluded two genomes with exceptionally high
353 heterozygosity (het > 0.02; > 5 standard deviations above the mean). Both were non-endangered,
354 and thus removing them made our determination of lower heterozygosity in endangered species
355 more conservative. Of the remaining 124, we excluded 19 genomes with allelic balance (ab)
356 values more than one standard deviation above the mean (>0.36). Abnormally high ab can
357 indicate sequencing biases with potential for artifacts in estimates of heterozygosity and/or SoH.
358 Our final dataset contains heterozygosity values for 105 genomes and SoH values for 98
359 genomes (Supplementary Table 3). For seven genomes, we were unable to estimate SoH because
360 the two components of the Gaussian Mixture Model overlapped completely. To ask about a
361 possible directional relationship between level of IUCN concern and overall heterozygosity or
362 SoH, we applied regression using IUCN category as an ordinal predictor. We also asked about
363 the relationship of diversity metrics to a set of species-level phenotypes for which correlations
364 were previously reported (Extended Data Table 3).

365 ***Alignment***

366 The alignment was generated using the progressive mode of Cactus^{37,56}. The topology used for
367 the guide-tree of the alignment was taken from TimeTree⁴⁵; the branch lengths of the guide-tree
368 were generated by a least-squares fit from a distance matrix. The distance matrix was based on
369 the UCSC 100-way phyloP fourfold-degenerate site tree⁵⁷ for those species which had
370 corresponding entries in the 100-way. For species not present in the 100-way, distance matrix
371 entries were more coarsely estimated using the distance estimated from Mash⁵⁸ to the closest
372 relative included in the 100-way data.

373 Cactus does not attempt to fully resolve the gene tree when multiple duplications take place
374 along a single branch, as there is an implicit restriction in Cactus that a duplication event be
375 represented as multiple regions in the child species aligned to a single region in the parent
376 species. This precludes representing discordance between gene tree and species tree that could
377 occur with either incomplete lineage-sorting or horizontal transfer. However, the guide tree has

378 minimal impact on the alignment, with little difference between alignments with different trees,
379 even when using a tree that is purposely wrong³⁷. Phenomena such as incomplete lineage sorting
380 that affect a subset of species are unlikely to substantially impact the detection of purifying
381 selection across the whole eutherian lineage described in Figure 3.

382 The alignment was generated in several steps on account of its large scale. First, a “backbone”
383 alignment of several long contiguity assemblies was generated. Next, separate clade alignments
384 were generated for each major clade in the alignment: Euarchonta, Glires, Laurasiatheria, and
385 Afrotheria/Xenarthra. The roots of these clade alignments were then aligned to the corresponding
386 ancestral genomes from the backbone, “stitching” these alignments together to create the final
387 alignment. (The process of aligning a genome to an existing ancestor is complex and further
388 described in the preprint introducing the progressive mode of Cactus³⁷).

389 We created a neutral model for the conservation analysis using ancestral repeats detected by
390 RepeatMasker⁵⁹ on the eutherian ancestral genome produced in the Cactus alignment (tRNA and
391 “low complexity” repeats were removed). To fit the neutral model, we used phyloFit from the
392 PHAST⁶⁰ package, using the REV (generalized reversible) model and EM optimization method.
393 The training input was a MAF exported on columns from the set of ancestral repeats mentioned
394 above. Since phyloFit does not support alignment columns containing duplicates, if a genome
395 had more than one sequence in a single alignment block, they were replaced with a single entry
396 representing the consensus base at each column.

397 We extracted initial conservation scores using PhyloP from the PHAST⁶⁰ package on a MAF
398 exported using human as a reference. We converted the PhyloP scores (which represent log-
399 scaled p-values of acceleration or conservation) into p-values, then into q-values using the FDR-
400 correction of Benjamini and Hochberg⁶¹. Any column with a resulting q-value less than 0.05 was
401 deemed significantly conserved or accelerated.

402 The alignment, as well as conservation annotations, are available on a UCSC Assembly Hub⁶²
403 hosted at broad.io/genomes (the link may be loaded into the “Track Hubs” section of the
404 browser) and at <https://alignment-output.s3.amazonaws.com/200m-v1.hal>.

405

406 *Acknowledgements*

407 We thank the many individuals who provided samples and advice, including Chris Adenyo,
408 Catherine Avila, Eric Baitchman, Richard Behringer, Adam Boyko, Matthew Breen, Kevin
409 Campbell, Polly Campbell, Chris J. Conroy, Kimberly Cooper, Liliana M. Dávalos, Frederic
410 Delsuc, Daniel Distel, Christopher Allan Emerling, Julie Fronczek, Neil Gemmel, Jeffrey Good,
411 Kai He, Kris Helgen, Allyson Hindle, Hopi Hoekstra, Rodney Honeycutt, Pavel Hulva, William
412 Israelsen, Boniface Kayang, Rosalind Kennerley, Marisa Korody, Danielle N. Lee, Edward
413 Louis, Matt MacManes, Ann Misuraca, Anna Mitelberg, Phillip Morin, Alice Mouton, Miho
414 Murayama, Michael Nachman, Asako Navarro, Rob Ogden, Bret Pasch, Sebastien Peuchmaille,
415 TJ Robinson, Stephen Rossiter, Manuel Ruedi, Ashley Seifert, Steven Thomas, Samuel Turvey,

416 Goedele Verbeylen, and the late Dr. R. J. Baker. We also thank the Broad Institute Genomics
417 Platform and SNP&SEQ Technology Platform (part of the National Genomics Infrastructure
418 (NGI) Sweden and Science for Life Laboratory).

419 This project was funded by NIH NHGRI R01HG008742 (EKK, BB, DPG, RS, JTM, JJ, HN, BP,
420 JA), Swedish Research Council Distinguished Professor Award (KLT, VDM, EM, JRSM), Knut
421 and Alice Wallenberg Foundation (KLT, VDM, EM, JRSM), Uppsala University (KLT, VDM,
422 EM, JRSM, JJ, JeA, LG), Broad Institute Next10 (LG), Gladstone Institutes (KSP), NIH NHGRI
423 5R01HG002939 (ArS, RH), NIH NHGRI 5U24HG010136 (ArS, RH), NIH NHGRI
424 5R01HG010485 (BP, MD), NIH NHGRI 2U41HG007234 (BP, MD, JA), NIH NIA
425 5PO1AG047200 (VNG), NIH NIA 1UH2AG064706 (VNG), BFU2017-86471-P
426 MINECO/FEDER, UE (TM), Secretaria d'Universitats i Recerca and CERCA Programme del
427 Departament d'Economia i Coneixement de la Generalitat de Catalunya GRC 2017 SGR 880
428 (TM), Howard Hughes International Early Career (TM), European Research Council Horizon
429 2020 #864203 (TM), Obra Social "La Caixa" (TM), BBSRC BBS/E/T/000PR9818, BBS/E/T/
430 000PR9783 (WH, WN), BBSRC Core Strategic Programme Grant BB/P016774/1 (WH, WN,
431 FD), Sir Henry Dale Fellowship 200517/Z/16/Z jointly funded by the Wellcome Trust and the
432 Royal Society to NRC (NRC), FJCI-2016-29558 MICINN (DJ), Prince Albert II Foundation of
433 Monaco and Canada, Global Genome Initiative, Smithsonian Institution (MN), European
434 Research Council Research Grant ERC-2012-StG311000 (ECT), UK Medical Research Council
435 MR/P026028/1 (WH, WN), National Science Foundation DEB-1457735 (MS), National Science
436 Foundation DEB-1753760 (WJM), National Science Foundation IOS-2029774 (EKK, DPG),
437 Robert and Rosabel Osborne Endowment (HL, JD), Swedish Research Council, FORMAS 221-
438 2012-1531 (JRSM), NSF RoL: FELLS: EAGER: DEB 1838283 (DR), Academy of Finland
439 Finnish Center of Excellence #312041 (TK).

440

441 **Consortium**

442 Diane P. Genereux¹, Aitor Serres^{2,3}, Joel Armstrong⁴, Jeremy Johnson¹, Voichita D. Marinescu⁵,
443 Eva Murén⁵, David Juan^{2,3}, Gill Bejerano^{6,7,8,9}, Nicholas R. Casewell¹⁰, Leona G. Chemnick¹¹,
444 Joana Damas¹², Federica Di Palma^{13,14}, Mark Diekhans⁴, Ian T. Fiddes⁴, Manuel Garber¹⁵,
445 Vadim N. Gladyshev^{1,16}, Linda Goodman^{1,17}, Wilfried Haerty¹⁴, Marlys L Houck¹¹, Robert
446 Hubley¹⁸, Teemu Kivioja^{19,20,21}, Klaus-Peter Koepfli²², Lukas F. K. Kuderna³, Eric S.
447 Lander^{1,23,24}, Jennifer R. S. Meadows⁵, William J. Murphy²⁵, Will Nash¹⁴, Hyun Ji Noh¹, Martin
448 Nweeia^{26,27,28}, Andreas R. Pfenning²⁹, Katherine S. Pollard^{30,31,32}, David Ray³³, Beth Shapiro^{34,35},
449 Arian Smit¹⁸, Mark Springer³⁶, Cynthia C. Steiner¹¹, Ross Swofford¹, Jussi Taipale^{19,21,37}, Emma
450 C. Teeling³⁸, Jason Turner-Maier¹, Jessica Alfoldi¹, Bruce Birren¹, Oliver A. Ryder¹¹, Harris
451 Lewin^{12,39}, Benedict Paten⁴, Tomas Marques-Bonet^{2,3,40,41}, Kerstin Lindblad-Toh^{1,5}, Elinor K.
452 Karlsson^{1,15,42}

453 (1) Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA; (2) CNAG-CRG,
454 Centre for Genomic Regulation (CRG), Barcelona Institute of Science and Technology (BIST),

455 Barcelona, Spain; (3) Institute of Evolutionary Biology (UPF-CSIC), PRBB, Barcelona,
456 Catalonia, Spain; (4) Center for Biomolecular Science and Engineering, University of California
457 Santa Cruz, Santa Cruz, CA, USA; (5) Science for Life Laboratory, Department of Medical
458 Biochemistry and Microbiology, Uppsala University, Uppsala, Sweden; (6) Department of
459 Biomedical Data Science, Stanford University, Stanford, CA, USA; (7) Department of Computer
460 Science, Stanford University, Stanford, CA, USA; (8) Department of Developmental Biology,
461 Stanford University, Stanford, CA, USA; (9) Department of Pediatrics, Stanford University,
462 Stanford, CA, USA; (10) Liverpool School of Tropical Medicine, Liverpool, UK; (11) San
463 Diego Zoo Institute for Conservation Research, San Diego, California, USA; (12) The UC Davis
464 Genome Center, University of California, Davis, California, USA; (13) Department of Biological
465 Sciences, University of East Anglia, Norwich, UK; (14) Earlham Institute, Norwich Research
466 Park, Norwich, UK; (15) Program in Bioinformatics and Integrative Biology, University of
467 Massachusetts Medical School, Worcester, Massachusetts, USA; (16) Brigham and Women's
468 Hospital, Harvard Medical School, Boston, Massachusetts, USA; (17) Stanford University,
469 Stanford, California, USA; (18) Institute for Systems Biology, Seattle, Washington, USA; (19)
470 Department of Biochemistry, University of Cambridge, Cambridge, United Kingdom; (20)
471 Department of Medical Biochemistry and Biophysics, Karolinska Institutet, Stockholm, Sweden;
472 (21) University of Helsinki, Helsinki, Finland; (22) Smithsonian Conservation Biology Institute,
473 Center for Species Survival, National Zoological Park, Front Royal, Virginia and Washington,
474 DC, USA; (23) Department of Biology, MIT, Cambridge, MA, USA; (24) Department of
475 Systems Biology, Harvard Medical School, Boston, MA, USA; (25) Veterinary Integrative
476 Biosciences, Texas A&M University, College Station, Texas, USA; (26) Marine Mammal
477 Program, Smithsonian Institution, Washington, DC, USA; (27) Restorative Dentistry and
478 Biomaterials Sciences, Harvard School of Dental Medicine, Boston, Massachusetts, USA; (28)
479 School of Dental Medicine, Case Western Reserve University, Cleveland, Ohio, USA; (29)
480 Carnegie Mellon University, School of Computer Science, Department of Computational
481 Biology, Pittsburgh, Pennsylvania, USA; (30) Chan-Zuckerberg Biohub, San Francisco,
482 California, USA; (31) Gladstone Institutes, San Francisco, California, USA; (32) University of
483 California, San Francisco, California, USA; (33) Department of Biological Sciences, Texas Tech
484 University, Lubbock, Texas, USA; (34) Ecology and Evolutionary Biology, University of
485 California, Santa Cruz, California, USA; (35) Howard Hughes Medical Institute, Chevy Chase,
486 Maryland, USA; (36) University of California, Riverside, California, USA; (37) Karolinska
487 Institute, Solna, Sweden; (38) School of Biology and Environmental Science, University College
488 Dublin, Dublin, Ireland; (39) Evolution and Ecology, University of California, Davis, California,
489 USA; (40) Catalan Institution of Research and Advanced Studies (ICREA), Barcelona, Spain;
490 (41) Institut Català de Paleontologia Miquel Crusafont, Universitat Autònoma de Barcelona,
491 Edifici ICTA-ICP, Barcelona, Spain; (42) Program in Molecular Medicine, University of
492 Massachusetts Medical School, Worcester, Massachusetts, USA

494 **Author contributions**

495 These authors contributed equally: Kerstin Lindblad-Toh and Elinor K. Karlsson

496 KLT, conceived the project, JJ, VDM, EM, NRC, LGC, JD, VNG, MLH, KPK, JRSM, WJM,
497 MN, DR, RS, ECT, JeA, OAR, HL, KLT and EKK, contributed to the acquisition of the
498 samples. JJ, VDM, EM, JD, LG, KPK, HN, CCS, RS, JTM, JeA, OAR, HL, KLT and EKK,
499 contributed to the production of the genome assemblies. DPG, AS, JA, JJ, DJ, ITF, LFKK, HL,
500 TM, KLT and EKK, contributed to the data analysis.. DPG, JJ, VDM, GB, FD, MD, ITF, MG,
501 VNG, WH, RH, TK, ESL, JRSM, ARP, KSP, ArS, MS, JT, JeA, BB, OAR, HL, BP, TM, KLT
502 and EKK contributed to the design and conduct of the project. DPG, ESL, WN, BS, OAR, KLT,
503 and EKK, wrote the manuscript, with input from all other authors

504

505 **Competing interests**

506 LG is a co-founder of, equity owner in and chief technical officer at Fauna Bio Incorporated.

507

508 **Additional Information:**

509 Supplementary Information is available for this paper.

510 Correspondence and requests for materials should be addressed to elinor@broadinstitute.org.

511 Reprints and permissions information is available at www.nature.com/reprints.

512

513 **Data Availability**

514 Details on each Zoonomia Project genome assembly, including NCBI Genbank⁶³ accession
515 numbers, are in **Supplementary Table 1**. Sequence data and genome assemblies are available at
516 <https://www.ncbi.nlm.nih.gov/>. Variant lists for each species are at <http://broad.io/variants>.
517 Source data for Figure 2 is provided and in the Zoonomia github repository (DOI
518 10.5281/zenodo.3887432). The Cactus alignment is at [https://alignment-](https://alignment-output.s3.amazonaws.com/200m-v1.hal)
519 [output.s3.amazonaws.com/200m-v1.hal](https://alignment-output.s3.amazonaws.com/200m-v1.hal). A visualization of the alignments and PhyloP data is
520 available by loading our assembly hub into the UCSC browser⁶² by copying the hub link
521 https://comparative-genomics-hubs.s3-us-west-2.amazonaws.com/200m_hub.txt into the "Track
522 Hubs" page. There are no restrictions on use.

523 **Code Availability**

524 The Discover *de novo* assembly code is available at
525 https://github.com/broadinstitute/discover_de_novo/releases/tag/v52488 (DOI
526 10.5281/zenodo.3870889), the Cactus pipeline is available at
527 <https://github.com/ComparativeGenomicsToolkit/cactus> (DOI 10.5281/zenodo.3873410) and

528 code for other analyses is available at <https://github.com/broadinstitute/Zoonomia/> (DOI
529 10.5281/zenodo.3887432).

530

531

532 **Extended Data Figure 1. Remarkable traits in non-human mammals.**

533 Sequences from species with remarkable phenotypes can inform human medicine, basic biology,
534 and biodiversity conservation, but sample collection can be challenging. **(A)** The Jamaican fruit
535 bat (*Artibeus jamaicensis*) maintains constant blood glucose across intervals of fruit-eating and
536 fasting⁶⁴, achieving homeostasis to a degree elusive in treatment of human diabetes. **(B)** The
537 North American beaver (*Castor canadensis*) avoids tooth decay by incorporating iron, rather
538 than magnesium, into tooth enamel, yielding an orange hue⁶⁵. **(C)** The thirteen-lined ground
539 squirrel (*Ictidomys tridecemlineatus*) prepares for hibernation by rapidly increasing the
540 thermogenic activity of brown fat⁶⁶, a process connected to improved glucose homeostasis and
541 insulin sensitivity in humans⁶⁷⁻⁶⁹. **(D)** The tiny bumblebee bat (*Craseonycteris thonglongyai*) is
542 among the smallest of mammals, making it a sparse source of DNA. **(E)**. The remote habitat of
543 the very rare Amazon River dolphin (*Inia geoffrensis*) precludes collection of the high-molecular
544 weight DNA. Image sources: (A) Merlin D. Tuttle/Science Source; (B) Stephen J.
545 Krasemann/Science Source; (C) Allyson Hindle; (D) Sébastien J. Puechmaille [CC BY-SA]; (E)
546 M. Watson/Science Source.

547

548 **Extended Data Figure 2. Sample collection can be challenging, and sequencing methods**

549 **must be selected to handle the sample quality.** To enable inclusion of species from across the
550 eutherian tree, including from the 50% of mammalian families not represented in existing
551 genome databases, the Zoonomia Project needed sequencing and assembly methods that produce
552 reliable data from DNA collected in remote locations, sometimes in only modest quantities, and
553 often without benefit of cold chains for transport. **(a)** For the marine species like the narwhal
554 (*Monodon monoceros*), simply accessing an individual in the wild can prove challenging. To
555 sample DNA from the near-threatened narwhal, for example, Martin Nweeia and Inuit guide
556 David Angnatsiak camped on an ice floe edge between Pond Inlet and Bylot Island, at the
557 northeastern tip of Baffin Island. After a narwhal was harvested by Inuit hunters as part of an
558 annual hunt, hours of flensing were necessary for collecting tissue samples. Shown, from left to
559 right: Frank McCann, Hans Christian Schmidt, Frederick Eichmiller, Martin Nweeia, James Orr
560 (facing backward), and Jack Orr (standing). **(b)** For endangered species like the Hispaniolan
561 solenodon (*Solenodon paradoxus*), sample collection must be designed to minimize stress to the
562 individual, limiting the amount of DNA that can be collected²². To collect DNA from the
563 endangered solenodon without imposing stress on individuals in the wild, Nicholas Casewell
564 turned to the world's only captive solenodons, housed off-exhibit at ZOODOM in the Dominican
565 Republic. With help from Zoo veterinarians, Casewell collected a small amount of blood from

566 the solenodon's rugged tail. Narwhal photograph by Gretchen Freund and courtesy of Martin
567 Nweeia. Solenodon photo courtesy of Lucy Emery.

568

569 **Extended Data Table 1. The Zoonomia Project data includes 132 genome assemblies.** These
570 assemblies include 131 different species, with two narwhals (male and female), and 10 genomes
571 upgraded to longer contiguity (including upgrade of an existing assembly for *Echinops telfairi*).
572 Species of concern on the IUCN Red List are indicated as Near Threatened (NT), Vulnerable(V),
573 Endangered(EN) or Critically Endangered (CR). * upgraded to longer contiguity; † upgraded to
574 longer contiguity using existing assembly.

575

576 **Extended Data Table 2: Power to detect constraint across data sets.** The expected number of
577 variants conserved by chance (false positives) was estimated for four genomic resources (the 29
578 Mammals Project⁷ dataset, the human only ExAC¹⁴ and gnomAD v3⁷⁰ datasets, and the
579 Zoonomia Project dataset) by applying a Poisson model of the distribution of substitution counts
580 in the genome. Branch length for gnomAD was estimated by dividing 526,001,545 single
581 nucleotide variants by 3.088 gigabases (human genome size). Branch length for Zoonomia was
582 measured as substitutions/site in the phyloP analysis of the Cactus alignment.

583

584 **Extended Data Table 3. Diversity statistics are not correlated with other species-level**
585 **phenotypes.** All phenotypes in the Pantheria database³⁰ for which at least 75% of the 75 “Least
586 Concern” species had a value were included in the analysis. For continuous phenotypes, values
587 were standardized to Z-scores prior to analysis (latitude was calculated as an absolute value) and
588 correlation measured by fitting a linear model using the core R function lm. For categorical
589 phenotypes with more than two categories, group means were compared using the core R
590 function aov to fit an analysis of variance model. None were significant after Bonferroni
591 correction for the number of traits considered (21).

592

593

594 **References**

- 595 1. Claussnitzer, M. *et al.* A brief history of human disease genetics. *Nature* **577**, 179–189 (2020).
- 596 2. Hiller, M. *et al.* A ‘forward genomics’ approach links genotype to phenotype using independent
597 phenotypic losses among related species. *Cell Rep.* **2**, 817–823 (2012).
- 598 3. Wasser, S. K. *et al.* Genetic assignment of large seizures of elephant ivory reveals Africa’s major
599 poaching hotspots. *Science* **349**, 84–87 (2015).
- 600 4. Wright, B. *et al.* Development of a SNP-based assay for measuring genetic diversity in the
601 Tasmanian devil insurance population. *BMC Genomics* **16**, 791 (2015).
- 602 5. Lappalainen, T., Scott, A. J., Brandt, M. & Hall, I. M. Genomic Analysis in the Age of Human
603 Genome Sequencing. *Cell* **177**, 70–84 (2019).
- 604 6. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic
605 variants. *Nat. Genet.* **46**, 310–315 (2014).
- 606 7. Lindblad-Toh, K. *et al.* A high-resolution map of human evolutionary constraint using 29 mammals.
607 *Nature* **478**, 476–482 (2011).
- 608 8. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide
609 association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
- 610 9. IUCN. The IUCN Red List of Threatened Species. Version 2019-2. <http://www.iucnredlist.org>
611 (2019).
- 612 10. Ryder, O. A. & Onuma, M. Viable Cell Culture Banking for Biodiversity Characterization and
613 Conservation. *Annu Rev Anim Biosci* **6**, 83–98 (2018).
- 614 11. Weisenfeld, N. I. *et al.* Comprehensive variation discovery in single human genomes. *Nat. Genet.*
615 **46**, 1350–1355 (2014).
- 616 12. Putnam, N. H. *et al.* Chromosome-scale shotgun assembly using an in vitro method for long-range
617 linkage. *Genome Res.* **26**, 342–350 (2016).
- 618 13. Kim, J. *et al.* Reconstruction and evolutionary history of eutherian chromosomes. *Proc. Natl. Acad.*
619 *Sci. U. S. A.* **114**, E5379–E5388 (2017).

- 620 14. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291
621 (2016).
- 622 15. Balasubramanian, S. *et al.* Using ALoFT to determine the impact of putative loss-of-function
623 variants in protein-coding genes. *Nat. Commun.* **8**, 382 (2017).
- 624 16. Meadows, J. R. S. & Lindblad-Toh, K. Dissecting evolution and disease using comparative
625 vertebrate genomics. *Nat. Rev. Genet.* **18**, 624–636 (2017).
- 626 17. Cooper, G. M. & Shendure, J. Needles in stacks of needles: finding disease-causal variants in a
627 wealth of genomic data. *Nat. Rev. Genet.* **12**, 628–640 (2011).
- 628 18. Baiz, M. D., Tucker, P. K. & Cortés-Ortiz, L. Multiple forms of selection shape reproductive
629 isolation in a primate hybrid zone. *Mol. Ecol.* (2018) doi:10.1111/mec.14966.
- 630 19. Dobzhansky, T. & Dobzhansky, T. G. *Genetics and the Origin of Species*. (Columbia University
631 Press, 1937).
- 632 20. Abegglen, L. M. *et al.* Potential Mechanisms for Cancer Resistance in Elephants and Comparative
633 Cellular Response to DNA Damage in Humans. *JAMA* **314**, 1850–1860 (2015).
- 634 21. Herrera-Álvarez, S., Karlsson, E., Ryder, O. A., Lindblad-Toh, K. & Crawford, A. J. How to make a
635 rodent giant: Genomic basis and tradeoffs of gigantism in the capybara, the world’s largest rodent.
636 *bioRxiv* 424606 (2018) doi:10.1101/424606.
- 637 22. Casewell, N. R. *et al.* Solenodon genome reveals convergent evolution of venom in eulipotyphlan
638 mammals. *Proc. Natl. Acad. Sci. U. S. A.* (2019) doi:10.1073/pnas.1906117.
- 639 23. Beichman, A. C. *et al.* Aquatic adaptation and depleted diversity: a deep dive into the genomes of
640 the sea otter and giant otter. *Mol. Biol. Evol.* (2019) doi:10.1093/molbev/msz101.
- 641 24. Damas, J. *et al.* Broad Host Range of SARS-CoV-2 Predicted by Comparative and Structural
642 Analysis of ACE2 in Vertebrates. *bioRxiv* 2020.04.16.045302 (2020)
643 doi:10.1101/2020.04.16.045302.
- 644 25. Xue, Y. *et al.* Mountain gorilla genomes reveal the impact of long-term population decline and
645 inbreeding. *Science* vol. 348 242–245 (2015).

- 646 26. Ceballos, F. C., Joshi, P. K., Clark, D. W., Ramsay, M. & Wilson, J. F. Runs of homozygosity:
647 windows into population history and trait architecture. *Nat. Rev. Genet.* **19**, 220–234 (2018).
- 648 27. Spielman, D., Brook, B. W. & Frankham, R. Most species are not driven to extinction before genetic
649 factors impact them. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 15261–15264 (2004).
- 650 28. Vinson, J. P. *et al.* Assembly of polymorphic genomes: algorithms and application to *Ciona*
651 *savignyi*. *Genome Res.* **15**, 1127–1135 (2005).
- 652 29. MacManes, M. D. & Lacey, E. A. The social brain: transcriptome assembly and characterization of
653 the hippocampus from a social subterranean rodent, the colonial tuco-tuco (*Ctenomys sociabilis*).
654 *PLoS One* **7**, e45524 (2012).
- 655 30. Jones, K. E. *et al.* PanTHERIA: a species-level database of life history, ecology, and geography of
656 extant and recently extinct mammals: Ecological Archives E090-184. *Ecology* **90**, 2648–2648
657 (2009).
- 658 31. Cardillo, M. Biological determinants of extinction risk: why are smaller species less vulnerable?
659 *Anim. Conserv.* **6**, 63–69 (2003).
- 660 32. Natesh, M. *et al.* Empowering conservation practice with efficient and economical genotyping from
661 poor quality samples. *Methods Ecol. Evol.* **10**, 853–859 (2019).
- 662 33. Lowry, D. B. *et al.* Breaking RAD: an evaluation of the utility of restriction site-associated DNA
663 sequencing for genome scans of adaptation. *Mol. Ecol. Resour.* **17**, 142–152 (2017).
- 664 34. Shapiro, B. Pathways to de-extinction: how close can we get to resurrection of an extinct species?
665 *Funct. Ecol.* **31**, 996–1002 (2017).
- 666 35. Benazzo, A. *et al.* Survival and divergence in a small group: The extraordinary genomic history of
667 the endangered Apennine brown bear stragglers. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E9589–E9597
668 (2017).
- 669 36. Saremi, N. F. *et al.* Puma genomes from North and South America provide insights into the genomic
670 consequences of inbreeding. *Nat. Commun.* **10**, 4769 (2019).
- 671 37. Armstrong, J. *et al.* Progressive alignment with Cactus: a multiple-genome aligner for the thousand-

- 672 genome era. (in submission).
- 673 38. Rands, C. M., Meader, S., Ponting, C. P. & Lunter, G. 8.2% of the Human genome is constrained:
674 variation in rates of turnover across functional element classes in the human lineage. *PLoS Genet.*
675 **10**, e1004525 (2014).
- 676 39. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome.
677 *Nature* **489**, 57–74 (2012).
- 678 40. GTEx Consortium *et al.* Genetic effects on gene expression across human tissues. *Nature* **550**, 204–
679 213 (2017).
- 680 41. Regev, A. *et al.* The Human Cell Atlas. *Elife* **6**, (2017).
- 681 42. Lewin, H. A. *et al.* Earth BioGenome Project: Sequencing life for the future of life. *Proc. Natl. Acad.*
682 *Sci. U. S. A.* **115**, 4325–4333 (2018).
- 683 43. Koepfli, K.-P., Paten, B., Genome 10K Community of Scientists & O’Brien, S. J. The Genome 10K
684 Project: a way forward. *Annu Rev Anim Biosci* **3**, 57–111 (2015).
- 685 44. Teeling, E. C. *et al.* Bat Biology, Genomes, and the Bat1K Project: To Generate Chromosome-Level
686 Genomes for All Living Bat Species. *Annu Rev Anim Biosci* **6**, 23–46 (2018).
- 687 45. Kumar, S., Stecher, G., Suleski, M. & Hedges, S. B. TimeTree: A Resource for Timelines,
688 Timetrees, and Divergence Times. *Mol. Biol. Evol.* **34**, 1812–1819 (2017).
- 689 46. *Mammal Species of the World. A Taxonomic and Geographic Reference (3rd ed).* (Johns Hopkins
690 University Press, 2005).
- 691 47. Vlieghe, D. *et al.* A new generation of JASPAR, the open-access repository for transcription factor
692 binding site profiles. *Nucleic Acids Res.* **34**, D95–7 (2006).

693

694

695

696

697

698 **Extended Data References**

- 699 48. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO:
700 assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*
701 **31**, 3210–3212 (2015).
- 702 49. Farré, M. *et al.* A near-chromosome-scale genome assembly of the gemsbok (*Oryx gazella*): an
703 iconic antelope of the Kalahari desert. *Gigascience* **8**, (2019).
- 704 50. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-
705 generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
- 706 51. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation
707 DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
- 708 52. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079
709 (2009).
- 710 53. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
711 *Bioinformatics* **25**, 1754–1760 (2009).
- 712 54. Benaglia, T., Chauveau, D., Hunter, D. & Young, D. mixtools: An R package for analyzing finite
713 mixture models. *J. Stat. Softw.* **32**, 1–29 (2009).
- 714 55. R Core Team. R: A Language and Environment for Statistical Computing. [https://www.R-](https://www.R-project.org/)
715 [project.org/](https://www.R-project.org/) (2019).
- 716 56. Paten, B. *et al.* Cactus: Algorithms for genome multiple sequence alignment. *Genome Res.* **21**,
717 1512–1528 (2011).
- 718 57. Haeussler, M. *et al.* The UCSC Genome Browser database: 2019 update. *Nucleic Acids Res.* **47**,
719 D853–D858 (2019).
- 720 58. Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using MinHash.
721 *Genome Biol.* **17**, 132 (2016).
- 722 59. Smit, A. F. A. and Hubley, R. and Green, P. RepeatMasker Open-4.0. (2013-2015).
- 723 60. Hubisz, M. J., Pollard, K. S. & Siepel, A. PHAST and RPHAST: phylogenetic analysis with

- 724 space/time models. *Brief. Bioinform.* **12**, 41–51 (2011).
- 725 61. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful
726 Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series B (Methodological)*
727 vol. 57 289–300 (1995).
- 728 62. Nguyen, N. *et al.* Comparative assembly hubs: web-accessible browsers for comparative genomics.
729 *Bioinformatics* **30**, 3293–3301 (2014).
- 730 63. Benson, D. A. *et al.* GenBank. *Nucleic Acids Res.* **41**, D36–42 (2013).
- 731 64. Effect of fasting on carbohydrate metabolism in frugivorous bats (*Artibeus lituratus* and *Artibeus*
732 *jamaicensis*). *Comp. Biochem. Physiol. B Biochem. Mol. Biol.* **143**, 279–284 (2006).
- 733 65. Gordon, L. M. *et al.* Dental materials. Amorphous intergranular phases control the properties of
734 rodent tooth enamel. *Science* **347**, 746–750 (2015).
- 735 66. Hindle, A. G. & Martin, S. L. Intrinsic circannual regulation of brown adipose tissue form and
736 function in tune with hibernation. *Am. J. Physiol. Endocrinol. Metab.* **306**, E284–99 (2014).
- 737 67. Stanford, K. I. *et al.* Brown adipose tissue regulates glucose homeostasis and insulin sensitivity. *J.*
738 *Clin. Invest.* **123**, 215–223 (2013).
- 739 68. Chondronikola, M. *et al.* Brown adipose tissue improves whole-body glucose homeostasis and
740 insulin sensitivity in humans. *Diabetes* **63**, 4089–4099 (2014).
- 741 69. Saito, M. *et al.* High incidence of metabolically active brown adipose tissue in healthy adult humans:
742 effects of cold exposure and adiposity. *Diabetes* **58**, 1526–1531 (2009).
- 743 70. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456
744 humans. *Nature* **581**, 434–443 (2020).

745