Genomic selection strategies to increase genetic gain in tea breeding programs

Nelson Lubanga¹, Gregor Gorjanc¹, Festo Massawe², Sean Mayes³, and Jon Bancic¹

¹Roslin Institute, The University of Edinburgh, Easter Bush Campus, Midlothian EH25 9RG, UK ²The University of Nottingham Malaysia, Jalan Broga, 43500 Semenyih, Selangor Darul Ehsan, Malaysia

³School of Biosciences, The University of Nottingham, Sutton Bonington Campus, Loughborough, Leicestershire LE12 5RD, UK

November 26, 2022

Abstract

Genomic selection (GS) can improve the efficiency of tea breeding compared to phenotypic selection (PS) by shortening the generation interval, increasing selection accuracy, and shortening the duration of the entire breeding program, especially at early stages. Tea (Camellia sinensis (L.) O. Kuntze) is mainly grown in low- to middle-income countries (LMIC) and is a global commodity. Breeding programs in these countries face the challenge of increasing genetic gain because the accuracy of selecting superior genotypes is low and resources are limited. Recurrent phenotypic selection has traditionally been the primary method for developing improved tea varieties and can take over 16 years. Therefore, the main objective of this study was to investigate the potential of implementing GS in tea breeding programs to speed up genetic progress despite the low labour costs in LMIC. We used stochastic simulations to compare three GS breeding programs with a commercial PS program over a 40-year breeding period. All GS breeding programs achieved higher genetic gains compared to PS. Seed-GSconst, in particular, proved to be the most cost-effective strategy for introducing GS into tea breeding programs. It introduces GS at the nursery stage, thereby increasing the predictive accuracy at the early stage of the breeding program. It also shortens the duration of the entire breeding program by three years and reduces the generation interval to two years. Our results indicate that GS is a promising strategy to improve genetic gain per unit time and cost in tea breeding programs.

1	Genomic selection strategies to increase genetic gain in tea breeding programs
2	Nelson Lubanga ¹ , Gregor Gorjanc ¹ , Festo Massawe ² , Sean Mayes ³ , Jon Bančič ¹
3	
4	¹ Roslin Institute, The University of Edinburgh, Easter Bush Campus, Midlothian EH25 9RG, UK
5	² The University of Nottingham Malaysia, Jalan Broga, 43500 Semenyih, Selangor Darul Ehsan,
6	Malaysia
7	³ School of Biosciences, The University of Nottingham, Sutton Bonington Campus, Loughborough,
8	Leicestershire LE12 5RD, UK
9	
10	Core ideas
11	• Genomic selection could result in 1.6 times greater genetic gain than phenotypic selection in
12	tea breeding programs.
13	• All genomic selection strategies achieved higher genetic gains per unit time and cost than
14	phenotypic selection method.
15	• Seed-GSconst is the most cost-effective strategy for introducing GS into tea breeding programs.
16	
17	
18	
19	
20	
21	
22	
23	
24	

25 Abstract

26 Genomic selection (GS) can improve the efficiency of tea breeding compared to phenotypic 27 selection (PS) by shortening the generation interval, increasing selection accuracy, and shortening the 28 duration of the entire breeding program, especially at early stages. Tea (Camellia sinensis (L.) O. 29 Kuntze) is mainly grown in low- to middle-income countries (LMIC) and is a global commodity. 30 Breeding programs in these countries face the challenge of increasing genetic gain because the 31 accuracy of selecting superior genotypes is low and resources are limited. Recurrent phenotypic 32 selection has traditionally been the primary method for developing improved tea varieties and can take 33 over 16 years. Therefore, the main objective of this study was to investigate the potential of 34 implementing GS in tea breeding programs to speed up genetic progress despite the low labour costs 35 in LMIC. We used stochastic simulations to compare three GS breeding programs with a commercial 36 PS program over a 40-year breeding period. All GS breeding programs achieved higher genetic gains 37 compared to PS. Seed-GSconst, in particular, proved to be the most cost-effective strategy for 38 introducing GS into tea breeding programs. It introduces GS at the nursery stage, thereby increasing 39 the predictive accuracy at the early stage of the breeding program. It also shortens the duration of the 40 entire breeding program by three years and reduces the generation interval to two years. Our results 41 indicate that GS is a promising strategy to improve genetic gain per unit time and cost in tea breeding 42 programs.

43

44 Abbreviations

ACT, Advanced Clonal Trial stage; ECT-GS, Elite Clonal Trials Genomic Selection Breeding
Program; ECT, Elite Clonal Trial stage; GS, Genomic selection; LMIC, low- to middle-income
countries; PS, Phenotypic selection; Seed-GSconst, Constrained Seedlings Genomic Selection
Breeding Program; Seed-GSunconst, Unconstrained Seedlings Genomic Selection Breeding Program

2

49 **1** Introduction

50 Tea (Camellia sinensis (L) O. Kuntze) is mainly grown in tropical and subtropical regions in 51 low- to middle-income countries (LMIC) (Han, Li, and Ahammed 2018). It is an important crop for 52 the economies of these countries as it provides a source of income for many smallholder farmers and 53 those employed in tea processing companies (Mukhtar and Ahmad 2000). Additionally, the tea 54 growing areas (mainly rural) have benefitted from improved social infrastructure such as good road 55 networks, schools and hospitals. All tea varieties currently grown in the world originated in India and 56 China and were either directly or indirectly imported from these two countries to other countries 57 (Meegahakumbura et al. 2016). The world population is steadily increasing and is expected to reach 9 58 billion people by 2050 (Perroy 2015) leading to an increase in demand for food and beverages (Valin 59 et al. 2014). Conventional tea breeding is well established in the major tea growing countries such as 60 India, China and Kenya and has led to the development of many superior varieties (Meegahakumbura 61 et al. 2016). Tea varieties developed through breeding have superior yield, quality and are resistant to 62 drought compared to seedling genotypes (Corley and Tuwei 2018). However, to sustain long-term tea 63 production and the increasing demand for tea, tea breeders need to continuously bring new improved 64 varieties to the market. The objectives of tea breeding vary among major tea-growing countries, 65 depending on local needs. However, the most important breeding objective is to develop varieties with 66 high yield and improved quality (colour, aroma, taste and mouthfeel) (Kamunya et al. 2012, Mondal 67 2014). Currently, tea productivity is seriously threatened by climate change, which is already causing 68 yield losses (Gunathilaka, Smart, and Fleming 2017, Sitienei, Juma, and Opere 2017), and decreased 69 quality (Han et al. 2017). Climate change has led to extreme and unpredictable weather patterns, 70 resulting in longer dry spells, heavy rainfall, more hail, and higher temperatures (Marx, Haunschild, 71 and Bornmann 2017, Batley and Edwards 2016). Additionally, the changing climate has led to 72 increased attacks of pests and diseases. Therefore, effective tea breeding strategies are needed to develop high-yielding and high-quality tea varieties that are also tolerant to biotic and abiotic stresses
(Mondal 2011, Muoki et al. 2020).

75 Tea breeding programs use recurrent phenotypic selection (PS) to select the best individuals 76 based on phenotypic values estimated from the *per se* performance of clones in clonal evaluation trials. 77 This involves the creation of genetic variation through crossing, followed by many years of recurrent 78 selection aimed at determining the genetic value of promising genotypes, leading to the identification 79 of new parents for crossing and the release of commercial varieties to farmers. In the initial phase of 80 the breeding program, new genotypes are first tested as seedlings in single bush (preliminary) trials. 81 Then, genetically identical teas (clones) are generated from selected seedlings through clonal 82 propagation (cuttings), allowing genotypes to be tested in clonal plots in multiple replications, at 83 multiple locations, and in multiple years (Carr 2018). PS has been quite successful in delivering 84 improved tea varieties over many years (Mondal 2014). However, this is a time-consuming process as 85 it takes 16 years to develop a tea variety from cross to release (Figure 1).

86 In modern times, plant breeding has started to move from complete reliance on PS to genomic-87 assisted selection due to improved molecular biology and high-throughput genotyping technologies 88 (Leng, Lübberstedt, and Xu 2017). Quantitative trait loci (QTL) mapping (Kamunya et al. 2010, 89 Malebe et al. 2021) and association mapping (Jin et al. 2016, Fang et al. 2021) have been tested in the 90 genetic improvement of tea, and several QTLs associated with yield and quality have been identified. 91 However, these QTLs have not been successfully applied in the genetic improvement of tea (Xia et al. 92 2020), since marker-assisted selection (MAS) methods do not account for the effects of minor QTLs 93 influencing quantitative traits. Only a few major QTLs have been identified in tea (Fang et al. 2021, 94 Yamashita et al. 2020), while the minor QTLs which influence important quantitative traits are ignored 95 (Heffner, Sorrells, and Jannink 2009). Many complex traits, including yield, quality and drought 96 tolerance are controlled by many genes with small effects, and therefore MAS is of limited use due to 97 low statistical power to detect individual QTLs (Bernardo and Yu 2007).

98 Genomic selection (GS) uses all available (genome-wide) markers to predict breeding values 99 (Meuwissen, Hayes, and Goddard 2001) and offers a great potential for identifying the best parents for 100 crossing and superior clones for variety development in tea breeding programs. Genomic estimated 101 breeding values (GEBVs) are calculated by summing marker effects that may or may not be in linkage 102 disequilibrium with one or more QTLs across the entire genome (Bernardo and Yu 2007). GS uses a 103 prediction model that is first trained using a population of genotyped and phenotyped individuals. The 104 trained model is then used to predict GEBVs of selection candidates with genotyping information but 105 no phenotypes. By then correlating estimated phenotypic values based on GEBVs with the actual 106 phenotypic data, it is possible assess the accuracy of the genomic selection model (Heffner, Sorrells, 107 and Jannink 2009, Meuwissen, Hayes, and Goddard 2001). For example, Lubanga, Massawe, and 108 Mayes (2021) investigated the potential use of GS to improve tea quality. They reported higher 109 prediction accuracies for all genomic prediction models compared to the pedigree model. Similar 110 findings were also reported by Yamashita et al. (2020), who also investigated the potential of GS for 111 improving tea quality. They found moderate prediction accuracies for the 6 GS models tested. In a tea 112 breeding program, GS can be used in four ways:

- 1131. to reduce the generation interval as new parents can be selected at the Seedlings stage.114Genotypes in the nursery can be genotyped. Superior genotypes can then be selected based on115GEBVs and planted in the germplasm garden for population improvement (Figure 2 and Figure
- 116

3),

- 117 2. to increase the accuracy of selecting superior tea genotypes at the Seedlings stage,
- 118 3. to increase the selection intensity. More seedlings can be genotyped at the nursery stage andpromising ones predicted accurately compared to PS,
- 4. to shorten the entire breeding program by eliminating some of the stages in a breeding program
 to enable faster release of varieties (Figure 2 and Figure 3).

5

However, the implementation of GS in LMIC faces limitations. In most of these countries, the cost of phenotypic selection is much lower compared to Europe and North America because the local population provides cheap labour. In addition, most breeding programs have limited investment budgets for conducting research. In addition, there is a lack of qualified personnel who are trained and understand the technique of GS and its practical implementation in breeding programs. The implementation of GS in tea breeding should therefore take into account the particular challenges of these programs.

129 Plant breeders have traditionally relied on field trial experiments to inform their decisions 130 (Rutkoski et al. 2015). However, evaluating these field trials takes a long time and is also expensive 131 (Wang and Wolfgang H 2007). Simulations are useful in determining the best breeding strategies and 132 can also be used to study the genetic gain, predictive accuracy, and cost-effectiveness of GS under 133 different scenarios (Gaynor, Gorjanc, and Hickey 2021). Stochastic simulations have been conducted 134 for many crops, including wheat (Gaynor et al. 2017), clonally propagated crops (Werner et al. 135 2020), maize (Powell et al. 2020), sorghum (Muleta, Pressoir, and Morris 2019), and trees (Iwata, 136 Hayashi, and Tsumura 2011). However, to our knowledge, no simulation studies have been published 137 integrating GS into tea breeding programs to investigate their feasibility and long-term outcomes. 138 This study aims to test the feasibility of implementing GS into a tea breeding program using 139 stochastic simulations. To this end, we used a PS breeding program as a baseline in which the number 140 of crosses, seedlings, replicates, and locations mimicked an actual commercial tea breeding program 141 (Unilever Tea Kenya). We estimated variance parameters from real field data. We developed three new 142 breeding programs based on the PS breeding program that integrated GS. Using simulations, we then 143 compared the baseline PS breeding program with three GS breeding programs. All simulations were 144 performed using AlphaSimR (Gaynor, Gorjanc, and Hickey 2021). Our objectives in this study were 145 (i) to investigate the potential of implementing GS in tea breeding programs despite the limited 146 breeding program resources and low labour costs in LMIC, (ii) to compare different strategies for 6

implementing GS in tea breeding programs at the same cost as PS, (iii) to investigate whether shortening the tea breeding generation interval and the entire breeding program duration by incorporating GS in breeding programs leads to higher genetic gains. In addition, we also evaluated two different strategies of parent selection.

151

2 MATERIALS AND METHODS

We used stochastic simulations to evaluate the possibility of implementing GS in tea breeding programs. We compared a PS breeding program and three breeding strategies incorporating GS. We subdivided the materials and methods section into simulation of the founder genotype population and simulation of the breeding programs.

- 156 We simulated the founder genotype population as follows:
- 157 i. Genome simulation: a genome sequence was simulated for a hypothetical diploid tea species
 158 (*Camellia sinensis* (L) O. Kuntze).
- 159 ii. Simulation of founder genotypes: the simulated genome sequences were used to generate a160 base population of 20 diploid founder genotypes.
- iii. Simulation of genetic values: a single trait representing yield was simulated for all founder
 genotypes by summing the additive effects with 2400 quantitative trait nucleotides (QTN).
- iv. Simulation of phenotypes: the phenotypes of all founder genotypes were simulated by adding
 random error to the total genetic value of the tea genotypes.
- 165 We simulated the breeding programs as follows:
- 166 i. Recent (burn-in) breeding phase: a PS breeding program for tea was simulated for a period of
 167 40 years (burn-in) to provide a shared starting point for the future breeding phase.
- 168 ii. Future breeding phase: three different GS breeding programs were simulated and compared to
- 169 the PS breeding program for an additional 40 years of breeding.

170 **2.1** Simulation of the founder genotype population

7

171 **2.1.1 Genome simulation**

172 We simulated a genome sequence with 15 pairs of chromosomes for a hypothetical diploid tea 173 species (Camellia sinensis (L) O. Kuntze). We then assigned a physical length of 10⁸ base pairs and a 174 genetic length of 1 Morgans to these chromosomes. We generated the chromosome sequences using 175 the Markovian coalescent simulator (MaCS) (Chen, Marjoram, and Wall 2009) implemented in 176 AlphaSimR (Gaynor, Gorjanc, and Hickey 2021). We estimated recombination rate as the ratio 177 between genetic length and physical genome length (i.e., 1 Morgans / 10^8 base pairs = 10^{-8}). We set the per-site mutation rate to 2.5 x 10^{-8} mutations per base pair. We set the effective population size (Ne) 178 179 to 100, as described by Werner et al. (2020).

180 **2.1.2 Simulation of founder genotypes**

We used the simulated genome sequences to generate a base population of 20 diploid founder genotypes in Hardy-Weinberg equilibrium. These genotypes were formed by randomly sampling 15 chromosome pairs per genotype. A set of 160 biallelic quantitative trait nucleotides (QTNs) and 600 single nucleotide polymorphisms (SNPs) were randomly selected along each chromosome, to simulate a quantitative trait that was controlled by 2400 QTN and an SNP marker array with 9000 genome-wide SNP markers.

187 **2.1.3 Simulation of genetic values**

We simulated genetic values for a single trait representing yield by summing the additive genetic effects at 2,400 randomly sampled QTN. We sampled additive genetic effects (*a*) from the standard normal distribution and scaled them to obtain an additive genetic variance of $\sigma_a^2 = 1$ in the founder population, as described in detail by Gaynor, Gorjanc, and Hickey (2021). The environmental effect represented the environmental component of the genotype by year (G x Y) interaction and was sampled for each year of the simulation from the standard normal distribution as described by Gaynor et al. (2017).

195 **2.1.4 Simulation of phenotypes**

We calculated the phenotypic values for yield by adding G x Y and random error to the additive genetic values. Therefore, the phenotypic value of genotype i grown in stage k of a breeding program in year j was calculated as;

$$y_{ij} = g_i + (gy)_{ij} + e_{ijk}$$

where g_i is the additive genetic value of genotype *i*; $(gy)_j$ is G x Y interaction effect associated with 200 genotype *i* and year *j*; and e_{ij} is error associated with genotype *i*, year *j*, and stage *k*. The random error 201 was sampled from the standard normal distribution with mean zero and an error variance σ_e^2 defined 202 203 by the target level of heritability at each testing stage of the tea breeding program. In the founder 204 population, we calculated the entry-mean values based on real data from Unilever Tea Kenya breeding program for narrow-sense heritability (h^2) at each of the breeding stages. The h^2 at the Seedling and 205 206 PT stages was 0.05, 0.45 in the advanced clonal testing stage (ACT) and 0.65 in the elite clonal testing 207 (ECT) stage. Narrow-sense heritabilities in later testing stages were higher because of the increased 208 number of replicates per genotype. We calculated narrow-sense heritability as

$$\frac{\sigma_a^2}{(\sigma_a^2 + \sigma_{gy}^2 / e + \sigma_e^2 / er)},$$

where σ_a^2 is the additive genetic variance, σ_{gy}^2 is the G x Y interaction variance, σ_e^2 is the residual variance and *e* and *r* are the number of environments and replicates within each environment, respectively.

213 **2.2 Recent (burn-in) breeding phase**

We simulated a PS breeding program over a 40-year period (burn-in) to establish a common baseline for the future tea breeding phase. The structure of the PS breeding program was based on the Unilever Tea Kenya breeding program (Figure 1). A description of the Unilever Tea Kenya breeding program can also be found in Corley and Tuwei (2018). To fill the breeding pipeline and provide a

218 starting point for the burn-in phase, we performed 16 crossing and selection cycles prior to the burn-in 219 phase. Each of these 16 cycles started with the same twenty founder genotypes in the crossing block 220 to perform 100 bi-parental crosses and 100 pollinations per cross (10,000 crosses in total). Based on 221 our experience, we assumed that 2,000 (approximately 20%) seedlings germinated and were grown in 222 the nursery for one year. In the third year, seedlings were planted in the field as preliminary trials (PT), 223 followed by a three-year evaluation period. Five hundred superior clones were selected and planted as 224 advanced clonal trials (ACT) and yield data were recorded for 5 years. Forty high yielding clones were 225 selected, advanced to the elite clonal trial (ECT) and yield data recorded for 6 years (Figure 1). 226 Selection of new parents and best clones in each testing stage were based on phenotypic records. In the 227 burn-in phase, the selection of new parents was done at the ECT stage in the year 16. Each year, we 228 replaced the 5 genotypes in the crossing block with the oldest *per se* performance with new high 229 yielding varieties from the ECT stage. The total duration of the PS breeding program was 16 years 230 (Figure 1).



231

Figure 1. Schematic overview of the phenotypic selection breeding program (PS). This program is based on the commercial breeding program (Unilever Tea Kenya). The solid line represents the stage at which the 5 or 20 new parents are selected based on phenotypic information. PT, Preliminary Trial stage; ACT, Advanced Clonal Trial stage and ECT, Elite Clonal Trial stage.

236 2.3 Future Breeding Phase

237 We used the future breeding phase to evaluate the PS breeding program and the three GS 238 breeding programs. We simulated each breeding program for an additional 40 years after the burn-in 239 breeding phase to evaluate each program with an equivalent starting point. The three GS strategies 240 were Seed-GSconst, Seed-GSunconst, and ECT-GS (see their description below). The GS programs 241 replaced PS with GS at different stages of the PS breeding program. The costs of the three GS strategies 242 were equalized to the estimated cost of the PS breeding program (\$71,880). Table 1 shows the sizes 243 and costs of the breeding programs. Equalization of operating costs in GS programs was done using 244 the estimated costs of genotyping and reducing program sizes at different breeding stages. Briefly, the 245 PT stage was eliminated for the Seed-GSconst and Seed-GSunconst programs, while the PT and ACT

- 246 stages were eliminated for the ECT-GS program. A summary of the key differences between the
- 247 breeding programs can be found in Table 2. A complete schematic description of the programs is shown
- 248 in Figures. 1, 2, 3. We assumed that the cost of genotyping per individual was \$15
- 249 (http://techservicespro.com/test-locations/). Phenotyping costs were estimated based on the Unilever
- 250 Tea Kenya breeding program.

251 Table 1. Summary of the tea breeding program sizes and annual costs of simulated breeding

252 programs.

Breeding program	Number of parents	Seedlings	РТ	ACT	ЕСТ	Cost (\$)
PS	20	2,000	2000	500	40	71,880
Seed-Gsconst	20	800	0	300	40	69,980
Seed- GSunconst	20	2,000	0	500	40	100,980
ECT-GS	20	800	0	0	90	72,970

253 PT, Preliminary Trial stage; ACT, Advanced Clonal Trial stage; and ECT, Elite Clonal Trial stage; PS, 254 phenotypic selection breeding program; GS, genomic selection; Seed-GSconst, seedlings GS breeding

- phenotypic selection breeding program, OS, genomic selection, seed-Osconst, seedings OS breeding
 program; Seed-GSunconst, seedlings GS breeding program with unconstrained budget; and ECT-GS,
- elite clonal GS breeding program.
- 257
- 258 We also compared two parent replacement methods for each strategy, namely:
- 1. replacing 25% of the parents after each breeding cycle,
- 260 2. replacing all the parents after each breeding cycle.

261 **Table 2. Summary of the key differences between the four breeding programs.**

Breeding Program	Parent selection stage	Number of parents	Generation interval / program duration (years)	Parent selection	Key features	
PS	ECT	15 old, 5 best parents	16/16	Phenotype	Conventional breeding	
Seed- GSconst	Seedlings	15 old, 5 new parents	2/13	GS	PT stage removed; 800 seedlings genotyped	
Seed- GSunconst	Seedlings	15 old, 5 new parents	2/13	GS	PT stage removed; 2,000 seedlings genotyped	
ECT-GS	Seedlings	15 old, 5 new parents	2/8	GS	PT and ACT stages removed; increased	

						number of clones
						tested in ECT stage
262	PT, Preliminar	y Trial stage;	ACT, Adva	nced Clonal Trial stage;	and ECT, Eli	te Clonal Trial stage; PS,
263	phenotypic sel	ection breeding	ng program;	GS, genomic selection;	Seed-GScon	st, seedlings GS breeding
264	program; Seed	-GSunconst,	seedlings G	S breeding program wit	h unconstrain	ed budget; and ECT-GS,
265	elite clonal GS	breeding pro	gram.			
266	2.3.1 Constrai	ined Seedling	gs Genomic	Selection Breeding Pr	ogram (Seed	-GSconst)
267	The Se	ed-GSconst p	orogram intr	oduced genotyping and	GS at the ear	liest Seedlings stage and
		Ĩ	C			
268	eliminated the	PT stage (Fig	gure 2). Eigl	ht hundred (8 per family) seedlings fr	om the 2,000 germinated
269	seeds in the nu	irsery were ra	indomly sel	ected for genotyping. G	enomic select	tion was used to advance
270	~~~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~			t the best 5 or 20 correts		the aldest nonents in the
270	genotypes to tr	ie next stage a	and to selec	t the best 5 or 20 genoty	pes to replace	e the oldest parents in the
271	crossing block	The PT stag	e was elimi	nated to avoid three ver	ors of field tes	sting and hence GS was
271	crossing brock	· · · · · · · · · · · · · · · · · · ·		indica to avoia tinee yet		sting, and nonce, ob was
272	used to advance	the best 30	0 genotypes	s from the Seedlings sta	ge to the AC	Γ stage. Yield trials were
			0 11	C .	C	0
273	recorded at the	ACT stage for	or 5 years. C	Genomic selection was u	sed to advance	e 40 promising clones to
274	the ECT stage.	The yield tria	als at the EC	CT stage were recorded t	for another 6	years. The Seed-GSconst
075		2		1 11 1 6		
275	program has a	2-year gener	ration interv	al and lasts a total of	13 years, whi	ch is three years shorter
276	a amount to the	DC handler	~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~ ~			
210	compared to th	ie PS breeding	g program.			



277

Figure 2. Schematic overview of the seedlings GS breeding program with constrained costs (Seed-GSconst) and unconstrained costs (Seed-GSunconst). The solid line represents the stage at which the 5 or 20 new parents are selected based on genomic prediction. The values outside parenthesis relate to Seed-GSconst and the values inside parenthesis relate to Seed-GSunconst. ACT, Advanced Clonal Trial stage; and ECT, Elite Clonal Trial stage.

283 2.3.2 Unconstrained Seedlings Genomic Selection Breeding Program (Seed-GSunconst)

The Seed-GSunconst program used a similar strategy to the Seed-GSconst program, only that it used an increased operating budget (Figure 2). There are two main differences between the two programs: (i) all 2,000 seedlings inside the Seedlings stage were genotyped and predicted with GS, and (ii) GS is used to advance a total of 500 genotypes (instead of 300) to the ACT stage. The Seed-GSunconst program also has a 2-year generation interval and lasts a total of 13 years, which is three years shorter compared to the PS breeding program.

290 **2.3.3 Elite Clonal Trials Genomic Selection Breeding Program (ECT-GS)**

291 The ECT-GS program introduced genotyping and GS at the earliest Seedlings stage and 292 eliminated the PT and ACT stages (Figure 3). The two stages were eliminated to avoid eight years of 293 testing and to reallocate the resources into genotyping of all 2,000 seedlings at the Seedlings stage as 294 in the PS program. Compared to the previous two GS programs, the seedlings were planted in the 295 nursery for an additional year to produce enough cuttings for direct planting at the ECT stage. Genomic 296 selection was then used to advance genotypes to the next stage and to select the best 5 or 20 genotypes 297 to replace the oldest parents in the crossing block. Ninety promising clones were advanced from the 298 Seedlings stage to the ECT stage, where they were evaluated for 6 years. The ECT-GS breeding has a 299 2-year generation interval and lasts a total of eight years, which is eight years shorter compared to the 300 PS breeding program.



Figure 3. Schematic overview of the elite clonal breeding program with GS (ECT-GS). The solid
line represents the stage at which the 5 or 20 new parents are selected based on genomic prediction.
ACT, Advanced Clonal Trial stage; and ECT, Elite Clonal Trial stage.

305 2.4 Training population & Genomic Selection Model

301

306 To initialize the training population, the last 6 years of the burn-in phase were used to collect 307 phenotypic data at the ACT stage for training the GS model (in year 41 and onwards). The initial 308 training population consisted of 6,000 phenotypic records. In subsequent years, new phenotypic data 309 were added to the training population as new yield trials were recorded. For the Seed-GSunconst and 310 Seed-GSconst programs, ACT and ECT data were used to update the training population, while only 311 ECT data were used to update the ECT-GS program. After 40 years of future breeding, the training 312 population records grew to 24,600 (Seed-GSunconst), 16,600 (Seed-GSconst), and 6,600 (ECT-GS) 313 records. We estimated genomic predictions using a ridge regression best linear unbiased prediction 314 model (RR-BLUP) (Meuwissen, Hayes, and Goddard 2001). In the model, we fitted year as a fixed 315 effect and allowed for heterogeneous error variance for each breeding stage. The predicted additive SNP effects at each marker locus were used to estimate the average effect of allele substitution for each
SNP. The allele substitution effects were then summed to estimate GEBVs. GEBVs of each genotype
were calculated by summing the predicted additive SNP effects at each marker locus.

319 **2.5** Evaluation and comparison of the tea breeding programs

320 We compared the efficacy of the three GS breeding programs with the PS program by measuring 321 the mean genetic values of the newly developed genotypes at the Seedlings stage over time. All 322 simulations for each strategy were replicated 10 times. We examined the genetic values of seedlings at 323 the Seedlings stage as this is the earliest stage at which all programs evaluate new crosses (F1 324 seedlings). We also evaluated genetic variance and selection accuracy in all breeding programs. We 325 evaluated genetic gain and genetic variance by plotting the mean and variance of seedling population 326 genetic values over time. We calculated the prediction accuracy for the GS breeding programs as the 327 correlation between the true genetic values and their GEBVs at the Seedlings stage. Conversion 328 efficiency was plotted as the change in genetic gain over genic variance.

329 **3 RESULTS**

We showed that all GS programs outperformed the PS breeding program. The Seed-Gsunconst program had the highest overall genetic gain. All GS breeding programs had higher selection accuracies compared to the PS breeding program. Although genetic variance decreased over time for all breeding programs, the GS programs had a large decrease in genetic variance compared to the PS breeding program. Replacement of all parents resulted in a slightly higher genetic gain, with a small decrease in genetic variance compared to replacement of 25% of parents. The PS breeding program had the highest conversion efficiency, but the lowest genetic gain compared to the GS programs.

337 **3.1 Genetic gain**

338 The GS breeding programs (Seed-GSunconst, Seed-GSconst, and ECT-GS) achieved greater 339 genetic gain compared to the PS breeding program, regardless of the number of parents replaced. This 340 is shown in Figure 4, where the population mean genetic value is plotted against the number of years 341 of breeding at the Seedlings stage. The first plot shows the trends for the mean genetic values of all 342 replicates for each of the tea breeding programs evaluated in the future breeding component when 25% 343 of the parents are replaced. The second plot shows the same trend when all parents are replaced after 344 each cycle for all breeding programs. Seed-GSunconst showed the greatest genetic gain compared to 345 all other programs. Both plots show that the overall ranking of the breeding programs in terms of total 346 genetic gain was consistent across the two proportions of parents replaced. The ranking of the breeding 347 programs from highest to lowest mean genetic gain was as follows: Seed-GSunconst, Seed-GSconst, 348 ECT-GS and PS.

349 The breeding programs where all parents were replaced showed slightly higher genetic gain than 350 when 25% of the parents were replaced. When all parents were replaced, the best program, Seed-351 GSunconst, generated 1.6 times the genetic gain of the PS breeding program. Seed-GSconst and ECT-352 GS generated 1.53 and 1.50 times the genetic gain of the PS breeding program, respectively. When 353 25% of the parents were replaced, Seed-GSUnconst generated 1.43 times the genetic gain of the PS 354 breeding program. On the other hand, Seed-GSconst and ECT-GS generated 1.39 and 1.36 times the 355 genetic gain of the PS breeding program, respectively. All GS breeding programs had a generation 356 interval of 2 years compared to 16 years generation interval of the PS breeding program (Figure 4).



358 Figure 4. Genetic gain over time for the four simulated breeding programs. Results are shown for 359 25% (left) and 100% (right) parents replaced separately. The lines within each of the two panels 360 represent the four breeding programs where each line represents population mean of genetic values for 361 the 10 simulated replicates and the shaded area showing 95% confidence intervals of the mean. Genetic 362 gain was measured at the seedlings stage. The black-coloured line is the phenotypic breeding program 363 (PS), the green-coloured line represent the seedlings GS breeding program (Seed-GSconst), the blue-364 coloured line represent the GS breeding program with increased funding (Seed-GSunconst) and the 365 red-coloured line represent the elite clonal GS breeding program (ECT-GS).

366 **3.2 Selection accuracy**

GS increased selection accuracy compared to PS, as shown in Figure 5, which plots the correlations between true and estimated genetic values for seedling entries over time. The first plot shows the selection accuracy for all breeding programs when 25% of all parents were replaced, and the second plot shows the selection accuracy when all parents were replaced. All GS breeding programs had higher selection accuracy compared to the PS breeding program, regardless of the number of parents replaced. Higher selection accuracy was observed when 25% of parents were replaced than when all parents were replaced. Selection accuracy was the same across years and the two methods of parent replacement. In the early years of future breeding, selection accuracy for the ECT-GS program was lower compared to the other GS breeding programs, but then gradually increased until it reached a plateau in year 20. Selection accuracy for the Seed-GSunconst, Seed-GSconst, and PS breeding programs remained constant over time. Both plots show that the ranking from highest to lowest mean selection accuracy is as follows: Seed-GSunconst, Seed-GSconst, ECT-GS, and PS (Figure 5).



379

380 Figure 5. Accuracy of selection over time for the four simulated breeding programs. Results are 381 shown for 25% (left) and 100% (right) parents replaced separately. The lines within each of the two 382 panels represent the four breeding programs where each line represents the mean accuracy of selection 383 for the 10 simulated replicates and the shaded area showing 95% confidence intervals of the mean. 384 Accuracy of selection was measured at the seedlings stage. The black-coloured line is the phenotypic 385 breeding program (PS), the green-coloured line represent the seedlings GS breeding program (Seed-386 GSconst), the blue-coloured line represent the GS breeding program with increased funding (Seed-387 GSunconst) and the red-coloured line represent the elite clonal GS breeding program (ECT-GS).

388 **3.3 Genetic variance**

389 The change in genetic variance over time is shown in Figure 6. The first plot shows the mean 390 genetic variance over time when 25% of the parents are replaced. The second plot shows the same 391 breeding programs when all parents are replaced. All breeding programs showed a decrease in genetic 392 variance over time. However, the rate of loss of genetic variance varied from breeding program to 393 breeding program. All GS breeding programs showed a tremendous decrease in genetic variance, while 394 the phenotypic selection program showed a slow and gradual decrease in genetic variance over time. 395 The difference in genetic variance when 25% and 100% of parents were replaced was small, except 396 during the transition period when GS was introduced and 25% of parents were replaced. Both plots 397 also show that the rank order from highest to lowest genetic variance was as follows: PS, Seed-398 GSconst, ECT-GS and Seed-GSunconst (Figure 6).



399

Figure 6. Genetic variance over time for the four simulated breeding programs. Results are shown for 25% (left) and 100% (right) parents replaced separately. The lines within each of the two panels represent the four breeding programs where each line represents mean genetic variance for the 10 simulated replicates and the shaded area showing 95% confidence intervals of the mean. Genetic variance was measured at the seedlings stage. The black-coloured line is the phenotypic breeding program (PS), the green-coloured line represent the seedlings GS breeding program (Seed-GSconst),

406 the blue-coloured line represent the GS breeding program with increased funding (Seed-GSunconst)

407 and the red-coloured line represent the elite clonal GS breeding program (ECT-GS).

408 **3.4 Conversion efficiency**

409 All breeding strategies had almost similar conversion efficiencies when all parents were replaced. 410 When 25% of the parents were replaced, the PS breeding program had more than twice the conversion 411 efficiency (56) but about 3.5 times less genetic gain than the GS breeding programs. This is illustrated 412 in Figure 7, which shows the long-term genetic gain in standard deviation units when all genic variance 413 is exhausted and is calculated by regressing realized genetic gain on lost genic variance over 40 years 414 of tea breeding. The first plot shows the change in genetic mean over genic standard deviation when 415 25% of the parents are replaced. The second plot shows the change in genetic mean over the genic 416 standard deviation when all parents are replaced. The slope of the change in genetic mean over the 417 change in the genic standard deviation quantifies the efficiency of converting genetic diversity into 418 genetic gain. The ranking of the breeding programs in terms of conversion efficiency when all parents 419 were replaced from highest to lowest was as follows: Seed-GSconst (34), PS (33), ECT-GS (32), and 420 Seed-GSunconst (31). The ranking of conversion efficiency of breeding programs when 25% of parents 421 were replaced from highest to lowest was as follows: PS (56), Seed-GSconst (28), Seed-GSunconst 422 (25), and ECT-GS (25).



423

Figure 7. Efficiency plot showing change of genetic mean and genic standard deviation over time for the four simulated breeding programs. Results are shown for 25% (left) and 100% (right) parents replaced separately. The black-coloured line is the phenotypic breeding program (PS), the greencoloured line represent the seedlings GS breeding program (Seed-GSconst), the blue-coloured line represent the GS breeding program with increased funding (Seed-GSunconst) and the red-coloured line represent the elite clonal GS breeding program (ECT-GS).

430 4 **DISCUSSION**

431 Tea breeding programs require the integration of efficient genomic-assisted breeding 432 approaches to increase the rate of genetic progress. However, currently, it is not clear how these could 433 be integrated into existing programs and whether the additional costs these approaches incur are worth 434 the effort. In our study, we used stochastic simulations for the first time to show that tea breeding 435 programs in LMIC can benefit from genomic selection despite low labour costs and limited research 436 budgets. We developed three different genomic selection breeding programs (Seed-GSconst, Seed-437 GSunconst and ECT-GS) and compared them with a phenotypic selection (PS) program based on the 438 commercial breeding program from Unilever Tea Kenya. To discuss our results, we examine the effects 439 of each breeding strategy on genetic gain, genetic variance, selection accuracy, and conversion efficiency. We also compared the effects of replacing all parents and 25% of parents after each breeding cycle. Our results confirm that the use of GS increases genetic gain compared to traditional PS in tea breeding by shortening the generation interval and increasing the accuracy of selection of superior parents in the Seedlings stage. GS also allowed shortening the entire breeding program by eliminating some of the stages. For example, in Seed-GSconst and Seed-GSunconst, we eliminated the PT stage, hence saving 3 years, while in ECT-GS we removed the PT and ACT stages, resulting in an 8-year reduction in the duration of the breeding program.

447 **4.1 Genetic gain**

Phenotypic selection is a very slow process in tea breeding as it takes 16 years to develop an improved tea variety (Figure 1). This is because tea is a perennial crop with a long generation interval – it takes between 3-6 years for tea bush to grow from seed to flower (Kamunya 2010). The multi-year testing of clones at many locations is also a time-consuming process. The selection accuracy at the Seedlings and PT stage is very low because selection is based on *per se* performance of single bush unreplicated seedlings. This contributes to the slow genetic progress in tea breeding programs and hence, genomic-assisted breeding approaches need to be considered.

455 Our results showed that all breeding programs using GS achieved greater genetic gain 456 compared to the PS breeding program, hence showing the potential of improving the rate of genetic 457 gain in tea breeding programs. This may be attributed to the improved prediction accuracy of selecting 458 superior parents by the GS model and the short generation interval (Cobb et al. 2019). All GS breeding 459 programs had a generation interval of two years, compared to 16 years for the PS breeding program. 460 This demonstrates the importance of reducing generation interval to increase genetic gain in tea 461 breeding programs. Our results showed a perfect rank correlation between generation interval and 462 genetic gain. Similar findings were reported by Bančič et al. (2021), who used stochastic simulations 463 to investigate the potential of using GS to improve yield during intercropping. They found that all

464 programs using GS produced significantly more yield than the PS breeding program mainly because 465 of reduced generation interval and increased prediction accuracy. Gaynor et al. (2017), using stochastic 466 simulations, also reported that both conventional GS breeding methods and the two-part GS strategy 467 significantly increased grain yield of inbred wheat compared to the PS breeding program. Reducing 468 cycle time has the advantage of increasing the frequency with which haplotypes are recombined and 469 exposed for selection, increasing the likelihood that a superior allele combination will emerge and be 470 selected (Atlin, Cairns, and Das 2017). Although Seed-GSunconst produced the most genetic gain, the 471 difference was not large when compared to the Seed-GSconst program, in which only a limited 472 proportion of seedlings were genotyped due to cost constraints. This was less expected and suggests 473 that genotyping more seedlings and testing more genotypes at ACT stage did not improve the genetic 474 gain. Replacing all the parents in the crossing block with new improved parents produced a higher 475 genetic gain, suggesting that substituting the poorest yielding and old parents with new higher yielding 476 parents increases the probability of combining favourable alleles in tea breeding programs.

477 Selection of superior seedlings is the major challenge in tea breeding programs due to high 478 selection intensity and extremely low selection accuracy at the seedlings stage. In the PS breeding 479 program, seedlings are selected based on per se performance, so the selection accuracy is very low. In 480 this study we estimated broad sense heritabilities at the PT, ACT and ECT stages from real data. The 481 heritability at the PT stage was 0.05 compared to the ACT (0.45) and ECT (0.65) stages. The PT stage 482 consists of a single bush (unreplicated) trial, and hence the low heritability. The ACT and ECT stages 483 are clonal trials with larger plots consisting of more clones of each seedling. The ACT stage is also a 484 multilocational trial. In clonal breeding programs, many seedlings are usually selected, thus the 485 increased selection intensity (Werner et al. 2020). For instance, at Unilever Tea Kenya, 500 seedlings 486 are selected each year in the PS breeding program to advance to the ACT stage (Corley and Tuwei 487 2018). The PS breeding program showed a cyclical pattern of genetic gain over time. This is because 488 the parents in the PS breeding program are updated after 16 years, whereas in the GS breeding

489 programs, they were updated every two years, making genetic progress more continuous. In practice, 490 this means that the GS breeding programs can deliver new varieties adapted to new biotic and abiotic 491 stress factors and market requirements more quickly.

492 **4.2** Selection accuracy

493 All GS breeding programs had higher selection accuracy compared to the PS program. The 494 prediction accuracy is high in the GS breeding program because selection at the Seedlings stage is 495 based on the predicted performance of the seedlings as clones, since the GS model is trained using data 496 from the clonal testing stages (Werner et al. 2020). In our simulation, the 500 clones from the ACT 497 stage were used as the initial training population using 6 years of historical phenotypic records (3,000 498 records). This improved the selection accuracy of parents because the phenotypic records at the ACT 499 stage used to train the prediction model had a higher heritability than the individual bushes at the 500 Seedlings stage in the PS breeding program. Additionally, there is increased accuracy of advancing the 501 promising genotypes from the Seedlings to ACT and from the ACT to ECT stages using GEBVs. There 502 is also a strong relationship between the training population and the selection candidates in the case of 503 GS prediction model (Neyhart et al. 2017). We also updated the training population each year with 504 new data from the previous cycle. As breeding cycles progress, the required linkage disequilibrium 505 (LD) between quantitative trait loci and markers is expected to change as a result of recombination, 506 selection and drift, leading to a decay in prediction accuracy (Lorenz et al. 2011). Therefore, the 507 training population should be regularly updated during recurrent selection to maintain the prediction 508 accuracy, as was also the case in our GS breeding programs. Previous research has identified the need 509 to update the training population using new data that may capture new LD generated over breeding 510 cycles. For instance, Neyhart et al. (2017) evaluated several methods for updating the training 511 population in a long-term GS. They reported that using a smaller but more recent training population 512 provided a slight advantage in prediction accuracy and genetic gain.

513 In the early years of the future breeding phase, selection accuracy was lower for the ECT-GS 514 program compared to the other two GS breeding programs. This could be because fewer clones were 515 used to update the training population each year. Each year, only 90 clones from the ECT stage in the 516 ECT-GS breeding program were used to update the training population. In the Seed-GSunconst and 517 Seed-GSconst breeding programs, clones from the ACT and ECT stages were used to update the 518 training population model. The number of clones used to update the Seed-GSunconst and Seed-519 GSconst breeding programs was 540 and 340 clones, respectively. The Seed-GSunconst breeding 520 program had the highest overall prediction accuracy, confirming that training population size is an 521 important factor in the development of GS breeding programs. This is consistent with the results of 522 previous studies that showed that a large training population is required to accurately estimate marker 523 effects (Zhang et al. 2017, Combs and Bernardo 2013). We also observed higher selection accuracy 524 when 25% of parents were replaced than when 100% parents were replaced. This could be because 525 when all the parents are replaced, there is a large shift in LD pattern leading to a decay in prediction 526 accuracy. Selection over time causes a change in LD between the quantitative trait loci (QTL) and the 527 markers. Shifts in the pattern of QTL-marker LD, if not captured, can result in decreased prediction 528 accuracy (Lorenz et al. 2011).

529 **4.3** Genetic variance and conversion efficiency

In our simulation, all breeding programs showed a decrease in genetic variance over time. However, all the GS breeding programs had a huge decrease in genetic variance compared to the PS breeding program, which showed a slow and gradual decrease over time. Similar results were reported by Tessema et al. (2020), who used stochastic simulations based on real data to quantify the increase in genetic gain by implementing GS in a traditional wheat-breeding program. They reported a significant decrease in genetic variance over a 25-year period for the breeding strategies using GS. The loss of genetic variance is due to increased selection accuracy in the early stages of a breeding program as well as shorter generation interval. Increased selection accuracy results in the Bulmer effect, which
decreases genetic variance under directional selection due to the build-up of negative linkage-disequilibrium between causal loci (Bulmer 1971). Selection causes changes in in genetic variances, allele
frequencies and LD relationships between markers and QTL (Muir 2007).

541 Our results showed that the GS breeding programs had lower conversion efficiency compared 542 to the PS breeding program when 25% of the parents were replaced. The conversion efficiency of the 543 PS breeding program decreased significantly when all parents were replaced, while it increased slightly 544 for the GS breeding programs when all parents were replaced. This implies that the PS program 545 converted the genetic gain over loss efficiently compared to the GS programs. A large reduction in 546 genetic variance limits long-term genetic gain in plant breeding because genetic variance is important 547 for continuous and sustained progress (Cobb et al. 2019). Genomic selection strategies that balance the 548 rates of genetic gain and loss of diversity and could be implemented in tea breeding programs include 549 optimal contribution selection (Sonesson, Woolliams, and Meuwissen 2012), optimal cross selection 550 (Gorjanc, Gaynor, and Hickey 2018), optimal contribution selection with branching (Santantonio and 551 Robbins 2020), optimal population value selection (Goiffon et al. 2017) and expected maximum 552 haploid breeding value selection (Müller, Schopp, and Melchinger 2018).

Interestingly, little difference in genetic variance was observed when 25% and 100% of the parents were replaced, except for the transition period when GS was introduced and 25% of parents were replaced. This is mostly because tea is a highly outcrossing and has an extremely highly heterozygous nuclear genome (Xia et al. 2020). Wang et al. (2020) reported that hybridization increased the heterozygosity and wide-ranging gene flow among tea populations with the spread of tea cultivation.

559 4.4 Simulation constrains and practical implementation of GS in tea breeding programs

560	We	e used a real commercial tea breeding program of Unilever Tea Kenya and its parameters as a
561	baselir	he to test the possible integration of GS. Our key constraints were the low operating cost of the
562	breedin	ng program due to cheap labour and limited resources for research. Kenya is one of the LMICs
563	where	labour costs are significantly lower compared to countries with advanced economies. For
564	examp	le, the average daily wage of a field worker at Unilever Tea Kenya is ~\$5 (based on the
565	ccollec	ctive bargaining aagreement (CBA) between the workers and tea companies in Kenya) while in
566	the UF	K the hourly wage is ~\$8-12 (based on salary reports from glassdoor.co.uk). Our simulations
567	showe	d that despite these constraints, it is possible to use GS in tea breeding. Our results provide
568	guidan	ce for several important decisions regarding resource allocation to increase genetic gain in tea
569	breedin	ng programs:
570	1.	Genotyping seedlings in nurseries and selecting the best parents based on GEBVs can increase
571		genetic gain by reducing generation interval and increasing selection accuracy.
572	2.	Elimination of preliminary evaluation (PT) stage reduced the cost of breeding and shortened
573		the breeding cycle by 3 years. The saved costs could be reallocated to genotyping more
574		seedlings in the nursery.
575	3.	Eliminating the PT and ACT stages in the ECT-GS program reduced the duration of the
576		breeding cycle by 8 years. However, genetic gains were lower compared to the other GS
577		strategies.
578	4.	If a breeding program has extra budget (\$30,000), genotyping all the seedlings and increasing
579		the number of genotypes evaluated at the ACT stage can increase genetic gain, however, only
580		slightly.
581	5.	GS is cost-effective for tea breeding programs with limited budgets when genotyping costs are
582		\$15. We expect the benefit of GS to increase in the future when costs of high-throughput
583		genotyping decrease even more - this will increase the selection intensity in a breeding program.

29

6. GS can also be used to advance superior genotypes for variety development, e.g., from the PT
to ACT and from ACT to ECT stages at a higher selection accuracy.

In addition, breeders will also need to consider whether all necessary facilities and equipment are available on site (e.g., freezers, sterile laboratories), train field technicians appropriately, determine whether genotyping can be done on site or else transportation costs should be considered, potentially collaborate with biometricians to optimize field trials, and develop GS pipelines for prediction using an appropriate modelling framework. Standardization and digitization of phenotyping protocols to ensure the best data quality will also be an important challenge, as this can greatly improve the predictive ability of GS.

593 5 Genomic selection for improvement of tea quality and practical implementation

594 Tea quality is an important attribute as it is the main determinant of price at the tea auction. It 595 is measured based on colour, aroma, taste and mouthfeel of tea liquor and the appearance of dry tea 596 (Zheng et al. 2016). These sensory attributes originate from biochemical compounds present in fresh 597 tea shoots such as catechins, alkaloids, amino acids and volatile compounds (Borse and Jagan Mohan 598 Rao 2012). Sensory evaluation using professional tasters is traditionally the main method used to 599 evaluate, grade and determine the price of tea (Liang et al. 2003). Although sensory evaluation is quick 600 and practical to use, it is limited since it requires identification and training to produce skilled and 601 experienced professional tasters (Stone and Sidel 2004), who are not easily found (Corley and Chomboi 602 2005). It is also time-consuming, the tasters sometimes get exhausted and the approach is susceptible 603 to many sources of variation because of individual tasters' preferences and moods (Sinija and Mishra 604 2011). Chemical and physical analytical methods have also been developed for identifying biochemical 605 components associated with tea quality (Liang et al. 2008) and they include; liquid chromatography 606 coupled with mass spectrometry (LC-MS), nuclear magnetic resonance (NMR), near infrared (NIR) 607 spectroscopy and chromatographic methods such as high-performance liquid chromatography (HPLC) and gas chromatography (GC) (Yashin et al. 2015, Zheng et al. 2016). Most of these techniques are
objective, repeatable and reproducible (Chen et al. 2015). However, most of these analytical techniques
are expensive to acquire and maintain and require specialized expertise to operate.

611 GS could be used to improve the selection of superior quality tea varieties (Lubanga, Massawe, 612 and Mayes 2021). The best way to implement GS in tea breeding programs is to integrate it into an 613 existing PS program. GS could be incorporated into an existing PS program by genotyping seedlings 614 at the nursery stage and predicting their genetic values using a GS model. Professional testers could be 615 used to obtain sensory data for training the GS model. Biochemical traits could also be measured at the 616 nursery stage using highly reproducible equipment such as HPLC, NMR or LC-MS. Samples (between 617 100-200) from ACT and ECT stages can be used for training the GS model using both genotypic and 618 phenotypic data. Low-cost genotyping platforms such as genotyping by sequencing (GBS) could be 619 used to obtain SNPs. Seedlings at the nursery stage could then be genotyped and predicted using the 620 trained model using only the genotypic data. Parents with high quality attributes could be selected for 621 crossing at the nursery stage while poor quality seedlings can be discarded at the seedlings stage. GS 622 could be implemented as described in Seed-GSconst, Seed-GSunconst and ECT-GS programs.

The use of GS for breeding of high-quality tea varieties can reduce the disadvantages associated with sensory evaluation methods and analytical techniques. For instance, high quality seedlings can be predicted at the nursery stage without the need to be tested by professional testers or analysed by analytical technical techniques. This increases the accuracy of predicting superior seedlings based on GEBVs, hence reducing the subjectivity associated with the professional testers. Additionally, the expensive costs associated with analytical equipment can be eliminated.

Our study used a practical breeding program at Unilever Tea Kenya and cost estimates to measure the cost-effectiveness of implementing genomic selection in tea breeding programs. All the programs in this study were constrained to equal operating costs, and therefore we can conclude that implementing genomic selection in tea breeding programs can increase the rate of genetic gain despite the challenges experienced in LMIC. Methods such as optimal cross selection (Gorjanc, Gaynor, and Hickey 2018) could be used to ensure that the newly introduced diversity is not quickly eliminated through genomic selection. It optimises the efficiency of converting genetic diversity into genetic gain through reducing the loss of genetic diversity and reducing the drop of genomic prediction accuracy with rapid cycling (Gorjanc, Gaynor, and Hickey 2018).

638 6 Conclusion

639 Our study provides excellent insights into the implementation of genomic selection in tea 640 breeding programs for yield in LMIC. The genomic selection scenarios and results will help tea 641 breeders with knowledge on how to design genomic selection strategies in breeding programs. We 642 show that incorporating GS in tea breeding programs can increase genetic gain up to 1.6 times more 643 than PS program, despite the low labour cost in LMIC. Moreover, the integration of GS does not 644 significantly change the structure of the existing tea breeding program. Rather, it can significantly 645 shorten its' duration. The increase in genetic gain in the GS breeding programs was due to higher 646 prediction accuracy and reduced generation interval. After 40 years of future breeding, the GS breeding 647 programs had lower genetic variance compared to PS, indicating the need to incorporate strategies that 648 balance genetic gain and genetic variance, such as the optimal contribution algorithm. We also 649 observed that replacing all parents resulted in higher genetic gain without significant loss of genetic 650 diversity. Tea quality is a very important attribute, but expensive and difficult trait to phenotype and 651 predict in breeding programs. We recommend further research to determine the most cost-effective 652 pipeline for implementing GS to improve tea quality and yield simultaneously.

653 7 Author contributions

N.L., J.B, F.M. S.M. and G.G designed the research and wrote the manuscript. N.L. and J.B contributed
to the writing of the scripts for the tea breeding simulations.

656 8 Acknowledgments

- 657 This research was supported by University of Edinburgh, University of Nottingham Malaysia, and
- 658 Unilever Tea Kenya.

659

660 **References**

661	Atlin, Gary N., Jill E. Cairns, and Biswanath Das. 2017. "Rapid breeding and varietal replacement
662	are critical to adaptation of cropping systems in the developing world to climate change."
663	Global Food Security 12:31-37. doi: https://doi.org/10.1016/j.gfs.2017.01.008.
664	Bančič, Jon, Christian R. Werner, R. Chris Gaynor, Gregor Gorjanc, Damaris A. Odeny, Henry F.
665	Ojulong, Ian K. Dawson, Stephen P. Hoad, and John M. Hickey. 2021. "Modeling Illustrates
666	That Genomic Selection Provides New Opportunities for Intercrop Breeding." Frontiers in
667	plant science 12:12.
668	Batley, Jacqueline, and David Edwards. 2016. "The application of genomics and bioinformatics to
669	accelerate crop improvement in a changing climate." Current opinion in plant biology 30:78-
670	81.
671	Bernardo, Rex, and Jianming Yu. 2007. "Prospects for genomewide selection for quantitative traits in
672	maize." Crop Science 47 (3):1082-1090.
673	Borse, B. B., and L. Jagan Mohan Rao. 2012. "Novel bio-chemical profiling of Indian black teas with
674	reference to quality parameters." Journal of Bioequivalence and Bioavailability 14:1-16.
675	Bulmer, M. G. 1971. "The effect of selection on genetic variability." The American Naturalist 105
676	(943):201-211.
677	Carr, Mike K. V. 2018. "Advances in tea agronomy."
678	Chen, Gary K., Paul Marjoram, and Jeffrey D. Wall. 2009. "Fast and flexible simulation of DNA
679	sequence data." Genome research 19 (1):136-142.
680	Chen, Quansheng, Dongliang Zhang, Wenxiu Pan, Qin Ouyang, Huanhuan Li, Khulal Urmila, and
681	Jiewen Zhao. 2015. "Recent developments of green analytical techniques in analysis of tea's
682	quality and nutrition." Trends in Food Science & Technology 43 (1):63-82.
	34

683	Cobb, Joshua N., Roselyne U. Juma, Partha S. Biswas, Juan D. Arbelaez, Jessica Rutkoski, Gary
684	Atlin, Tom Hagen, Michael Quinn, and Eng Hwa Ng. 2019. "Enhancing the rate of genetic
685	gain in public-sector plant breeding programs: lessons from the breeder's equation."
686	Theoretical and Applied Genetics 132 (3):627-645. doi: 10.1007/s00122-019-03317-0.
687	Combs, Emily, and Rex Bernardo. 2013. "Accuracy of genomewide selection for different traits with
688	constant population size, heritability, and number of markers." The Plant Genome 6
689	(1):plantgenome2012-11.
690	Corley, R. H. V., and K. C. Chomboi. 2005. "Tea tasting-a statistical evaluation of tasters' skills."
691	<i>Tea</i> 26 (1):10-18.
692	Corley, R. H. V., and G. Tuwei. 2018. "The well-Bred tea Bush." Advances in tea agronomy.
693	Cambridge University Press, Cambridge:106-136.
694	Fang, Kaixing, Zhiqiang Xia, Hongjian Li, Xiaohui Jiang, Dandan Qin, Qiushuang Wang, Qing
695	Wang, Chendong Pan, Bo Li, and Hualing Wu. 2021. "Genome-wide association analysis
696	identified molecular markers associated with important tea flavor-related metabolites."
697	Horticulture Research 8 (1):1-17.
698	Gaynor, R. Chris, Gregor Gorjanc, Alison R. Bentley, Eric S. Ober, Phil Howell, Robert Jackson, Ian
699	J. Mackay, and John M. Hickey. 2017. "A two-part strategy for using genomic selection to
700	develop inbred lines." Crop Science 57 (5):2372-2386.
701	Gaynor, R. Chris, Gregor Gorjanc, and John M. Hickey. 2021. "AlphaSimR: an R package for
702	breeding program simulations." $G3$ 11 (2):jkaa017.
703	Goiffon, Matthew, Aaron Kusmec, Lizhi Wang, Guiping Hu, and Patrick S. Schnable. 2017.
704	"Improving response in genomic selection with a population-based selection strategy: optimal
705	population value selection." Genetics 206 (3):1675-1682.
	35

706	Gorjanc, Gregor, R. Chris Gaynor, and John M. Hickey. 2018. "Optimal cross selection for long-term
707	genetic gain in two-part programs with rapid recurrent genomic selection." Theoretical and
708	applied genetics 131 (9):1953-1966.
709	Gunathilaka, R. P. Dayani, James C. R. Smart, and Christopher M. Fleming. 2017. "The impact of
710	changing climate on perennial crops: the case of tea production in Sri Lanka." Climatic
711	<i>change</i> 140 (3-4):577-592.
712	Han, Wen-Yan, Ji-Gang Huang, Xin Li, Zhi-Xin Li, Golam Jalal Ahammed, Peng Yan, and John
713	Richard Stepp. 2017. "Altitudinal effects on the quality of green tea in east China: a climate
714	change perspective." European Food Research and Technology 243 (2):323-330.
715	Han, Wen-Yan, Xin Li, and Golam Jalal Ahammed. 2018. Stress physiology of tea in the face of
716	climate change: Springer.
717	Heffner, Elliot L., Mark E. Sorrells, and Jean-Luc Jannink. 2009. "Genomic selection for crop
718	improvement." Crop Science 49 (1):1-12.
719	Iwata, Hiroyoshi, Takeshi Hayashi, and Yoshihiko Tsumura. 2011. "Prospects for genomic selection
720	in conifer breeding: a simulation study of Cryptomeria japonica." Tree genetics & genomes 7
721	(4):747-758.
722	Jin, Ji-Qiang, Ming-Zhe Yao, Chun-Lei Ma, Jian-Qiang Ma, and Liang Chen. 2016. "Association
723	mapping of caffeine content with TCS1 in tea plant and its related species." Plant physiology
724	and biochemistry 105:251-259.
725	Kamunya, S. M., F. N. Wachira, R. S. Pathak, R. Korir, V. Sharma, R. Kumar, P. Bhardwaj, R.
726	Chalo, P. S. Ahuja, and R. K. Sharma. 2010. "Genomic mapping and testing for quantitative
727	trait loci in tea (Camellia sinensis (L.) O. Kuntze)." Tree Genetics & Genomes 6 (6):915-929.
728	doi: 10.1007/s11295-010-0301-2.

729	Kamunya, Samson M., Francis N. Wachira, Ram S. Pathak, Richard C. Muoki, and Ram K. Sharma.
730	2012. "Tea improvement in Kenya." In Global Tea Breeding, 177-226. Springer.
731	Kamunya, Samson Machohi. 2010. "Genetic parameters and quantitative trait loci mapping in tea,
732	Camellia sinensis (L.) O. Kuntze."
733	Leng, Peng-fei, Thomas Lübberstedt, and Ming-liang Xu. 2017. "Genomics-assisted breeding – A
734	revolutionary strategy for crop improvement." Journal of Integrative Agriculture 16
735	(12):2674-2685. doi: 10.1016/s2095-3119(17)61813-6.
736	Liang, Y. R., Q. Ye, J. Jin, H. Liang, J. L. Lu, Y. Y. Du, and J. J. Dong. 2008. "Chemical and
737	instrumental assessment of green tea sensory preference." International Journal of Food
738	Properties 11 (2):258-272.
739	Liang, Yuerong, Jianliang Lu, Lingyun Zhang, Shan Wu, and Ying Wu. 2003. "Estimation of black
740	tea quality by analysis of chemical composition and colour difference of tea infusions." Food
741	chemistry 80 (2):283-290.
742	Lorenz, Aaron J., Shiaoman Chao, Franco G. Asoro, Elliot L. Heffner, Takeshi Hayashi, Hiroyoshi
743	Iwata, Kevin P. Smith, Mark E. Sorrells, and Jean-Luc Jannink. 2011. "Genomic selection in
744	plant breeding: knowledge and prospects." Advances in agronomy 110:77-123.
745	Lubanga, Nelson, Festo Massawe, and Sean Mayes. 2021. "Genomic and pedigree-based predictive
746	ability for quality traits in tea (Camellia sinensis (L.) O. Kuntze)." Euphytica 217 (3):32. doi:
747	10.1007/s10681-021-02774-3.
748	Malebe, M. P., R. K. Koech, E. G. N. Mbanjo, S. M. Kamunya, A. A. Myburg, and Z. Apostolides.
749	2021. "Construction of a DArT-seq marker-based genetic linkage map and identification
750	of QTLs for yield in tea (Camellia sinensis (L.) O. Kuntze)." Tree Genetics & Genomes 17
751	(1):9. doi: 10.1007/s11295-021-01491-1.

752	Marx, Werner, Robin Haunschild, and Lutz Bornmann. 2017. "Global warming and tea production-
753	The bibliometric view on a newly emerging research topic." <i>Climate</i> 5 (3):46.
754	Meegahakumbura, M. K., Moses C. Wambulwa, K. K. Thapa, M. M. Li, M. Möller, J. C. Xu, J. B.
755	Yang, B. Y. Liu, S. Ranjitkar, and J. Liu. 2016. "Indications for three independent
756	domestication events for the tea plant (Camellia sinensis (L.) O. Kuntze) and new insights
757	into the origin of tea germplasm in China and India revealed by nuclear microsatellites." PloS
758	one 11 (5):e0155369.
759	Meuwissen, Theo H. E., Ben J. Hayes, and Michael E. Goddard. 2001. "Prediction of total genetic
760	value using genome-wide dense marker maps." Genetics 157 (4):1819-1829.
761	Mondal, Tapan Kumar. 2011. "Camellia." In Wild crop relatives: genomic and breeding resources,
762	15-39. Springer.
763	Mondal, Tapan Kumar. 2014. Breeding and biotechnology of tea and its wild species: Springer
764	Science & Business Media.
765	Muir, W. M. 2007. "Comparison of genomic and traditional BLUP-estimated breeding value
766	accuracy and selection response under alternative trait and genomic parameters." Journal of
767	Animal Breeding and Genetics 124 (6):342-355.
768	Mukhtar, Hasan, and Nihal Ahmad. 2000. "Tea polyphenols: prevention of cancer and optimizing
769	health" The American journal of clinical nutrition 71 (6):1698S-1702S.
770	Muleta, Kebede T., Gael Pressoir, and Geoffrey P. Morris. 2019. "Optimizing Genomic Selection for
771	a Sorghum Breeding Program in Haiti: A Simulation Study." G3: Genes/Genomes/Genetics 9
772	(2):391. doi: 10.1534/g3.118.200932.

773	Müller, Dominik, Pa	scal Schopp, and Al	brecht E. Melchinger.	. 2018. "Selectio	on on expected

- 774 maximum haploid breeding values can increase genetic gain in recurrent genomic selection."
- 775 *G3: Genes, Genomes, Genetics* 8 (4):1173-1181.
- 776 Muoki, Chalo Richard, Tony Kipkoech Maritim, Wyclife Agumba Oluoch, Samson Machohi
- 777 Kamunya, and John Kipkoech Bore. 2020. Combating Climate Change in the Kenyan Tea
- Industry. *Frontiers in plant science* 11: 339. Accessed 2020. doi:10.3389/fpls.2020.00339.
- 779 Neyhart, Jeffrey L., Tyler Tiede, Aaron J. Lorenz, and Kevin P. Smith. 2017. "Evaluating Methods of
- 780 Updating Training Data in Long-Term Genomewide Selection." *G3:*

781 *Genes/Genetics* 7 (5):1499. doi: 10.1534/g3.117.040550.

- 782 Perroy, R. 2015. World Population Prospects. United Nations, 1 (6042), 587–92.
- Powell, Owen M., R. Chris Gaynor, Gregor M. Gorjanc, Christian R. Werner, and John M. Hickey.
 2020. "A Two-Part Strategy using Genomic Selection in Hybrid Crop Breeding Programs."
 bioRxiv.
- Rutkoski, J., R. P. Singh, J. Huerta-Espino, S. Bhavani, J. Poland, J. L. Jannink, and M. E. Sorrells.
 2015. "Genetic gain from phenotypic and genomic selection for quantitative resistance to
 stem rust of wheat."
- Santantonio, Nicholas, and Kelly Robbins. 2020. "A hybrid optimal contribution approach to drive
 short-term gains while maintaining long-term sustainability in a modern plant breeding
 program." *bioRxiv*.
- Sinija, V. R., and Hari Niwas Mishra. 2011. "Fuzzy analysis of sensory data for quality evaluation
 and ranking of instant green tea powder and granules." *Food and bioprocess technology* 4
 (3):408-416.

- 795 Sitienei, Betty J., Shem G. Juma, and Everline Opere. 2017. "On the Use of Regression Models to
- Predict Tea Crop Yield Responses to Climate Change: A Case of Nandi East, Sub-County of
 Nandi County, Kenya." *Climate* 5 (3):54.
- 798 Sonesson, Anna K., John A. Woolliams, and Theo H. E. Meuwissen. 2012. "Genomic selection
- requires genomic control of inbreeding." *Genetics Selection Evolution* 44 (1):1-10.
- Stone, H., and J. L. Sidel. 2004. "Sensory Evaluation Practices, Elsevier Academic Press." *California, USA*.
- 802 Tessema, Biructawit Bekele, Huiming Liu, Anders Christian Sørensen, Jeppe Reitan Andersen, and
- Just Jensen. 2020. "Strategies Using Genomic Selection to Increase Genetic Gain in Breeding
 Programs for Wheat." *Frontiers in Genetics* 11 (1538). doi: 10.3389/fgene.2020.578123.
- 805 Valin, Hugo, Ronald D. Sands, Dominique Van der Mensbrugghe, Gerald C. Nelson, Helal
- 806 Ahammad, Elodie Blanc, Benjamin Bodirsky, Shinichiro Fujimori, Tomoko Hasegawa, and
- 807 Petr Havlik. 2014. "The future of food demand: understanding differences in global economic
- 808 models." *Agricultural Economics* 45 (1):51-67.
- 809 Wang, Jian-kang, and Pfeiffer Wolfgang H. 2007. "Simulation Modeling in Plant Breeding:
- 810 Principles and Applications." *Agricultural Sciences in China* 6 (8):908-921. doi:
- 811 https://doi.org/10.1016/S1671-2927(07)60129-1.

812 Wang, Xinchao, Hu Feng, Yuxiao Chang, Chunlei Ma, Liyuan Wang, Xinyuan Hao, A'lun Li, Hao

- 813 Cheng, Lu Wang, Peng Cui, Jiqiang Jin, Xiaobo Wang, Kang Wei, Cheng Ai, Sheng Zhao,
- 814 Zhichao Wu, Youyong Li, Benying Liu, Guo-Dong Wang, Liang Chen, Jue Ruan, and Yajun
- 815 Yang. 2020. "Population sequencing enhances understanding of tea plant evolution." *Nature*
- 816 *communications* 11 (1):4447-4447. doi: 10.1038/s41467-020-18228-8.

817	Werner, Christian R., R. Chris Gaynor, Daniel J. Sargent, Alessandra Lillo, Gregor Gorjanc, and
818	John M. Hickey. 2020. "Genomic selection strategies for clonally propagated crops."
819	bioRxiv.
820	Xia, En-Hua, Wei Tong, Qiong Wu, Shu Wei, Jian Zhao, Zheng-Zhu Zhang, Chao-Ling Wei, and
821	Xiao-Chun Wan. 2020. "Tea plant genomics: achievements, challenges and perspectives."
822	Horticulture Research 7 (1):7. doi: 10.1038/s41438-019-0225-4.
823	Yamashita, Hiroto, Tomoki Uchida, Yasuno Tanaka, Hideyuki Katai, Atsushi J. Nagano, Akio
824	Morita, and Takashi Ikka. 2020. "Genomic predictions and genome-wide association studies
825	based on RAD-seq of quality-related metabolites for the genomics-assisted breeding of tea
826	plants." Scientific Reports 10 (1):17480. doi: 10.1038/s41598-020-74623-7.
827	Yashin, Alexandr Ya, Boris V. Nemzer, Emilie Combet, and Yakov I. Yashin. 2015. "Determination
828	of the chemical composition of tea by chromatographic methods: a review." Journal of Food
829	<i>Research</i> 4 (3):56-87.
830	Zhang, Ao, Hongwu Wang, Yoseph Beyene, Kassa Semagn, Yubo Liu, Shiliang Cao, Zhenhai Cui,
831	Yanye Ruan, Juan Burgueño, and Felix San Vicente. 2017. "Effect of trait heritability,
832	training population size and marker density on genomic prediction accuracy estimation in 22
833	bi-parental tropical maize populations." Frontiers in Plant Science 8:1916.
834	Zheng, Xin-Qiang, Qing-Sheng Li, Li-Ping Xiang, and Yue-Rong Liang. 2016. "Recent advances in
835	volatiles of teas." Molecules 21 (3):338.

836