1    *Type: Letter/Methods*

2

3    # selscan 2.0: scanning for sweeps in unphased data

4

5    *Zachary A. Szpiech[1,2,*]*

6

7    *[1] Department of Biology, Pennsylvania State University, University Park, PA 16802, USA*
8    *[2] Institute for Computational and Data Sciences, Pennsylvania State University, University Park,*
9    *PA 16802, USA*
10   *[*] Correspondence: szpiech@psu.edu*

11

12   **Abstract**

13      Haplotype-based scans to identify recent and ongoing positive selection have become

14   commonplace in evolutionary genomics studies of numerous species across the tree of life.

15   However, the most widely adopted approaches require phased haplotypes to compute the key

16   statistics. Here we release a major update to the selscan software that re-defines popular

17   haplotype-based statistics for use with unphased "multi-locus genotype" data. We provide

18   unphased implementations of iHS, nSL, XP-EHH, and XP-nSL and evaluate their performance

19   across a range of important parameters in a generic demographic history. Source code and

20   executables are available at https://www.github.com/szpiech/selscan.

21   **1 Introduction**

22      Haplotype-based summary statistics—such as iHS (Voight, et al. 2006), nSL (Ferrer-

23   Admetlla, et al. 2014), XP-EHH (Sabeti, et al. 2007), and XP-nSL (Szpiech, et al. 2021)—have

24   become commonplace in evolutionary genomics studies to identify recent and ongoing positive

25   selection in populations (e.g.,Colonna, et al. 2014; Zoledziewska, et al. 2015; Nedelec, et al.

26   2016; Crawford, et al. 2017; Meier, et al. 2018; Lu, et al. 2019; Zhang, et al. 2020; Salmon, et al.

27   2021). When an adaptive allele sweeps through a population, it leaves a characteristic pattern

28   of long high-frequency haplotypes and low genetic diversity in the vicinity of the allele. These

29   statistics aim to capture these signals by summarizing the decay of haplotype homozygosity as

30   a function of distance from a putatively selected region, either within a single population (iHS

31    and nSL) or between two populations (XP-EHH and XP-nSL). However, each of these statistics

32    presumes that haplotype phase is known.

33        Recent work has shown that converting haplotype data into multi-locus genotype data is

34    an effective approach for using haplotype-based selection statistics such as G12, LASSI, and

35    saltiLASSI (Harris, et al. 2018; Harris and DeGiorgio 2020; DeGiorgio and Szpiech 2021) in

36    unphased data. Recognizing this, we have reformulated the iHS, nSL, XP-EHH, and XP-nSL

37    statistics to use multi-locus genotypes and provided an easy-to-use implementation in selscan

38    2.0 (Szpiech and Hernandez 2014). We also evaluate the performance of these unphased

39    statistics under various generic demographic models.

## 2 New Approaches

41        When the --unphased flag is set in selscan v2.0+, biallelic genotype data is collapsed

42    into multi-locus genotype data by representing the genotype as either 0, 1, or 2—the number of

43    derived alleles observed. In this case, selscan v2.0+ will then compute iHS, nSL, XP-EHH, and

44    XP-nSL as described below. We follow the notation conventions of Szpiech and Hernandez

45    (2014).

### 2.1 Extended Haplotype Homozygosity

47        In a sample of $n$ diploid individuals, let $\mathcal{C}$ denote the set of all possible genotypes at

48    locus $x_0$. For multi-locus genotypes, $\mathcal{C} := \{0,1,2\}$, representing the total counts of a derived

49    allele. Let $\mathcal{C}(x_i)$ be the set of all unique haplotypes extending from site $x_0$ to site $x_i$ either

50    upstream or downstream of $x_0$. If $x_1$ is a site immediately adjacent to $x_0$, then $\mathcal{C}(x_1) :=$

51    $\{00,01,02,10,11,12,20,21,22\}$, representing all possible two-site multi-locus genotypes. We can

52    then compute the extended haplotype homozygosity (EHH) of a set of multi-locus genotypes as

53
$$EHH(x_i) \sum_{h \in \mathcal{C}(x_i)} \frac{\binom{n_h}{2}}{\binom{n}{2}},$$

54    where $n_h$ is the number of observed haplotypes of type $h$.

55        If we wish to compute the EHH of a subset of observed haplotypes that all contain the

56    same 'core' multi-locus genotype, let $\mathcal{H}_c(x_i)$ be the partition of $\mathcal{C}(x_i)$ containing genotype $c \in \mathcal{C}$

57    at $x_0$. For example, choosing a homozygous derived genotype ($c = 2$) as the core, $\mathcal{H}_2 :=$

58    $\{20,21,22\}$. Thus, we can compute the EHH of all individuals carrying a given genotype at site $x_0$

59    extending out to site $x_i$ as

60
$$EHH_c(x_i) = \sum_{h \in \mathcal{H}_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}},$$

61    where $n_h$ is the number of observed haplotypes of type $h$ and $n_c$ is the number of observed

62    multi-locus genotypes with core genotype of $c$. Finally, we can compute the complement EHH of

63    a sample of multi-locus genotypes as

64
$$cEHH_c(x_i) = \sum_{h \in \mathcal{C}(x_i) \setminus \mathcal{H}_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_{c'}}{2}},$$

65    where $n_{c'}$ is the number of observed multi-locus genotypes with a core genotype of not $c$.

66    **2.2 iHS and nSL**

67        Unphased iHS and nSL are calculated using the equations above. First, we compute the

68    integrated haplotype homozygosity (iHH) for the homozygous ancestral ($c = 0$) and derived ($c =$

69    2) core genotypes as

70
$$iHH_c = \sum_{i=1}^{|\mathcal{D}|} \frac{1}{2}\big(EHH_c(x_{i-1}) + EHH_c(x_i)\big)\mathrm{d}(x_{i-1}, x_i) + \sum_{i=1}^{|\mathcal{U}|} \frac{1}{2}\big(EHH_c(x_{i-1}) + EHH_c(x_i)\big)\mathrm{d}(x_{i-1}, x_i),$$

71    where $\mathcal{D}$ is the set of downstream sites from the core locus and $\mathcal{U}$ is the set of upstream sites.

72    $d(x_{i-1}, x_i)$ is a measure of genomic distance between to markers and is the genetic distance in

73    centimorgans or physical distance in basepairs for iHS (Voight, et al. 2006) or the number of

74    sites observed for nSL (Ferrer-Admetlla, et al. 2014). We similarly compute the complement

75    integrated haplotype homozygosity (ciHH) for both homozygous core genotypes as

76
$$ciHH_c = \sum_{i=1}^{|\mathcal{D}|} \frac{1}{2}\big(cEHH_c(x_{i-1}) + cEHH_c(x_i)\big)\mathrm{d}(x_{i-1}, x_i)$$

77
$$+ \sum_{i=1}^{|\mathcal{U}|} \frac{1}{2}\big(cEHH_c(x_{i-1}) + cEHH_c(x_i)\big)\mathrm{d}(x_{i-1}, x_i).$$

78   The (unstandardized) unphased iHS is then calculated as

79
$$iHS = \begin{cases} iHS_2, & \text{if } iHS_2 > iHS_0 \\ \text{-}iHS_0, & \text{otherwise} \end{cases},$$

80   where $iHS_2 = \log_{10}\left(\frac{iHH_2}{ciHH_2}\right)$ and $iHS_0 = \log_{10}\left(\frac{iHH_0}{ciHH_0}\right)$. Unstandardized iHS scores are then

81   normalized in frequency bins, as previously described (Voight, et al. 2006; Ferrer-Admetlla, et

82   al. 2014). Unstandardized unphased nSL is computed similarly with the appropriate distance

83   measure. Large positive scores indicate long high-frequency haplotypes with a homozygous

84   derived core genotype, and large negative scores indicate long high-frequency haplotypes with

85   a homozygous ancestral core genotype. Clusters of extreme scores in both directions indicate

86   evidence for a sweep.

87   **2.3 XP-EHH and XP-nSL**

88       Unphased XP-EHH and XP-nSL are calculated by comparing the iHH between

89   populations $A$ and $B$, using the entire sample in each population. iHH in a population P is

90   computed as

91
$$iHH_P = \sum_{i=1}^{|\mathcal{D}|} \frac{1}{2}\big(EHH(x_{i-1}) + EHH(x_i)\big)\mathrm{d}(x_{i-1}, x_i) + \sum_{i=1}^{|\mathcal{U}|} \frac{1}{2}\big(EHH(x_{i-1}) + EHH(x_i)\big)\mathrm{d}(x_{i-1}, x_i),$$

92   where the distance measure is given as centimorgans or basepairs for XP-EHH (Sabeti, et al.

93   2007) and number of sites observed for XP-nSL (Szpiech, et al. 2021). The XP statistics

94   between population $A$ and $B$ are then computed as $XP = \log_{10}\left(\frac{iHH_A}{iHH_B}\right)$ and are normalized

95   genome wide in a single bin. Large positive scores indicate long high-frequency haplotypes in

96    population $A$, and large negative scores indicate long high-frequency haplotypes in population

97    $B$. Clusters of extreme scores in one direction indicate evidence for a sweep in that population.

## 3 Methods

### 3.1 Simulations

100    We evaluate the performance of the unphased versions of iHS, nSL, XP-EHH, and XP-

101    nSL under a generic two-population divergence model using the coalescent simulation program

102    discoal (Kern and Schrider 2016). We explore five versions of this generic model and name

103    them Demo 1 through Demo 5 (Table 1). Let $N_0$ and $N_1$ be the effective population sizes of

104    Population 0 and Population 1 after the split from their ancestral population (of size $N_A$). For

105    Demo 1, we keep a constant population size post-split and let $N_0 = N_1 = 10,000$. For Demo 2,

106    we keep a constant population size post-split and let $N_0 = 2N_1 = 10,000$. For Demo 3, we keep

107    a constant population size post-split and let $2N_0 = N_1 = 10,000$. For Demo 4, we initially set

108    $N_0 = N_1 = 10,000$ and let $N_0$ grow stepwise exponentially every 50 generations starting at 2,000

109    generations ago until $N_0 = 5N_1 = 50,000$. For Demo 5, we initially set $N_0 = N_1 = 10,000$ and let

110    $N_1$ grow stepwise exponentially every 50 generations starting at 2,000 generations ago until

111    $5N_0 = N_1 = 50,000$.

112    For each demographic history we vary the population divergence time $t_d \in$

113    $\{2000, 4000, 8000\}$ generations ago. For non-neutral simulations, we simulate a sweep in

114    Population 0 in the middle of the simulated region across a range of selection coefficients $s \in$

115    $\{0.005, 0.01, 0.02\}$. We vary the frequency at which the adaptive allele starts sweeping as $e \in$

116    $\{0, 0.01, 0.02, 0.05, 0.10\}$, where $e = 0$ indicates a hard sweep and $e > 0$ indicates a soft sweep,

117    and we also vary the frequency of the selected allele at time of sampling $f \in \{0.7, 0.8, 0.9, 1.0\}$

118    as well as $g \in \{50, 100\}$ representing fixation of the sweeping allele $g$ generations ago. For all

119    simulations we set the genome length to be $L = 500,000$ basepairs, the ancestral effective

120    population size to be $N_A = 10,000$, the per site per generation mutation rate at $\mu = 2.35 \times 10^{-8}$,

121    and the per site per generation recombination rate at $r = 1.2 \times 10^{-8}$. For neutral simulations, we

122    simulate 1,000 replicates for each parameter set, and for non-neutral simulations we simulate

123    100 replicates for each parameter set. As iHS and nSL are single population statistics, we only

124    analyze Demo 1, Demo 3, and Demo 4 with these statistics, as Demo 2 and Demo 5 have a

125    constant size history identical to Demo 1 for Population 0, where the sweeps are simulated.

126         For all simulations, we compute the relevant statistics (--ihs, --nsl, --xpehh, or --xpnsl)

127    with selscan v2.0, using the --unphased and --trunc-ok flags. For iHS and XP-EHH, we also use

128    the --pmap flag in order to use physical distance instead of a recombination map.

129    **3.2 Power and False Positive Rate**

130         To compute power for iHS and nSL, we follow the approach of Voight et al. (2006). For

131    these statistics, each non-neutral replicate is individually normalized jointly with all matching

132    neutral replicates in 1% allele frequency bins. Because extreme values of the statistic are likely

133    to be clustered along the genome (Voight, et al. 2006), we then compute the proportion of

134    extreme scores ($|iHS| > 2$ or $|nSL| > 2$) within 100kbp non-overlapping windows. We then bin

135    these windows into 10 quantile bins based on the number of scores observed in each window

136    and call the top 1% of these windows as putatively under selection. We calculate the proportion

137    of non-neutral replicates that fall in this top 1% as the power. To compute the false positive rate,

138    we compute the proportion of neutral simulations that fall within the top 1%.

139         To compute power for XP-EHH and XP-nSL, we follow the approach of Szpiech et al.

140    (2021). For these statistics, each non-neutral replicate is individually normalized jointly with all

141    matching neutral replicates. Because extreme values of the statistic are likely to be clustered

142    along the genome (Szpiech, et al. 2021), we then compute the proportion of extreme scores

143    (XP-EHH $> 2$ or XP-nSL $> 2$) within 100kbp non-overlapping windows. We then bin these

144    windows into 10 quantile bins based on the number of scores observed in each window and call

145    the top 1% of these windows as putatively under selection. We calculate the proportion of non-

146    neutral replicates that fall in this top 1% as the power. To compute the false positive rate, we

147    compute the proportion of neutral simulations that fall within the top 1%.

## 4 Results

149        We find that the unphased versions of iHS and nSL have good power (Figures 1, S1-S4,

150    S13-16, and S21-24) to detect selection prior to fixation of the allele, with nSL generally

151    outperforming iHS. In smaller populations (Figure 1C and 1D), power does suffer relative to

152    larger populations (Figure 1A, 1B, 1E, 1F). Each of these statistics also have low false positive

153    rates hovering around 1% (Table S1).

154        Similarly, we find that the unphased versions of XP-EHH and XP-nSL have good power

155    as well (Figures 2, 3, S5-S12, S17-S20, and S25-S32). When the sweep takes place in the

156    smaller of the two populations (Figure 2C and 2D), we see a similar decrease in power. When

157    one population is undergoing exponential growth (Figure 3) performance is generally quite

158    good, likely the result of a larger effective selection coefficient in large populations. These two-

159    population statistics generally outperform their single-population counterparts, especially for

160    sweeps that have reached fixation recently. Each of these statistics also have low false positive

161    rates hovering around 1% (Table S1).

## 5 Discussion

163        We introduce multi-locus genotype versions of four popular haplotype-based selection

164    statistics—iHS (Voight, et al. 2006), nSL (Ferrer-Admetlla, et al. 2014), XP-EHH (Sabeti, et al.

165    2007), and XP-nSL (Szpiech, et al. 2021)—that can be used when the phase of genotypes is

166    unknown. We implement these updates in the latest v2.0 update of the program selscan

167    (Szpiech and Hernandez 2014), with source code and pre-compiled binaries available at

168    https://www.github.com/szpiech/selscan.

## 6 Acknowledgements

173



174
175    **Figure 1**. Power curves for unphased implementations of iHS (A, C, and E) and nSL (B, D, and

176    F) under demographic histories Demo 1 (A and B), Demo 3 (C and D), and Demo 4 (E and F). $s$

177    is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the

178    number of generations at time of sampling since fixation, $e$ is the frequency at which selection

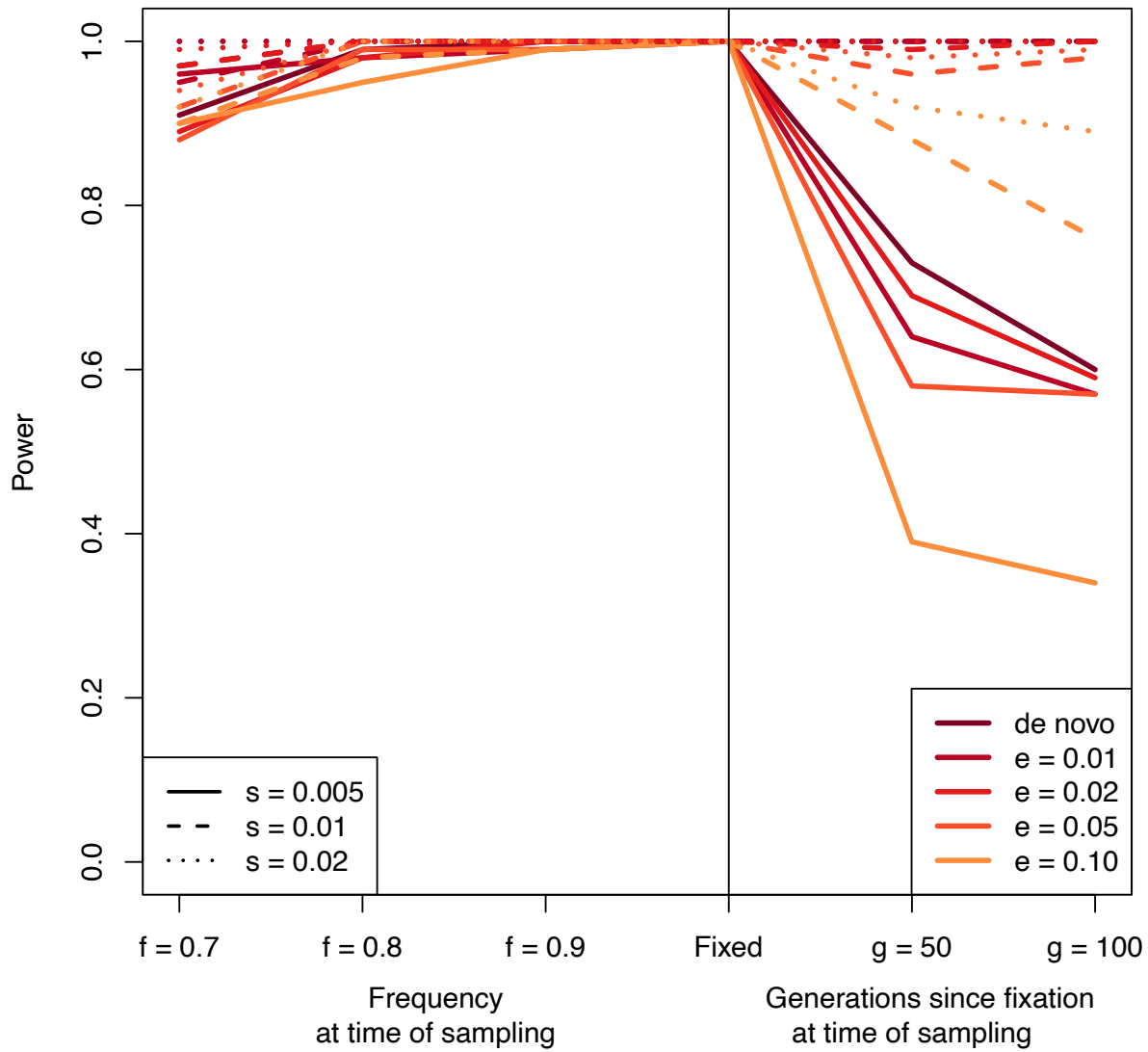179    began, and $t_d$ is the time in generations since the two populations diverged.

**Figure 2**. Power curves for unphased implementations of XP-EHH (A, C, and E) and XP-nSL (B, D, and F) under demographic histories Demo 1 (A and B), Demo 2 (C and D), and Demo 3 (E and F). $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
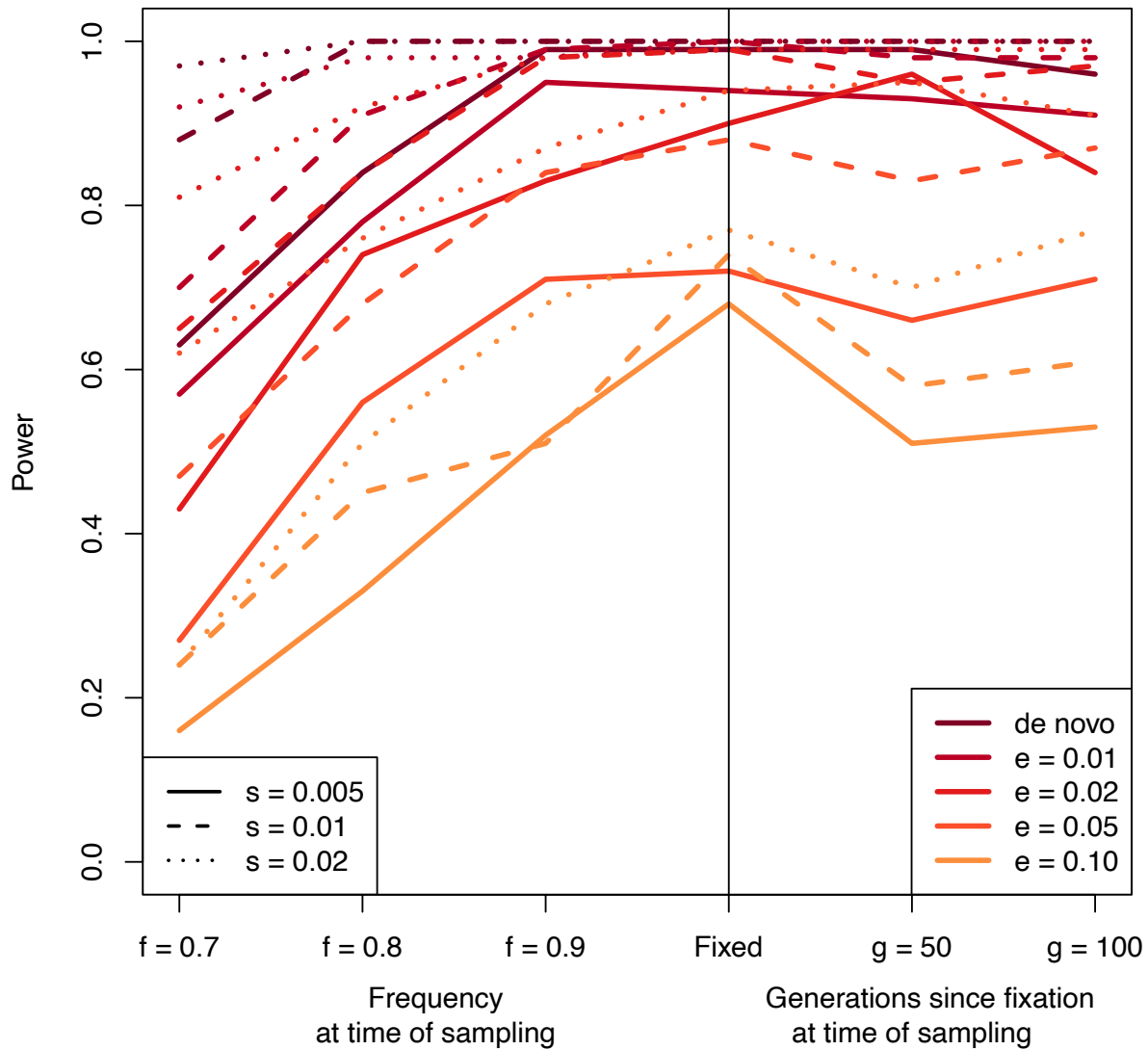
**Figure 3**. Power curves for unphased implementations of XP-EHH (A and C) and XP-nSL (B and D) under demographic histories Demo 4 (A and B), and Demo 5 (C and D). $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
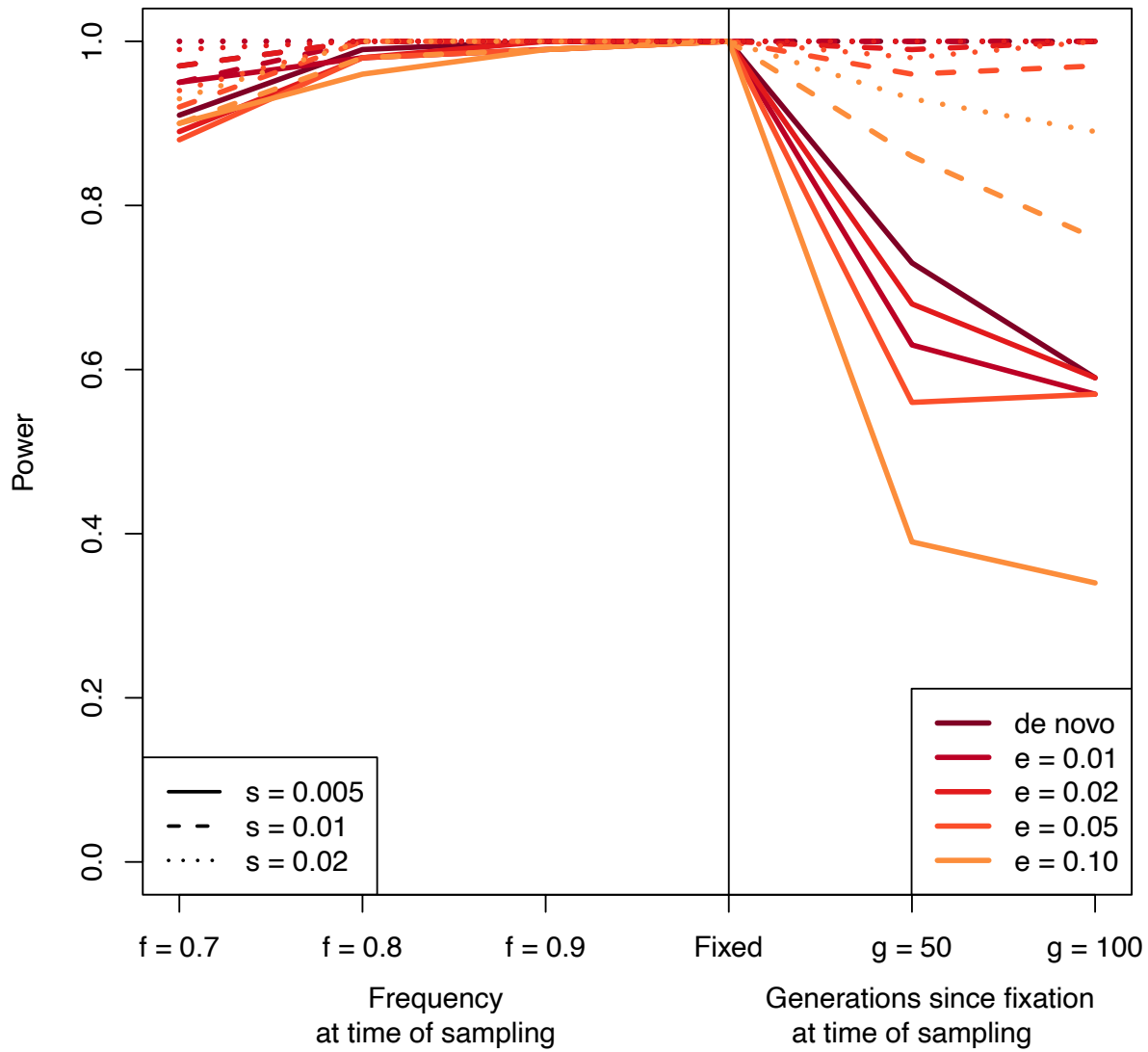
**Figure S1**. Demo 1 iHS $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

**Figure S2**. Demo 1 iHS $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
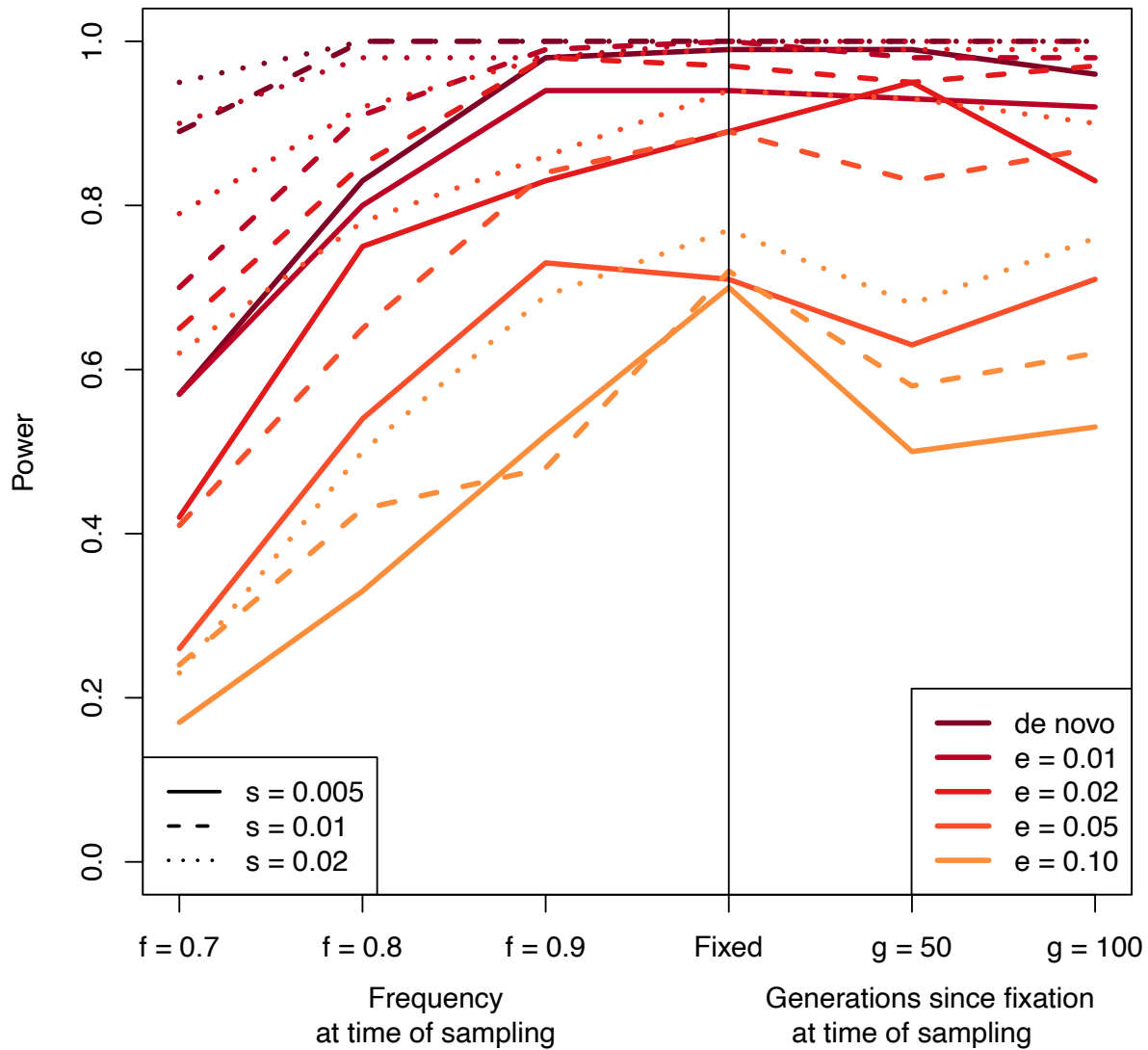
**Figure S3**. Demo 1 nSL $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
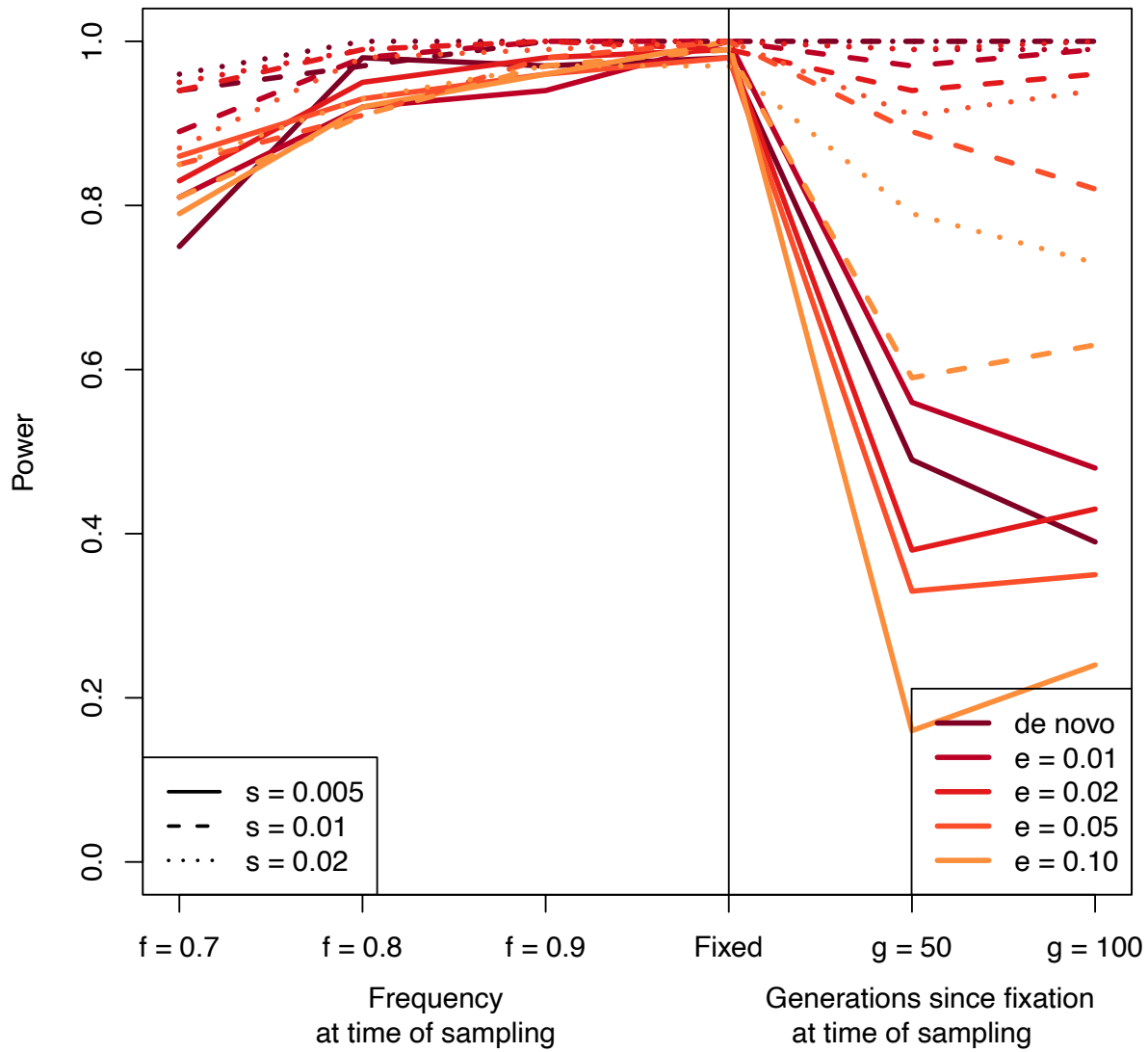
**Figure S4**. Demo 1 nSL $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
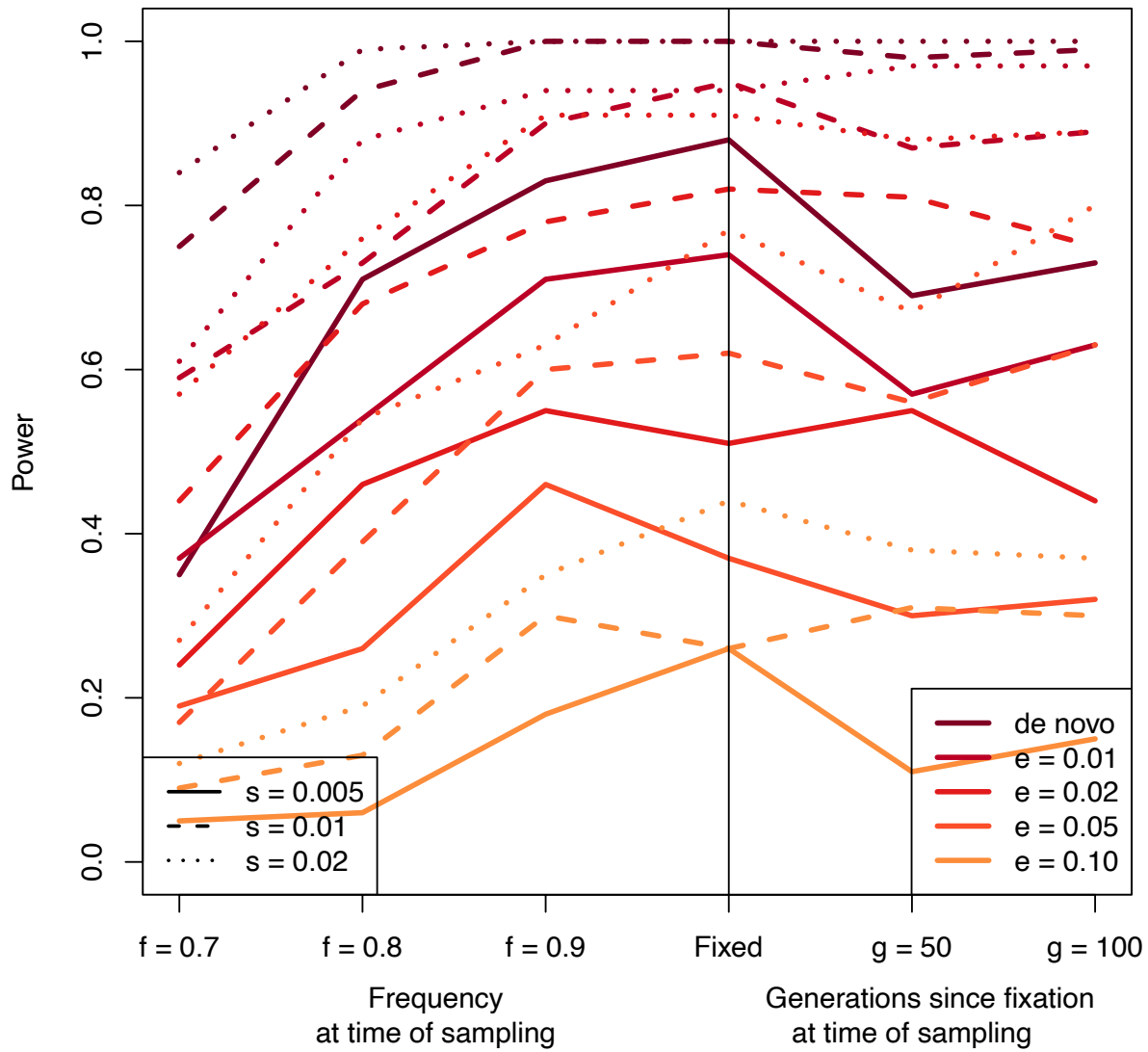
**Figure S5**. Demo 1 XP-EHH $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

**Figure S6**. Demo 1 XP-EHH $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

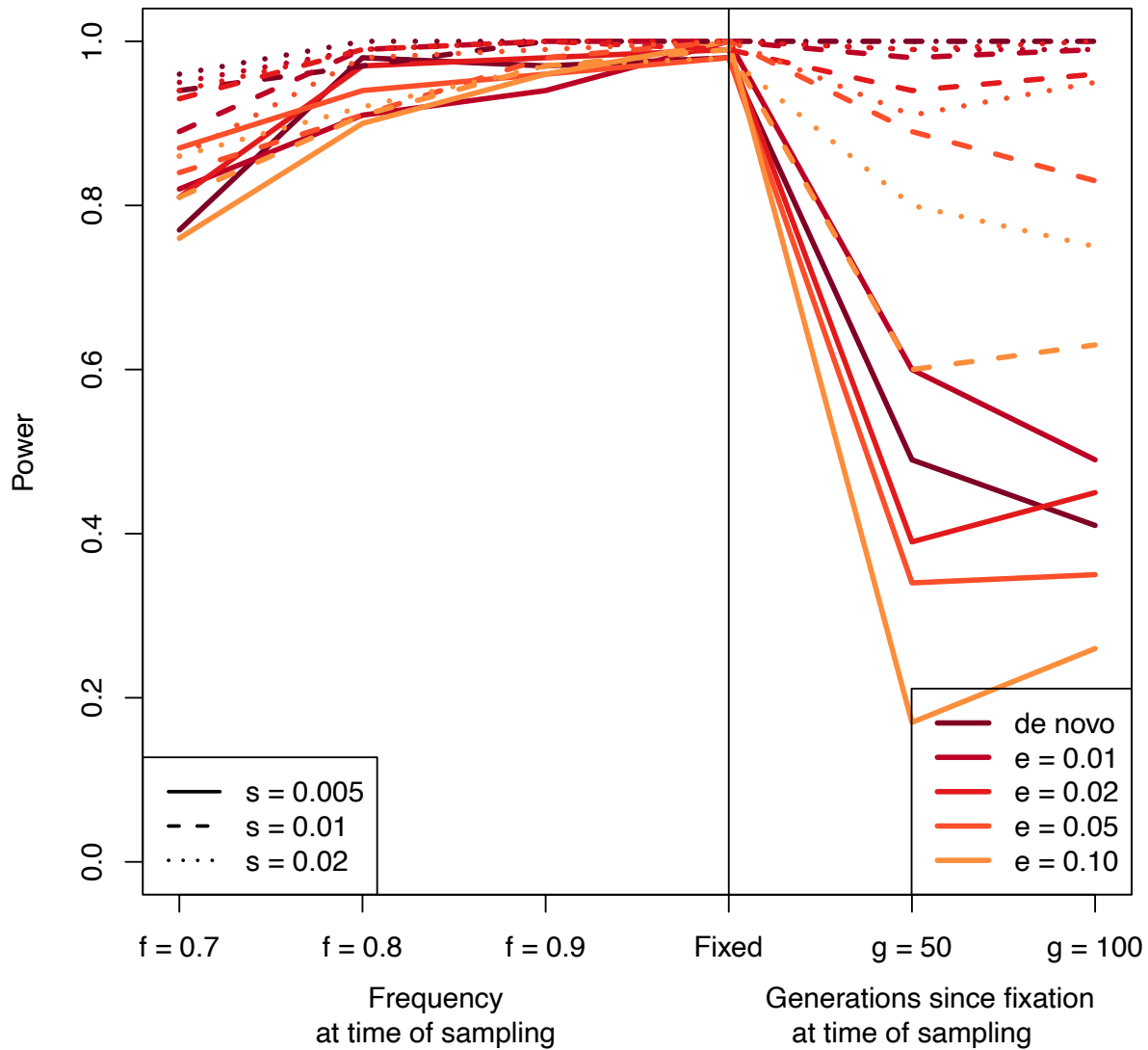**Figure S7**. Demo 1 XP-nSL $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

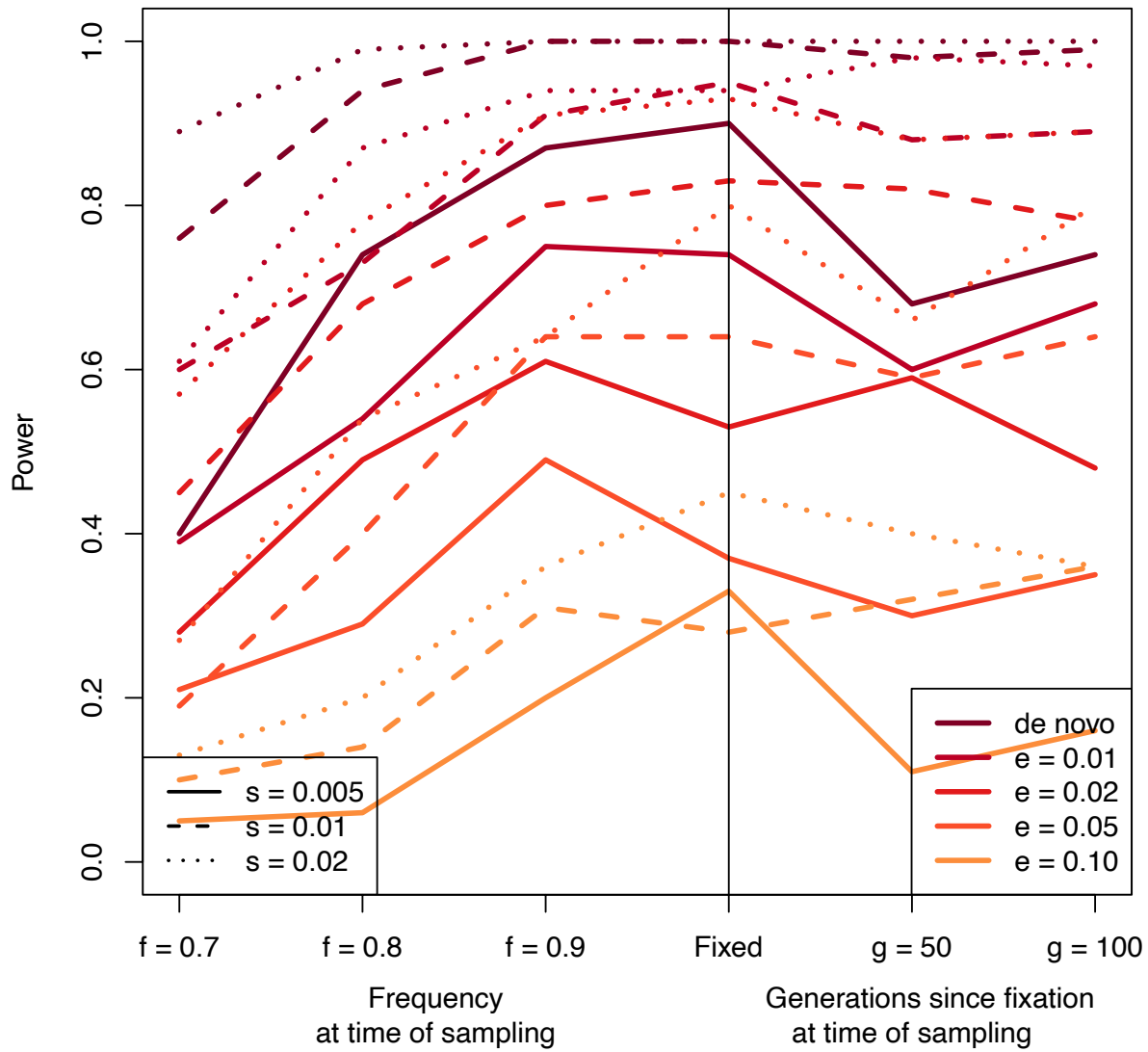XP−nSL; $t_d$ = 8000



**Figure S8**. Demo 1 XP-nSL $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
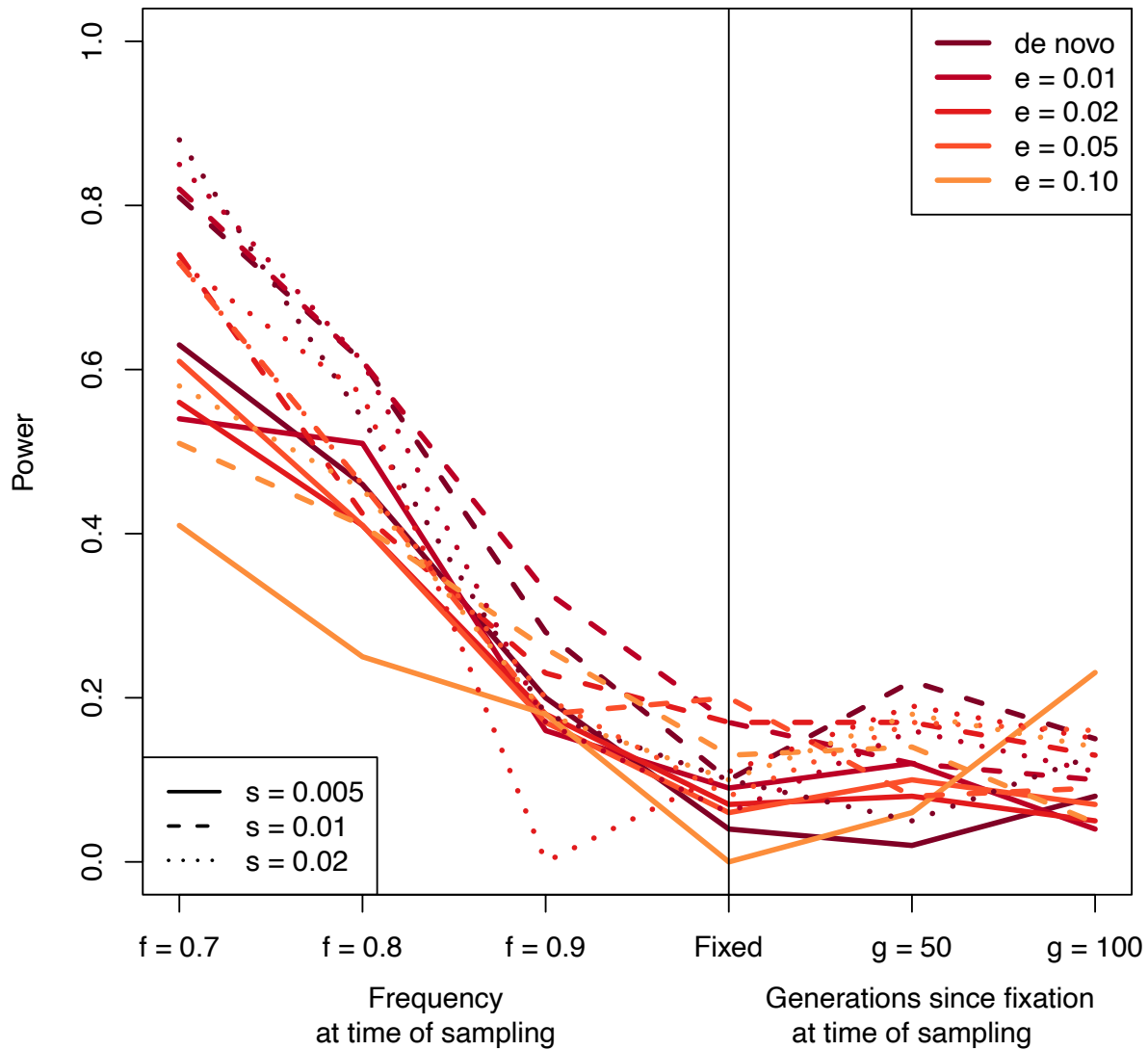
**Figure S9**. Demo 2 XP-EHH $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

**Figure S10**. Demo 2 XP-EHH $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

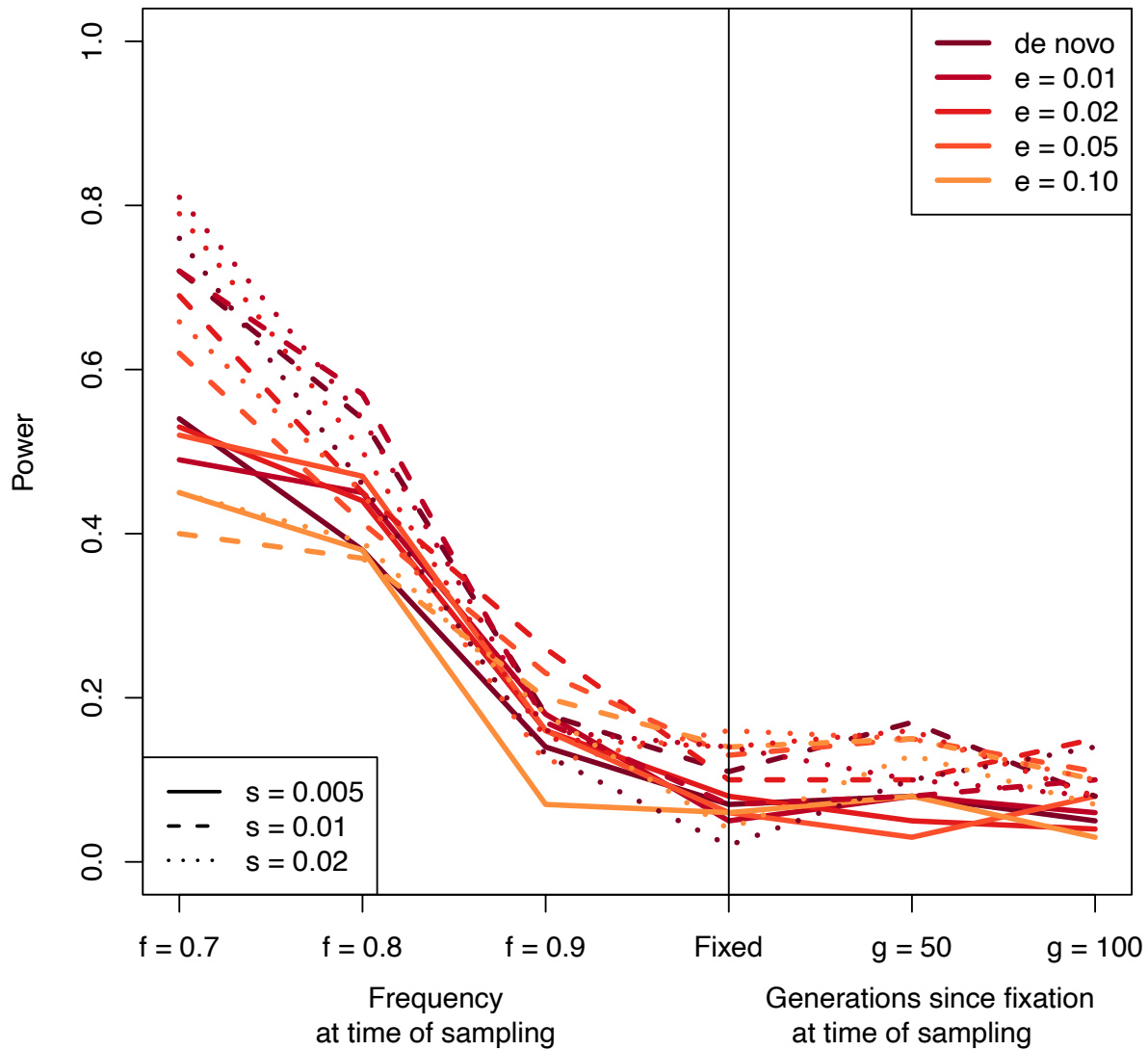**Figure S11**. Demo 2 XP-nSL $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
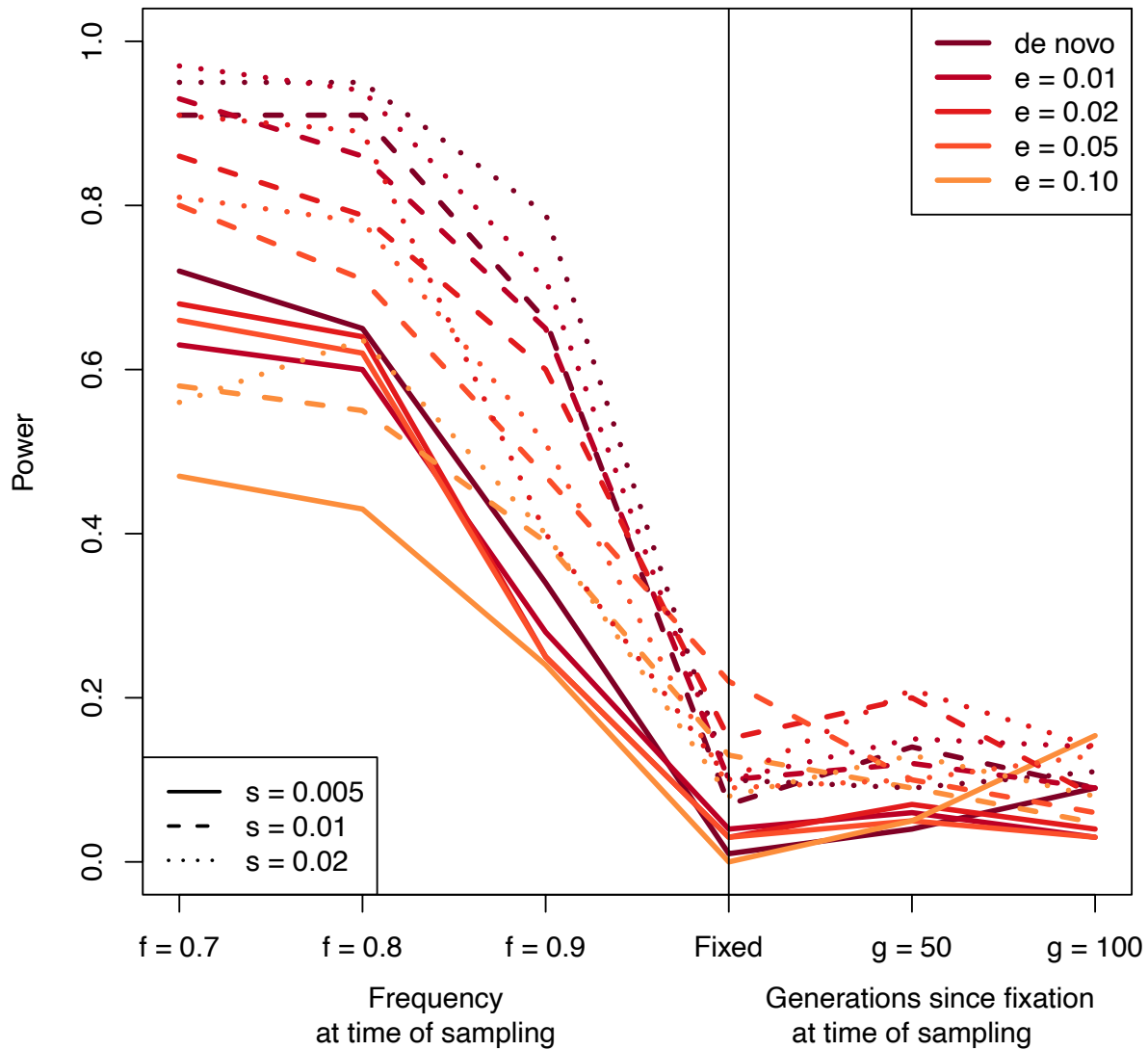
**Figure S12**. Demo 2 XP-nSL $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

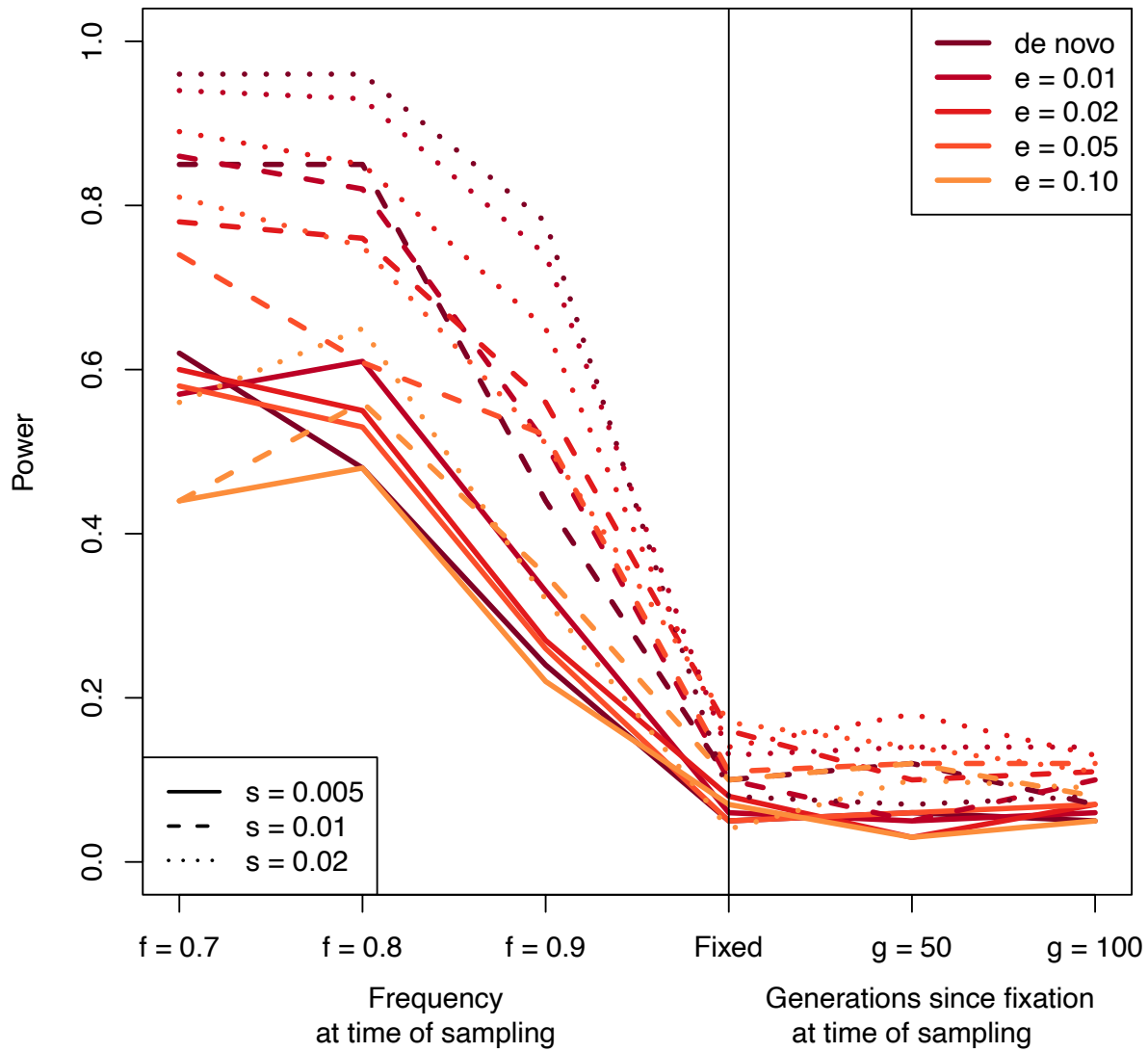**Figure S13**. Demo 3 iHS $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

**Figure S14**. Demo 3 iHS $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

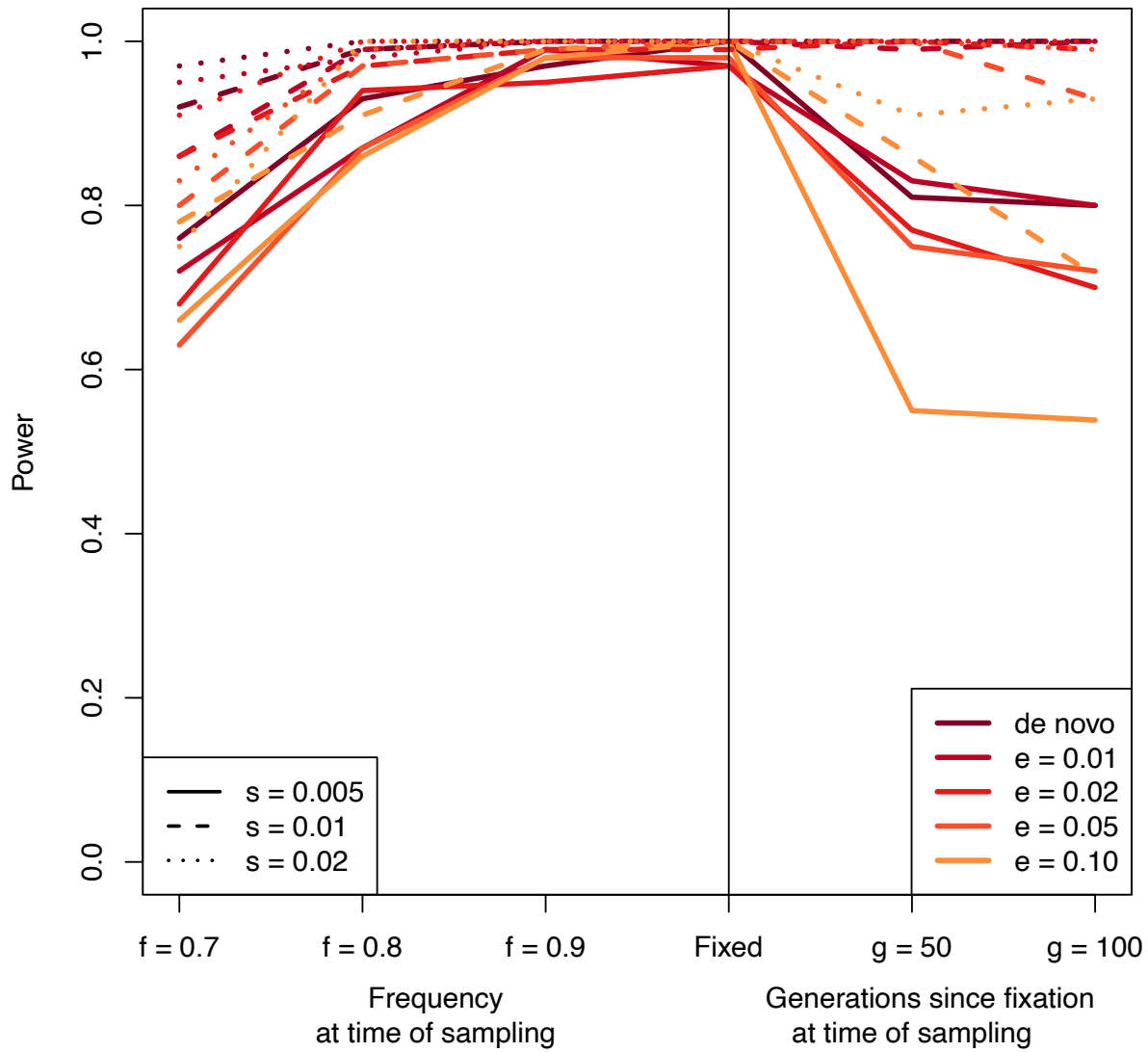**Figure S15**. Demo 3 nSL $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

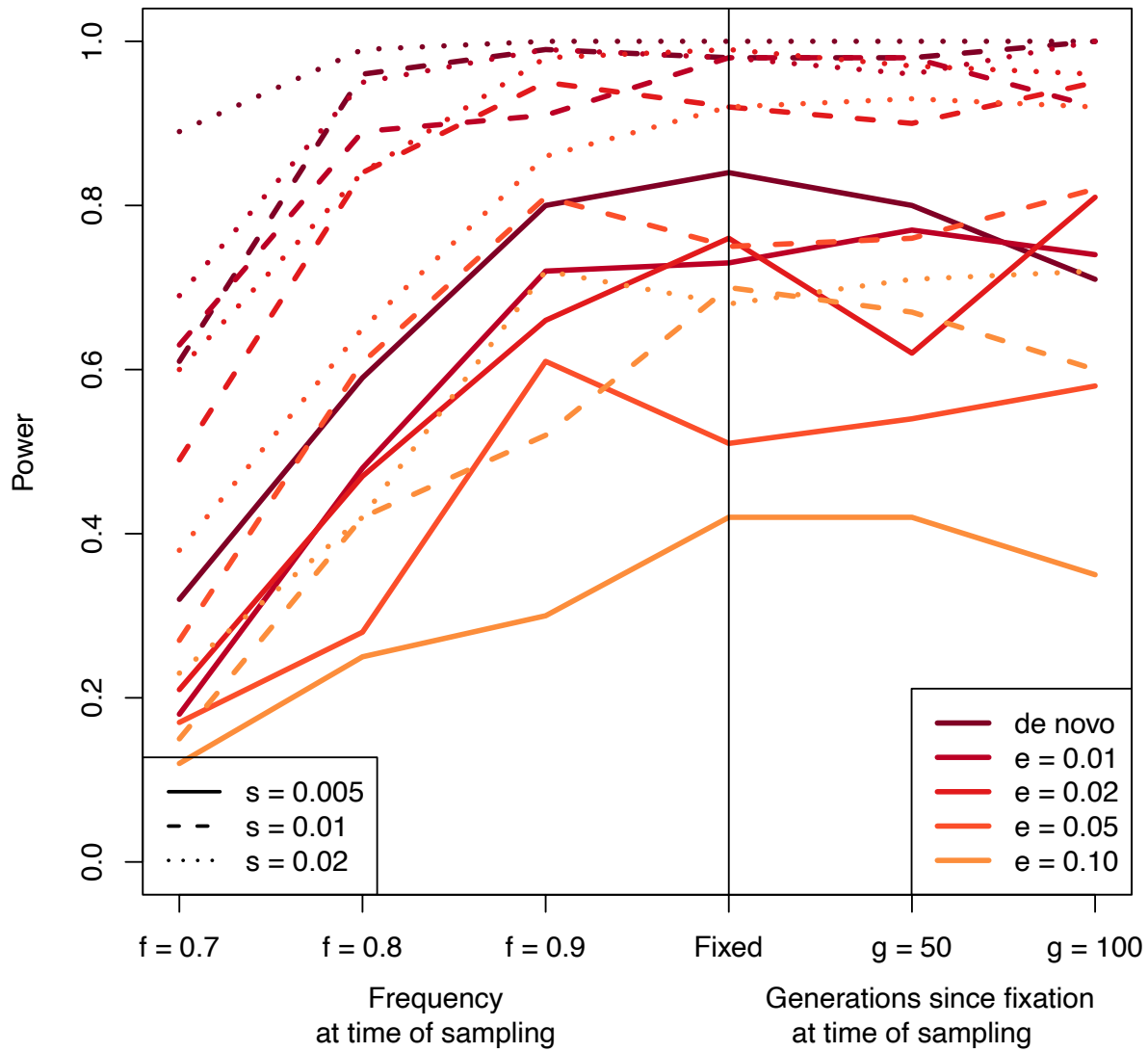**Figure S16**. Demo 3 nSL $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

**Figure S17.** Demo 3 XP-EHH $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

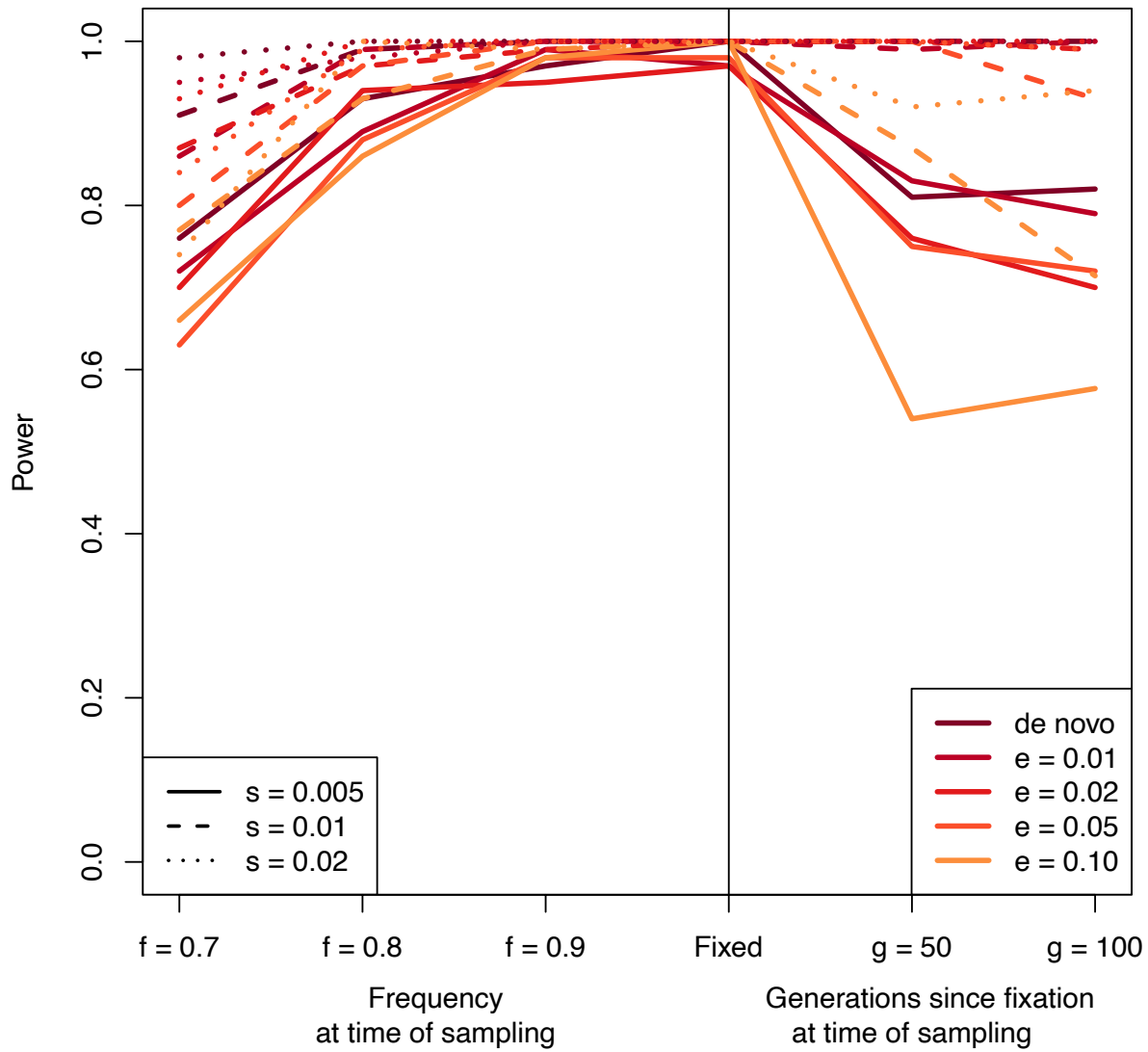**Figure S18**. Demo 3 XP-EHH $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

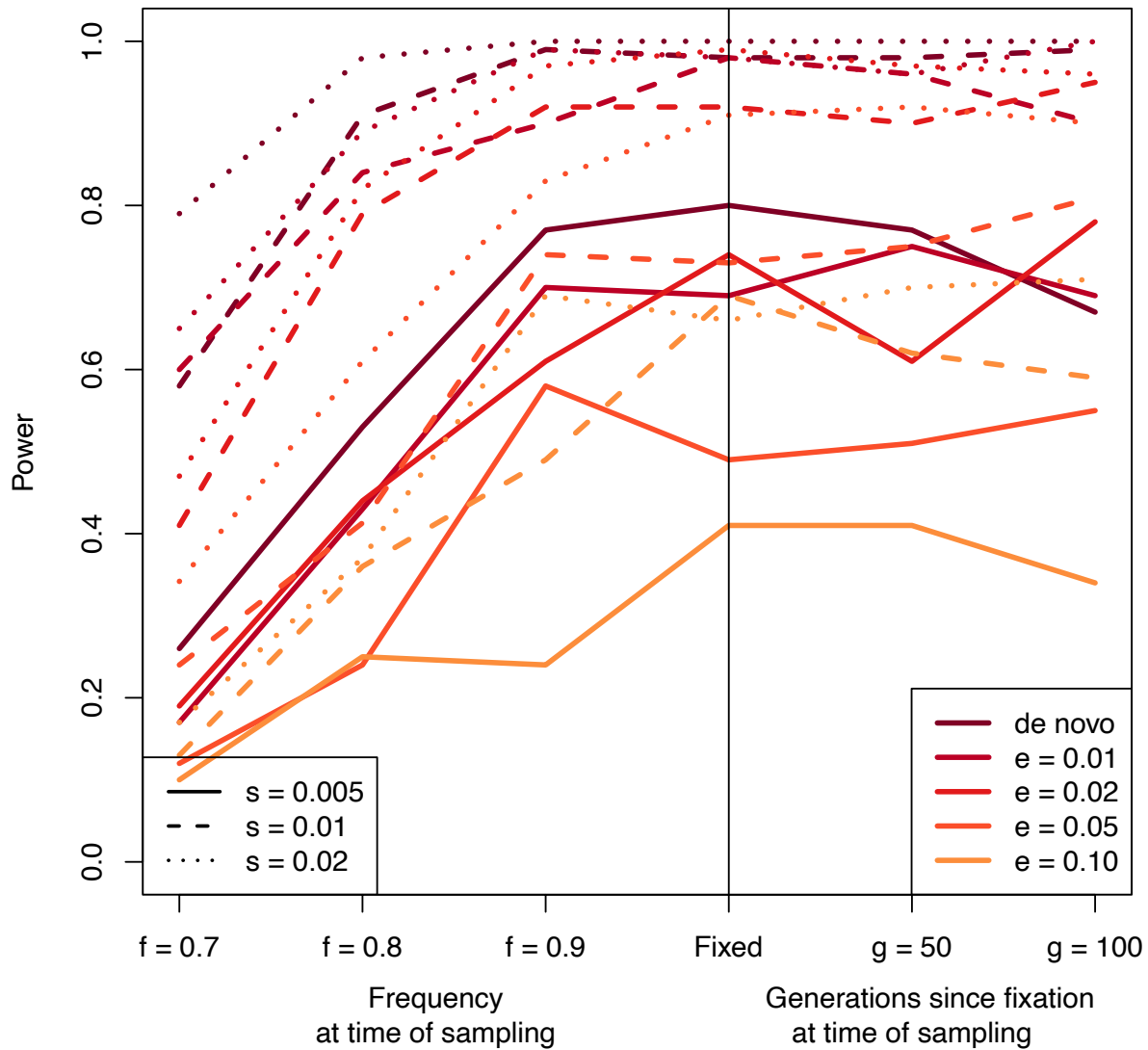**Figure S19**. Demo 3 XP-nSL $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

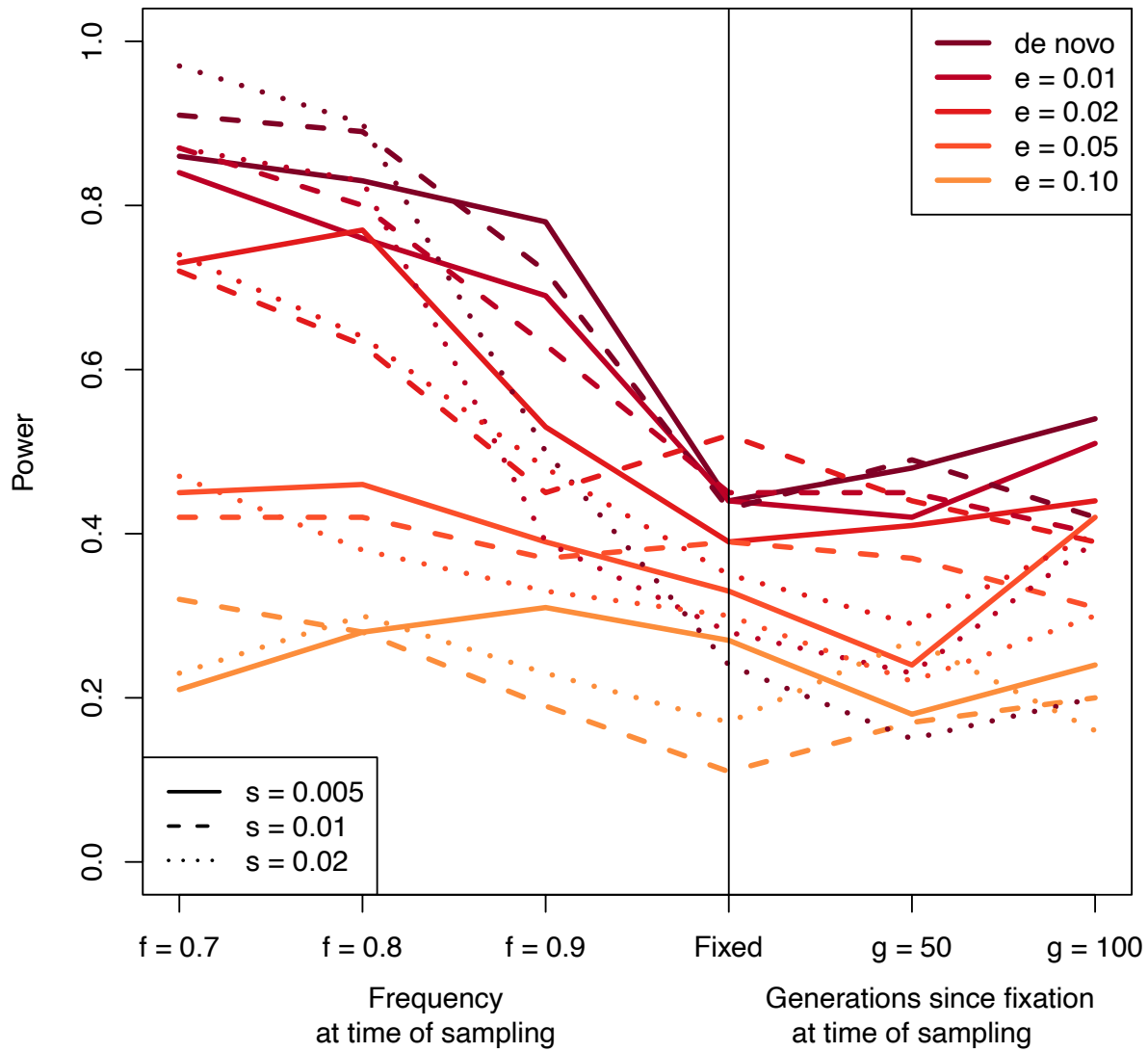**Figure S20**. Demo 3 XP-nSL $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

**Figure S21**. Demo 4 iHS $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the

frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of

sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in

generations since the two populations diverged.

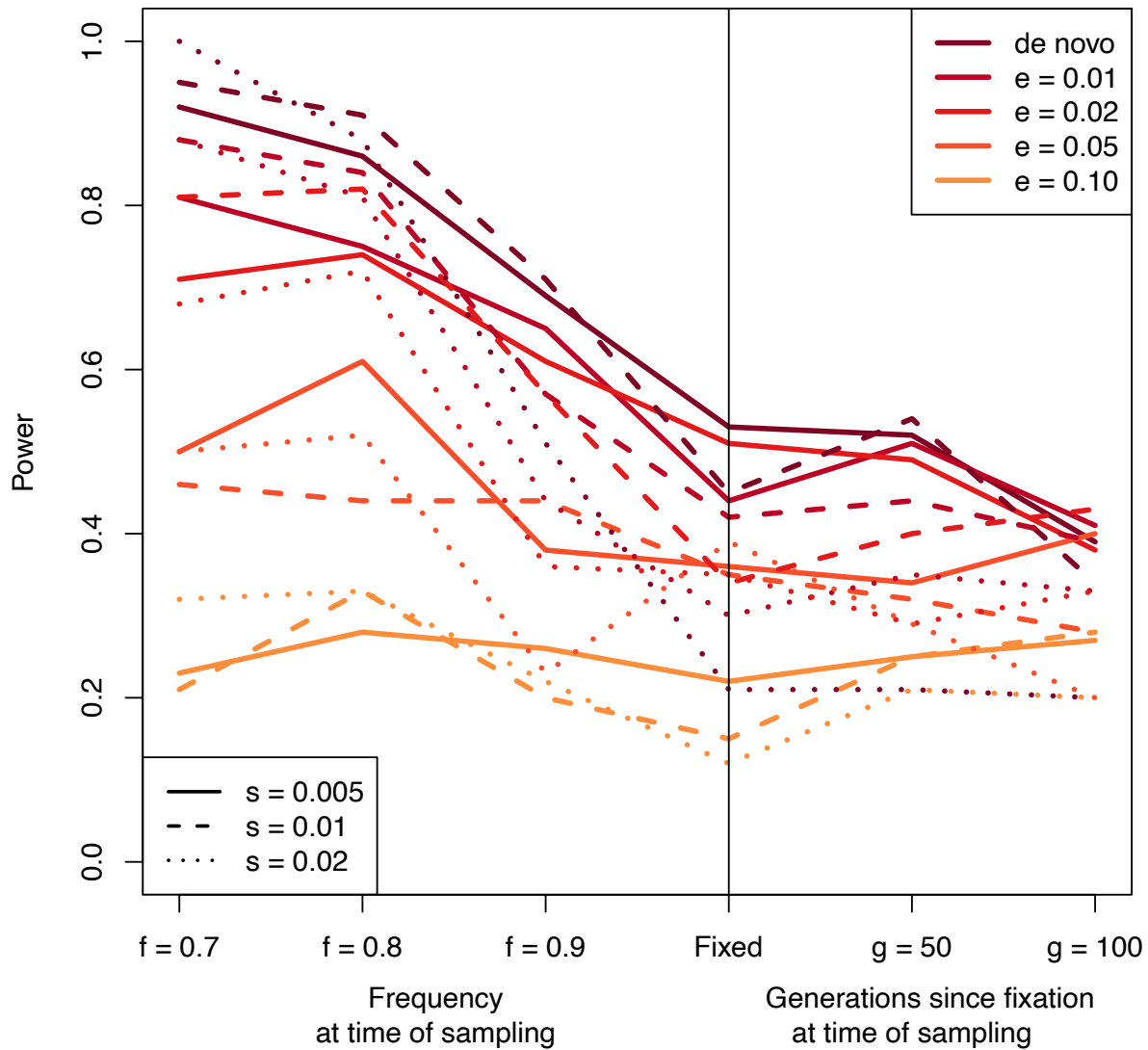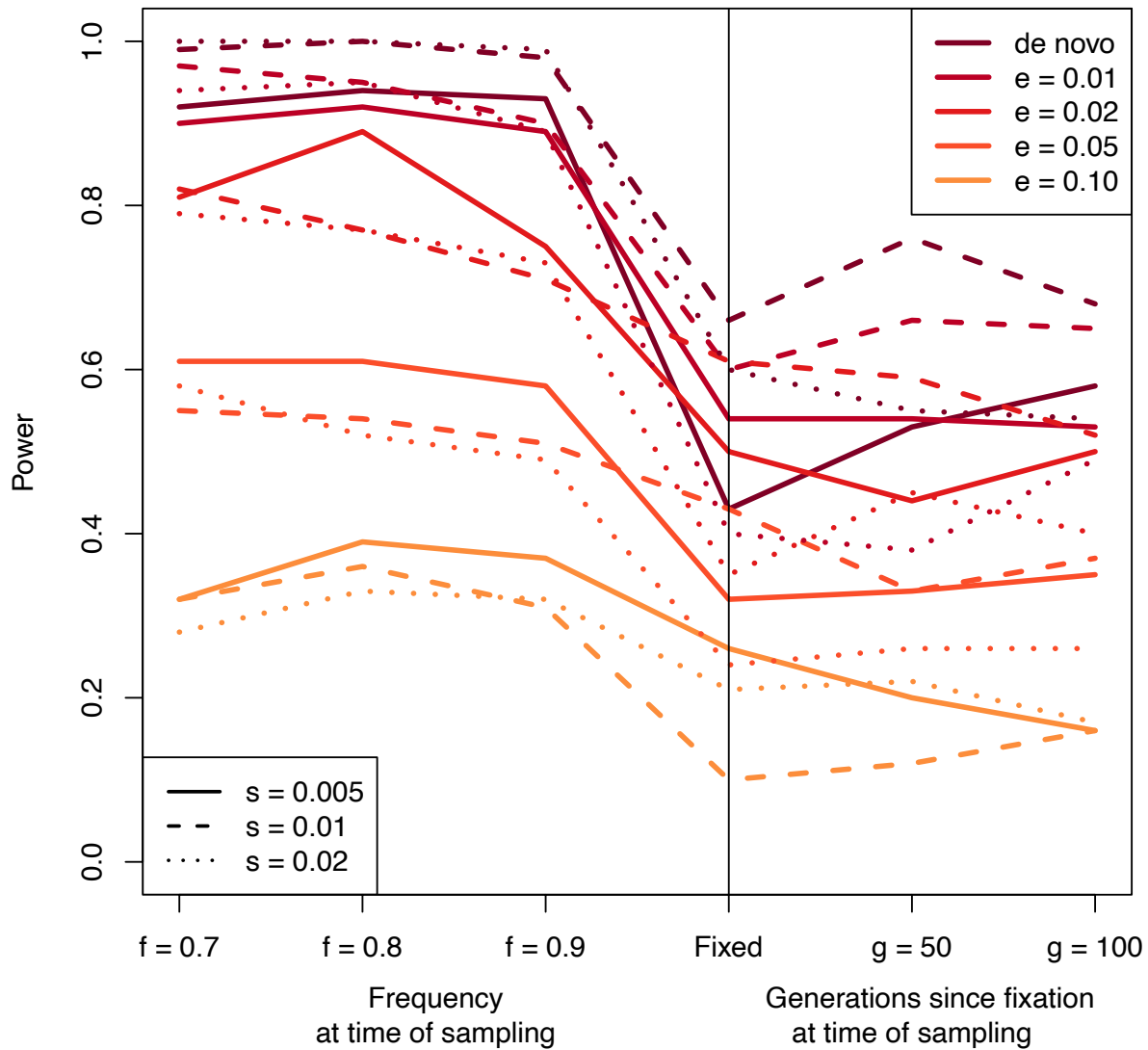**Figure S22**. Demo 4 iHS $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

**Figure S23**. Demo 4 nSL $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

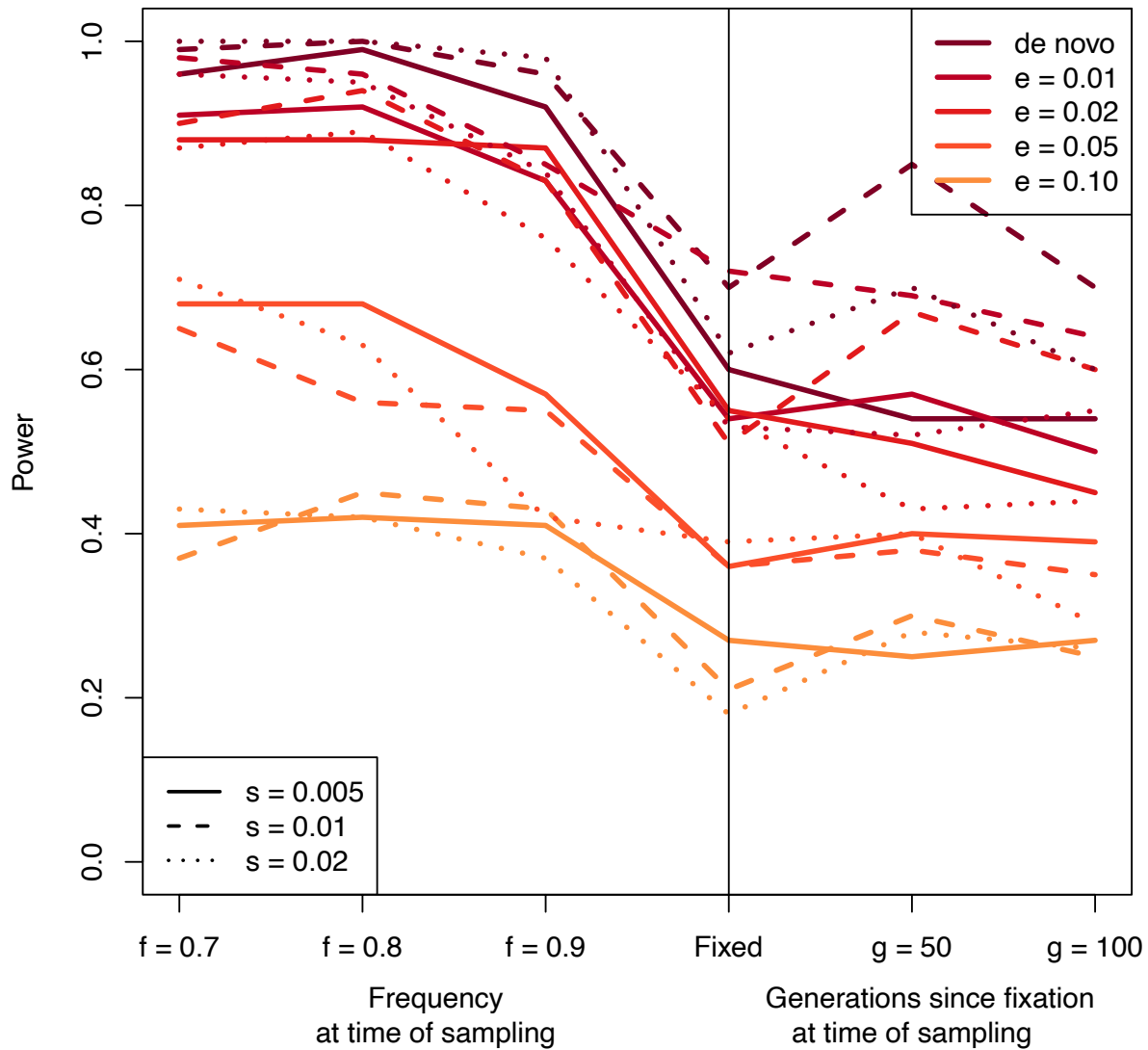**Figure S24**. Demo 4 nSL $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

XP−EHH; $t_d$ = 2000

**Figure S25**. Demo 4 XP-EHH $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
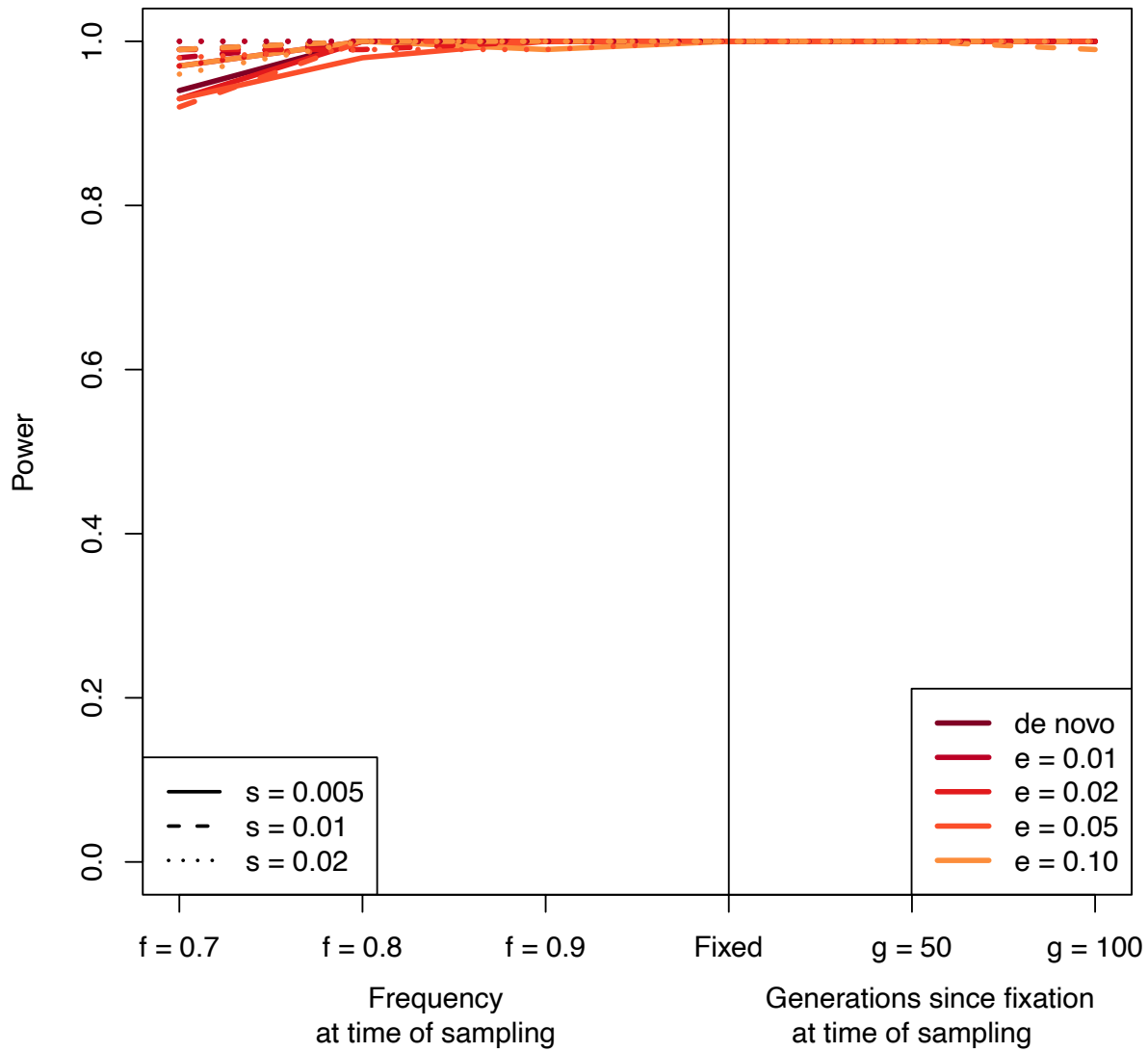
**Figure S26**. Demo 4 XP-EHH $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

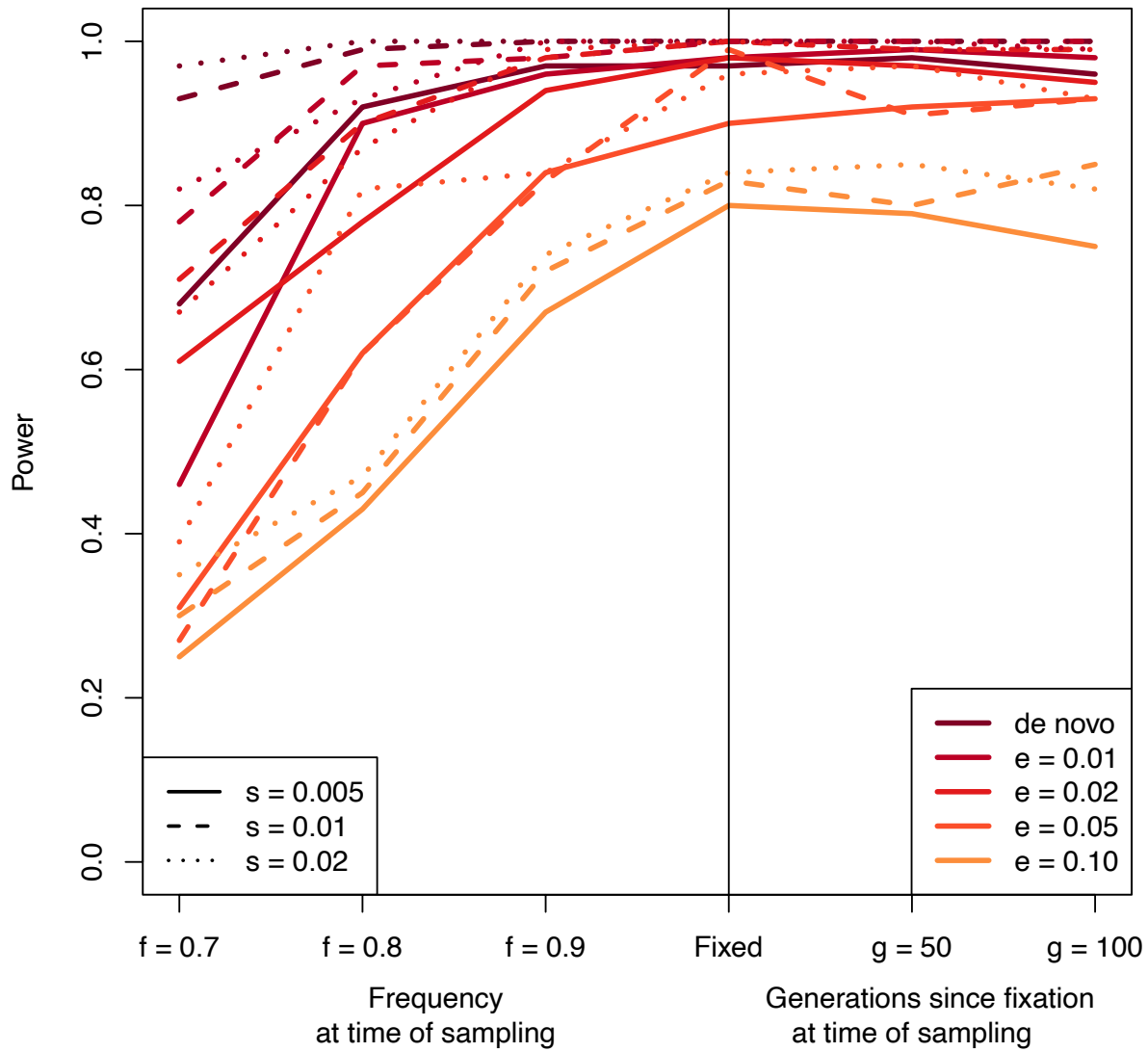XP−nSL; $t_d = 2000$



**Figure S27**. Demo 4 XP-nSL $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
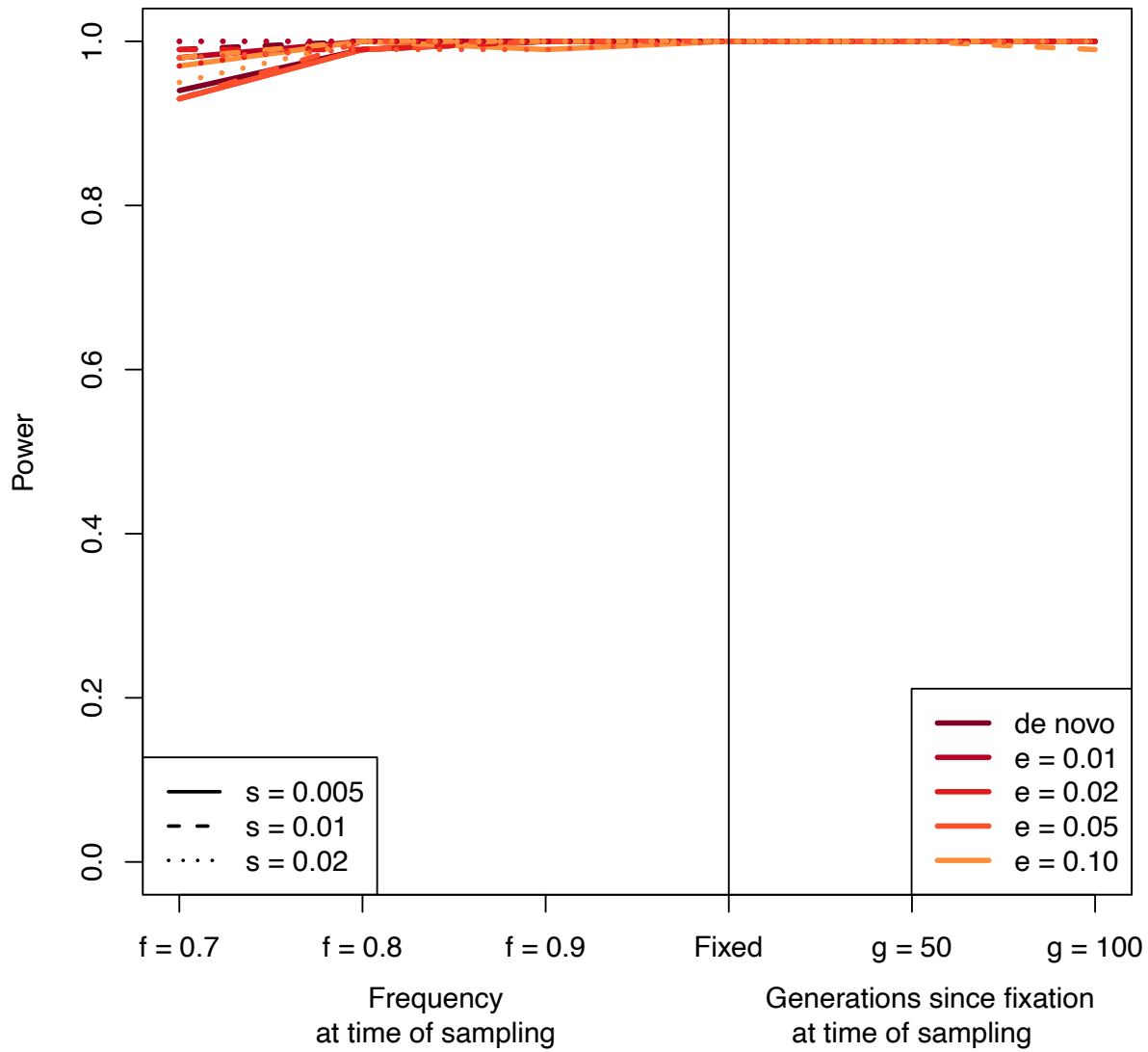
**Figure S28**. Demo 4 XP-nSL $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

**Figure S29**. Demo 5 XP-EHH $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.

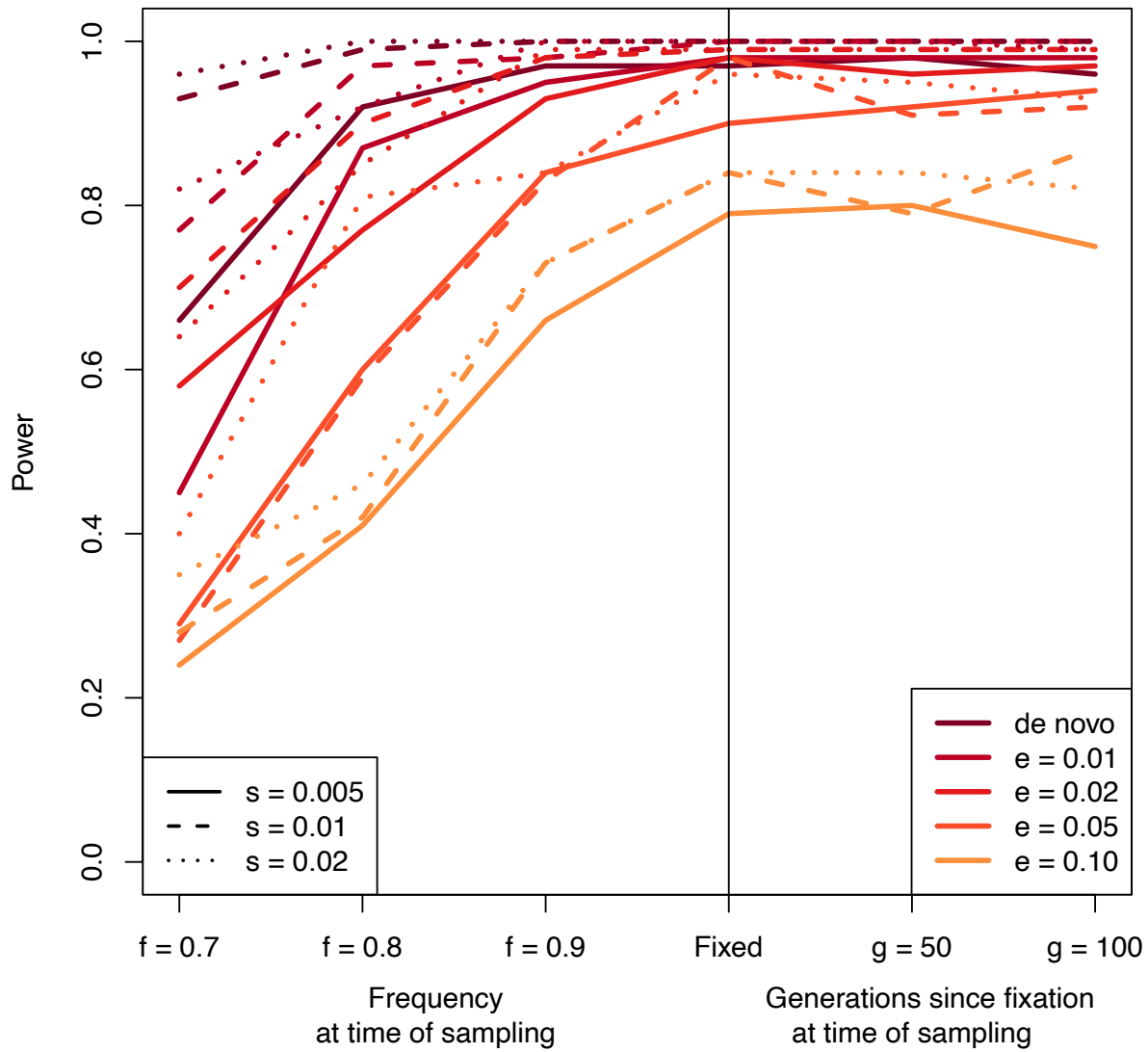XP–EHH; $t_d = 8000$

**Figure S30**. Demo 5 XP-EHH $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
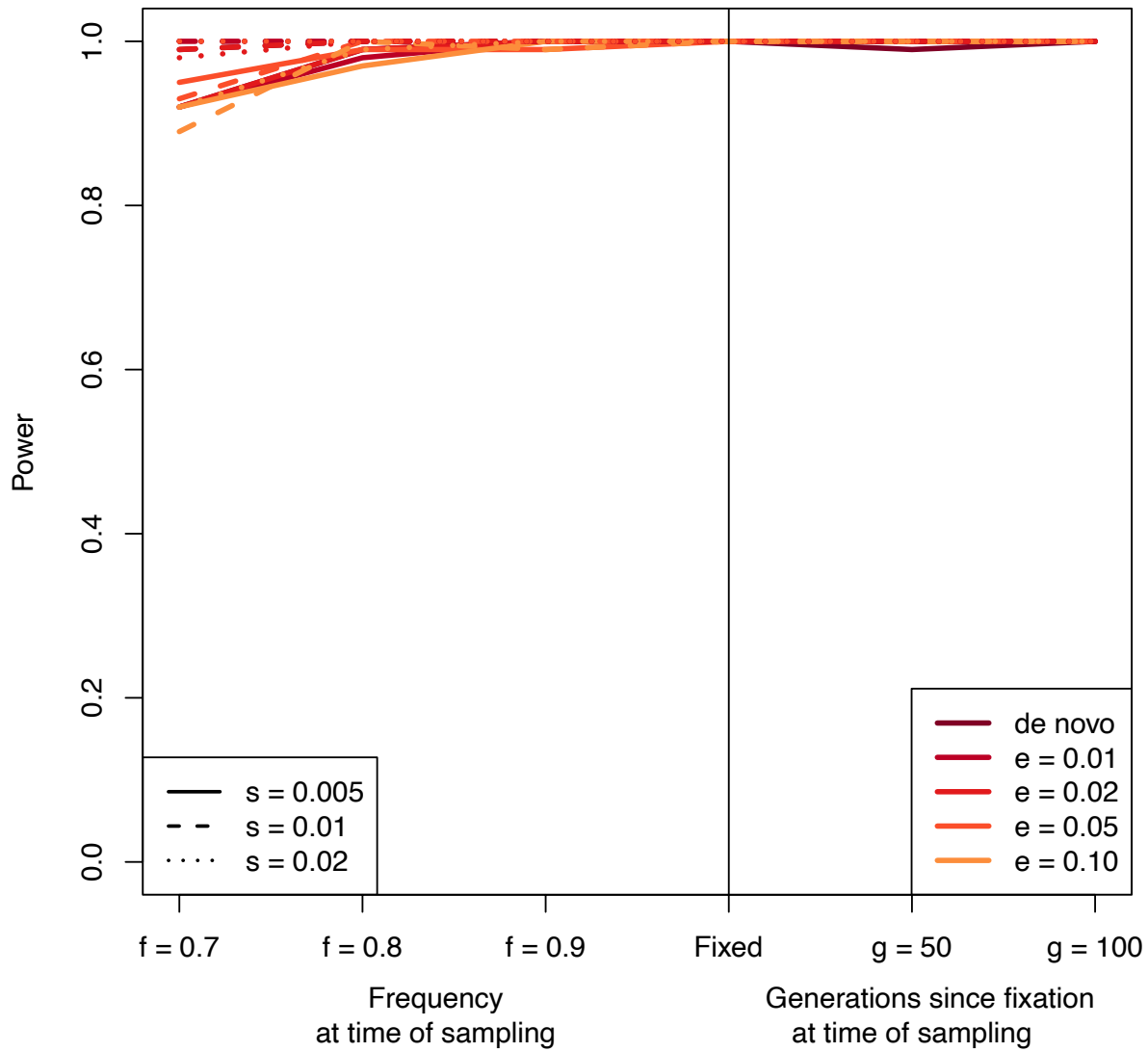
346

**Figure S31**. Demo 5 XP-nSL $t_d = 2000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
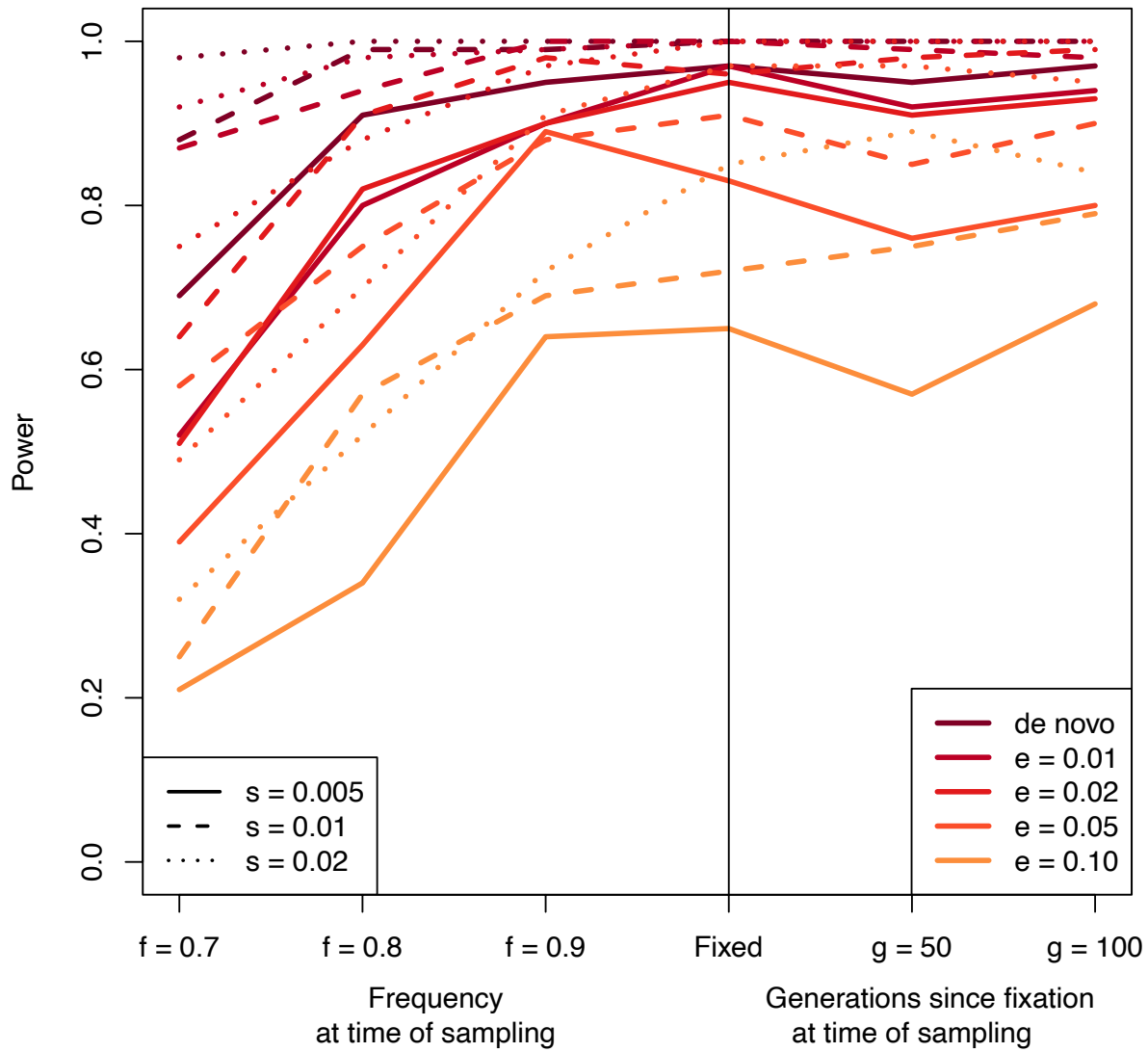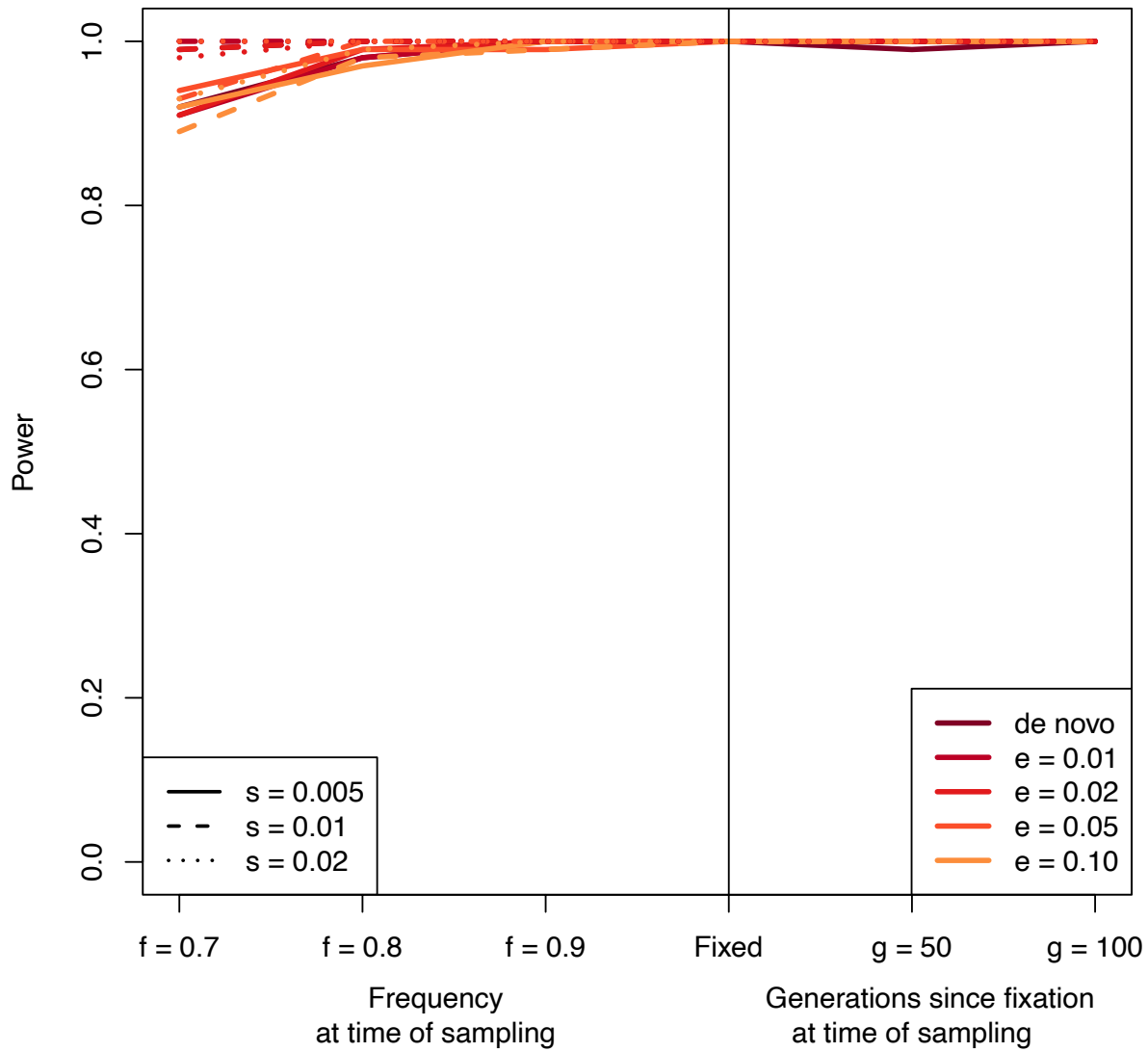
**Figure S32**. Demo 5 XP-nSL $t_d = 8000$ power curves. $s$ is the selection coefficient, $f$ is the frequency of the adaptive allele at time of sampling, $g$ is the number of generations at time of sampling since fixation, $e$ is the frequency at which selection began, and $t_d$ is the time in generations since the two populations diverged.
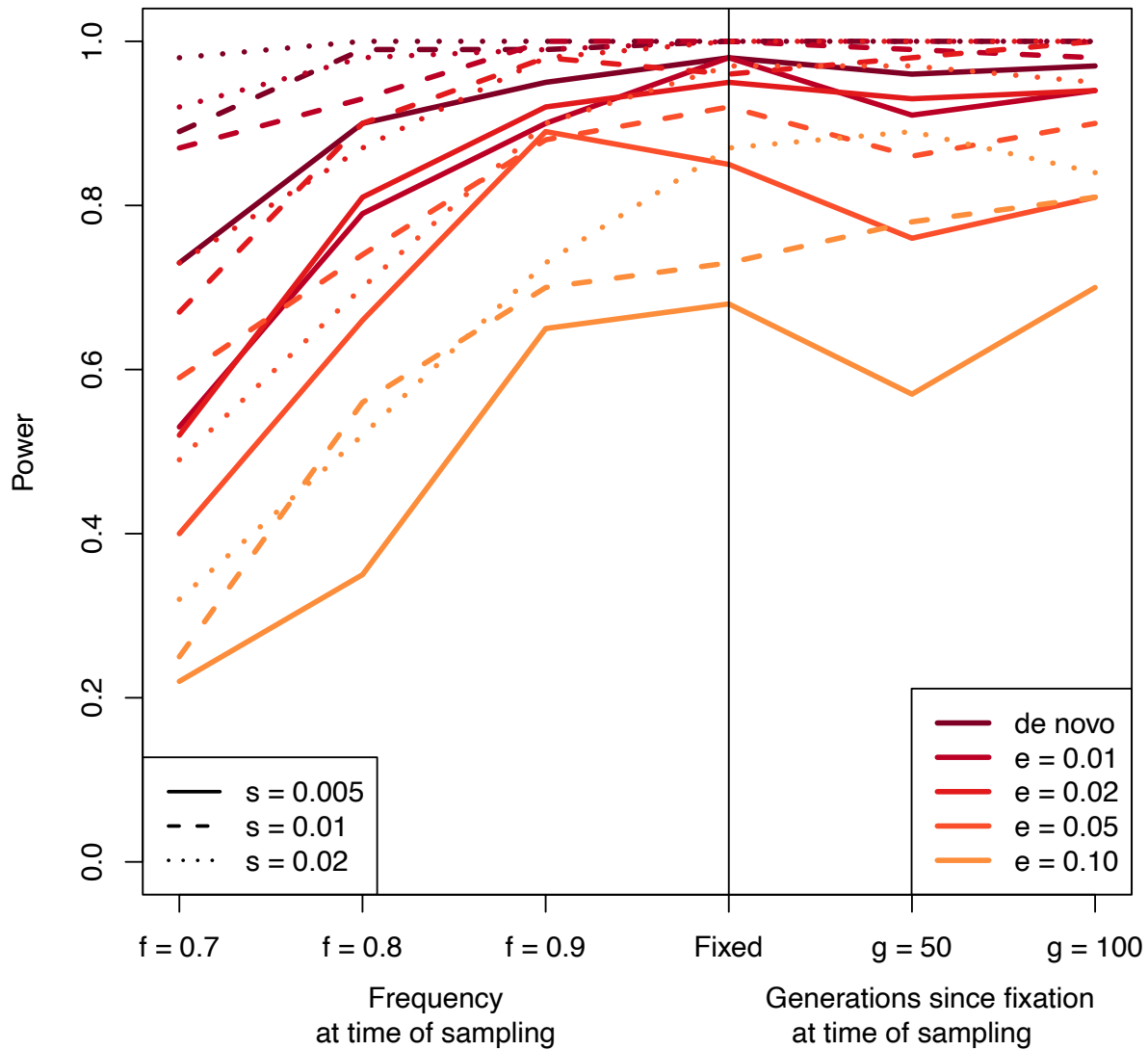
361 **Table 1**. Demographic history parameters for simulations.

| | $N_A$ | $N_0$ at split | $N_0$ at present | $N_1$ at split | $N_1$ at present | $t_d$ |
|---|---|---|---|---|---|---|
| Demo 1 | 10,000 | 10,000 | 10,000 | 10,000 | 10,000 | 2,000/4,000/8,000 |
| Demo 2 | 10,000 | 10,000 | 10,000 | 5,000 | 5,000 | 2,000/4,000/8,000 |
| Demo 3 | 10,000 | 5,000 | 5,000 | 10,000 | 10,000 | 2,000/4,000/8,000 |
| Demo 4 | 10,000 | 10,000 | 50,000$^\dagger$ | 10,000 | 10,000 | 2,000/4,000/8,000 |
| Demo 5 | 10,000 | 10,000 | 10,000 | 10,000 | 50,000$^\dagger$ | 2,000/4,000/8,000 |

362 $^\dagger$The reached via exponential growth starting 2,000 generations ago.
363
364 **Table S1**. False positive rate computed from neutral simulations for varying $t_d$ and demographic
365 history.

| | | $t_d = 2000$ | $t_d = 4000$ | $t_d = 8000$ |
|---|---|---|---|---|
| iHS | Demo 1 | 0.013 | 0.1 | 0.009 |
| | Demo 3 | 0.007 | 0.013 | 0.007 |
| | Demo 4 | 0.015 | 0.018 | 0.008 |
| nSL | Demo 1 | 0.01 | 0.015 | 0.008 |
| | Demo 3 | 0.008 | 0.011 | 0.007 |
| | Demo 4 | 0.014 | 0.021 | 0.014 |
| XP-EHH | Demo 1 | 0.013 | 0.013 | 0.016 |
| | Demo 2 | 0.017 | 0.009 | 0.015 |
| | Demo 3 | 0.01 | 0.011 | 0.012 |
| | Demo 4 | 0.012 | 0.014 | 0.014 |
| | Demo 5 | 0.011 | 0.012 | 0.013 |
| XP-nSL | Demo 1 | 0.014 | 0.011 | 0.013 |
| | Demo 2 | 0.019 | 0.011 | 0.012 |
| | Demo 3 | 0.011 | 0.011 | 0.012 |
| | Demo 4 | 0.012 | 0.012 | 0.014 |
| | Demo 5 | 0.011 | 0.012 | 0.014 |

366
367

368 **References**
369

370 Colonna V, Ayub Q, Chen Y, Pagani L, Luisi P, Pybus M, Garrison E, Xue Y, Tyler-Smith C,
371     Genomes Project C, et al. 2014. Human genomic regions with exceptionally high levels
372     of population differentiation identified from 911 whole-genome sequences. Genome
373     Biol 15:R88.
374 Crawford NG, Kelly DE, Hansen MEB, Beltrame MH, Fan S, Bowman SL, Jewett E, Ranciaro A,
375     Thompson S, Lo Y, et al. 2017. Loci associated with skin pigmentation identified in
376     African populations. Science 358.
377 DeGiorgio M, Szpiech ZA. 2021. A spatially aware likelihood test to detect sweeps from
378     haplotype distributions. bioRxiv:2021.2005.2012.443825.
379 Ferrer-Admetlla A, Liang M, Korneliussen T, Nielsen R. 2014. On detecting incomplete soft
380     or hard selective sweeps using haplotype structure. Mol Biol Evol 31:1275-1291.
381 Harris AM, DeGiorgio M. 2020. A likelihood approach for uncovering selective sweep
382     signatures from haplotype data. Mol Biol Evol.

383    Harris AM, Garud NR, DeGiorgio M. 2018. Detection and Classification of Hard and Soft
384        Sweeps from Unphased Genotypes by Multilocus Genotype Identity. Genetics 210:1429-
385        1452.
386    Kern AD, Schrider DR. 2016. Discoal: flexible coalescent simulations with selection.
387        Bioinformatics 32:3839-3841.
388    Lu K, Wei L, Li X, Wang Y, Wu J, Liu M, Zhang C, Chen Z, Xiao Z, Jian H, et al. 2019. Whole-
389        genome resequencing reveals Brassica napus origin and genetic loci involved in its
390        improvement. Nat Commun 10:1154.
391    Meier JI, Marques DA, Wagner CE, Excoffier L, Seehausen O. 2018. Genomics of Parallel
392        Ecological Speciation in Lake Victoria Cichlids. Mol Biol Evol 35:1489-1506.
393    Nedelec Y, Sanz J, Baharian G, Szpiech ZA, Pacis A, Dumaine A, Grenier JC, Freiman A, Sams
394        AJ, Hebert S, et al. 2016. Genetic Ancestry and Natural Selection Drive Population
395        Differences in Immune Responses to Pathogens. Cell 167:657-669 e621.
396    Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH, McCarroll
397        SA, Gaudet R, et al. 2007. Genome-wide detection and characterization of positive
398        selection in human populations. Nature 449:913-918.
399    Salmon P, Jacobs A, Ahren D, Biard C, Dingemanse NJ, Dominoni DM, Helm B, Lundberg M,
400        Senar JC, Sprau P, et al. 2021. Continent-wide genomic signatures of adaptation to
401        urbanisation in a songbird across Europe. Nat Commun 12:2983.
402    Szpiech ZA, Hernandez RD. 2014. selscan: an efficient multithreaded program to perform
403        EHH-based scans for positive selection. Mol Biol Evol 31:2824-2827.
404    Szpiech ZA, Novak TE, Bailey NP, Stevison LS. 2021. Application of a novel haplotype-based
405        scan for local adaptation to study high-altitude adaptation in rhesus macaques. Evol
406        Lett 5:408-421.
407    Voight BF, Kudaravalli S, Wen X, Pritchard JK. 2006. A map of recent positive selection in
408        the human genome. Plos Biology 4:e72.
409    Zhang SJ, Wang GD, Ma P, Zhang LL, Yin TT, Liu YH, Otecko NO, Wang M, Ma YP, Wang L, et
410        al. 2020. Genomic regions under selection in the feralization of the dingoes. Nat
411        Commun 11:671.
412    Zoledziewska M, Sidore C, Chiang CWK, Sanna S, Mulas A, Steri M, Busonero F, Marcus JH,
413        Marongiu M, Maschio A, et al. 2015. Height-reducing variants and selection for short
414        stature in Sardinia. Nat Genet 47:1352-1356.
415