

Theor Appl Genet (2002) 104:1325–1334
DOI 10.1007/s00122-001-0854-4

S. Chandra · Z. Huaman · S. Hari Krishna · R. Ortiz

Optimal sampling strategy and core collection size of Andean tetraploid potato based on isozyme data – a simulation study

Received: 5 June 2001 / Accepted: 8 November 2001 / Published online: 26 April 2002
© Springer-Verlag 2002

Abstract Selection of an appropriate sampling strategy is an important prerequisite to establish core collections of appropriate size in order to adequately represent the genetic spectrum and maximally capture the genetic diversity in available crop collections. We developed a simulation approach to identify an optimal sampling strategy and core-collection size, using isozyme data from a CIP germplasm collection on an Andean tetraploid potato. Five sampling strategies, constant (C), proportional (P), logarithmic (L), square-root (S) and random (R), were tested on isozyme data from 9,396 Andean tetraploid potato accessions characterized for nine isozyme loci having a total of 38 alleles. The 9,396 accessions, though comprising 2,379 morphologically distinct accessions, were found to represent 1,910 genetically distinct groups of accessions for the nine isozyme loci using a sort-and-duplicate-search algorithm. From each group, one accession was randomly selected to form a genetically refined entire collection (GREC) of size 1,910. The GREC was used to test the five sampling strategies. To assess the behavior of the results in repeated sampling, $k = 1,500$ and $5,000$ independent random samples (without replacement) of admissible sizes $n = 50(50)1,000$ for each strategy were drawn from GREC. Allele frequencies (AF) for the 38 alleles and lo-

cus heterozygosity (LH) for the nine loci were estimated for each sample. The goodness of fit of samples AF and LH with those from GREC was tested using the χ^2 test. A core collection of size $n = 600$, selected using either the P or the R sampling strategy, was found adequately to represent the GREC for both AF and LH. As similar results were obtained at $k = 1,500$ and $5,000$, it seems adequate to draw 1,500 independent random samples of different sizes to test the behavior of different sampling strategies in order to identify an appropriate sampling approach, as well as to determine an optimal core collection size.

Keywords Andean tetraploid potato · Core collection · Sampling strategy · Simulation

Introduction

Frankel and Brown (1984) proposed the concept of a core collection to enable efficient and cost-effective management and utilization of crop genetic resources. They defined a core collection as a limited subset of accessions from an existing germplasm collection that adequately represents the genetic spectrum of, and captures maximal genetic diversity in, a collection held in a genebank. An ideal core collection should include entries that are also ecologically and genetically distinct from one another (Brown 1989). A core collection that meets these requirements acts as a representative entry point to the whole collection in order to facilitate the processes of crop genetic improvement and research.

The gene bank at Centro Internacional de la Papa (CIP, Lima, Perú) maintains one of the largest collections of tetrasomic Andean potatoes (*Solanum tuberosum* subsp. *andigena*) (Huaman 1998). These accessions have been characterized both at morphological and genetic levels, the latter using isozyme markers. Isozymes markers in Andean potatoes have been employed for assessing genetic variation (Zimmerer and Douches 1991; Quiros et al. 1992), determining rates of out-crossing be-

Communicated by H.C. Becker

S. Chandra · S. Hari Krishna
International Crops Research Institute
for the Semi-Arid Tropics (ICRISAT),
Patancheru, 502 324 AP, India

Z. Huaman
Centro Internacional de la Papa, Apartado 1558, Lima 100, Perú

R. Ortiz (✉)
International Institute of Tropical Agriculture (IITA),
c/o Lambourn & Co., Carolyn House, 26 Dingwall Road,
Croydon, CR9 3EE, UK
e-mail: r.ortiz@cgiar.org

Present address:

Z. Huaman, Pro Biodiversity of the Andes,
Av. Raul Ferrero # 1354, Lima 12, Perú

tween primitive cultivated potatoes (Rabinowitz et al. 1990), characterizing North American tetraploid potato cultivars (Douches et al. 1991; Douches and Ludlam 1991) and for determining how human selection affects genetic diversity in tetraploid potatoes (Ortiz and Huaman 2001). Some of these isozymes are associated with important agronomic characters in potato-segregating populations (Ortiz et al. 1993; Freyre and Douches 1994; Freyre et al. 1994).

Recently, Huaman et al. (2000a) developed a core collection of 306 Andean tetraploid potatoes from a subset of morphologically distinct 2,379 accessions. The latter were selected from an existing whole collection of 10,722 accessions held in the CIP gene-bank after removing from it 8,343 duplicate accessions based on several morphological traits. A square-root sampling approach was used to select the core of 306 entries from each geographical division of Latin American countries, from which the 2,379 accessions were collected. Data on nine isozyme markers were subsequently used to investigate the genetic structure of the 2,379 accessions and to assess the genetic representativeness of the core of 306 entries in terms of allele frequencies and locus heterozygosity (Huaman et al. 2000b).

The objective of this research was to develop and apply a simulation approach to determine an optimal sampling strategy and core collection size for Andean tetraploid potato accessions using only the isozyme data.

Materials and methods

Genetic materials

The original Andean tetraploid collection at CIP consisted of 10,722 accessions from eight Latin American countries. Of these, only 9,396 accessions, characterized for morphology and nine isozymes, were included in this study. These 9,396 accessions represented 2,379 morphologically distinct genotypes (Huaman et al. 2000b).

Genetic markers

Allozyme diversity was determined using horizontal gel-electrophoresis and two buffer systems. The procedures for tissue processing, electrophoresis, gel staining and allozyme scoring were those of Douches and Quiros (1988) and Huaman et al. (2000b). These nine isozyme loci covering a total of 38 alleles were: isocitric acid dehydrogenase 1 (*Idh-1* in chromosome I), malate dehydrogenase 1 (*Mdh-1* in chromosome VII), malate dehydrogenase 2 (*Mdh-2*), and phosphoglucose isomerase 1 (*Pgi-1* in chromosome XII) for histidine-citrate at pH 5.7, and Diaphorase 1 (*Dia-1* in chromosome V), glutamate oxaloacetate transaminase 1 (*Got-1* in chromosome VIII), glutamate oxaloacetate transaminase 2 (*Got-2* in chromosome VII), phosphoglucose mutase 1 (*Pgm-1* in chromosome III), and phosphoglucose mutase 2 (*Pgm-2* in chromosome IV) for tris-borate at pH 8.3.

Creation of a genetically refined entire collection

The data consist of counts Y_{ijk} of allele $k = 1, \dots, a_j \in [3, 6]$ at locus $j = 1, \dots, n_l (= 9)$ for accession $i = 1, \dots, N (= 9,396)$ with the property $\sum_k Y_{ijk} = 4 \forall (i, j)$. The allele counts Y_{ijk} were transformed to allele frequencies $P_{ijk} = (1/4)Y_{ijk}$ with $\sum_k P_{ijk} = 1 \forall (i, j)$.

A sort-and-duplicate-search algorithm found the N accessions to fall into $K = 1,910$ distinct allelic-configuration/genotype classes, with $N_i \in [1, 198]$ duplicate genotypes present in class $t = 1, \dots, K, \sum_t N_t = N$. The original entire collection (OEC) of N accessions was therefore first reduced to a genetically refined entire collection (GREC) of $K = 1,910$ distinct tetraploid genotypes by randomly selecting one accession from each of the K genotype classes. The GREC, rather than the OEC, was used to investigate the suitability of different sampling strategies and to determine the optimal core collection size. Use of GREC ensures that the core contains genetically distinct entries.

Sampling strategies

Five sampling strategies were investigated, random (R), constant (C), proportional (P), logarithmic (L) and square root (S). For the R strategy, accessions were randomly selected from the GREC using simple random sampling without replacement (SRSWOR), in keeping with the fact that a core should include distinct entries. For C, P, L and S strategies, the 1,910 accessions in the GREC were first grouped into eight clusters according to the country of their collection as follows: Argentina 73, Bolivia 258, Colombia 105, Ecuador 131, Guatemala 24, Mexico 16, Peru 1,276 and Venezuela 27. From each of these eight clusters, the number of accessions n_u ($u = 1, \dots, 8$), to obtain a specified core sample of $n = \sum_u n_u$ accessions, was selected using intra-cluster SRSWOR as follows:

Strategy	Intra-cluster sample-size n_u	Admissible/tested core sample size n
C	$n_u = n/8$	$n = 50(50)150^*$
L	$n_u = n[\log(K_u)/\sum_u \log(K_u)]$	$n = 50(50)250$
P	$n_u = K_u (n/K)$	$n = 50(50)1,000$
S	$n_u = n[\sqrt{K_u}/\sum_u \sqrt{K_u}]$	$n = 50(50)400$
R	–	$n = 50(50)1,000$

K_u = size of cluster u ; *sample sizes varied between 50 and 150 with an increment of 50

Estimation of allele frequencies and locus heterozygosity

The allele frequencies (AF) P_{jk} for allele k at locus j , and locus heterozygosity (LH) H_j for locus j in the GREC were computed as follows:

$$P_{jk} = (a_{jk}/a_j) = (\sum_t a_{tjk} / \sum_t \sum_k a_{tjk}) = [\sum_t a_{tjk} / (4 \times K)] \forall j \quad (1)$$

$$H_j = \{1 - [\sum_k \#(P_{jk} = 1)/K]\}, \quad (2)$$

where a_{jk} : total count of the k -th-type allele at locus j across K genotypes, a_j : total count of all allele types at locus j across K genotypes, a_{tjk} : count of the k -th-type allele at locus j for genotype $t = 1, \dots, K, \#(P_{jk} = 1)$: number of genotypes homozygous for allele k at locus j .

Sample estimates p_{jk} and h_j of P_{jk} and H_j respectively for a sample of size n drawn from using any sampling strategy were obtained from equations (1) and (2) respectively, with K replaced by the sample size n .

Chi-square tests of goodness-of-fit

Goodness-of-fit of the sample estimates to the population values was tested using χ^2 tests as follows:

$$\chi^2(\text{AF}) = \sum_j \chi_j^2(\text{AF})df = \sum_j (a_j - 1) = 28$$

at a level of significance (LOS) α for across-the-loci (genome-wide) fit of AF ($H_0: p_{jk} = P_{jk}, j=1, \dots, n_1, k=1, \dots, a_j$),

$$\chi_j^2(\text{AF}) = 4n \sum_k [(p_{jk} - P_{jk})^2 P_{jk}] k = 1, \dots, a_j \quad df = a_j - 1$$

at LOS $\alpha_j = \alpha/8$ for an individual locus-wise fit of AF ($H_0: p_{(j)k} = P_{(j)k}, k = 1, \dots, a_j$), $\chi^2(\text{LH}) = \sum_j \chi_j^2(\text{LH}) df = n_1 = 9$ at LOS α for across-the-loci (genome-wide) fit of LH ($H_0: h_j = H_j, j = 1, \dots, 9$), and $\chi_j^2(\text{LH}) = (h_j - H_j)^2 [H_j(1 - H_j)/n] df = 1$ at LOS $\alpha_j = \alpha/8$ for an individual locus-wise fit of LH ($H_0: h_j = H_j$).

Simulations

Inferences based on just one sample of a particular size n could be misleading as this does not give an idea of the likely variation in the results had we drawn more samples of that size. Repeated samples provide an objective assessment of the degree of consistency, stability and reproducibility of results. Therefore, $k = 1,500$ and 5,000 independent random samples of a particular size $n = 50(50)1,000$, as admissible for a given sampling strategy, were drawn according to the afore-stated five sampling strategies. Two values of k were chosen to determine the adequate number of random samples to be simulated.

A sample size and a strategy that consistently do not reject H_0 at a chosen level of significance α across all k repeated samples are the safest sample size and strategy to use. This is practically unlikely to happen as long as $n < N$. However, for a given sam-

pling strategy, a sample size n for which, under H_0 , the k -observed χ^2 -values follow the corresponding theoretical χ^2 distribution, provides a lower bound, if that exists, on optimal sample size. We used the Kolmogorov-Smirnov (K-S) test (Sokal and Rohlf 1981) to identify this lower bound on the optimal sample size for each sampling strategy. Having identified the lower bound on an optimal n , the optimal n for a given sampling strategy can be determined from a suitably chosen characteristic of the frequency distribution of the k -observed χ^2 -values. Some possible candidate-characteristics are the maximum, upper-0.05-quantile, and a median of the observed distribution of the k values of χ^2 . The maximum is obviously the safest to use as it covers the maximum possible risk in terms of the largest possible discrepancy between GREC and sample values. However, since χ^2 can theoretically assume a maximum value of infinity, it is likely that, with increasing n , the observed maximum χ^2 values may show an erratic pattern, which they did, (see Tables 1 and 2). that This situation will make it difficult to clearly identify an optimal sample size and strategy. Use of the median, compared to using the observed upper-0.05-quantile χ^2 , on the other hand, covers much-less risk. We therefore chose to use the upper-0.05-quantile of the observed distribution of k χ^2 -values to judge the suitability of a sample size and strategy. Any upper-0.05-quantile χ^2 -value that is non-significant at a chosen level of significance α implies that, for the corresponding sample size and strategy, all samples of that size will consistently deliver non-significant χ^2 values 95% of the time, and hence provide a good fit to the GREC. Also, the more the P -value of the observed upper-0.05-quantile χ^2 -value exceeds the specified α , less is the discrepancy between GREC and sample values. From this perspective, one could choose an α -value more than the conven-

Table 1 Quantiles of 1,500 observed χ^2 values for allele frequencies for different sample sizes (n) under proportional strategy

n	Min	0.95-uq ^a	0.75-uq	0.50-uq	0.25-uq	0.05-uq	Max	D ^b
50	12.70 <i>0.9941^c</i>	21.24 <i>0.8152</i>	28.51 <i>0.4377</i>	35.32 <i>0.1607</i>	44.45 <i>0.0251</i>	78.58 <i>0.0000</i>	132.98 <i>0.0000</i>	0.3603
100	12.24 <i>0.9957</i>	21.94 <i>0.7841</i>	28.90 <i>0.4176</i>	35.40 <i>0.1585</i>	43.60 <i>0.0304</i>	62.66 <i>0.0002</i>	129.76 <i>0.0000</i>	0.3583
150	12.30 <i>0.9955</i>	20.94 <i>0.8278</i>	28.20 <i>0.4539</i>	34.86 <i>0.1740</i>	42.48 <i>0.0390</i>	56.49 <i>0.0011</i>	108.00 <i>0.0000</i>	0.3421
200	10.16 <i>0.9992</i>	20.40 <i>0.8494</i>	28.00 <i>0.4644</i>	34.48 <i>0.1855</i>	42.68 <i>0.0374</i>	55.60 <i>0.0014</i>	92.96 <i>0.0000</i>	0.3296
250	11.60 <i>0.9973</i>	20.20 <i>0.8571</i>	27.60 <i>0.4858</i>	34.00 <i>0.2009</i>	40.90 <i>0.0548</i>	52.45 <i>0.0034</i>	75.50 <i>0.0000</i>	0.3085
300	9.60 <i>0.9995</i>	20.04 <i>0.8630</i>	27.00 <i>0.5182</i>	32.52 <i>0.2539</i>	39.06 <i>0.0800</i>	51.12 <i>0.0048</i>	80.52 <i>0.0000</i>	0.2598
350	10.22 <i>0.9991</i>	19.88 <i>0.8688</i>	26.32 <i>0.5555</i>	32.06 <i>0.2721</i>	38.85 <i>0.0834</i>	49.84 <i>0.0067</i>	76.02 <i>0.0000</i>	0.2355
400	8.64 <i>0.9998</i>	19.04 <i>0.8969</i>	25.60 <i>0.5950</i>	30.88 <i>0.3224</i>	36.96 <i>0.1197</i>	47.52 <i>0.0121</i>	74.08 <i>0.0000</i>	0.1837
450	10.44 <i>0.9990</i>	18.36 <i>0.9167</i>	24.48 <i>0.6560</i>	29.88 <i>0.3690</i>	35.82 <i>0.1472</i>	46.26 <i>0.0164</i>	75.60 <i>0.0000</i>	0.1399
500	9.60 <i>0.9995</i>	17.80 <i>0.9311</i>	23.80 <i>0.6920</i>	29.00 <i>0.4125</i>	35.20 <i>0.1641</i>	44.00 <i>0.0278</i>	69.40 <i>0.0000</i>	0.0983
550	10.34 <i>0.9990</i>	17.16 <i>0.9454</i>	22.88 <i>0.7390</i>	27.72 <i>0.4794</i>	33.66 <i>0.2123</i>	42.68 <i>0.0374</i>	81.40 <i>0.0000</i>	0.0420
600	8.88 <i>0.9998</i>	16.56 <i>0.9568</i>	22.08 <i>0.7776</i>	26.88 <i>0.5248</i>	31.92 <i>0.2778</i>	41.04 <i>0.0533</i>	60.24 <i>0.0004</i>	0.0403
650	9.36 <i>0.9996</i>	16.12 <i>0.9640</i>	21.58 <i>0.8004</i>	26.00 <i>0.5730</i>	30.68 <i>0.3315</i>	39.26 <i>0.0768</i>	62.66 <i>0.0002</i>	0.0952
700	8.68 <i>0.9998</i>	15.68 <i>0.9703</i>	20.72 <i>0.8368</i>	24.78 <i>0.6398</i>	29.96 <i>0.3651</i>	38.08 <i>0.0969</i>	61.88 <i>0.0002</i>	0.1555
750	7.50 <i>1.0000</i>	15.30 <i>0.9751</i>	20.10 <i>0.8608</i>	24.00 <i>0.6815</i>	28.80 <i>0.4227</i>	36.90 <i>0.1211</i>	48.90 <i>0.0086</i>	0.2020
800	9.60 <i>0.9995</i>	14.24 <i>0.9854</i>	18.88 <i>0.9018</i>	23.04 <i>0.7310</i>	27.52 <i>0.4901</i>	35.04 <i>0.1687</i>	52.16 <i>0.0037</i>	0.2667

^a Upper quantile (uq)

^b Kolmogorov-Smirnov test-statistic value ($D_{0.05} = 0.035, D_{0.01} = 0.042$ based on $k = 1,500$)

^c P -value (*italics*) of the above observed χ^2 values

tional values of 0.05 and 0.01 to further minimize the risk of picking up an inappropriate sample size and strategy. The P -values corresponding to the observed upper-0.05-quantile χ^2 -values, summarized in a tabular or graphical form, provide an objective probabilistic basis to compare the suitability of different sampling strategies to help determine the optimal sample size and strategy, with α chosen according to the risk one wants to cover.

Our strategy in determining an optimal sample size for a given sampling strategy was to adopt the approach of the preceding paragraph to first check the overall genome-wide fit. Having identified the genome-wide optimal sample size for a chosen α , the suitability of that sample size at individual loci was determined using a locus-wise level of significance $\alpha_j = \alpha/n_L$ based on the Bonferroni correction, where the denominator $n_L \leq n_i$ represents the number of independent linkage groups on which the n_i isozyme loci are located.

An optimal sampling strategy is defined as one that, for the observed upper-0.05-quantile χ^2 , provides a smaller genome-wide optimal sample size with a P -value \geq to the chosen level of significance, α . It is anticipated that, due to the difference in the way AF and LH are estimated, different sample sizes may turn out to be optimal for AF and LH. We took the larger of the two optimal sample sizes as the optimal sample size for both AF and LH.

Results

The results were similar for $k = 1,500$ and $5,000$. Accordingly, we will subsequently report $k = 1,500$ in pre-

senting and discussing the results. The K-S test showed that, for the admissible values of n , a lower bound on optimal n did not exist for the C, L and S strategies. At the same time, for all three strategies, the P -value corresponding to the upper-0.05-quantile χ^2 (AF) and the upper-0.05-quantile χ^2 (LH) never exceeded $\alpha = 0.05$ for any of the sample sizes. These three strategies, regarded as non-optimal because of the above reasons, are therefore not discussed further.

Allele frequencies

The genome-wide frequency distribution of the 1,500 observed χ^2 (AF)-values for the P and R strategies for different sample sizes n are summarized in Tables 1 and 2 respectively. As expected from the law of large numbers, the χ^2 values show a generally decreasing trend as the sample size n increases. Figure 1 depicts the observed upper-0.05-quantile χ^2 (AF)-values and their corresponding P -values for the P and R strategies. Figure 2 provides the values of the upper-0.05-quantile χ^2_j (AF)-values and their corresponding P -values for individual loci. The K-S test-statistic for the P-strategy (Table 1) is non-significant (at $\alpha = 0.01$) at $n = 550$,

Table 2 Quantiles of 1,500 observed χ^2 values for allele frequencies for different sample sizes (n) under random strategy

n	Min	0.95- uq^a	0.75- uq	0.50- uq	0.25- uq	0.05- uq	Max	D^b
50	12.96 <i>0.9931^c</i>	21.81 <i>0.7900</i>	29.39 <i>0.3930</i>	36.52 <i>0.1298</i>	45.61 <i>0.0191</i>	66.23 <i>0.0001</i>	169.92 <i>0.0000</i>	0.3951
100	12.16 <i>0.9959</i>	21.48 <i>0.8048</i>	29.32 <i>0.3964</i>	36.10 <i>0.1401</i>	44.96 <i>0.0223</i>	64.82 <i>0.0001</i>	108.92 <i>0.0000</i>	0.3842
150	10.32 <i>0.9991</i>	20.91 <i>0.8291</i>	29.25 <i>0.3999</i>	36.30 <i>0.1351</i>	44.58 <i>0.0243</i>	59.79 <i>0.0004</i>	100.32 <i>0.0000</i>	0.3799
200	12.88 <i>0.9934</i>	21.00 <i>0.8253</i>	29.16 <i>0.4044</i>	35.76 <i>0.1488</i>	44.00 <i>0.0278</i>	57.84 <i>0.0008</i>	103.92 <i>0.0000</i>	0.3811
250	13.30 <i>0.9915</i>	20.30 <i>0.8533</i>	28.65 <i>0.4304</i>	35.20 <i>0.1641</i>	42.05 <i>0.0429</i>	54.95 <i>0.0017</i>	83.90 <i>0.0000</i>	0.3535
300	12.24 <i>0.9957</i>	20.28 <i>0.8540</i>	27.36 <i>0.4987</i>	33.72 <i>0.2102</i>	41.16 <i>0.0519</i>	51.96 <i>0.0039</i>	88.44 <i>0.0000</i>	0.2918
350	13.02 <i>0.9928</i>	20.09 <i>0.8612</i>	26.88 <i>0.5248</i>	32.62 <i>0.2500</i>	39.06 <i>0.0800</i>	50.82 <i>0.0052</i>	82.88 <i>0.0000</i>	0.2607
400	11.20 <i>0.9980</i>	19.84 <i>0.8702</i>	26.16 <i>0.5643</i>	32.00 <i>0.2745</i>	38.88 <i>0.0829</i>	49.60 <i>0.0072</i>	78.56 <i>0.0000</i>	0.2339
450	11.52 <i>0.9975</i>	18.72 <i>0.9066</i>	24.84 <i>0.6365</i>	29.88 <i>0.3690</i>	36.27 <i>0.1359</i>	47.52 <i>0.0121</i>	70.56 <i>0.0000</i>	0.1419
500	10.60 <i>0.9988</i>	18.20 <i>0.9210</i>	24.20 <i>0.6709</i>	29.40 <i>0.3925</i>	35.60 <i>0.1531</i>	45.50 <i>0.0196</i>	65.80 <i>0.0001</i>	0.1182
550	8.36 <i>0.9999</i>	17.38 <i>0.9407</i>	23.10 <i>0.7280</i>	27.94 <i>0.4676</i>	33.88 <i>0.2049</i>	43.34 <i>0.0323</i>	64.46 <i>0.0001</i>	0.0616
600	5.52 <i>1.0000</i>	16.56 <i>0.9568</i>	22.80 <i>0.7430</i>	27.60 <i>0.4858</i>	33.24 <i>0.2270</i>	42.48 <i>0.0390</i>	68.64 <i>0.0000</i>	0.0415
650	10.14 <i>0.9992</i>	15.86 <i>0.9678</i>	21.58 <i>0.8004</i>	26.52 <i>0.5445</i>	31.72 <i>0.2861</i>	40.43 <i>0.0605</i>	67.60 <i>0.0000</i>	0.0681
700	7.28 <i>1.0000</i>	15.96 <i>0.9664</i>	20.72 <i>0.8368</i>	25.20 <i>0.6169</i>	30.24 <i>0.3518</i>	38.64 <i>0.0869</i>	55.44 <i>0.0015</i>	0.1494
750	7.20 <i>1.0000</i>	15.60 <i>0.9714</i>	20.10 <i>0.8608</i>	24.60 <i>0.6495</i>	29.10 <i>0.4075</i>	37.50 <i>0.1082</i>	58.80 <i>0.0006</i>	0.1743
800	9.60 <i>0.9995</i>	14.40 <i>0.9841</i>	19.52 <i>0.8813</i>	23.68 <i>0.6983</i>	27.84 <i>0.4730</i>	36.48 <i>0.1308</i>	57.92 <i>0.0007</i>	0.2452

^a Upper quantile (uq)

^b Kolmogorov-Smirnov test-statistic value ($D_{0.05} = 0.035$, $D_{0.01} = 0.042$ based on $k = 1,500$)

^c P -value (*italics*) of the above observed χ^2 values

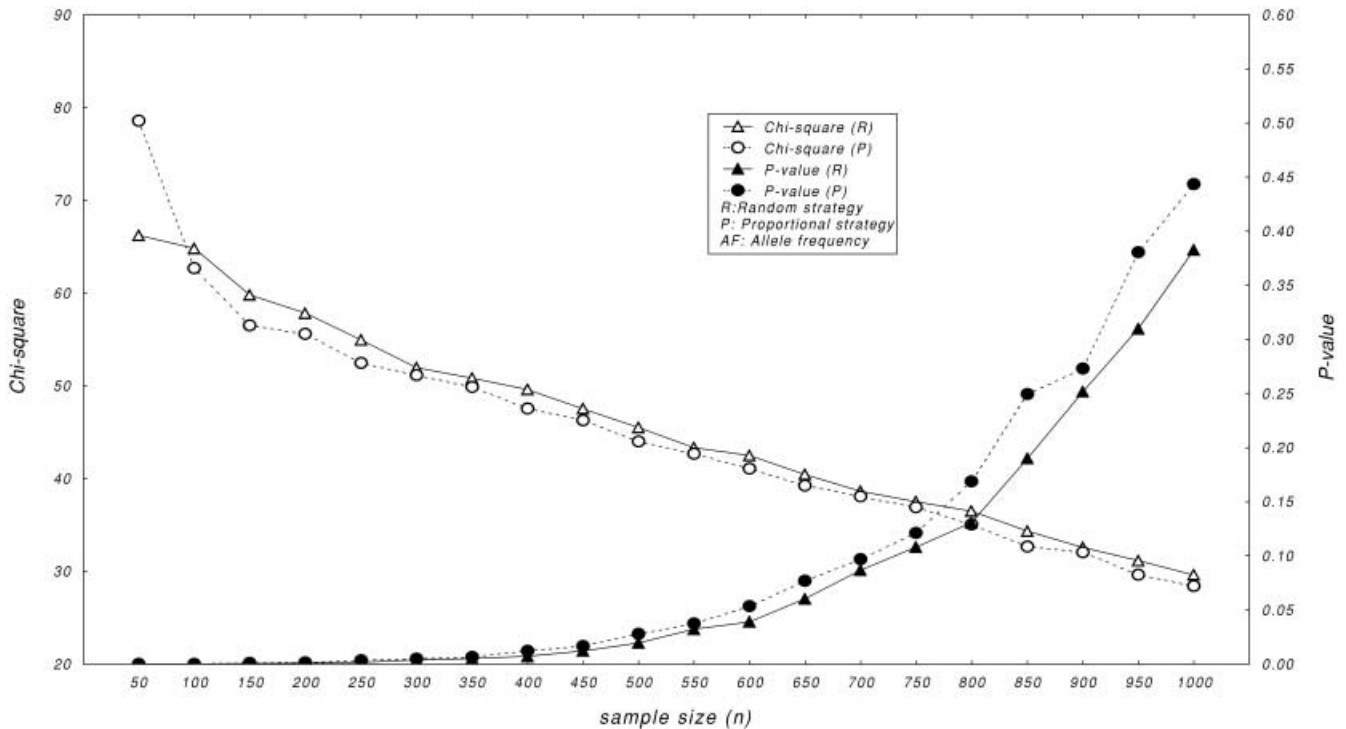


Fig. 1 Observed upper-0.05-quantile χ^2 and P -values for allele frequency under proportional and random strategies

which serves as a lower bound on an optimal n under the P strategy. However, for $\alpha = 0.05$, the corresponding upper-0.05-quantile χ^2 is significant at $n = 550$ having a P -value of 0.0374. At $n = 600$, the P -value ($=0.0533$) of the upper-0.05-quantile χ^2 exceeds $\alpha = 0.05$. At $\alpha = 0.05$ and $n = 600$, for each individual locus the P -value always exceeds $\alpha_j = 0.05/8 = 0.00625$ (Fig. 2). Therefore, $n = 600$ is the optimal n under the P strategy at $\alpha = 0.05$. For more risk to be covered by choosing say, e.g. $\alpha = 0.10$, the optimal n needs to be about 750 (Table 1, Fig. 1). Results for the R-strategy (Table 2, Figs. 1, 2) were similar to that of the P strategy with the difference that the K-S test-statistic was non-significant (at $\alpha = 0.01$) at $n = 600$, with $n = 650$ being the optimal n , which relative to the P strategy exceeds it by 50.

Locus heterozygosity

Tables 3 and 4 list the genome-wide frequency distributions of the 1,500 observed $\chi^2(\text{LH})$ -values for the P and R strategies respectively. Figure 3 depicts the observed upper-0.05-quantile $\chi^2(\text{LH})$ -values and their corresponding P -values for the P and R strategies. Figure 4 shows the values of the upper-0.05-quantile $\chi^2_j(\text{LH})$ -values and their corresponding P -values for individual loci. For the P strategy, the K-S test-statistic was always significant (at $\alpha = 0.05$) for all sample sizes n . However, for all n , the P -value corresponding to the

upper-0.05-quantile χ^2 was always greater than $\alpha = 0.05$. Thus $n = 50$ could be taken as the minimum sample size for a genome-wide fit at $\alpha = 0.05$. Also, for $\alpha = 0.05$ at $n = 50$, Fig. 4 shows that for each individual locus the P -value always much exceeded $\alpha_j = 0.05/8 = 0.00625$. Therefore, $n = 50$ is the optimal n under the P strategy at $\alpha = 0.05$. For the R strategy (Table 4), the K-S test identified $n = 50$ as the lower bound on an optimal n , this n also being the optimal n as the P -value ($= 0.0582$) corresponding to the upper-0.05-quantile χ^2 exceeded $\alpha = 0.05$. The locus-wise results for the R-strategy (Figs. 3, 4) were similar to that of the P strategy. Accordingly, $n = 50$ is also the optimal n for the R strategy at $\alpha = 0.05$.

Optimal sampling strategy and core collection size

Results from the preceding two paragraphs indicate that, for AF and LH considered simultaneously, there is little difference in performance of the P and R strategy, with P performing slightly better than R. A core collection size of about 600 entries selected using either the P or the R strategy is optimal to adequately represent the genetic spectrum of, and to maximally capture the genetic diversity (in terms of LH) in, the GREC.

As evident from the results reported above, LH requires a much-smaller optimal sample size than AF. An optimal sample size chosen solely on the basis of LH is, therefore, not likely to adequately represent the genetic spectrum of the population. A safer approach in arriving at an optimal sample size therefore seems to be to consider the (larger) optimal sample size for AF as the optimal sample size.

Fig. 2 Locus-wise observed upper-0.05-quantile χ^2 and P -values for AF under proportional (P) and random (R) strategies

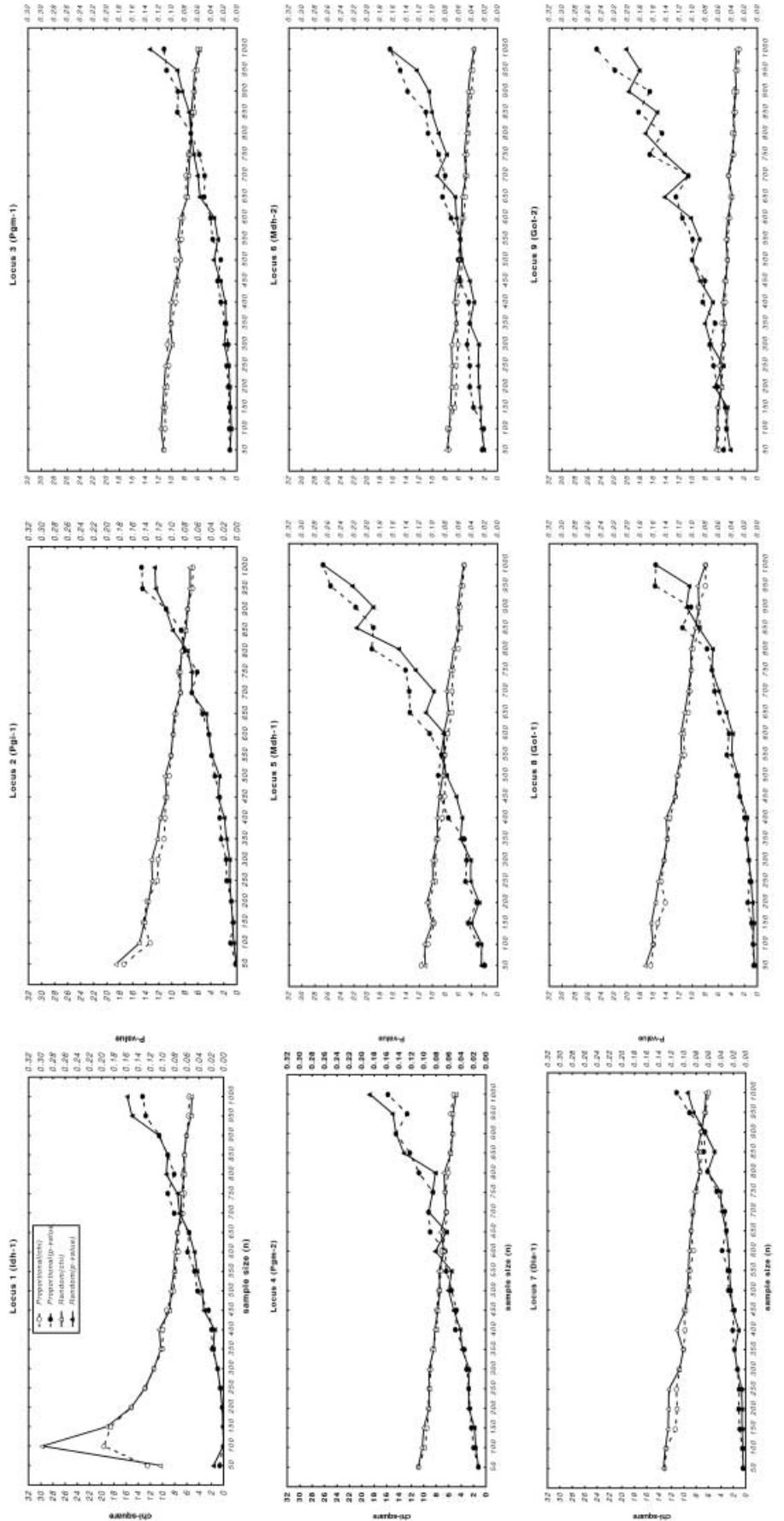


Table 3 Quantiles of 1,500 observed χ^2 values for locus heterozygosity for different sample sizes (n) under proportional strategy

n	Min	0.95- uq^a	0.75- uq	0.50- uq	0.25- uq	0.05- uq	Max	D^b
50	1.36 <i>0.9980^c</i>	3.29 <i>0.9516</i>	5.60 <i>0.7794</i>	7.98 <i>0.5360</i>	10.92 <i>0.2810</i>	15.69 <i>0.0737</i>	30.78 <i>0.0003</i>	0.0485
100	1.22 <i>0.9988</i>	2.92 <i>0.9674</i>	5.36 <i>0.8018</i>	7.77 <i>0.5579</i>	10.64 <i>0.3011</i>	15.54 <i>0.0771</i>	30.61 <i>0.0003</i>	0.0685
150	1.08 <i>0.9992</i>	3.05 <i>0.9621</i>	5.43 <i>0.7958</i>	7.54 <i>0.5811</i>	10.31 <i>0.3261</i>	16.11 <i>0.0646</i>	24.95 <i>0.0030</i>	0.0915
200	0.72 <i>0.9999</i>	2.83 <i>0.9705</i>	5.21 <i>0.8155</i>	7.50 <i>0.5848</i>	10.20 <i>0.3348</i>	14.89 <i>0.0941</i>	26.58 <i>0.0016</i>	0.1033
250	0.95 <i>0.9995</i>	2.98 <i>0.965</i>	5.14 <i>0.8216</i>	7.13 <i>0.6231</i>	9.85 <i>0.3631</i>	14.41 <i>0.1083</i>	32.20 <i>0.0002</i>	0.1311
300	0.38 <i>1.0000</i>	2.74 <i>0.9736</i>	4.99 <i>0.8355</i>	7.05 <i>0.6318</i>	9.59 <i>0.385</i>	13.80 <i>0.1294</i>	25.57 <i>0.0024</i>	0.1429
350	0.87 <i>0.9997</i>	2.72 <i>0.9744</i>	4.77 <i>0.8542</i>	6.78 <i>0.6598</i>	9.39 <i>0.4025</i>	14.16 <i>0.1168</i>	26.36 <i>0.0018</i>	0.1673
400	0.95 <i>0.9995</i>	2.47 <i>0.9817</i>	4.72 <i>0.8577</i>	6.93 <i>0.6443</i>	9.51 <i>0.3915</i>	13.81 <i>0.1292</i>	21.33 <i>0.0113</i>	0.1565
450	0.93 <i>0.9996</i>	2.55 <i>0.9794</i>	4.69 <i>0.8607</i>	6.51 <i>0.6877</i>	8.94 <i>0.4424</i>	13.25 <i>0.1518</i>	21.72 <i>0.0098</i>	0.2058
500	0.61 <i>0.9999</i>	2.36 <i>0.9843</i>	4.56 <i>0.8705</i>	6.28 <i>0.7119</i>	8.65 <i>0.4706</i>	13.00 <i>0.1628</i>	25.70 <i>0.0023</i>	0.2297
550	0.61 <i>0.9999</i>	2.45 <i>0.9822</i>	4.34 <i>0.8878</i>	6.20 <i>0.7193</i>	8.35 <i>0.499</i>	12.76 <i>0.1738</i>	20.08 <i>0.0174</i>	0.2528
600	0.64 <i>0.9999</i>	2.24 <i>0.9871</i>	4.00 <i>0.9115</i>	5.71 <i>0.7686</i>	8.07 <i>0.5274</i>	12.12 <i>0.2069</i>	24.11 <i>0.0041</i>	0.2946
650	0.63 <i>0.9999</i>	2.27 <i>0.9865</i>	4.01 <i>0.9106</i>	5.79 <i>0.761</i>	8.11 <i>0.5229</i>	11.26 <i>0.2581</i>	21.17 <i>0.0119</i>	0.2848
700	0.35 <i>1.0000</i>	2.05 <i>0.9907</i>	3.83 <i>0.9219</i>	5.54 <i>0.7848</i>	7.67 <i>0.5679</i>	11.79 <i>0.2254</i>	19.85 <i>0.0189</i>	0.3250
750	0.51 <i>1.0000</i>	2.22 <i>0.9874</i>	3.77 <i>0.9259</i>	5.26 <i>0.811</i>	7.29 <i>0.6071</i>	10.61 <i>0.3035</i>	17.35 <i>0.0435</i>	0.3616
800	0.88 <i>0.9997</i>	2.16 <i>0.9886</i>	3.75 <i>0.9273</i>	5.30 <i>0.8077</i>	7.13 <i>0.6236</i>	10.82 <i>0.2882</i>	19.21 <i>0.0235</i>	0.3799

^a Upper quantile (uq)

^c P -value (*italics*) of the above observed χ^2 values

^b Kolmogorov-Smirnov test-statistic value ($D_{0.05} = 0.035$, $D_{0.01} = 0.042$ based on $k = 1,500$)

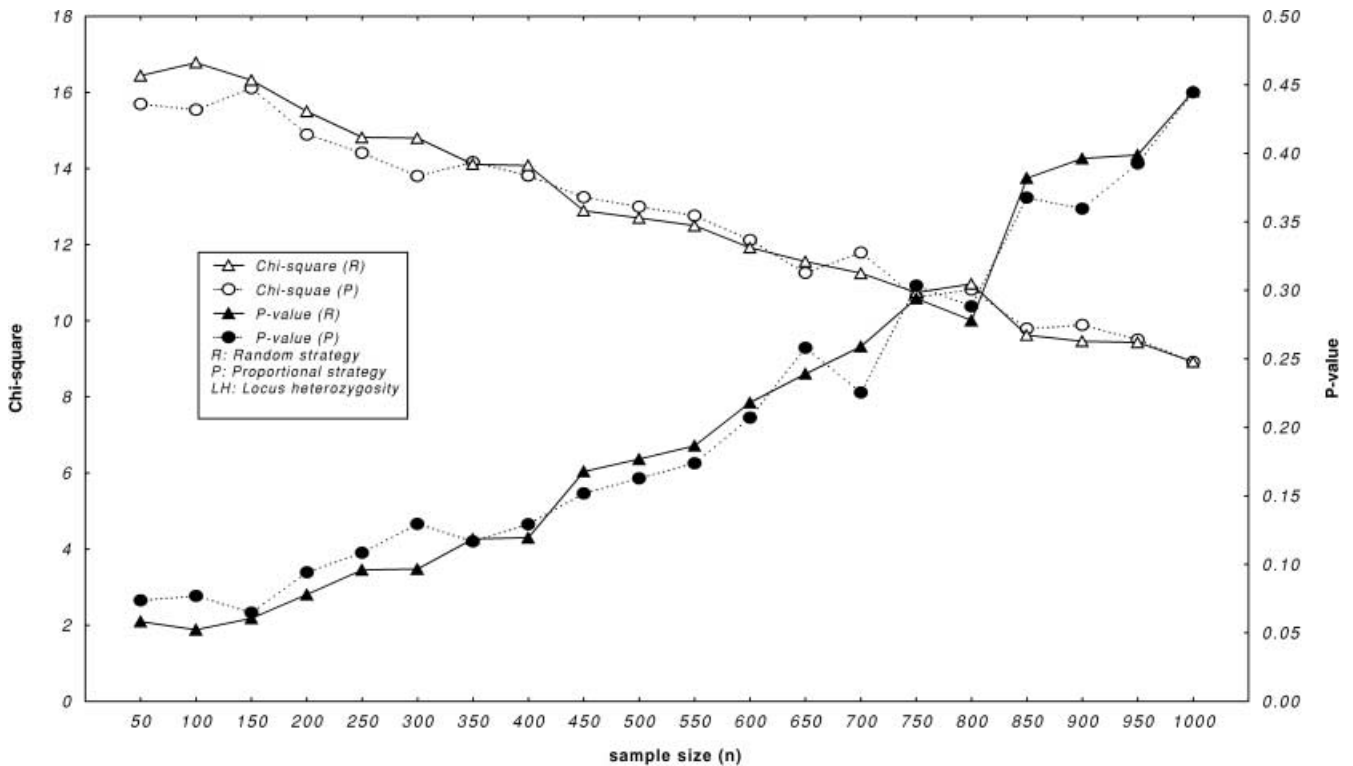


Fig. 3 Observed upper-0.05-quantile χ^2 and P -values for locus heterozygosity under proportional and random strategies

Table 4 Quantiles of 1,500 observed χ^2 values for locus heterozygosity for different sample sizes (n) under random strategy

n	Min	.95- uq^a	0.75- uq	0.50- uq	0.25- uq	0.05- uq	Max	D^b
50	0.85 <i>0.9997^c</i>	3.17 <i>0.9572</i>	5.88 <i>0.7515</i>	8.36 <i>0.4988</i>	11.16 <i>0.2651</i>	16.44 <i>0.0582</i>	30.53 <i>0.0004</i>	0.0226
100	0.74 <i>0.9998</i>	3.00 <i>0.9643</i>	5.57 <i>0.7821</i>	7.95 <i>0.5392</i>	10.86 <i>0.2856</i>	16.78 <i>0.0523</i>	29.62 <i>0.0005</i>	0.0609
150	0.86 <i>0.9997</i>	3.20 <i>0.9559</i>	5.53 <i>0.7863</i>	7.86 <i>0.5485</i>	10.42 <i>0.3177</i>	16.32 <i>0.0605</i>	29.44 <i>0.0005</i>	0.0726
200	0.97 <i>0.9995</i>	2.79 <i>0.9722</i>	5.29 <i>0.8085</i>	7.50 <i>0.5854</i>	10.21 <i>0.3335</i>	15.50 <i>0.0780</i>	28.81 <i>0.0007</i>	0.1050
250	0.83 <i>0.9997</i>	2.97 <i>0.9653</i>	5.10 <i>0.8259</i>	7.40 <i>0.5953</i>	10.12 <i>0.3407</i>	14.82 <i>0.0959</i>	24.11 <i>0.0041</i>	0.1120
300	0.98 <i>0.9995</i>	2.82 <i>0.9711</i>	5.03 <i>0.8319</i>	7.08 <i>0.6286</i>	9.99 <i>0.3516</i>	14.8 <i>0.0966</i>	25.21 <i>0.0027</i>	0.1384
350	0.98 <i>0.9995</i>	2.71 <i>0.9745</i>	4.92 <i>0.8413</i>	6.92 <i>0.6457</i>	9.53 <i>0.3895</i>	14.11 <i>0.1185</i>	27.61 <i>0.0011</i>	0.1618
400	0.87 <i>0.9997</i>	2.74 <i>0.9737</i>	4.90 <i>0.8430</i>	6.91 <i>0.6464</i>	9.58 <i>0.3860</i>	14.08 <i>0.1194</i>	31.87 <i>0.0002</i>	0.1568
450	1.00 <i>0.9994</i>	2.60 <i>0.9780</i>	4.51 <i>0.8744</i>	6.37 <i>0.7028</i>	8.66 <i>0.4689</i>	12.89 <i>0.1677</i>	27.19 <i>0.0013</i>	0.2333
500	0.87 <i>0.9997</i>	2.50 <i>0.9808</i>	4.49 <i>0.8767</i>	6.31 <i>0.7085</i>	8.60 <i>0.4748</i>	12.70 <i>0.1768</i>	23.47 <i>0.0052</i>	0.2318
550	0.99 <i>0.9995</i>	2.39 <i>0.9837</i>	4.35 <i>0.8872</i>	6.29 <i>0.7107</i>	8.45 <i>0.4895</i>	12.50 <i>0.1865</i>	27.9 <i>0.001</i>	0.2560
600	0.81 <i>0.9998</i>	2.26 <i>0.9866</i>	4.16 <i>0.9007</i>	5.86 <i>0.754</i>	8.14 <i>0.5202</i>	11.92 <i>0.2181</i>	24.09 <i>0.0042</i>	0.2933
650	0.91 <i>0.9996</i>	2.42 <i>0.983</i>	4.07 <i>0.9067</i>	5.74 <i>0.7653</i>	7.93 <i>0.5413</i>	11.56 <i>0.2390</i>	21.03 <i>0.0125</i>	0.3071
700	0.83 <i>0.9997</i>	2.26 <i>0.9866</i>	3.91 <i>0.9170</i>	5.48 <i>0.7906</i>	7.54 <i>0.5812</i>	11.25 <i>0.2589</i>	20.67 <i>0.0142</i>	0.3338
750	0.65 <i>0.9999</i>	2.23 <i>0.9872</i>	3.78 <i>0.9251</i>	5.48 <i>0.7907</i>	7.52 <i>0.5834</i>	10.74 <i>0.2941</i>	18.24 <i>0.0325</i>	0.3394
800	0.27 <i>1.0000</i>	2.18 <i>0.9883</i>	3.71 <i>0.9295</i>	5.22 <i>0.8149</i>	7.19 <i>0.6176</i>	10.97 <i>0.2780</i>	18.98 <i>0.0254</i>	0.3739

^a Upper quantile (uq)

^b Kolmogorov-Smirnov test-statistic value ($D_{0.05} = 0.035$, $D_{0.01} = 0.042$ based on $k = 1,500$)

^c P -value (*italics*) of the above observed χ^2 values

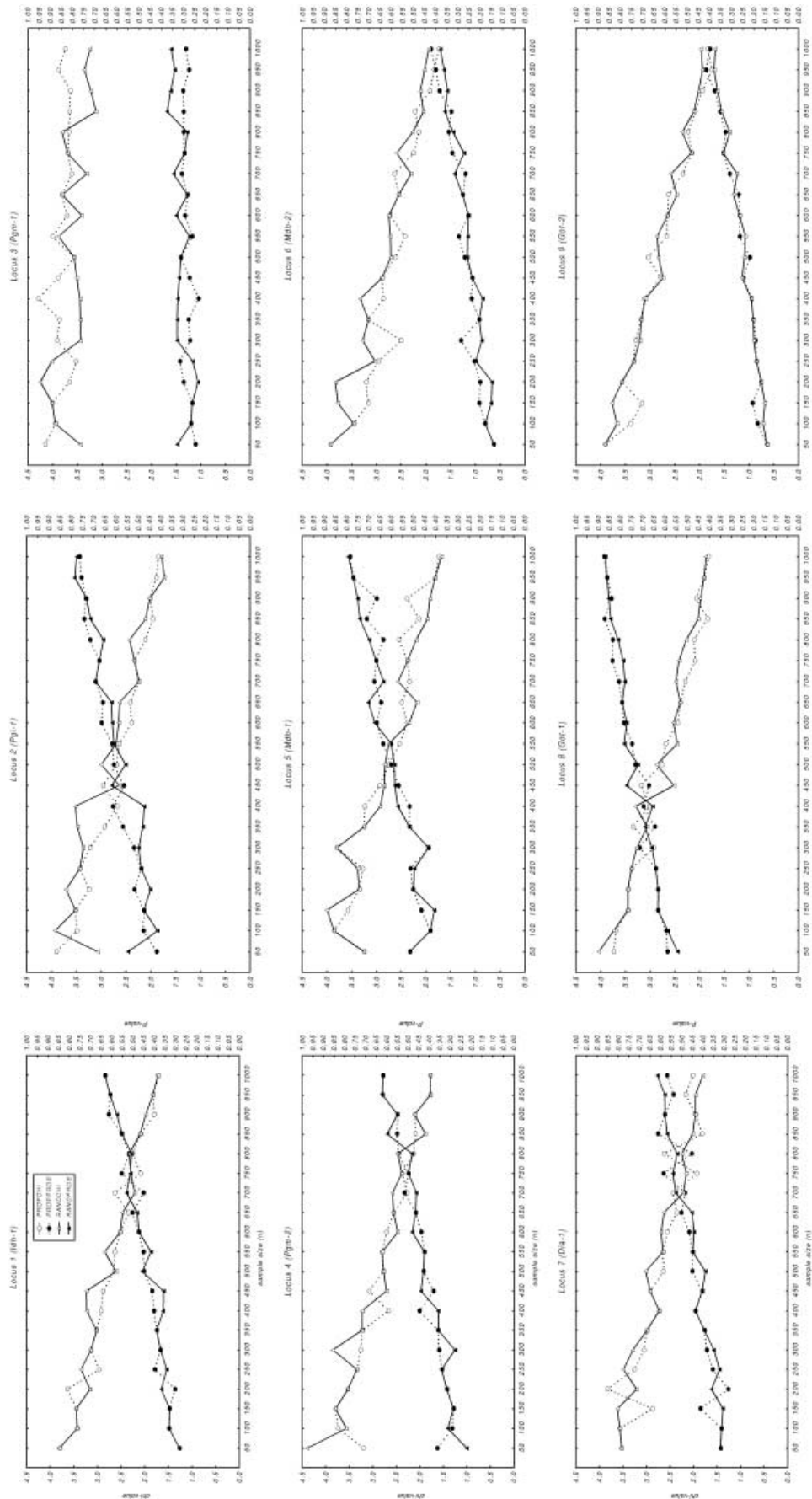
Discussion

Brown (1989) used the sampling theory of Ewens (1972) to propose a fraction $n/N_e = 0.10$ as an optimal sampling fraction for randomly sampling n_r core entries from a germplasm collection of effective population size N_e . By doing so, Brown expected that at least 70% of existent alleles could be retained with 95% certainty. Ewens' sampling theory assumed that the finite germplasm collection contained selectively neutral alleles whose frequencies were in Hardy-Weinberg equilibrium. However, randomly sampling n_r core entries from a finite population of effective size N_e is not genetically equivalent to sampling n_r core entries from N accessions unless genetic duplicates are first removed. We achieved this by choosing to work with GREC, rather than with the original entire collection. However, the assumption of selectively neutral alleles may not hold for many genes that control adaptive traits since these are products of long-term natural and artificial selection. In fact, as pointed out by Yonezawa et al. (1995), the neutrality principle may not hold for some isozymes. The assumption of Hardy-Weinberg equilibrium may also not be valid since accessions in the collection do not interbreed with one another.

The major theoretical argument for core collections in seed crops is that a small number of samples may be efficient in retaining alleles at single loci (Brown 1989). This leads one to presume that the breeders would assemble alleles into genotypes at will in crossing programs. The relative efficiency of a few samples (approximately 10% of N) is attributable to the expectation that the number of alleles increases in proportion to the logarithm of the number N of available samples in the entire collection. However, in clonal crops like potato, much more interest surrounds the whole genotype; specific combinations of genes in highly heterozygous combinations could be worth preserving, and the number of genotypes (genets) preserved increases in direct proportion to the number of samples, assuming duplicates are removed. Realizing these specific features in clonal crops, Brown (1995) suggested that the proportion of entries in the core, rather than fixing at 10%, might have to be higher or lower than 10%. The findings of this research, giving an optimal sampling fraction of $600/1,910 = 0.31$, agrees with Brown's views.

Huaman et al. (2000b) found that a core collection of 306 entries adequately represented their morphologically duplicate-free collection of 2,379 accessions. However, an examination of their Table 2 shows that, with $n = 306$,

Fig. 4 Locus-wise observed upper-0.05-quantile χ^2 and P -values for locus heterozygosity under proportional and random strategies



two loci (Got-1 and Pgi-1) fail to be adequately represented in the population. The sum of individual-locus $\chi^2(\text{AF})$ values in their Table 2 comes to $\chi^2(\text{AF}) = 55.385$ ($df = 28$) with a P -value of 0.0015. This value of $\chi^2(\text{AF})$, corresponding to $n = 306$, is included in the range of 1,500 $\chi^2(\text{AF})$ values for $n = 300$ for both P and R strategies (Tables 1 and 2). This result provides validity to, and confidence in, the simulation approach employed in this study. Table 3 in Huaman *et al.* (2000b) also needs correction in the value of $\chi^2(\text{LH})$, which should have been computed according to the $\chi^2(\text{LH})$ formula given in Materials and methods and should have been 15.647 ($df = 9$; $P = 0.075$) in place of 5.729 ($df = 8$; $P = 0.678$) as reported. Simulation results clearly establish that Huaman *et al.* (2000b) need to revise their optimal core collection size from 306 to about 600 using either the P or the R strategy.

The conclusions regarding optimal core sample size and strategy arrived at for the potato collection obviously hold for the nine isozyme loci for which the accessions in the available collection were characterized. These may change when additional loci are used to characterize the collection.

The simulation approach, developed here using potato isozyme data, could be generally applied on genetic or molecular data of any crop species for identifying the optimal sampling strategy and core collection size, with suitable minor modifications as necessary.

References

- Brown AHD (1989) Core collections: a practical approach to genetic resource management. *Genome* 31:818–824
- Brown AHD (1995) The core collection at the crossroads. In: Hodgkin T, Brown AHD, van Hintum ThJL, Morales EAV (eds) Core collections of plant genetic resources. John Wiley and Sons, New York, pp 3–19
- Douches DS, Ludlam K (1991) Electrophoretic characterization of North American potato cultivars. *Am Potato J* 68:767–780
- Douches DS, Quiros CF (1988) Additional loci in tuber-bearing solanums: inheritance and linkage relationships. *J Hered* 79: 377–384
- Douches DS, Ludlam K, Freyre R (1991) Isozyme and plastid DNA assessment of pedigrees of nineteenth-century potato cultivars. *Theor Appl Genet* 82:192–200
- Ewens WJ (1972) The sampling theory of selectively neutral alleles. *Theor Pop Biol* 3:87–112
- Frankel OH, Brown AHD (1984) Current plant genetic resources – a critical appraisal. In: Genetics: new frontiers (vol IV). Oxford and IBH Publishing, New Delhi, India, pp 1–13
- Freyre R, Douches DS (1994) Development of a model for marker-assisted selection of specific gravity in diploid potato across environments. *Crop Sci* 34:1361–1368
- Freyre R, Warnke S, Sosinski B, Douches DS (1994) Quantitative trait locus analysis of tuber dormancy in diploid potato (*Solanum* spp.). *Theor Appl Genet* 89:474–480
- Huaman Z (1998) Collection, maintenance and evaluation of potato genetic resources. *Plant Var Seeds* 11:29–38
- Huaman Z, Ortiz R, Gomez R (2000a) Selecting a *Solanum tuberosum* ssp. *andigena* core collection using morphological, geographical, disease and pest descriptors. *Am J Potato Res* 77:183–90
- Huaman Z, Ortiz R, Zhang D, Rodriguez F (2000b) Isozyme analysis of entire and core collections of *Solanum tuberosum* subsp. *Andigena* potato cultivars. *Crop Sci* 40:273–276
- Ortiz R, Huaman Z (2001) Allozyme polymorphism in tetraploid potato gene pools and the effect of human selection. *Theor Appl Genet* (in press)
- Ortiz R, Douches DS, Kotch GP, Peloquin SJ (1993) Use of haploids and isozyme markers for genetic analysis in the polyploid potato. *J Genet Breed* 47:283–288
- Quiros CF, Ortega R, van Raamsdock L, Herrera-Montoya M, Cisneros P, Schmidt E, Brush S (1992) Increase of potato genetic resources in their center of diversity: the role of natural outcrossing and selection by the Andean farmers. *Genet Res Crop Evol* 39:107–112
- Rabinowitz D, Linder CR, Ortega R, Begazo D, Murguía H, Douches DS, Quiros CF (1990) High levels of interspecific hybridization between *Solanum sparsipilum* and *S. stenotomum* in experimental plots in the Andes. *Am Potato J* 67:73–81
- Sokal RR, Rohlf FJ (1981) Biometry. W. H. Freeman and Co, New York
- Yonezawa K, Nomura T, Morisima H (1995) Sampling strategies for use in stratified germplasm collections. In: Hodgkin T, Brown AHD, van Hintum ThJL, Morales EAV (eds) Core collections of plant genetic resources. John Wiley and Sons, New York, pp 35–53
- Zimmerer KS, Douches DS (1991) Geographical approaches to crop conservation: the partitioning of genetic diversity in Andean potatoes. *Econ Bot* 45:176–189