

1 **Unbiased whole genomes from mammalian feces using fluorescence-activated cell sorting**

2

3

4 Joseph D. Orkin^{1,2*}, Marc de Manuel³, Roman Krawetz⁴, Javier del Campo⁵, Claudia Fontserè³,
5 Lukas F. K. Kuderna³, Ester Lizano³, Jia Tang⁶, Tomas Marques-Bonet³, Amanda D. Melin^{1,2}

6

7

8 1. Department of Anthropology & Archaeology, University of Calgary, Calgary, Alberta, Canada

9

10 2. Alberta Children's Hospital Research Institute, University of Calgary, Calgary, Alberta, Canada

11

12 3. Institut de Biologia Evolutiva (Consejo Superior de Investigaciones Científicas–Universitat
13 Pompeu Fabra), Barcelona Biomedical Research Park, Doctor Aiguader 88, Barcelona, Catalonia
14 08003, Spain.

15

16 4. Department of Cell Biology and Anatomy, University of Calgary, Calgary, Alberta, Canada

17

18 5. Department of Marine Biology and Oceanography, Institut de Ciències del Mar (Consejo
19 Superior de Investigaciones Científicas), Barcelona, Catalonia, Spain

20

21 6. Department of Biological Sciences, University of Calgary, Calgary, Alberta, Canada

22

23 ***Corresponding Author:**

24 Joseph D. Orkin (joseph.orkin@ucalgary.ca)

25

26

27
28

ABSTRACT

29 Non-invasive genomic research on free-ranging mammals typically relies on the use of
30 fecal DNA. This requires the isolation and enrichment of endogenous DNA, given its small
31 proportion compared to bacterial DNA. Current approaches for acquiring large-scale genomic
32 data from feces rely on bait-and-capture techniques. While this technique has greatly improved
33 our understanding of mammalian population genomics, it is limited by biases inherent to the
34 capture process, including allele dropout, low mapping rates, PCR duplication artifacts, and
35 structural biases. We report here a new method for generating whole mammalian genomes from
36 feces using fluorescence-activated cell sorting (FACS). Instead of enriching endogenous DNA
37 from extracted fecal DNA, we isolated mammalian cells directly from feces. We then built
38 fragment libraries with low input material from commercially available kits, which we
39 sequenced at high and low coverage. We validated this method on feces collected from primates
40 stored in RNAlater for up to three years. We sequenced one fecal genome at high coverage
41 (12X) and 15 additional fecal genomes at low coverage (0.1X - 4X). For comparative purposes,
42 we also sequenced DNA from nine blood or tissue samples opportunistically collected from
43 capuchin monkeys that died of natural causes or were treated in a local rehabilitation center.
44 Across all fecal samples, we achieved median mapping and duplication rates of 82% and 6%,
45 respectively. Our high-depth fecal genome did not differ in the distribution of coverage,
46 heterozygosity, or GC content from those derived from blood or tissue origin. As a practical
47 application of our new approach with low coverage fecal genomes, we were able to resolve the
48 population genetic structure of capuchin monkeys from four sites in Costa Rica.

49

INTRODUCTION

Advances in DNA sequencing technology have allowed for great strides to be made in comparative genomics (Arandjelovic & Vigilant, 2018; Corlett, 2017; Perry, 2014). It is now commonplace in a single study of a non-model organism to sequence partial genomes from multiple individuals. If fresh tissue or blood samples can be acquired from a handful of individuals, sequencing a *de novo* reference genome or generating a panel of single nucleotide variants is relatively straightforward. However, answering population level questions typically depends upon the non-invasive collection of fecal samples from free-ranging animals. Unfortunately, less than 5% of the extracted DNA from a fecal sample typically originates from an endogenous source (i.e. the host animal) (Hernandez-Rodriguez et al., 2017; Snyder-Mackler et al., 2016), while the remaining 95% comes from microorganisms and dietary items. For many species, this combination of factors makes it unfeasible to sequence whole genomes at high coverage from a large number of individuals. The resulting dearth of population-wide high-coverage genomes limits the scope of questions that can be asked and answered in genomics, ecology, and conservation.

Thanks to recent advances in non-invasive genomics, it has become possible to sequence partial genomes by enriching the proportion of endogenous DNA in feces (Chiou & Bergey, 2018; Perry, Marioni, Melsted, & Gilad, 2010; Snyder-Mackler et al., 2016). Through the use of targeted bait-and-capture and reduced representation libraries, this approach allows for the sequencing of single nucleotide variant (SNV) sets, which has begun to provide important insights into population structure and local adaptation of free-ranging mammals (Chiou, 2017; de Manuel et al., 2016; Wall et al., 2016). Despite the promise of this approach, DNA enrichment suffers from biases and impracticalities that limit its ability to uniformly cover a genome. Current bait-and-capture techniques are subject to inherent biases in the type of DNA captured (e.g. non-repetitive elements, GC content, reduced representation libraries, inconsistent hybridization); requires the costly and time consuming generation of RNA or DNA baits; have limited ability to enrich endogenous DNA (mean: ~57% of mapped reads (Snyder-Mackler et al., 2016)); and have high average PCR duplication rates (mean: ~38% of mapped reads (Snyder-Mackler et al., 2016)). Methylation-based enrichment offers a promising and cost-effective alternative to bait-and-capture for SNV generation, although it also suffers from inherent bias in the composition of the enriched libraries, and has limited enrichment capacity (mean: <50% of mapped reads (Chiou & Bergey, 2018)). While both approaches are viable for partial data, neither offers the realistic possibility of truly unbiased, cost-effective whole genome sequencing.

Through a novel application of fluorescence-activated cell sorting (FACS), we present a rapid, cost-effective method of isolating a host animal's intestinal epithelial cells for DNA extraction and genome sequencing. With this approach, we have routinely mapped more than 80% of reads to the host genome, strongly suggesting they are from endogenous DNA. This method requires no targeted enrichment of DNA, RNA baits, or methylation. It allows for DNA to be extracted and libraries built with commercially available kits, removing many of the challenges of enrichment-based techniques. Furthermore, our method allows for the long-term, room-temperature stabilization of samples, making it possible for field workers to collect samples with ease from remote areas without the need for temperature sensitive storage.

Here, we propose a novel protocol to isolate intestinal epithelial cells from the feces of white-faced capuchin monkeys (*Cebus imitator*) up to three years after initial collection. From these cells, we generated low coverage genomes from 17 fecal samples and selected one of them

96 for deeper sequencing (targeting ~10X - 15X coverage). In so doing, we have generated the first
97 uniformly-distributed, high-coverage, whole genome of a mammal from its feces. To
98 demonstrate the breadth of fecal FACS, we also conducted an analysis of population genetic
99 structure in two Costa Rican forest reserves using DNA derived from both fecal FACS and
100 traditional blood/tissue extractions.

101

102

METHODS

103 2.1 Sample Collection

104 We collected fecal samples from free-ranging white-faced capuchin monkeys (*Cebus*
105 *imitator*) at Sector Santa Rosa (SSR), part of the Área de Conservación Guanacaste in
106 northwestern Costa Rica, which is a 163,000 hectare tropical dry forest nature reserve (Figure 1).
107 Behavioral research of free-ranging white-faced capuchins has been ongoing at SSR since the
108 1980's which allows for the reliable identification of known individuals from facial features and
109 bodily scars (Fedigan & Rose-Wiles, 1996). We collected 14 fresh fecal samples from 12 white-
110 faced capuchin monkeys immediately following defecation (Table 1). We placed 1 mL of feces
111 into conical 15 mL tubes pre-filled with 5 mL of RNAlater. RNAlater preserved fecal samples
112 were sent to the University of Calgary, where they were stored at room temperature for up to
113 three years. To evaluate other preservation methods, we also collected two additional capuchin
114 monkey fecal samples (SSR-FL and a section of SSR-ML) and one spider monkey (*Ateles*
115 *geoffroyi*) fecal sample, which we stored in 1X PBS buffer and then froze in liquid nitrogen with
116 a betaine cryopreservative (Rinke et al., 2014). Given the logistical challenges of carrying liquid
117 nitrogen to remote field sites, we prioritized evaluation of samples stored in RNAlater.

118 Finally, we took tissue and blood samples opportunistically. During the course of our
119 study, 4 individual capuchin monkeys died of natural causes at SSR, from whom we were able to
120 collect tissue samples, which we stored in RNAlater. By collaborating with *Kids Saving the*
121 *Rainforest* veterinary rehabilitation clinic in Quepos, Costa Rica, we acquired blood samples
122 from 5 more Costa Rican white-faced capuchins who were undergoing treatment at the facility
123 (although we were unable to collect paired fecal samples). Samples were collected with
124 permission from the Area de Conservacion Guanacaste (ACG-PI-033-2016) and CONAGEBIO
125 (R-025-2014-OT-CONAGEBIO). Samples were exported from Costa Rica under permits from
126 CITES and Area de Conservacion Guanacaste (2016-CR2392/SJ #S 2477, 2016-CR2393/SJ #S
127 2477, DGVS-030-2016-ACG-PI-002-2016; 012706) and imported with permission from the
128 Canadian Food and Inspection agency (A-2016-03992-4).

129

130 2.2 FACS

131 Before isolating cells by Fluorescence-activated cell sorting (FACS), fecal samples were
132 prepared using a series of washes and filtration steps. Fecal samples were vortexed for 30 s and
133 centrifuged for 30 s at 2,500 g. Then the supernatant was passed through a 70 um filter into a 50
134 mL tube and washed with DPBS. After transferring the resultant filtrate to a 15 mL tube, it was
135 centrifuged at 1,500 RPM for 5 minutes to pellet the cells. Then we twice washed the cells with
136 13 mL of DPBS. We added 500 uL of DPBS to the pellet and re-filtered through a 35 um filter
137 into a 5 mL FACS tube. We prepared a negative control (to control for auto-fluorescence) with
138 500 uL of DPBS and one drop of the cell solution. To the remaining solution, we added 1 uL of
139 AE1/AE3 Anti-Pan Cytokeratin Alexa Fluor® 488 antibody or TOTO-3 DNA stain, which we
140 allowed to incubate at 4°C for at least 30 minutes.

141 We isolated cells using a BD FACSAria™ Fusion (BD Biosciences) flow cytometer at the
142 University of Calgary Flow Cytometry Core. To sterilize the cytometer's fluidics before
143 processing each sample, we ran a 3% bleach solution through the system for four minutes at
144 maximum pressure. We assessed background fluorescence and cellular integrity, by processing
145 the negative control sample prior to all prepared fecal samples. For each sample we first gated
146 our target population by forward and side scatter characteristics that were likely to minimize
147 bacteria and cellular debris (Figure 2). Secondary and tertiary gates were implemented to remove
148 cellular agglomerations. Finally, we selected cells with antibody or DNA fluorescence greater
149 than background levels. In cases when staining was not effective, we sorted solely on the first
150 three gates. Cells were pelleted and frozen at -20°C.

151

152 2.3 DNA Extraction and Shotgun Sequencing

153 We extracted fecal DNA (fDNA) with the QIAGEN DNA Micro kit, following the
154 "Small volumes of blood" protocol. To improve DNA yield, we increased the lysis time to three
155 hours, and incubated 50 mL of 56°C elution buffer on the spin column membrane for 10
156 minutes. DNA concentration was measured with a Qubit fluorometer. Additionally, to calculate
157 endogenous DNA enrichment, we extracted DNA directly from five fecal samples prior to their
158 having undergone FACS. We extracted DNA from the nine tissue and blood samples using the
159 QIAGEN Genra Puregene Tissue kit and DNeasy blood and tissue kit, respectively.

160 For the fecal samples, DNA was fragmented to 350 bp with a Covaris sonicator. We built
161 whole genomic sequencing libraries with the NEB Next Ultra 2 kit using 10-11 PCR cycles.
162 Fecal genomic libraries were sequenced on an Illumina NextSeq (2x150 PE) at the University of
163 Calgary genome sequencing core. We resequenced one fecal sample at high coverage on an
164 Illumina HighSeq 4000 at the McDonnell Genome Institute at Washington University in St.
165 Louis (MGI). High-coverage, whole genomic shotgun libraries were prepared for the blood and
166 tissue DNA samples and sequenced on an Illumina X-10 at MGI.

167

168 2.3 Mapping and SNV Generation

169 Reads were trimmed of sequencing adaptors with Trimmomatic (Bolger, Lohse, &
170 Usadel, 2014). Subsequently, we mapped the *Cebus* reads to the *Cebus imitator* 1.0 reference
171 genome ([GCF_001604975.1](https://doi.org/10.1093/gbe/abz011)) with BWA mem (Li & Durbin, 2009) and removed duplicates with
172 Picard Tools (<http://broadinstitute.github.io/picard/>). We called SNVs for each sample
173 independently using the *Cebus* genome and the GATK UnifiedGenotyper pipeline (-out_mode
174 EMIT_ALL_SITES) (McKenna et al., 2010). Genomic VCFs were then combined using GATK's
175 CombineVariants restricting to positions with a depth of coverage between 3 and 100, mapping
176 quality above 30, no reads with mapping quality zero and variant PHRED scores above 30.
177 Sequencing reads from one of the high coverage fecal samples (SSR-FL) bore a strong signature
178 of human contamination (16%), and were thus excluded from SNV generation. We included
179 reads from nine tissue/blood samples and one frozen fecal sample with high coverage (SSR-ML).
180 In total, we generated 4,184,363 SNVs for downstream analyses.

181 To remove potential human contamination from sequenced libraries, we mapped trimmed
182 reads to the *Cebus imitator* 1.0 and human (hg38) genomes simultaneously with BBsplit
183 (Bushnell, 2016). Using default BBsplit parameters, we binned separately reads that mapped
184 unambiguously to either genome. Ambiguously mapping reads (i.e. those mapping equally well
185 to both genomes) were assigned to both genomic bins, and unmapped reads were assigned to a
186 third bin. We calculated the amount of human genomic contamination as the percentage of total

187 reads unambiguously mapping to the human genome (Table 2). After removing contaminant
188 reads, all libraries with at least 0.5X genomic coverage were used for population analysis.

189 In order to test the effect of fecal FACS on mapping rates, we selected five samples at
190 random (SSR-CH, SSR-NM, SSR-LE, SSR-PR, SSR-SN) to compare pre- and post-FACS
191 mapping rates. To test for an increase in mapping percentage, we ran a one-sample paired
192 Wilcoxon signed-rank test on the percentages of reads that mapped exclusively to the *Cebus*
193 genome before and after flow FACS. Additionally, we ran Pearson's product moment
194 correlations to test for an effect of the number of cells (log10 transformed) on rates of mapping,
195 read duplication, and ng of input DNA. The above tests were all performed in R.

196

197 2.5 High coverage fecal genome comparison

198 We made several comparisons between our high-coverage feces-derived genome and the
199 blood/tissue-derived genomes using window-based approaches. For each test, the feces-derived
200 genome should fall within the range of variation for members of its population of origin (SSR).
201 Deviations from this, for examples all fecal genomes clustering together, would indicate biases
202 in our DNA isolation methods. To assess this, we constructed 10 KB / 4KB sliding windows
203 along the largest scaffold (21,314,911 bp) in the *C. imitator* reference genome. From these
204 windows, we constructed plots of coverage density and the distribution of window coverage
205 along the scaffold. Secondly, we assessed the level of heterozygosity in 1 MB / 200 KB sliding
206 windows throughout the ten largest scaffolds. For each high-coverage genome, we plotted the
207 density distribution of window heterozygosity. We measured genome-wide GC content with the
208 Picard Tools CollectGcBiasMetrics function. The percentage of GC content was assessed against
209 the distribution of normalized coverage and the number of reads in 100 bp windows per the
210 number reads aligned to the windows.

211

212 2.6 Population genomic analysis

213 Given the large degree of difference in coverage among our samples, (less than 1X to
214 greater than 50X), performed pseudodiploid allele calling on all samples using custom scripts.
215 For each library, at each position in the SNV set, we selected a single, random read from the
216 sequenced library. From that read, we called the variant information at the respective SNV site
217 for the given library. In so doing, we generated a VCF with a representative degree of variation
218 and error for all samples.

219 To assess population structure and infer splits between northern and southern groups of
220 Costa Rican white-faced capuchins, we constructed principal components plots with
221 EIGENSTRAT (Price et al., 2006) and built population trees with TreeMix (Pickrell & Pritchard,
222 2012). Because we ascertained variants predominantly with libraries that were of tissue/blood
223 origin, we built principal components solely with SNVs from these libraries and projected the
224 remaining fecal libraries onto the principal components. For our maximum likelihood trees, we
225 used three outgroups (*Ateles geoffroyi*, *Saimiri sciureus*, and *Cebus albifrons*), with *A. geoffroyi*
226 serving as the root of the tree. Given the geographic distance and anthropogenic deforestation
227 between northern and southern populations, we assumed no migration. To account for linkage
228 disequilibrium, we grouped SNVs into windows of 1,000 SNVs.

229

230

230 RESULTS

231

232 3.1 Isolation of intestinal epithelial cells using Fluorescence-activated cell sorting (FACS)

233 Flow cytometry can be used to discriminate among categories of cells by examining the
234 manner in which light scatters in response to cellular properties. We interpreted forward scatter
235 (FSC) and side scatter (SSC) as measures of cellular size and granularity (complexity),
236 respectively. When cells are intact, free of agglomerations, and of limited variety, they form
237 easily identifiable clusters, particularly when bound with fluorescently labeled antibodies. In
238 contrast to this idealized schema, abundant cellular debris prevented us from observing distinct
239 cellular populations when assessing the relationship between FSC and SSC (Figure 2) of fecal
240 samples. The vast majority of events were usually clustered in lower range of FSC, and likely of
241 bacterial origin. To exclude bacteria insofar as possible, we implemented a FSC gate that only
242 included events above this cluster, typically the top $\frac{1}{2}$ to $\frac{2}{3}$ of the FSC range. From the 14
243 RNAlater preserved capuchin fecal samples, we isolated a median of 1,739 cells, with a range of
244 129 - 62,201 (Table 2). Typically, we collected a few hundred or thousand cells, but in two cases
245 of poor fluorescent staining (SSR-FN and the RNAlater preserved SSR-ML sample), we sorted
246 the larger gated populations, irrespective of fluorescent intensity. From the frozen samples, SSR-
247 FL and SSR-ML, we collected 4,405 and 2,546 cells, respectively. Similarly, from the spider
248 monkey sample, which we split into two separate FACS runs, we isolated 4,026 and 602 cells.
249

250 3.2 Mapping of genomic libraries

251 From each cellular population, we successfully extracted DNA and prepared sequencing
252 libraries. Among the RNAlater preserved capuchin samples, the total amount of DNA per sample
253 was low, ranging from 2.96 to 21.50 ng, with a median value of 7.85 ng (Table 2). A relationship
254 between the number of cells was not significantly correlated with the amount of extracted DNA
255 ($R=0.227$; 95% CI (-0.345, 0.676); $t=0.808$, $p > 0.05$) or mapping rate ($R = -0.204$; 95% CI (-
256 0.663, 0.367); $t = -0.721$; $p > 0.05$). Median mapping rates reached 93% (range: 55 - 98%) with
257 BWA-MEM and 82% (range: 11 - 95%) with the more stringent BBsplit settings (Figures 3A,
258 3B, Table 2). Read duplication levels were low, with a median value of 9% (range: 2 - 40%)
259 resulting in 63% (range: 8 - 92%) of reads being unique and mapping to the *Cebus imitator* 1.0
260 genome. The amount of duplicate reads was distributed bimodally across individuals, with reads
261 from five samples having substantially higher duplication rates than the remaining nine. The rate
262 of duplication was significantly correlated ($R = -0.751$; 95% CI (-0.917, -0.366); $t = -3.94$; $p <$
263 0.01) with the number of cells (log₁₀ transformed), decreasing sharply above a threshold of
264 about 1,000 cells (Figure 3D).

265 The samples frozen in liquid nitrogen mapped at comparable rates to those preserved in
266 RNAlater. From the two frozen capuchin samples, SSR-ML and SSR-FL, respectively, we
267 extracted 10.50 and 6.72 ng of DNA. These two samples mapped at 96% and 80.4% with BWA-
268 MEM and 90% and 42% with BBsplit (5% and 3% duplicates), respectively. We extracted 6.96
269 and 4.50 ng of DNA from the two runs of the spider monkey sample, which mapped at a
270 substantially lower rate of 54% and 49% with BWA-MEM and 12% with BBsplit for both (1%
271 duplicates for both).

272 We observed little to no human contamination in the RNAlater preserved samples. For
273 nine of the 14 samples, BBsplit mapped between 0.61 and 1.25% of reads to hg38 (median
274 0.96%); however, in four cases 2.86 - 5.80% of reads were binned to the human genome. Human
275 mapped reads were also low for the frozen SSR-ML (1.25%) and spider monkey (2.83% and
276 1.82%) samples. However, SSR-FL appeared to have substantial human contamination (15.77%
277 of reads). This may be due to initial processing of these three samples, which were stored using
278 the cryopreservation method, at the field site. We conducted the initial vortexing, centrifugation,

279 and collection of supernatant (see section 2.2) at the SSR field station, which is likely where
280 SSR-FL was contaminated. Due to this, we examine the mapping rates using only the RNAlater
281 preserved samples. However, we were able to decontaminate reads bioinformatically, and
282 include the decontaminated reads in downstream analyses where appropriate.

283 By sorting fecal samples with FACS, we substantially increased the percentage of reads
284 mapping to the target genome. We selected five samples at random (SSR-CH, SSR-NM, SSR-
285 LE, SSR-PR, SSR-SN) to compare pre- and post-FACS mapping rates. The mapping rates of
286 unsorted feces ranged from 10 - 42%, with a median of 14% (Figure 3C). After flow sorting
287 aliquots of these fecal samples, we obtained significantly higher mapping rates ($V = 15$, $p <$
288 0.05) for each sample, ranging from 64 - 95%, with a median of 85%, resulting in a median 6.07
289 fold enrichment.

290

291 3.3 High coverage fecal genome

292 Given that the sample SSR-ML had a high mapping percentage, a low rate of duplication,
293 and was effectively free of human-specific mapping, we selected it for sequencing at high
294 coverage. Using $\frac{1}{2}$ of one HiSeq 4000 lane, we achieved an average coverage of $\sim 12X$ across the
295 *Cebus imitator* 1.0 genome.

296 When comparing the high coverage fecal and tissue genomes from the Santa Rosa site,
297 we observed no substantial difference in quality, coverage, heterozygosity, or GC content
298 (Figures 3 and 4). For each genome, the distribution of per site coverage followed a roughly
299 normal distribution with a small number of positions uncovered ($\sim 2\%$) (Figure 3A). Coverage
300 along the largest scaffold from the *Cebus* genome was uniform in both tissue and fecal samples
301 (Figure 3B). No obvious area of excessively high or low coverages is apparent in the fecal
302 genome compared to that of the tissue derived genomes. Importantly, the fecal genome does not
303 have any obvious gaps in coverage. Likewise, levels of heterozygosity were comparable between
304 fecal and tissue genomes (Figure 3C, D). The fluctuating levels of heterozygosity across the
305 largest genomic scaffold in 100 KB windows is highly similar for SSR-ML and SSR-CR (Figure
306 3D), indicative of their close familial relationship. Finally, the distribution of GC content across
307 the genome does not suffer from substantial bias (Figure 5B). Although the normalized coverage
308 at the extremes of the GC distribution is on the higher end of the capuchin samples (Figure 5A),
309 it falls well within the range of other samples for the vast majority of the genome where GC
310 content ranges from $\sim 20 - 75\%$ (Figure 5B).

311

312 3.4 Population structure

313 We observed likely population subdivision between the northern and southern groups of
314 white-faced capuchins in our SNV set. This separation corresponds to the ecological division of
315 the season tropical dry forests in the north from the non-seasonal tropical wet forests in the
316 south. Given the limitations of the available sampling sites, it is possible that the appearance of
317 an ecological divide is actually evidence of isolation by distance.

318 All individuals from the north and the south are sharply discriminated by the first
319 principle component of the PCA (Figure 6A). The second component indicates a higher degree
320 of genetic variation within the southern individuals. All the northern individuals form a tight
321 cluster on the PCA plot, in contrast to those from the south, which are more widely dispersed
322 along PC 2. Furthermore, the single individual from the northern site of Cañas clusters closely
323 with the individuals from Santa Rosa, despite a geographic distance of more than 100 km, which
324 suggests that isolation by distance might not be the sole reason for population differentiation. No

325 clustering was observed within the four individuals from the southern sites of Manuel Antonio
326 and Quepos, apart from their separation from the northern individuals along PC 1. Because we
327 generated the principal components with samples from the primary SNV set and projected the
328 remaining samples (fecal flow FACS and tissue-based outgroups), the outgroup taxa are
329 expected to fall in between the two main sampling clusters of white-faced capuchins. As
330 expected, the three outgroup taxa (*C. albifrons*, *S. sciureus*, and *A. geoffroyi*) fall in the center of
331 the PCA plot.

332 The pattern of clustering generated by our maximum likelihood SNV tree recapitulates
333 the expected patterns of geographic distance and ecological separation in our sample (Figure 6).
334 Among the white-faced capuchin monkeys, the northern and southern clades represent the main
335 split in the tree. Each clade is subdivided according to the two sampling sites within the
336 geographic/ecological regions. Furthermore, the three outgroup taxa split by the expected degree
337 of evolutionary distance. These relationships are not perturbed by the fact that samples were a
338 mixture of traditional tissue-based genomic libraries and libraries generated by fecal flow-FACS.
339 This pattern is evident both within the northern sites and outgroup taxa. Additionally, depth of
340 coverage does not appear to affect the pattern of clustering. Our sample ranged in coverage from
341 less than 1X to greater than 50X. In spite of this, the pattern of geographic/ecological subdivision
342 held.

343 DISCUSSION

344
345
346 In this manuscript, we describe a novel use of FACS to isolate cells from the feces of
347 free-ranging mammals for population and comparative genomics. We have demonstrated that
348 fecal FACS is an effective means for: 1) the enrichment of endogenous DNA from non-invasive
349 primate samples; 2) the generation of unbiased whole genomes at high coverage or low coverage
350 sequencing libraries suitable for population genomic analysis. Isolating genome-scale
351 information from non-invasively collected samples remains a major challenge in molecular
352 ecology. Although DNA can be extracted readily from museum specimens and captive
353 individuals (Guschanski et al., 2013; Prado-Martinez et al., 2013; van der Valk, Lona Durazo,
354 Dalén, & Guschanski, 2017), the vast majority of the world's mammalian genomic diversity
355 remains in free-ranging individuals. Our results indicate that fecal FACS has the potential for
356 widespread application in molecular ecology and the broadening of non-invasive genomics for
357 threatened and cryptic mammals.

358 4.1 Performance and cost-effectiveness

359
360 Current techniques to isolate whole genomic information from fecal samples depend
361 upon the enrichment of endogenous DNA from extracted fDNA (Chiou & Bergey, 2018; Perry et
362 al., 2010; Snyder-Mackler et al., 2016). While these methods have proven effective for SNV
363 analyses, particularly at low coverage (Chiou, 2017; de Manuel et al., 2016; Wall et al., 2016),
364 they remain of limited genomic scope. The total mapping rate of endogenous reads from the
365 highest performing enrichment protocol is 57%, with a non-duplicate mapping rate of 38%
366 (Snyder-Mackler et al., 2016). The median non-duplicate rate that we generated through FACS is
367 63% (82% when including duplicates), substantially outperforming that of enrichment-based
368 approaches. While sequencing costs have fallen dramatically in recent years, maximizing the
369 proportion of non-duplicate reads in sequencing libraries remains a critical factor in determining
370 the feasibility of sampling schemes. Studies that aim to sequence tens or hundreds of fecal

371 individuals at high coverage are simply not practical for most labs, given the current cost
372 structure. We were able to isolate primate cells from feces for roughly \$40 per sample. Given
373 that each sample required about 30 minutes of FACS time and three hours of wet lab preparation
374 time (per batch of samples), a trained lab worker could prepare five to ten samples per day,
375 presuming the availability of FACS resources. Although these costs of time and money are not
376 negligible, this may be a justifiable expense for projects where the increased mapping rate and
377 genomic coverage are desired.

378 While our fecal FACS method is effective in white-faced capuchin monkeys and
379 Geoffroy's spider monkeys, we acknowledge that further validation in other species is warranted.
380 Given the disparity in mapping rates between the capuchin and spider monkey samples, it is
381 possible that cytometry protocols would need to be optimized toward the particularities of a
382 given species' feces and conditions. Consistent with this notion is the fact that the fecal sample
383 (SSR-SB1) with low mapping success, was substantially darker than the other capuchin samples,
384 which, depending on the dietary items consumed, typically have a green, brown, or rust
385 coloration. Mapping was substantially improved in the replicate sample (SSR-SB2), which was
386 collected on a different day. Curiously, we did not observe a relationship between the number of
387 sorted cells and the concentration of extracted DNA. However, this is likely explained by
388 residual intercalating dyes used in FACS process remaining in the sorted cells and interfering
389 with Qubit quantification (Kuderna et al., 2018). Additionally, it is peculiar that the mapping
390 rates of the libraries we built from unsorted fDNA were so high (median 14%). Typically, less
391 than 5% of fDNA is of an endogenous source, although some chimpanzee samples have been
392 reported to have up to 25% endogenous reads (Hernandez-Rodriguez et al., 2017). Further
393 testing of capuchin fecal samples with lower endogenous DNA concentration is worth pursuing,
394 as mapping rates for endogenous DNA from unprocessed and enriched fDNA are often
395 correlated (Chiou & Bergey, 2018; Hernandez-Rodriguez et al., 2017; Snyder-Mackler et al.,
396 2016). However, because cell sorting is not a targeted DNA enrichment process, we find it
397 unlikely that post-FACS mapping rates should depend on the concentration of endogenous
398 fDNA; accordingly, we did not observe any such relationship among the five samples we tested
399 for enrichment (Figure 3C). Furthermore, we did not observe a correlation between the number
400 of isolated cells and the mapping rate; in one case, we obtained a 94% mapping rate with only
401 140 cells. Presuming that the flow cytometer is sorting cells correctly, and that those cells
402 contain viable DNA, the mapping rate should only be contingent upon the accuracy of the cell
403 sorting process.

404 We have demonstrated that RNAlater is an effective, long-term, room-temperature
405 cellular storage medium for fecal FACS. In the great majority of cases, FACS involves the
406 sorting of living cellular populations, and attempts to sort dead cells are often met with
407 skepticism (Sasaki, Dumas, & Engleman, 1987). Dead cells are typically distorted and
408 fragmented, yielding populations that are difficult to discriminate. We attempted to freeze fresh
409 feces with liquid nitrogen and a betaine cryopreservative, following the single-cell protocol of
410 Rinke et al. (2014). Unfortunately, many of these samples contained extremely large amounts of
411 cellular debris, likely from improper cryopreservation in field conditions. Additionally,
412 cryopreservation of samples required a non-trivial amount of laboratory preparation in non-
413 sterile field conditions that we believe introduced substantial human contamination to SSR-FL.
414 While we were able to sequence one of the frozen samples (SSR-ML) at high coverage and
415 replicate it with RNAlater, we cannot presently recommend in-field cryopreservation of fecal
416 samples for FACS. RNAlater is often commonly used in molecular field primatology, because it

417 offers long-term, stable preservation of host DNA at room temperature. For our purposes, it also
418 offered the distinct advantage of not requiring any in-field laboratory preparation, which
419 minimized human contamination of our cellular populations. Attempts to flow sort RNAlater
420 preserved cells of any origin are extremely scant, and we only found two such studies in the
421 literature (Barrett et al., 2002; Zaitoun, Erickson, Schell, & Epstein, 2010). We observed a
422 substantial improvement in the cellular integrity in the RNAlater preserved samples. Although
423 cells preserved in RNAlater are dead, they suffer minimal histological disruption, and maintain
424 cellular epitopes critical for antibody binding (Florell et al., 2001). Given this array of benefits,
425 we recommend preservation of fresh fecal samples in RNAlater when collected in field
426 conditions.

427

428 4.2 Quality and feasibility of high-coverage fecal genomics

429 We have presented the first high-coverage, unbiased mammalian genome, derived
430 exclusively from feces. While traditional bait-and-capture approaches to non-invasive genomics
431 have allowed for broad sampling of the mammalian genome from feces, such methods remain
432 limited by genomic bias. When compared to tissue-derived capuchin genomes, our FACS-
433 derived fecal genome indicates no such biases. SSR-ML consistently fell within or immediately
434 adjacent to the observable range of variation of the other tissue samples collected from Sector
435 Santa Rosa. While we acknowledge that it would have been optimal to compare high-coverage
436 whole genomes generated from the blood and feces of the same individual, this was not possible,
437 because of our non-invasive sampling strategy. Nonetheless, we are able to infer such a
438 comparison through the use of pedigree data in our SSR samples. The scaffold-wide pattern of
439 heterozygosity observable in SSR-ML (Figure 4D) is nearly identical to that of SSR-CR (tissue),
440 who was his sibling. This relationship is further supported by the population clustering results
441 (section 4.3). Furthermore, the SSR-ML sample we used in SNV calling did not bear any
442 indication of human contamination. In order to remain consistent in comparison with the tissue
443 and blood derived samples, we did not remove reads mapping to the hg38 with BBsplit during
444 SNV calling. Because the *Cebus* genome is less complete than hg38, it is likely that the majority
445 of human-specific mapping from this and other samples is artifactual. Given the consistency
446 similarity of the SSR-ML sample to the others from SSR, we and suggest that FACS is a viable
447 approach to expand the horizons of non-invasive population and conservation genomics.

448 Prior to selecting libraries for high-coverage sequencing, we suggest that multiple
449 libraries should be run on a lower throughput sequencing platform (e.g. MiSeq). Given the
450 variability in sequencing outcomes inherent in our technique, it would be prudent to avoid
451 wasting sequencing capacity on libraries that lack the requisite diversity for high-depth
452 sequencing. Working with extremely low numbers of cells, which is sometimes the result of the
453 FACS process, can result non-trivial duplication rates and the potential for the introduction of
454 human contaminants. Given that our FACS protocol only requires a small amount of fecal slurry,
455 processing two or three aliquots from the same fecal sample would increase the number of cells
456 and, presumably, the available diversity in cases where it was deemed necessary.

457

458 4.3 Population structure of white-faced capuchin monkeys

459 By successfully discriminating among two populations of white-faced capuchins in Costa
460 Rica, we have demonstrated that fecal FACS is effective for low-coverage applications of
461 population and conservation genomics. While bait-and-capture approaches remain a valuable

462 tool for the assessment of population genetic structure from real-world distributions of free-
463 ranging mammals, fecal FACS provides a simple alternative approach.

464 The clustering patterns in our trees and PCA plots do not reveal any samples that deviate
465 from their expected geographic or ecological origin. These relationships are robust to both the
466 coverage levels (< 1X to > 50X) and biological origins (feces, tissue, and blood) of the samples.
467 The tight geographic clustering of individuals within the SSR sampling locale provides
468 reasonable evidence that there is no substantial effect from fecal FACS on population structure.
469 Were it the case that fecal FACS introduced substantial bias, we would have expected the fecal
470 samples from SSR to plot in a separate cluster from those of tissue origin. As fecal and tissue
471 samples fall in the same general cluster, this is no evidence of such an effect. Furthermore,
472 known pedigree information from SSR corresponds to the genetic relationships observed in our
473 SNV tree. SSR-ML (fecal) and SSR-CR (tissue) form an internal clade in the tree (sixth and
474 seventh points from the top). These two individuals also cluster adjacent to each other on the
475 PCA plot.

476

477 4.4 Summary

478 Through a novel use of flow cytometry/FACS, we have developed a new method for the
479 isolation of epithelial cells from mammalian feces for population genomics. We generated the
480 first high-coverage, unbiased mammalian genome solely from feces. Additionally, we have
481 demonstrated that fecal FACS can be used to generate low coverage SNP datasets that function
482 well in population assignment and clustering algorithms. Fecal FACS is cost-effective and free
483 of the biases that commonly occur in traditional bait-and-capture approaches to the enrichment
484 of endogenous DNA from feces. Furthermore, fecal FACS does not require costly impractical
485 preservation of biomaterial in liquid nitrogen; rather, we rely on room-temperature stable storage
486 in RNAlater. Fecal FACS offers great benefits to the field of mammalian conservation and
487 population genomics.

488

489 **AUTHOR CONTRIBUTIONS**

490

491 Research was designed by JDO, ADM, RK and JC; performed by JDO, RK, CF, LFKK, EL, and
492 JT; analyzed by JDO and MM; and written by JDO, ADM, and TMB.

493

494 **DATA ACCESSIBILITY STATEMENT**

495

496 Sequencing data will be archived on NCBI SRA and made publically available.

497

498 **ACKNOWLEDGEMENTS**

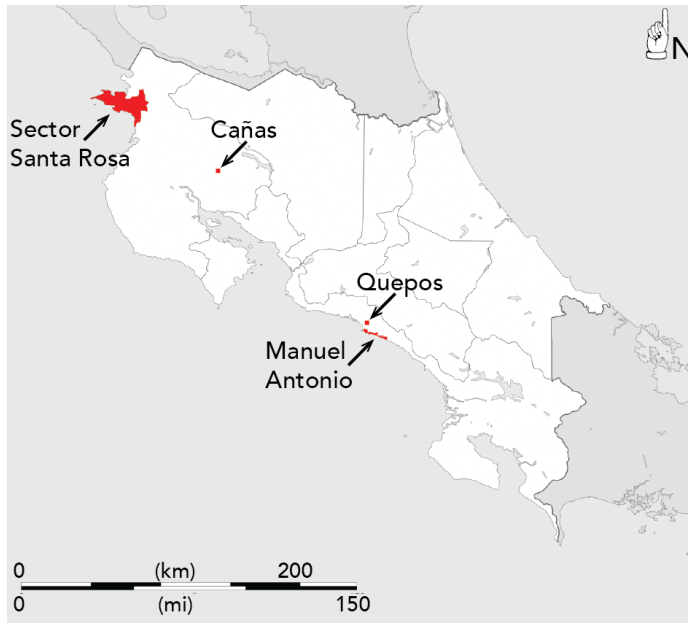
499

500 We would like to thank Laurie Kennedy, Yiping Liu from the University of Calgary Flow
501 Cytometry Core for their patience and assistance with developing this protocol. Additionally. We
502 thank Jene Weatherhead, Shelley Wegener, Frank Visser, Gwen Duytschaever for molecular
503 laboratory assistance. We acknowledge Oscar Fornas for helpful discussion. Thanks to Wes
504 Warren, Pat Minx, Mike Montague, Shoji Kawamura, and J. Pedro Magalhaes for their
505 involvement with the development of the *Cebus imitator* reference genome. PJ Perry, Shasta
506 Webb, Rachel Williamson, and Saul Cheves Hernandez assisted with sample acquisition. This
507 research was supported by the National Sciences and Engineering Research Council of Canada

508 (NSERC), and the Canada Research Chairs program to ADM. TMB is supported by BFU2017-
509 86471-P (MINECO/FEDER, UE), Howard Hughes International Early Career, Obra Social "La
510 Caixa" and Secretaria d'Universitats i Recerca del Departament d'Economia i Coneixement de la
511 Generalitat de Catalunya. JDO is supported by the Alberta Children's Hospital Research Institute
512 (ACHRI). CF is supported by "La Caixa" PhD fellowship. L.F.K.K. is supported by an FPI
513 fellowship associated with BFU2014-55090-P (MINECO/FEDER, UE)
514
515

516
517

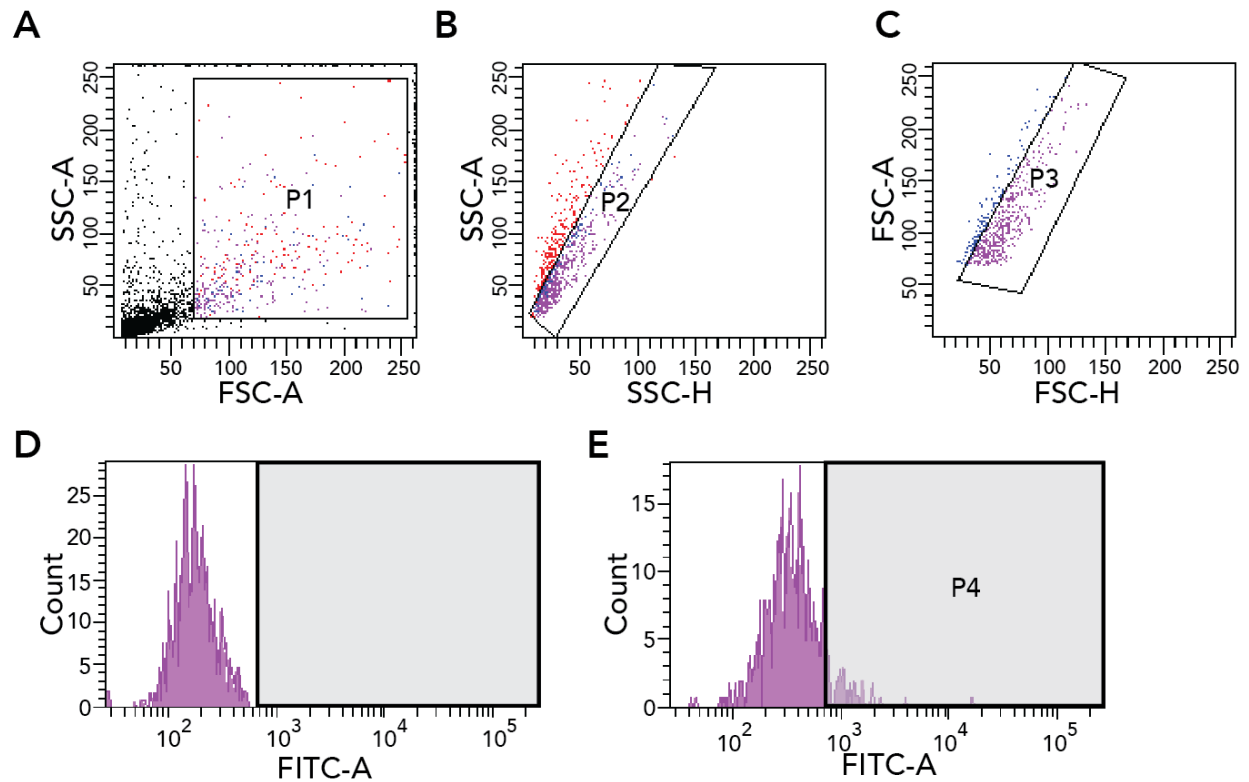
FIGURES



518
519

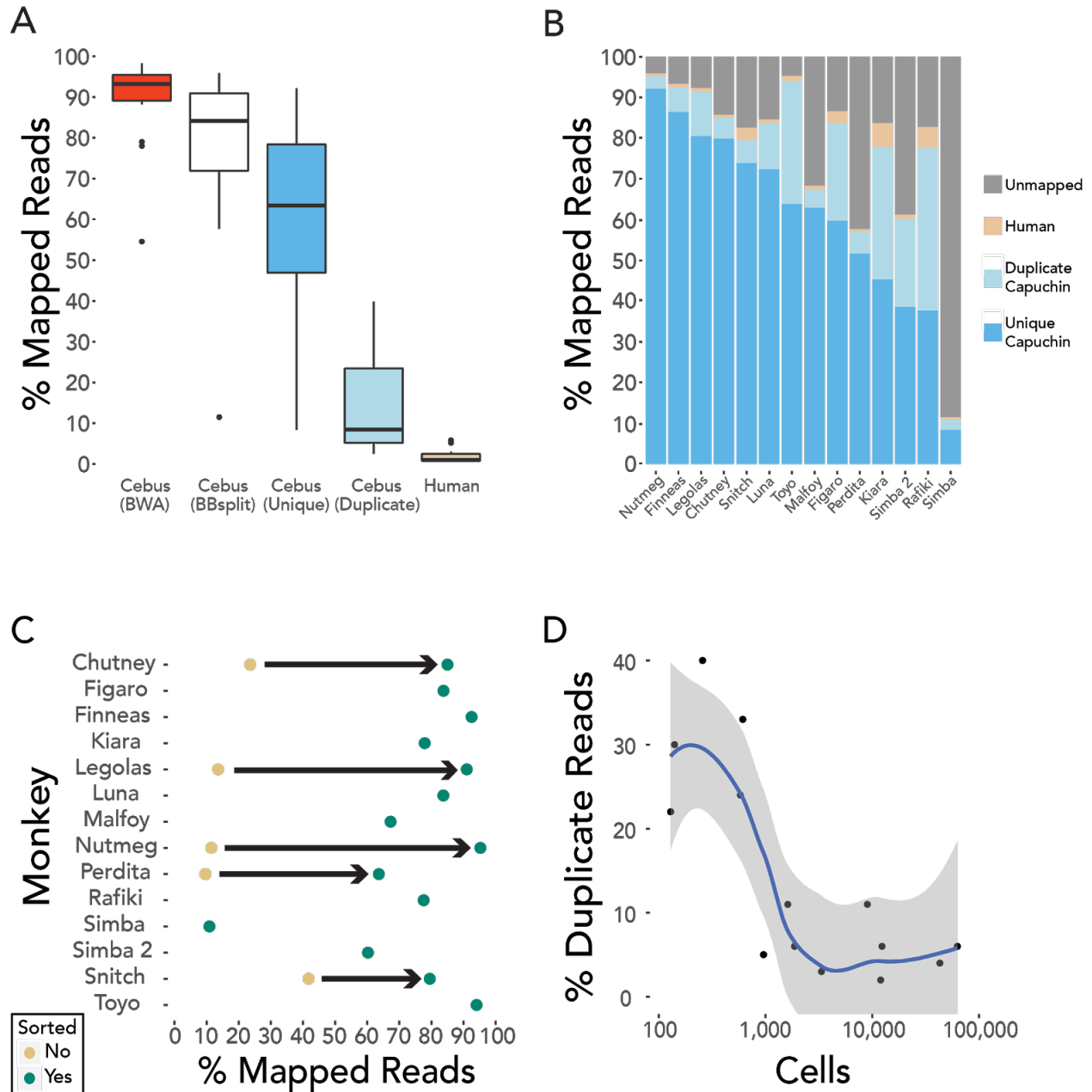
520 **Figure 1:** Map of sampling sites. Sector Santa Rosa (SSR) and Cañas are situated in the northern
521 dry forest and samples from Quepos and Manuel Antonio are from the southern wet forest. Map
522 courtesy of Eric Gaba—Wikimedia Commons user: Sting.

523
524



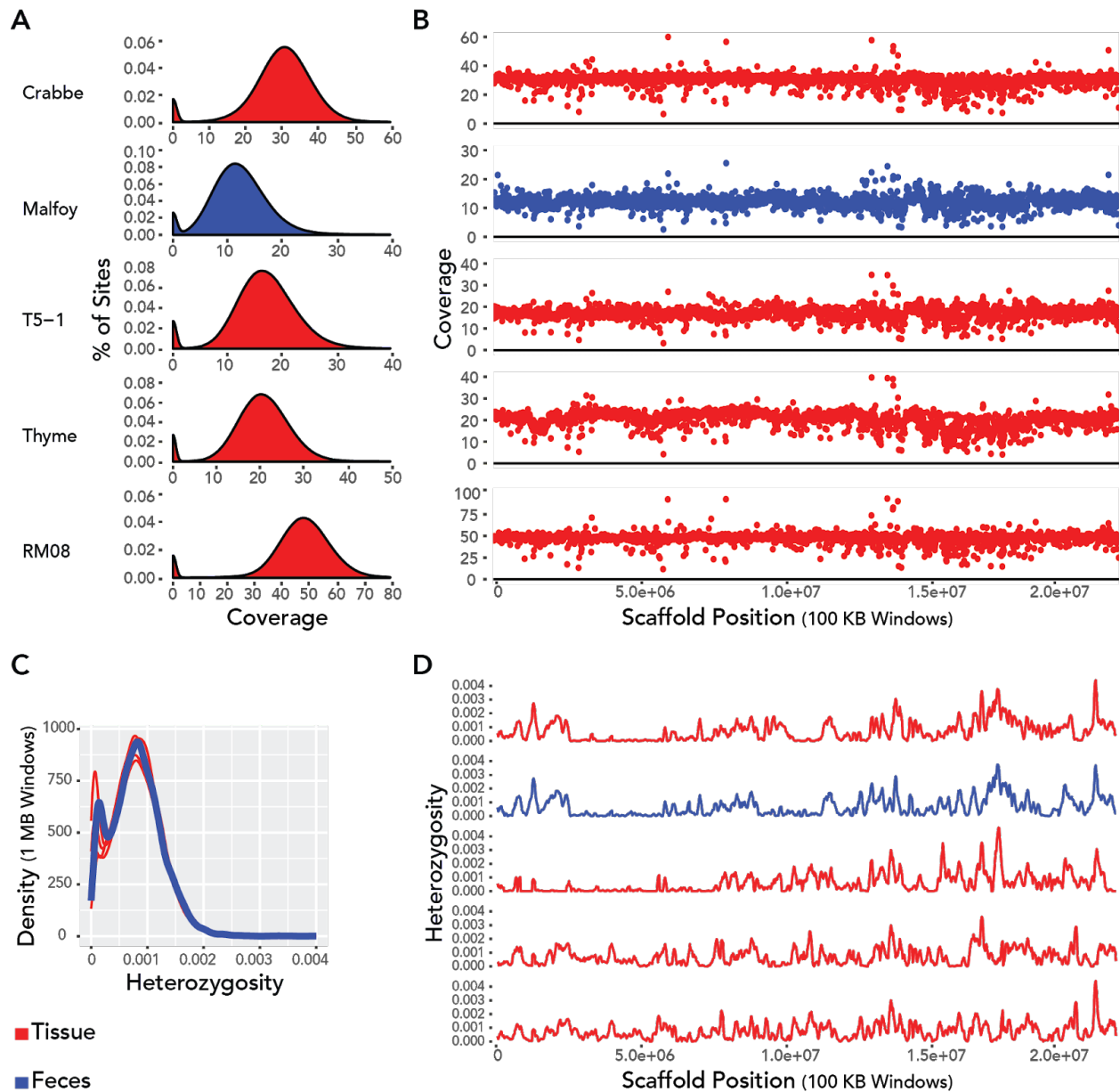
525
526

527 **Figure 2:** FACS gating strategy. Cells were gated first by size and complexity to avoid bacteria
528 and cellular debris (A), followed by discrimination of cellular agglomerations (B and C).
529 Fluorescence of AE1/AE3 Anti-Pan Cytokeratin Alexa Fluor® 488 antibody (FITC-A) is
530 depicted in unstained (D) and stained (E) cellular populations. Epithelial cells were identified as
531 those fluorescing beyond background levels, as depicted in the P4 gate.
532



533
534
535
536
537
538
539
540

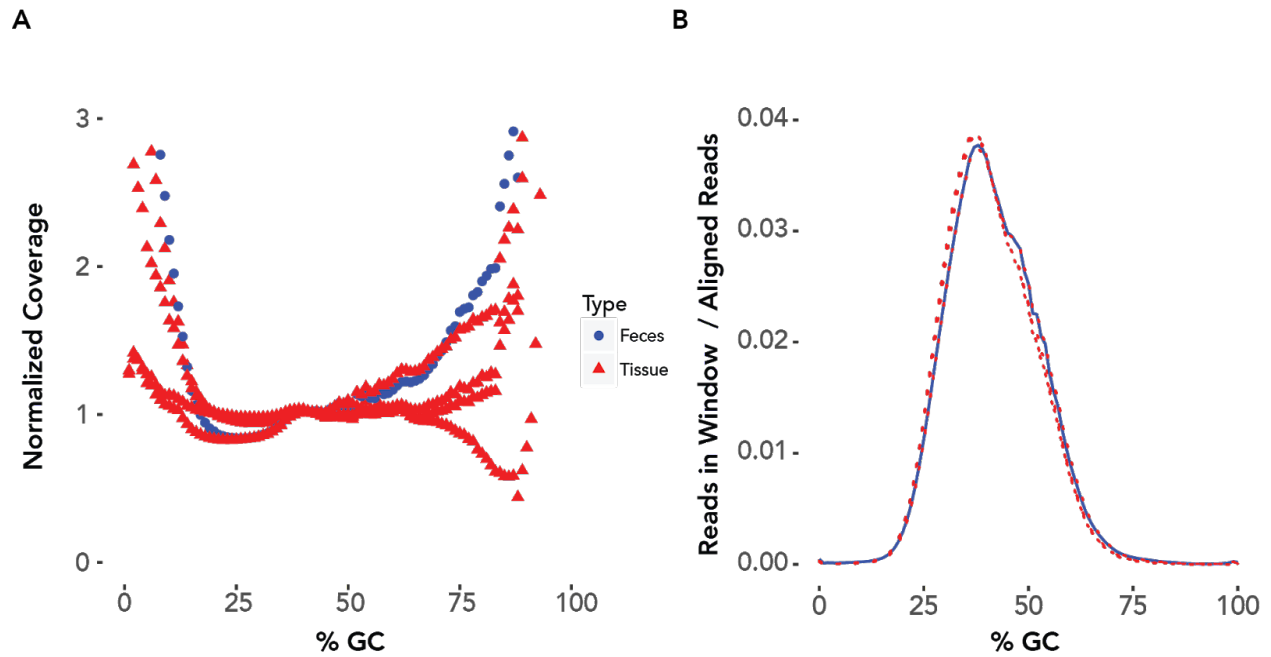
Figure 3: Mapping percentages of sequencing reads from RNAlater preserved fDNA libraries prepared with FACS for A) all samples, and B) individual libraries. C) Increase in mapping rate for RNAlater preserved samples. D) Relationship between mapped read duplication and number of cells with LOESS smoothing. The duplicate rate decreases sharply once a threshold of about 1,000 cells is reached.



541
542

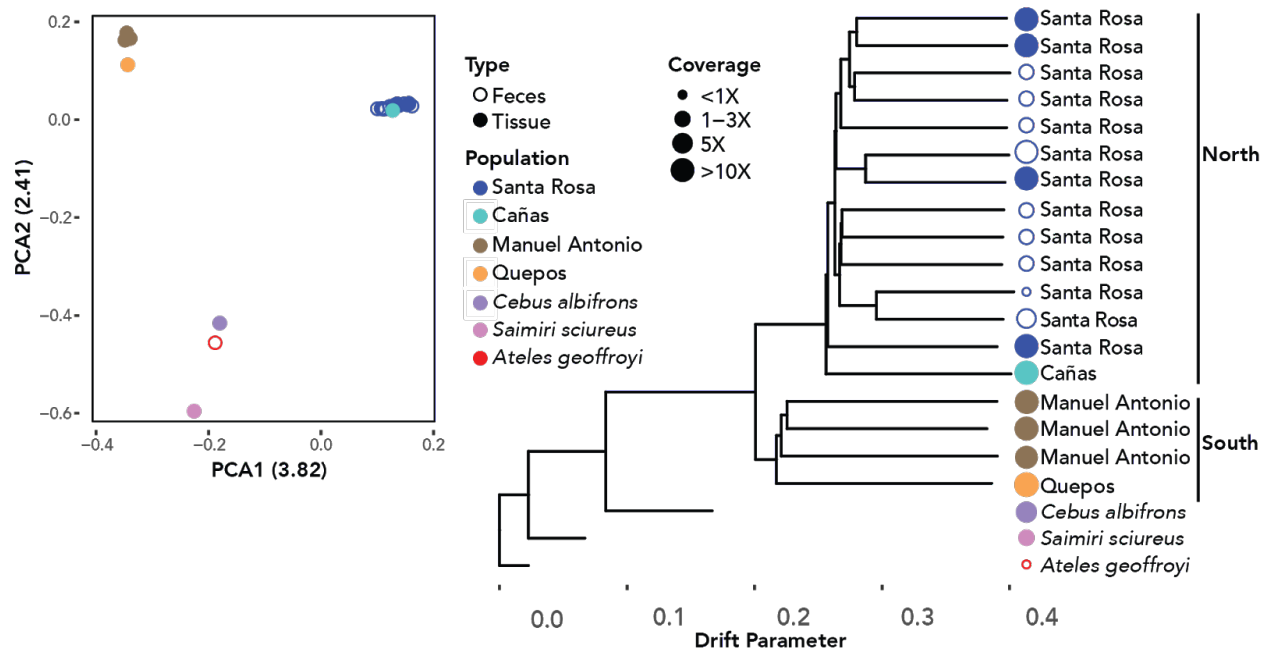
543 **Figure 4:** A) Density of genomic coverage of high coverage genomes from Santa Rosa. B)
544 Average coverage per 100 KB window along the largest scaffold of the *C. imitator* 1.0 reference
545 genome. C) Density of 1 MB windows at varying levels of heterozygosity along the entire
546 genome. D) Heterozygosity of 100 KB windows along the largest scaffold of the *C. imitator* 1.0
547 reference genome. The top two genomes (SSR-CR and SSR-ML) are from siblings. The order of
548 individuals in figures B and D correspond to that of figure A.

549



550
551
552
553
554
555
556
557

Figure 5: Percent of GC content across the genome for the four tissue (red) and one fecal (blue) samples from Sector Santa Rosa. GC content does not substantially differ for each type of sample. A) Average normalized coverage at each percentage of GC. B) Number of reads per 100 bp window (scaled by the number aligned reads) at each percentage of GC.



558
559

560 **Figure 6:** Left: Principal components of 14 fecal and 10 blood/tissue libraries from white faced
561 capuchin and three outgroups. Right: Maximum likelihood tree of 9 fecal and 10 blood/tissue
562 libraries. Samples with less 0.5X coverage were excluded. Among the white-faced capuchin
563 samples, individuals from northern (dry forest) and southern (wet forest) regions form the
564 primary split; secondary splits reflect the individuals from different sites within regions.
565

566
567
568
569

TABLES:

Table 1: Origins and preservation information for *Cebus imitator* samples.

Sample	Region	Site	Sample Type	Preservation
SSR-NM	North	Sector Santa Rosa	Feces	RNAlater
SSR-TY	North	Sector Santa Rosa	Feces	RNAlater
SSR-FN	North	Sector Santa Rosa	Feces	RNAlater
SSR-LE	North	Sector Santa Rosa	Feces	RNAlater
SSR-CH	North	Sector Santa Rosa	Feces	RNAlater
SSR-FG	North	Sector Santa Rosa	Feces	RNAlater
SSR-LU	North	Sector Santa Rosa	Feces	RNAlater
SSR-SN	North	Sector Santa Rosa	Feces	RNAlater
SSR-KI	North	Sector Santa Rosa	Feces	RNAlater
SSR-RF	North	Sector Santa Rosa	Feces	RNAlater
SSR-ML	North	Sector Santa Rosa	Feces	RNAlater
SSR-PR	North	Sector Santa Rosa	Feces	RNAlater
SSR-SB1	North	Sector Santa Rosa	Feces	RNAlater
SSR-SB2	North	Sector Santa Rosa	Feces	RNAlater
SSR-ML	North	Sector Santa Rosa	Feces	Frozen
SSR-FL	North	Sector Santa Rosa	Feces	Frozen
SSR-CR	North	Sector Santa Rosa	Tissue	Frozen
SSR-FL	North	Sector Santa Rosa	Tissue	Frozen
SSR-TH	North	Sector Santa Rosa	Tissue	Frozen
SSR-T5-1	North	Sector Santa Rosa	Tissue	Frozen
SSR-RM08	North	Sector Santa Rosa	Tissue	Frozen
CNS-HE	North	Cañas	Blood	Frozen
KSTR29	South	Manuel Antonio	Blood	Frozen
KSTR116	South	Manuel Antonio	Blood	Frozen
KSTR159	South	Manuel Antonio	Blood	Frozen
KSTR64	South	Quepos	Blood	Frozen

570
571

572 **Table 2:** FACS and mapping results from *Cebus* and *Ateles* fecal samples
 573

Monkey	Library	Cells	PCR Cycles	Total DNA (ng)	% Mapping					X Coverage
					BWA mem	BBsplit <i>Cebus</i>	Unique <i>Cebus</i>	Duplicate <i>Cebus</i>	BBsplit Human	
SSR-ML	SSR-ML Frozen	2546	11	10.50	96	90	85	5	1.25	11.7
SSR-FL	SSR-FL	4405	12	6.72	80	42	40	3	15.77	4.4
SSR-FN	SSR-FN	62601	8	21.50	97	93	86	6	0.81	2.8
SSR-FG	SSR-FG	580	10	9.75	94	84	60	24	2.86	2.0
SSR-LU	SSR-LU	8998	10	8.00	93	84	72	11	0.89	2.0
SSR-ML	SSR-ML RNAlater	42837	10	8.26	88	67	63	4	1.08	1.9
SSR-TY	SSR-TY	140	10	7.70	98	94	64	30	1.24	1.5
SSR-SB	SSR-SB 2	129	10	9.00	79	60	39	22	1.00	1.1
SSR-SB	SSR-SB 1	11944	10	6.25	55	11	8	2	0.61	
SSR-KI	SSR-KI	612	10	9.00	93	78	45	33	5.80	1.0
SSR-RF	SSR-RF	257	10	10.00	92	78	38	40	5.18	0.7
SSR-NM	SSR-NM	3336	11	3.38	98	95	92	3	0.66	0.4
SSR-CH	SSR-CH	957	11	4.06	93	85	80	5	0.74	0.4
SSR-LE	SSR-LE	1612	11	2.96	96	91	81	11	0.91	0.3
SSR-SN	SSR-SN	1866	11	3.96	92	79	74	6	3.07	0.2
SSR-PR	SSR-PR	12316	11	3.13	78	64	58	6	0.68	0.1
Spider	Spider 1	4026	12	6.96	54	12	11	1	2.83	0.4
	Spider 2	602	11	4.50	49	12	10	1	1.82	
	Median (<i>Cebus</i>)*	2079		7.85	93	82	63	6	1	

574
 575
 576
 577

578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600

REFERENCES

- Arandjelovic, M., & Vigilant, L. (2018). Non-invasive genetic censusing and monitoring of primate populations. *American Journal of Primatology*. doi:10.1002/ajp.22743
- Barrett, M. T., Glogovac, J., Prevo, L. J., Reid, B. J., Porter, P., & Rabinovitch, P. S. (2002). High-quality RNA and DNA from flow cytometrically sorted human epithelial cells and tissues. *BioTechniques*, 32(4), 888–90, 892, 894, 896.
- Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*, 30(15), 2114–2120.
- Bushnell, B. (2016). BBMap short read aligner. *University of California, Berkeley, California*. URL [Http://sourceforge.net/projects/bbmap](http://sourceforge.net/projects/bbmap).
- Chiou, K. L. (2017). *Population Genomics of a Baboon Hybrid Zone in Zambia* (PhD Thesis). Washington University in St. Louis. Retrieved from <https://doi.org/10.7936/K7348HS3>
- Chiou, K. L., & Bergey, C. M. (2018). Methylation-based enrichment facilitates low-cost, noninvasive genomic scale sequencing of populations from feces. *Scientific Reports*, 8(1), 1975.
- Corlett, R. T. (2017). A Bigger Toolbox: Biotechnology in Biodiversity Conservation. *Trends in Biotechnology*, 35(1), 55–65.
- de Manuel, M., Kuhlwilm, M., Frandsen, P., Sousa, V. C., Desai, T., Prado-Martinez, J., ... Marques-Bonet, T. (2016). Chimpanzee genomic diversity reveals ancient admixture with bonobos. *Science*, 354(6311), 477–481.
- Fedigan, L., & Rose-Wiles, L. (1996). See how they grow: Tracking capuchin monkey populations in a regenerating Costa Rican dry forest. In M. A. Norconk, A. L. Rosenberger, & P. A. Garber (Eds.), *Adaptive radiations of Neotropical primates* (pp. 289–307). Springer.

- 601 Florell, S. R., Coffin, C. M., Holden, J. A., Zimmermann, J. W., Gerwels, J. W., Summers, B.
602 K., ... Leachman, S. A. (2001). Preservation of RNA for functional genomic studies: a
603 multidisciplinary tumor bank protocol. *Modern Pathology: An Official Journal of the United*
604 *States and Canadian Academy of Pathology, Inc*, 14(2), 116–128.
- 605 Guschanski, K., Krause, J., Sawyer, S., Valente, L. M., Bailey, S., Finstermeier, K., ...
606 Savolainen, V. (2013). Next-generation museomics disentangles one of the largest primate
607 radiations. *Systematic Biology*, 62(4), 539–554.
- 608 Hernandez-Rodriguez, J., Arandjelovic, M., Lester, J., de Filippo, C., Weihmann, A., Meyer,
609 M., ... Marques-Bonet, T. (2017). The impact of endogenous content, replicates and pooling
610 on genome capture from faecal samples. *Molecular Ecology Resources*. doi:10.1111/1755-
611 0998.12728
- 612 Kuderna, L. F. K., Lizano, E., Julia, E., Gomez-Garrido, J., Serres-Armero, A., Kuhlwilm, M., ...
613 Marques-Bonet, T. (2018, June 13). *Selective single molecule sequencing and assembly of a*
614 *human Y chromosome of African origin*. *bioRxiv*. doi:10.1101/342667
- 615 Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler
616 transform. *Bioinformatics*, 25(14), 1754–1760.
- 617 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., ... DePristo,
618 M. A. (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-
619 generation DNA sequencing data. *Genome Research*, 20(9), 1297–1303.
- 620 Perry, G. H. (2014). The Promise and Practicality of Population Genomics Research with
621 Endangered Species. *International Journal of Primatology*, 35(1), 55–70.
- 622 Perry, G. H., Marioni, J. C., Melsted, P., & Gilad, Y. (2010). Genomic-scale capture and
623 sequencing of endogenous DNA from feces. *Molecular Ecology*, 19(24), 5332–5344.

- 624 Pickrell, J. K., & Pritchard, J. K. (2012). Inference of population splits and mixtures from
625 genome-wide allele frequency data. *PLoS Genetics*, *8*(11), e1002967.
- 626 Prado-Martinez, J., Sudmant, P. H., Kidd, J. M., Li, H., Kelley, J. L., Lorente-Galdos, B., ...
627 Marques-Bonet, T. (2013). Great ape genetic diversity and population history. *Nature*,
628 *499*(7459), 471–475.
- 629 Price, A. L., Patterson, N. J., Plenge, R. M., Weinblatt, M. E., Shadick, N. A., & Reich, D.
630 (2006). Principal components analysis corrects for stratification in genome-wide association
631 studies. *Nature Genetics*, *38*(8), 904–909.
- 632 Rinke, C., Lee, J., Nath, N., Goudeau, D., Thompson, B., Poulton, N., ... Woyke, T. (2014).
633 Obtaining genomes from uncultivated environmental microorganisms using FACS-based
634 single-cell genomics. *Nature Protocols*, *9*(5), 1038–1048.
- 635 Sasaki, D. T., Dumas, S. E., & Engleman, E. G. (1987). Discrimination of Viable and Non-
636 Viable Cells Using Propidium Iodide in Two Color Immunofluorescence. *Alan R. Liss, Inc.*
637 *Cytometry*, *8*, 413–420.
- 638 Snyder-Mackler, N., Majoros, W. H., Yuan, M. L., Shaver, A. O., Gordon, J. B., Kopp, G. H., ...
639 Tung, J. (2016). Efficient Genome-Wide Sequencing and Low-Coverage Pedigree Analysis
640 from Noninvasively Collected Samples. *Genetics*, *203*(2), 699–714.
- 641 van der Valk, T., Lona Durazo, F., Dalén, L., & Guschanski, K. (2017). Whole mitochondrial
642 genome capture from faecal samples and museum-preserved specimens. *Molecular Ecology*
643 *Resources*, *17*(6), e111–e121.
- 644 Wall, J. D., Schlebusch, S. A., Alberts, S. C., Cox, L. A., Snyder-Mackler, N., Nevenon, K.
645 A., ... Tung, J. (2016). Genomewide ancestry and divergence patterns from low-coverage

646 sequencing data reveal a complex history of admixture in wild baboons. *Molecular Ecology*,
647 25(14), 3469–3483.

648 Zaitoun, I., Erickson, C. S., Schell, K., & Epstein, M. L. (2010). Use of RNAlater in
649 fluorescence-activated cell sorting (FACS) reduces the fluorescence from GFP but not from
650 DsRed. *BMC Research Notes*, 3, 328.