

Title: Whole-genome characterization in pedigreed non-human primates using Genotyping-By-Sequencing (GBS) and imputation.

Authors: Ben N Bimber¹, Michael J Raboin¹, John Letaw¹, Kimberly Nevonen¹, Jennifer E Spindel², Susan R McCouch², Rita Cervera-Juanes¹, Eliot Spindel¹, Lucia Carbone¹, Betsy Ferguson¹, Amanda Vinson^{1*}

Institutional affiliations: ¹Primate Genetics Section, Oregon National Primate Research Center, Beaverton, Oregon, and the Oregon Health & Science University, Portland, Oregon

²Section of Plant Breeding and Genetics, School of Integrative Plant Sciences, Cornell University, Ithaca, New York.

*Corresponding author. Amanda Vinson, Asst. Scientist, Primate Genetics Section, Oregon National Primate Research Center, Beaverton, Oregon

Author email addresses: Ben Bimber: bimber@ohsu.edu
Michael Raboin: raboin@ohsu.edu
John Letaw: letaw@ohsu.edu
Kimberly Nevonen: nevonen@ohsu.edu
Jennifer Spindel: jes46@cornell.edu
Susan McCouch: srm4@cornell.edu
Rita Cervera-Juanes: cerveraj@ohsu.edu
Eliot Spindel: spindele@ohsu.edu
Lucia Carbone: carbone@ohsu.edu
Betsy Ferguson: fergusob@ohsu.edu
Amanda Vinson: vinsona@ohsu.edu

1 Title: Whole-genome characterization in pedigreed non-human primates
2 using Genotyping-By-Sequencing (GBS) and imputation.
3

4 Authors: Ben N Bimber¹, Michael J Raboin¹, John Letaw¹, Kimberly
5 Nevonen¹, Jennifer E Spindel², Susan R McCouch², Rita Cervera-
6 Juanes¹, Eliot Spindel¹, Lucia Carbone¹, Betsy Ferguson¹,
7 Amanda Vinson^{1*}
8

9 Institutional affiliations: ¹Primate Genetics Section, Oregon National Primate Research
10 Center, Beaverton, Oregon, and the Oregon Health & Science
11 University, Portland, Oregon
12
13 ²Section of Plant Breeding and Genetics, School of Integrative
14 Plant Sciences, Cornell University, Ithaca, New York.
15

16 *Corresponding author. Amanda Vinson, Asst. Scientist, Primate Genetics Section,
17 Oregon National Primate Research Center, Beaverton, Oregon
18

19 Author email addresses: Ben Bimber: bimber@ohsu.edu
20
21 Michael Raboin: raboin@ohsu.edu
22
23 John Letaw: letaw@ohsu.edu
24
25 Kimberly Nevonen: nevonen@ohsu.edu
26
27 Jennifer Spindel: jes46@cornell.edu
28
29 Susan McCouch: srm4@cornell.edu
30
31 Rita Cervera-Juanes: cerveraj@ohsu.edu
32
33 Eliot Spindel: spindele@ohsu.edu
34
35 Lucia Carbone: carbone@ohsu.edu
36
37 Betsy Ferguson: fergusob@ohsu.edu
38
39 Amanda Vinson: vinsona@ohsu.edu
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

34 ABSTRACT:

35 *Background:* Rhesus macaques are widely used in biomedical research, but the application of
36 genomic information in this species to better understand human disease is still undeveloped.
37 Whole-genome sequence (WGS) data in pedigreed macaque colonies could provide substantial
38 experimental power, but the collection of WGS data in large cohorts remains a formidable
39 expense. Here, we describe a cost-effective approach that selects the most informative
40 macaques in a pedigree for whole-genome sequencing, and imputes these dense marker data
41 into all remaining individuals having sparse marker data, obtained using Genotyping-By-
42 Sequencing (GBS).

43 *Results:* We developed GBS for the macaque genome using a single digest with *Pst*I, followed
44 by sequencing to 30X coverage. From GBS sequence data collected on all individuals in a 16-
45 member pedigree, we characterized an optimal 22,455 sparse markers spaced ~125 kb apart.
46 To characterize dense markers for imputation, we performed WGS at 30X coverage on 9 of the
47 16 individuals, yielding ~10.2 million high-confidence variants. Using the approach of
48 “Genotype Imputation Given Inheritance” (GIGI), we imputed alleles at an optimized dense set
49 of 4,920 variants on chromosome 19, using 490 sparse markers from GBS. We assessed
50 changes in accuracy of imputed alleles, 1) across 3 different strategies for selecting individuals
51 for WGS, i.e., a) using “GIGI-Pick” to select informative individuals, b) sequencing the most
52 recent generation, or c) sequencing founders only; and 2) when using from 1-9 WGS individuals
53 for imputation. We found that accuracy of imputed alleles was highest using the GIGI-Pick
54 selection strategy (median 92%), and improved very little when using >4 individuals with WGS
55 for imputation. We used this ratio of 4 WGS to 12 GBS individuals to impute an expanded set of
56 ~14.4 million variants across all 20 macaque autosomes, achieving ~85-88% accuracy per
57 chromosome.

58 *Conclusions:* We conclude that an optimal tradeoff exists at the ratio of 1 individual selected for
59 WGS using the GIGI-Pick algorithm, per 3-5 relatives selected for GBS, a cost savings of ~67-

60 83% over WGS of all individuals. This approach makes feasible the collection of accurate,
61 dense genome-wide sequence data in large pedigreed macaque cohorts without the need for
62 expensive WGS data on all individuals.

63

64 KEYWORDS: whole-genome sequencing, genotyping-by-sequencing, imputation, macaque,
65 pedigree

66

67 BACKGROUND:

68

69 The analysis of whole-genome sequence data in non-human primates (NHPs) can play a
70 significant role in advancing the application of genomic medicine to human disease. Potential
71 uses of these data include the identification of novel genetic variants that influence conserved
72 pathways of disease pathology, the development of novel therapeutics that target these
73 variants, and the characterization of variants that influence efficacy and response to
74 therapeutics. Given their high degree of genetic and physiological similarity to humans, and their
75 ubiquity in biomedical research, it is surprising that the use of the rhesus macaque for these
76 purposes has been so slow to develop. One likely reason for this delay is the dearth of genome-
77 wide sequence information on sufficient numbers of animals to support such studies, which
78 typically require large numbers of phenotyped and genotyped subjects. However, the collection
79 of dense sequence data in large cohorts remains a formidable expense, and a cost-effective
80 solution to this problem is needed if we are to reap the full benefit of non-human primate (NHP)
81 models in both basic science and preclinical research.

82

83 Whole-genome sequencing (WGS) of large cohorts remains a very expensive undertaking, both
84 now and likely long after we achieve the \$1,000 per genome benchmark. Several sequencing
85 strategies have been developed to address this problem, each of which strikes a different

86 balance between sequencing costs, sequence depth, and coverage across the genome.
87 Although deep WGS is the most unbiased and comprehensive method for surveying genetic
88 variants [1], at approximately \$3000/genome for standard 30X coverage, its cost remains
89 prohibitive in the foreseeable future for large cohort studies. A second strategy aims to cover the
90 whole genome but at greatly reduced depth (i.e., “low coverage sequencing”), which lowers
91 costs to \$200-600/genome. However, this approach reduces the accuracy of resulting genotype
92 data, particularly for smaller studies of rare or low-frequency variants [2]. Another strategy is to
93 sequence only a portion of the genome, i.e., “reduced representation” approaches, which offers
94 a reasonable compromise between sequencing depth and breadth of coverage. The most
95 common of the reduced representation approaches is whole-exome sequencing, currently
96 ~\$300/genome for 30X coverage, in which a commercial hybridization kit is used to capture
97 genomic fragments enriched for exons in protein-coding genes. While this approach produces
98 coverage of genomic regions that are of interest to many Mendelian diseases, coverage of
99 regulatory elements or other non-coding regions is sacrificed. Moreover, most commercial
100 exome capture tools are designed for humans or rodents, and thus will miss some portion of the
101 NHP exome.

102
103 More recently, a reduced representation approach called genotyping-by-sequencing (GBS) has
104 lowered the cost per genome dramatically, by taking advantage of classical molecular biology
105 methods that capture a more evenly distributed subset of the genome. In the GBS method,
106 standard restriction enzymes target conserved cut sites that span the genome, and the resulting
107 fragments are sequenced to the desired coverage. While these fragments still represent only a
108 small portion of the genome, they can be distributed more evenly than in other methods such as
109 exome capture. Importantly, because the GBS approach does not require proprietary capture
110 technology and can be highly multiplexed, costs can be reduced to as little as \$50 per genome
111 and this approach can be applied to species that lack available commercial arrays. This

112 approach has been applied to many agricultural and other economically important species to
113 construct dense genetic linkage maps and identify QTLs [3-8], to improve genome assemblies
114 [9], and to investigate population structure, diversity, and evolutionary history [10-12].

115

116 Further gains in the amount of sequence information obtained at the lowest possible cost could
117 be achieved by combining GBS data with imputation, particularly for NHP cohorts with pedigree
118 information. In this approach, whole-genome sequence data collected in selected individuals
119 within the pedigree are used to impute dense genotypes into their many relatives, in which only
120 sparse genotype data (e.g., obtained by GBS) has been collected. These sparse data from GBS
121 are used to anchor the imputation of genotypes at intervening and more densely spaced loci
122 across the genome, by leveraging information on expected allele-sharing among relatives. This
123 strategy is appealing for many captive NHP breeding colonies, where deep and well-defined
124 pedigrees could permit extremely cost-effective, whole-genome characterization. However, the
125 selection of the most informative animals in the pedigree for whole-genome sequencing is
126 expected to have a large impact on the success of this approach, and studies addressing
127 optimal selection strategies have only been published for human pedigrees, which are typically
128 much smaller and less complex than those characterized for NHP cohorts.

129

130 Although whole-genome sequencing combined with GBS and imputation presents a significant
131 opportunity for obtaining dense sequence data at minimal cost, this approach has not yet been
132 applied to a pedigreed NHP cohort. Thus, our objectives were to, 1) develop a reliable GBS
133 method in the macaque genome to support pedigree-based imputation; 2) to assess the extent
134 and accuracy of imputed dense marker data from WGS using sparse marker data from GBS;
135 and 3) to compare the accuracy of imputed dense marker data among different strategies
136 commonly used to select individuals for WGS, and among different ratios of WGS to GBS
137 individuals in the pedigree. Here, we show that a *Pst*I digest in the macaque genome produces

138 >22,000 high-quality sparse variants that are suitable for use in imputation. We further show
139 that the pedigree-based “Genotype Imputation Given Inheritance” imputation approach (i.e.,
140 “GIGI”; [13]), combined with the GIGI-Pick method [14] of selecting individuals for WGS, allowed
141 us to impute >14 million variants throughout a 16-member pedigree with ~85-88% accuracy,
142 using only 4 individuals with WGS and 12 individuals with GBS. This strategy represents a
143 reasonable tradeoff between sequencing costs, and the amount and quality of dense sequence
144 data obtained on as many individuals as possible.

145

146 METHODS:

147

148 *Animal care and welfare:* All macaque samples used in this study were collected during routine
149 veterinary care procedures approved by the Institutional Animal Care and Use Committee of the
150 Oregon Health & Science University (Protocol Number: IS00002621); these samples are part of
151 the much larger ONPRC DNA Biobank. Animal care personnel and staff veterinarians of the
152 ONPRC provide routine and emergency health care to all animals in accordance with the Guide
153 for the Care and Use of Laboratory Animals, and the ONPRC is certified by the Association for
154 Assessment and Accreditation of Laboratory Animal Care International.

155

156 *Pedigree configuration and validation:* We selected 16 closely related animals from the larger
157 Oregon National Primate Research Center (ONPRC) colony pedigree as the focus of this study
158 (see Fig. 1). These animals were selected to represent the most common relationships in the
159 colony, including parent/offspring, half-sibling, half-avuncular, half-cousin, and
160 grandparent/grandchild relationships. Because assumed pedigree relationships may prove to be
161 incorrect when comprehensive genotype data are examined, we explored the accuracy of our
162 focal 16-member pedigree using a set of ~5,000 markers on chromosome 19 generated from
163 our GBS sequencing experiments (see *Imputation accuracy on chromosome 19 across*

164 *selection strategies*, below), employing algorithms that assess Mendelian consistent error both
165 pairwise between relatives and within families, as implemented in PedCheck [15] and GIGI-
166 Check [16] software. No significant departures from expected patterns of allele-sharing were
167 noted, confirming the validity of the pedigree configuration depicted in Fig. 1. Nine animals were
168 selected using an *ad hoc* approach for whole-genome sequencing in this study, based on their
169 position within the pedigree.

170

171 *Genomic DNA Isolation and Quantification:* Genomic DNA (gDNA) was extracted from 3 ml of
172 whole blood using the ArchivePure DNA Blood Kit (5 Prime, Inc.), following the manufacturer's
173 recommendations. Genomic DNA was quantified with the Qubit High Sensitivity dsDNA Assay
174 (Life Technologies, CA).

175

176 *Genotyping-By-Sequencing (GBS):* To determine the optimal restriction enzymes for conducting
177 GBS in rhesus macaques, we first performed *in silico* prediction of cut sites using the most
178 recent build of the macaque genome [17] that would be expected to produce 60,000-100,000
179 DNA fragments in the 200-500 bp size range, while also minimizing the presence of repeat
180 sequences (e.g. retrotransposons, DNA satellites). We initially tested the enzymes *ApeKI*, *BgIII*,
181 *EcoRI*, *HindIII*, *PspXI*, *PstI*, and *Sall*, ultimately selecting *BgIII* and *PstI* as the two enzymes
182 most likely to meet these criteria. We then generated GBS libraries based on these 2 enzymes
183 using a modified version of the method described by Elshire et al. [18]. Specifically, to create the
184 adaptors, oligonucleotides for the top and bottom strands for each barcoded adaptor and for the
185 two common adaptors (one for *BgIII* and one for *PstI*) were paired and annealed in 1X
186 Annealing Buffer (20mM NaCl, 10mM Tris-HCl pH 7.5, 2mM MgCl₂) using a thermal cycler (3
187 min at 95C, ramp down 1.6C/min for 44 cycles, cool to 4C). All adaptors were quantified with the
188 Qubit Broad Range dsDNA Assay (Invitrogen). Each of the 32 barcoded adaptors was then
189 paired with a common adaptor at a 1:1 ratio. Each of the 16 genomic DNAs was digested with

190 *BglII* and *PstI* in separate reactions. All 32 reactions (500ng DNA, 10U enzyme, in 20uL
191 volume) were incubated for 2 hours at 37°C, and digests were ligated (400U T4 DNA Ligase
192 (New England Biolabs) to adaptor mixes (4.5ng *BglII*, 15ng *PstI*, in 50ul volume) for 1 hour at
193 22°C. Four (4) ul from each ligation reaction was combined into two separate pools, one per
194 enzyme. Both pools were cleaned with DNA Clean and Concentrator (Zymo Research) and
195 eluted in 50uL. Following amplification parameters in Elshire et al [18], PCR was performed on
196 10 ul of each pool (Q5 High Fidelity 2X MM (New England Biolabs), 25 pmol of each primer, in
197 50 ul volume) using Primers A and B, to extend and complete the sequencing adaptors.
198 Libraries were purified using the Qiaquick PCR Purification Kit (Qiagen), quantified with the
199 Qubit High Sensitivity dsDNA Assay (Invitrogen) and validated with the Bioanalyzer High
200 Sensitivity Assay (Agilent). A one-sided 0.8X size selection with AMPure XP beads (Beckman-
201 Coulter) was used to enrich larger size fragments. Libraries were sequenced on an Illumina
202 NextSeq at the Oregon Health & Science University Integrated Genomics Laboratory to produce
203 30X coverage.

204

205 *Whole Genome Sequencing:* Per sample, 1 µg of gDNA was sheared using a Bioruptor
206 UCD200 (Diagenode, Denville, NJ), generating fragments around 300 bp. Libraries were
207 constructed using the NEXTflex DNA Sequencing Kit and NEXTflex DNA barcodes (BIOO
208 Scientific, Austin,TX) following the manufacturer's instructions. Briefly, the ends of the sheared
209 gDNA were repaired and adenylated, then ligated to barcoded adaptors using the reagents
210 provided. Next, fragments of 200-400 bp were excised from a 1% agarose gel. The products
211 were amplified by PCR using 8 cycles, then purified using 1X AMPure XP beads (Beckman-
212 Coulter). The final libraries were quantified with the Qubit High Sensitivity dsDNA Assay (Life
213 Technologies, CA) and validated using the 2100 Bioanalyzer High Sensitivity Assay (Agilent
214 Technologies, Santa Clara, CA). Libraries were sequenced on a HiSeq3000 at the Oregon State

215 University Center for Genome Research and Biocomputing, to produce 30x coverage with
216 paired-end, 150 bp reads.

217

218 *Analysis of Sequence Data:* Both whole genome and GBS data were processed using the best
219 practice recommendations from the Broad Institute's Genome Analysis Toolkit (GATK; [19, 20]),
220 adapted for rhesus macaque. Briefly, paired-end reads were trimmed using Trimmomatic [21],
221 and aligned to the most current rhesus macaque reference genome, "MacaM" [17], using
222 Burrows-Wheeler Aligner [22]. GATK's HaplotypeCaller was used to produce VCF files,
223 followed by genotype calling using GenotypeGVCFs. The resulting VCF was filtered for quality,
224 strand bias, and proximity to the read end. Additional hard filters include removal of, 1) single-
225 nucleotide variant (SNV) clusters within a 20 bp span, 2) SNVs located within regions with
226 greater than twice the mean coverage (potential CNV or mismapped reads), 3) SNVs that
227 display non-Mendelian inheritance, and 4) SNVs located within repetitive regions. The analyses
228 also employed Picard tools [23] and FASTQC [24] for quality control of the raw data, JBrowse
229 [25] to visualize data, and BEDTools [26] to evaluate SNV and imputation marker distribution.
230 Sequence data were managed and analyzed using DISCVR-Seq [27], a LabKey Server-based
231 system [28].

232

233 *Sequencing and imputation strategy:* We focused on chromosome 19 as a test case in order to
234 develop an analytical pipeline that could be applied to all the remaining chromosomes. We
235 performed imputation using the method of GIGI ("Genotype Imputation Given Inheritance"; [13]),
236 as this method has been successfully used to impute genotypes with high accuracy in extended
237 human pedigrees. This approach infers inheritance vectors (IVs, representing shared
238 chromosomal segments) at sparse marker locations conditioned on observed sparse marker
239 genotypes, and then infers IVs at dense marker locations conditioned on the sparse marker IVs,
240 together with the genetic map. The probability distribution is then estimated for each missing

241 genotype at a dense marker position, conditioned on observed genotypes at all dense marker
242 positions, corresponding allele frequencies, and IVs corresponding to dense markers. In the
243 last step, genotypes may be called using these estimated probabilities, based on user-defined
244 thresholds. We estimated inheritance vectors according to the algorithm of [29], as
245 implemented using a Markov-Chain Monte-Carlo (MCMC) sampler in the `gl_auto` function in the
246 software package for genetic epidemiology MORGAN 3; available at [30]. The GIGI approach
247 has been implemented in a software package of the same name, and is available at [31].

248

249 To characterize the sparse set of markers needed to facilitate imputation of dense marker
250 genotypes on chromosome 19, we identified a set of markers that could be detected reliably by
251 GBS for as many macaques as possible. Accordingly, we selected a set of high-quality SNVs
252 that, 1) were sequenced to at least 20X depth across the majority of GBS libraries, 2) were
253 spaced evenly across the genome, 3) had minor allele frequencies (MAF) >0.25 , and 4) were in
254 excess of what was needed to meet the desired goal of ~ 0.5 - 1.0 cM average marker spacing.
255 We refer to these as “framework” markers, as discussed in Cheung et al., 2013 [13]. Using this
256 approach, the desired spacing can be maintained in an approximate fashion, even when
257 individuals are missing a substantial amount of genotype data, an outcome characteristic of the
258 GBS method [32]. Second, we selected a non-overlapping set of SNVs from our WGS data that
259 were more densely distributed than the framework markers, designated as our “dense” markers,
260 and which we attempted to impute into animals having only sparse framework marker
261 genotypes from GBS. These dense markers were selected from the set of all high-confidence
262 SNVs identified in our cohort.

263

264 To determine the success of imputing dense marker data into animals having only sparse
265 framework marker data, we evaluated the accuracy of imputed alleles for each recipient based
266 on their framework marker data obtained by GBS. Here, we define accuracy as the proportion

267 of alleles imputed correctly among all attempted allele calls at that position, such that a correctly
268 imputed allele is concordant with the allele call from either WGS or GBS sequence data. We
269 additionally define rare variants as those having only 1 copy among a total of 30 chromosomes
270 (i.e., singletons, present at ~3% frequency) with WGS data that we have studied to date,
271 including the 9 individuals discussed in this paper and an additional 6 unrelated rhesus
272 macaques (unpublished data). We define accuracy of imputation for rare variants as the
273 proportion of rare alleles imputed correctly among all rare variant heterozygotes called from
274 either WGS or GBS sequence data.

275

276 We assessed accuracy of imputed variants over a range of 1-9 animals with dense marker data
277 from WGS, imputed into the remaining pedigree members using only sparse data from GBS.
278 Individuals with WGS data were added consecutively in the following order: B, H, J, F, M, K, P,
279 C, D (SEE FIG. 1). Thus, the first scenario used only the most informative animal (B) with WGS
280 to impute genotypes into the remaining 15 animals within the pedigree. Subsequent scenarios
281 retained the previous animal(s), added the next most informative animal, and imputed
282 genotypes into the remaining animals within the pedigree. This procedure was conducted
283 iteratively, until all 9 animals with WGS were used to impute genotypes into the remaining 7
284 animals in the pedigree. We used the GIGI-Pick algorithm [14] to rank our 9 animals with WGS.
285 This algorithm calculates a metric of coverage, defined as the expected percentage of allele
286 copies called for a variant at a random locus, conditional on fixed IVs for a specific choice of
287 individual(s), and then iteratively selects those individuals with the highest coverage, calculated
288 by integrating over all possible genotype configurations within a given pedigree. This algorithm
289 is implemented in the suite of software based on the GIGI approach, and is available at [33]. We
290 evaluated accuracy of imputed genotypes for each of our recipient macaques by masking all
291 non-framework genotypes in recipient animals, and comparing imputed genotypes to masked
292 genotypes obtained from either WGS or GBS data, depending on the data available for each

293 recipient. Specifically, imputed genotypes were compared to genotypes from WGS where
294 available, but for recipients with only GBS data available, imputed genotypes were compared to
295 genotypes at any SNVs with coverage by GBS that were not designated as framework markers.
296 Imputed genotypes were called as the most probable genotype at that position, using allele
297 frequencies established from all Indian-origin rhesus macaques sequenced to date at the
298 ONPRC as a reference (n=15, including the 9 animals from this study and 6 additional unrelated
299 animals).

300

301 To evaluate differences in imputation success associated with different sequencing strategies,
302 we compared the accuracy of genotypes imputed by GIGI on chromosome 19, among 3
303 different sequencing strategies, including GIGI-Pick and two common heuristic methods. These
304 2 heuristic methods include whole-genome sequencing of, 1) pedigree founders only, or 2) the
305 most recent generation (i.e., individuals typically located at the bottom of the pedigree). To
306 compare the different selection strategies, we examined accuracy for the scenario in which
307 dense markers from 3 animals selected for WGS are imputed into the remaining 13 pedigree
308 members with GBS data, based on using individuals B, H, and J (GIGI-Pick selections), B, C,
309 and D (“Founders”), and M, P, and K (“Pedigree bottom”) strategies (see Fig. 1). Imputed
310 genotypes were called using 2 different strategies available in the GIGI algorithm, i.e.,
311 genotypes were either above the default probability threshold, or simply as the most probable
312 genotype at that position (“Threshold” vs. “Most Likely”, respectively), using allele frequencies
313 established from all Indian-origin rhesus macaques sequenced to date at the ONPRC as a
314 reference. To assess the utility of the GIGI imputation approach for imputing low-frequency
315 variants, we further evaluated genotype accuracy within the GIGI-Pick strategy of 3 WGS into
316 13 GBS described above.

317

318

319 RESULTS:

320 *Whole-genome Sequencing and Variant Calling:* We obtained an average 566,035,688 read
321 pairs per sample (range 495,617,772-735,313,000) for each of the 9 individuals with WGS.
322 These reads were aligned to MacaM, the most recent macaque genome build [17], to produce
323 an average 27X coverage across the genome (range 24-33X). From these reads, a total of
324 10,193,425 high-confidence SNVs were identified across all 9 individuals, with an average of
325 5,037,341 variants detected per individual. The transition/transversion ratio (Ti/Tv) observed in
326 this study was 2.17, consistent with observations in larger macaque cohorts (unpublished data).
327 This set of sites served as the source of our optimal dense marker set, as described in Methods.

328
329 *Genotyping-By-Sequencing and Variant Calling:* For each of the 16 pedigree members, we
330 prepared and sequenced GBS libraries based on individual digests for both *BglIII* and *PstI*.
331 Among *BglIII* libraries, we obtained an average 3,754,352 reads per sample, resulting in an
332 average 4,851,004 base-pairs (bp) from 42,919 fragments with at least 20X coverage per
333 sample (equivalent to 0.17% of the genome). In contrast, among *PstI* libraries, we obtained an
334 average 5,686,709 reads per sample, resulting in an average 10,682,162 bp from 130,247
335 fragments with >20X coverage per sample (equivalent to 0.38% of the genome) (Fig. 2A-2B).
336 Notably, although the *PstI* libraries originally had ~1.5X more reads than *BglIII* libraries, they had
337 >3-fold the number of fragments with high-depth coverage. Virtually all GBS fragments were
338 adjacent to their predicted restriction sites, but a small number appeared to be distant from
339 these sites (Fig. 2C-D). While these results may reflect off-target sequencing, it is also possible
340 that they reflect restriction sites in the genome of one or more individuals not predicted by the
341 current macaque reference genome.

342
343 We next determined the number of high-quality SNVs available from each digest that would be
344 suitable for use as framework markers in imputation. We first restricted SNVs to those with

345 >20X coverage, leaving a total of 52,500 high-confidence variants in the *BglII* samples, and
346 239,428 variants in *PstI* samples. We further restricted these SNVs to only those that were
347 concordant between WGS and GBS data, which included 96.3% of SNVs for *BglII* (range 95.9-
348 97.2% among individuals) and 96.6% of SNVs for *PstI* (range 96.1-97.2% among individuals).
349 To maximize the probability of the variant being present in as many animals as possible, we
350 also restricted SNVs to only those with MAF >0.25, which further reduced these numbers to
351 7,399 variants for *BglII* and 22,455 variants for *PstI*, with an average distance between SNVs of
352 376,818 bp and 125,280 bp, respectively (see Fig. 2E-F). We were not able to call genotypes in
353 all individuals for all SNV sites, due to variation among individuals in sequence quality at each
354 site. From the *PstI* data, all individuals had sufficient data to call genotypes at an average
355 15,516 (~69% of 22,455) of these SNVs, but only an average of 4,418 (~60% of 7,399) of these
356 sites could be called for all individuals from the *BglII* data. Based on the significantly greater
357 numbers of high-quality SNVs across the macaque genome available from *PstI* sequence data,
358 we chose this enzyme for all final imputation analyses.

359
360 *Imputation accuracy on chromosome 19 across 3 different selection strategies:* Our initial tests
361 of imputation focused on chromosome 19. To characterize the set of framework markers
362 needed to facilitate imputation on this chromosome, we selected 490 variants spaced ~100 kb
363 apart, from the set of 981 high-MAF variants on this chromosome, as described above (Fig.
364 2G). To characterize the set of dense markers to be imputed on this chromosome, we selected
365 4,920 variants spaced ~10 kb apart and which were not framework markers, from the set of
366 260,000 SNVs located on this chromosome (see above). We additionally removed a set of 578
367 markers that consistently performed poorly in imputation. This reduced and optimized set of
368 dense markers on chromosome 19 was used in all imputation analyses on this chromosome, in
369 order to limit the computational time required.

370

371 Using these chromosome-specific framework and dense marker sets, we evaluated the “Bottom
372 of Pedigree”, the “Founders”, and the “GIGI-Pick” strategies for selecting the 3 most informative
373 of the 9 individuals with WGS data, followed by imputation of dense markers into the remaining
374 13 individuals, based on their framework marker data from GBS. For both genotype-calling
375 methods, the GIGI-Pick selection strategy produced slightly higher median accuracy of imputed
376 genotypes than either of the other strategies, at 89.5% (“Most Likely” method, ML) or 90.1%
377 accuracy (“Threshold” method, THR), compared to median accuracy of 88.4% (ML) or 88.8%
378 (THR) in the “Bottom of Pedigree” strategy, and 88.6% (ML) or 87.6% (THR) in the “Founders”
379 strategy (Fig. 3A-B). However, more individuals in the “GIGI-Pick” strategy displayed greater
380 genotype accuracy than in either of the other 2 strategies (interquartile range (IQR) from 89.5%-
381 92.6% for GIGI-Pick, compared to 84.2-90.1% for “Bottom of Pedigree”, and 85.9-89.5% for
382 “Founders”). While the difference in median accuracy estimated by both genotype calling
383 methods was <1%, 100% of genotypes were called in the ML method, while only 47.7% were
384 called using the THR method across all strategies. Thus, the ML method called genotypes at an
385 average 2,350 more markers per subject, while maintaining nearly identical accuracy. The GIGI-
386 Pick strategy did produce one individual that consistently displayed much lower genotype
387 accuracy than all other individuals; individual D had only 72.9-75.6% of genotypes accurately
388 imputed, depending on calling method. Neither the “Bottom of Pedigree” nor the “Founders”
389 strategy produced any individuals with substantially lower genotype accuracy than other
390 pedigree members. Finally, we used the GIGI-Pick results to examine accuracy of imputed
391 singleton alleles, as defined in Methods. These rare alleles were imputed a total of 405 times,
392 representing 178 distinct alleles, with an accuracy of 93%.

393

394 *Imputation accuracy on chromosome 19 across 9 sequencing ratios:* Using the same
395 chromosome 19 framework and dense marker sets, we used the GIGI algorithm with the ML
396 genotype-calling method to impute our dense markers from 1-9 individuals with WGS into the

397 remaining 7-15 pedigree members, as described in Methods. Among all imputation scenarios
398 that added consecutively from 1 to 9 individuals with WGS, median accuracy among recipients
399 increased from 86.3% to 93.2%, while variation in accuracy decreased (IQR from 5% at N=1
400 WGS, to 1% at N=9 WGS). Variability in accuracy improved in a stepwise fashion; greater
401 variability tended to occur in parallel with an increase in median accuracy, but would improve
402 during the following scenario in which the greater accuracy was retained or further increased.
403 Individual D was consistently an outlier across 8 out of 9 scenarios, from 1 through 5 at 76%
404 accuracy, then increasing to 84% accuracy in scenarios 6-8, due to the inclusion of WGS data
405 from K, the child of D. Median accuracy surpassed 90% beginning at N=4 individuals with
406 WGS; although variation in accuracy continued to decrease across all remaining scenarios, only
407 slight gains in median accuracy were observed beyond this sequencing ratio (i.e., imputing from
408 4 individuals with WGS into 12 individuals with GBS) (Fig. 4).

409
410 *Imputation of dense markers across the genome:* To evaluate our imputation strategy across
411 the whole genome, we employed the same criteria outlined above to generate framework
412 marker sets for each of the 20 macaque autosomes. The number of framework markers per
413 chromosome ranged from 350 to 636, with mean spacing between framework markers among
414 all chromosomes of ~273 kb (108 kb-467 kb). At this stage, we expanded our dense marker set
415 to a total 14,384,988 SNVs across the macaque genome, by including SNVs discovered
416 previously from whole-genome sequence data in an additional 6 unrelated ONPRC animals (in
417 prep). Using the strategy identified by our analyses as the one most likely to maximize
418 genotype accuracy while minimizing overall costs, we used the first 4 individuals among our 9
419 with WGS ranked by the GIGI-Pick algorithm, and imputed our expanded dense marker data
420 into the remaining 12 pedigree members, based on the ML call method (Fig 5). Per
421 chromosome, median accuracy ranged from 85-88%; this accuracy is somewhat lower than the
422 ~92% accuracy achieved in our original experiment using chromosome 19. However, unlike our

423 previous analysis that imputed a much smaller set of dense markers on chromosome 19, our
424 final genome-wide imputation included all known variants. As before, D had consistently lower
425 genotype accuracy at 73-77%. Individual E was also an outlier on multiple chromosomes, with
426 accuracy ranging from 79-84%.

427

428 DISCUSSION:

429

430 The rhesus macaque is widely used in academic biomedical research, primarily due to its utility
431 as a model of human HIV infection and pathology. Although this species is well-known for the
432 susceptibility to HIV that it shares with humans [34, 35], it is not widely appreciated that
433 macaques naturally display variation in susceptibility to a broad spectrum of diseases and
434 disorders that mimic those found in humans, e.g., dyslipidemia, alcohol abuse, macular
435 degeneration, and anxiety [36-42]. While the macaque was identified early as a high priority for
436 assembly of a reference genome, and a draft genome subsequently published in 2007[43], the
437 systematic application of genome-wide data in the macaque to the study of human health and
438 disease has yet to materialize. Since 2007, next-generation sequencing technology has
439 speeded the collection of genomic data at steadily decreasing cost, but only a relatively small
440 number of additional macaque genomes have been explored for variation, and none have yet
441 been systematically applied to the problem of human disease. This is unfortunate, given that
442 rhesus macaque colonies at many of the national primate research centers constitute a powerful
443 resource for genetic analysis of disease that is on par with many human genetic isolates, due to
444 their maintenance in large outbred and pedigreed colonies within a homogeneous environment.
445 Here, in order to catalyze the application of genomic data in macaques to the study of human
446 disease, we present an approach that will make feasible the collection of accurate, dense
447 genome-wide sequence data in large numbers of pedigreed macaques without the need for
448 expensive whole-genome sequence data on all individuals.

449

450 Our approach is based on using a low-cost reduced representation sequencing method
451 (genotyping-by-sequencing, GBS), to facilitate pedigree-based imputation of dense marker
452 genotypes from selected individuals with whole genome sequence data. In this study, we
453 evaluated the ability of 2 candidate restriction enzymes (*BglII* and *PstI*) to produce genomic
454 fragments for GBS, using both *in silico* and empirical methods. When compared to *BglII*, we
455 show that *PstI* produces substantially larger numbers of high-quality SNVs that are supported by
456 greater sequence read depth. Further, we found that *PstI* libraries provided sufficient coverage
457 over more than twice the number of high-quality variants needed to generate the sparse
458 “framework” markers required to support imputation. This is important because fluctuations in
459 sequencing coverage among individuals are a known characteristic of the GBS method [32],
460 resulting in the frequent inability to call genotypes at all sites in all individuals. Thus, *PstI*
461 produces far more high-quality SNVs than are actually needed, which increases the likelihood of
462 observing a minimum number of framework markers in every individual. Ultimately, the extent
463 of coverage provided by *PstI* allowed us to impute genotypes at ~14.4 million SNVs over all 20
464 autosomes, using only 4 individuals with WGS and 12 individuals with GBS, at a median 85-
465 88% accuracy throughout our 16-member pedigree. This approach could be applied at a
466 relatively reasonable cost to other managed or natural colonies of Indian-origin rhesus
467 macaques with pedigree information, and potentially to similar groups of other macaque
468 subspecies.

469

470 The selection of individuals for WGS that will maximize the accuracy of imputed genotypes
471 throughout the pedigree is a critical component of this approach. We compared GIGI-Pick [14],
472 a pedigree-based statistical approach to prioritizing subjects for WGS, to two other common
473 heuristic methods for selecting individuals for WGS, including sequencing only the most recent
474 generation of the pedigree (“Bottom of Pedigree”), and sequencing only pedigree founders

475 (“Founders”). Due to the small size of our sample pedigree, we used 3 individuals with WGS to
476 impute genotypes at our streamlined set of dense markers, into 13 remaining individuals based
477 on their framework marker genotypes from GBS data. We show that while median accuracy is
478 only somewhat greater for the GIGI-Pick approach compared to the other two approaches,
479 many more individuals overall displayed higher accuracy using the GIGI-Pick selection method,
480 as reflected in the strong upward shift of the interquartile range. Importantly, rare alleles were
481 imputed with exceptionally high accuracy using the GIGI-Pick selection strategy, suggesting that
482 this strategy offers powerful support for downstream analysis of rare variant effects on complex
483 traits. It is possible that these 3 selection methods may perform differently for alternative
484 pedigree configurations, e.g., in a more shallow pedigree, sequencing founders or the most
485 recent pedigree members may provide information equivalent to the more formal strategy
486 implemented in GIGI-Pick. However, we note that the GIGI-Pick approach results in a clear
487 advantage even in our small pedigree that extends to only 2 generations, but which includes
488 many of the most common relationship types found in NHP colonies. Our results are consistent
489 with those of Cheung et al. [14], in that the GIGI-Pick selection approach substantially
490 outperformed both the “Bottom of Pedigree” and “Founders” (i.e., “PRIMUS” in [14]) approaches
491 in the ability to impute common alleles, although our results indicate more consistent accuracy
492 with the “Founders” approach than with the “Bottom of Pedigree” approach.

493
494 Using the WGS individuals ranked in order of priority by the GIGI-Pick algorithm, we examined
495 the gain in accuracy of imputed genotypes throughout the pedigree achieved by the consecutive
496 addition of from 1-9 whole-genome sequences. Our results demonstrate that there are
497 excellent compromises available that balance sequencing costs and the ability to obtain dense
498 and accurate marker data. While the accuracy of imputed genotypes was greatest when using
499 all 9 individuals with WGS, most of this accuracy was achieved using the first 4 WGS
500 individuals, i.e., at 4 WGS individuals, median accuracy is at ~92.4% but increased only another

501 0.8% with the addition of the remaining 5 genomes. These results suggest that an optimal
502 tradeoff between the animals selected for WGS and GBS exists at the ratio of 1 individual
503 selected for WGS, per 3-5 relatives selected for GBS, a cost savings of ~67-83% over WGS of
504 all 6 individuals. We note that these estimates of accuracy were based on careful selection of an
505 optimal, and thus reduced, set of dense markers that were available on chromosome 19. While
506 this strategy was used deliberately to reduce the total computational time required for this study,
507 in our subsequent imputation of the full set of dense markers across the genome, median
508 accuracy only decreased by 5-8% for all chromosomes.

509

510 The increase in overall accuracy observed with additional WGS individuals was not shared
511 uniformly among all individuals in the pedigree. For example, while there was a large increase
512 in accuracy between 3 and 4 individuals with WGS, all of this change is due to increased
513 accuracy in A, a founder (see Fig. 1). In this scenario, this improvement is almost certainly due
514 to the inclusion of WGS data from F, a child of A. We note that D remained an outlier in the
515 distribution of genotype accuracy throughout virtually all imputation analyses based on the
516 ranking of WGS individuals by the GIGI-Pick algorithm. This may be due to the limited initial
517 selection of WGS individuals located in the far right lineage, i.e., only when K, P, and C are
518 added to J and used for imputation does accuracy rise for D. This result is consistent with the
519 GIGI-Pick approach, which balances the selection of closely related individuals within the
520 pedigree to facilitate phasing of genotypes, with the selection of more distant relatives to
521 increase the chance of observing unique founder alleles [14]. Because of this compromise, we
522 note that while sequencing only pedigree founders is not the best strategy for maximizing
523 accuracy, founders may remain unselected using the GIGI-Pick approach when the ability to
524 phase genotypes produces more expected allele calls than does the observation of unique
525 founder alleles, for a fixed number of selected individuals. This result also highlights the
526 importance that prior knowledge of phenotypes plays in selecting individuals for WGS. If traits of

527 interest are known to segregate in a particular lineage within the larger pedigree, it may be
528 advisable to manually assign either founders or a close descendant in that lineage for WGS, if
529 neither individual is selected using a more unbiased approach.

530

531 The imputation of dense, genome-wide genotypes with high accuracy will allow the unbiased
532 mapping of genetic variants in the macaque genome to disease traits, using either linkage or
533 association approaches. Both of these approaches are important tools in translational research,
534 and should further advance the understanding of human disease already made possible by
535 research in this species. Large pedigreed colonies of macaques, such as the ~4,500 macaques
536 found at the Oregon National Primate Research Center, provide an almost unequaled resource
537 for translational genetic research, due to their multi-generational pedigree structure and the
538 excess of rare and low-frequency variants expected to segregate within this pedigree. Rare and
539 low-frequency variants are expected to play a significant role in human disease [44-47], and we
540 have shown that we can impute these variants in the macaque genome with high accuracy and
541 at a reasonable cost, using the approach we outline here. Moreover, our findings suggest that
542 this approach can be modified to support specific research goals. For example, it may be
543 beneficial to take advantage of the less accurate but greater amount of information provided by
544 the full set of imputed dense markers during initial discovery of variants either linked to or
545 associated with a disease trait, while fine-mapping or replication of a putative trait locus might
546 employ a reduced, optimal set of dense markers likely to provide greater genotype accuracy
547 over a smaller region of interest.

548

549 In this study, we demonstrated that it is feasible to obtain comprehensive genome-wide variation
550 at a fraction of the cost of whole-genome sequencing using the GBS method and pedigree-
551 based imputation, which we describe for the first time in a non-human primate genome.
552 However, imputed genotypes will only be as accurate as the underlying whole-genome

553 sequence data and the reference genome to which it is compared. The macaque genome has
554 undergone multiple revisions; however, it remains in draft form and is less complete than many
555 other common model organisms [17, 43]. There are also extremely limited available data on
556 genetic variants in macaques, and no databases with comprehensive information on known
557 variants or population-level allele frequency information are publicly accessible. Both of these
558 factors present obstacles for accurate whole-genome variant calling in macaques, and will thus
559 reduce accuracy for any genotyping approach. Improvements in both of these areas are
560 urgently needed in order to fully realize the value of the macaque as a genetic model of human
561 disease.

562

563 DECLARATIONS:

564

565 1) List of abbreviations: WGS, whole-genome sequencing; GBS, Genotyping-By-Sequencing;
566 SNV, single-nucleotide variant; ONPRC, Oregon National Primate Research Center; MAF,
567 minor allele frequency; GIGI, Genotype Imputation Given Inheritance; CNV, copy number
568 variant; BWA, Burrows-Wheeler Aligner; VCF, variant call format; GATK, Genome Analyzer
569 ToolKit; MCMC, Markov Chain Monte Carlo; ML, most likely genotype calling method; THR,
570 threshold genotype calling method.

571

572 2) Ethics approval and consent to participate:

573 *Animal care and welfare:* All macaque samples used in this study were collected during routine
574 veterinary care procedures approved by the Institutional Animal Care and Use Committee of the
575 Oregon Health & Science University (Protocol Number: IS00002621); these samples are part of
576 the much larger ONPRC DNA Biobank. Animal care personnel and staff veterinarians of the
577 ONPRC provide routine and emergency health care to all animals in accordance with the Guide

578 for the Care and Use of Laboratory Animals, and the ONPRC is certified by the Association for
579 Assessment and Accreditation of Laboratory Animal Care International.

580

581 3) Consent for publication: Not applicable.

582

583 4) Availability of data and material: The GBS and WGS sequence data files supporting the
584 conclusions of this article are in the process of submission to the NCBI Sequence Read Archive,
585 <http://www.ncbi.nlm.nih.gov/sra>.

586

587 5) Competing interests: A. Vinson receives compensation as an external consultant to Novo-
588 Nordisk, USA. No other authors declare competing interests.

589

590 6) Funding sources: This project was supported by the Office of the Director/Office of
591 Research Infrastructure Programs (OD/ORIP) of the NIH (grant no. OD011092). This funding
592 body played no part in the design of the study, in the collection, analysis, and interpretation of
593 data, or in the writing of this manuscript.

594

595 7) Authors contributions: All authors made substantial contributions to the conception and
596 design of this study. AV, BB, BF, LC, and RCJ wrote or critically revised the manuscript. BB,
597 AV, JL, KN, LC, BF, MR, and ES analyzed or interpreted data. All authors gave approval for the
598 submission of this manuscript.

599

600 8) Acknowledgements: We owe particular thanks to Dr. Laura Cox of the Texas Biomedical
601 Research Institute for initial discussions of GBS methods. We also thank Dr. Ellen Wijsman for
602 helpful discussion of key issues around imputation in pedigrees, and Dr. Charles Cheung for
603 support with the GIGI suite of software. We thank the ONPRC Bioinformatics service core for

604 initial processing and analysis of sequence data, and the ONPRC DNA Bank for access to NHP
605 samples used in this project. This project was supported by the Office of the Director/Office of
606 Research Infrastructure Programs (OD/ORIP) of the NIH (grant no. OD011092).

607

608 9) Authors' information: Not applicable.

609

610 10) Endnotes: Not applicable.

611

612 11) References:

- 613 1. Meynert AM, Ansari M, FitzPatrick DR, Taylor MS: **Variant detection sensitivity and**
614 **biases in whole genome and exome sequencing.** *BMC Bioinformatics* 2014, **15**:247.
- 615 2. Li Y, Sidore C, Kang HM, Boehnke M, Abecasis GR: **Low-coverage sequencing:**
616 **implications for design of complex trait association studies.** *Genome Res* 2011,
617 **21**(6):940-951.
- 618 3. Bielenberg DG, Rauh B, Fan S, Gasic K, Abbott AG, Reighard GL, Okie WR, Wells CE:
619 **Genotyping by Sequencing for SNP-Based Linkage Map Construction and QTL**
620 **Analysis of Chilling Requirement and Bloom Date in Peach [*Prunus persica* (L.)**
621 **Batsch].** *PloS one* 2015, **10**(10):e0139406.
- 622 4. De Donato M, Peters SO, Mitchell SE, Hussain T, Imumorin IG: **Genotyping-by-**
623 **sequencing (GBS): a novel, efficient and cost-effective genotyping method for**
624 **cattle using next-generation sequencing.** *PloS one* 2013, **8**(5):e62137.
- 625 5. Palti Y, Vallejo RL, Gao G, Liu S, Hernandez AG, Rexroad CE, 3rd, Wiens GD:
626 **Detection and Validation of QTL Affecting Bacterial Cold Water Disease**
627 **Resistance in Rainbow Trout Using Restriction-Site Associated DNA Sequencing.**
628 *PloS one* 2015, **10**(9):e0138435.
- 629 6. Pootakham W, Ruang-Areerate P, Jomchai N, Sonthirod C, Sangsrakru D, Yoocha T,
630 Theerawattanasuk K, Nirapathpongporn K, Romruensukharom P, Tragoonrung S *et al.*
631 **Construction of a high-density integrated genetic linkage map of rubber tree**
632 **(*Hevea brasiliensis*) using genotyping-by-sequencing (GBS).** *Front Plant Sci* 2015,
633 **6**:367.
- 634 7. Spindel J, Wright M, Chen C, Cobb J, Gage J, Harrington S, Lorieux M, Ahmadi N,
635 McCouch S: **Bridging the genotyping gap: using genotyping by sequencing (GBS)**
636 **to add high-density SNP markers and new value to traditional bi-parental mapping**
637 **and breeding populations.** *Theor Appl Genet* 2013, **126**(11):2699-2716.
- 638 8. Xu Y, Huang L, Ji D, Chen C, Zheng H, Xie C: **Construction of a dense genetic**
639 **linkage map and mapping quantitative trait loci for economic traits of a doubled**
640 **haploid population of *Pyropia haitanensis* (Bangiales, Rhodophyta).** *BMC Plant Biol*
641 2015, **15**:228.
- 642 9. Zhou Z, Liu S, Dong Y, Gao S, Chen Z, Jiang J, Yang A, Sun H, Guan X, Jiang B *et al.*
643 **High-Density Genetic Mapping with Interspecific Hybrids of Two Sea Urchins,**
644 ***Strongylocentrotus nudus* and *S. intermedius*, by RAD Sequencing.** *PloS one* 2015,
645 **10**(9):e0138585.

- 646 10. Burrell AM, Pepper AE, Hodnett G, Goolsby JA, Overholt WA, Racelis AE, Diaz R, Klein
647 PE: **Exploring origins, invasion history and genetic diversity of Imperata cylindrica**
648 **(L.) P. Beauv. (Cogongrass) in the United States using genotyping by sequencing.**
649 *Mol Ecol* 2015, **24**(9):2177-2193.
- 650 11. Escudero M, Eaton DA, Hahn M, Hipp AL: **Genotyping-by-sequencing as a tool to**
651 **infer phylogeny and ancestral hybridization: a case study in Carex (Cyperaceae).**
652 *Mol Phylogenet Evol* 2014, **79**:359-367.
- 653 12. Johnson JL, Wittgenstein H, Mitchell SE, Hyma KE, Temnykh SV, Kharlamova AV,
654 Gulevich RG, Vladimirova AV, Fong HW, Acland GM *et al.* **Genotyping-By-Sequencing**
655 **(GBS) Detects Genetic Structure and Confirms Behavioral QTL in Tame and**
656 **Aggressive Foxes (Vulpes vulpes).** *PLoS one* 2015, **10**(6):e0127013.
- 657 13. Cheung CY, Thompson EA, Wijsman EM: **GIGI: an approach to effective imputation**
658 **of dense genotypes on large pedigrees.** *Am J Hum Genet* 2013, **92**(4):504-516.
- 659 14. Cheung CY, Marchani Blue E, Wijsman EM: **A statistical framework to guide**
660 **sequencing choices in pedigrees.** *Am J Hum Genet* 2014, **94**(2):257-267.
- 661 15. O'Connell JR, Weeks DE: **PedCheck: a program for identification of genotype**
662 **incompatibilities in linkage analysis.** *Am J Hum Genet* 1998, **63**(1):259-266.
- 663 16. Cheung CY, Thompson EA, Wijsman EM: **Detection of Mendelian consistent**
664 **genotyping errors in pedigrees.** *Genet Epidemiol* 2014, **38**(4):291-299.
- 665 17. Zimin AV, Cornish AS, Maudhoo MD, Gibbs RM, Zhang X, Pandey S, Meehan DT,
666 Wipfler K, Bosinger SE, Johnson ZP *et al.* **A new rhesus macaque assembly and**
667 **annotation for next-generation sequencing analyses.** *Biol Direct* 2014, **9**(1):20.
- 668 18. Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, Mitchell SE: **A**
669 **robust, simple genotyping-by-sequencing (GBS) approach for high diversity**
670 **species.** *PLoS one* 2011, **6**(5):e19379.
- 671 19. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,
672 Altshuler D, Gabriel S, Daly M *et al.* **The Genome Analysis Toolkit: a MapReduce**
673 **framework for analyzing next-generation DNA sequencing data.** *Genome Res* 2010,
674 **20**(9):1297-1303.
- 675 20. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A,
676 Jordan T, Shakir K, Roazen D, Thibault J *et al.* **From FastQ data to high confidence**
677 **variant calls: the Genome Analysis Toolkit best practices pipeline.** *Curr Protoc*
678 *Bioinformatics* 2013, **11**(1110):11 10 11-11 10 33.
- 679 21. Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina**
680 **sequence data.** *Bioinformatics* 2014, **30**(15):2114-2120.
- 681 22. Li H, Durbin R: **Fast and accurate long-read alignment with Burrows-Wheeler**
682 **transform.** *Bioinformatics* 2010, **26**(5):589-595.
- 683 23. **Picard Tools** [<http://broadinstitute.github.io/picard/>]. 30 December 2015.
- 684 24. **FASTQC** [<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>]. 30 December
685 2015.
- 686 25. Skinner ME, Uzilov AV, Stein LD, Mungall CJ, Holmes IH: **JBrowse: a next-generation**
687 **genome browser.** *Genome Res* 2009, **19**(9):1630-1638.
- 688 26. Quinlan AR: **BEDTools: The Swiss-Army Tool for Genome Feature Analysis.** *Curr*
689 *Protoc Bioinformatics* 2014, **47**:11 12 11-11 12 34.
- 690 27. **DISCVR-Seq** [<https://github.com/bbimber/discvr-seq/>]. 30 December 2015.
- 691 28. Nelson EK, Piehler B, Eckels J, Rauch A, Bellew M, Hussey P, Ramsay S, Nathe C,
692 Lum K, Krouse K *et al.* **LabKey Server: an open source platform for scientific data**
693 **integration, analysis and collaboration.** *BMC Bioinformatics* 2011, **12**:71.
- 694 29. Tong L, Thompson E: **Multilocus lod scores in large pedigrees: combination of**
695 **exact and approximate calculations.** *Hum Hered* 2008, **65**(3):142-153.

- 696 30. **Morgan** [<http://www.stat.washington.edu/thompson/Genepi/MORGAN/Morgan.shtml>]. 30
697 December 2015.
- 698 31. **GIGI: Genotype Imputation Given Inheritance**
699 [<https://faculty.washington.edu/wiisman/progdists/gigi/software/GIGI/GIGI.html>]. 30
700 December 2015.
- 701 32. Torkamaneh D, Belzile F: **Scanning and Filling: Ultra-Dense SNP Genotyping**
702 **Combining Genotyping-By-Sequencing, SNP Array and Whole-Genome**
703 **Resequencing Data.** *PloS one* 2015, **10**(7):e0131533.
- 704 33. **GIGI-Pick: Subject Selection for Sequencing in Pedigrees**
705 [[https://faculty.washington.edu/wiisman/progdists/gigi/software/GIGI-Pick/GIGI-](https://faculty.washington.edu/wiisman/progdists/gigi/software/GIGI-Pick/GIGI-Pick.html)
706 [Pick.html](https://faculty.washington.edu/wiisman/progdists/gigi/software/GIGI-Pick/GIGI-Pick.html)]. 30 December 2015.
- 707 34. Letvin NL, Rao SS, Montefiori DC, Seaman MS, Sun Y, Lim SY, Yeh WW, Asmal M,
708 Gelman RS, Shen L *et al*: **Immune and Genetic Correlates of Vaccine Protection**
709 **Against Mucosal Infection by SIV in Monkeys.** *Sci Transl Med* 2011, **3**(81):81ra36.
- 710 35. Nomura T, Matano T: **Association of MHC-I genotypes with disease progression in**
711 **HIV/SIV infections.** *Front Microbiol* 2012, **3**:234.
- 712 36. Ferguson B, Hunter JE, Luty J, Street SL, Woodall A, Grant KA: **Genetic load is**
713 **associated with hypothalamic-pituitary-adrenal axis dysregulation in macaques.**
714 *Genes Brain Behav* 2012, **11**(8):949-957.
- 715 37. Francis PJ, Appukuttan B, Simmons E, Landauer N, Stoddard J, Hamon S, Ott J,
716 Ferguson B, Klein M, Stout JT *et al*: **Rhesus monkeys and humans share common**
717 **susceptibility genes for age-related macular disease.** *Hum Mol Genet* 2008,
718 **17**(17):2673-2680.
- 719 38. Hartig W, Goldhammer S, Bauer U, Wegner F, Wirths O, Bayer TA, Grosche J:
720 **Concomitant detection of beta-amyloid peptides with N-terminal truncation and**
721 **different C-terminal endings in cortical plaques from cases with Alzheimer's**
722 **disease, senile monkeys and triple transgenic mice.** *J Chem Neuroanat* 2010,
723 **40**(1):82-92.
- 724 39. Lindell SG, Schwandt ML, Sun H, Sparenborg JD, Bjork K, Kasckow JW, Sommer WH,
725 Goldman D, Higley JD, Suomi SJ *et al*: **Functional NPY variation as a factor in stress**
726 **resilience and alcohol consumption in rhesus macaques.** *Arch Gen Psychiatry*
727 2010, **67**(4):423-431.
- 728 40. Spinelli S, Schwandt ML, Lindell SG, Heilig M, Suomi SJ, Higley JD, Goldman D, Barr
729 CS: **The serotonin transporter gene linked polymorphic region is associated with**
730 **the behavioral response to repeated stress exposure in infant rhesus macaques.**
731 *Dev Psychopathol* 2012, **24**(1):157-165.
- 732 41. Vallender EJ, Ruedi-Bettschen D, Miller GM, Platt DM: **A pharmacogenetic model of**
733 **naltrexone-induced attenuation of alcohol consumption in rhesus monkeys.** *Drug*
734 *Alcohol Depend* 2010, **109**(1-3):252-256.
- 735 42. Vinson A, Mitchell AD, Toffey D, Silver J, Raboin MJ: **Sex-specific heritability of**
736 **spontaneous lipid levels in an extended pedigree of Indian-origin rhesus**
737 **macaques (*Macaca mulatta*).** *PloS one* 2013, **8**(8):e72241.
- 738 43. Rhesus Macaque Genome S, Analysis C, Gibbs RA, Rogers J, Katze MG, Bumgarner
739 R, Weinstock GM, Mardis ER, Remington KA, Strausberg RL *et al*: **Evolutionary and**
740 **biomedical insights from the rhesus macaque genome.** *Science* 2007,
741 **316**(5822):222-234.
- 742 44. Fritsche LG, Igl W, Bailey JN, Grassmann F, Sengupta S, Bragg-Gresham JL, Burdon
743 KP, Hebring SJ, Wen C, Gorski M *et al*: **A large genome-wide association study of**
744 **age-related macular degeneration highlights contributions of rare and common**
745 **variants.** *Nat Genet* 2015.

- 746 45. Mancuso N, Rohland N, Rand KA, Tandon A, Allen A, Quinque D, Mallick S, Li H, Stram
747 A, Sheng X *et al*: **The contribution of rare variation to prostate cancer heritability.**
748 *Nat Genet* 2015.
- 749 46. Sasaki MM, Skol AD, Hungate EA, Bao R, Huang L, Kahn SA, Allan JM, Brant SR,
750 McGovern DP, Peter I *et al*: **Whole-exome Sequence Analysis Implicates Rare**
751 **IL17REL Variants in Familial and Sporadic Inflammatory Bowel Disease.** *Inflamm*
752 *Bowel Dis* 2016, **22**(1):20-27.
- 753 47. Yu B, Pulit SL, Hwang SJ, Brody JA, Amin N, Auer PL, Bis JC, Boerwinkle E, Burke GL,
754 Chakravarti A *et al*: **Rare Exome Sequence Variants in CLCN6 Reduce Blood**
755 **Pressure Levels and Hypertension Risk.** *Circ Cardiovasc Genet* 2015.
- 756

757 12) Figure titles and captions: Figures 1-5, uploaded separately. Figure titles and captions are
758 as follows:

759

760 **Figure 1. Pedigree diagram.** Pedigree diagram of the 16 subjects included in this study.

761 Subjects with whole genome sequence data are shown in gray; all subjects have GBS data.

762

763 **Figure 2. Evaluation of GBS library coverage and SNVs.** a) The number of positions with at
764 least 20X coverage; b) Total contiguous fragments with >20X coverage; c) Distance between
765 each GBS fragment and nearest predicted cut site for the *BglII* libraries (all fragments >400 bp
766 are grouped into a single bin); d) Distance between each GBS fragment and the nearest
767 predicted cut site for *PstI* libraries; e) Distance between high MAF (>0.25) SNVs in *BglII*; f)
768 Distance between high MAF SNVs in *PstI*; g) Total SNVs detected per enzyme, total that were
769 concordant with WGS data, and total SNVs with MAF >0.25.

770

771 **Figure 3. Imputation accuracy on chromosome 19, among different strategies for selecting 3**

772 individuals for WGS. Comparison of imputation accuracy among 3 different strategies for

773 selecting 3 individuals for WGS within the 16-member pedigree: “Bottom of Pedigree” (subjects

774 M, P, K), “Founders” (subjects B, C, D), and “GIGI-Pick” (B, H, J). Imputation of an optimal set

775 of dense markers was conducted for chromosome 19 from the 3 individuals with WGS, into the

776 13 recipient individuals with GBS data, using the GIGI imputation algorithm with the “Most-

777 Likely” (A) and “Threshold” (B) genotype calling methods. Circles indicate accuracy for each of
778 the 13 individuals; triangles indicate outlier individuals.

779

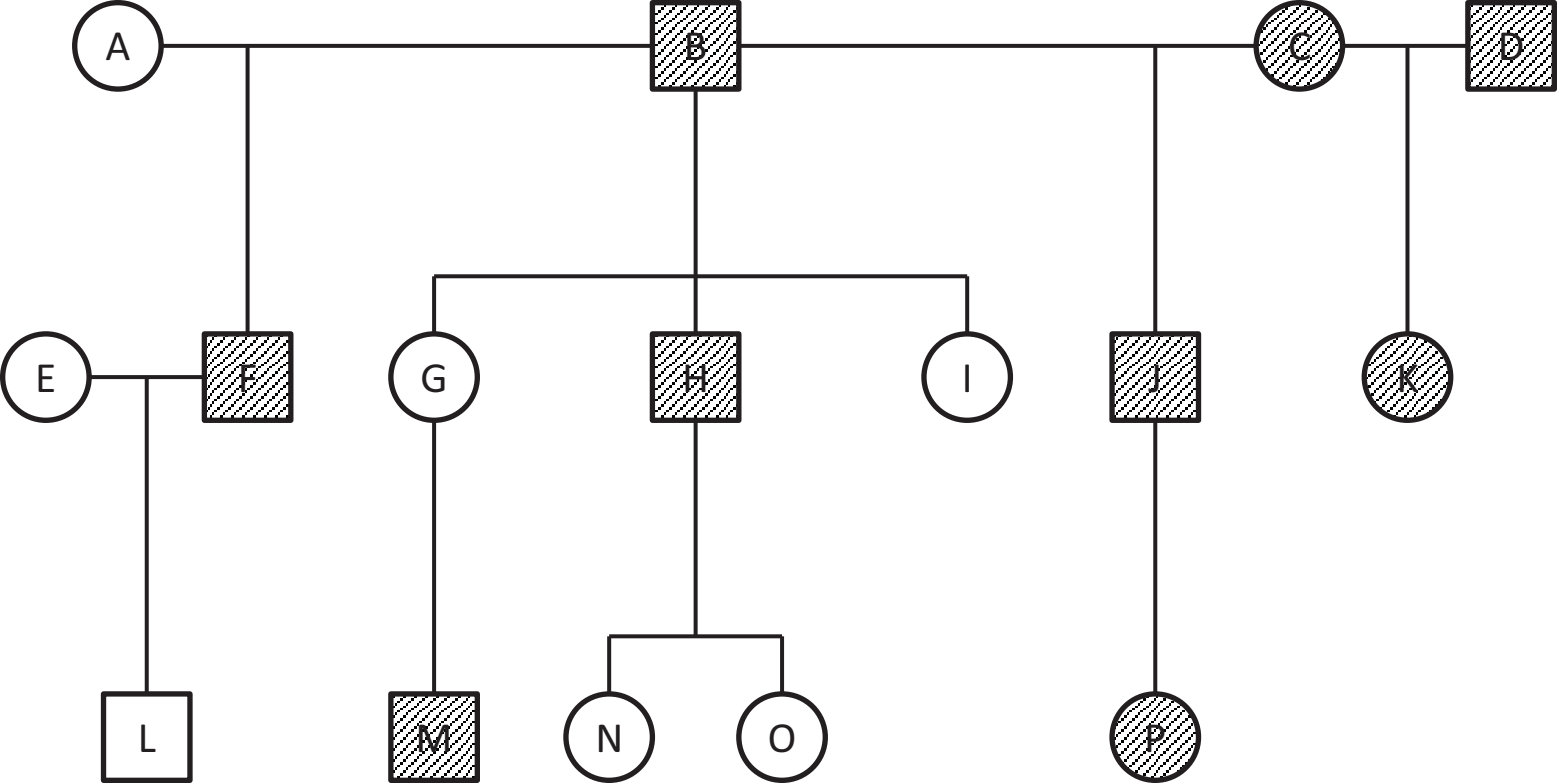
780 **Figure 4. Imputation accuracy on chromosome 19 by total number of individuals selected for**
781 **WGS.** Accuracy for an optimal set of dense markers on chromosome 19, using 1-9 individuals
782 with WGS data, imputed into all remaining pedigree members with GBS data , using the “most
783 likely” genotype calling method in the GIGI algorithm [13]. Circles indicate accuracy for each
784 individual; triangles indicate outliers. Individuals with WGS data were selected by the GIGI-Pick
785 algorithm [14] and used for imputation in the following order: B, H, J, F, M, K, P, C, D (see Fig.
786 1).

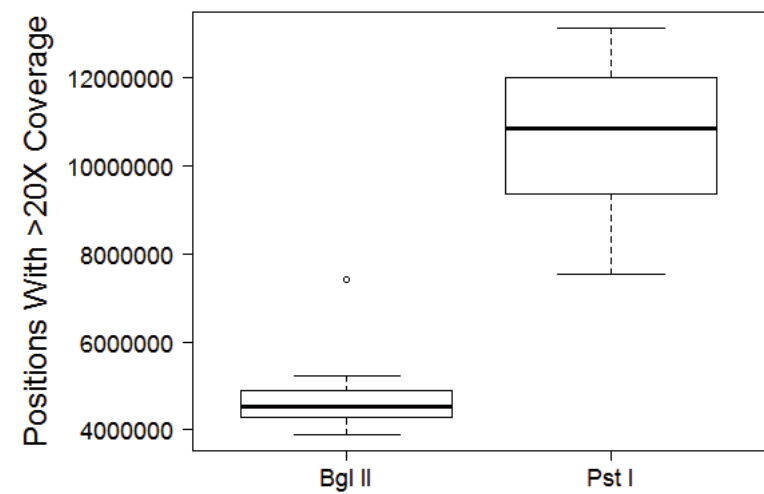
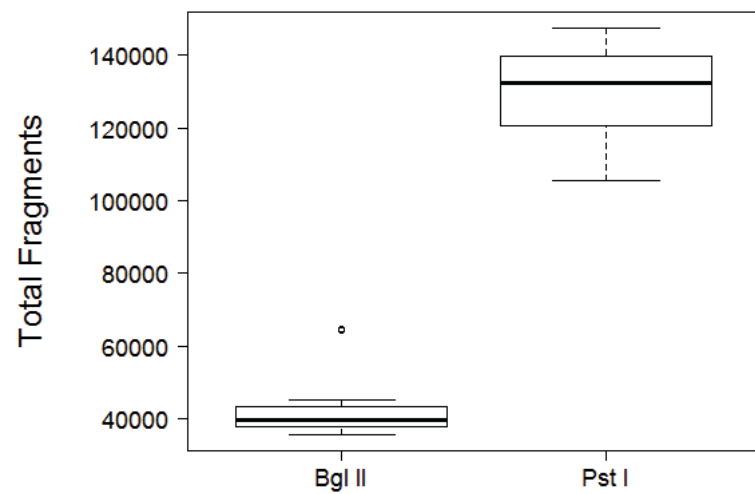
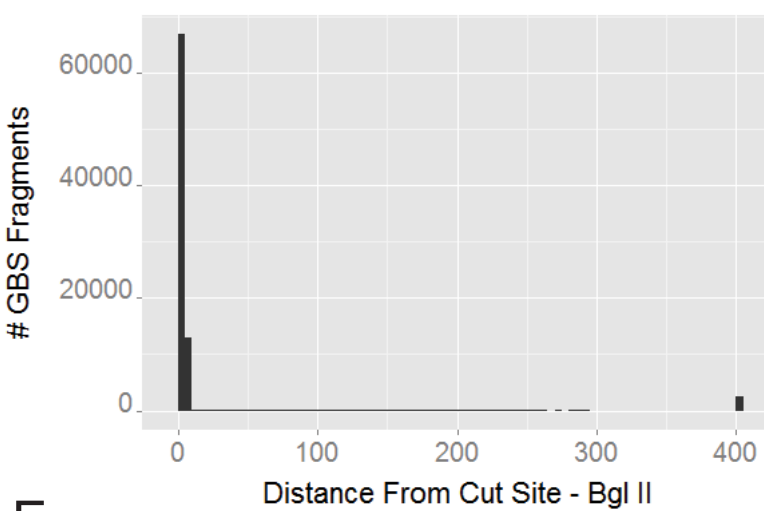
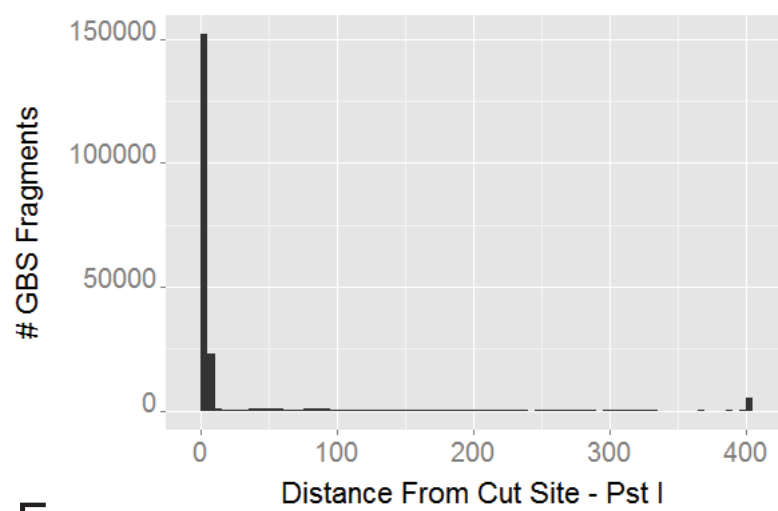
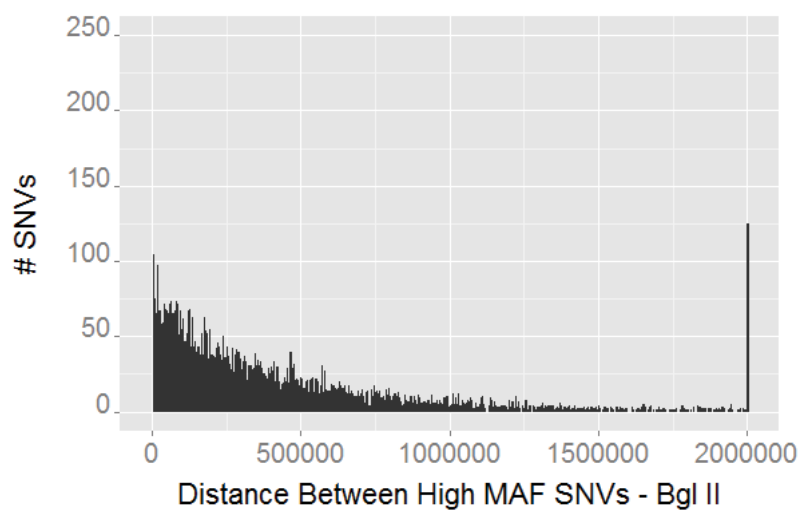
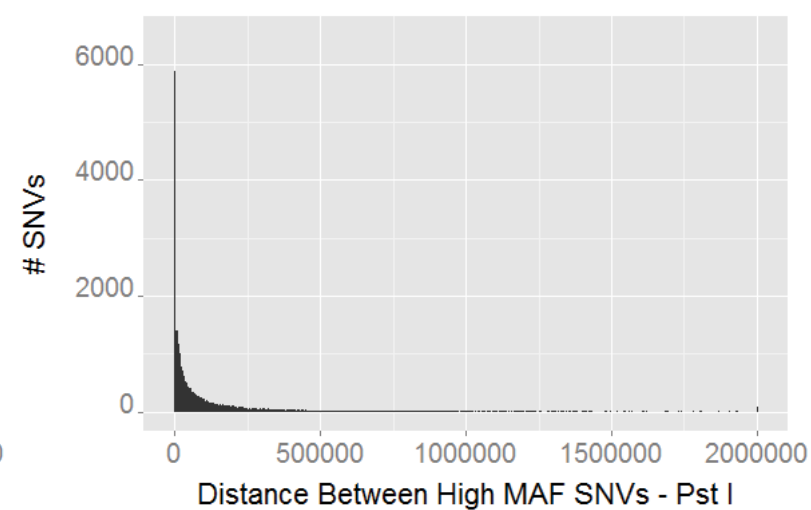
787

788 **Figure 5. Imputation accuracy across the genome at an expanded set of dense markers.**
789 Accuracy of alleles imputed across all autosomes at 14,384,988 dense markers. Data
790 represent accuracy of alleles imputed at dense markers for 12 pedigree members with GBS
791 data, imputed from individuals B, H, J, and F. Alleles were called using the “most likely”
792 genotype calling method in the GIGI algorithm [13]. Triangles indicate outliers.

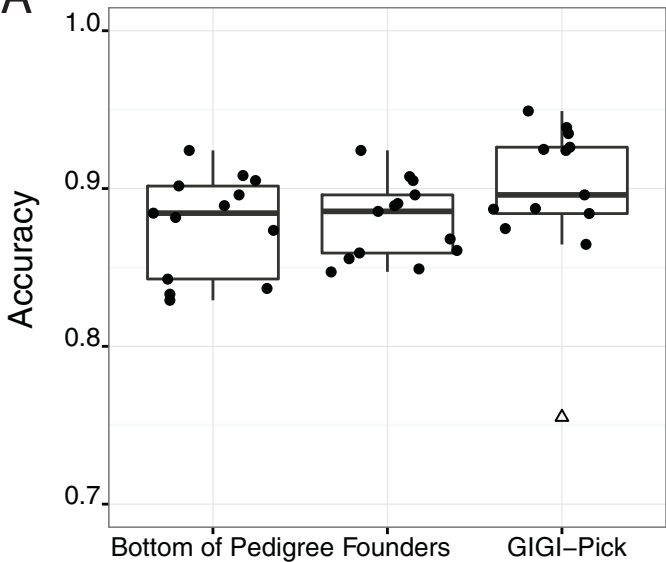
793

794 13) Tables and captions: Not applicable.



A**B****C****D****E****F****G**

	Bgl II	Pst I
Total SNVs	52,500	239,428
Concordant With WGS	50,557	231,287
High MAF SNVs	7,399	22,455

A**B**