

# UC San Diego

## UC San Diego Previously Published Works

### Title

Adapting Genotyping-by-Sequencing and Variant Calling for Heterogeneous Stock Rats

### Permalink

<https://escholarship.org/uc/item/6gq5k467>

### Journal

G3: Genes, Genomes, Genetics, 10(7)

### ISSN

2160-1836

### Authors

Gileta, Alexander F  
Gao, Jianjun  
Chitre, Apurva S  
[et al.](#)

### Publication Date

2020-07-01

### DOI

10.1534/g3.120.401325

Peer reviewed

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20

**Adapting genotyping-by-sequencing and variant calling for heterogeneous stock rats**

Alexander F. Gileta\* †<sup>¶</sup>, Jianjun Gao\*<sup>¶</sup>, Apurva S. Chitre\*, Hannah V. Bimschleger\*, Celine L. St. Pierre\*, Shyam Gopalakrishnan#, Abraham A. Palmer\*<sup>‡</sup>

\* Department of Psychiatry, University of California San Diego, La Jolla, California, 92093

† Department of Human Genetics, University of Chicago, Chicago, Illinois, 60637

# Natural History Museum of Denmark, University of Copenhagen, 2200 København N, Denmark

‡ Institute for Genomic Medicine, University of California San Diego, La Jolla, California, 92093

<sup>¶</sup> These authors contributed equally to this work.

21

22 **Running title:** GBS and variant calling in HS rats

23

24 **Key words:** genotyping-by-sequencing, heterogeneous stock, rat,  
25 imputation

26

27 **Corresponding author:** Abraham A. Palmer

28 **Mailing address:** 9500 Gilman Drive #0667, La Jolla, CA, 92093

29 **Phone number:** 858-534-2093

30 **Email:** [aapalmer@ucsd.edu](mailto:aapalmer@ucsd.edu)

31

## ABSTRACT

32 The heterogeneous stock (**HS**) is an outbred rat population derived from  
33 eight inbred rat strains. HS rats are ideally suited for genome wide  
34 association studies; however, only a few genotyping microarrays have ever  
35 been designed for rats and none of them are currently in production. To  
36 address the need for an efficient and cost effective method of genotyping HS  
37 rats, we have adapted genotype-by-sequencing (**GBS**) to obtain genotype  
38 information at large numbers of single nucleotide polymorphisms (**SNPs**). In  
39 this paper, we have outlined the laboratory and computational steps we took  
40 to optimize double digest genotype-by-sequencing (**ddGBS**) for use in rats.  
41 We also evaluate multiple existing computational tools and explain the  
42 workflow we have used to call and impute over 3.7 million SNPs. We also  
43 compared various rat genetic maps, which are necessary for imputation,  
44 including a recently developed map specific to the HS. Using our approach,  
45 we obtained concordance rates of 99% with data obtained using data from a  
46 genotyping array. The principles and computational pipeline that we describe  
47 could easily be adapted for use in other species for which reliable reference  
48 genome sets are available.

49

## INTRODUCTION

50 Advances in next-generation sequencing technology over the past decade  
51 have enabled the discovery of high-density, genome-wide single nucleotide  
52 polymorphisms (**SNPs**) in model systems. Comprehensive assays of the

53 standing genetic variation in these organisms has allowed for the  
54 identification of quantitative trait loci (**QTL**) and the application of numerous  
55 population genetic and phylogenetic methods. However, due to the high  
56 degree of linkage disequilibrium (**LD**) in QTL mapping populations,  
57 sequencing whole genomes is not necessary. Many populations are the result  
58 of numerous generations of interbreeding inbred strains, allowing for  
59 recombination to produce an admixed population with known founder  
60 haplotypes. Due to the relatively slow rate of accumulation of recombination  
61 events, these populations contain large chunks of the genome derived from  
62 the same founder haplotype. Nearby SNPs are therefore often in strong LD  
63 with physically adjacent loci, effectively 'tagging' nearby variation and  
64 thereby reducing the number of sites that need to be directly genotyped.  
65 Several reduced-representation sequencing approaches that take advantage  
66 of LD structure have been previously described (Miller *et al.* 2007; van  
67 Orsouw *et al.* 2007; Van Tassell *et al.* 2008; Baird *et al.* 2008; Huang *et al.*  
68 2009; Andolfatto *et al.* 2011; Elshire *et al.* 2011; Davey *et al.* 2011; Poland  
69 and Rife 2012; Peterson *et al.* 2012; Sun *et al.* 2013; Scheben *et al.* 2017).  
70 Using these methods, thousands of SNPs can be identified in large numbers  
71 of samples for a fraction of the price of whole-genome sequencing (Chen *et*  
72 *al.* 2013; He *et al.* 2014). The advantages of these methods are especially  
73 attractive when applied to less commonly utilized species or strains for which  
74 genotyping microarrays are not available.

75           Of the existing reduced-representation protocols, the genotyping-by-  
76 sequencing (**GBS**) approach developed by Elshire et al. (Elshire *et al.* 2011)  
77 has been frequently modified to accommodate other species: soybean  
78 (Sonah *et al.* 2013), rice (Furuta *et al.* 2017), oat (Fu and Yang 2017),  
79 chicken (Pértille *et al.* 2016; Wang *et al.* 2017), mouse (Parker *et al.* 2016),  
80 fox (Johnson *et al.* 2015), and cattle (De Donato *et al.* 2013), among others.  
81 The greatly varying genomic composition among organisms necessitates a  
82 diverse and customized set of approaches for obtaining high-quality  
83 genotypes. As such, both the GBS protocol and computational pipeline  
84 require modifications when applied to a new species. Recent work from our  
85 group showed that GBS can be effectively applied to outbred mice (Parker *et*  
86 *al.* 2016; Zhou *et al.* 2018; Gonzales *et al.* 2018) and rats (Fitzpatrick *et al.*  
87 2013). However, those publications used protocols that had not been  
88 optimized, leaving significant room for improvement in genotype quality and  
89 marker density. Additionally, although several tools and workflows for the  
90 analysis of GBS data have been described, including Stacks (Catchen *et al.*  
91 2013), IGST-GBS (Sonah *et al.* 2013), TASSEL-GBS (Glaubitz *et al.* 2014),  
92 Fast-GBS (Torkamaneh *et al.* 2017), and GB-eaSy (Wickland *et al.* 2017), the  
93 majority were developed and optimized for use in plant species. Given the  
94 lack of well-developed genomic resources in these species, they do not  
95 leverage the wealth of genomic data available for model organisms such as  
96 rats. Here we describe the customized computational and laboratory  
97 protocols for applying GBS to HS rats.

98           The HS is an outbred rat population created in 1984 using eight inbred  
99 strains and has been maintained since then with the goal of minimizing  
100 inbreeding and maximizing the genetic diversity of the colony (Johannesson  
101 *et al.* 2008; Woods and Mott 2017). After more than 80 generations of  
102 accumulated recombination events, their genome has become a fine-scale  
103 mosaic of the inbred founders' haplotypes. The breeding scheme and the  
104 number of accumulated generations has made the HS colony attractive for  
105 genetic studies. Additionally, extensive deep sequencing data exists for  
106 many inbred rat strains, including the eight progenitor strains (Rat Genome  
107 Sequencing and Mapping Consortium *et al.* 2013; Hermsen *et al.* 2015;  
108 Ramdas *et al.* 2019), allowing for accurate imputation to millions of  
109 additional SNPs following direct genotyping of only a subset.

110           Detailed here are the steps we have taken to optimize a rat GBS  
111 protocol and computational pipeline. Drawing on existing protocols (Elshire  
112 *et al.* 2011; Poland *et al.* 2012; Peterson *et al.* 2012; Parker *et al.* 2016) as  
113 templates, we redesigned our previous GBS approach (Parker *et al.* 2016;  
114 Gonzales *et al.* 2018) and have developed a novel, reference-based, high-  
115 throughput workflow to accurately and cost-effectively call and impute  
116 variants from low-coverage double digest GBS (**ddGBS**) data in HS rats. This  
117 publication is intended as a resource for others who might wish to perform  
118 GBS in rats and should provide a roadmap for adapting GBS for use in new  
119 species. We demonstrate that with a suitable reference panel, applying

120 reduced representation approaches and imputation in model systems can  
121 provide high-confidence genotypes on millions of genome-wide markers.

## 122 MATERIALS AND METHODS

### 123 **Tissue samples and DNA extraction**

124 Samples for this study originated from three sources: an in house advanced  
125 intercross line (**AIL**) derived from LG/J and SM/J mice (Gonzales *et al.* 2018),  
126 Sprague Dawley (**SD**) rats from Charles River Laboratories and Harlan  
127 Sprague Dawley, Inc. (Gileta *et al.* 2018), and an HS rat colony (Woods and  
128 Mott 2017; Chitre *et al.* 2018). Early stages of ddGBS optimization utilized  
129 AIL genomic DNA extracted from spleen by a standard salting-out protocol.  
130 Later optimization steps were performed using genomic DNA from SD rats  
131 extracted from tail tissue using the PureLink Genomic DNA Mini Kit (Thermo  
132 Fisher Scientific, Waltham, MA). HS rat DNA was extracted from spleen tissue  
133 using the Agencourt DNAdvance Kit (Beckman Coulter Life Sciences,  
134 Indianapolis, IN). All genomic DNA quality and purity was assessed by  
135 NanoDrop 8000 (Thermo Fisher Scientific, Waltham, MA). Interestingly, we  
136 observed that rat genomic DNA derived from either spleen or tail tissue  
137 appears to degrade faster than mouse genomic DNA following extraction by  
138 either of the above protocols; therefore, we recommend storing rat genomic  
139 DNA at -20° and using it within weeks of extraction whenever possible.

### 140 ***In silico* digest of rat genome**



141 We used *in silico* digests to aid in the selection of restriction enzymes, with  
142 the goal of maximizing the proportion of the genome captured at sufficient  
143 depth to make confident genotype calls (Kent *et al.* 2002). We used the  
144 *restrict* function in EMBOSS (version 6.6.0) (Rice *et al.* 2000) in conjunction  
145 with the REBASE database published by New England BioLabs (NEB; version  
146 808) (Roberts and Macelis 1999) to perform *in silico* digest of the current  
147 release of the Norway brown rat reference genome, designated rn6. For the  
148 primary restriction enzyme, we chose PstI, which had been successfully used  
149 in numerous project (Fitzpatrick *et al.* 2013; Parker *et al.* 2016; Gonzales *et*  
150 *al.* 2018). We performed the digest with PstI alone and then with PstI paired  
151 with each of 7 secondary enzymes: AluI, BfaI, DpnI, HaeIII, MluCI, MspI, and  
152 NlaIII. We only considered fragments with one PstI cut site and one cut site  
153 from the secondary enzyme because the adapter and primer sets are  
154 designed to only allow these fragments to be amplified.

### 155 **Restriction enzyme selection**

156 Initial criteria for selecting a secondary restriction enzyme were a 4bp  
157 recognition sequence, no ambiguity in the recognition sequence (i.e. N's),  
158 compatibility with the NEB CutSmart Buffer, and an incubation temperature  
159 of 37°C. The list of enzymes meeting these criteria at the time included AluI,  
160 BfaI, DpnI, HaeIII, MluCI, MspI, and NlaIII. Using the *in silico* digest data, we  
161 looked to maximize the portion of the genome contained within a fragment  
162 size range of 125-275bp (250-400bp with annealed adapters and primers)  
163 (Figure 2; Table 1). We excluded enzymes that produced blunt ends, both

164 because it would be more difficult to anneal adapters to blunt ended  
165 fragments and because our adapters would then also anneal to blunt ends  
166 produced by DNA shearing. We also excluded methylation-sensitive  
167 enzymes, as we did not want to limit the breadth of our sequencing efforts,  
168 accepting the possibility of read pileup in repetitive regions. Based on these  
169 criteria, as well as maximizing the percent of the genome captured, NlaIII,  
170 BfaI, and MluCI were selected for further testing. The final choice of enzyme  
171 (NlaIII) was determined empirically and is detailed in the Results.

172

### 173 **ddGBS library preparation and sequencing**

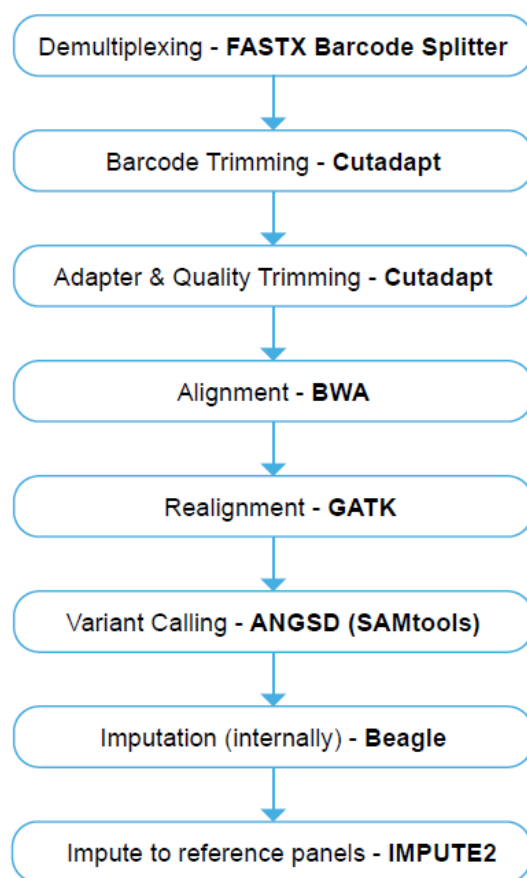
174 The full ddGBS protocol is available in File S1. In brief, approximately 1 $\mu$ g of  
175 DNA was used per sample. Sample DNA, PstI barcoded adapters, and NlaIII Y-  
176 adapter were combined in a 96-well plate and allowed to evaporate at 37°C  
177 overnight. The PstI adapter barcode is “in-line” such that each sequencing  
178 read from a given sample contains both the PstI overhang sequence (4bps)  
179 and a unique adapter sequence (4-8bps) prior to the beginning of the insert  
180 sequence. Sample DNA and adapters were re-eluted on day two with a  
181 PstI/NlaIII digestion mix and incubated at 37°C for two hours to allow for  
182 complete digestion. Ligation reagents were then added and incubated at  
183 16°C for one hour to anneal the adapters to the DNA fragments, followed by  
184 a 30-minute incubation at 80°C to inactivate the restriction enzymes. Sample  
185 libraries were purified using a plate from a MinElute 96 UF PCR Purification

186 Kit (QIAGEN Inc., Hilden, Germany), vacuum manifold, and ddH<sub>2</sub>O. Once re-  
187 eluted, libraries were quantified in duplicate with Quant-IT PicoGreen  
188 (Thermo Fisher Scientific, Waltham, MA) and pooled to the desired level of  
189 multiplexing (i.e. 12, 24, or 48 samples per library). Each pooled library was  
190 then concentrated by splitting the pooled volume across 2-3 wells of the  
191 MinElute vacuum plate and resuspending the library at desired volume for  
192 use in the Pippin Prep. The concentrated pool was quantified to ensure the  
193 gel cassette was not overloaded with DNA (>5µg). The pool was then loaded  
194 into the Pippin Prep for size selection (300-450bps) using a 2% agarose gel  
195 cassette on a Pippin Prep (Sage Science, Beverly, MA). Size-selected libraries  
196 were then PCR amplified for 12 cycles to add the Illumina sequencing  
197 primers and increase the quantity of DNA, concentrated again, and checked  
198 for quality on an Agilent 2100 Bioanalyzer with a DNA 1000 Series II chip  
199 (Agilent Technologies, Santa Clara, CA)., Bioanalyzer results were used to  
200 assure sufficient DNA concentration and to identify excessive primer dimer  
201 peaks.

202         As a pilot, an initial 96 HS samples were sequenced, 12 samples per  
203 library, at Beckman Coulter Genomics (now GENEWIZ) on an Illumina HiSeq  
204 2500 with v4 chemistry and 125bp single-end reads. Subsequently, we  
205 began using a set of 48 unique barcoded adapters (File S2) to multiplex 48  
206 HS samples per ddGBS library. Each library thereafter was run on a single  
207 flow cell lane on an Illumina HiSeq 4000 with 100bp single-end reads at the  
208 IGM Genomics Center (University of California San Diego, La Jolla, CA). We

209 also obtained ddGBS data on the HiSeq 4000 for a select set of 96 samples  
210 that had been previously genotyped on a custom Affymetrix Axiom MiRat  
211 625k microarray (Part#: 550572), providing us with a “gold standard” array-  
212 based dataset with which to compare to our ddGBS data.

213 **Figure 1. ddGBS sequencing data analysis workflow.** Each step of the  
214 workflow is described in the text.



215

---

## 217 **Evaluation of ddGBS pipeline performance**

218 We present the steps required to call and impute genotypes from raw ddGBS  
219 sequencing data in Figure 1. During optimization of the pipeline,  
220 performance was assessed by two primary metrics: (1) the number of  
221 variants called and (2) genotype concordance rates for calls made in 96 HS

222 rats that had both ddGBS genotypes and array genotypes from a custom  
223 Affymetrix Axiom MiRat 625k microarray. There were two checkpoints in the  
224 GBS pipeline where genotype quality (measured by concordance rate) was  
225 assessed. The first was after “internal” imputation with Beagle (Browning  
226 and Browning 2009, 2016), whereby we leverage information from samples  
227 that had sufficient read depth to make a confident genotype call at a given  
228 locus in order to impute the genotype of other samples that had lower read  
229 depths at that locus. The second checkpoint was after “external” imputation,  
230 meaning imputation to our reference panel with IMPUTE2 to obtain genotype  
231 calls at loci we did not directly capture by our GBS method. (Howie *et al.*  
232 2009, 2012). A third, additional metric we checked was the transition to  
233 transversion ratio ( $T_S T_V$ ), which is expected to be  $\sim 2$  for intergenic regions.  
234 The steps as outlined in the following sections reflect the final version of the  
235 pipeline. Variant calling and imputation steps utilized all available samples  
236 run on the HiSeq 4000 (3,000+ rats), though genotype concordance rates  
237 could only be calculated for the set of 96 HS samples for which we had array  
238 genotype calls.

### 239 **Demultiplexing**

240 The PstI adapter barcodes were used to demultiplex FASTQ files into  
241 individual sample files. Three demultiplexing software packages were tested:  
242 FASTX Barcode Splitter v0.0.13 [RRID: SCR\_005534] (Hannon Lab 2010),  
243 GBSX v1.3 (Herten *et al.* 2015), and an in-house Python script (Parker *et al.*  
244 2016). Reads that could not be matched with any barcode (maximum of 1

245 mismatch allowed), or that lacked the appropriate enzyme cut site, were  
246 discarded. Samples with less than two million reads after demultiplexing  
247 were discarded as these appeared to be outliers (Figure S4) and were  
248 observed to have high rates of missingness in their genotype calls. Data  
249 concerning demultiplexing are shown in Table S1 and are from a single HS  
250 rat sequenced in a 12-sample library on one lane after demultiplexing and  
251 adapter/quality trimming.

## 252 **Barcode, adapter, and quality trimming**

253 Read quality was assessed using FastQC v0.11.6 (Andrews 2017). We  
254 compared the efficacy of two rapid, lightweight software options for trimming  
255 barcodes, adapters, and low-quality bases from the NGS reads: Cutadapt  
256 v1.9.1 (Martin 2011) and the FASTX Clipper/Trimmer/Quality Trimmer tools  
257 v0.0.13 (Hannon Lab 2010) (Table S2). A base quality threshold of 20 was  
258 used and reads shorter than 25bp were discarded.

## 259 **Read alignment and indel realignment**

260 *Rattus norvegicus* genome assembly rn6 was used as the reference genome  
261 for read alignment with the Burrows-Wheeler Aligner v0.7.5a (BWA) [RRID:  
262 SCR\_010910] (Li and Durbin 2009) using the *mem* algorithm. We then used  
263 GATK IndelRealigner v3.5 [RRID: SCR001876] (McKenna *et al.* 2010) to  
264 improve alignment quality by locally realigning reads around a reference set  
265 of known indels in 42 whole-genome sequenced inbred rat strains, including  
266 the eight HS progenitor strains (Hermsen *et al.* 2015).

## 267 **Variant calling**

268 Variants were called, and genotype likelihoods were computed at variant  
269 sites using ANGSD v0.911, under the SAMtools model for genotype  
270 likelihoods (ANGSD-SAMtools) (Korneliussen *et al.* 2014; Durvasula *et al.*  
271 2016). Further, using ANGSD-SAMtools, we inferred the major and minor  
272 alleles (*-domajorminor* 1) from the genotype likelihoods, retaining only high  
273 confidence polymorphic sites (*-snp\_pval* 1e-6), and estimated the allele  
274 frequencies based on the inferred alleles (*-domaf* 1). We discarded sites  
275 missing read data in more than 4% of samples (*-minInd*). Additionally, we  
276 tested multiple thresholds for minimum base (*-minQ*) and mapping (*-*  
277 *minMapQ*) qualities.

## 278 **Internal imputation**

279 Beagle v4.1 (Browning and Browning 2009, 2016) was used to improve the  
280 genotyping within the samples without the use of an external reference  
281 panel. Missing and low quality genotypes were imputed by borrowing  
282 information from other individuals in the dataset with high quality  
283 information at these same variant sites. Before settling on the combination  
284 of ANGSD-SAMtools and Beagle for genotype calling and internal imputation,  
285 we also experimented with GATK's HaplotypeCaller (McKenna *et al.* 2010)  
286 with various parameter settings, but with unsatisfactory results (Figure 3).

## 287 **Quality Control for genotypes before imputation using an external** 288 **reference panel**

289 To verify the quality of the “internally” imputed genotypes prior to imputing  
290 SNPs from the 42 inbred strain reference panel (Hermsen *et al.* 2015), we  
291 checked concordance rates for the 96 HS animals with array genotypes,  
292 examined the  $T_sT_v$  ratio, and assessed whether the sex as recorded in the  
293 pedigree records agreed with the sex empirically determined by the  
294 proportion of reads on the X-chromosome out of the total number of reads  
295 (Figure S1). We also identified Mendelian errors using the --mendel option in  
296 *plink* and known pedigree information for 1,136 trios from 214 families within  
297 the HS sample. Using the fraction of the trios that were informative for a  
298 given SNP and the formula  $1-(1-2p(1-p))^3$ , where  $p$  represents the minor  
299 allele frequency of the allele, we formed curves for the distributions of the  
300 expected number of Mendelian errors for both SNPs and samples and chose  
301 the inflection points as thresholds for the number of Mendelian errors  
302 allowed.

### 303 **Data preparation for phasing with external reference panel**

304 First, in our study sample of 96 samples, we only retained variants  
305 previously identified in the 8 HS founder strains because we expected the  
306 polymorphisms in our samples to be limited to the variation present in the  
307 founders (Hermsen *et al.* 2015; Ramdas *et al.* 2019). Further, to improve  
308 imputation efficiency, we employed a pre-phasing step with IMPUTE2  
309 (*prephase\_g*) (Howie *et al.* 2012) prior to imputation. Pre-phasing only needs  
310 to be performed once, allowing us to reuse the estimated haplotypes from



311 our dataset for imputation with multiple different reference panels. A  
312 flowchart outlining the pre-phasing protocol is presented in Figure S2.

### 313 **Genetic maps**

314 Genetic maps are required for phasing and imputation with IMPUTE2. When  
315 we began this project, no strain-specific recombination map was available for  
316 HS rats. Thus, we considered a sparse genetic map for SHRSPxBN (Steen *et*  
317 *al.* 1999). We also tested two types of linearly interpolated genetic maps,  
318 with recombination rates set at either 1cM/Mb or the chromosome specific  
319 averages for rats, as reported by Jensen-Seaman *et al.* (Jensen-Seaman  
320 2004). Lastly, late in the course of this project, we experimented with an HS-  
321 specific genetic map developed by Littrell *et al.* (Littrell *et al.* 2018).

### 322 **Imputation to reference panel**

323 We used a combination of existing sequencing and array genotyping data  
324 from the HS rat founder and other inbred laboratory rat strains (Hermsen *et*  
325 *al.* 2015) as reference panel for imputation. Genotype data underwent QC  
326 and were phased by Beagle into single chromosome haplotype files.  
327 Haplotype files were then created using the workflow detailed in Figure S2.  
328 Imputation by IMPUTE2 was performed in 5Mb windows using the  
329 aforementioned reference panels and genetic maps.

### 330 **Data availability**

331 The ddGBS protocol and adapter sequences used to generate the data  
332 presented in this paper are available at  
333 <https://doi.org/10.6084/m9.figshare.12284432.v1>. All supplementary figures  
334 are available at <https://doi.org/10.6084/m9.figshare.12280814.v1>.  
335 Supplementary tables can be found at  
336 <https://doi.org/10.6084/m9.figshare.12284444.v1>. Genotype data will be  
337 available at <https://dx.doi.org/10.6084/m9.figshare.8243222>. The code  
338 necessary to run the steps of the computational pipeline outlined in this  
339 publication is available at <https://dx.doi.org/10.6084/m9.figshare.8243156>.  
340 Supplementary Files are available at  
341 <https://dx.doi.org/10.6084/m9.figshare.8243129>. Remaining files necessary  
342 for imputation (genetic maps, reference data, etc.) can be found with the  
343 following links: <https://dx.doi.org/10.6084/m9.figshare.11919615>,  
344 <https://dx.doi.org/10.6084/m9.figshare.11919573>, <https://dx.doi.org/10.6084/m9.figshare.11919597>.

## 346 RESULTS

### 347 **ddGBS optimization**

348 Previous projects utilizing GBS in mice and rats (Fitzpatrick *et al.* 2013;  
349 Parker *et al.* 2016; Gonzales *et al.* 2018) often encountered an issue where  
350 certain regions of the genome experienced high pileups of reads per sample  
351 (>100x), while other regions were covered by just 1-2 reads. This read  
352 distribution imbalance can be caused in part by PCR amplification bias,

353 where a subset of fragments are preferentially amplified until they dominate  
354 the final library (Kanagawa 2003; Aird *et al.* 2011). These previous protocols  
355 utilized 18 cycles of amplification. We tested reducing this to 8, 10, 12, or 14  
356 cycles and found that below 12 cycles, there was insufficient PCR product to  
357 accurately quantify and pool for sequencing. The reduction in the number of  
358 PCR cycles was expected to reduce PCR bias, though this was not explicitly  
359 tested.

360 Another concern regarding previous sequencing results was an excess  
361 of long fragments (>700bps as determined by *in silico* digest). We observed  
362 that longer sequencing fragments often do not provide sufficient reads to  
363 make confident genotype calls (< 5 reads per sample), putatively due to  
364 inefficient bridge amplification and clustering on Illumina flow cells.  
365 Sequencing these long fragments is therefore wasteful. We tested three  
366 methods of combating this issue, including increasing the digestion time or  
367 enzyme concentration, performing size selection on the libraries, and using a  
368 two-enzyme restriction digest.

369 We considered the possibility that the restriction enzyme digests might  
370 not be running to completion. To address this possibility, we increased the  
371 duration of the digestion from 2 hours to 3 or 4 hours. We also tried  
372 increasing the number of units of PstI enzyme added, to ensure complete  
373 digest. Neither of these modifications impacted the final fragment length  
374 distribution of the library, indicating that the digest was reaching completion  
375 after 2 hours using the original concentration of PstI (File S3 - wells 1-6).

376 Our previous GBS protocol did not have an explicit library fragment  
377 size selection step. The final library was purified using a MinElute PCR  
378 Purification Kit (QIAGEN Inc., Hilden, Germany), which isolates PCR products  
379 70bp-4kb in length, leaving a wide range of fragment sizes in the final  
380 library, under the assumption that only shorter fragments would bridge  
381 amplify on the flow cell. This method was imprecise and had low  
382 reproducibility, negatively impacting our ability obtain reads at consistent  
383 sites across libraries. Rather than attempt size selection by gel extraction,  
384 we chose to utilize a Pippin Prep, which automates the elution of DNA  
385 libraries of desired fragment size ranges. By using this automated size  
386 selection, we reduced the proportion of the genome targeted for sequencing,  
387 Additionally, since restriction enzymes make predominantly consistent cuts  
388 across samples (barring the presence of polymorphisms in RE recognition  
389 sites), it is ensured that highly similar sets of genomic fragments will be  
390 sequenced across sample libraries. Since the clustering process involves a  
391 bridge amplification step that preferentially amplifies library fragments with  
392 shorter insert sizes (Illumina, Inc. 2014), we kept the size selection window  
393 narrow (250-400bps) to avoid introducing a bias in which fragments were  
394 sequenced. A comparison of the fragment size distributions for the protocols  
395 before and after introduction of the Pippin Prep is shown in File S4.

396 To increase the proportion of the genome captured within the fragment  
397 size window, we pursued a double digest of the genome using a secondary  
398 enzyme with a more frequently occurring recognition sequence. When used

399 alone, *in silico* digest of the rn6 reference genome by PstI (Figure 2; Table 1)  
400 showed that only ~0.5% of the genome would have fallen within a 150bp  
401 fragment size window selected on the Pippin Prep. Previously, we performed  
402 GBS in CFW mice using the single-enzyme approach and observed that large  
403 regions of the genome that were not covered by sequencing reads (Parker *et*  
404 *al.* 2016). Therefore, we sought to increase the fraction of the genome that  
405 was accessible to GBS, so that there would be sufficient SNPs to tag the  
406 majority of the variation in the rat genome. Additionally, we were concerned  
407 about potential biases in coverage, heterozygosity, and the minor allele  
408 frequency (**MAF**) spectrum that may be introduced by a less complete  
409 capture of the genome. Flanagan and Jones have performed an empirical  
410 study comparing single- to double- digest RAD-seq and found that double-  
411 digest RAD-seq had lower rates of allelic dropout, decreased variance in  
412 between-sample per SNP coverage, less allele frequency inflation due to PCR  
413 bias, and reduced batch effects (Flanagan and Jones 2018).

414       The number of fragments with one of each of the cut sites was  
415 summed for all observed lengths and the results summarized in Figure 2 and  
416 Table 1. Bfal, MluCI, and NlaIII were chosen for further testing due to their  
417 compatibility with PstI digestion reagents and temperatures, sticky ends, and  
418 the proportion of the genome falling in the size selection window in the *in*  
419 *silico* analysis. We ruled out Bfal because it only had a 2bp overhang after  
420 cleavage, which we empirically showed leads to a high concentration of  
421 adapter dimer in the sequencing libraries (S5 File). NlaIII was chosen over

422 MluCI because it contained the greatest portion of the genome in the desired  
 423 size selection window.

424 **Table 1. Restriction enzyme options for double digest.**

Restriction Enzyme(s)	Recognition sequence	Length of Overhang (bp)	% Genome in 250-400bp Region <sup>+</sup>	% Genome in 300-450bp Region <sup>+</sup>
PstI	CTGCA <sup>^</sup> G	4	0.48%	0.56%
PstI + AluI	AG <sup>^</sup> CT	0	3.06%	2.88%
PstI + Bfal	C <sup>^</sup> TAG	2	3.10%	3.25%
PstI + DpnI*	GA <sup>^</sup> TC	0	2.69%	3.00%
PstI + HaeIII	GG <sup>^</sup> CC	0	2.71%	2.79%
PstI + MluCI	<sup>^</sup> AATT	4	3.32%	3.21%
PstI + MspI	C <sup>^</sup> CGG	2	1.16%	1.24%
PstI + NlaIII	CATG <sup>^</sup>	4	3.45%	3.31%

425 The percent genome in region columns indicate the percentage of the  
 426 genome that falls within the provided fragment size ranges and can  
 427 therefore be captured by GBS.

428

429 \* Restriction enzyme is methylation sensitive.

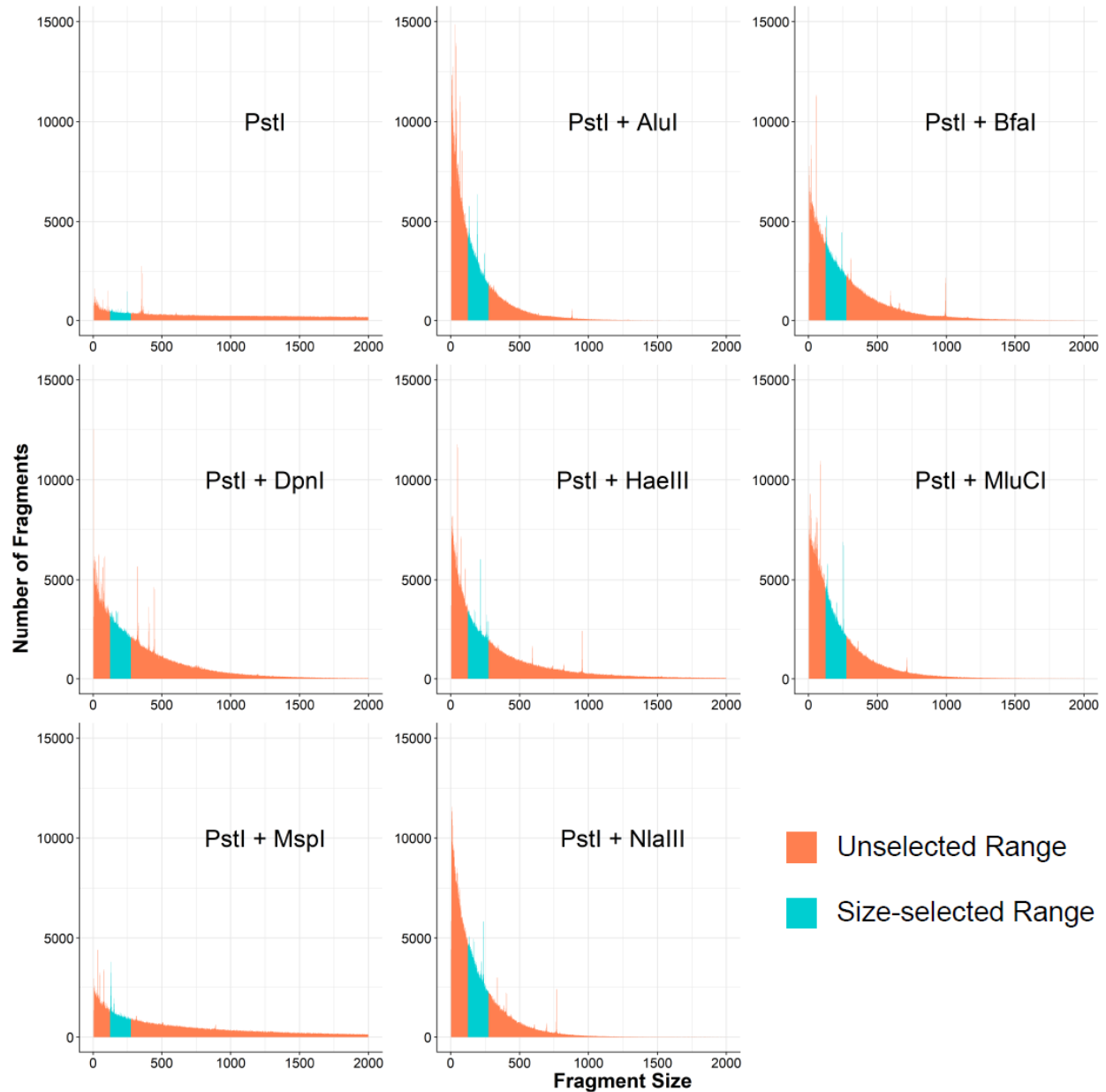
431 <sup>+</sup> Calculated using rn6 genome length of 2,870,182,909bps.

432

433

434 **Figure 2. *In silico* digest fragment distributions for PstI and**  
 435 **potential secondary restriction enzymes.**

436



437

438

439 Each panel represents an independent digest of rn6 with the listed  
 440 enzyme(s). Regions highlighted in blue are fragments that would be selected  
 441 by the Pippin Prep (125-275bp) after annealing adapters and primers. These  
 442 regions are quantified in Table 1 by multiplying the length of the fragments  
 443 by the number of fragments to estimate the portion of the genome captured.

444

445 In our previous GBS protocol, all fragments were cut on both ends by  
 446 PstI. By using a substantially lower concentration of the barcoded PstI

447 adapter than the common PstI adapter, we ensured the barcoded adapter  
448 would be the limiting reagent and the majority of fragments with an  
449 annealed barcoded adapter would have a common adapter on the other end.  
450 This is crucial, as having one of each of the adapters is required for proper  
451 amplification of the fragments on the flow cell. However, when using both  
452 PstI and NlaIII, the library is predominantly composed of fragments cut on  
453 both sides by NlaIII (File S6), which will amplify during PCR with a common  
454 adapter, but not on the flow cell. Therefore, we employed a Y-adapter  
455 (Poland *et al.* 2012) to control the direction of the first round of PCR and  
456 prevent two-sided NlaIII fragments from dominating the final sequencing  
457 library (File S2).

458         We tested numerous quantities of PstI and NlaIII adapters in an  
459 attempt minimize the amount used and avoid adapter dimers in the final  
460 libraries. For the barcoded PstI adapters, we tested 120pmol, 60pmol,  
461 20pmol, 4.0pmol, 2.67pmol, 1.60pmol, 0.53pmol, and 0.20pmol; for the NlaIII  
462 Y-adapter, 30pmol, 10pmol, 5.0pmol, 4.0pmol, and 1.0pmol (Files S7 & S8).  
463 We found that 0.20pmol of PstI adapter and 4pmol of NlaIII Y-adapter yielded  
464 sufficient library and minimized the presence of adapter dimers.

465         We sequenced a trial flow cell with 8 pooled ddGBS libraries of 12 SD  
466 rat samples each (96 total) on a HiSeq 2500 (Illumina, San Diego, CA) with  
467 125bp reads and v3 chemistry, obtaining an average of 15.3 million reads  
468 per sample. Given the NlaIII *in silico* digest results suggested we were  
469 capturing ~3.4% of the genome and that we were using 125bp reads, this



470 corresponded to approximately 20x coverage of captured sites. We  
471 subsequently increased the number of samples to 48 per library for the HS  
472 rats because we hypothesized 5x would be sufficient coverage per sample  
473 when utilizing imputation to a reference panel. We also discovered that a  
474 portion of the reads contained sequence fragments of the NlaIII adapter  
475 sequence, indicating there were fragments with insert sizes smaller than  
476 125bps in the final library. To avoid this, we increased the fragment size  
477 range to 300-450bps (Table 1), which corresponds to a 175-325bp insert size  
478 once the adapters and primers are accounted for. We noted however that  
479 the library size distribution obtained from the Pippin Prep was uniformly  
480 shifted towards higher fragment lengths (Figure S3). This is a result of the  
481 high concentrations of our libraries after pooling and loading the gel  
482 cartridge near the upper limit of the recommend number of micrograms of  
483 DNA, which can cause slower migration of the DNA across the gel matrix.

484 The final ddGBS protocol can be found in File S1 and the necessary  
485 primer and adapter sequences in File S2. This protocol was used for the  
486 sequencing of all HS rats included in subsequent computational optimization  
487 steps.

## 488 **Demultiplexing**

489 The number of base pairs of sequencing data retained after demultiplexing  
490 was fairly consistent across demultiplexing software (Table S1). We  
491 ultimately decided to use FASTX Barcode Splitter because it yielded the

492 greatest number of reads after quality/adaptor trimming and had faster run  
493 times. An average of 330 million 100bp reads were obtained per library,  
494 resulting in ~7 million reads per sample. Figure S4 shows the distribution of  
495 reads counts for all samples after demultiplexing.

#### 496 **Adapter and quality trimming**

497 Read quality was substantially improved after trimming the barcode and  
498 adaptor sequences and low-quality base pairs at the ends of reads (Figure  
499 S5). Overall read counts were only marginally reduced by quality trimming  
500 (Table S1). We observed that the number of called variant sites and the  
501 genotyping rate were both greater when using reads initially processed by  
502 cutadapt (Martin 2011) than reads processed by the FASTX\_Toolkit (Table  
503 S2). Importantly, a large portion of the additional identified variants were  
504 known variant sites from the 42 inbred strains reference set (Figure S6),  
505 indicating the elevated call rate was at least in part due to capturing more  
506 true variant sites. We viewed this as sufficient support for proceeding with  
507 cutadapt for adaptor removal and quality trimming.

#### 508 **Mapping quality**

509 The number of called variants and genotype call rates were identical at read  
510 mapping quality (mapQ) thresholds of either 20 or 30 (Table S3) within  
511 ANGSD. As the ANGSD mapQ threshold was raised to 45, there was a small  
512 reduction in the number of called variants, and then much greater losses at  
513 thresholds of 60 or 90. Fortunately, discordance rates between ddGBS and

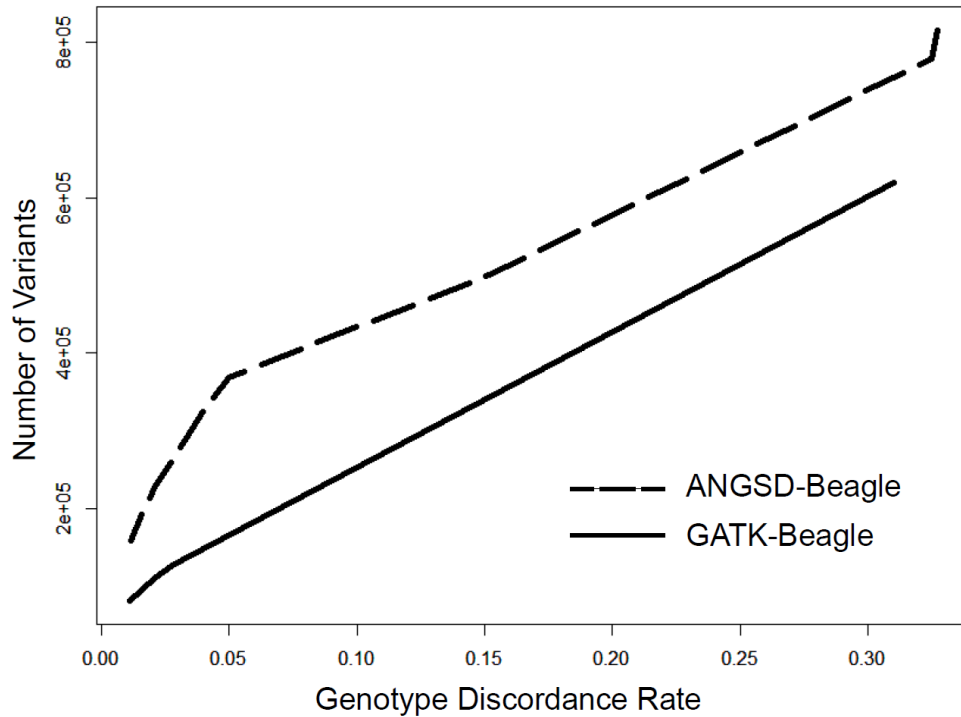
514 array genotypes were stable at both low and high mapQ thresholds, despite  
515 the putatively higher quality of the alignments (Figure S7). This permitted us  
516 to select a lower mapQ threshold (mapQ = 20), maximizing the number of  
517 variants called without sacrificing genotyping accuracy.

### 518 **Variant calling**

519 Figure 3 shows that across all levels of genotype discordance rates  
520 (comparing ddGBS with the array genotyping data), the combination of the  
521 ANGSD-SAMtools with BEAGLE produced more SNPs at various discordance  
522 thresholds than GATK's HaplotypeCaller (McKenna *et al.* 2010; DePristo *et al.*  
523 2011). This observation held when variants were limited only to biallelic sites  
524 and SNPs with an MAF > 0.05 (Figure S8). We speculate that the poorer  
525 performance of HaplotypeCaller may be due in part to the sparsity and non-  
526 uniform distribution of GBS genotype data across the genome and the high  
527 level of genotype call missingness across samples prior to imputation.

---

528  
529 **Figure 3. Genotype discordance rates between array data and**  
530 **variants called by GATK/Beagle or ANGSD-SAMtools/Beagle.**



531

532 The figure compares the number of variants called by combination of  
 533 ANGSD-SAMtools and Beagle or GATK HaplotypeCaller and Beagle at various  
 534 thresholds of genotype discordance with array data. Calls were made using  
 535 the 96 HS rats with array data. The x-axis represents the genotype  
 536 discordance rate thresholds and the y-axis is the number of variants that  
 537 surpass that threshold for each genotype calling method.

538

539 ANGSD supports four different models for estimating genotype  
 540 likelihoods: SAMtools, GATK, SOAPsnp and SYK. We compared the methods  
 541 to determine which produced the most SNPs with the lowest error rates. The  
 542 SOAPsnp model demonstrated an advantage in genotype accuracy and  
 543 number of variants called post-imputation with Beagle (Figure S9). However,  
 544 SOAPsnp requires considerably more time (1.7x for 96 samples) and memory  
 545 and scales poorly with sample size. With greater than 2,000 samples, we  
 546 were unable to allocate sufficient memory for the SOAPsnp model to

---

547 successfully run, even after dividing the chromosomes into several, smaller  
548 chunks. The marginal benefits of SOAPsnp in accuracy and number of  
549 variants were far outweighed by its limitations when applied to a large  
550 sample set. The GATK model showed a greater number of variants for more  
551 lenient genotype discordance rate threshold. This is in contrast with what  
552 was observed in Figures 3 and Figure S8 because ANGSD utilizes the direct  
553 genotype likelihood method from the first implementation of GATK's Unified  
554 Genotyper, whereas we had previously tested GATK's HaplotypeCaller.  
555 Interestingly though, as the stringency for discordance rate increased, the  
556 number of variants converged across the SAMtools, GATK, and SOAPsnp  
557 models. We proceeded with the SAMtools model for genotype likelihood  
558 estimation due to its previous support in the GBS literature (Torkamaneh *et*  
559 *al.* 2017), accepting a nominal decrease in highly concordant variants (Figure  
560 S9) for a large reduction in run time and memory usage.

### 561 **Imputation to reference panel**

562 Imputation is used in two complimentary ways in our protocol. As described  
563 earlier, after ddGBS, not all samples will have sufficient sequencing coverage  
564 at captured polymorphic loci to make a confident genotype call. Therefore,  
565 we first use imputation from other well-covered samples to “fill in the  
566 blanks” and assign genotypes to SNP loci in the subset of individuals that  
567 lacked confident calls at these sites. After these missing genotypes have  
568 been imputed in all samples, we then use the genotype information we have  
569 for the SNPs captured by ddGBS along with the reference panel data on the

570 original 8 HS founders (Hermsen *et al.* 2015; Ramdas *et al.* 2019) to impute  
571 genotype calls at sites that were inaccessible to ddGBS sequencing. Thus,  
572 our second application of imputation is similar to the human genetics  
573 application in which imputation using 1000 Genomes (1000 Genomes Project  
574 Consortium *et al.* 2010) increases the number of SNPs beyond those included  
575 on a given microarray platform. IMPUTE2 was selected over Beagle for this  
576 application because it has been shown to perform better with smaller  
577 reference panels from populations with substantial LD (Frischknecht *et al.*  
578 2014; Friedenberg and Meurs 2016)

579         Before starting this imputation step, we observed an inflated transition/  
580 transversion ratio (Table S4) in our ANGSD-SAMtools/Beagle SNPs. This issue  
581 was ameliorated when the SNP set was filtered for only “known” variants  
582 that were previously identified in either the 42 inbred strains (Hermsen *et al.*  
583 2015) or the 8 deep-sequenced HS founders (Ramdas *et al.* 2019). For  
584 imputation, we therefore only provided IMPUTE2 with previously identified  
585 variant sites from our ANGSD-SAMtools/Beagle output. Prior to running  
586 IMPUTE2, we also filtered the variants for biallelic sites with a genotype call  
587 in at least two individuals. Using pedigree data for the HS rats, we further  
588 removed samples showing an order of magnitude higher level of Mendelian  
589 error than the sample mean. We further removed SNPs that had an error rate  
590 surpassing a threshold of  $\sim 0.005$  (Figure S10; inflection point). There were 4  
591 samples and 4,179 SNPs removed from subsequent analyses. Lastly, we  
592 removed any samples where the X chromosome read ratio (reads mapped to

593 the X chromosome divided by total reads) was incompatible with their  
594 reported sex. We used hard threshold of 3% of total reads (empirically  
595 determined), where individuals with more than 3% X-mapped reads were  
596 determined to be female and below 3%, male (Figure S1).

597         There were three major genomic reference datasets available for the  
598 HS rats. The first reference set was obtained from Baud et al. (Rat Genome  
599 Sequencing and Mapping Consortium *et al.* 2013) and contained sequence  
600 data and genotype calls for the 8 founders of the HS. The second came from  
601 Hermsen et al. (Hermsen *et al.* 2015) which contains sequence and genotype  
602 data on 42 distinct laboratory rats strains and substrains, 8 of which were  
603 the founders of the HS from Baud et al., but analyzed alongside a new set of  
604 strains. The third reference set came from Ramdas et al. (Ramdas *et al.*  
605 2019), who independently performed whole-genome sequencing and made  
606 genotype calls on the 8 HS founder strains. It was unclear which set of  
607 genotypes would provide the best reference for imputation from our ddGBS  
608 data, so we tested five different possible subsets of available data (Table 2).  
609 From Hermsen et al., we used (1) all 42 inbred strains, (2) only the 34 strains  
610 that were not the HS founders, and (3) only the 8 HS founder strains. Then  
611 from Baud et al. and Ramdas et al., we tested the 8 HS founders only from  
612 each study. The most accurate imputation was observed for the reference  
613 set containing only the 8 deep-sequenced HS founder strains (Ramdas *et al.*  
614 2019); however, imputation to this set had the lowest genotyping rate of all  
615 panels. In contrast, using the 42 rat inbred strains displayed a balance of

616 high accuracy and low missingness, leading us to choose this as our  
 617 reference set. To better understand the role of the 8 founder strains, which  
 618 were part of the 42 strain reference panel, we created a reference panel that  
 619 included only the 34 non-HS founder strains. As expected, discordance rates  
 620 were much higher when only considering non-founders. However, the  
 621 genotype missingness was lower for the 34 than the 8 founders alone,  
 622 suggesting a combination of the two was the optimal set.

623

624

625

626

627

628

629

630 **Table 2. Imputation accuracy based on different variant reference**  
 631 **panels for IMPUTE2.**

632 The table includes five different possible reference panels for imputation.  
 633 The 42 inbred strains, 34 non-founder inbred strains, and 8 HS founders from  
 634 the 42 inbred strains all were derived from Hermsen et al. 2015 (Hermsen et  
 635 al. 2015). The UMich 8 HS founders were obtained from Ramdas et al. 2019  
 636 (Ramdas et al. 2019). The final set of 8 HS founder was taken from Baud et  
 637 al. 2013 (Rat Genome Sequencing and Mapping Consortium et al. 2013).

		<b>Chr1</b>	<b>Chr2</b>
<b>42 inbred strains</b>	<b>Discordance rate</b>	0.011	0.010
	<b># Variants</b>	790,659	882,993



	<b>Genotyping Rate</b>	0.85	0.81
<b>All 34 non-founder inbred strains</b>	<b>Discordance rate</b>	0.035	0.030
	<b># Variants</b>	812,550	912,749
	<b>Genotyping Rate</b>	0.84	0.80
<b>8 HS founders only from the 42 inbred strains</b>	<b>Discordance rate</b>	0.012	0.011
	<b># Variants</b>	805,424	902,061
	<b>Genotyping Rate</b>	0.57	0.53
<b>UMich 8 HS founders only</b>	<b>Discordance rate</b>	0.0059	0.008
	<b># Variants</b>	865,514	898,621
	<b>Genotyping Rate</b>	0.42	0.41
<b>Baud et. al 2013 8 HS founders only</b>	<b>Discordance rate</b>	0.0095	0.0096
	<b># Variants</b>	507,909	540,844
	<b>Genotyping Rate</b>	0.43	0.40

638

639 IMPUTE2 requires a genetic map. As described in the methods section,  
640 we considered four different genetic maps, two that were empirically derived  
641 and two that were linear extrapolations based on the physical map (Figure  
642 S11). All genetic maps performed similarly (Table S5). Surprisingly, the  
643 linear genetic maps performed just as well as the HS-specific map (Littrell *et*  
644 *al.* 2018). Thus, for simplicity, we chose to use the chromosome-specific  
645 values initially published by Jensen-Seaman (Jensen-Seaman 2004).

646 To obtain our final set of ~3.7 million variants, a final round of variant  
647 filtering was performed after imputation to the 42 strain reference panel. We  
648 removed SNPs with MAF < 0.005, a post-imputation genotyping rate < 90%,  
649 and SNPs that violated HWE with  $p < 1 \times 10^{-10}$ .

650

## DISCUSSION

651 The use of microarrays and WGS for genotyping large samples in model  
652 organisms remains cost-prohibitive. There is therefore an urgent and wide-  
653 spread need for high-performance and economical methods of obtaining  
654 genome-wide genotype data. While reduced-representation approaches have  
655 been utilized in numerous species of plants and animals, including rodents  
656 (Peterson *et al.* 2012; Fitzpatrick *et al.* 2013; Parker *et al.* 2016; Zhou *et al.*  
657 2018; Gonzales *et al.* 2018), there has yet to be a published protocol  
658 optimized specifically for rats. Prior to sequencing thousands of HS samples  
659 with GBS for our mapping efforts, we wanted to ensure we were capturing  
660 the greatest possible number of high-quality variants at the lowest possible  
661 cost. The protocol we present here is the culmination of careful testing and  
662 optimization of each step of the GBS protocol for rats. We have now applied  
663 the approach to 4,973 HS rats, as well as 4,608 Sprague Dawley rats (Gileta  
664 *et al.* 2018).

665 Our previous GBS protocol (Parker *et al.* 2016), which was designed for  
666 use with CFW mice, was unsuitable for our current genotyping efforts in HS  
667 rats, due to the much higher levels of genetic diversity in the HS population.  
668 There are multiple reasons we chose to develop our own computational  
669 pipeline for GBS rather than using existing workflows. Foremost, the  
670 prominent GBS analysis pipelines were developed and optimized for use in  
671 crop species (Sonah *et al.* 2013; Catchen *et al.* 2013; Glaubitz *et al.* 2014;  
672 Torkamaneh *et al.* 2017; Wickland *et al.* 2017), some of which are polyploid

673 and have differing levels of variation and LD than outbred rodent  
674 populations. Additionally, there were elements of each pipeline that did not  
675 meet our needs or lacked customizability. For instance, TASSEL-GBS v2  
676 (Glaubitz *et al.* 2014) trims all reads to 92 base pairs; however, other  
677 projects underway in our lab utilized up to 125bp reads, leading to a ~20%  
678 reduction in data. TASSEL-GBS also ignores read base quality scores, which  
679 are informative in probabilistic frameworks for estimating uncertainty in  
680 alignments and variant calls (Li *et al.* 2008; DePristo *et al.* 2011; Nielsen *et*  
681 *al.* 2011), and uses a naïve binomial likelihood ratio method for calling SNPs.  
682 Stacks has previously shown poor performance in demultiplexing (Herten *et*  
683 *al.* 2015; Torkamaneh *et al.* 2017) and does not make use of the reference  
684 genome for priors when calling SNPs (Catchen *et al.* 2013). Fast-GBS relies  
685 on Platypus (Rimmer *et al.* 2014) for variant calling (WGS500 Consortium *et*  
686 *al.* 2014; Torkamaneh *et al.* 2017), which employs a Bayesian method of  
687 constructing candidate haplotypes that works poorly with low-pass  
688 sequencing data and does not scale well to large sample sizes (Li *et al.*  
689 2018). Lastly, none of these pipelines included an imputation step, which is  
690 crucial for filling in missing genotypes in GBS data and can provide millions  
691 of additional SNPs given an appropriate composite reference panel (Howie *et*  
692 *al.* 2011; Huang and Tseng 2014).

693         Though we have not explicitly tested each alternate GBS pipeline for  
694 the purposes of this publication, this has been recently done by Wickland *et*  
695 *al.* (Wickland *et al.* 2017). Their pipeline GB-eaSy, which ours most closely

696 resembles, was found to be superior by a number of metrics to Stacks,  
697 TASSEL-GBS, IGST, and Fast-GBS. Similar to GB-eaSy, our pipeline utilizes a  
698 double-digest GBS protocol, aligns reads to the reference genome with *bwa*  
699 *mem*, and uses the SAMtools genotype likelihood model for calling SNPs (Li  
700 2011). The combination of *bwa mem* and SAMtools algorithm was  
701 independently shown to have the best performance for calling SNPs from  
702 Illumina data (Hwang *et al.* 2015), further supporting our choice of these  
703 programs for read alignment and variant calling. Additionally, using the  
704 ANGSD wrapper provided us with the ability to convert the posterior  
705 genotype probabilities into genotype dosages for mapping studies  
706 (Korneliussen *et al.* 2014).

707         A minor difference between GB-eaSy and our pipeline is the use of  
708 cutadapt (Martin 2011) rather than GBSX (Herten *et al.* 2015) for  
709 demultiplexing, though both performed equally well (Table S1). The primary  
710 improvement is our extension of the pipeline with the implementation of  
711 effective internal and reference-based imputation steps using the 42 inbred  
712 rat genomes (Hermsen *et al.* 2015) and 8 deep-sequenced HS founders from  
713 UMich (Ramdas *et al.* 2019). There are two stages of imputation in our  
714 pipeline. The first one is accomplished by Beagle and aims to fill in missing  
715 genotypes at called variants using information from other samples. This  
716 raises the genotype call rate to 100%, but it may also introduce errors due to  
717 insufficient information, emphasizing the need for careful filtering steps. The  
718 second stage of imputation made use of IMPUTE2 and an external reference

719 panels of variants called from WGS data on the 8 inbred HS founders, as well  
720 as 34 additional inbred rat strains. We decided to include the 34 additional  
721 strains because of the elevated genotyping rate we observed upon their  
722 inclusion in the IMPUTE2 reference panel. We attribute this to the presence  
723 of haplotypes that exist in both the 8 the HS founder strains and a subset of  
724 the 34 additional strains in this panel. The benefits of using a composite  
725 reference panel have been previously noted (Zhang *et al.* 2013; Huang and  
726 Tseng 2014); there is increased accuracy and decreased missingness in the  
727 imputed genotype data.

728         In summary, we have adapted a GBS protocol and genotyping and  
729 imputation pipeline to obtain dense genotypes on genome-wide markers in  
730 highly-multiplexed HS rats. After quality filtering on the level of SNP and  
731 sample, over 3.7 million SNPs were called with a concordance rate of 99%.  
732 The ddGBS protocol and bioinformatic methods used to produce this data are  
733 publicly available, easy to handle, and cost-effective. The presented  
734 workflow could be feasibly followed with marginal modifications for  
735 application in other species. The steps taken toward optimizing the wet lab  
736 protocols are easily applied to novel organisms, as is the computational  
737 pipeline so long as there are reliable reference genome sets available for use  
738 in alignment and imputation.

739

740

741

742

743

744

745

746

747

748

749

750

751

752

ACKNOWLEDGEMENTS:

753 This work was supported by P50 DA037844, R21 DA036672, T32 GM07197,

754 and F31 DA039638.

755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779

## LITERATURE CITED

1000 Genomes Project Consortium, G. R. Abecasis, D. Altshuler, A. Auton, L. D. Brooks *et al.*, 2010 A map of human genome variation from population-scale sequencing. *Nature* 467: 1061-1073.

Aird, D., M. G. Ross, W.-S. Chen, M. Danielsson, T. Fennell *et al.*, 2011 Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries. *Genome Biology* 12: R18.

Andolfatto, P., D. Davison, D. Erezyilmaz, T. T. Hu, J. Mast *et al.*, 2011 Multiplexed shotgun genotyping for rapid and efficient genetic mapping. *Genome Res.* 21: 610-617.

Andrews, S., 2017 *FastQC*.

Baird, N. A., P. D. Etter, T. S. Atwood, M. C. Currey, A. L. Shiver *et al.*, 2008 Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers (J. C. Fay, Ed.). *PLoS ONE* 3: e3376.

Browning, B. L., and S. R. Browning, 2009 A Unified Approach to Genotype Imputation and Haplotype-Phase Inference for Large Data Sets of Trios and Unrelated Individuals. *The American Journal of Human Genetics* 84: 210-223.

Browning, B. L., and S. R. Browning, 2016 Genotype Imputation with Millions of Reference Samples. *The American Journal of Human Genetics* 98: 116-126.

Catchen, J., P. A. Hohenlohe, S. Bassham, A. Amores, and W. A. Cresko, 2013 Stacks: an analysis tool set for population genomics. *Molecular Ecology* 22: 3124-3140.

Chen, Q., Y. Ma, Y. Yang, Z. Chen, R. Liao *et al.*, 2013 Genotyping by Genome Reducing and Sequencing for Outbred Animals (S. Zhao, Ed.). *PLoS ONE* 8: e67500.

780 Chitre, A. S., O. Polesskaya, K. Holl, J. Gao, R. Cheng *et al.*, 2018 Genome wide  
781 association study of body weight, body mass index, adiposity, and fasting  
782 glucose in 3,173 outbred rats.

783 Davey, J. W., P. A. Hohenlohe, P. D. Etter, J. Q. Boone, J. M. Catchen *et al.*, 2011  
784 Genome-wide genetic marker discovery and genotyping using next-  
785 generation sequencing. *Nature Reviews Genetics* 12: 499–510.

786 De Donato, M., S. O. Peters, S. E. Mitchell, T. Hussain, and I. G. Imumorin, 2013  
787 Genotyping-by-Sequencing (GBS): A Novel, Efficient and Cost-Effective  
788 Genotyping Method for Cattle Using Next-Generation Sequencing (J. C.  
789 Nelson, Ed.). *PLoS ONE* 8: e62137.

790 DePristo, M. A., E. Banks, R. Poplin, K. V. Garimella, J. R. Maguire *et al.*, 2011 A  
791 framework for variation discovery and genotyping using next-generation DNA  
792 sequencing data. *Nature Genetics* 43: 491–498.

793 Durvasula, A., P. J. Hoffman, T. V. Kent, C. Liu, T. J. Y. Kono *et al.*, 2016 ANGSD -  
794 wrapper: utilities for analysing next-generation sequencing data. *Mol Ecol*  
795 *Resour* 16: 1449–1454.

796 Elshire, R. J., J. C. Glaubitz, Q. Sun, J. A. Poland, K. Kawamoto *et al.*, 2011 A Robust,  
797 Simple Genotyping-by-Sequencing (GBS) Approach for High Diversity Species  
798 (L. Orban, Ed.). *PLoS ONE* 6: e19379.

799 Fitzpatrick, C. J., S. Gopalakrishnan, E. S. Cogan, L. M. Yager, P. J. Meyer *et al.*, 2013  
800 Variation in the Form of Pavlovian Conditioned Approach Behavior among  
801 Outbred Male Sprague-Dawley Rats from Different Vendors and Colonies:  
802 Sign-Tracking vs. Goal-Tracking (P. Campolongo, Ed.). *PLoS ONE* 8: e75042.



803 Flanagan, S. P., and A. G. Jones, 2018 Substantial differences in bias between  
804 single-digest and double-digest RAD-seq libraries: A case study. *Molecular*  
805 *Ecology Resources* 18: 264–280.

806 Friedenber, S. G., and K. M. Meurs, 2016 Genotype imputation in the domestic dog.  
807 *Mamm. Genome* 27: 485–494.

808 Frischknecht, M., M. Neuditschko, V. Jagannathan, C. Drögemüller, J. Tetens *et al.*,  
809 2014 Imputation of sequence level genotypes in the Franches-Montagnes  
810 horse breed. *Genet. Sel. Evol.* 46: 63.

811 Fu, Y.-B., and M.-H. Yang, 2017 Genotyping-by-Sequencing and Its Application to  
812 Oat Genomic Research, pp. 169–187 in *Oat*, edited by S. Gasparis. Springer  
813 New York, New York, NY.

814 Furuta, T., M. Ashikari, K. K. Jena, K. Doi, and S. Reuscher, 2017 Adapting  
815 Genotyping-by-Sequencing for Rice F2 Populations. *Genes|*  
816 *Genomes|Genetics* 7: 881–893.

817 Gileta, A. F., C. J. Fitzpatrick, A. S. Chitre, C. L. St. Pierre, E. V. Joyce *et al.*, 2018  
818 Genetic characterization of outbred Sprague Dawley rats and utility for  
819 genome-wide association studies.

820 Glaubitz, J. C., T. M. Casstevens, F. Lu, J. Harriman, R. J. Elshire *et al.*, 2014 TASSEL-  
821 GBS: A High Capacity Genotyping by Sequencing Analysis Pipeline (N. A.  
822 Tinker, Ed.). *PLoS ONE* 9: e90346.

823 Gonzales, N. M., J. Seo, A. I. Hernandez-Cordero, C. L. St. Pierre, J. S. Gregory *et al.*,  
824 2018 Genome wide association analysis in a mouse advanced intercross line.  
825 Hannon Lab, 2010 *FASTX-Toolkit*.

826 He, J., X. Zhao, A. Laroche, Z.-X. Lu, H. Liu *et al.*, 2014 Genotyping-by-sequencing  
827 (GBS), an ultimate marker-assisted selection (MAS) tool to accelerate plant  
828 breeding. *Frontiers in Plant Science* 5:.

829 Hermsen, R., J. de Ligt, W. Spee, F. Blokzijl, S. Schäfer *et al.*, 2015 Genomic  
830 landscape of rat strain and substrain variation. *BMC Genomics* 16:.

831 Herten, K., M. S. Hestand, J. R. Vermeesch, and J. K. Van Houdt, 2015 GBSX: a toolkit  
832 for experimental design and demultiplexing genotyping by sequencing  
833 experiments. *BMC Bioinformatics* 16:.

834 Howie, B. N., P. Donnelly, and J. Marchini, 2009 A Flexible and Accurate Genotype  
835 Imputation Method for the Next Generation of Genome-Wide Association  
836 Studies (N. J. Schork, Ed.). *PLoS Genetics* 5: e1000529.

837 Howie, B., C. Fuchsberger, M. Stephens, J. Marchini, and G. R. Abecasis, 2012 Fast  
838 and accurate genotype imputation in genome-wide association studies  
839 through pre-phasing. *Nature Genetics* 44: 955.

840 Howie, B., J. Marchini, and M. Stephens, 2011 Genotype Imputation with Thousands  
841 of Genomes. *G3&#58; Genes|Genomes|Genetics* 1: 457-470.

842 Huang, X., Q. Feng, Q. Qian, Q. Zhao, L. Wang *et al.*, 2009 High-throughput  
843 genotyping by whole-genome resequencing. *Genome Research* 19: 1068-  
844 1076.

845 Huang, G.-H., and Y.-C. Tseng, 2014 Genotype imputation accuracy with different  
846 reference panels in admixed populations. *BMC Proceedings* 8: S64.

847 Hwang, S., E. Kim, I. Lee, and E. M. Marcotte, 2015 Systematic comparison of  
848 variant calling pipelines using gold standard personal exome variants.  
849 *Scientific Reports* 5: 17875.

850 Illumina, Inc., 2014 Nextera(R) Library Validation and Cluster Density Optimization:  
851 Guidelines for generating high-quality data with Nextera library preparation  
852 kits.

853 Jensen-Seaman, M. I., 2004 Comparative Recombination Rates in the Rat, Mouse,  
854 and Human Genomes. *Genome Research* 14: 528-538.

855 Johannesson, M., R. Lopez-Aumatell, P. Stridh, M. Diez, J. Tuncel *et al.*, 2008 A  
856 resource for the simultaneous high-resolution mapping of multiple  
857 quantitative trait loci in rats: The NIH heterogeneous stock. *Genome*  
858 *Research* 19: 150-158.

859 Johnson, J. L., H. Wittgenstein, S. E. Mitchell, K. E. Hyma, S. V. Temnykh *et al.*, 2015  
860 Genotyping-By-Sequencing (GBS) Detects Genetic Structure and Confirms  
861 Behavioral QTL in Tame and Aggressive Foxes (*Vulpes vulpes*) (W. J. Murphy,  
862 Ed.). *PLOS ONE* 10: e0127013.

863 Kanagawa, T., 2003 Bias and artifacts in multitemplate polymerase chain reactions  
864 (PCR). *Journal of Bioscience and Bioengineering* 96: 317-323.

865 Kent, W. J., C. W. Sugnet, T. S. Furey, K. M. Roskin, T. H. Pringle *et al.*, 2002 The  
866 Human Genome Browser at UCSC. *Genome Research* 12: 996-1006.

867 Korneliussen, T. S., A. Albrechtsen, and R. Nielsen, 2014 ANGSD: Analysis of Next  
868 Generation Sequencing Data. *BMC Bioinformatics* 15:.

869 Li, H., 2011 A statistical framework for SNP calling, mutation discovery, association  
870 mapping and population genetical parameter estimation from sequencing  
871 data. *Bioinformatics* 27: 2987-2993.

872 Li, H., and R. Durbin, 2009 Fast and accurate short read alignment with Burrows-  
873 Wheeler transform. *Bioinformatics* 25: 1754-1760.

874 Li, H., J. Ruan, and R. Durbin, 2008 Mapping short DNA sequencing reads and calling  
875 variants using mapping quality scores. *Genome Research* 18: 1851-1858.

876 Li, Z., Y. Wang, and F. Wang, 2018 A study on fast calling variants from next-  
877 generation sequencing data using decision tree. *BMC Bioinformatics* 19:.

878 Littrell, J., S.-W. Tsaih, A. Baud, P. Rastas, L. Solberg-Woods *et al.*, 2018 A High-  
879 Resolution Genetic Map for the Laboratory Rat. *G3&#58; Genes|*  
880 *Genomes|Genetics* g3.200187.2018.

881 Martin, M., 2011 Cutadapt removes adapter sequences from high-throughput  
882 sequencing reads. *EMBnet.journal* 17: 10.

883 McKenna, A., M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis *et al.*, 2010 The  
884 Genome Analysis Toolkit: A MapReduce framework for analyzing next-  
885 generation DNA sequencing data. *Genome Research* 20: 1297-1303.

886 Miller, M. R., J. P. Dunham, A. Amores, W. A. Cresko, and E. A. Johnson, 2007 Rapid  
887 and cost-effective polymorphism identification and genotyping using  
888 restriction site associated DNA (RAD) markers. *Genome Research* 17: 240-  
889 248.

890 Nielsen, R., J. S. Paul, A. Albrechtsen, and Y. S. Song, 2011 Genotype and SNP  
891 calling from next-generation sequencing data. *Nature Reviews Genetics* 12:  
892 443-451.

893 van Orsouw, N. J., R. C. J. Hogers, A. Janssen, F. Yalcin, S. Snoeijers *et al.*, 2007  
894 Complexity Reduction of Polymorphic Sequences (CRoPS™): A Novel  
895 Approach for Large-Scale Polymorphism Discovery in Complex Genomes (I.  
896 Baxter, Ed.). *PLoS ONE* 2: e1172.

897 Parker, C. C., S. Gopalakrishnan, P. Carbonetto, N. M. Gonzales, E. Leung *et al.*,  
898 2016 Genome-wide association study of behavioral, physiological and gene  
899 expression traits in outbred CFW mice. *Nature Genetics* 48: 919–926.

900 Pértille, F., C. Guerrero-Bosagna, V. H. da Silva, C. Boschiero, J. de R. da S. Nunes *et*  
901 *al.*, 2016 High-throughput and Cost-effective Chicken Genotyping Using Next-  
902 Generation Sequencing. *Scientific Reports* 6: 26929.

903 Peterson, B. K., J. N. Weber, E. H. Kay, H. S. Fisher, and H. E. Hoekstra, 2012 Double  
904 Digest RADseq: An Inexpensive Method for De Novo SNP Discovery and  
905 Genotyping in Model and Non-Model Species (L. Orlando, Ed.). *PLoS ONE* 7:  
906 e37135.

907 Poland, J. A., P. J. Brown, M. E. Sorrells, and J.-L. Jannink, 2012 Development of High-  
908 Density Genetic Maps for Barley and Wheat Using a Novel Two-Enzyme  
909 Genotyping-by-Sequencing Approach (T. Yin, Ed.). *PLoS ONE* 7: e32253.

910 Poland, J. A., and T. W. Rife, 2012 Genotyping-by-Sequencing for Plant Breeding and  
911 Genetics. *The Plant Genome Journal* 5: 92.

912 Ramdas, S., A. B. Ozel, M. K. Treutelaar, K. Holl, M. Mandel *et al.*, 2019 Extended  
913 regions of suspected mis-assembly in the rat reference genome. *Sci Data* 6:  
914 39.

915 Rat Genome Sequencing and Mapping Consortium, A. Baud, R. Hermsen, V. Guryev,  
916 P. Stridh *et al.*, 2013 Combined sequence-based and genetic mapping  
917 analysis of complex traits in outbred rats. *Nature Genetics* 45: 767–775.

918 Rice, P., I. Longden, and A. Bleasby, 2000 EMBOSS: The European Molecular Biology  
919 Open Software Suite. *Trends in Genetics* 16: 276–277.

920 Rimmer, A., H. Phan, I. Mathieson, Z. Iqbal, S. R. F. Twigg *et al.*, 2014 Integrating  
921 mapping-, assembly- and haplotype-based approaches for calling variants in  
922 clinical sequencing applications. *Nat. Genet.* 46: 912–918.

923 Roberts, R. J., and D. Macelis, 1999 REBASE--restriction enzymes and methylases.  
924 *Nucleic Acids Research* 27: 312–313.

925 Scheben, A., J. Batley, and D. Edwards, 2017 Genotyping-by-sequencing approaches  
926 to characterize crop genomes: choosing the right tool for the right  
927 application. *Plant Biotechnology Journal* 15: 149–161.

928 Sonah, H., M. Bastien, E. Iquira, A. Tardivel, G. Légaré *et al.*, 2013 An Improved  
929 Genotyping by Sequencing (GBS) Approach Offering Increased Versatility and  
930 Efficiency of SNP Discovery and Genotyping (Z. Liu, Ed.). *PLoS ONE* 8:  
931 e54603.

932 Steen, R. G., A. E. Kwikteck-Black, C. Glenn, J. Gullings-Handley, W. Van Etten *et al.*,  
933 1999 A high-density integrated genetic linkage and radiation hybrid map of  
934 the laboratory rat. *Genome Res.* 9: AP1-8, insert.

935 Sun, X., D. Liu, X. Zhang, W. Li, H. Liu *et al.*, 2013 SLAF-seq: An Efficient Method of  
936 Large-Scale De Novo SNP Discovery and Genotyping Using High-Throughput  
937 Sequencing (J. Aerts, Ed.). *PLoS ONE* 8: e58700.

938 Torkamaneh, D., J. Laroche, M. Bastien, A. Abed, and F. Belzile, 2017 Fast-GBS: a  
939 new pipeline for the efficient and highly accurate calling of SNPs from  
940 genotyping-by-sequencing data. *BMC Bioinformatics* 18:.

941 Van Tassell, C. P., T. P. L. Smith, L. K. Matukumalli, J. F. Taylor, R. D. Schnabel *et al.*,  
942 2008 SNP discovery and allele frequency estimation by deep sequencing of  
943 reduced representation libraries. *Nature Methods* 5: 247–252.

944 Wang, Y., X. Cao, Y. Zhao, J. Fei, X. Hu *et al.*, 2017 Optimized double-digest  
945 genotyping by sequencing (ddGBS) method with high-density SNP markers  
946 and high genotyping accuracy for chickens (P. Xu, Ed.). PLOS ONE 12:  
947 e0179073.

948 WGS500 Consortium, A. Rimmer, H. Phan, I. Mathieson, Z. Iqbal *et al.*, 2014  
949 Integrating mapping-, assembly- and haplotype-based approaches for calling  
950 variants in clinical sequencing applications. Nature Genetics 46: 912–918.

951 Wickland, D. P., G. Battu, K. A. Hudson, B. W. Diers, and M. E. Hudson, 2017 A  
952 comparison of genotyping-by-sequencing analysis methods on low-coverage  
953 crop datasets shows advantages of a new workflow, GB-eaSy. BMC  
954 Bioinformatics 18:.

955 Woods, L. C. S., and R. Mott, 2017 Heterogeneous Stock Populations for Analysis of  
956 Complex Traits, pp. 31–44 in *Systems Genetics*, edited by K. Schughart and R.  
957 W. Williams. Springer New York, New York, NY.

958 Zhang, P., X. Zhan, N. A. Rosenberg, and S. Zöllner, 2013 Genotype imputation  
959 reference panel selection using maximal phylogenetic diversity. Genetics  
960 195: 319–330.

961 Zhou, X., C. L. St. Pierre, N. M. Gonzales, R. Cheng, A. S. Chitre *et al.*, 2018 Genome-  
962 wide association study, replication, and mega-analysis using a dense marker  
963 panel in a multi-generational mouse advanced intercross line.

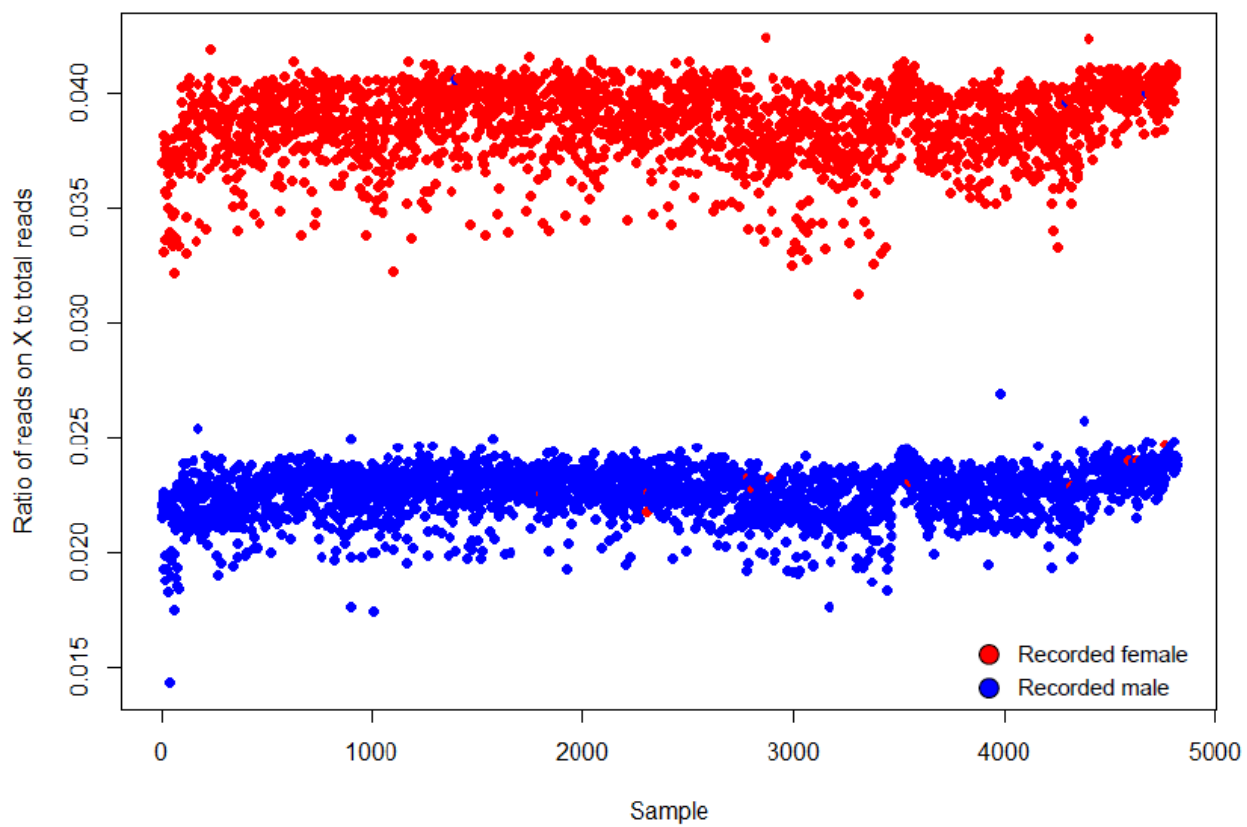
964

---

965 **Figure S1. Ratio of reads on X-chromosome to total sequencing**  
966 **reads.**

967 The color of the points indicates the pedigree-recorded sex of the samples.  
968 Females are expected to have approximately twice as many reads for the X-  
969 chromosome. Samples that did not cluster with their pedigree-recorded sex  
970 were removed from the study for possible sample mix-up.

971



972

973

974

975

976

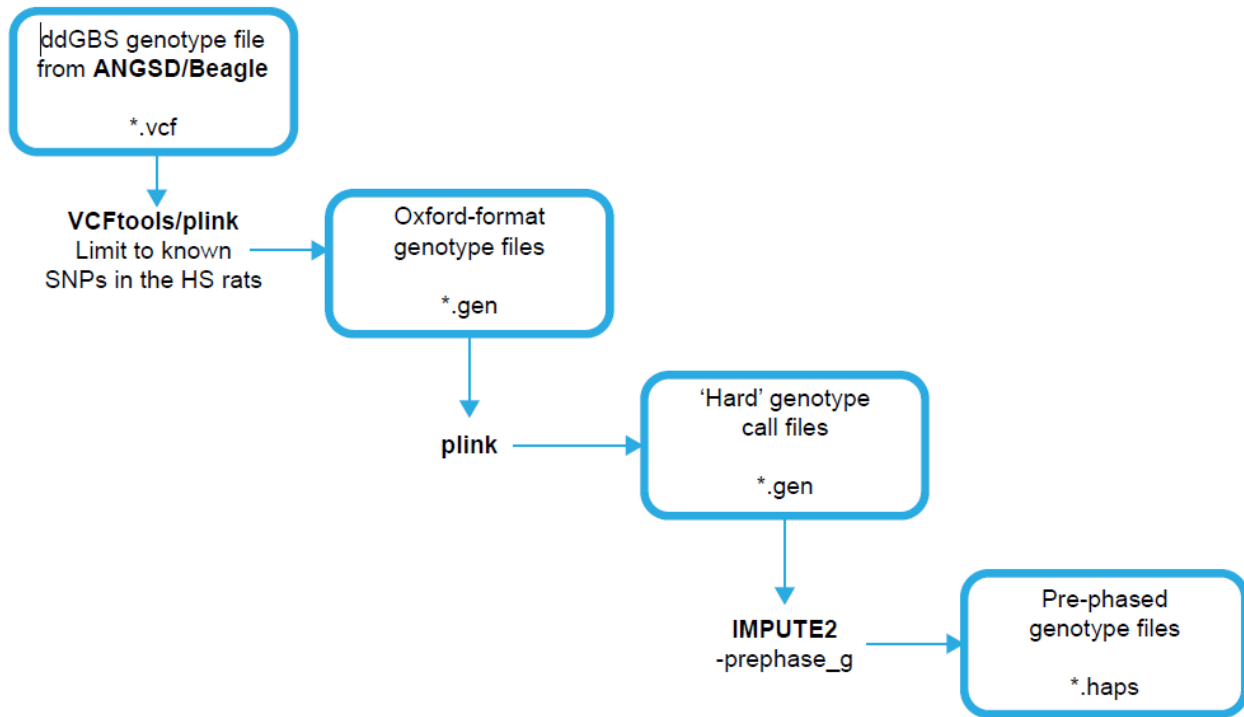
977

978



980 **Figure S2. Data preparation workflow for imputation with IMPUTE2.**

981



982

983

984

985

986

987

988

989

990

991

992

993

994

995

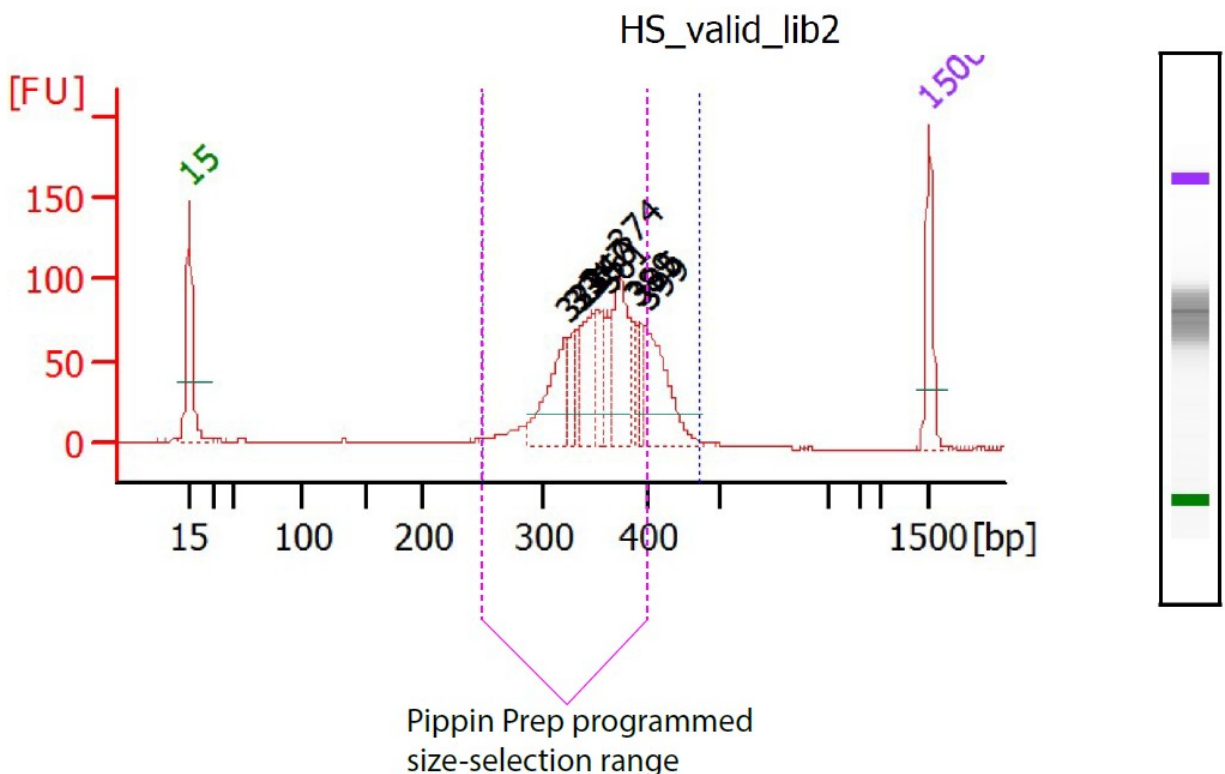
996

997  
998  
999

**Figure S3. Programmed vs. empirical Pippin Prep fragment size range.**

1000 This plot comes from the Bioanalyzer output for a pooled HS library. The x-  
1001 axis shows the library fragment sizes in base pairs, and the y-axis is in  
1002 fluorescent units, which represent the quantity of the fragments on the gel  
1003 chip. There is approximately a 50-75bp shift in the empirical library  
1004 distribution compared to expectation due to the high quantity of fragments  
1005 loaded into the Pippin Prep gel cassette.

1006



1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015

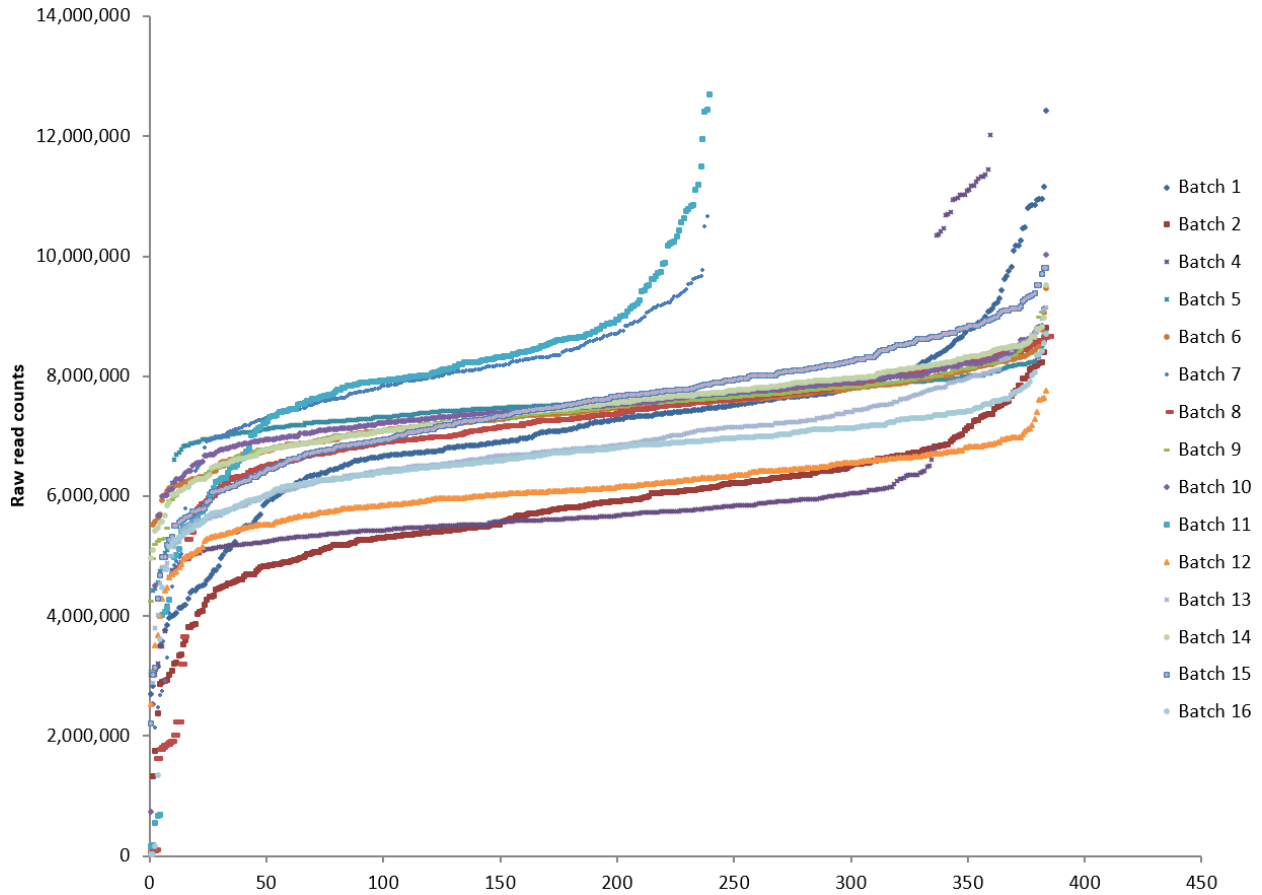
1016

1017

1018 **Figure S4. Raw read counts grouped by shipment batch.**

1019 Raw read counts are on a per-sample basis after demultiplexing FASTQ files  
1020 with FASTX Barcode Splitter. Each batch represents a set of samples from a  
1021 given shipment.

1022



1023

1024

1025

1026

1027

1028

1029

1030

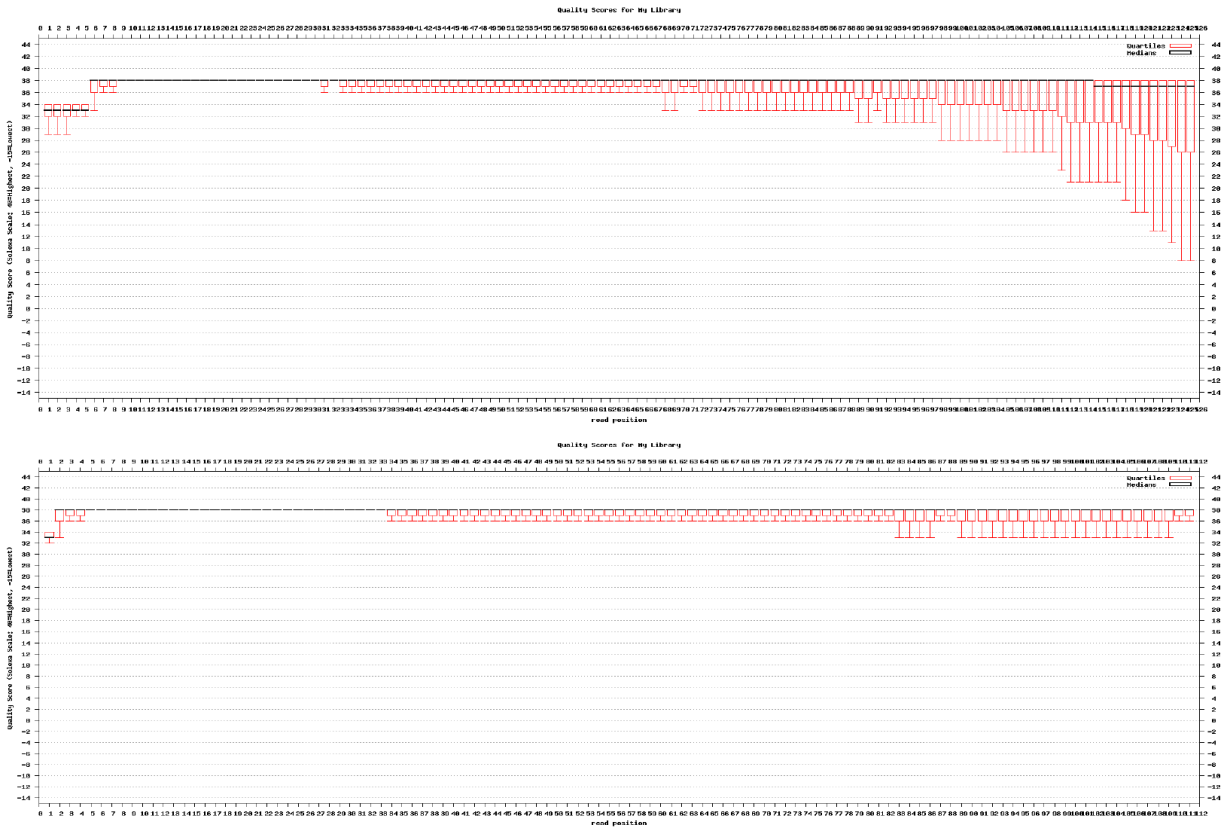
1031

1032

1033

1034 **Figure S5. FASTQC results pre- and post-filtering with Cutadapt.**

1035 FASTQC results are from a single sample from the original set of 96 HS  
1036 samples prepared in 12-plex and sequenced on the Illumina HiSeq 2500 with  
1037 125bp reads.  
1038



1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

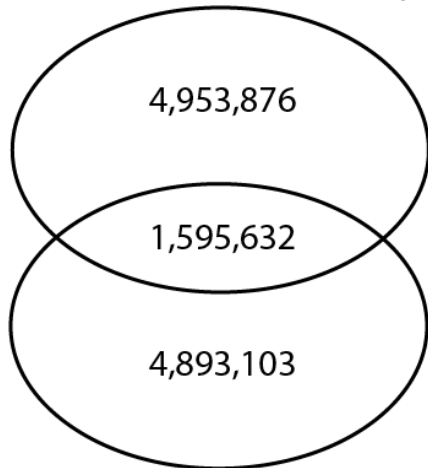
1050

1051

1052 **Figure S6. Overlap of called SNPs with known variants after read**  
1053 **trimming with FASTX or Cutadapt.**

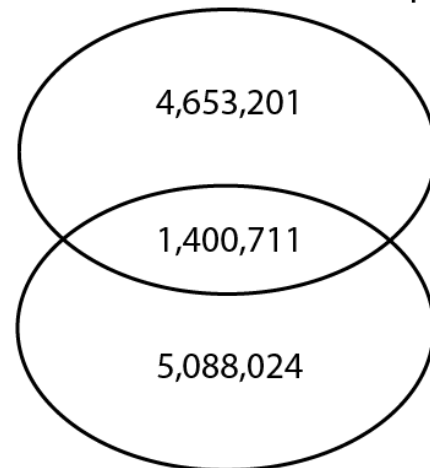
1054

GBS variants after Cutadapt



Known variants from  
42 inbred strains

GBS variants after FASTX Clipper



Known variants from  
42 inbred strains

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

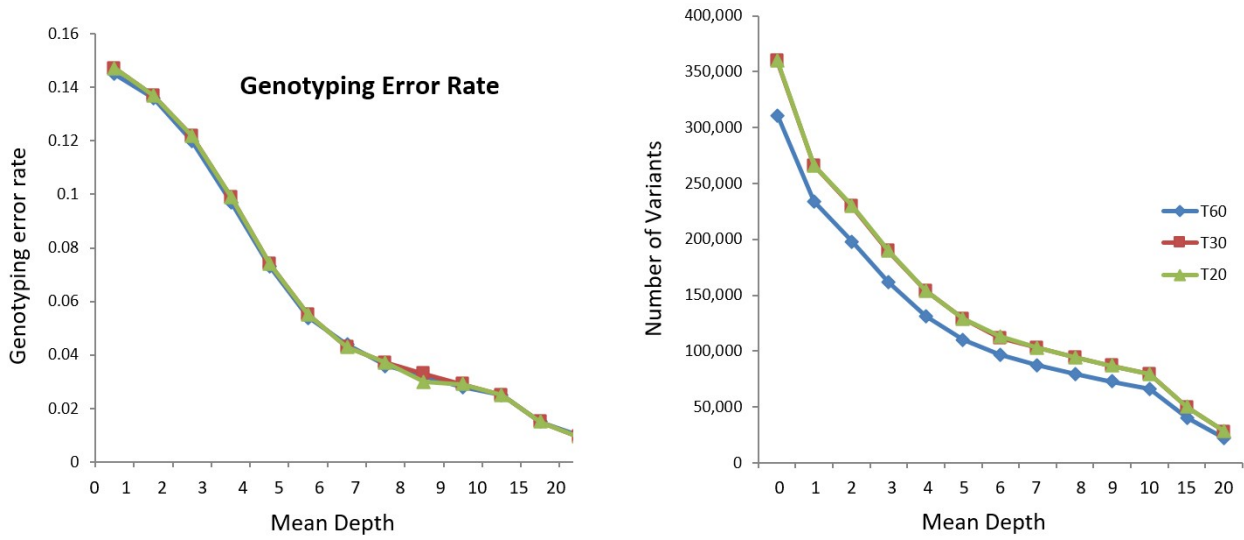
1070

1071

1072 **Figure S7. Mapping quality thresholds.**

1073 Genotyping error rate and number of variants by mean depth per sample per  
1074 variant site for mapping quality thresholds of 20, 30, and 60.

1075



1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

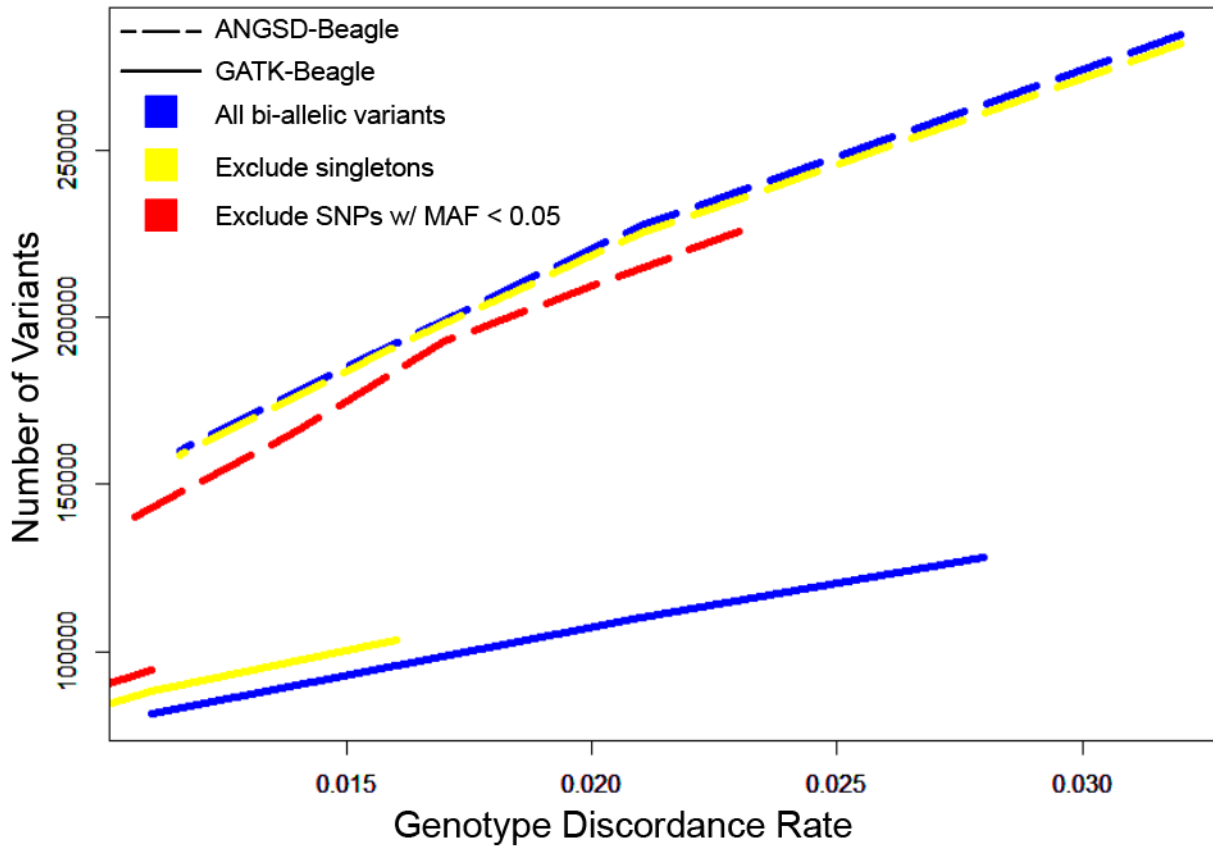
1090

1091

1092 **Figure S8. ANGSD-SAMtools vs GATK HaplotypeCaller, filtered calls.**

1093 The panel compares the number variants called by combination of ANGSD-  
1094 SAMtools and Beagle or GATK HaplotypeCaller and Beagle at various  
1095 thresholds of genotype discordance with array data. Calls were made using  
1096 the 96 HS rats with array data. The x-axis represents the genotype  
1097 discordance rate thresholds and the y-axis is the number of variants that  
1098 surpass that threshold for each genotype calling method. Additional filters  
1099 were applied to the original SNP sets and the plot zooms in on a smaller  
1100 range of acceptable discordance rates compared to Figure 3. Blue lines  
1101 represent the unfiltered SNP set. Yellow lines have been filtered for  
1102 singletons. Red lines have further excluded SNPs with an MAF < 0.05. Each  
1103 line contains the same number of points.

1104



1105

1106

1107

1108

1109

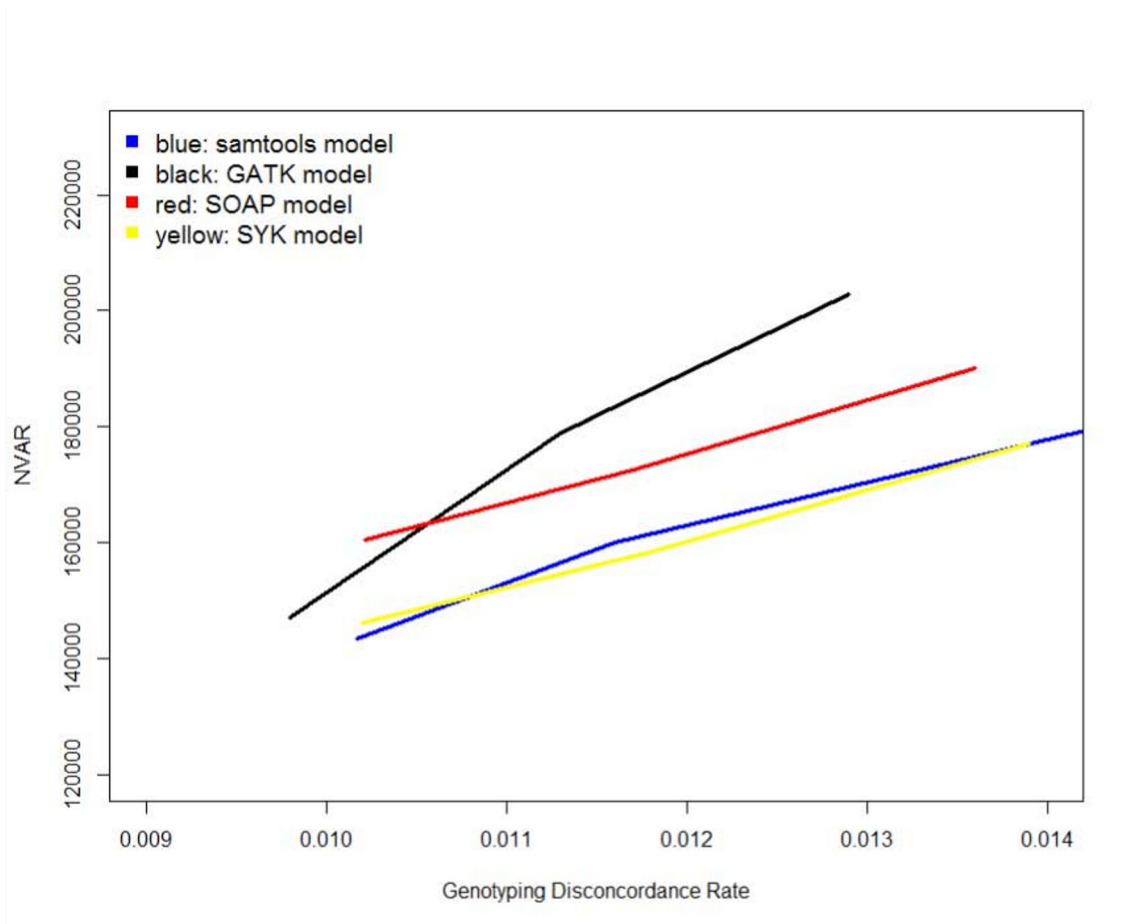
1110

1111

1112

1113 **Figure S9. Number of variants by genotype discordance rates for 4**

1114 **ANGSD genotype likelihood models.**



1115

1116

1117

1118

1120

1121

1122

1123

1124



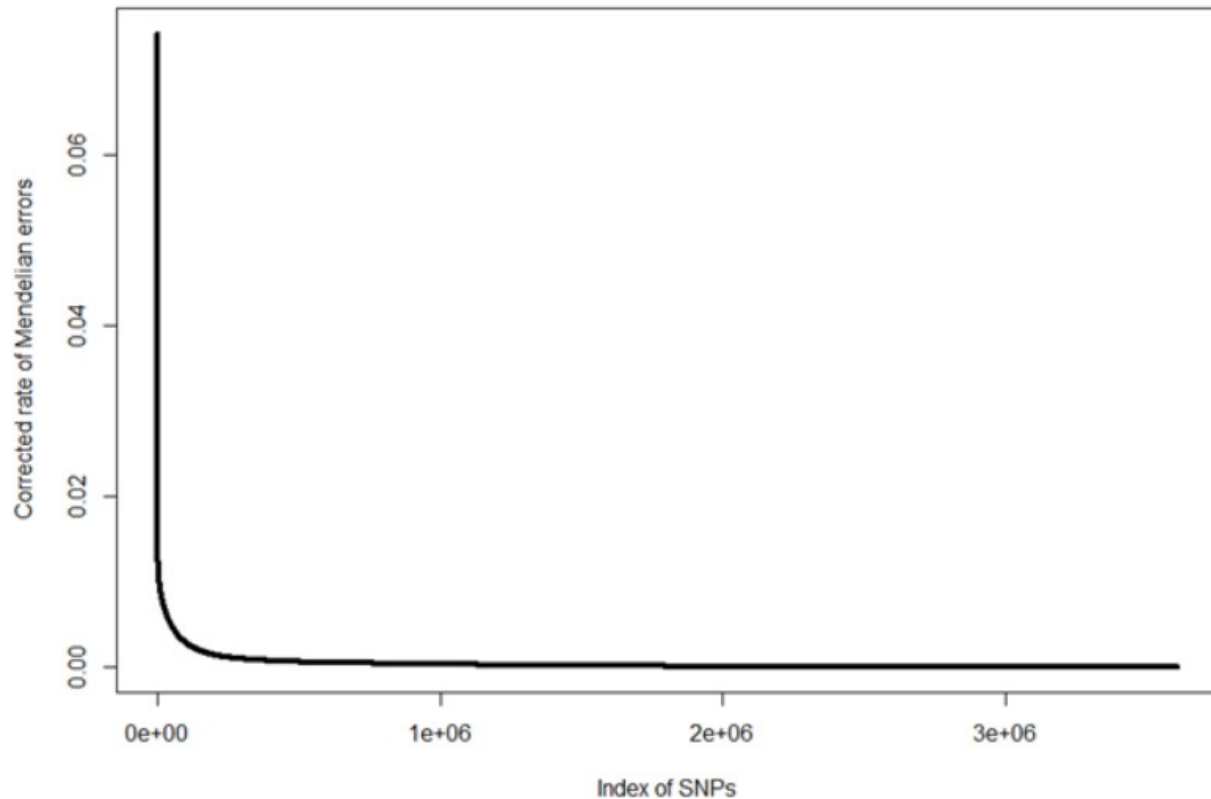
1125

1126

1127

1128 **Figure S10. Mendelian error rates.**

1129 The plot shows the Mendelian error rate for all SNPs. A threshold was set at  
1130 the inflection point of the curve ( $\sim 0.005$ ) and all SNPs above that threshold  
1131 were removed from the data set.



1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

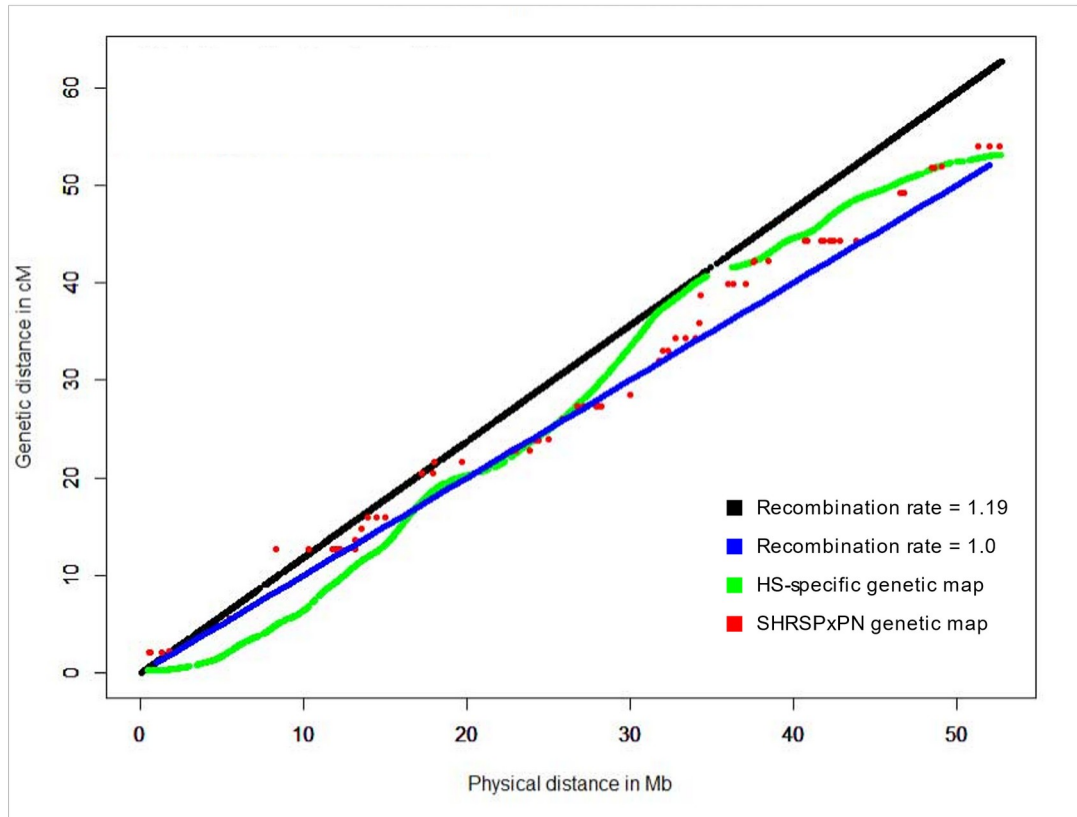
1142

1143

1144

1145 **Figure S11. Available rat genetic maps.**

1146 Plotted physical and genetic distances are for chromosome 12.



1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158 **Table S1. Demultiplexing performance.**

1159 All methods began with the same number of reads from the original FASTQ.  
1160 Final read and base pair counts are from after the reads have been trimmed  
1161 of adapter, barcode, and restriction site sequences, as well as low-quality  
1162 base pairs (< Q20).

1163

	<b>In-house Python Script</b>	<b>GBSX</b>	<b>FASTX Barcode Splitter</b>
<b>Reads with Nlalll adapter sequence</b>	545,177 (3.07%)	475,581 (2.67%)	547,697 (3.07%)
<b>Total bps processed</b>	2,061,523,464	2,116,436,361	2,227,542,500
<b>Total bps written to file</b>	2,059,714,312	2,114,841,934	2,225,724,833
<b>Proportion of bps retained</b>	99.91%	99.92%	99.92%
<b>Reads post-processing</b>	17,771,754	17,786,280	17,820,340

1164

1165

1166

1168

1169 **Table S2. Comparison of variants calls after filtering with FASTX vs**  
1170 **Cutadapt.**

1171 Data shown comes from the original set of 96 HS samples prepared in 12-  
1172 plex and sequenced on the Illumina HiSeq 2500. At this early step of pipeline  
1173 optimization, variants were called utilizing GATK UnifiedGenotyper.

1174

	<b>FASTX Clipper</b>	<b>Cutadapt</b>
<b>Number of variants</b>	6,075,821	6,581,115
<b>Genotyping call rate</b>	0.17	0.19
<b>Mean minor allele count</b>	3.96	4.25
<b>Mean minor allele frequency</b>	0.15	0.15
<b>Number of singletons</b>	433,960	548,975
<b>Number monomorphic sites</b>	807,453	773,074
<b>Transition/ transversion ratio</b>	2.32	2.40
<b>T<sub>I</sub>T<sub>V</sub> ratio for singletons</b>	3.23	3.40
<b>Mean variant read depth</b>	109.56	126.35
<b>Mean quality score</b>	601.79	715.56

1175

1176

1177

1178

1179

1180

1182

1183 **Table S3. Variant metrics resulting from reads filtered at different**  
 1184 **mapping quality thresholds.**

1185 Data shown comes from the original set of 96 HS samples prepared in 12-  
 1186 plex and sequenced on the Illumina HiSeq 2500. Variants were called  
 1187 utilizing the SAMtools model and the -minMapQ filter in ANGSD. Calls were  
 1188 unfiltered.

	<b>MAPQ = 20</b>	<b>MAPQ = 30</b>	<b>MAPQ = 45</b>	<b>MAPQ = 60</b>	<b>MAPQ = 90</b>
<b>Number of variants</b>	372,860	372,330	363,790	316,949	233,322
<b>Genotyping call rate</b>	0.64	0.64	0.64	0.61	0.75
<b>Mean minor allele count</b>	5.96	5.96	6.06	5.86	7.36
<b>Mean minor allele frequency</b>	0.18	0.18	0.18	0.18	0.19
<b>Number of singletons</b>	16,781 (4.50%)	16,732 (4.49%)	16,550 (4.55%)	17,352 (5.47%)	11,773 (5.05%)
<b>Number of monomorphic sites</b>	122,478 (32.85%)	122,188 (32.82%)	116,738 (32.09%)	100,074 (31.57%)	56,179 (24.08%)
<b>Transition/transversion ratio</b>	1.23	1.24	1.26	1.31	1.41
<b>T<sub>i</sub>T<sub>v</sub> ratio for singletons</b>	1.27	1.28	1.28	1.31	1.38
<b>Mean variant read depth</b>	157.78	157.73	159.25	152.48	188.80
<b>Mean quality score</b>	2,547	2,548	2,556	2,461	2,954

1189

1190

1192

1193 **Table S4. Transition/transversion ratio before and after known sites**  
1194 **filtering.**

1195 The presented data comes from ANGSD-SAMtools/Beagle variant calls for  
1196 3,601 HS samples, prior to imputation with IMPUTE2. Known SNPs came from  
1197 both the 42 inbred genomes from Hermsen et. al 2015 (Hermsen et al. 2015)  
1198 and the 8 inbred HS founder strains sequenced by the University of Michigan  
1199 (Ramdas et al. 2019).

1200

	<b>Unfiltered SNPs</b>	<b>Filtered for known SNPs</b>
<b>AC</b>	15,157	9,166
<b>AG</b>	888,657	42,275
<b>AT</b>	15,432	7,610
<b>CG</b>	18,043	8,061
<b>CT</b>	893,653	41,938
<b>GT</b>	15,118	9,177
<b>T<sub>s</sub></b>	1,782,310	84,213
<b>T<sub>v</sub></b>	63,750	34,014
<b>T<sub>s</sub>T<sub>v</sub></b>	27.96	2.48
<b>Total # SNPs</b>	1,846,060	118,227

1201

1202

1204

1205 **Table S5. Imputation accuracy for chromosome 12 across different**  
1206 **genetic maps.**

1207 The number of variants used for the concordance check is dependent on the  
1208 overlap of the imputed variants with array data for the 96 HS rats with array  
1209 genotypes. The MAF filter only removes monomorphic sites within the 96 HS  
1210 rat sample used for the concordance check.

1211

	<b>cM/Mb = 1.00</b>	<b>cM/Mb = 1.16</b>	<b>SHRSPxP N</b>	<b>HS- specific</b>
<b>Number of variants before QC</b>	158,452	158,452	158,452	158,452
<b>Genotyping rate before QC</b>	0.94	0.92	0.92	0.92
<b>Variant removed for missingness &gt; 10%</b>	22,217	28,959	28,356	28,858
<b>Variants removed for MAF &lt; 0.005</b>	50,380	61,270	61,592	59,812
<b>Variants removed for HWE &lt; 1x10<sup>-10</sup></b>	53	56	57	56
<b>Number of variants after QC</b>	85,802	68,167	68,447	69,726
<b>Genotyping rate after QC</b>	0.93	0.91	0.92	0.91
<b>Number of variants in concordance check</b>	5,912	5,590	5,594	5,646
<b>Discordance rate</b>	0.095	0.011	0.011	0.010

1212

1213

1214