

1 Incorporating selfing to purge deleterious alleles in a cassava genomic selection program

2

3 Mohamed Somo¹ and Jean-Luc Jannink^{1,2*}

4

5 ¹Dept. of Plant Breeding and Genetics, Cornell University, Ithaca, NY, USA

6 ²USDA-ARS, R.W. Holley Center for Agriculture and Health, Ithaca, NY, USA

7

8

9 *Corresponding author: Jean-Luc Jannink, jeanluc.jannink@usda.gov

10

11 **Abbreviations:** BSL, Breeding Scheme Language, CET, clonal evaluation trial, GEBV,
12 genomic estimated breeding value, GS, genomic selection, PYT, preliminary yield trial

13

14

15

16 **Abstract**

17 Cassava has been found to carry high levels of recessive deleterious mutations and it is known to
18 suffer from inbreeding depression. Breeders therefore consider specific approaches to decrease
19 cassava's genetic load. Using self fertilization to unmask deleterious recessive alleles and
20 therefore accelerate their purging is one possibility. Before implementation of this approach we
21 sought to understand better its consequences through simulation. Founder populations with high
22 directional dominance were simulated using a natural selection forward simulator. The founder
23 population was then subjected to five generations of genomic selection in schemes that did or did
24 not include a generation of phenotypic selection on selfed progeny. We found that genomic
25 selection was less effective under the directional dominance model than under the additive
26 models that have commonly been used in simulations. While selection did increase favorable
27 allele frequencies, increased inbreeding during selection caused decreased gain in genotypic
28 values under the directional dominance. While purging selection on selfed individuals was
29 effective in the first breeding cycle, it was not effective in later cycles, an effect we attributed to
30 the fact that the generation of selfing decreased the relatedness of the genomic prediction training
31 population from selection candidates. That decreased relatedness caused genomic prediction
32 accuracy to be lower in schemes incorporating selfing. We found that selection on individuals
33 partially inbred by one generation of selfing did increase mean genetic value of the partially
34 inbred population, but that this gain was accompanied by a relatively small increase in favorable
35 allele frequencies such that improvement in the outbred population was lower than might have
36 been intuited.

37 **Introduction**

38 Over the evolution of wild ancestors of crop plants, during domestication, and during breeding,
39 crop populations can accumulate deleterious alleles (Valluru et al., 2019; Ramu et al., 2017).
40 This problem may be more severe in clonally-propagated crops, such as cassava, because of the
41 many propagation cycles they undergo between meiotic recombination events. The propagation
42 cycles mean more rounds of DNA replication leading to opportunities for mutation whereas
43 recombination facilitates purging by decreasing the association of favorable and deleterious
44 alleles across loci. Cassava has been found to carry high levels of deleterious mutations,
45 distributed throughout the genome (Ramu et al., 2017). Ramu et al. (2017) further found that
46 over the course of modern breeding, deleterious alleles have become more common, but that
47 performance among cultivars has been maintained by the masking of these alleles in
48 heterozygous state. Given the recessive mode of action of these alleles, bringing them to
49 homozygous state could be a way to detect and purge them.

50 Selfing and subsequent inbreeding depression leading to fitness loss can cause deleterious
51 recessive allele purging under natural selection (Charlesworth & Willis, 2009; Lande et al.,
52 1994). Although rounds of selfing can help purge recessive deleterious alleles, it is less efficient
53 against mildly deleterious alleles. This means that while selfing could help eliminate strongly
54 deleterious alleles it might also accelerate the fixation of mildly deleterious alleles (Boakes &
55 Wang, 2005).

56
57 Selfing in cassava is associated with inbreeding depression and loss of fitness (Rojas et al., 2009;
58 Nuwamanya et al., 2011; Ceballos et al., 2004). The expression of recessive alleles has been used
59 for gene discovery. For example, genotypes carrying waxy genes associated with starch were S₁
60 derived from AM206-5 clone (Ceballos et al., 2007). Similarly, Prochnik et al. (2012) used the
61 predominantly homozygous AM560-2 cassava clone (S₃ derived from MCOL-1505) for genome
62 sequencing. Studies have reported varied inbreeding depression on different cassava traits. Both
63 Rojas et al. (2009) and Nuwamanya et al. (2011) examined the effect of a single round of selfing
64 on fresh root yield and found 60% reduction in yield. Kawuki et al. (2011) reported a 36%
65 decrease of dry matter following a cycle of selfing. Finally, de Freitas et al. (2016) found that
66 inbreeding depression varied widely between clones, ranging from 2% to 55% for fresh root
67 yield, 0% to 9% for height, and 0% to 2% for dry matter content.

68

69 With GS breeding schemes being adopted in cassava breeding (Wolfe et al., 2017), there is a
70 higher risk of inbreeding because cycles of selection are more frequent and each cycle entails a
71 bottleneck (Ozimati et al., 2019). Ozimati et al. (2019) compared inbreeding between C0 and C1
72 populations using the average of diagonal elements of a kinship matrix, as a measure of the
73 inbreeding coefficient. While they found less inbreeding in the C1 than in the C0 population,
74 they attributed that to a purposeful crossing design that prevented genomically similar
75 individuals from being crossed. Other GS studies, as well as simulation, have found a potential
76 for rapid loss of diversity (Rutkoski et al., 2015; Jannink, 2010).

77

78 Empirically testing the impact of selfing cycles for cassava with a relatively long life cycle
79 would be slow and expensive so that breeders are reluctant to explore selfing as a component of
80 the GS scheme. To our knowledge, no GS breeding studies with selfing have been conducted to
81 estimate the effect that selfing might have on gain from selection. In addition to empirical tests,
82 stochastic simulation can be used to test breeding schemes and gain insight into the impact of
83 specific changes such as incorporating selfing. In this study our objectives were to (1) compare
84 breeding schemes with and without selfing in scenarios with different training population sizes
85 and selection intensities under directional dominance and additive modes of gene action, (2)
86 identify mechanisms affecting the observed responses to selection in terms of changes in the
87 frequency of favorable alleles and deleterious homozygote genotypes.

88

89 **Materials and methods**

90 **Generation of founder haplotypes using forward simulation**

91 Because the cassava genome has been found to harbor many deleterious recessive mutations
92 (Ramu et al., 2017), we chose to start breeding schemes from founder haplotypes with evolved
93 directional dominance, induced both by the mode of gene action and historical natural selection.
94 We generated starting populations using the forward simulator SLiM (Haller & Messer, 2017).
95 SLiM is a stochastic simulator enabling flexible simulation at the base pair level. Our
96 simulations involved evolutionarily constrained segments that we call here “genes” though
97 superficial evaluation shows that these segments resemble the genes of molecular geneticists
98 poorly. Likewise, “base pairs” in our genes underwent stochastic mutations with different
99 probabilities and severity of deleteriousness. Our objective was to generate haplotypes with
100 directional dominance including some level of complexity in their gametic phase disequilibrium
101 and allele frequency spectrum. We simulated genomes with five chromosomes of 30 Mbp and 90
102 cM length each. These chromosomes were interspersed with genes every 100,000 bp (300 genes
103 per chromosome). Genes were 20,000 bp in length, with characteristic mutation effects and mode
104 of action dependent on position within the gene (Table 1). The simulation gave genes rather long
105 regions subject to mutations with milder effects where the deleterious allele was only partially
106 recessive, and more constrained regions where mutations had stronger effects, were more
107 frequently deleterious, and where the deleterious alleles were more completely recessive. The
108 overall mutation rate was $1e-7$ per base pair per generation. The forward simulation was run for
109 12,000 generations with a census population size of 500.
110 We recognize that the cassava genome is much different from the genome that we simulated,
111 with 770 million base pairs of sequence, 18 chromosomes and 33,033 annotated genes. In
112 recombination space, the genome spans 2412 cM (International Cassava Genetic Map
113 Consortium, 2015). We nevertheless believe that the genome we simulated can provide
114 instructive insights.

115 **Table 1** Mutation characteristics as a function of region in genes. The selection coefficient s
116 gives the fitness of the homozygote mutant ($W = 1 + s$) relative to the non-mutant ($W = 1$). The
117 dominance coefficient h gives the fitness of the heterozygote ($W = 1 + hs$). The Deleterious
118 Frequency indicates the percentage of mutants in the region that are deleterious.

Base pair position of region	Deleterious Frequency (%)	Deleterious		Favorable	
		s	h	s	h
1 – 10,000	90	-0.004	0.20	0.004	0.80
10,001 – 14,000	95	-0.008	0.10	0.004	0.80
14,001 – 16,000	98	-0.016	0.05	0.008	0.90
Two region types alternating every 200bp	98	-0.030	0.02	0.008	0.90
16,001 – 20,000	95	-0.008	0.10	0.004	0.80

119

120 **Breeding scheme simulation**

121 The BreedingSchemeLanguage (BSL) R Package (Yabe et al., 2017) was used to simulate
122 breeding schemes. We assumed there were no new mutations during the breeding program.
123 Simulations in SLiM parameterized mutations in terms of their effects on fitness. The BSL
124 environment, however, parameterized mutations in terms of their effects on a quantitative trait.
125 In the BSL, then, the effect of the trait was proportional to the s value shown in Table 1, and its
126 degree of dominance was the same as in SLiM, such that the favorable (high-valued) allele was
127 always dominant over the deleterious allele. The initial genetic variance was set to 1. The
128 founders were considered to have been well phenotyped, thus initial error variance was set to 1,
129 leading to heritability 0.5 for founders. Error variances associated with subsequent evaluations at
130 seedling (as is common in cassava, we refer to any plant grown from botanical seed as a
131 “seedling” regardless of how old and big it gets), clonal evaluation trial (CET), and preliminary
132 yield trial (PYT) plot types for model updating were set to 16, 9, and 4, respectively. The error
133 variances for different evaluation types were held constant in each cycle of the breeding
134 schemes.

135

136 We implemented a five-cycle scheme after initial evaluation and selection of founder
137 populations. The standard scheme involving selfing was as follows. Founders' genomic
138 estimated breeding values (GEBVs) were estimated using ridge regression (Endelman, 2011) and
139 sixty individuals with high GEBVs were selected as parents for the next generation. Those
140 parents were selfed, creating 10 S_1 progeny per parent. The resulting 600 partial inbreds were
141 evaluated as seedlings. The five progenies with highest values were chosen from each of the 60
142 families for the crossing nursery. To generate 1,000 outbred progenies, the 300 individuals (60
143 selected parents x 5 S_1 per parent) were randomly mated, excluding within-family matings.
144 Outbred progenies were genotyped and their GEBVs estimated using all prior phenotypic data,
145 again selecting the best 60 individuals. This scheme was repeated in subsequent cycles.

146
147 The genomic prediction training population started with the founders (cycle C0). After the creation
148 of C3 progeny (but before they were selected), the training population was updated with
149 phenotypic records of 440 non-parental cycle C1 individuals and the 60 C1 parents of C2 progeny
150 using CET plots. After the creation of C4 progeny, the training population was updated with
151 phenotypic records of 440 non-parental cycle C2 individuals and the 60 C2 parents of C3 progeny
152 using CET plots as well as 190 non-parental C1 individuals and the 60 C1 parents of C2 progeny
153 using PYT plots. The number of C5 progeny (the final generation) was set to 200. All breeding
154 schemes were replicated 24 times.

155

156 **Simulation scenarios**

157 In addition to the standard scheme with selfing, we simulated a no selfing scheme. There, the 60
158 selected parents were intermated at random leading directly to the next generation. We
159 considered schemes with or without a selfing generation to have the same overall breeding cycle
160 time. We also simulated a "selfing-plus" (Self+) scheme used to update the training population
161 more rapidly. In the Self+ scheme, the phenotypic data from evaluating S_1 progeny of selected
162 parents were incorporated back into the training population by attributing to the parent the mean
163 value of its ten S_1 progeny. Standard schemes involved 600 founders as the initial training
164 population. We also tested founder population sizes of 200 and 1500. Standard schemes involved
165 1,200 markers per chromosome. We also tested schemes with 3,600 markers per chromosome.
166 Standard schemes involved a selection intensity of 10% among founders and 6% thereafter. We

167 also tested selection intensities of 5% among founders and 3% thereafter. Finally, we simulated
168 the standard schemes with an additive gene model to contrast with the directional dominance
169 gene model.

170

171 **Analysis of simulation results**

172 We tracked change in genotypic value across breeding cycles. The QTL could have effects of
173 four different sizes (Table 1). We tracked mean favorable allele frequency across all loci within
174 an effect class, so that “allele frequency” was actually calculated over many loci. We also
175 tracked a composite favorable allele frequency as a weighted mean, with weights being the effect
176 of the locus class (i.e., weights of 4, 8, 16, and 30, Table 1). Similarly, we tracked homozygote
177 deleterious genotype frequency for which the frequency is over individuals rather than over
178 gametes. We measured the impact of selection on inbreeding depression by calculating the loss
179 of genetic value due to one generation of self-fertilization in the founder population as compared
180 to the Cycle 5 population. For all of these measures, we performed simple ANOVA or t-tests for
181 relevant comparisons across cycles or across simulation scenarios to determine significance of
182 differences.

183

184 We called selection on GEBVs of outcrossed individuals as the “main selection” and on selfed S_1
185 individuals as “purging selection.” Populations in schemes with no selfing were subject only to
186 main selection while those with selfing underwent main selection and purging selection.

187 Comparison of gain due to main and purging selection were done for additive and dominance
188 models for the Self, NoSelf, and Self+ schemes, and across the three founder population sizes
189 (200, 600, and 1500).

190

191 **Data availability**

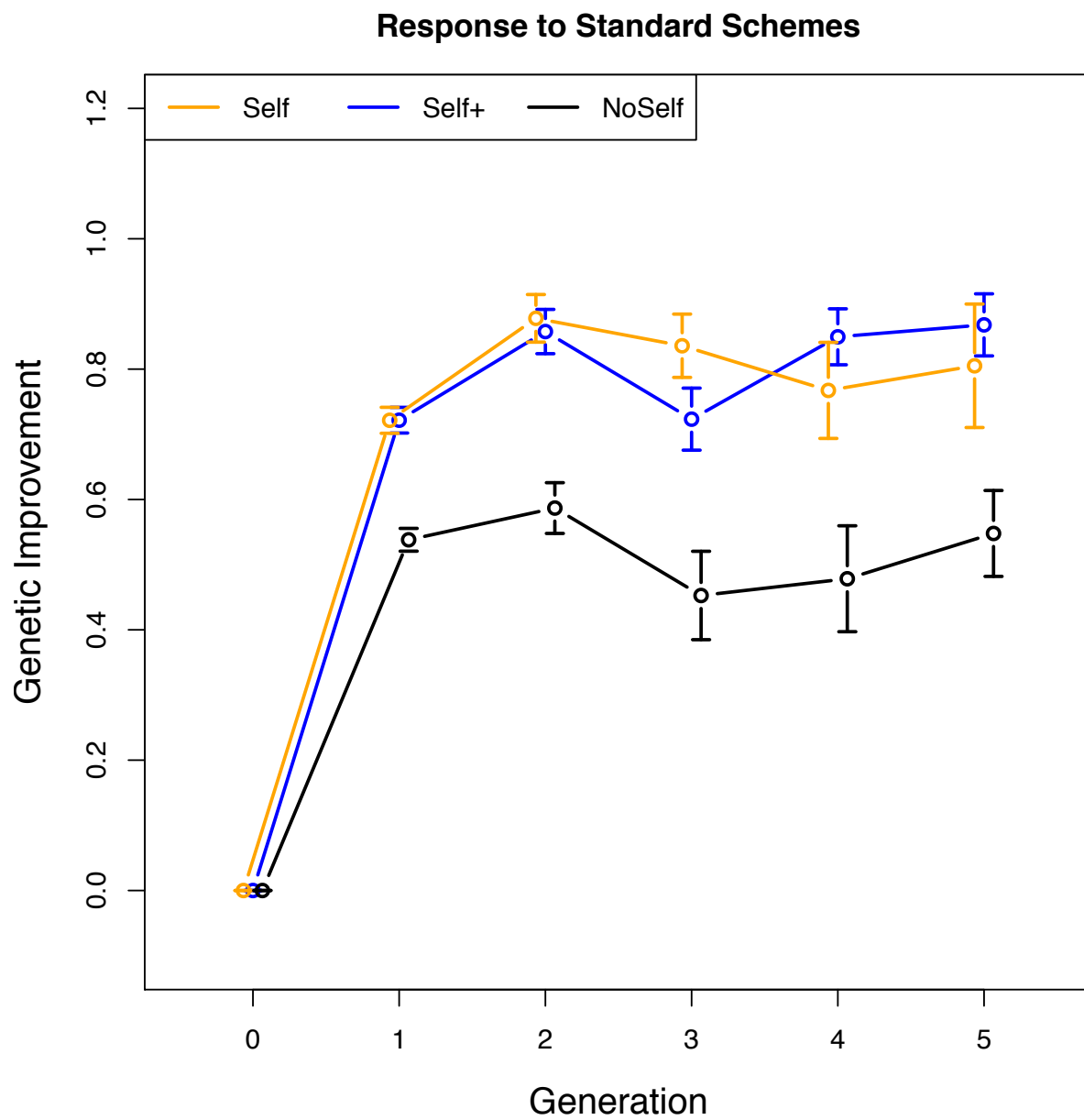
192 No empirical data were used or generated. The simulation and analysis scripts are deposited on
193 github at <https://github.com/jeanlucj/PurgingBySelfingSimulations>.

194

195 **Results**

196 **Genomic selection under directional dominance**

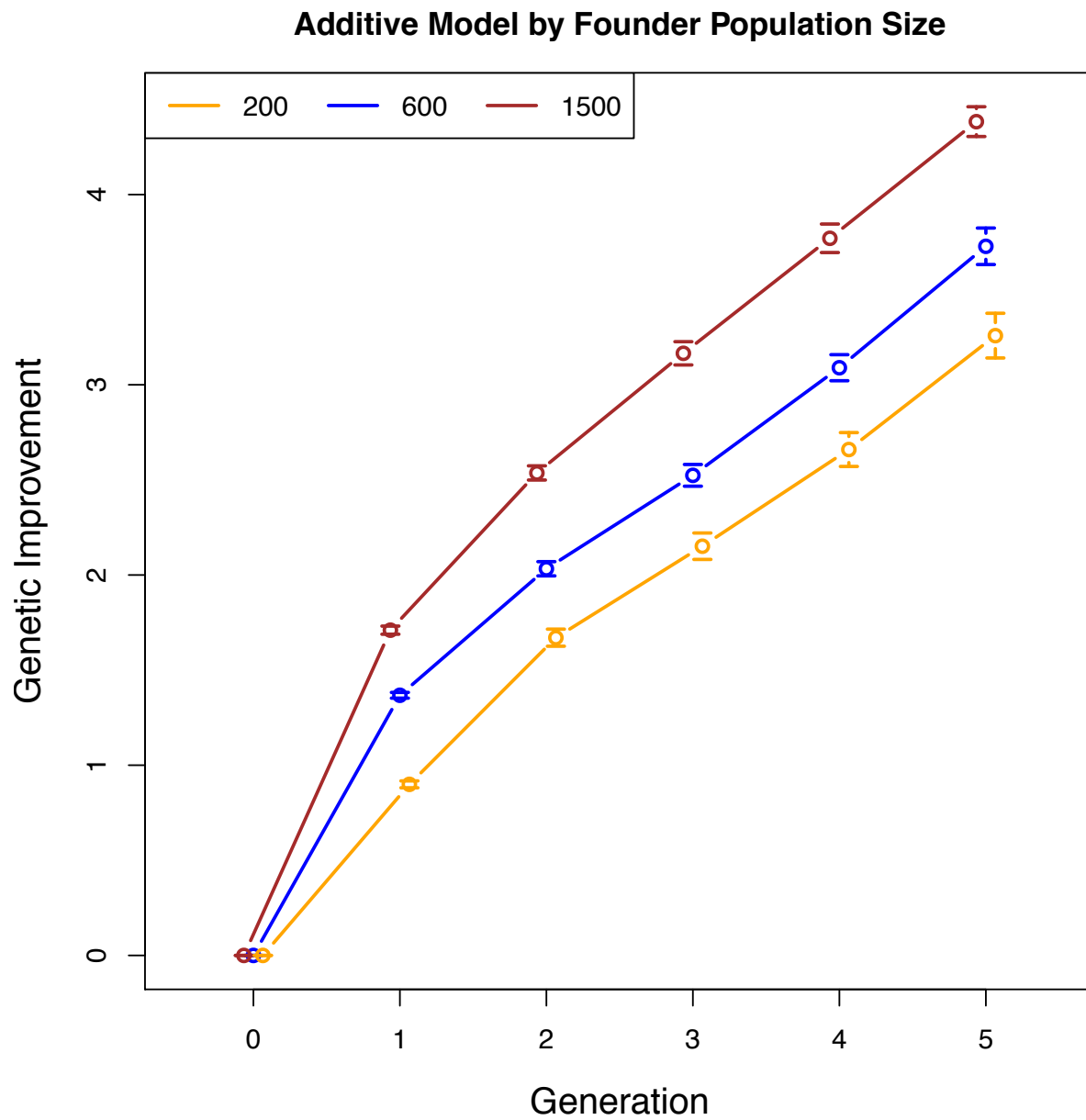
197 We observed significant differences among simulated breeding schemes (Fig. 1). The Self and
198 Self+ schemes achieved higher means than the NoSelf scheme. Qualitatively, the responses of all
199 schemes followed a similar pattern: rapid gain in the first cycle, diminished gain in the second
200 cycle, loss in the third cycle, and recovery in the fourth and fifth cycles (Fig. 1). This pattern
201 meant that gain in genetic mean from generations 1 to 5 was not significant. Genomic prediction
202 models lose accuracy after selection (e.g., Muir, 2007). To determine if the pattern observed here
203 was simply caused by that mechanism or was caused by the mode of gene action, we simulated
204 selection under additive gene action (Fig. 2). Under additivity, the pattern of response was very
205 different, with steady gain observed across different founder population sizes (Fig. 2). Note that
206 gain under the additive model was three to four times greater than under the directional
207 dominance model, despite all schemes starting from populations with the same genetic variance.
208 Variation in gain among replicate simulations was also lower under the additive than the
209 dominance gene models, as reflected by the smaller standard errors of population means (Fig. 2).
210 To better understand the response under directional dominance, we examined allele and
211 genotype frequencies across the four classes of loci in the simulations. We wanted to distinguish
212 two hypotheses. First, loss of accuracy of the prediction model would be expected to occur more
213 quickly under directional dominance than additivity (Duenk et al., 2019). So, despite little loss of
214 accuracy under the additive model, the question of a loss of accuracy for breeding value
215 prediction under dominance was still relevant. Second, under directional dominance the
216 genotypic value depends strongly on the frequency of the deleterious allele homozygote
217 genotype frequency. Consequently, decreased response could be due to increased frequency of
218 that genotype. Favorable allele and deleterious homozygote genotype frequencies supported the
219 second hypothesis (Fig. 3). In particular, the breeding values of individuals was a function of the
220 favorable allele frequencies (Falconer & Mackay, 1996), which trended upward for all classes of
221 loci (Fig. 3), though certainly more so for the larger effect alleles. Favorable allele frequency
222 changes were small, as is to be expected under an infinitesimal model. The trend upward shows
223 the effectiveness of genomic selection to increase breeding values, even under directional
224 dominance.



225

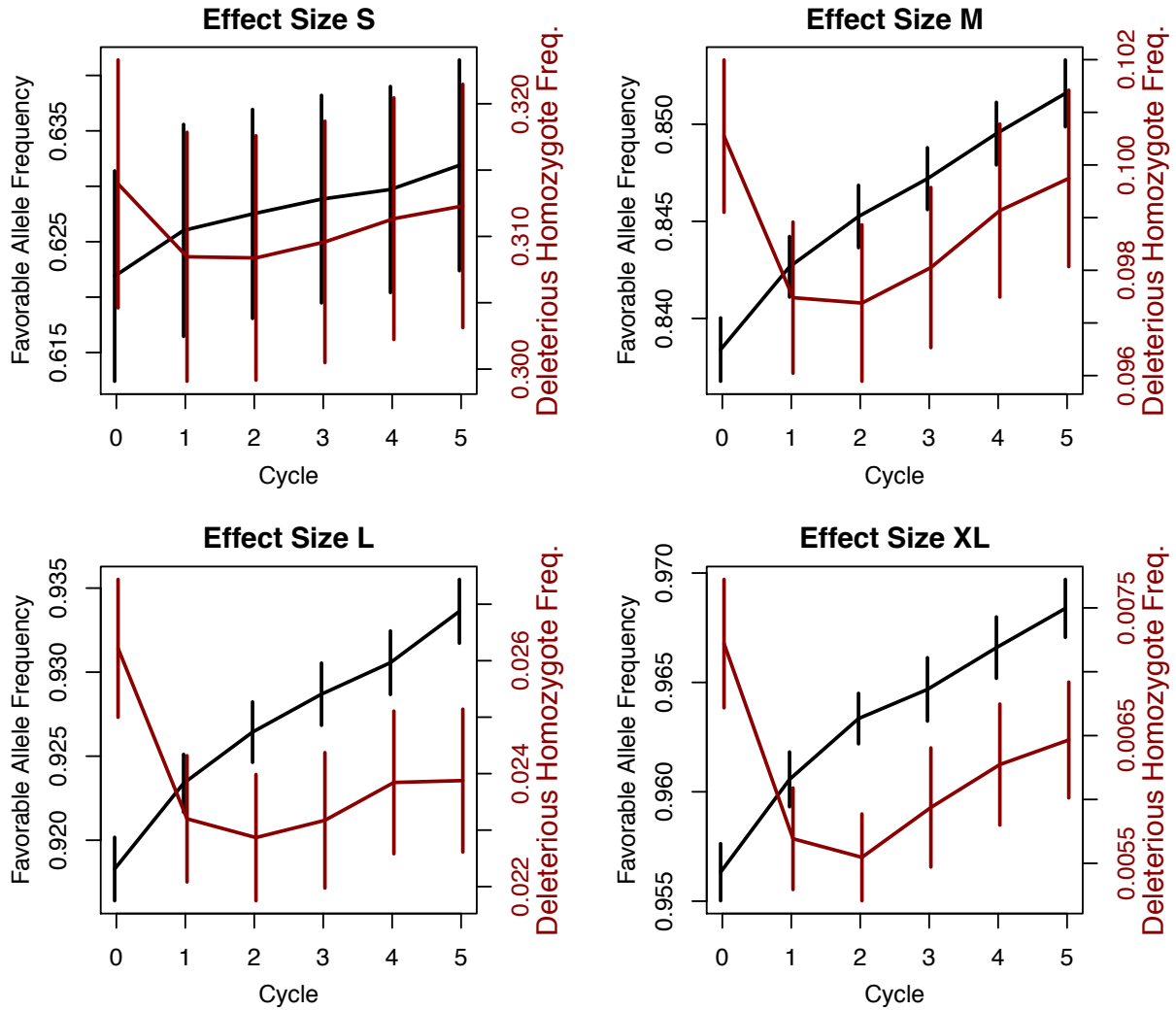
226 **Fig. 1.** Gain from selection using two schemes with selfing (Self and Self+) compared to one without selfing
227 (NoSelf) over five generations. The number of founders was 600.

228



229

230 **Fig. 2.** Gain from the Self breeding scheme under an additive gene action model, assuming different
231 founder population sizes. Other than mode of gene action, all parameters were the same as for the
232 breeding schemes shown in Fig. 1.



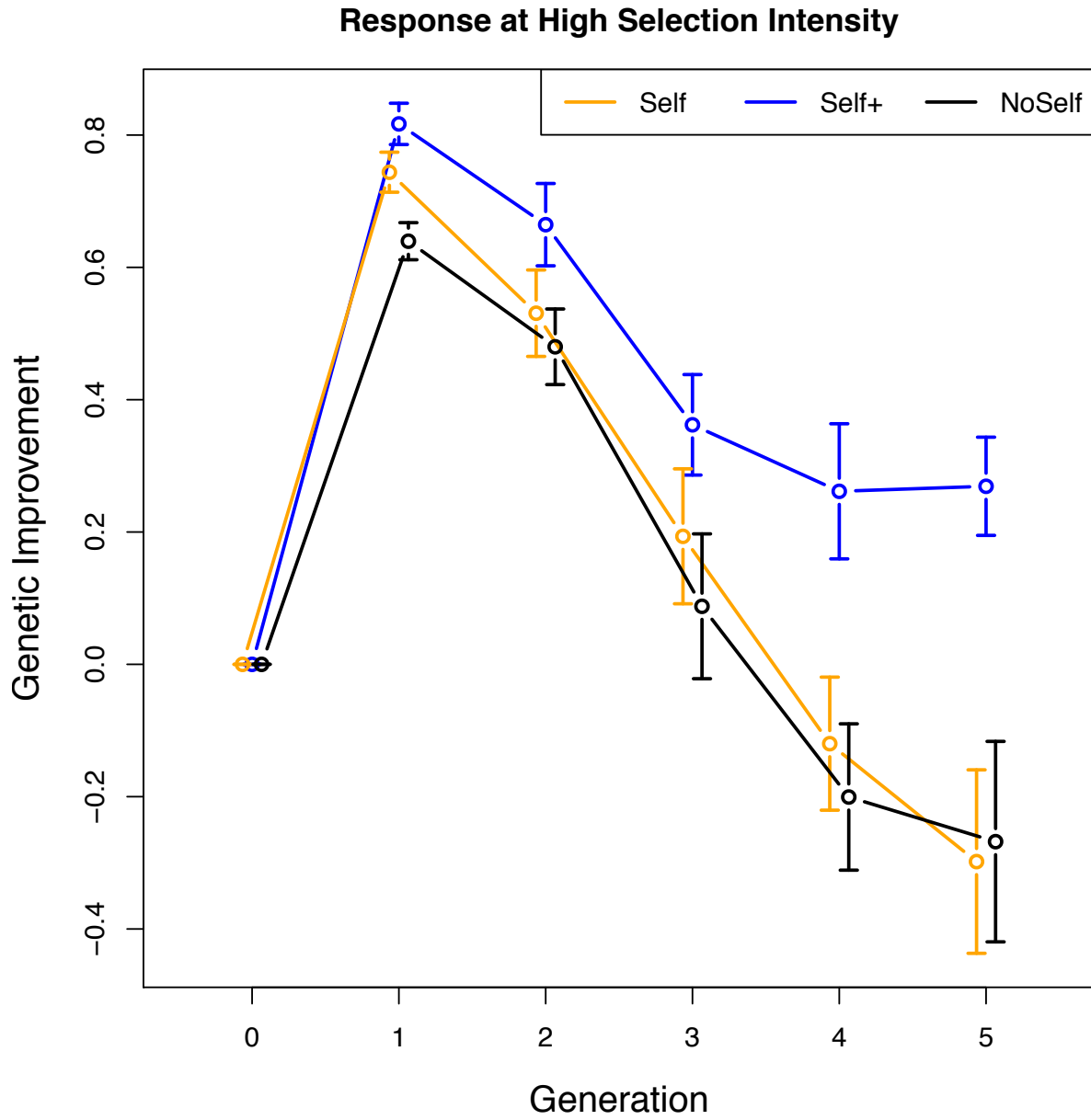
233

234 **Fig. 3.** Favorable allele frequencies (black, left scale) and deleterious allele homozygote frequencies (red,
235 right scale, note differences in scales) over five cycles of selection in the standard setting. Effect sizes of
236 0.04, 0.08, 0.16, and 0.30 labeled S, M, L, and XL, respectively. Error bars are +/- one standard error
237 across the 24 repeated simulations.

238 Apparently paradoxically, however, the frequency of the deleterious allele homozygote genotype
239 also increased. The trend for increase of favorable allele frequency with concomitant increase in
240 the frequency of deleterious allele homozygotes was observed across all allele effect sizes. As a
241 possible explanation for these opposing trends, we highlight the change in effective population
242 size between the founders (N_e of about 500) and during the breeding cycles (N_e of about 60).
243 That change meant that even as favorable alleles increased in frequency, variation in allele
244 frequency across loci also occurred, leading to fixation at some loci. The hypothesis of fixation
245 leading to higher frequencies of deleterious homozygote genotypes was supported by response to
246 selection under high selection intensity (Fig. 4) and given a small founder population size (Fig.
247 5). In the case of high selection intensity, the number of individuals selected was lower, leading
248 to more severe bottlenecks. Comparing Fig. 4 to Fig. 1, the higher selection intensity leads to
249 greater gain in the first selection event than the standard scenario, but in subsequent generations
250 the populations crash, leading to negative responses to selection (Fig. 4). Interestingly, the Self+
251 scheme, with its more rapid (though limited) updating of the training population to some extent
252 rescues the populations from this crash.

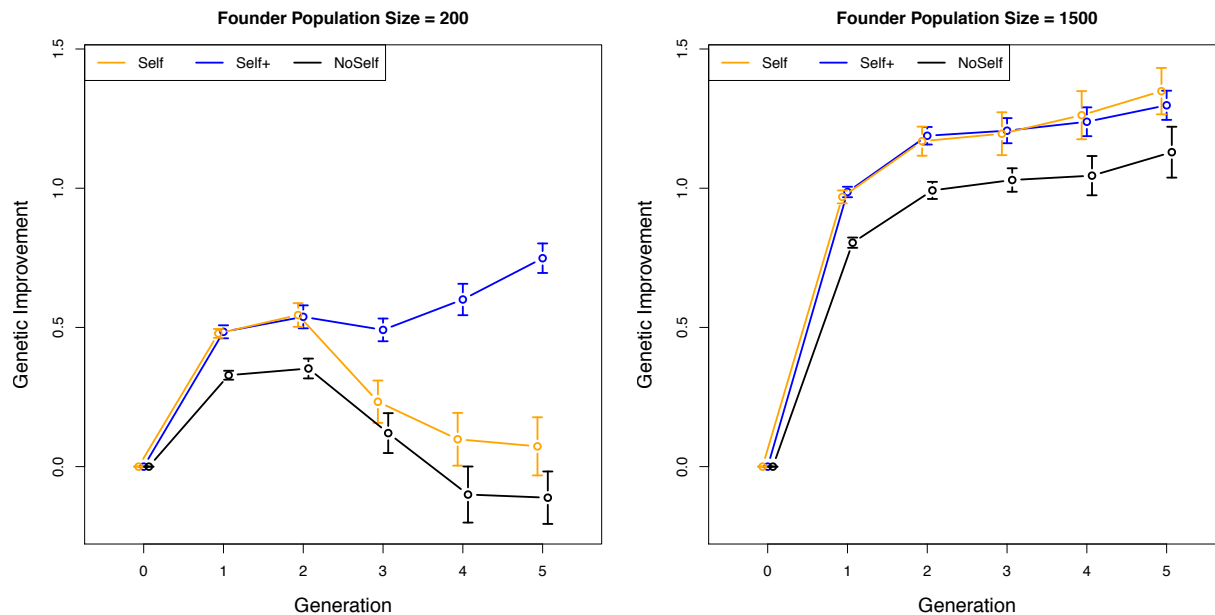
253
254 In the case of small founder population size, it is known that with a smaller training population,
255 genomic selection causes higher co-selection of relatives and therefore more rapid loss of
256 diversity / inbreeding (Jannink et al., 2010). With the small founder population (and therefore
257 smaller training population) we again saw a population crash and negative responses to selection
258 (Fig. 5). The Self+ scheme again rescued these populations from a total crash. In contrast, with
259 the large founder population of 1500 individuals, there was never a decline in the population
260 mean, and the Self and Self+ schemes perform equally well (Fig. 5). As for the standard breeding
261 scheme, the mechanism generating dramatic drops in genotypic value had to do with the
262 accumulation of deleterious homozygotes in the population. For both the high intensity selection
263 and the small founder population breeding scenarios, the Self scheme had higher rates of
264 deleterious homozygote accumulation than the Self+ scheme (Supp. Fig. 1), even though the
265 increase in the favorable allele frequency was similar across schemes. For the large founder
266 population breeding scenario, Self and Self+ schemes did not differ, with low increase in
267 deleterious homozygote frequency and greater increase in favorable allele frequency.

268



269

270 **Fig. 4.** Gain from selection using two schemes with selfing (Self and Self+) compared to one without
271 selfing (NoSelf) over five generations. All parameters were the same as for the breeding schemes shown
272 in Fig. 1 with the exception that 30 parents were selected in each cycle as opposed to 60.
273



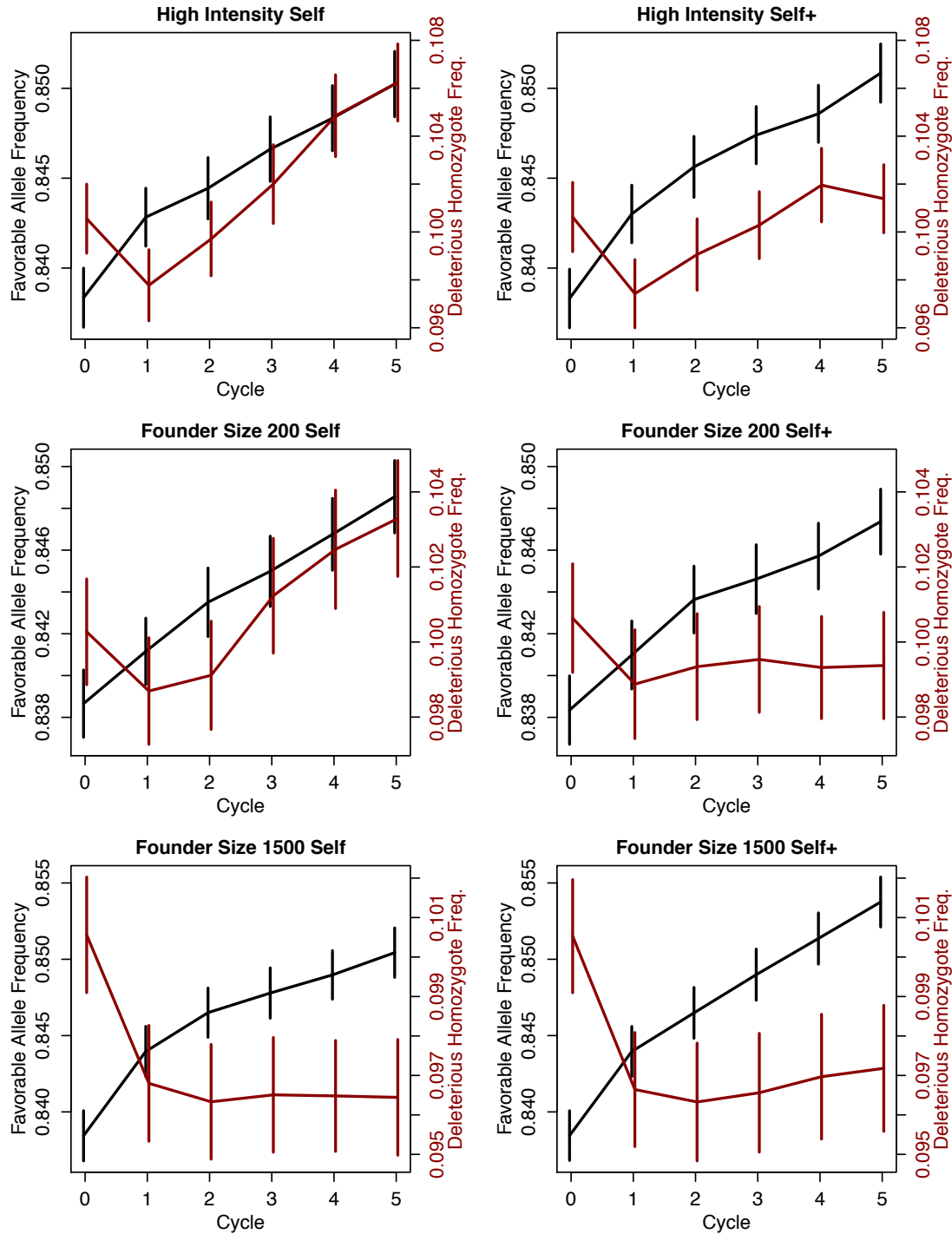
274

275 **Fig. 5.** Gain from selection using two schemes with selfing (Self and Self+) compared to one without
276 selfing (NoSelf) over five generations. All parameters were the same as for the breeding schemes shown
277 in Fig. 1 with the exception that founder population sizes were 200 and 1500 as opposed to 600.
278

279 Selection efficiency after the first cycle with and without selfing

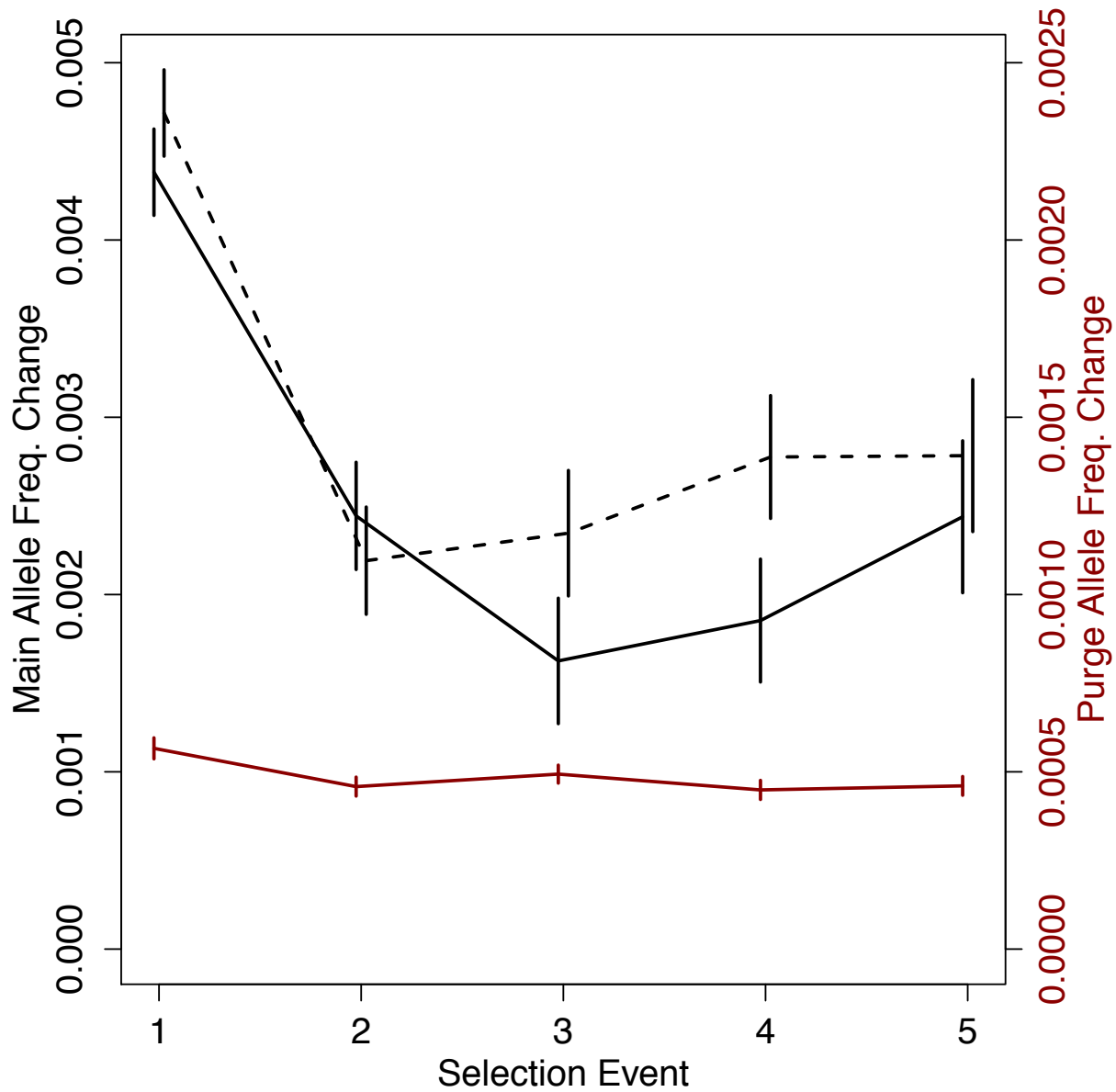
280 Another observation from the standard scheme was that all of the increased gain in the Self and
281 Self+ schemes relative to the NoSelf scheme occurred in the first two cycles of selection (Fig. 1).
282 Thereafter, the gains of the Self and Self+ schemes were not significantly different from that of
283 the NoSelf scheme. The selfing schemes benefit from two selection events leading to genetic
284 gain, the main event from genomic selection on outcrossed individuals and the purging event on
285 selfed individuals. For the selfing schemes not to do better than the NoSelf scheme requires that
286 one or both of these events must decrease in effectiveness in the selfing schemes after the first
287 two cycles of selection. In fact, only the main event decreased in effectiveness during selection
288 events 3 and 4 (Fig. 6). We hypothesize that the reason for the drop in genomic prediction
289 accuracy in the selfing but not the NoSelf scheme is that, by adding a generation between the
290 training and selection candidate populations, the degree of relationship between those
291 populations was lower in the selfing schemes than in the NoSelf scheme. The degree of
292 relationship has a large impact on prediction accuracy (Clark et al., 2012).

293



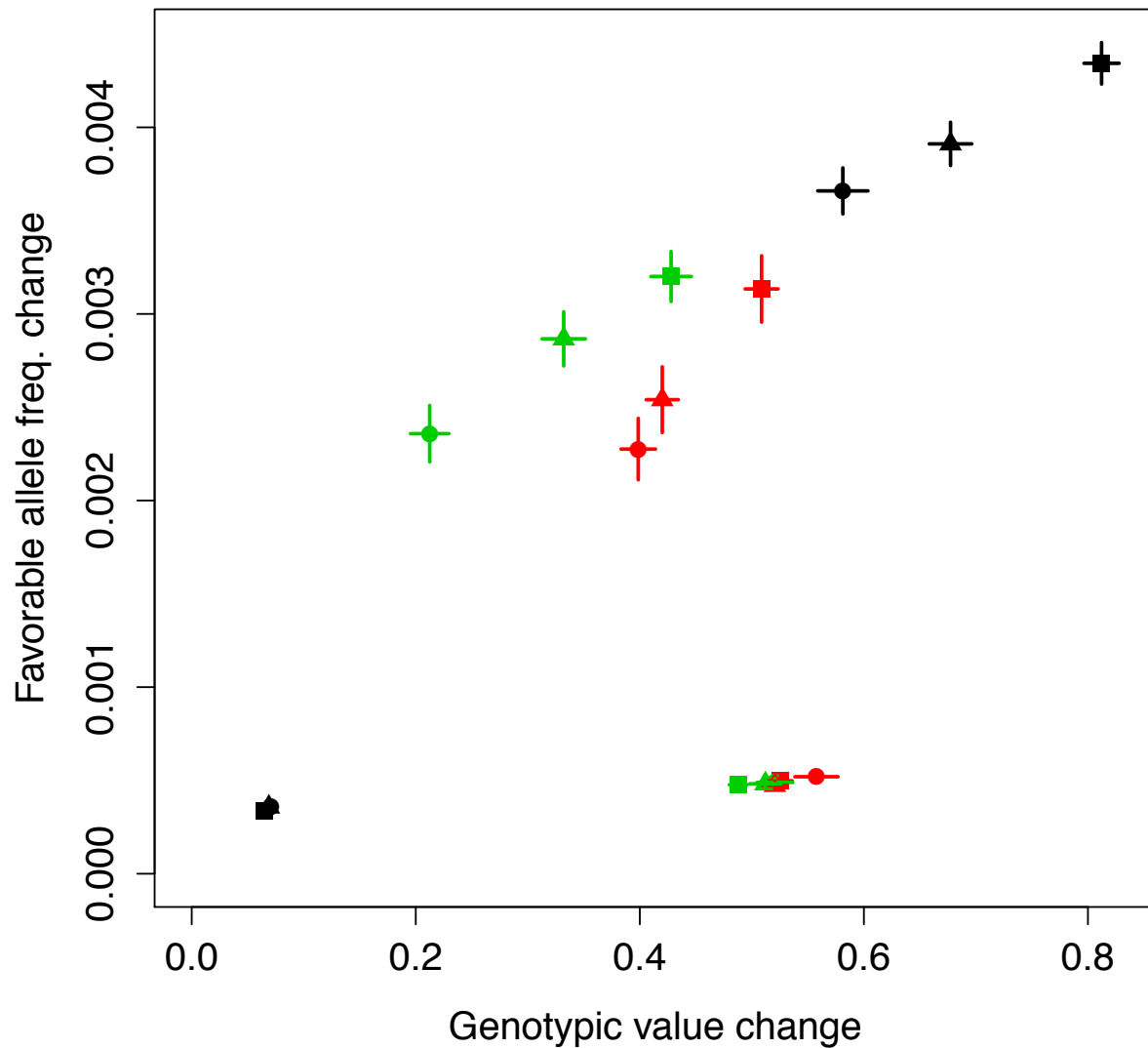
294

295 **Supp. Fig. 1.** Weighted mean of favorable allele frequencies (black, left scale) and deleterious allele homozygote
296 frequencies (red, right scale, note differences in scales) over five cycles of selection. Weights were 0.04, 0.08, 0.16,
297 and 0.30 for the four locus classes, S, M, L, and XL, respectively. Error bars are +/- one standard error across the 24
298 repeated simulations. Left and right columns show the Self and Self+ breeding schemes, respectively. Top, middle,
299 and bottom rows show founder population size of 600 with high selection intensity, and founder population size of
300 200 and 1500 with standard selection intensity, respectively.



301

302 **Fig. 6.** Weighted mean of favorable allele frequencies during the main (black, left scale) and purging (red,
303 right scale, note differences in scales) selection events in the standard selection scheme. Weights were 4,
304 8, 16, and 30 for the four locus classes, S, M, L, and XL, respectively. Error bars are +/- one standard
305 error across the 24 repeated simulations.



306

307 **Fig. 7.** Comparison of genotypic value change versus weighted favorable allele frequency change caused
308 by selection events. Black, red, and green symbols are, respectively, for additive gene action, and
309 dominance with Self and Self+ selection. Round, triangle, and square symbols are, respectively, for founder
310 populations sizes of 200, 600, and 1500. Favorable allele frequency changes under 0.001 occurred during
311 purge selection, while changes above 0.002 occurred during main selection. Whiskers show one standard
312 error above and below the observed value. An absence of whiskers indicates a standard error smaller than
313 the size of the symbol.

314

315 **Inbreeding depression**

316 We compared inbreeding depression before versus after the five cycles of selection by comparing
317 the loss of genetic value caused by one generation of self-fertilization, either on the founder
318 population or on the Cycle 5 population. As expected, inbreeding among founders for all breeding
319 scenarios were not significantly different (given that no breeding treatment had yet been applied).
320 One generation of selfing on a founder population caused a mean decrease of 5.63 in the genetic
321 value, with a standard deviation of 0.38 across replications. After five cycles of selection, there
322 were significant differences among selection schemes, though, surprisingly, the difference was
323 that inbreeding depression was higher under the Self+ scheme while the Self and NoSelf schemes
324 showed similar but lower inbreeding depression (Table 2). Equally surprisingly, the number of
325 founders was not significantly related to decrease in inbreeding depression (Table 2). We
326 hypothesize that the level of inbreeding depression was due not only to the extent of deleterious
327 allele purging, but also to the level of inbreeding previously attained in the population. Thus, there
328 was greater inbreeding depression under the Self+ than the Self scheme not because it was less
329 effective at purging but because the Self+ population was less inbred than the Self population (Fig.
330 3, Supp. Fig. 1). This hypothesis could also explain why inbreeding depression was as great given
331 a 1500 versus a 200 founder population: the former was more effective at purging but also was
332 less inbred than the latter (Supp. Fig. 1).

333

334 **Table 2.** Inbreeding depression from one generation of selfing after five cycles of selection for
335 different breeding schemes. Means across 24 replications are reported and have a standard error
336 of the mean of 0.08.

Num. Founders	Selection Scheme		
	Self+	Self	NoSelf
200	5.24	4.83	4.84
600	5.12	4.76	4.89
1500	5.18	4.73	4.79

337

338 **Discussion**

339 This study investigated the effect of purging recessive deleterious alleles by selection on selfed
340 individuals during five generations of genomic-assisted breeding. Perhaps most obviously, we
341 showed that response under directional dominance differs dramatically from that under an additive
342 model (contrast Fig. 1 and 2). By analyzing changes in favorable allele frequency and deleterious
343 homozygote genotype frequency, we showed that the sudden decline in response under directional
344 dominance was caused by drift fixing deleterious alleles at some loci, even as the favorable allele
345 frequency, averaged over all loci, increased (Fig. 3). These opposing effects led to stagnating or
346 declining genotypic values. Support for this interpretation also came from simulation scenarios
347 with higher selection intensity (hence more drift, Fig. 4) and smaller training population size
348 (hence more co-selection of relatives, Fig. 5).

349 On the strict question of the value of selecting among selfed progeny, our results did not lead to
350 an unambiguous answer. First, selection among selfed progeny was effective in the sense that it
351 led to a large shift in genotypic value among those progeny (Fig. 7). This shift, however, did not
352 lead to a large shift in the favorable allele frequency (Fig. 7). The change in favorable allele
353 frequency during purging selection was only about one fourth of that during main selection (Fig.
354 6). The change in genotypic value was of course affected by the simulation parameters we
355 assumed: relatively high error variance among seedlings and low selection intensity during purging
356 selection. Empirical estimates of error variance or broad sense heritability for yield traits in cassava
357 seedlings have not been reported (e.g., Ozimati et al. 2019). Such estimates would be difficult to
358 obtain given that all seedlings are genetically unique and cannot be replicated. An estimate of the
359 residual from additive effects, which would be straightforward to estimate, would include a
360 component of non-additive genetic variation. In the presence of dominance, that component could
361 be considerable. Our choice of a high error variance on the seedlings was also influenced by the
362 fact that cassava storage root properties differ between seedling and clonal plants (Ceballos et al.,
363 2004), so that lack of genetic correlation between seedling and clonal performance would also
364 contribute to the error of seedling measures to predict clonal performance. Thus, it seemed
365 reasonable to assume that the error variance on a single plant would be substantially greater than
366 on a clonal evaluation performed on multiple plants. We justified the low selection intensity (five
367 plants chosen out of ten grown) only as a matter of labor and cost savings. Obtaining many selfed

368 progeny and planting them all for a large number of selected parents would be logistically
369 challenging.

370

371 While the change in allele frequency from purging selection was low, the increase in genetic
372 value among S_1 progeny from that selection event was similar to that observed in the mean
373 selection (Fig. 7). Under additive gene action (black symbols in Fig. 7), there was a direct
374 relationship between the genotypic value and the favorable allele frequency changes, regardless
375 of main versus purge selection. Under additivity, however, purging selection caused a smaller
376 genotypic value change and thus a smaller favorable allele frequency change (Fig. 7). Under
377 directional dominance gene action, while purging selection caused greater genotypic value
378 change than main selection, it caused little favorable allele frequency change (Fig. 7). Empirical
379 selection studies using selfing have sometimes also observed disappointing gains (Wardyn et al.,
380 2009). Under purging selection, individuals' genotypic values depend on how many deleterious
381 homozygote loci they carry. But there is not a perfect correlation between the number of
382 homozygous loci and the total number of deleterious alleles, which may also be in the
383 heterozygous state. In other words, genotypic value offers an uncertain guide to favorable allele
384 content, such that the strong shift of mean genotypic value under purging did not correspond to a
385 strong shift in favorable allele frequency. We also found that while selection on selfed
386 individuals did lead to a reliable increment in breeding value in each cycle, it also caused the
387 genomic prediction model accuracy to decrease (Fig. 6). We assume that the effect was due only
388 to the addition of a generation between the training population and the selection candidates,
389 making them more distantly related (Clark et al. 2012). That effect erased the benefit of the
390 purging selection. It is unclear whether the accuracy drop would have persisted over further
391 cycles of selection, as the training population began to be updated (Fig. 6). We opted to simulate
392 only five cycles of selection because five cycles actually represent a long period of time, a
393 minimum of ten years. In that period of time, it is unlikely that any single breeding scheme
394 would persist, given technological and breeding priority changes. The value of simulating longer
395 time periods is therefore unclear. Finally, we note that we did not impose any penalty on the
396 selfing schemes for the amount of time that adding the generation of selfing would require. That

397 time would decrease the number of breeding cycles possible per unit time, an effect that has
398 strong negative repercussions (Cobb et al., 2019).

399

400 **Conclusions**

401 Cassava carries a high genetic load (Ramu et al., 2017) and is known to suffer from inbreeding
402 depression (Ceballos et al., 2004). These observations suggest attempting to purge deleterious load
403 by selection on partially inbred individuals could be worthwhile. We simulated one approach to
404 implement purging selection in the context of a breeding program using genomic selection. We
405 observed a favorable initial response (cycles 1 and 2) to adding a generation of selfing to the
406 breeding scheme, but this benefit did not persist beyond those cycles. Over subsequent cycles
407 (cycles 3 to 5), the increase in favorable allele frequency occurring during purging selection only
408 compensated for the loss in accuracy a selfing cycle caused during the main genomic prediction
409 cycle. Thus, despite added cost and overall breeding cycle length from including selfing, no net
410 gain was observed relative to a scheme without selfing.

411 It is difficult to extrapolate from our results to what a cassava breeder may experience empirically
412 because the results do depend both on questions of underlying genetic architecture and on breeding
413 scheme details. Thus, it is difficult to make strong recommendations. We did find, however, a few
414 somewhat non-intuitive results that we believe breeders should keep in mind as they consider
415 breeding scheme modifications. First, we found that the breeding schemes increased favorable
416 allele frequencies even though genotypic values were stagnant (Fig. 1 and Fig. 3). While we
417 believe that this phenomenon can occur generally, it depended on a specific aspect of the
418 simulation that will differ across real breeding programs. In particular, our simulated founders
419 came from a population with an effective size of 500, whereas the standard breeding scheme had
420 an upper bound to the effective population size of 60 due to the number of parents randomly mated
421 in each cycle. This rapid downward change led to greater drift and therefore fixation of the
422 deleterious allele at some loci. Thus, the lack of gain from the first to the fifth cycles (Fig. 1) in
423 our simulation may not be generalizable to practicing breeding programs.

424 An observation that we believe will be generalizable was that introducing a generation of selfing
425 into the scheme will have conflicting impacts: on the one hand, it will provide some genetic gain
426 itself. On the other hand, it will distance selection candidates from the genomic prediction training

427 population, thereby causing a decrease in accuracy. In our simulation, these conflicting impacts
428 balanced out beyond the first selection cycle so that the selfing schemes generated about equal
429 gain to the scheme without selfing (Fig. 1, 4, and 5). A related observation was that rapid additions
430 to the training population in the Self+ scheme caused decreased inbreeding (Fig. 4, 5, and Supp.
431 Fig. 1). We also think that this result is generalizable. We do not have an explanation as to why
432 this effect did not lead to greater gains per cycle in the Self+ than the Self scheme. We do think
433 the effect lead to greater potential for inbreeding depression at the end of the five selection cycles
434 (Table 2), simply because the standing population was less inbred to begin with.
435 Finally, we observed that selection on partially inbred individuals led to a large gain in genotypic
436 value among those individuals, but that that large gain was not reflected in a large gain in favorable
437 allele frequency (Fig. 7). Under strong directional dominance, genotypic value may not be a good
438 guide to breeding value. It is worth considering this effect more carefully. Changes in genotypic
439 value are what breeders observe directly and it therefore guides their intuition. Here we show
440 through simulation that selection on partially inbred individuals may not lead to as strong gain in
441 their outcrossed progeny as we might intuit.

442

443 **Acknowledgments**

444 The research conducted at Cornell University was part of NextGen cassava project funded by
445 Bill & Melinda Gates foundation and UKaid (Grant 1048542). The research was conducted using
446 the resources of the Cornell University Institute of Biotechnology Bioinformatics Facility
447 (BioHPC). The authors gratefully acknowledge the support of BioHPC staff.

448 **References**

- 449 Boakes, E.H & Wang, J. 2005. A simulation study on detecting purging of inbreeding depression
450 in captive populations. *Genetical Research*, 86(02):139-148.
- 451 Ceballos, H., Iglesias, C. A., Pérez, J. C., & Dixon, A. G. (2004). Cassava breeding: opportunities
452 and challenges. *Plant molecular biology*, 56(4), 503-516
- 453 Ceballos, H., Sánchez, T., Morante, N., Fregene, M., Dufour, D., Smith, A. M., ... & Mestres, C.
454 (2007). Discovery of an amylose-free starch mutant in cassava (*Manihot esculenta* Crantz).
455 *Journal of Agricultural and Food Chemistry*, 55(18), 7469-7476.
- 456 Charlesworth, D., & Willis, J. H. 2009. The genetics of inbreeding depression. *Nature Reviews*
457 *Genetics*, 10(11):783-796.
- 458 Clark, S.A., J.M. Hickey, H.D. Daetwyler, and J.H.J. van der Werf. 2012. The importance of
459 information on relatives for the prediction of genomic breeding values and the implications
460 for the makeup of reference data sets in livestock breeding schemes. *Genet. Sel. Evol.* 44:
461 4.
- 462 Cobb, J.N., R.U. Juma, P.S. Biswas, J.D. Arbelaez, J. Rutkoski, et al. 2019. Enhancing the rate
463 of genetic gain in public-sector plant breeding programs: lessons from the breeder's
464 equation. *Theor. Appl. Genet.* 132:627–645.
- 465 de Freitas, J.P.X., V. da Silva Santos, and E.J. de Oliveira. 2016. Inbreeding depression in
466 cassava for productive traits. *Euphytica*, 209(1):137-145.
- 467 Duenk, P., P. Bijma, M.P.L. Calus, Y.C.J. Wientjes, and J.H.J. van der Werf. 2020. The Impact
468 of Non-additive Effects on the Genetic Correlation Between Populations. *G3* 10(2): 783–
469 795.
- 470 Endelman, J.B. 2011. Ridge regression and other kernels for genomic selection with R package
471 rrBLUP. *The Plant Genome*, 4(3):250-255.
- 472 Falconer, D.S., and T.F.C. Mackay. 1996. Introduction to Quantitative Genetics. Longman Sci.
473 and Tech., Harlow, U.K.
- 474 Haller, B. C., and P. W. Messer. 2017. SLiM 2: Flexible." *Interactive Forward Genetic*.
- 475 International Cassava Genetic Map Consortium (ICGMC). 2014. High-resolution linkage map
476 and chromosome-scale genome assembly for cassava (*Manihot esculenta* Crantz) from 10
477 populations. *G3* 5(1): 133–144.
- 478 Jannink, J.-L. 2010. Dynamics of long-term genomic selection. *Genet. Sel. Evol.* 42(1): 35.
- 479 Jannink, J.L., Lorenz, A.J. and Iwata, H., 2010. Genomic selection in plant breeding: from theory
480 to practice. *Briefings in functional genomics*, 9(2), pp.166-177.
- 481 Kawuki, R. S, Nuwamanya E, Labuschagne MT, Herselman L, Ferguson M (2011) Segregation
482 of selected agronomic traits in six S1 cassava families. *J Plant Breed Crop Sci* 3:154–160
- 483 Lande, R., Schemske, D. W., & Schultz, S. T. (1994). High inbreeding depression, selective
484 interference among loci, and the threshold selfing rate for purging recessive lethal
485 mutations. *Evolution*, 965-978.
- 486 Lorenz, A.J., Chao, S., Asoro, F.G., Heffner, E.L., Hayashi, T., Iwata, H., Smith, K.P., Sorrells,
487 M.E. and Jannink, J.L. 2011. Genomic selection in plant breeding: knowledge and
488 prospects. In *Advances in agronomy* 110:77-123.
- 489 Lorenz, A. J. 2013. Resource allocation for maximizing prediction accuracy and genetic gain of
490 genomic selection in plant breeding: a simulation experiment. *G3*, 3:481–91
- 491 Muir, W.M. 2007. Comparison of genomic and traditional BLUP-estimated breeding value
492 accuracy and selection response under alternative trait and genomic parameters. *J. Anim.*
493 *Breed. Genet.* 124(6): 342–355

- 494 Nuwamanya, E., Herselman, L. and Ferguson, M. 2011. Segregation of selected agronomic traits
495 in six S1 cassava families. *Journal of Plant Breeding and Crop Science* **3**, 154-160.
- 496 Ozimati, A., Kawuki, R., Esuma, W., Kayondo, S.I., Pariyo, A., Wolfe, M. and Jannink, J.L.,
497 2019. Genetic Variation and Trait Correlations in an East African Cassava Breeding
498 Population for Genomic Selection. *Crop Science* **59**: 460–473
- 499 Prochnik, S., Marri, P. R., Desany, B., Rabinowicz, P. D., Kodira, C., Mohiuddin, M., and
500 Rokhsar, D. S. (2012). The cassava genome: current progress, future directions. *Tropical*
501 *plant biology*, *5*(1), 88-94
- 502 Ramu, Punna, Williams Esuma, Robert Kawuki, Ismail Y. Rabbi, Chiedozi Egesi, Jessen V.
503 Bredeson, Rebecca S. Bart, Janu Verma, Edward S. Buckler, and Fei Lu. "Cassava
504 haplotype map highlights fixation of deleterious mutations during clonal propagation."
505 *Nature Genetics* (2017).
- 506 Rojas, M. C., Pérez, J. C., Ceballos, H., Baena, D., Morante, N., & Calle, F. (2009). Analysis of
507 Inbreeding Depression in Eight S Cassava Families. *Crop Science*, *49*(2), 543-548.
- 508 Rutkoski, J., R.P. Singh, J. Huerta-Espino, S. Bhavani, J. Poland, et al. 2015. Genetic Gain from
509 Phenotypic and Genomic Selection for Quantitative Resistance to Stem Rust of Wheat.
510 *Plant Genome* **8**. doi: 10.3835/plantgenome2014.10.0074.
- 511 Valluru, R., Gazave, E.E., Fernandes, S.B., Ferguson, J.N., Lozano, R., Hirannaiah, P., Zuo, T.,
512 Brown, P.J., Leakey, A.D., Gore, M.A. and Buckler, E.S., 2019. Deleterious Mutation
513 Burden and Its Association with Complex Traits in Sorghum (*Sorghum bicolor*).
514 *Genetics*, *211*(3):1075-1087.
- 515 Wardyn, B.M., J.W. Edwards, and K.R. Lamkey. 2009. Inbred-progeny selection is predicted to
516 be inferior to half-sib selection for three maize populations. *Crop Sci.* **49**:443-450.
- 517 Wolfe, M.D., D.P. Del Carpio, O. Alabi, L.C. Ezenwaka, U.N. Ikeogu, et al. 2017. Prospects for
518 Genomic Selection in Cassava Breeding. *Plant Genome*. doi:
519 10.3835/plantgenome2017.03.0015.
- 520 Yabe, S., Iwata, H. and Jannink, J.L., 2017. A simple package to script and simulate breeding
521 schemes: the breeding scheme language. *Crop Science*, *57*(3), pp.1347-13
522