# Maximizing the acquisition of unique reads in non-invasive capture sequencing experiments

Running title: Increasing coverage of captured fecal DNA

Claudia Fontsere[1], Marina Alvarez-Estape[1], Jack Lester[2], Mimi Arandjelovic[2], Martin Kuhlwilm[1], Paula Dieguez[2], Anthony Agbor[2], Samuel Angedakin[2], Emmanuel Ayuk Ayimisin[2], Mattia Bessone[2], Gregory Brazzola[2], Tobias Deschner[2], Manasseh Eno-Nku[3], Anne-Céline Granjon[2], Josephine Head[2], Parag Kadam[4], Ammie K. Kalan[2], Mohamed Kambi[2], Kevin Langergraber[5,6], Juan Lapuente[2,7], Giovanna Maretti[2], Lucy Jayne Ormsby[2], Alex Piel[4], Martha Robbins[2], Fiona Stewart[4], Virginie Vergnes[8], Roman M. Wittig[2,9], Hjalmar S. Kühl[2,10], Tomas Marques-Bonet[1,11,12,13] †, David A. Hughes[14,15] * and Esther Lizano[1,13] † *

[1] Institut de Biologia Evolutiva, (CSIC-Universitat Pompeu Fabra), PRBB, Doctor Aiguader 88, Barcelona, 08003, Spain.

[2] Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany.

[3] WWF Cameroon Country Programme Office, BP6776; Yaoundé, Cameroon.

[4] School of Biological and Environmental Sciences, Liverpool John Moores University, James Parsons Building, Byrom street, Liverpool, L3 3AF, UK.

[5] School of Human Evolution and Social Change, Arizona State University, 900 Cady Mall, Tempe, AZ 85287 Arizona State University, PO Box 872402, Tempe, AZ 85287-2402 USA.

[6] Institute of Human Origins, Arizona State University, 900 Cady Mall, Tempe, AZ 85287 Arizona State University, PO Box 872402, Tempe, AZ 85287-2402 USA.

[7] Comoé Chimpanzee Conservation Project, Kakpin, Comoé National Park, Ivory Coast.

[8] Wild Chimpanzee Foundation (WCF) 23BP238 Abidjan, Côte d'Ivoire 23.

28    [9] Taï Chimpanzee Project, Centre Suisse de Recherches Scientifiques, BP 1301, Abidjan 01, CI, Côte d'Ivoire.

30    [10] German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Deutscher Platz 5e, 04103 Leipzig.

32    [11] CNAG-CRG, Centre for Genomic Regulation (CRG), Barcelona Institute of Science and Technology (BIST), Baldiri i Reixac 4, 08028 Barcelona, Spain.

34    [12] Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Catalonia 08010, Spain.

       [13] Institut Català de Paleontologia Miquel Crusafont, Universitat Autònoma de Barcelona, Columnes s/n,
36    08193 Cerdanyola del Vallès, Spain.

       [14] MRC Integrative Epidemiology Unit at University of Bristol, Bristol, BS8 2BN, UK.

38    [15] Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, BS8 2BN, UK.

40    * Esther Lizano and David A. Hughes should be considered joint senior author.

       [†]**Corresponding author:** Esther Lizano and Tomas Marques-Bonet

42

# Abstract

44    Non-invasive samples as a source of DNA are gaining interest in genomic studies of endangered species. However, their complex nature and low endogenous DNA

46    content hamper the recovery of good quality data. Target capture has become a productive method to enrich the endogenous fraction of non-invasive samples, such

48    as feces, but its sensitivity has not yet been extensively studied. Coping with fecal samples with an endogenous DNA content below 1% is a common problem when prior

50    selection of samples from a large collection is not possible. However, samples classified as unfavorable for target capture sequencing might be the only

52    representatives of unique specific geographical locations or to answer the question of interest.

54    To explore how library complexity may be increased without repeating DNA extractions and generating new libraries, here we have captured the exome of 60 chimpanzees

56    (*Pan troglodytes*) using fecal samples with very low proportions of endogenous content (< 1%).

58    Our results indicate that by performing additional hybridizations of the same libraries, the molecular complexity can be maintained to achieve higher coverage. Also,

60    whenever possible the starting DNA material for capture should be increased. Lastly, we have specifically calculated the sequencing effort needed to avoid exhausting the

62    library complexity of enriched fecal samples with low endogenous DNA content.

This study provides guidelines, schemes and tools for laboratories facing the

64    challenges of working with non-invasive samples containing extremely low amounts of endogenous DNA.

66 **Keywords**: Non-invasive samples, fecal samples, target capture, molecular complexity, conservation genomics, chimpanzees.

68 ## Introduction

Studies of wild populations that are unamenable to invasive sampling (eg: trapping or
70 darting) often rely on the usage of low quality and/or quantity DNA samples (Schwartz, Luikart, & Waples, 2007; Vigilant & Guschanski, 2009), traditionally restricting the
72 analysis to neutral markers or genetic loci such as microsatellites (Arandjelovic et al., 2011; Inoue et al., 2013; Mengüllüoğlu, Fickel, Hofer, & Förster, 2019; Orkin, Yang,
74 Yang, Yu, & Jiang, 2016), autosomal regions (Fischer, Wiebe, Pääbo, & Przeworski, 2004) and the mitochondrial genome (Fickel, Lieckfeldt, Ratanakorn, & Pitra, 2007;
76 Thalmann, Hebler, Poinar, Pääbo, & Vigilant, 2004). Depending on the researcher's question, these neutral genetic markers may continue to be the most economical and
78 efficient method (Shafer et al., 2015). However, for other questions such as cataloging genetic diversity, assessing kinship, making fine inferences of demographic history, or
80 evaluating disease susceptibility it becomes increasingly relevant to acquire a more representative view of the genome (Ouborg, Pertoldi, Loeschcke, Bijlsma, & Hedrick,
82 2010; Primmer, 2009; Shafer et al., 2015; Städele & Vigilant, 2016; Steiner, Putnam, Hoeck, & Ryder, 2013).

84 Conservation genomics of ecologically-crucial, non-model organisms, and especially threatened species such as great apes, have largely benefited from the current
86 advances in next-generation sequencing (NGS) technologies (Gordon et al., 2016; Locke et al., 2011; Mikkelsen et al., 2005; Scally et al., 2012). The ability to
88 simultaneously interrogate hundreds of thousands of genetic markers across an entire

4

genome allows greater resolution on inferences of demographic parameters, genetic

90 variation, gene flow, inbreeding, natural selection, local adaptation and the

evolutionary history of the studied species (De Manuel et al., 2016; Prado-Martinez et

92 al., 2013; Xue et al., 2015).

The major impediment to the study of wild, threatened, natural populations continues

94 to be the difficulties in acquiring samples of known location from a large number of

individuals. To avoid disturbing and negatively influencing endangered species

96 (alteration of social group dynamics, infections and stress) (Morin, Wallis, Moore,

Chakraborty, & Woodruff, 1993; Taberlet, Luikart, & Waits, 1999), but also to track

98 cryptic or monitor reintroduced species (De Barba et al., 2010; Ferreira et al., 2018;

Reiners, Encarnação, & Wolters, 2011; Stenglein, Waits, Ausband, Zager, & Mack,

100 2010), sampling often relies on non-invasive (NI) sources of DNA such as feces and

hair, rather than invasive samples such as blood or other tissues, which yield better

102 DNA quality and quantity.

NI samples have a complex nature: they are typically composed of low proportions of

104 host or endogenous DNA (eDNA), are highly degraded (Perry, Marioni, Melsted, &

Gilad, 2010; Taberlet et al., 1999), and contain genetic material from the host's

106 microbiota and from species living in the environment where the sample was collected

(i.e., exogenous DNA) (Hicks et al., 2018). The proportion of endogenous versus

108 exogenous DNA can be highly variable (Hernandez-Rodriguez et al., 2018) and as

previous literature has proposed, may depend on the environmental conditions, with

110 humidity and ambient temperature having the highest influence (Goossens, Chikhi,

Utami, De Ruiter, & Bruford, 2000; Harestad & Bunnell, 1987; King, Schoenecker, Fike,

112 & Oyler-McCance, 2018; Nsubuga et al., 2004). Because of this, the employment of

techniques that generate sequences of the whole genomic content of the samples,

114  such as NGS, has not been economically feasible until recently. Target enrichment

technologies, also known as capture, have become a common and successful

116  methodology in ancient DNA studies (Carpenter et al., 2013) and have allowed for a

more cost-effective use of NGS on NI samples, as the endogenous to exogenous DNA

118  ratio greatly improves, thus reducing the sequencing effort (Perry et al., 2010; Snyder-

Mackler et al., 2016; van der Valk, Lona Durazo, Dalén, & Guschanski, 2017). Capture

120  methods reduce the relative cost of sequencing and improve the quality of the data by

building DNA libraries that are hybridized to complementary baits for selected target

122  regions (partial genomic regions, a chromosome(s), the exome, or the whole genome)

increasing the proportion of the targeted eDNA to be sequenced.

124  Despite the existence of technical studies describing the use of NI samples for the

study of wild chimpanzees (*Pan troglodytes*) (Hernandez-Rodriguez et al., 2018; White

126  et al., 2019) many aspects remain to be investigated. For instance, in Hernandez-

Rodriguez et al., samples were selected to cover the entire range of observed average

128  fragmentation lengths and percentage of eDNA, in order to be as representative as

possible. As a result, they observed a sequencing bias due to the different percentage

130  of endogenous content in captured samples. To avoid that outcome, they proposed

performing equi-endogenous pools instead of the standard pooling of libraries

132  according to molarity. White et al. followed this recommendation and yielded a more

balanced representation across samples. However, their experiments were limited to

134  only those samples with a proportion of eDNA above 2% (White et al., 2019). As shown

by Hernandez-Rodriguez et al. there is a positive association between endogenous

136  content and the amount of data acquired from a sample, such that when possible, one

should use those samples with higher endogenous content. However, the proportion

138     of chimpanzee fecal samples with eDNA above 2% is often very low (<20%) (White et al., 2019).

140     The NI chimpanzee samples used in this study were collected from 15 different geographic sites across the whole species' ecological habitat in Africa and included all

142     four subspecies, thus representing a wide variety of sampling and environmental conditions. With this screening approach we were able to examine how the proportion

144     of eDNA content varies between each site, revealing that the majority of collected samples in some sites have low proportions of eDNA (<1%). Therefore, when prior

146     selection of samples from a large collection is not possible, the only ones to represent a specific location or relevant to the scientific question might be those with extremely

148     low proportions of endogenous content. Because of that, we have focused our efforts on developing approaches to retrieve the maximum data possible from challenging

150     samples.

In that regard, we sought to capture the exome of 60 chimpanzee fecal samples as

152     part of the Pan African Programme: The Cultured Chimpanzee (PanAf) (http://panafrican.eva.mpg.de/) (Kühl et al., 2019) with eDNA estimates below 1%. We

154     used a commercial human exome to evaluate how the coverage of targeted genomic regions may be increased in a collection of samples that may be regarded as

156     unfavorable for target capture sequencing. We confirmed the importance of the correct estimation of eDNA and the pooling of libraries accordingly to avoid sequencing bias

158     across samples (Hernandez-Rodriguez et al., 2018). We also expanded on previously explored and unexplored guidelines to ensure the maintenance of the captured

160     molecule diversity or library complexity such as the number of libraries in a pool, the

performance of additional hybridizations and increasing the total DNA starting material

162   for capture (Hernandez-Rodriguez et al., 2018; Perry et al., 2010; Snyder-Mackler et al., 2016; White et al., 2019).

164   Our results provide the most comprehensive exploration to date of target enrichment efficiency in very low eDNA fecal samples, and guidelines to improve the quality of the

166   data without re-extracting DNA and preparing new libraries. These findings could greatly benefit the conservation effort on great apes, as well as any other species with

168   similar DNA sampling limitations.

# Material and Methods

170   ## Samples and Library Preparation

Chimpanzee fecal samples from 15 different sites in Africa were collected as part of

172   the PanAf (Figure 1A). Approximately 5g ("hazelnut-size") of feces were collected from each chimpanzee fecal sample and stored in the field using a two-step ethanol-silica

174   preservation method (Nsubuga et al., 2004). Depending on the density of the sample, between 10 and 80 mg of dry fecal sample were extracted using a Qiagen robot with

176   the QIAamp Fast DNA Stool Mini Kit (Qiagen) with modifications (Lester et al, in review). The extractions were screened using a microsatellite genotyping assay

178   (Arandjelovic et al., 2009; Arandjelovic et al., 2011) and up to 20 samples from each PanAf field site were selected as follows (1) those that amplified at the most loci of the

180   15 loci panel, (2) represented unique individuals, and (3) were ascertained to be non-first degree relatives (Csilléry et al., 2006) (302 samples) (Supporting Information

182   Table S1). To ensure sufficient template DNA for library preparation, the 302 samples were re-extracted using the same QIAamp kit and between 100 and 200 mg of dry

184    fecal sample. Total DNA concentration and fragmentation were measured on a

Fragment Analyzer using a Genomic DNA 50Kb Analysis kit (Advanced Analytical) and

186    the fragmentation level was calculated with PROSize Data Analysis Software (Agilent

Technologies). Endogenous DNA content (fraction of mammalian DNA, relative to gut

188    microbial and other environmental genetic material) was estimated by qPCR (Morin,

Chambers, Boesch, & Vigilant, 2001). Finally, percentage of endogenous content for

190    each sample was calculated by dividing the chimpanzee eDNA concentration by the

total DNA concentration. We selected 60 samples with an intermediate percentage of

192    eDNA (0.41-0.85%, average 0.61%) from the 302 screened samples (range of

endogenous distribution: 0-47.57%, average 1.49%) (Supporting Information S1 and

194    Table S2).

A single library was prepared for each of the 60 samples following the BEST protocol

196    (Carøe et al., 2018) starting with 200 ng total DNA (from a sample) with minor

modifications. Specifically, double in-line barcoded adapters were used, barcoding

198    each sample at both ends of its library to allow for its unique identification within a pool

(Rohland & Reich, 2012). Library concentration was calculated using Agilent 2100

200    BioAnalyzer and DNA7500 assay kit. A detailed protocol for library construction can

be found in Supplementary Information.

202

Pooling and Capture

204    Endogenous DNA content is a key factor in target-capture experiments directly

influencing the yield of on-target reads and molecule diversity (Hernandez-Rodriguez

206    et al., 2018). Our equi-endogenous sample pooling strategy follows two criteria. First,

samples belonging to a pool have similar eDNA proportions according to a 1:2 ratio

rule: the sample with highest proportion of eDNA cannot double the sample with the lowest. Second, each sample within a pool contributes the same total amount of eDNA (μg) to the final pool, creating an equi-endogenous pool. So, the sample with the lowest percentage of eDNA will contribute more total DNA to the final pool compared to the sample with the highest, but the amount of eDNA per sample will be equivalent.

According to the estimates of eDNA, we pooled the 60 libraries into three primary pools (see graphical representation in Figure 2). The first pool (P1) with 2 μg total DNA (in the pool) consisted of 10 samples with an average endogenous content of 0.81% (range 0.69-0.85%). The second pool (P2) had 4 μg total DNA and consisted of 20 samples and an average endogenous content of 0.69% (range 0.58-0.80%). The 30 remaining libraries were pooled into the third pool (P3) of 6 μg total DNA with an average endogenous content of 0.49% (range 0.41-0.66%) (Table 1 and Figure 3A, Supporting Information Table S2). Subsequently, each initial primary pool was subdivided into two (P1E1, P1E2), four (P2E1, P2E2, P2E3, P2E4) and six (P3E1, P3E2, P3E3, P3E4, P3E5, P3E6) exome capture (E) replicates each consisting of 1 μg of total DNA.

Independently, we repeated the construction of the primary pools (P1, P2 and P3), but with each having 4 μg total DNA. Each of these new primary pools was then divided into two replicates of 2 μg each (P1E3, P1E4, P2E5, P2E6, P3E7, P3E8). As a consequence of generating replicate primary pools, six of the 60 libraries were exhausted and are not present in these replicate primary pools. As a result, across all 60 samples and 18 hybridizations there are a total of 388 individual hybridization experiments (Figure 2). All details are provided in Table 1.

Each exome capture experiment consisted of two consecutive hybridizations, or dual-capture reactions as previously recommended (Hernandez-Rodriguez et al., 2018) using the SureSelect Human All Exon V6 RNA library baits from Agilent Technologies and was performed following the manufacturer's protocol with some modifications (full protocol is available in Supporting Information), and started with either 1 µg or 2 µg total DNA (Table 1 and Figure 2). After the first hybridization reaction and the subsequent PCR enrichment, we performed the second hybridization reaction with all available material. The final captured pool was amplified with indexed primers (Kircher, Sawyer, & Meyer, 2012), double-indexing each library within a pool, thereby tagging each library to a specific hybridization experiment. Double inline barcoded (sample specific) and double indexed (pool specific) libraries allow for multiplexing many libraries into a single pool and sequencing many pools into a single sequencing lane, even when the same sample library is present in multiple hybridization reactions. This permits the tracking of unique experiments.

For the reminder of the article when we use the word "capture" or "hybridization", we will always be referring to the dual-capture or two consecutive rounds of capture hybridizations that are described above.

## Sequencing and Mapping

Captured libraries were pooled into 3 sequencing batches and sequenced on a total of 3.75 lanes of a HiSeq 4000 with 2x100 paired-end reads: SeqBatch1 (P1E1, P2E1, P2E2, P3E1, P3E2, P3E3), SeqBatch2 (P1E2, P2E3, P2E4, P3E4, P3E5, P3E6) and SeqBatch3 (P1E3, P1E4, P2E5, P2E6, P3E6, P3E7, P3E8) (Table 1).

Demultiplexed FASTQ files were trimmed with Trimmomatic (version 0.36) (Bolger, Lohse, & Usadel, 2014) to remove the first 7 nucleotides corresponding to the in-line barcode (HEADCROP: 7), the Illumina adapters (ILLUMINACLIP:2:30:10), and bases with an average quality less than 20 (SLIDINGWINDOW:5:20). Paired-end reads were aligned to human genome Hg19 (GRCh37, Feb.2009 (GCA_000001405.1)) using BWA (version 0.7.12) (Li & Durbin, 2009). Duplicates were removed using PicardTools (version 1.95) (http://broadinstitute.github.io/picard/) with MarkDuplicates option. Further filtering of the reads was carried out to discard secondary alignments and reads with mapping quality lower than 30 using samtools (version 1.5) (Li et al., 2009). From now on, we will refer to those reads remaining after filtering as "reliable reads". To retrieve the reliable reads on-target we used intersectBed from BEDTOOLS package (version 2.22.1) (Quinlan & Hall, 2010) using exome target regions provided by Agilent. In cases where we combined sequencing data, we merged filtered bam files from different hybridizations using MergeSamFiles option from PicardTools (version 1.95) (http://broadinstitute.github.io/picard/). Since the merged bam files can still contain duplicates generated during library preparation, we removed duplicates and then retrieved the reliable reads on-target using the same methodology as above.

For all previous steps, the total number of reads were counted using PicardTools (version 1.95) (http://broadinstitute.github.io/picard/) with CollectAlignmentSummaryMetrics option. The percentage of human contamination was estimated by using positions where modern humans and chimpanzees consistently differ. We used previously published diversity data on high-coverage genomes from the *Pan* species (chimpanzee and bonobos) (De Manuel et al., 2016) and human diversity data from the 1000 Genomes Project (Auton et al., 2015),

selecting positions where the human allele is observed at more than 98% frequency,

278 and a different allele is observed in almost all *Pan* individuals (136 out of 138

chromosomes). Genome-wide, 5,646,707 chimpanzee-specific positions were

280 identified. Using samtools mpileup (Li et al., 2009), we retrieved the number of

observations of human-like alleles at these positions in the mapped reads, and

282 estimated the human contamination as the fraction of observations for the human-like allele

across all positions.

284

## Capture performance

286 Capture performance was evaluated by calculating the enrichment factor (EF), capture

specificity (CSp), library complexity (LC), and capture sensitivity (CS) as described in

288 Hernandez-Rodriguez *et al* (2018). EF is calculated as the ratio of the number of

reliable reads on-target to the total reads sequenced divided by the fraction of the

290 target space (64Mb) to the genome size (~3Gb). CSp is defined as the ratio of reliable

on-target reads to the total number of reliable reads. LC is defined as the number of

292 reliable reads divided by the total number of mapped reads (containing duplicated

reads). Capture sensitivity (CS) is defined as the number of target regions with an

294 average coverage of at least one (DP1) - but also four (DP4), ten (DP10), twenty

(DP20) or fifty (DP50) - divided by the total number of target regions provided by the

296 manufacturer (n = 243,190). To calculate the average coverage of the target regions

we used samtools (version 1.5) with the option bedcov (Li et al., 2009).

298 To generate molecular complexity or library complexity curves (MC), we used the

subsampling without replacement strategy implemented in Preseq software (version

300 2.0.7) with c_curve option (http://smithlabresearch.org/software/preseq/) from the bam

files without removing duplicates. MCs were sequentially estimated by adding the

302     production reads, i.e. raw reads produced by sequencing, from additional

hybridizations, one at a time until all hybridizations from the same library were merged

304     (schematic representation in Figure S2).

Correlation coefficients among all pairs of study variables were estimated. Spearman's

306     rho (cor.test(, method = "sp") from R stats package) was estimated when comparing

two numeric variables. Among two categorical variables we estimated Cramér's V,

308     derived from a chi-squared test (chisq.test() from R stats package). When comparing

a numeric and categorical variable we took the square root of the R-squared statistic

310     derived from a univariate linear model (lm() from R stats package) with a rank normal

transformation (rntransform() modified from the GenABEL package to randomly split

312     tied values) on the dependent, numerical values. In addition, univariate and

multivariate type I hierarchical analysis of variances (ANOVA; anova() from R stats

314     package) were performed to estimate the variance explained (or eta-squared) each

experimental variable has on performance summary statistics (number of unique

316     reads, reliable reads, EF, LC, CS and CSp). We down-sampled libraries to 1,500,000

reads (n=274) to remove production reads as a confounding factor. Each performance

318     statistic was rank normal transformed with ties being randomly split to ensure normality

of the dependent variable. Univariate analysis focused on the effect that subspecies,

320     geographic sampling site, total DNA concentration, endogenous DNA concentration,

percent endogenous DNA, average fragment length, pool, amount of DNA in a

322     hybridization, hybridization and sequencing batch had on each performance statistic.

A multivariate model was built to conform with experimental (hierarchical) order, such

324     that each dependent variable (performance summary statistic, CS at DP1) was

explained by ~ subspecies + site + % eDNA + average fragment size + pool + amount

326 of DNA + hybridization + sequencing batch + error. Again, the variance explained by

each independent variable was summarized by computing the eta-square statistic

328 derived from the sums of squares for each variable using a type I hierarchical ANOVA.

All statistical analyses were performed in R (version 3.5.2) (R Core Team, 2018).

330

## Results

332 ### Sample Description

Samples were collected from 15 different PanAf sites distributed across the entire

334 range of chimpanzees in Africa (Figure 1A and Supporting Information Table S1). The

302 screened samples had an average eDNA of 1.49%, ranging from 0 to 47.75%

336 (Figure 1B, Supporting Information Figure S1A and Table S1) with 70.2% of the

samples below 1% eDNA, according to qPCR estimates (Figure 1C). The average

338 fragment length for screened samples was 3,479.94 bp (ranging from 72 to 17,966 bp)

(Supporting Information Figure S1B and Table S1).

340 We observe variation on the average endogenous content among geographical sites

(Figure 1B), and also variation on fragment length among geographical sites

342 (Supporting Information Figure S1B). For instance, samples collected in a specific

location such as Campo Ma'an (Cameroon) have an average eDNA of 0.02%, an

344 extremely low value compared to the average of all sites of 1.49%. On the other hand,

some sites such as Ngogo (Uganda) have samples with higher than average eDNA

346 (6.95%) (Supporting Information Table S3). This might be explained by the influence

of weather, humidity and temperature on DNA preservation and bacterial growth in the

348     fecal sample before collection as well as a product of sample age and quality of

sampling conditions (Brinkman, Schwartz, Person, Pilgrim, & Hundertmark, 2010;

350     Goossens et al., 2000; Harestad & Bunnell, 1987; King et al., 2018; Nsubuga et al.,

2004; Wedrowicz, Karsa, Mosse, & Hogan, 2013).

352     A total of 60 samples with a mean percent endogenous content of 0.58% and range

from 0.41% to 0.85%, and with a median human contamination of 0.0875% from all

354     four chimpanzee subspecies and 14 geographic sites were carried forward into target

capture enrichment experiments (Table S2). After double-inline-barcoded library

356     production, the 60 samples were placed into 3 pools with 10, 20 and 30 samples each.

Samples were divided into pools based on their percent endogenous content, such

358     that those samples with higher levels of percent endogenous content were in P1 with

10 samples (mean = 0.81) and those with the smallest were in P3 with 30 samples

360     (mean = 0.49; P2 mean = 0.69) (Figure 3A). As such the percent endogenous DNA is

highly structured among the three pools, explaining 81% of the variation in eDNA

362     (univariate linear model using rank normal transformed % eDNA; p-value = $2.05 \times 10^{-91}$) (Supporting Information Figure S4A).

364 ## Read Summary Statistics and Capture Performance

As illustrated in Figure 3B across a total of 18 hybridization experiments sequenced

366     we obtained ~1.40 billion reads distributed among 3 pools. Of those, ~1.19 billion were

mapped reads (85.19%), with ~203 million reads being considered duplicate-free,

368     reliable reads (14.6%). After removing off-target reads, we obtained a total of ~174

million on-target-reliable reads (12.48%) (Supporting Information Table S4, Supporting

370     Information Figure S3A). However, on average each hybridization experiment yielded

an average of 17.35% on-target-reliable reads, with a range of 4.15% in our earliest

372 experiments to 34.85% in our later experiments (Supporting information Table S5). The observed high levels of duplicates are a consequence of the low endogenous

374 content of the samples and the exhaustion of library complexity during sequencing; we will elaborate on outcome and improvements below.

376 The ~1.40 billion reads were not equally distributed among the 3 pools (production reads explained by pools; $r^2 = 0.41$, p-value = $3.24\text{x}10^{-16}$) or 18 hybridizations ($r^2 = $

378 0.62, p-value = $2.59\text{x}10^{-30}$). In fact, two hybridizations of P1 (P1E1, P1E2) were sequenced to an average depth of 18 million reads, while all other hybridizations had

380 an average depth of 3 million reads (Figure 3C). This very deep sequencing, in P1E1 and P1E2, led to a point where the library complexity was exhausted, leading to the

382 sequencing of a high number of PCR duplicates (Supporting Information Figure S3A, S3B and Supporting Information Table S5). We therefore reduced subsequent

384 sequencing efforts, as discussed in section "Optimization of required production reads", for the remaining replicate hybridizations.

386 All capture performance summary statistics (Supporting Information Table S4), to the exception of capture specificity (CSp), are strongly correlated with the number of

388 production reads acquired (median correlation coefficient = 0.422, CI = 0.03 to 0.93; Supporting information Figure S4A, Table S6). Given this, and also because of the

390 distinct difference in the number of production reads between P1E1 and P1E2 and all other hybridizations we down-sampled all experiments to 1.5 million production reads,

392 retaining only those 274 sample/hybridization experiments with 1.5 million production reads, and re-estimated all capture performance summary statistics (Supporting

394 Information Figure S4B, Table S7 and S8). The effect each experimental variable has

on performance was estimated in a univariate linear model after rank normal transforming each summary statistic (Figure 4A). We observed a near uniformity in the variance explained by each experimental variable across each performance statistics. In short, the average, ranked order of variance explained by each explanatory variable are sample (86.50%), hybridization (38.72%), sequencing batch (28.78%), site (20.5%), pool (13%), % endogenous DNA (11%), subspecies (8.85%), starting DNA amount (7.35%), endogenous DNA concentration (5.14%), average fragmentation size (2.12%,), and total DNA concentration (2.07%). Given these observations we may conclude that variation in hybridization and sequencing are crucial to performance. However, sample quality and starting material varies among our hybridizations and sequencing batches. These tendencies can be observed in Figure 5A-C. We account for this in a multivariate linear model followed by a decomposition of the variance in a type I hierarchical analysis of variance (ANOVA). To do so we fit a linear model ordered by experimental choices, as described in materials and methods, to explain Capture Sensitivity (CS) at DP1 which is being used here as an example of capture performance. This model indicates that hybridization explains, on average, an attenuated 17.80% of the variation in performance, followed by percent endogenous content (17.11%), site (9.62%), subspecies (9.26%), pool (3.92%) and then the amount of DNA in the hybridization (3.58 %) (Figure 4B). Results for all other performance summary statistics mirror those for CS at DP1 and can be seen in Figure S5.

## Relevance of Equi-Endogenous Pools

The observations of Hernandez-Rodriguez et al. and White et al. suggest that pooling libraries by eDNA concentration (in equi-endogenous pools) prior to hybridization

418    capture should reduce or remove the effect of variation in eDNA across samples on

targeted capture sequencing performance. Indeed, eDNA did not have a major

420    influence on production reads or on-target reads, although a slightly positive trend can

be observed in some hybridizations of P2 (Supporting Information Figure S6). Without

422    equi-endogenous pooling, it is expected that samples with higher eDNA would

accumulate more on-target reads than other samples with lower eDNA as observed by

424    Hernandez-Rodriguez et al. The reason why in P2 we find some outliers might be

traced to both pipetting variations and inaccurate endogenous measurements from

426    qPCR values due to the presence of inhibitors (Morin et al., 2001). Avoiding outliers is

extremely important in limiting variability within a pool. For example, sample N183-5

428    accumulated 29.4% of total raw reads in P2, when a value 5% (1/20 of 100%) was

expected (Supporting Information Figure S7).

430    ## Impact of Amount of Starting DNA for Capture on Library Complexity

One major decision when performing capture experiments is the amount of starting

432    DNA in the pool. In twelve hybridizations we used the manufacturer's suggested

amount of starting material, 1 µg for each pool. For the last two hybridizations of each

434    pool (a total of six hybridizations) we doubled the starting material, up to 2 µg of pooled

libraries (Table 1). With this approach we aimed to test the effect on the final LC when

436    doubling the amount of DNA and to determine how much DNA should be used for fecal

capture experiments. We observed an average increase of 2.8-fold in LC for

438    experiments using 2 µg of total DNA in the hybridization relative to those using 1 µg

(Supporting Information Figure S3B). However, given that production reads also vary

440    between these two conditions, we down-sampled the data to 1,500,000 reads per

library. After this correction we still observed 2-fold higher LC when starting the

442 experiments with 2 µg of total DNA in all pools (Figure 5D).

Molecular complexity, as influenced by the amount total DNA in a hybridization, was

444 further investigated by evaluating the relationship between MC and production reads

in a MC curve analysis. The MC curve for each hybridization was obtained by

446 subsampling without replacement their reads. The results supported the conclusion

above: increasing the amount of total DNA in the hybridization increased the MC

448 (Supporting Information Figure S8). Therefore, whenever there is sufficient library

available, it is advisable to start with 2 µg rather than 1 µg.

450 ## Molecular Complexity and Capture Sensitivity

One of the critical aspects to increase coverage is to acquire as many unique on-target

452 reads as possible without exhausting the library's molecular complexity. We applied a

subsampling without replacement method to assess how many mapped reads are

454 unique after incrementally adding production reads from replicate hybridizations. In

principle, molecular complexity curves that plateau quickly are derived from low

456 complexity libraries, and conversely high complexity libraries may not reach plateau.

Thereby the plateau indicates when there are no new unique reads to be sampled or

458 sequenced (see Supporting Information Figure S2 for a schematic representation).

We performed the analysis of molecular complexity in libraries belonging to P3 since

460 more hybridization replicates were available (8 in total) for 30 libraries. We found that

for the majority of the libraries, performing additional hybridizations increased the

462 number of unique reads retrieved (Supporting Information Figure S9, example library

N259-5). However, there were libraries that quickly hit exhaustion where performing

additional hybridizations would add little extra information (Supporting Information Figure S9, example library Kay2-32). Overall, by performing additional hybridizations, it was possible to retrieve new unique reads and thus increase the final coverage (Figure 6A), because libraries themselves were not exhausted but merely their hybridization-captured molecules reached exhaustion.

Following the same strategy, we calculated the sensitivity in P1, P2 and P3 (4, 6 and 8 replicates respectively). After cumulatively adding data from replicate hybridizations we covered 85.57% in P1 (95% CI: 74.78-96.36%), 76.23% in P2 (95% CI: 64.55-87.91%) and 79.83% in P3 (95% CI: 74.44-85.22%) on average of the target space, with at least 1 read (Supporting Information Figure S10). Interestingly, no sample covered 100% of target space. Looking carefully into this, we observed that precisely the same 3,804 regions (1.54%) were never covered in any replicate hybridizations, suggesting that some regions are either difficult to capture (Kong, Lee, Liu, Hirschhorn, & Mandl, 2018) or are too divergent between *Homo* and *Pan* to either capture or map these particular sequences (Supporting Information Figure S11).

For deeper coverage of at least 4 or 10 reads, we still observed a positive progression, with each additional hybridization increasing coverage, indicating that additional hybridizations would result in an increase of the proportion of the genome covered at these depths as well (Supporting Information Figure S10).

## Optimization of Required Production Reads

Assessing the amount of sequencing needed is one of the major decisions when planning an experiment. As a result of the low eDNA content of most fecal samples, derived libraries can easily reach saturation (i.e., high levels of duplicated reads).

Therefore, sequencing depth should be carefully calculated. Without previous

488    knowledge, we sequenced the first 2 hybridizations for P1, the first 4 hybridizations for

P2, and the first 6 hybridizations for P3 in three lanes of a HiSeq 4000. For P1 only

490    ~6% and for P2 and P3 only ~13% of production reads were unique reads (Supporting

Information Table S5), indicative of high levels of PCR duplicates due to library

492    exhaustion. To avoid over-sequencing in our next experiments, we set an arbitrary

threshold to recover approximately 20% of the "informative" data (unique reads)

494    available in a hybridization experiment. Using the data from SeqBatch 1 and 2, we

estimated that on average, for samples with less than 1% eDNA, we would sequence

496    at most 2 million mapped reads per library (Figure S12). Given that 80% of reads

mapped to the genome in these experiments, we estimated that we would need to

498    sequence at most 2.5 million production reads per library (Supporting Information

Table S5).

500    To test these estimates, we sequenced the remaining hybridizations (P1E3, P1E4,

P2E5, P2E6, P3E7, P3E8) in three-fourths of a HiSeq 4000 lane. The number of

502    average production reads obtained were 3.5, 2.0 and 1.5 million for libraries in

hybridizations from P1, P2, and P3, respectively. On average ~38% (range: 8.09-

504    50.81%) of reads were unique reads in all pools (Supporting Information Figure S13).

We note that these values exceeded what we observed in the previous hybridization

506    experiments. An outcome we attribute to the increase in starting material (2 µg), also

used in these experiments, as noted above.

## Pooling Strategy

Choosing how many samples to pool is a difficult decision, since little is known on how the pool size will affect the final molecular complexity. Taking advantage of our pooling strategy (Figure 2), we assessed the effect of size on the average library complexity for all samples within each hybridization with a subsampling without replacement strategy.

When only a single hybridization was performed, a single library within a pool of 10, 20 or 30 would, on average, result in a similar number of unique molecules (Figure 6B, Supporting Information Figure S14). However, there is a tendency for samples in smaller pools (P1) to perform better than those in larger pools. This could be explained by our experimental design, where samples with higher eDNA content are in smaller pools. However, let us address this possibility here. Using CS as an example summary statistic, we observed that CS is higher for pools with smaller numbers of samples in them (Figure 5C). Given median estimates, a pool of 10 libraries (median CS = 0.46) had 1.44-fold higher CS than a pool of 20 libraries (median CS = 0.32), and 1.92-fold higher than a pool of 30 libraries (median CS = 0.24). Between a pool of 20 and a pool of 30, the ratio was 1.33-fold (Figure 5C and Supporting Information Figure S15). If we remove the effect of having a variable number of production reads across experiments by down-sampling, this observation still remains (Supporting Information Figure S16). That is, smaller pools do have higher CS estimates, and pools linearly account for 18% of the variation in CS (univariate ANOVA, p-value=$3.47 \times 10^{-12}$ (Figure 4A)). Finally, if we correct for all experimental variables with a multivariate analysis, as done above, we show that 'Pool' only accounts for 4% of the variation in CS (Figure 4B), but the effect of pool size remains significant (multivariate ANOVA, p-value = $2.7 \times 10^{-4}$;

532 Supporting Information Figure S16). However, this effect on CS attenuates with

additional hybridizations (4, 6 and 8, for P1, P2 and P3 respectively) for the same pool

534 (Supporting Information Figure S17). Moreover, a similar outcome can be observed

when comparing the effect of pool size on LC. After sequentially adding data from

536 replicate hybridizations in each pool (see Supporting Information Figure S2 for a

schematic representation), we can acquire the same number of unique reliable reads

538 (Figure 6C, Supporting Information S16).


540 # Discussion


Capturing host DNA from fecal samples is a challenging endeavor. Previous work has

542 shown that the retrieval of genomic data from fecal samples by target enrichment

methodologies is a feasible and powerful tool for conservation and evolutionary studies

544 (Perry, 2014; Snyder-Mackler et al., 2016). However, obtaining good quality and

quantity DNA from fecal samples is not always possible. Because of that, many studies

546 have characterized the technical difficulties of capturing DNA from non-invasive

samples and proposed different strategies (Hernandez-Rodriguez et al., 2018; van der

548 Valk et al., 2017; White et al., 2019). Van der Valk et al. (2017) captured the whole

mitochondrial genome but no autosomal regions, and describe the biases introduced

550 during capture such as DNA fragment size, jumping PCR and divergence between bait

and target species. The study performed by Hernandez-Rodriguez et al. (2018)

552 systematically analyzed the capture performance and library complexity. While they

described that pooling different libraries into the same hybridization is feasible, they

554 did not discuss how many of them should be pooled. Also, they concluded that

performing multiple libraries from the same extract or even from different extracts from

556  the same sample can increase the final complexity. Finally, they recommended performing two capture rounds for the same library. On the other hand, White et al.

558  (2019) suggested to do only one capture round, at least when eDNA is higher than 2-3%, stressing the importance of pooling libraries as well as taking into consideration

560  the eDNA content, as first proposed by Hernandez-Rodriguez et al.

The present study addresses these gaps left unexplored by the previous studies. We

562  focused our analysis on a representative set of samples with very low proportions of endogenous content (< 1%) as are often found in the field. After screening 302

564  samples, we found that up to 70% of samples are below this threshold, similar to what was already described (White et al., 2019). Hence, if time and economic reasons

566  hinder the ability to collect and select the best samples, the only available one(s) might have low eDNA. This may be a common situation when using historical samples,

568  aiming for a large sample size, or if an interesting sampling location is particularly challenging in terms of low eDNA (such as Campo Ma'an, Figure 1B).

570  For these reasons, it is of utmost importance to characterize ways to maximize the amount of data to be recovered from these types of samples. In this regard, we have

572  extensively evaluated how to increase library complexity without doing more extractions or library preparations from the same sample, how many libraries to pool

574  together, and how much starting amount of DNA should be used in a capture, as well as the impact of endogenous content for pooling.

576  Consistent with previous findings (Hernandez-Rodriguez et al., 2018; White et al., 2019), we determined that assessing the endogenous content of fecal samples and

578  pooling them equi-endogenously is a practical way to equally distribute raw reads between samples. Importantly, the correct estimation of the proportion of eDNA is key

580 for the success of this method. Thus, we recommend the usage of shotgun sequencing

(Hernandez-Rodriguez et al., 2018) rather than qPCR estimates, since the later can

582 easily fluctuate due to the presence of inhibitors (Morin et al., 2001).

In regard to the performance of target capture sequencing experiments, gaining new

584 unique reads is crucial to reach higher sensitivity, which is a good predictor of capture

success. Here, we have established an approach to obtain new unique reads using

586 the same prepared libraries. Since it is mainly during capture experiments when the

molecular diversity is reduced, we propose to perform additional hybridizations from

588 the same library so the final coverage can reach higher values. If the library complexity

is already very low, the only solution is to re-extract DNA or prepare a new library from

590 the same sample (Hernandez-Rodriguez et al., 2018).

We observed a better performance (MC and CS) in small pools, when evaluating initial

592 results derived from the entire dataset. However, after correcting for other variables

that differ among pools, the effect is attenuated and can only explain ~4% of the

594 variance, an effect that may be largely negligible for most studies. Moreover,

performing additional hybridizations can also compensate for this effect. Therefore, we

596 do not conclude, based on this data, that pool size is a major contributor to

performance. However, in cases where libraries have small proportions of eDNA, we

598 would advocate for the reduction of the number of samples per pool so that pipetting

volumes may remain larger, and as a consequence variability due to pipetting error

600 may be reduced. Otherwise when the eDNA proportion is not a limiting factor, pooling

more libraries together and performing additional hybridizations can be a good

602 strategy.

It is worth noting that without taking into consideration individual sample quality and the amount of starting material used, one of the most influential variables on the performance of target capture enrichment experiments is the hybridization experiment itself. After accounting for all other variables, it still explains 18% of the variation. This is due to the technical complexity and variability inherent to these experiments. Careful equipment optimization, material selection, preparation and experience will aid in minimizing this variation, although it is likely to remain a sensitive experiment that requires diligence.

Finally, we have illustrated that a sequencing effort of exome-captured fecal samples with low eDNA (< 1%) should be set at ~3 million reads per library in a pool to avoid exhausting the molecular complexity. We have benefited from the usage of double-barcoded and double-indexed libraries to multiplex many samples in a single sequencing lane. This becomes a great advantage because we can utilize high throughput sequencing technologies at a lower price per read.

To summarize, when starting a project involving fecal samples, we recommend screening your set of samples based on quantity and quality of the DNA extracted. If having related individuals in the study should be avoided, microsatellite genotyping could be an option, helping as well to discard samples with high amount of PCR inhibitors. Further selection of samples should be based on the proportion of eDNA; we recommend using shotgun sequencing from the prepared libraries. Performing re-extractions of the most valuable samples and preparing replicate libraries from each extract can help increase the final molecular complexity. As we have shown here, another approach to achieve higher molecular complexity is based on conducting additional hybridizations of the captured libraries, always pooling libraries in an equi-

endogenous manner, and starting with more library material than the standard protocol

628 suggests. Finally, we suggest not sequencing the captured libraries very deeply, since

their molecular complexity is already very low and over-sequencing can result in rapidly

630 depleting the economic feasibility of the experiment.

In the study presented here we have thoroughly explored approaches to increase the

632 molecular diversity and capture sensitivity and hence the final coverage of exome

captured fecal samples with extremely low endogenous content in an attempt to help

634 laboratories facing the challenges of working with non-invasive samples.


# Acknowledgments

# References

676  Arandjelovic, M., Guschanski, K., Schubert, G., Harris, T. R., Thalmann, O., Siedel,
        H., & Vigilant, L. (2009). Two-step multiplex polymerase chain reaction improves
678     the speed and accuracy of genotyping using DNA from noninvasive and
        museum samples. *Molecular Ecology Resources*, *9*(1), 28–36. doi:
680     10.1111/j.1755-0998.2008.02387.x

     Arandjelovic, M., Head, J., Rabanal, L. I., Schubert, G., Mettke, E., Boesch, C., …
682     Vigilant, L. (2011). Non-invasive genetic monitoring of wild central chimpanzees.
        *PLoS ONE*, *6*(3), e14761. doi: 10.1371/journal.pone.0014761

684  Auton, A., Abecasis, G. R., Altshuler, D. M., Durbin, R. M., Bentley, D. R.,
        Chakravarti, A., … Schloss, J. A. (2015). A global reference for human genetic
686     variation. *Nature*, Vol. 526, pp. 68–74. doi: 10.1038/nature15393

     Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: A flexible trimmer for
688     Illumina sequence data. *Bioinformatics*, *30*(15), 2114–2120. doi:
        10.1093/bioinformatics/btu170

690  Brinkman, T. J., Schwartz, M. K., Person, D. K., Pilgrim, K. L., & Hundertmark, K. J.
        (2010). Effects of time and rainfall on PCR success using DNA extracted from
692     deer fecal pellets. *Conservation Genetics*, *11*(4), 1547–1552. doi:
        10.1007/s10592-009-9928-7

694  Carøe, C., Gopalakrishnan, S., Vinner, L., Mak, S. S. T., Sinding, M. H. S.,
        Samaniego, J. A., … Gilbert, M. T. P. (2018). Single-tube library preparation for
696     degraded DNA. *Methods in Ecology and Evolution*, *9*(2), 410–419. doi:
        10.1111/2041-210X.12871

698  Carpenter, M. L., Buenrostro, J. D., Valdiosera, C., Schroeder, H., Allentoft, M. E.,
        Sikora, M., … Bustamante, C. D. (2013). Pulling out the 1%: Whole-Genome
700     capture for the targeted enrichment of ancient dna sequencing libraries.
        *American Journal of Human Genetics*, *93*(5), 852–864. doi:
702     10.1016/j.ajhg.2013.10.002

     Csilléry, K., Johnson, T., Beraldi, D., Clutton-Brock, T., Coltman, D., Hansson, B., …
704     Pemberton, J. M. (2006). Performance of marker-based relatedness estimators
        in natural populations of outbred vertebrates. *Genetics*, *173*(4), 2091–2101. doi:
706     10.1534/genetics.106.057331

     De Barba, M., Waits, L. P., Genovesi, P., Randi, E., Chirichella, R., & Cetto, E.
708     (2010). Comparing opportunistic and systematic sampling methods for non-
        invasive genetic monitoring of a small translocated brown bear population.
710     *Journal of Applied Ecology*, *47*(1), 172–181. doi: 10.1111/j.1365-
        2664.2009.01752.x

712  De Manuel, M., Kuhlwilm, M., Frandsen, P., Sousa, V. C., Desai, T., Prado-Martinez,
        J., … Marques-Bonet, T. (2016). Chimpanzee genomic diversity reveals ancient
714     admixture with bonobos. *Science*, *354*(6311), 477–481. doi:
        10.1126/science.aag2602

716  Ferreira, C. M., Sabino-Marques, H., Barbosa, S., Costa, P., Encarnação, C., Alpizar-
        Jara, R., … Alves, P. C. (2018). Genetic non-invasive sampling (gNIS) as a cost-
718     effective tool for monitoring elusive small mammals. *European Journal of Wildlife
        Research*, *64*(4). doi: 10.1007/s10344-018-1188-8

720  Fickel, J., Lieckfeldt, D., Ratanakorn, P., & Pitra, C. (2007). Distribution of haplotypes

722     and microsatellite alleles among Asian elephants (Elephas maximus) in Thailand. *European Journal of Wildlife Research*, *53*(4), 298–303. doi: 10.1007/s10344-007-0099-x

724 Fischer, A., Wiebe, V., Pääbo, S., & Przeworski, M. (2004). Evidence for a Complex Demographic History of Chimpanzees. *Molecular Biology and Evolution*, *21*(5), 726     799–808. doi: 10.1093/molbev/msh083

Goossens, B., Chikhi, L., Utami, S. S., De Ruiter, J., & Bruford, M. W. (2000). A multi-728     samples, multi-extracts approach for microsatellite analysis of faecal samples in an arboreal ape. *Conservation Genetics*, *1*(2), 157–162. doi: 730     10.1023/A:1026535006318

Gordon, D., Huddleston, J., Chaisson, M. J. P., Hill, C. M., Kronenberg, Z. N., 732     Munson, K. M., … Eichler, E. E. (2016). Long-read sequence assembly of the gorilla genome. *Science, 352*(6281), aae0344. doi: 10.1126/science.aae0344

734 Harestad, A. S., & Bunnell, F. L. (1987). Persistence of Black-Tailed Deer Fecal Pellets in Coastal Habitats. *The Journal of Wildlife Management*, *51*(1), 33. doi: 736     10.2307/3801624

Hernandez-Rodriguez, J., Arandjelovic, M., Lester, J., de Filippo, C., Weihmann, A., 738     Meyer, M., … Marques-Bonet, T. (2018). The impact of endogenous content, replicates and pooling on genome capture from faecal samples. *Molecular* 740     *Ecology Resources*, *18*(2), 319–333. doi: 10.1111/1755-0998.12728

Hicks, A. L., Lee, K. J., Couto-Rodriguez, M., Patel, J., Sinha, R., Guo, C., … 742     Williams, B. L. (2018). Gut microbiomes of wild great apes fluctuate seasonally in response to diet. *Nature Communications*, *9*(1), 1786. doi: 10.1038/s41467-744     018-04204-w

Inoue, E., Akomo-Okoue, E. F., Ando, C., Iwata, Y., Judai, M., Fujita, S., … 746     Yamagiwa, J. (2013). Male genetic structure and paternity in western lowland gorillas (Gorilla gorilla gorilla). *American Journal of Physical Anthropology*, 748     *151*(4), 583–588. doi: 10.1002/ajpa.22312

King, S. R. B., Schoenecker, K. A., Fike, J. A., & Oyler-McCance, S. J. (2018). Long-750     term persistence of horse fecal DNA in the environment makes equids particularly good candidates for noninvasive sampling. *Ecology and Evolution*, 752     *8*(8), 4053–4064. doi: 10.1002/ece3.3956

Kircher, M., Sawyer, S., & Meyer, M. (2012). Double indexing overcomes 754     inaccuracies in multiplex sequencing on the Illumina platform. *Nucleic Acids Research*, *40*(1), 1–8. doi: 10.1093/nar/gkr771

756 Kong, S. W., Lee, I. H., Liu, X., Hirschhorn, J. N., & Mandl, K. D. (2018). Measuring coverage and accuracy of whole-exome sequencing in clinical context. *Genetics* 758     *in Medicine, 20*(12), 1617–1626. doi: 10.1038/gim.2018.51

Kühl, H. S., Boesch, C., Kulik, L., Haas, F., Arandjelovic, M., Dieguez, P., … Kalan, 760     A. K. (2019). Human impact erodes chimpanzee behavioral diversity. *Science (New York, N.Y.)*, *363*(6434), 1453–1455. doi: 10.1126/science.aau4532

762 Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, *25*(14), 1754–1760. doi: 764     10.1093/bioinformatics/btp324

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., … Durbin, R. 766     (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, *25*(16), 2078–2079. doi: 10.1093/bioinformatics/btp352

768 Locke, D. P., Hillier, L. W., Warren, W. C., Worley, K. C., Nazareth, L. V., Muzny, D.

770  M., … Wilson, R. K. (2011). Comparative and demographic analysis of orang-
utan genomes. *Nature*, *469*(7331), 529–533. doi: 10.1038/nature09687

Mengüllüoğlu, D., Fickel, J., Hofer, H., & Förster, D. W. (2019). Non-invasive faecal
772  sampling reveals spatial organization and improves measures of genetic
diversity for the conservation assessment of territorial species: Caucasian lynx
774  as a case species. *PLoS ONE*, *14*(5). doi: 10.1371/journal.pone.0216549

Mikkelsen, T. S., Hillier, L. W., Eichler, E. E., Zody, M. C., Jaffe, D. B., Yang, S. P., …
776  Waterston, R. H. (2005). Initial sequence of the chimpanzee genome and
comparison with the human genome. *Nature*, *437*(7055), 69–87. doi:
778  10.1038/nature04072

Morin, P. A., Chambers, K. E., Boesch, C., & Vigilant, L. (2001). Quantitative PCR
780  analysis of DNA from noninvasive samples fro accurate microsatellite genotyping
of wild chimpanzees. *Molecular Ecology*, 1835–1844.

782  Morin, P. A., Wallis, J., Moore, J. J., Chakraborty, R., & Woodruff, D. S. (1993). Non-
invasive sampling and DNA amplification for paternity exclusion, community
784  structure, and phylogeography in wild chimpanzees. *Primates*, *34*(3), 347–356.
doi: 10.1007/BF02382630

786  Nsubuga, A. M., Robbins, M. M., Roeder, A. D., Morin, P. A., Boesch, C., & Vigilant,
L. (2004). Factors affecting the amount of genomic DNA extracted from ape
788  faeces and the identification of an improved sample storage method. *Molecular
Ecology*, *13*(7), 2089–2094. doi: 10.1111/j.1365-294X.2004.02207.x

790  Orkin, J. D., Yang, Y., Yang, C., Yu, D. W., & Jiang, X. (2016). Cost-effective scat-
detection dogs: Unleashing a powerful new tool for international mammalian
792  conservation biology. *Scientific Reports*, *6*(1), 34758. doi: 10.1038/srep34758

Ouborg, N. J., Pertoldi, C., Loeschcke, V., Bijlsma, R. K., & Hedrick, P. W. (2010).
794  Conservation genetics in transition to conservation genomics. *Trends in
Genetics*, *26*(4), 177–187. doi: 10.1016/j.tig.2010.01.001

796  Perry, G. H. (2014). The Promise and Practicality of Population Genomics Research
with Endangered Species. *International Journal of Primatology*, *35*(1), 55–70.
798  doi: 10.1007/s10764-013-9702-z

Perry, G. H., Marioni, J. C., Melsted, P., & Gilad, Y. (2010). Genomic-scale capture
800  and sequencing of endogenous DNA from feces. *Molecular Ecology*, *19*(24),
5332–5344. doi: 10.1111/j.1365-294X.2010.04888.x

802  Prado-Martinez, J., Sudmant, P. H., Kidd, J. M., Li, H., Kelley, J. L., Lorente-Galdos,
B., … Marques-Bonet, T. (2013). Great ape genetic diversity and population
804  history. *Nature*, *499*(7459), 471–475. doi: 10.1038/nature12228

Primmer, C. R. (2009). From conservation genetics to conservation genomics.
806  *Annals of the New York Academy of Sciences*, Vol. 1162, pp. 357–368. doi:
10.1111/j.1749-6632.2009.04444.x

808  Quinlan, A. R., & Hall, I. M. (2010). BEDTools: A flexible suite of utilities for
comparing genomic features. *Bioinformatics*, *26*(6), 841–842. doi:
810  10.1093/bioinformatics/btq033

Reiners, T. E., Encarnação, J. A., & Wolters, V. (2011). An optimized hair trap for
812  non-invasive genetic studies of small cryptic mammals. *European Journal of
Wildlife Research*, *57*(4), 991–995. doi: 10.1007/s10344-011-0543-9

814  Rohland, N., & Reich, D. (2012). Cost-effective, high-throughput DNA sequencing
libraries for multiplexed target capture. *Genome Research*, *22*(5), 939–946. doi:
816  10.1101/gr.128124.111

Scally, A., Dutheil, J. Y., Hillier, L. W., Jordan, G. E., Goodhead, I., Herrero, J., …
818      Durbin, R. (2012). Insights into hominid evolution from the gorilla genome
         sequence. *Nature*, *483*(7388), 169–175. doi: 10.1038/nature10842
820  Schwartz, M. K., Luikart, G., & Waples, R. S. (2007). Genetic monitoring as a
         promising tool for conservation and management. *Trends in Ecology and*
822      *Evolution*, Vol. 22, pp. 25–33. doi: 10.1016/j.tree.2006.08.009
     Shafer, A. B., Wolf, J. B., Alves, P. C., Bergströ, L., Bruford, M. W., Brä nnströ, I., …
824      Zielin, P. (2015). Genomics and the challenging translation into conservation
         practice. *Trends in Ecology & Evolution*, *30*(2), 78–87. doi:
826      10.1016/j.tree.2014.11.009
     Snyder-Mackler, N., Majoros, W. H., Yuan, M. L., Shaver, A. O., Gordon, J. B., Kopp,
828      G. H., … Tung, J. (2016). Efficient genome-wide sequencing and low-coverage
         pedigree analysis from noninvasively collected samples. *Genetics*, *203*(2), 699–
830      714. doi: 10.1534/genetics.116.187492
     Städele, V., & Vigilant, L. (2016). Strategies for determining kinship in wild
832      populations using genetic data. *Ecology and Evolution*, *6*(17), 6107–6120. doi:
         10.1002/ece3.2346
834  Steiner, C. C., Putnam, A. S., Hoeck, P. E. A., & Ryder, O. A. (2013). Conservation
         Genomics of Threatened Animal Species. *Annual Review of Animal Biosciences*,
836      *1*(1), 261–281. doi: 10.1146/annurev-animal-031412-103636
     Stenglein, J. L., Waits, L. P., Ausband, D. E., Zager, P., & Mack, C. M. (2010).
838      Efficient, Noninvasive Genetic Sampling for Monitoring Reintroduced Wolves.
         *Journal of Wildlife Management*, *74*(5), 1050–1058. doi: 10.2193/2009-305
840  Taberlet, P., Luikart, G., & Waits, L. P. (1999). Noninvasive genetic sampling: Look
         before you leap. *Trends in Ecology and Evolution*, *14*(8), 323–327. doi:
842      10.1016/S0169-5347(99)01637-7
     Thalmann, O., Hebler, J., Poinar, H. N., Pääbo, S., & Vigilant, L. (2004). Unreliable
844      mtDNA data due to nuclear insertions: A cautionary tale from analysis of humans
         and other great apes. *Molecular Ecology*, *13*(2), 321–335. doi: 10.1046/j.1365-
846      294X.2003.02070.x
     van der Valk, T., Lona Durazo, F., Dalén, L., & Guschanski, K. (2017). Whole
848      mitochondrial genome capture from faecal samples and museum-preserved
         specimens. *Molecular Ecology Resources*, *17*(6), e111–e121. doi:
850      10.1111/1755-0998.12699
     Vigilant, L., & Guschanski, K. (2009). Using genetics to understand the dynamics of
852      wild primate populations. *Primates*, *50*(2), 105–120. doi: 10.1007/s10329-008-
         0124-z
854  Wedrowicz, F., Karsa, M., Mosse, J., & Hogan, F. E. (2013). Reliable genotyping of
         the koala (Phascolarctos cinereus) using DNA isolated from a single faecal
856      pellet. *Molecular Ecology Resources*, *13*(4), 634–641. doi: 10.1111/1755-
         0998.12101
858  White, L. C., Fontsere, C., Lizano, E., Hughes, D. A., Angedakin, S., Arandjelovic, M.,
         … Vigilant, L. (2019). A roadmap for high-throughput sequencing studies of wild
860      animal populations using noninvasive samples and hybridization capture.
         *Molecular Ecology Resources*, *19*(3), 609–622. doi: 10.1111/1755-0998.12993
862  Xue, Y., Prado-Martinez, J., Sudmant, P. H., Narasimhan, V., Ayub, Q., Szpak, M., …
         Scally, A. (2015). Mountain gorilla genomes reveal the impact of long-term
864      population decline and inbreeding. *Science*, *348*(6231), 242–245. doi:

866
# Data Accessibility

868 All raw sequencing data have been deposited at ENA and are available under the accession code PRJEB37173 (http://www.ebi.ac.uk/ena/data/view/PRJEB37173).

870
# Author Contributions

CF, TMB, DAH and EL designed the study. MA and HSK direct the Pan African

872 Programme: The Cultured Chimpanzee. MA and HSK obtained funding for the project.

MA, PD, AA, SA, EAA, MB, GB, TD, MEN, ACG, JH, PK, AKK, MK, KL, JL, GM, LJO,

874 AP, MR, FS, VV and RMW supervised, conducted field work and collected samples.

CF, MAE, EL, JL, MA performed experiments. CF and DAH performed the analysis.

876 MAE, MK, DAH, TMB, EL provided analytical support. CF wrote the manuscript with

input from all authors.

878

# Supporting Information

880 Additional supporting information with extended methods and supplementary figures

and tables can be found online in the Supporting information section at the end of the

882 article.

# Conflict of Interest

884 Authors declare no conflict of interest.

**FIGURE 1.** Sample description. (a) Geographical location of the 15 sites from the Pan African Programme: The Cultured Chimpanzee (PanAf). (b) Endogenous DNA (eDNA) content for all screened samples according to geographic origin. The maximum value of the x-axis has been set to 10% eDNA for visual purposes. (c) eDNA distribution for all screened samples. Samples with > 10% eDNA are excluded (N=5). In the boxplot, lower and upper hinges correspond to first and third quartiles and the lower and upper whiskers extend to the smallest or largest value no further than 1.5 times the interquartile range (distance between the 1st and 3rd quartile).


**FIGURE 2.** Pooling strategy illustration. P1 has 10 libraries with average endogenous of 0.81%. We performed two primary pools of 2 μg and 4 μg each that were further divided into four hybridization pools, two at 1 μg and two at 2 μg. P2 has 20 libraries with average endogenous of 0.69%. Two primary pools of 4 μg were divided into four hybridization pools of 1 μg each and two hybridizations pools of 2 μg. P3 has 30 libraries and an average endogenous of 0.49%. Two primary pools of 6 μg and 4 μg were distributed into six hybridization pools of 1μg and two hybridization pools of 2 μg each. Colors represent the sequencing batch.


**FIGURE 3.** Capture performance and sequencing. (a) Percentage of eDNA among hybridizations, structured by pools (P1, P2 and P3). (b) Sequencing stats across all samples for the 18 hybridizations in 3,75 HiSeq 4000 lanes. (c) Distribution of production reads across 18 hybridizations. The colors red, blue and yellow found in the box plots for figure (a) and (c) denote the sequencing batch to which each hybridization was assigned. In the boxplots, lower and upper hinges correspond to first and third quartiles and the lower and upper whiskers extend to the smallest or largest value no further than 1.5 times the interquartile range (distance between the 1st and 3rd quartile).


**FIGURE 4**. Analysis of variance. (a) Estimated variance explained from univariate linear models after rank normal transforming each performance summary statistic (columns). LC stands for library complexity and DP describes read depth at different cutoffs (1, 4, 10, 20 and 50 reads) (b) Multivariate type I ANOVA of the experimental variables affecting Capture Sensitivity (CS) at depth 1. Both models are built down-sampling libraries to 1,500,000 reads.


**FIGURE 5.** Summary stats after down-sampling to 1,500,000 reads: (a) Enrichment factor and (d) Capture Specificity (c) Capture Sensitivity at depth 1 for the 18 hybridizations in P1, P2 and P3; colors illustrate sequencing batch. (d) Library complexity contrasting the amount of starting DNA (1 μg or 2 μg) in down-sampled data and structured by pools (P1=Pool1, P2=Pool2, P3=Pool3). See Figure 2 for more details on pools. In the boxplots, lower and upper hinges correspond to first and third quartiles and the lower and upper whiskers extend to the smallest or largest value no further than 1.5 times the interquartile range (distance between the 1st and 3rd quartile).


**FIGURE 6**. Analysis of coverage and LC with hybridizations done with 1 μg. (a) Coverage after merging data from additional hybridizations with up to 2, 4 and 6 for P1, P2 and P3. (b) Comparison of average LC curves of individual hybridizations belonging to pools with different size. Each line is the average of libraries within each hybridization and the surrounding area is the standard deviation. (c) Two examples comparing the effect of pool size on the average LC curves from merged hybridization: P1 (10 samples) - 1 hybridization, P2 (20 samples) – 2 hybridizations and P3 (30 samples) – 3 hybridizations; and P1 (10 samples) - 2 hybridizations, P2 (20 samples) – 4 hybridizations and P3 (30 samples) – 6 hybridizations. Sample Lib1-6D in P2 was removed from the analysis due to low coverage.

| Pool | Average eDNA content (range) | Hybridization ID | Number of pooled libraries | Total DNA | Sequencing Batch |
|---|---|---|---|---|---|
| Pool 1 (**P1**) | 0.81% (0.60% - 0.85%) | P1E1 | 10 | 1 μg | SeqBatch1 |
| | | P1E2 | 10 | 1 μg | SeqBatch2 |
| | | P1E3 | 9 | 2 μg | SeqBatch3 |
| | | P1E4 | 9 | 2 μg | SeqBatch3 |
| Pool 2 (**P2**) | 0.69% (0.58% - 0.80%) | P2E1 | 20 | 1 μg | SeqBatch1 |
| | | P2E2 | 20 | 1 μg | SeqBatch1 |
| | | P2E3 | 20 | 1 μg | SeqBatch2 |
| | | P2E4 | 20 | 1 μg | SeqBatch2 |
| | | P2E5 | 19 | 2 μg | SeqBatch3 |
| | | P2E6 | 19 | 2 μg | SeqBatch3 |
| Pool 3 (**P3**) | 0.49% (0.41% - 0.66%) | P3E1 | 30 | 1 μg | SeqBatch1 |
| | | P3E2 | 30 | 1 μg | SeqBatch1 |
| | | P3E3 | 30 | 1 μg | SeqBatch1 |
| | | P3E4 | 30 | 1 μg | SeqBatch2 |
| | | P3E5 | 30 | 1 μg | SeqBatch2 |
| | | P3E6 | 30 | 1 μg | SeqBatch2 |
| | | P3E7 | 26 | 2 μg | SeqBatch3 |
| | | P3E8 | 26 | 2 μg | SeqBatch3 |

936

938 **TABLE 1.** Pooling Strategy. Sixty libraries were divided into 3 pools for capture hybridization experiments in 4 replicates for P1, 6 replicates for P2 and 8 replicates for P3. Total DNA
940 represents the starting material for each capture hybridization.