

1 **Title:** Genetic characterization of a captive marmoset colony using genotype-by-sequencing

2

3 **Authors:** Cole SA^{1,2}, Lyke MM^{1,2}, Christensen C^{1,2}, Newman D^{1,2}, Bagwell A^{1,2}, Galindo S^{1,2},
4 Glenn J², Layne-Colon D^{1,2}, Sayers K^{1,2}, Tardif SD¹, Cox LA³, Ross CN^{1,2}, Cheeseman IH^{1,2}.

5

6 **Author Affiliations:**

7 ¹Southwest National Primate Research Center, San Antonio, TX;

8 ²Texas Biomedical Research Institute, San Antonio, TX;

9 ³Center for Precision Medicine, Wake Forest University School of Medicine, Winston-Salem, NC

10

11 **ABSTRACT:**

12 The marmoset is a fundamental non-human primate model for the study of aging, neurobiology,
13 and many other topics. Genetic management of captive marmoset colonies is complicated by
14 frequent chimerism in the blood and other tissues, a lack of tools to enable cost-effective, genome-
15 wide interrogation of variation, and historic mergers and migrations of animals between colonies.
16 We implemented genotype-by-sequencing (GBS) of hair follicle derived DNA (a minimally
17 chimeric DNA source) of 82 marmosets housed at the Southwest National Primate Research
18 Center (SNPRC). Our primary goals were the genetic characterization of our marmoset population
19 for pedigree verification and colony management and to inform the scientific community of the
20 functional genetic makeup of this valuable resource. We used the GBS data to reconstruct the
21 genetic legacy of recent mergers between colonies, to identify genetically related animals whose
22 relationships were previously unknown due to incomplete pedigree information, and to show that
23 animals in the SNPRC colony appear to exhibit low levels of inbreeding. Of the >99,000 single-
24 nucleotide variants (SNVs) that we characterized, >9,800 are located within gene regions known
25 to harbor pathogenic variants of clinical significance in humans. Overall, we show the combination
26 of low-resolution (sparse) genotyping using hair follicle DNA is a powerful strategy for the genetic

27 management of captive marmoset colonies and for identifying potential SNVs for the development
28 of biomedical research models.

29

30 **Keywords:** *Callithrix jacchus*; hair follicle DNA; pedigree; genetic ancestry; biomedical
31 research; captive non-human primates

32

33 **INTRODUCTION:**

34 The common marmoset (*Callithrix jacchus*) is an important non-human primate (NHP) model
35 for biomedical research (Miller, 2017; Ross & Salmon, 2019; Servick, 2018). The marmoset
36 provides unique practical advantages relative to other NHPs, including small body size, rapid
37 reproductive maturation, high fecundity, compressed life cycle, and comparative ease of handling
38 (Abbott et al., 2003; Kishi et al., 2014; Tardif & Ross, 2019; Tardif et al., 2003). Additionally,
39 aspects of their social behaviors and communication closely resemble those observed in humans
40 (Miller et al., 2016). These advantages have contributed to the recent use of marmosets for
41 developing new disease models using gene editing techniques and translational studies focused
42 on gene therapy (Kishi et al., 2014; National Academies of Sciences & Medicine, 2019).
43 Marmosets have also been utilized as a model species for a number of research areas, including
44 aging, neuroscience, and infectious disease (National Academies of Sciences & Medicine, 2019),
45 and new applications continue to be explored.

46 Given this increasing importance, there has been recent interest in genetic characterization
47 of existing captive colonies to develop colony management strategies that maintain their genetic
48 diversity and long-term value (reviewed in Harding 2017). Several approaches have been used
49 for population genetic characterization of NHPs, including whole genome sequencing (WGS),
50 whole exome sequencing (WES), and genotype-by-sequencing (GBS) (Bimber et al., 2016). The
51 first whole genome sequence arising from a common marmoset was generated from the genome
52 of an animal housed at the Southwest National Primate Research Center (SNPRC) (MGSAC,

53 2014). The genome was sequenced using Sanger sequencing (6x) and a whole-genome shotgun
54 approach (MGSAC, 2014). The marmoset genome assembly was then improved through deep
55 sequencing of a marmoset housed in Japan using next generation sequencing (NGS) (Sato et
56 al., 2015), and subsequent efforts continue to increase the quality and coverage (Rogers & del
57 Rosario, 2019).

58 An important aspect of the genetic characterization of NHPs as models for human health and
59 disease is the identification of functional single-nucleotide variants (SNVs) (Cline & Karchin, 2011;
60 Xue et al., 2016). SNVs represent common genetic variation that can be useful for identifying
61 variants that influence gene functions (Jasinska, 2020). Of particular interest are SNVs located in
62 genes with known phenotypic consequences associated with health and disease in humans.
63 Potentially functional SNVs have been identified and cataloged for several NHP species, including
64 rhesus macaques (*Macaca mulatta*; Bimber et al. 2017, Xue et al. 2016) and vervet monkeys
65 (*Chlorocebus sabaues*; Huang et al. 2015) with the goal of identifying sites relevant to studies of
66 human health, though this has not yet been extensively done for marmosets. Once identified,
67 these variants can be explored as potential orthologs to human disease genes, further facilitating
68 the development of effective marmoset models for biomedical research.

69 A special consideration when interpreting marmoset genomic data is that marmoset
70 littermates are chimeric (Benirschke et al., 1962). Marmoset tissues have a range of chimerism,
71 with hematopoietic-derived tissues showing more extensive chimerism (Ross et al., 2007;
72 Sweeney et al., 2012; Takabayashi & Katoh, 2015). To minimize the effect of chimerism when
73 assessing genetic diversity based upon sequence analyses (sequencing or genotyping),
74 investigators choose tissues with minimal chimerism such as hair, skin, or nail (Silva et al., 2017;
75 Takabayashi & Katoh, 2015). Due to the low yield of DNA derived from some of these sources,
76 cultured fibroblasts have been used for more comprehensive molecular approaches such as WGS
77 (NPRC Marmoset Genomics Working Group). To minimize chimerism in our samples, we

78 extracted DNA from hair follicles, which have far lower levels of chimerism than blood (Ross et
79 al., 2007) and yield more DNA of higher quality than finger nails.

80 The marmoset colony housed at the SNPRC is particularly valuable because of its genetic
81 history. It is an outbred population that can be traced for up to 12 generations, with a founder
82 population of 120 and a current population of \approx 400. Based upon pedigree, the SNPRC population
83 is particularly diverse, with founders from the University of Zurich, the University of London,
84 Marmoset Research Center-Oak Ridge (MARCOR), Osage Research Primates (ORP), Harlan,
85 Wisconsin National Primate Research Center (WNPRC), NIHCHD (National Institutes of Health-
86 Child Health and Development), New England NPRC (NEPRC), and Worldwide Primates (WWP).
87 The SNPRC colony manager uses pedigree analyses and relatedness tools to evaluate and
88 create breeding pairs that protect the genetic diversity of the colony. The current research will
89 increase our knowledge of the genetic diversity and relatedness of these important primate
90 resources.

91 The goal of this study was to develop a genome-wide genetic resource for the SNPRC
92 pedigreed marmoset colony that can be used in a cost-effective manner for: (1) pedigree
93 verification; (2) assisting with colony management; and (3) informing the scientific community
94 regarding the functional genetic makeup of this valuable colony resource. Although WGS costs
95 continue to decrease and low coverage WGS can be done relatively inexpensively, confidently
96 calling variants in potentially chimeric samples requires at least moderate coverage. For initial
97 development of this resource, we thought it was important to provide high resolution single
98 nucleotide variant data in exons, which are typically the highest priority SNVs for investigators.
99 While WES shows promise in characterizing NHP genomic variation (Chan et al., 2021)
100 commercially available WES currently uses human DNA sequences for exon capture, and exons
101 for some marmoset genes would be missed. For example, Chan and colleagues (2021) showed
102 that 6.5% of the coding exons, 32% of the 5' UTR exons, and 9.8% of the 3' UTR exons were not
103 captured using the human exon capture kit. Therefore, we chose to use a species agnostic GBS

104 approach where we can enrich for coding regions of the genome by restriction enzyme selection
105 and generate moderate coverage for each read at about one-fifth the cost per sample of WES.

106

107 **METHODS:**

108 SNPRC marmoset colony:

109 The SNPRC marmoset breeding colony was originally established in 1993 by Dr. Suzette
110 Tardif at the University of Tennessee-Knoxville with ten founding animals imported from Zurich,
111 Switzerland (University of Zurich) and 35 from the United Kingdom (UK; University of London).
112 The colony added animals from MARCOR, ORP, and the NICHD colony prior to moving the
113 colony to the SNPRC in 2001. Additional animals were imported for projects and production
114 starting in 2005 including WNPRC, Harlan, NEPRC, NIH and WWP resulting in 120 founding
115 animals. In 2015 the colony was further expanded through a merger between the SNPRC and
116 NEPRC colonies, introducing more than 180 new animals to the SNPRC colony. Several
117 marmosets have also returned to the SNPRC colony following sales to outside institutions. Most
118 of the marmoset colony belongs to one pedigree, which includes >3,000 animals (~400 living)
119 spanning 12 generations in depth as of April 2023. Hair follicle DNA was collected from 82 animals
120 for GBS, and used for sequencing and evaluation. This research protocol was approved by the
121 Texas Biomedical Research Institutional Animal Care and Use Committee #927CJ.

122

123 Isolation of DNA

124 For this study, hair follicle samples were opportunistically collected from 82 living animals.
125 Hair and follicles were collected at the time of a physical exam or of euthanasia. A site on the
126 body was chosen that was visually free of scent marking secretions, such as the tail or lower back.
127 A clump of about 50-100 hairs was grasped with clean forceps and pulled to remove hair with the
128 follicles intact. The clump was placed follicle end first into a sterile screw cap tube and frozen at -
129 80C. For DNA isolation, the hair and follicles were rinsed in 2mls of PBS (2x) prior to the start of

130 the DNA isolation procedure, which used the QIAmp DNA Mini Kit (Qiagen) following the
131 manufacturer's protocol for tissues, including the optional RNase step.

132

133 Genotype-by-sequencing (GBS)

134 GBS was completed using the DNA from 82 individual marmosets. GBS libraries were
135 prepared following the GBS method developed by Elshire et al. (2011) and optimized for the
136 marmoset genome. In brief, 50 ng of genomic DNA was digested in 20 µl reactions with 4U Avall.
137 Adapters with barcodes from 4-9 bp in length, designed using the GBS Barcode Generator
138 (deenabio.com/services/gbs-adapters) were then ligated to the digestion products in 50 µl
139 reactions using 400 cohesive end units of T4 DNA ligase (NEB). Following ligation, samples were
140 pooled and purified using the QIAquick PCR Purification Kit (Qiagen). Pooled DNA libraries were
141 amplified in 50 µl reactions with 1x NEBNext High Fidelity PCR Master Mix (NEB) and 12.5 pmol
142 PCR primers containing complementary sequences for adapter-ligated DNA. Temperature
143 cycling consisted of 72° for 5 min, 98° for 30 s, followed by 18 cycles of 98° for 10 s, 65° for 30 s,
144 and 72° for 30 s, with a final extension at 72° for 5 min. Amplified libraries were purified using the
145 QIAquick PCR Purification Kit (Qiagen) and the quality and quantity of each library assessed
146 using the Agilent DNA 1000 chip (Agilent Technologies) and KAPA Library Quantification Kit
147 (Kapa Biosystems), respectively. The final DNA libraries were hybridized to Rapid Run Flow Cells
148 (Illumina) for cluster generation using the TruSeq™ PE Cluster Kit (Illumina) and sequenced with
149 the TruSeq™ SBS Kit (Illumina) on an HiSeq 2500 (Illumina) using a 150-cycle paired-end
150 sequencing run.

151 Sequence reads were demultiplexed with GBSX v1.0.1 (Herten et al. 2015;
152 github.com/GenomicsCoreLeuven/GBSX). Sequence data were imported into Partek Flow
153 (Partek, Inc.). Sequence reads were trimmed and filtered based on a minimum quality score
154 (Phred) of 30 and a minimum read length of 25. The filtered reads were aligned to the marmoset
155 C_jacchus3.2.1 assembly using BWA-MEM alignment tool (Li, 2013) Single nucleotide variants

156 (SNVs) were detected using a minimum Phred quality score of 30 and filtered further using GATK
157 to include only those that had a minimum read depth of 5 and a minimum log-odds ratio of 3.05
158 to ensure high-quality variants. SNVs were annotated and effects predicted using Ensembl
159 Variant Effect Predictor (VEP, McLaren et al. 2016; [useast.ensembl.org/info/docs/tools/vep/
160 index.html](http://useast.ensembl.org/info/docs/tools/vep/index.html)). Sequences have been deposited in NCBI with Accession numbers SRR18101012-
161 SRR18100931.

162

163 Generating a high confidence genetic database

164 We retained autosomal, biallelic SNV data from the hair follicle DNA samples, and loaded the
165 data into SOLAR (Almasy and Blangero 1998) for initial evaluation of allele frequencies and
166 descriptive tallies. A modified pedigree based on the 82 original samples, their parents, and their
167 descendants was developed for use with all pedigree-based analyses (n=822). One of these
168 animals was unrelated to the others and was therefore not included in pedigree and kinship
169 estimation, though included in all other analyses. INFER, which is a program within the PEDSYS
170 software system (Dyke, 1996), examines each offspring-father-mother triplet and, when possible,
171 adds missing alleles and genotypes according to the Mendelian laws of transmission. The
172 program iterates through the pedigree as many times as necessary until no more assignments
173 can be made. The inferred data were then combined with the pedigree data for further analyses.

174 Following the inference of new genotypes, we used SimWalk2 (Sobel et al., 2001) to identify
175 and remove genotypes inconsistent with Mendelian properties within families including the grand-
176 parental generation, as well as distributions of alleles within entire sibships. SimWalk2 reports the
177 *overall* probability of mistyping at each observed genotype (in fact, at each observed allele). When
178 genotypes were flagged with a significant probability of mistyping, they were removed from the
179 dataset. The resultant data file was again evaluated with SOLAR (Almasy & Blangero, 1998) to
180 create a summary list of all variant loci. This summary provided the number of samples counted
181 per variant, the SNV major and minor allele frequencies and the associated p-values of a test of

182 whether the Hardy-Weinberg equilibrium (HWE) holds. From the summary list of all variant loci,
183 we selected and removed all SNVs that had low call rates keeping only those SNVs where 95%
184 of samples (78 of initial 82) were typed. The remaining set of annotated variants represent the
185 high confidence database carried forward for further analyses.

186

187 Statistical and Genetic Analysis

188 We performed inference of historical admixture in the colony with ADMIXTURE v1.3.0
189 (Alexander et al., 2015) using high quality SNVs from the 82 animals directly genotyped in this
190 study. In addition, we have substantial genotype data from 48 animals where genotypes were
191 inferred by inheritance in INFER, which were also included to increase coverage of the pedigree.
192 ADMIXTURE analysis was run unsupervised on these 130 animals using 1-5 theoretical ancestral
193 populations (K). Values of $K > 3$ failed cross validation and were discarded. We then integrated
194 publicly available WGS data from 9 animals from WNPRC ($n=2$), NEPRC ($n=2$), and SNPRC
195 ($n=5$) (MGSAC, 2014). We ran a supervised analysis ($K=3$) using the WGS samples incorporating
196 the population labels from each primate center.

197 We estimated kinship from genetic data using lcMLkin (Lipatov et al., 2015), using high quality
198 genotypes from 81 animals with GBS data (excluding the individual that was unrelated to the rest
199 of the pedigree). lcMLkin estimates kinship while accounting for incomplete data and the reduced
200 ability to capture both alleles at a given locus. Additionally, empirical kinship was directly inferred
201 from the pedigree records of 3,232 animals using the kinship2 package in R v4.2.2 (Sinnwell et
202 al., 2014). We found a strong correlation between mean read depth and the proportion of
203 heterozygous base calls. To account for this dependency in our analysis we performed linear
204 regression of proportion of heterozygous base calls using mean read depth and generation as
205 covariates. This was implemented in the lm function in R.

206 To investigate functionally significant genetic variation, we assessed the overlap of the GBS
207 SNVs with gene regions associated with immune function and inflammation, neurological traits,

208 aging, obesity, and diabetes using the NCBI “Gene” database
209 (<https://www.ncbi.nlm.nih.gov/gene>). The NCBI Gene database contains the known functions of
210 genes of many different species, including humans and NHPs. We also merged our GBS-
211 identified SNVs with the human ClinVar database ([ncbi.nlm.nih.gov/clinvar](https://www.ncbi.nlm.nih.gov/clinvar)) to identify SNVs that
212 were located in genes with identified human variants that result in phenotypes with potentially
213 clinical outcomes.

214

215 **RESULTS:**

216 Genotype-by-sequencing identifies SNV with potential functional significance in marmosets.

217 Forty-five percent of reads aligned to the transcriptome. As the transcriptome is 3.3% of the
218 marmoset genome, this indicates that the choice of restriction enzyme to target coding regions
219 was appropriate for generating gene-centric GBS data. The GBS data from the 82 marmoset hair
220 follicle DNA samples yielded 231,317 biallelic SNV loci. Excluding loci from chromosomes X and
221 Y yielded 216,015 SNVs. After implementing our quality control procedures which include
222 Mendelian error cleaning and including only SNV with high call rates, we obtained a high quality
223 set of SNV for further analysis. Table 1 lists the types of potential functional effects identified for
224 this marmoset SNV dataset. We assessed the overlap of the GBS SNVs with gene regions
225 associated with immune function and inflammation, neurological traits, aging, obesity, and
226 diabetes in humans using the NCBI Gene database. We identified 7,738 SNVs associated with
227 immunity/inflammation, 289 associated with neurological traits, 2,544 with aging, 5,715 with
228 obesity, and 5,402 with diabetes. We also merged our GBS-identified SNV with the human
229 ClinVar database ([ncbi.nlm.nih.gov/clinvar](https://www.ncbi.nlm.nih.gov/clinvar)) and identified 9,897 SNVs that were located in genes
230 with identified human variants that result in phenotypes with potentially clinical outcomes (Table
231 1).

232

Number of called loci (all chromosomes, including multi-allelic)	263,575
Number of biallelic SNVs prior to Mendelian inheritance check (chr1-22 only)	216,015
Number SNVs after INFER and Mendelian Inheritance checking (chr1-22 only)	201,892
Number SNVs (inferred, Mendelian-cleaned) present in at least 95% of samples	99,439
Summaries for High Quality SNVs Data Set (n=99,439):	
Intergenic variant	44,622
Intron variant	36,199
Missense variant (coding missense variant)	1,704
SNVs in marmoset genes	54,817
SNVs in human genes	45,663
SNVs in gene regions harboring human clinical variants	20,242
Gene regions with clinical significance of "pathogenic or likely pathogenic" (ClinVar)	9,897
SNVs where major allele doesn't equal ref allele	13,018
SNVs with HIGH Impact Score	86
SNVs in genes associated with Immunity or Inflammation	7,738
SNVs in genes associated with Neurological Traits	289
SNVs in genes associated with Aging	2,544
SNVs in genes associated with Obesity	5,715
SNVs in genes associated with Diabetes	5,402

233

234

235 Admixture analyses reveal the impact of past mergers on colony structure.

236 To assess the impact of the colony mergers, we combined cleaned, imputed GBS data from

237 the 82 animals sequenced here with previously reported whole genome sequencing data and

238 imputed genotypes inferred solely from pedigree structure. For admixture analyses, we lowered

239 our threshold for SNV filtering to include 131,648 biallelic SNVs typed in >50% of individuals. We

240 estimated the number of ancestral populations giving rise to genetic diversity in the SNPRC

241 colony and the prevalence of ancestry components in the joint set (GBS and inferred) of 130

242 animals using ADMIXTURE. Cross-validation of cluster numbers $K > 3$ did not converge and were

243 excluded. At $K=2$ (the best supported number of clusters), SNPRC animals showed a clear

244 bimodal distribution of ancestry with 60.8% of animals deriving >95% of their ancestry from a

245 single component (Fig. 1A). These results were also supported at $K=3$ (Fig. 1B). Results of the
246 supervised admixture analysis with population labels for each primate center also indicate that
247 the ancestry of most animals is derived from a single population (Fig. 1C). Figure 1D shows the
248 known pedigree with each animal included in the admixture analysis shaded by the proportion of
249 NEPRC ancestry. There is a clear subdivision in the pedigree which recapitulates the colony
250 merging. Our data allows us to assess admixture in the SNPRC colony and capture the impact of
251 past colony mergers on population structure.

252

253 GBS data augments colony pedigree data.

254 For colony management, genetic data directly captures inbreeding and relatedness, and can
255 be used to inform breeding strategies and correct or enhance pedigree records. We estimated
256 kinship directly from the genetic data for the 82 animals with GBS data using lcMLkin. We found
257 this to have strong correlation to pedigree records ($r^2=0.78$, $p<2.2\times 10^{-16}$, Fig. 2), with a moderate
258 global inflation of kinship estimate (intercept from a linear model=0.036, Fig. 2A). We identified
259 12 animals who showed no observed kinship (<0.01) from pedigree records, though were close
260 relatives from genetic data ($\pi\hat{r}>0.1$). Subsequent inspection of the pedigree and animal transfer
261 records revealed the relationships of 11 of these animals (Fig. 2B). We generated a heat map to
262 compare the pedigree and GBS based estimates of relatedness (Fig. 2C). While both estimates
263 capture close familial relationships, GBS based estimates of relatedness show generally higher
264 estimates of relatedness (Fig. 2C).

265

266 Genetic diversity is not declining in the SNPRC colony.

267 A major concern in captive pedigrees is the erosion of diversity over time due to inbreeding.
268 We tested if there was a decline in the proportion of heterozygotes in successive generations in
269 our data. As there is a significant concern that read depth will influence the ability to accurately
270 identify heterozygous sites, we fit a linear model on the percentage of heterozygous sites in a

271 sample against the mean read depth and the pedigree generation. Both read depth and pedigree
272 generation were significant predictors of percentage of heterozygous sites (Table 2). Notably, a
273 more complex model including ancestry components did not show ancestry to be a significant
274 predictor suggesting the proportion of heterozygous sites was not variable between founding
275 populations. In both cases (increased read depth and more contemporary generations) there was
276 a positive relationship with percentage of heterozygous sites. While we do not explicitly derive an
277 estimate of the heterozygosity here, the average proportion of heterozygous sites in this
278 population is not decreasing, suggesting this is a genetically healthy breeding population
279 minimally impacted by inbreeding.

280
281
282
283

Table 2. Linear regression of proportion of heterozygous sites against pedigree generation and mean depth.

	Estimate (std error)	<i>t</i> -value	<i>p</i> -value
Generation	0.03 (0.007)	4.708	1.05x10 ⁻⁵
Mean Depth	1.95 (0.16)	12.297	4.92x10 ⁻²⁰

284
285
286

DISCUSSION:

287 Marmosets have become important resources for biomedical research (Marini et al., 2018).
288 They have been used as models for studying a range of human conditions, from aging to
289 infectious disease (Miller, 2017; Ross & Salmon, 2019). Marmosets are phylogenetically more
290 distant from humans than catarrhine primates, but there are several traits that make them better
291 suited for specific types of research. Their relatively shorter lifespans make them ideal models for
292 studies of primate aging (Ross & Salmon, 2019). Their higher fecundity and tendency towards
293 twinning are useful for studying aspects of pregnancy and for maintaining colony numbers (Abbott
294 et al., 2003; Tardif & Ross, 2019; Tardif et al., 2003). Additionally, their small body size and
295 relative ease of handling make housing and caring for marmosets more practical than other NHPs
296 (Abbott et al., 2003; Miller, 2017). The marmosets at the SNPRC comprise a controlled breeding
297 population with multi-generational pedigree data, making them ideal subjects for genetic

298 characterization. With the recent increase in genome sequencing of NHPs commonly used in
299 biomedical applications, the identification of SNVs corresponding to genes with known human
300 health outcomes has been highly beneficial. Using GBS on hair follicle DNA, we identified 216,015
301 high quality SNVs in 82 marmosets from the SNPRC colony (Table 1). After quality control
302 including Mendelian error correction and including only SNV with a high call rate and further
303 filtering using public databases, we were able to identify variants associated with immune function
304 and inflammation (n=7,738), neurological traits (n=289), aging (n=2,544), obesity (n=5,715), and
305 diabetes (n=5,402) in humans (Table 1). Identifying potentially functional health related SNVs
306 allows researchers to focus on specific regions of the genome and model genetic mechanisms of
307 human disease.

308 One of the key factors when considering animals as subjects for translational research is the
309 overall genetic health of the population (Harding, 2017; Haus et al., 2014). It is important to
310 minimize potential adverse health effects related to inbreeding that could confound experimental
311 design and results (Haus et al., 2014; Honess et al., 2010). Inbreeding and decreased genetic
312 diversity are concerns when breeding any captive population as gene flow is generally limited
313 (Harding, 2017; Haus et al., 2014). Early genetic studies suggested that overall genetic diversity
314 of marmosets and other callitrichids was historically low (Dixson et al., 1988; Faulkes et al., 2003;
315 Forman et al., 1986; Nievergelt et al., 2000; Watkins et al., 1991), so developing a low cost yet
316 effective way to measure and track genetic diversity in our population has been a priority. Using
317 the GBS data and the proportion of heterozygous sites as a measure of diversity, here we have
318 shown that the individuals in the primary pedigree at SNPRC are genetically healthy and in
319 general the colony appears to exhibit low levels of inbreeding. This is likely due to the diverse
320 provenances of our founding population and the ongoing efforts of our colony manager and others
321 to reduce the loss of diversity. As shown in our data, both read depth and pedigree generation
322 were predictors of the percentage of heterozygous sites, and both showed positive relationships
323 (Table 2). Our GBS approach and data QC ensured that we had a high confidence set of SNVs

324 for analyses, and caution that low pass sequencing methods may not reveal all informative
325 genetic variation. Regarding the effect of pedigree generation on heterozygosity, this was likely
326 due in part to recent colony mergers that introduced new individuals into the population. In
327 general, these migrations increase genetic diversity and decrease inbreeding, as was the case
328 with the marmoset colony mergers discussed above.

329 While migrations and animal transfers can help maintain or increase genetic diversity, when
330 animals are moved between NPRCs or other organizations, potential relationships between
331 individuals can be missing from the pedigree record. When new animals are imported, they are
332 considered to be unrelated to others in the colony, as there is no known relationship. Our GBS
333 data uncovered previously unknown relationships in twelve individuals in our colony (Figure 2).
334 Subsequent inspection of pedigree records allowed us to uncover some of those relationships.
335 For example, a male marmoset with several offspring in the SNPRC colony was transferred to
336 another location, where he then had offspring. One of his male offspring (CJTXGBS00082) was
337 transferred back into the SNPRC colony, where he had unknown half-siblings. Additionally, one
338 or both parents of three of the individuals (CJTXGBS00060, CJTXGBS00062, CJTXGBS00066)
339 were part of the founding population, and their relationships to each other were only revealed with
340 the GBS data. Prior to generating the GBS data, estimations of kinship and relatedness were
341 based solely on pedigree, which is extensive and is a valuable tool, but with some inherent
342 limitations due to potentially missing data. We have demonstrated that genetic data can play a
343 pivotal role in identifying potential errors in the pedigree for further exploration. These related
344 individuals, or their close relatives, might otherwise have been picked as breeding pairs,
345 inadvertently increasing levels of inbreeding in the population.

346 Another outcome of importing animals from other colonies is that genetic admixture occurs.
347 Geographically isolated populations tend to have genetic signatures due to changes in allele
348 frequencies based on local adaptations (Cheng et al., 2022). Genetic signatures can be
349 developed over long evolutionary periods and detected at the species level, but they can also be

350 driven by short term population isolation, such as in captive breeding colonies (MGSAC, 2014;
351 Schoener, 2011). When discrete populations are brought into contact and interbreed, genetic
352 admixture occurs. The extent of admixture can be measured over subsequent generations using
353 genetic data, such as GBS (Alexander et al., 2015). The ability to track admixture and ancestry is
354 beneficial when analyzing potential variation in disease risk and response when developing
355 animal models for biomedical research (Shriner, 2017). Here we were able to successfully identify
356 genetic signatures corresponding to colony of origin using admixture analyses of our GBS data.
357 We ran both supervised and unsupervised analyses and both reflected population structure at the
358 colony level. Because of their relatively short evolutionary histories and rapidly shifting allele
359 frequencies, migration of animals among NPRCs is driving population structure.

360 An important component of NHP colony management is the assessment and maintenance of
361 genetic health and diversity. This is especially critical for animals used in biomedical research,
362 where certain genetic traits may influence disease susceptibility or outcomes (Haus et al., 2014;
363 Honess et al., 2010). It is also critical to maintain genetic diversity for the endurance and
364 expansion of the colony itself as a long-term research resource. For marmosets, the choice of
365 tissue used for genetic analyses requires special consideration due to their chimeric nature. Hair
366 follicles are among the least chimeric tissues, are collected non-invasively, and yield high quality
367 DNA suitable for advanced genetic sequencing (Ross et al., 2007). Recent years have seen an
368 increase in commercially available sequencing options with decreasing costs, yet some methods,
369 such as WGS, are still cost prohibitive for population-level genetic screening, especially for those
370 colonies with limited resources. We have demonstrated that GBS provides high quality, affordable
371 data using sparse genetic characterization for population management and for assessing
372 ancestry and colony genetic health. The combination of hair follicle DNA with GBS represents a
373 successful, non-invasive, cost-effective approach for colony management and for understanding
374 evolutionary diversity of captive marmosets being used in biomedical applications.

375

376 **ACKNOWLEDGEMENTS:**

377 **Financial Support:** This research was supported by NIH-NCRR grant P51 RR013986 to the
378 Southwest National Primate Research Center.

379

380 **Conflict of interest:** The authors have declared no conflict of interest.

381

382 **REFERENCES:**

383

384 Abbott, D. H., Barnett, D. K., Colman, R. J., Yamamoto, M. E., & Schultz-Darken, N. J. (2003). Aspects of
385 common marmoset basic biology and life history important for biomedical research.
386 *Comparative medicine*, 53(4), 339-350.

387 Alexander, D. H., Shringarpure, S. S., Novembre, J., & Lange, K. (2015). Admixture 1.3 software manual.
388 *Los Angeles: UCLA Human Genetics Software Distribution.*

389 Almasy, L., & Blangero, J. (1998). Multipoint quantitative-trait linkage analysis in general pedigrees. *The*
390 *American Journal of Human Genetics*, 62(5), 1198-1211.

391 Benirschke, K., Anderson, J. M., & Brownhill, L. E. (1962). Marrow chimerism in marmosets. *Science*,
392 138(3539), 513-515.

393 Bimber, B. N., Raboin, M. J., Letaw, J., Nevenon, K. A., Spindel, J. E., McCouch, S. R., Cervera-Juanes, R.,
394 Spindel, E., Carbone, L., & Ferguson, B. (2016). Whole-genome characterization in pedigreed
395 non-human primates using genotyping-by-sequencing (GBS) and imputation. *BMC genomics*,
396 17(1), 1-16.

397 Chan, J., Yao, W., Howard, T. D., Hawkins, G. A., Olivier, M., Jorgensen, M. J., Cheeseman, I. H., Cole, S.
398 A., & Cox, L. A. (2021). Efficiency of whole-exome sequencing in old world and new world
399 primates using human capture reagents. *Journal of medical primatology*, 50(3), 176-181.

400 Cheng, J. Y., Stern, A. J., Racimo, F., & Nielsen, R. (2022). Detecting selection in multiple populations by
401 modeling ancestral admixture components. *Molecular biology and evolution*, 39(1), msab294.

402 Cline, M. S., & Karchin, R. (2011). Using bioinformatics to predict the functional impact of SNVs.
403 *Bioinformatics*, 27(4), 441-448.

404 Dixon, A., Hastie, N., Patel, I., & Jeffreys, A. (1988). DNA 'fingerprinting' of captive family groups of
405 common marmosets (*Callithrix jacchus*). *Folia Primatologica*, 51(1), 52-55.

406 Dyke, B. (1996). PEDSYS: a pedigree data management system. *Population Genetics Laboratory,*
407 *Department of Genetics, Southwest Foundation for Biomedical Research, San Antonio.*

408 Elshire, R. J., Glaubitz, J. C., Sun, Q., Poland, J. A., Kawamoto, K., Buckler, E. S., & Mitchell, S. E. (2011). A
409 robust, simple genotyping-by-sequencing (GBS) approach for high diversity species. *PLoS One*,
410 6(5), e19379.

411 Faulkes, C., Arruda, M., & Monteiro da Cruz, M. (2003). Matrilineal genetic structure within and among
412 populations of the cooperatively breeding common marmoset, *Callithrix jacchus*. *Molecular*
413 *Ecology*, 12(4), 1101-1108.

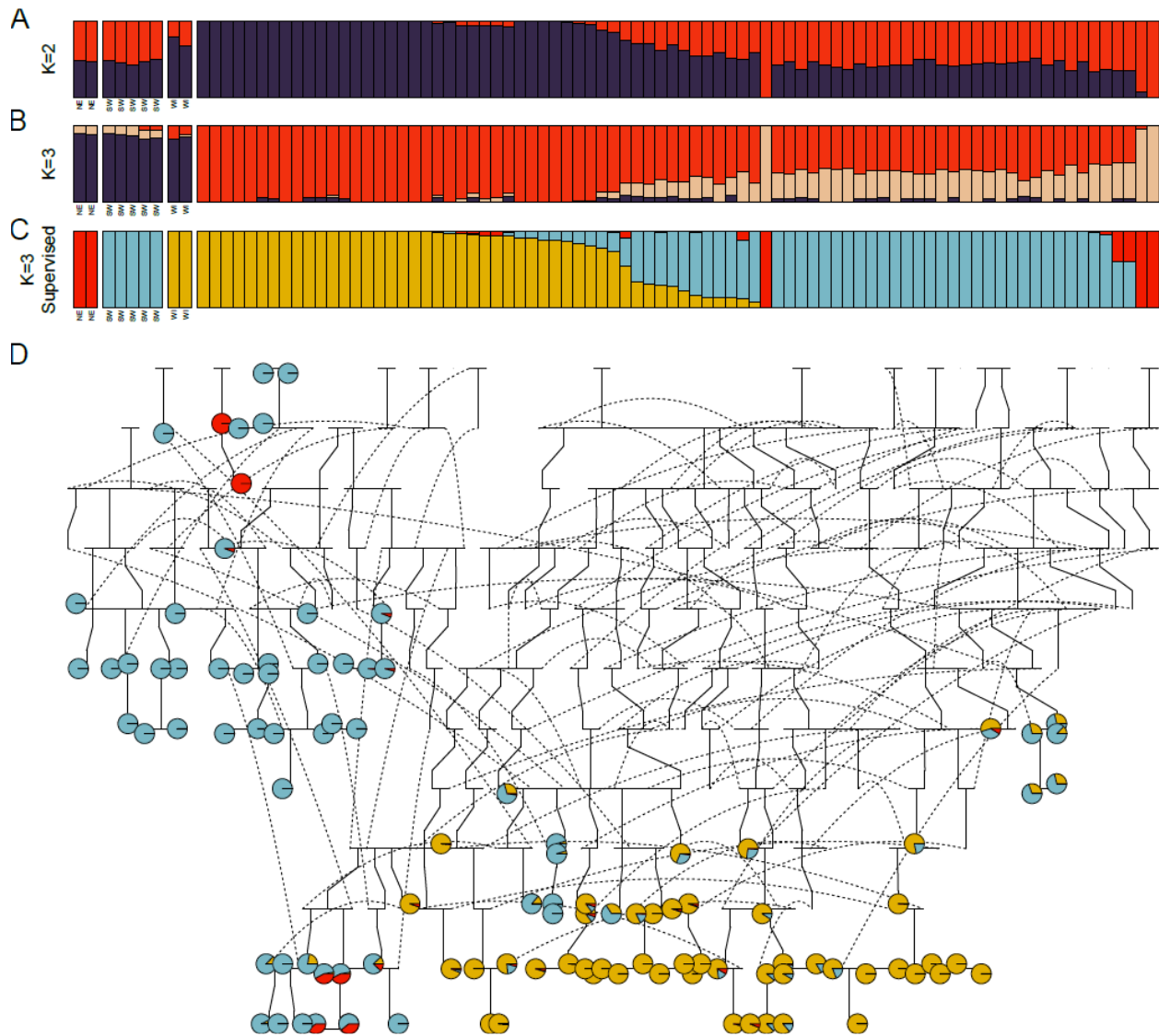
414 Forman, L., Kleiman, D. G., Bush, R. M., Dietz, J. M., Ballou, J. D., Phillips, L. G., Coimbra-Filho, A. F., &
415 O'Brien, S. J. (1986). Genetic variation within and among lion tamarins. *American Journal of*
416 *Physical Anthropology*, 71(1), 1-11.

417 Harding, J. D. (2017). Nonhuman primates and translational research: progress, opportunities, and
418 challenges. *ILAR journal*, 58(2), 141-150.

- 419 Haus, T., Ferguson, B., Rogers, J., Doxiadis, G., Certa, U., Rose, N. J., Teepe, R., Weinbauer, G. F., & Roos,
420 C. (2014). Genome typing of nonhuman primate models: implications for biomedical research.
421 *Trends in Genetics*, *30*(11), 482-487.
- 422 Herten, K., Hestand, M. S., Vermeesch, J. R., & Van Houdt, J. K. (2015). GBSX: a toolkit for experimental
423 design and demultiplexing genotyping by sequencing experiments. *BMC bioinformatics*, *16*(1), 1-
424 6.
- 425 Honess, P., Stanley-Griffiths, M., Narainapoulle, S., Naiken, S., & Andrianjalahatra, T. (2010). Selective
426 breeding of primates for use in research: consequences and challenges. *Animal welfare*, *19*(S1),
427 57-65.
- 428 Huang, Y. S., Ramensky, V., Jasinska, A. J., Jung, Y., Choi, O.-W., Cantor, R. M., Juretic, N., Wasserscheid,
429 J., Kaplan, J. R., & Jorgensen, M. J. (2015). Sequencing strategies and characterization of 721
430 vervet monkey genomes for future genetic analyses of medically relevant traits. *BMC biology*,
431 *13*(1), 1-10.
- 432 Jasinska, A. J. (2020). Resources for functional genomic studies of health and development in nonhuman
433 primates. *American Journal of Physical Anthropology*, *171*, 174-194.
- 434 Kishi, N., Sato, K., Sasaki, E., & Okano, H. (2014). Common marmoset as a new model animal for
435 neuroscience research and genome editing technology. *Development, growth & differentiation*,
436 *56*(1), 53-62.
- 437 Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv*
438 *preprint arXiv:1303.3997*.
- 439 Lipatov, M., Sanjeev, K., Patro, R., & Veeramah, K. R. (2015). Maximum likelihood estimation of
440 biological relatedness from low coverage sequencing data. *BioRxiv*, 023374.
- 441 Marini, R. P., Wachtman, L. M., Tardif, S. D., Mansfield, K., & Fox, J. G. (2018). *The common marmoset in*
442 *captivity and biomedical research*. Academic Press.
- 443 McLaren, W., Gil, L., Hunt, S. E., Riat, H. S., Ritchie, G. R., Thormann, A., Flicek, P., & Cunningham, F.
444 (2016). The ensembl variant effect predictor. *Genome biology*, *17*(1), 1-14.
- 445 MGSAC. (2014). The Marmoset Genome Sequencing and Analysis Consortium: The common marmoset
446 genome provides insight into primate biology and evolution. *Nature genetics*, *46*, 850-857.
447 <https://doi.org/doi.org/10.1038/ng.3042>
- 448 Miller, C. T. (2017). Why marmosets? *Developmental Neurobiology*.
- 449 Miller, C. T., Freiwald, W. A., Leopold, D. A., Mitchell, J. F., Silva, A. C., & Wang, X. (2016). Marmosets: a
450 neuroscientific model of human social behavior. *Neuron*, *90*(2), 219-233.
- 451 National Academies of Sciences, E., & Medicine. (2019). Care, Use, and Welfare of Marmosets as Animal
452 Models for Gene Editing-Based Biomedical Research: Proceedings of a Workshop.
- 453 Nievergelt, C. M., Digby, L. J., Ramakrishnan, U., & Woodruff, D. S. (2000). Genetic analysis of group
454 composition and breeding system in a wild common marmoset (*Callithrix jacchus*) population.
455 *International Journal of Primatology*, *21*, 1-20.
- 456 Rogers, J., & del Rosario, R. (2019). Marmoset Genomics and Genetic Diversity. Care, Use, and Welfare
457 of Marmosets as Animal Models for Gene Editing-Based Biomedical Research: Proceedings of a
458 Workshop,
- 459 Ross, C. N., French, J. A., & Orti, G. (2007). Germ-line chimerism and paternal care in marmosets
460 (*Callithrix kuhlii*). *Proceedings of the National Academy of Sciences*, *104*(15), 6278-6282.
- 461 Ross, C. N., & Salmon, A. B. (2019). Aging research using the common marmoset: Focus on aging
462 interventions. *Nutrition and healthy aging*, *5*(2), 97-109.
- 463 Sato, K., Kuroki, Y., Kumita, W., Fujiyama, A., Toyoda, A., Kawai, J., Iriki, A., Sasaki, E., Okano, H., &
464 Sakakibara, Y. (2015). Resequencing of the common marmoset genome improves genome
465 assemblies and gene-coding sequence analysis. *Scientific reports*, *5*(1), 1-8.

- 466 Schoener, T. W. (2011). The newest synthesis: understanding the interplay of evolutionary and
467 ecological dynamics. *science*, 331(6016), 426-429.
- 468 Servick, K. (2018). US labs clamor for marmosets. In: American Association for the Advancement of
469 Science.
- 470 Shriner, D. (2017). Overview of admixture mapping. *Current protocols in human genetics*, 94(1), 1.23. 21-
471 21.23. 28.
- 472 Silva, M. O., ARMADA, J. L. A., Verona, C. E. S., Heliodoro, G., & Nogueira, D. M. (2017). Cytogenetics and
473 molecular genetic analysis of chimerism in marmosets (*Callithrix*: Primates). *Anais da Academia*
474 *Brasileira de Ciências*, 89, 2793-2804.
- 475 Sinnwell, J. P., Therneau, T. M., & Schaid, D. J. (2014). The kinship2 R package for pedigree data. *Human*
476 *heredity*, 78(2), 91-93.
- 477 Sobel, E., Sengul, H., & Weeks, D. E. (2001). Multipoint estimation of identity-by-descent probabilities at
478 arbitrary positions among marker loci on general pedigrees. *Human heredity*, 52(3), 121-131.
- 479 Sweeney, C. G., Curran, E., Westmoreland, S. V., Mansfield, K. G., & Vallender, E. J. (2012). Quantitative
480 molecular assessment of chimerism across tissues in marmosets and tamarins. *BMC genomics*,
481 13(1), 1-7.
- 482 Takabayashi, S., & Katoh, H. (2015). Noninvasive genotyping of common marmoset (*Callithrix jacchus*) by
483 fingernail PCR. *Primates*, 56(3), 235-240.
- 484 Tardif, S. D., & Ross, C. N. (2019). Reproduction, Growth, and Development. In *The Common Marmoset*
485 *in Captivity and Biomedical Research* (pp. 119-132). Elsevier.
- 486 Tardif, S. D., Smucny, D. A., Abbott, D. H., Mansfield, K., Schultz-Darken, N., & Yamamoto, M. E. (2003).
487 Reproduction in captive common marmosets (*Callithrix jacchus*). *Comparative medicine*, 53(4),
488 364-368.
- 489 Watkins, D. I., Garber, T. L., Chen, Z. W., Toukatly, G., Hughes, A. L., & Letvin, N. L. (1991). Unusually
490 limited nucleotide sequence variation of the expressed major histocompatibility complex class I
491 genes of a New World primate species (*Saguinus oedipus*). *Immunogenetics*, 33, 79-89.
- 492 Xue, C., Raveendran, M., Harris, R. A., Fawcett, G. L., Liu, X., White, S., Dahdouli, M., Deiros, D. R., Below,
493 J. E., & Salerno, W. (2016). The population genomics of rhesus macaques (*Macaca mulatta*)
494 based on whole-genome sequences. *Genome research*, 26(12), 1651-1662.
- 495
- 496

497 **FIGURES:**
498



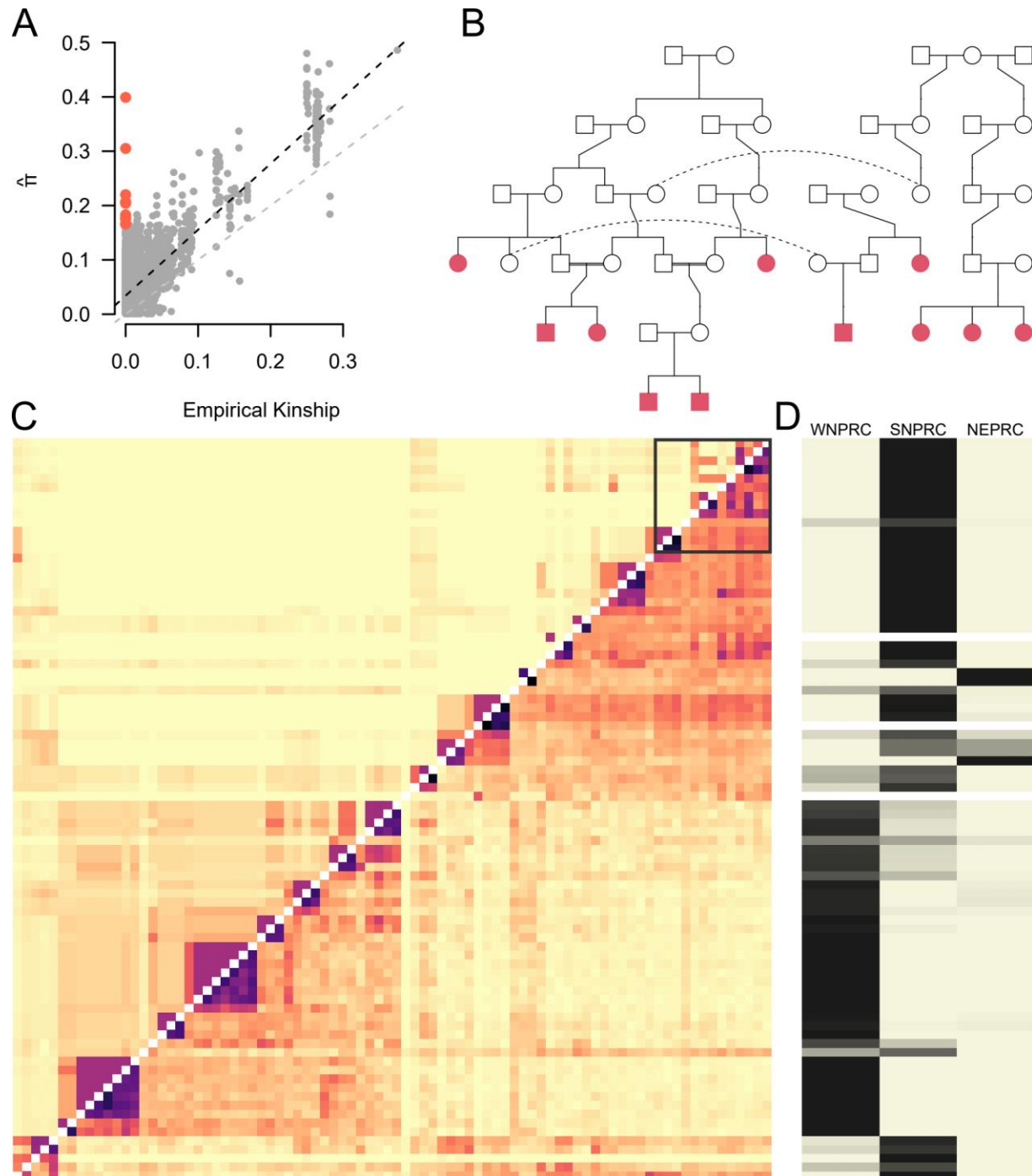
499
500

501 **Figure 1. Admixture in the SNPRC colony.** High quality SNVs from GBS from 82 marmosets
502 were integrated with whole genome sequencing from NEPRC (columns 1-2, n=2), SNPRC
503 (columns 3-7, n=5) and WNPRC (columns 8-9, n=2) to infer recent ancestry of pedigreed animals
504 in the USA. Each bar shows the inferred ancestry proportions of each animal specifying either Fig
505 1A, 2 ancestral populations (K=2, top panel and best cross-validation score) or Fig 1B, 3 (K=3,
506 middle panel). Cross validation for K>3 ancestral populations failed. In Fig 1C, ancestry was
507 independently inferred in a supervised analysis using the population labels shown below the first
508 9 bars and corresponding to each NPRC. Each analysis showed a subdivided population, where

509 most animals have a major ancestry component from a single population. Fig 1D shows the known
510 pedigree with each animal included in the ADMIXTURE analysis shaded to represent the
511 proportion of NEPRC, SNPRC, and WNPRC ancestry. Dashed lines show where an individual is
512 placed multiple times in the figure. Given the complex nature of large pedigrees this is an artifact
513 of plotting data.

514

515



516
517

518 **Figure 2. Inference of relatedness from GBS data.** (A) a comparison of empirical kinship (x
519 axis) and relatedness inferred from GBS ($\hat{\pi}$, y axis). The grey dashed line shows the expectation
520 from a perfect correlation between approaches, the black dashed line shows the observed
521 relationship between approaches. Highlighted in the red dots are comparisons between 12

522 individuals which are discordant between approaches. (B) Pedigree of 11 of the 12 individuals
523 uncovered using GBS (C) A heatmap showing the pedigree (upper triangle) and GBS (lower
524 triangle) based estimates of relatedness. Darker colors denote a higher degree of relatedness
525 between individuals. Highlighted by a box in the top right corner are the individuals shown in red
526 in (A). The large difference in estimates of relatedness between the GBS and pedigree based
527 estimates is driven by incomplete pedigree information from a single individual subject to
528 migration between colonies. (D) Inferred ancestry components for each individual, darker colors
529 represent the proportion of ancestry from Fig 1C.

530
531
532
533