

# Estimation of Pairwise Genetic Distances Under Independent Sampling of Segregating Sites vs. Haplotype Sampling

Max Shpak<sup>1,2,3</sup>, Yang Ni<sup>4</sup>, Jie Lu<sup>5</sup> and Peter Mueller<sup>4,6</sup>

<sup>1</sup> Sarah Cannon Research Institute, Austin TX 78705, USA

<sup>2</sup> Center for Systems and Synthetic Biology, University of Texas, Austin TX 78712, USA

<sup>3</sup> Fresh Pond Research Institute, Cambridge MA 02140, USA

<sup>4</sup> Department of Statistics and Data Science, University of Texas, Austin TX 78712, USA

<sup>5</sup> Genetics Division, Fisher Scientific, Austin TX 78744, USA

<sup>6</sup> Department of Mathematics, University of Texas, Austin TX 78712, USA

February 14, 2017

## Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>The sampling models</b>	<b>8</b>
2.1	Case 1: Independent Locus Sampling (ILS) . . . . .	12
2.2	Case 2: Whole Haplotype Sampling (WHS) . . . . .	13
2.3	Difference and independence . . . . .	15
2.4	Implications . . . . .	17
<b>3</b>	<b>Comparison to individual-based simulations</b>	<b>20</b>
<b>4</b>	<b>Analysis of cancer sequence data</b>	<b>23</b>
<b>5</b>	<b>Discussion</b>	<b>27</b>
<b>6</b>	<b>Acknowledgments</b>	<b>30</b>

<b>7</b>	<b>Statement of Effort</b>	<b>30</b>
<b>8</b>	<b>Figures and Tables</b>	<b>31</b>
<b>9</b>	<b>Appendix A1: Ordered Pairs of Pairs</b>	<b>36</b>

*Running Title* : Sample variance of genetic distance

*keywords* : genetic distance, linkage disequilibrium, Tajima's estimator, cancer genomics, next-generation sequencing

*corresponding author* : Max Shpak, St. David's Medical Center,  
1015 E. 32nd St, Suite 414, Austin TX 78705. Ph: 512-544-8077,  
Email: shpak.max@gmail.com

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

## Abstract

Genetic distance is a standard measure of variation in populations. When sequencing genomes individually, genetic distances are computed over all pairs of multilocus haplotypes in a sample. However, when next-generation sequencing methods obtain reads from heterogeneous assemblages of genomes (e.g. for microbial samples in a biofilm or cells from a tumor), individual reads are often drawn from different genomes. This means that pairwise genetic distances are calculated across independently sampled sites rather than across haplotype pairs. In this paper, we show that while the expected pairwise distance under whole haplotype sampling (WHS) is the same as with independent locus sampling (ILS), the sample variances of pairwise distance differ and depend on the direction and magnitude of linkage disequilibrium (LD) among polymorphic sites. We derive a weighted LD value that, when positive, predicts higher sample variance in estimated genetic distance for WHS. Weighted LD is positive when on average, the most common alleles at two loci are in positive LD. Using individual-based simulations of an infinite sites model under Fisher-Wright genetic drift, variances of estimated genetic distance are found to be almost always higher under WHS than under ILS, suggesting a reduction in estimation error when sites are sampled independently. We apply these results to haplotype frequencies from a lung cancer tumor to compute weighted LD and the variances in estimated genetic distance under ILS vs. WHS, and find that the the relative magnitudes of variances under WHS vs. ILS are sensitive to sampled allele frequencies.

## 26 **1 Introduction**

27 Genetic variation is the raw material for evolutionary change, consequently,  
28 one of the defining empirical questions in evolutionary and population genet-  
29 ics is the measurement of genetic heterogeneity in natural and experimental  
30 populations (Lewontin et al., 1974; Ellegren and Galtier, 2016). In addition  
31 to its importance to furthering our basic understanding of the evolutionary  
32 process (Hansson and Westerberg, 2002), characterization of genetic varia-  
33 tion has applied significance in many endeavors relevant to human welfare,  
34 including biomedical research. For example, the extent of genetic variation  
35 in populations of pathogens can be predictive of their ability to adapt to an-  
36 tibiotic treatment (Martinez and Baquero, 2000; MacLean et al., 2010) and  
37 immune response, while genetic heterogeneity among populations of cancer  
38 cells is predictive of their potential for metastatic disease and of a tumor's  
39 ability to develop resistance to chemotherapies (Dexter and Leith, 1986;  
40 Burrell et al., 2013; Sun and Yu, 2015). Estimates of genetic variation are  
41 equally relevant to maintaining diversity in crop and livestock strains (Na-  
42 tional Research Council, 1993; Fu, 2015) and to the maintenance of viable  
43 populations in biological conservation (Van Dyke, 2008).

44 The recent advent of high-throughput technologies for DNA sequencing  
45 allows researchers to measure genetic variation within and among popula-  
46 tions with very large sample sizes and high statistical power. The methods  
47 developed for characterizing genetic variation in studies of multicellular, usu-  
48 ally sexually reproducing model organisms can now be applied to genomic  
49 studies of typically clonal unicellular organisms such as microbes growing on  
50 biofilms, to populations of genetically heterogeneous cancer cells in a tumor,  
51 or to viruses in serums. In many cases, both the underlying population ge-  
52 netic models and the descriptive statistics used to measure genetic variation

53 in microbial or tumor samples must be adjusted to take into consideration  
54 the biological characteristics of the populations under study (such as the  
55 absence of meiotic recombination), as well as for the statistical properties of  
56 what are often different methods of sampling.

57 One of the most widely used measures of genetic variation in a population  
58 is the mean pairwise genetic distance among genomes, which is an estimate  
59 of the total heterozygosity across all polymorphic sites. For a sample of  $n$   
60 genotypes, the mean pairwise distance is calculated as:

$$\hat{\pi}_1 = 2 \sum_{i,j} \pi_{ij} / n(n-1), \quad (1)$$

61 where  $\pi_{i,j}$  is the Hamming distance for the haplotype pair  $z_i, z_j$ , summed  
62 over all polymorphic sites in haplotypes  $i$  and  $j$ , i.e.  $\pi_{ij} = \sum_s f(z_{is}, z_{js})$  for  
63  $f(z_{is}, z_{js}) = 1$  if site  $s$  has different nucleotides in haplotypes  $z_i, z_j$  and 0  
64 otherwise.

65 The parameter  $\hat{\pi}$  is of importance not only as a summary statistic of  
66 genetic variation, but as an estimator of key population genetic parameters.  
67 Under neutral evolution in an infinite sites model (Kimura, 1969; Tajima,  
68 1996),  $\hat{\pi}$  estimates the population mutation rate (Tajima, 1989), i.e. for a  
69 diploid population with  $N$  individuals and a per-generation genomic muta-  
70 tion rate  $u$ ,

$$E[\hat{\pi}] = 4Nu = \theta_\pi,$$

71 where  $\theta_\pi$  represents the distance-based Tajima estimator for this parameter  
72 (in a population of  $N$  haploids,  $\theta = 2Nu$ ). As a result,  $\hat{\pi}$  provides an esti-  
73 mate of neutral effective population size in populations when the mutation  
74 rate  $u$  is known approximately. Additionally, comparisons of  $\theta_\pi$  estimated

75 from genetic distances to the population mutation rate estimated from the  
76 number of segregating sites  $S_n$  in a sample of  $n$  genotypes

$$S_n / \left( \sum_{i=1}^{n-1} 1/i \right) = \theta_S,$$

77 (Watterson, 1975) is the basis for the Tajima D test for selection. Values of  
78  $\theta_\pi$  that are inflated relative to  $\theta_S$  may be the result of diversifying selection  
79 or recent population bottlenecks, while values of  $\hat{\pi}$  that are smaller than  
80 the value expected from the number of polymorphic sites indicate a history  
81 of selective sweeps or, alternatively, a recent population expansion with a  
82 relative large number of recent, rare variants. Understanding the error in  
83 estimates of  $\hat{\pi}$  relates directly to the error inherent to estimates of population  
84 mutation rate  $\theta$  and tests for neutral evolution derived from this parameter.

85 In studies of multicellular organisms,  $\hat{\pi}$  is estimated directly from com-  
86 plete haplotypes sampled from  $n$  different individuals, each of which has  
87 been sequenced over the region(s) containing the segregating sites of inter-  
88 est, i.e.  $\hat{\pi}$  is computed from the Hamming distances among actual haplotype  
89 pairs. This pairwise comparison of genotypes across multiple loci is possible  
90 because the co-occurring genotypes across sites in the genome are known for  
91 individually sequenced genomes. In contrast, for most samples of microbes  
92 or of cancer cells, the application of next generation sequencing (NGS) meth-  
93 ods (Goodwin et al., 2016) entail sampling reads from an unknown number  
94 of different genomes (in contrast to multicellular tissue samples from indi-  
95 vidual organisms, which are assumed to be genetically homogeneous). In  
96 the limiting case, if the read coverage depth at each segregating site is suf-  
97 ficiently small relative to the number of individual genomes in a sample,  
98 every read is likely to be drawn from a different individual cell and geno-

99 type (assuming non-adjacent segregating sites that occur on separate reads).  
100 Consequently, sampling in this way for a read depth of  $n$  is not statistically  
101 equivalent to sequencing  $n$  individuals at the same number of sites. The  
102 estimated mean pairwise genetic distance for independent sampling of loci  
103 from different genomes is:

$$\hat{\pi}_2 = 2 \sum_s \sum_{i_s, j_s} f(z_{i_s, s}, z_{j_s, s}) / n(n-1), \quad (2)$$

104 where  $z_{i_s, s}$  is the identity of the  $i$ th allele sampled at locus  $s$ , which is as-  
105 sumed to be from a different genome (distinct cell or organism) with respect  
106 to the  $i$ th sample at some other site (in contrast to  $z_{i_s}$  in Eqn. (1), which  
107 represents site  $s$  in haplotype  $i$ ). We include the subindex  $s$  in  $i_s, j_s$  to high-  
108 light this. As in Eqn. (1),  $f(x, y)$  is an indicator function equal to 1 if the  
109 nucleotide pair is not identical and 0 otherwise.

110 Throughout this paper, we will refer to these two modes of genotype  
111 sampling as as whole haplotype sampling (WHS) and as independent locus  
112 sampling (ILS), respectively. The difference between WHS and ILS is illus-  
113 trated schematically in Figure 1.

114

115 FIGURE 1 HERE

116

117 Although WHS is usually used to estimate genetic distances when se-  
118 quencing multicellular organisms while ILS is standard for assemblages of  
119 microbes or tumor cells, one can apply single cell sequencing (corresponding  
120 to WHS) to microbe and tumor cells (Navin, 2015; Gawad et al., 2016). It  
121 is also possible to sample loci via independent reads from different genomes  
122 (ILS) in multicellular organisms if one sequences sufficiently many individual

123 organisms (i.e. more individuals genotyped than there are reads), although  
124 it usually isn't practical to do so. Therefore, it is instructive to compare  
125 the sampling distributions of  $\hat{\pi}$  obtained for WHS and ILS. Even in cases  
126 where only either ILS or WHS are practically feasible, it is important to  
127 understand the potential sources of error in estimates of genetic distance  
128 given the type of sampling being used. The sample variances of  $\hat{\pi}$  under ILS  
129 vs. WHS are of particular significance, as they determine the expected error  
130 in our point estimates of genetic distance in a population, and by extension,  
131 the reliability of test statistics for the consequences of natural selection or  
132 population dynamics such as Tajima's D.

133 Below, we will derive the expectations and sample variances of pairwise  
134 genetic distance under the two modes of sampling with the fewest possi-  
135 ble a priori assumptions about the number and distribution of mutations.  
136 We test these analytical predictions against samples from simulated popula-  
137 tions undergoing neutral evolution via random mutation and Fisher-Wright  
138 genetic drift under an infinite-sites model. We also apply these results to  
139 estimating the variances in genetic distances using single nucleotide variant  
140 (SNV) frequency data from lung cancer tumors.

## 141 **2 The sampling models**

142 Consider a population of  $N$  organisms with some distribution of mutations  
143 over  $S$  segregating sites (in the population, as opposed to  $S_n \ll S$  in a  
144 sample of  $n$ ). We wish to estimate the mean genetic distance  $\hat{\pi}$  for the  
145 population and its sample variance  $var(\hat{\pi})$  under the WHS and ILS models  
146 of sampling. For WHS, we draw  $n \ll N$  individual organisms (or cells) from  
147 the population and sequence their entire genomes, exomes, or any regions  
148 containing the polymorphic sites of interest. For simplification but without

149 loss of generality, assume that the sample consists of  $n$  haploid genotypes  
150 or known/phased haplotypes, regardless of how they were sequenced or the  
151 number of reads (we note that if we were working with diploid genotypes,  
152 phasing would not matter if pairwise distances are computed with respect  
153 to the per-site genotype).

154 For an idealized model of ILS in an aggregate sample of microbes or  
155 cells, we assume that the number of individual genomes (i.e. from different  
156 tumor cells or microbes) that contribute reads to a sample is much larger  
157 than the sequencing read depth (mean coverage depth)  $n$ . If this is the  
158 case, we can assume (approximately) that the majority of reads are sampled  
159 different individual genomes. If we make the further assumption that reads  
160 are short, the majority of reads will contain at most a single polymorphic  
161 site. Together, these conditions imply that the majority of polymorphic sites  
162 will be sampled from different genomes, or, more precisely, each polymorphic  
163 site is sampled independently of other polymorphic sites with respect to  
164 their genome of origin (in the second panel of Figure 1, several sites are  
165 sampled from the same genome simply because there are very few genomes to  
166 draw this random sample from). When computing average pairwise genetic  
167 distance, WHS sums over the Hamming distances of all haplotype pairs,  
168 while ILS is the sum over all pairs for each of the  $S_n$  segregating sites  
169 sampled from different individuals.

170 Without loss of generality, we also assume an infinite sites model so that  
171 there only two alleles per segregating site. This allows an unambiguous  
172 binary classification of alleles, with mutations as ancestral "wildtype" vs.  
173 "reference" genotype (in the case of tumors, the reference corresponds to  
174 the normal germline genotype, with somatic mutations defining the variant  
175 genotypes of the clonal lineages), and to specify the direction of linkage dis-

176 equilibrium. We note, however, that the results derived below are applicable  
177 to multiallelic states provided that some allele (usually the most common,  
178 or, in the case of cancer genomics, the germline allele) is designated as a  
179 reference and all other alleles are pooled together to create an aggregate  
180 biallelic state.

181

182 **Definitions.** In this subsection and throughout the manuscript, we will  
183 make use of the following definitions and terminology as a formal way of  
184 characterizing and distinguishing between Eqns (1) and (2) in the introduc-  
185 tion:

186

187 *Variables:* Let  $z$  denote a genotype, at either single locus  $s$  or across multiple  
188 loci. We define the frequency distribution of  $z$  over samples  $i$  as  $z_i \sim p(z)$ ,  
189 which are iid among  $i = 1..n$ . As above, we use  $z_{is}$  to denote site  $s$  in  
190 haplotype  $i$  (for WHS), and  $z_{i,s}$  to denote sample  $i$  at site  $s$  when sites are  
191 sampled independently (ILS).

192

193 *Pairs:* In both cases, that is, for WHS and ILS, respectively, the esti-  
194 mators  $\hat{\pi}_1$  and  $\hat{\pi}_2$  include an average  $\sum_{i<j} \phi_{ij}/n(n-1)$  of some function  
195  $\phi_{ij} = \phi(x_i, x_j)$  of pairs of i.i.d. random variables  $x_i, i = 1, \dots, n$ . In the  
196 case of ILS  $x_i = z_{is}$  and  $\phi(x_i, x_j) = f_{ijs}$  with  $f_{ijs} = I(z_{is} \neq z_{js})$  (and an  
197 additional sum over  $s$ , outside the average). In the case of WHS the random  
198 variables are  $x_i = z_i$  and  $\phi(x_i, x_j) = g_{ij} = \sum_s f_{ijs}$ . Importantly, while the  
199 r.v.'s  $x_i$  are independent, pairs  $(x_i, x_j)$  and  $(x_i, x_k)$  that share a common  
200 element are not.

201

202 *Moments of  $\phi_{ij}$ :* We define  $E(\phi_{ij}) = \mu$ ,  $var(\phi_{ij}) = \sigma^2$ . We also define an

203 expectation for an indicator function on pairs of pairs with a shared element  
204 as  $E(\phi_{ij}, \phi_{jk}) = \kappa$ .

205

206 *Pairs of pairs*: Let  $P$  denote the set of all ordered pairs of pairs, with  $P_3 \subset P$   
207 defining the subset of ordered pairs of pairs with a single shared element,

$$P = \{[(i, j), (k, \ell)] : i < j, k < \ell \text{ and } (i, j) < (k, \ell)\}$$

$$P_3 = \{[(i, j), (k, \ell)] : i < j, k < \ell \text{ and } (i, j) < (k, \ell) \text{ and } |\{i, j, k, \ell\}| = 3\}$$

208 Numbers of pairs: The number of ordered pairs, and the number of ordered  
209 pairs of pairs with a shared element are, respectively

$$N_2 = n(n-1)/2$$

$$N_3 = n(n-1)(n-2)/2$$

210 The value of  $N_3$  follows from the fact that there are  $n(n-1)(n-2)/6$  ways  
211 to select a triplet  $i, j, k$ , and three ways to select a shared element from this  
212 triplet. In Appendix A1, we cover some of the properties of ordered pairs  
213 of pairs, including the derivation of the following relation which we will use  
214 below to compute  $\text{var}(\hat{\pi})$  under ILS and WHS,

$$\text{var}(\hat{\phi}_n) = \frac{\sigma^2}{N_2} + 2\frac{N_3}{N_2^2}(\kappa - \mu^2), \quad (3)$$

215 where  $\hat{\phi}_n = \frac{1}{N_2} \sum_{i < j} \phi_{ij}$  is a sample estimate of  $E(\phi_{ij}) = \mu$ . We will use this  
216 result twice, once for ILS with  $\phi_{ij} = f_{ijs}$ , and once for WHS with  $\phi_{ij} = g_{ij}$ .

## 217 **2.1 Case 1: Independent Locus Sampling (ILS)**

218 For ILS, we use the indicator function at a single site  $s$ ,  $f_{ij,s} = I(z_{i_s,s} \neq z_{j_s,s})$ ,

219 where  $z_{i_s,s} \sim \text{Bern}(p_s)$ , i.e.  $p(z_{i_s,s}) = p_s$  for  $z_{i_s,s} \in 0, 1$  such that

$$\mu_s = E(f_{ij,s}) = h_s = 2p_s(1 - p_s)$$

$$\sigma_s^2 = \text{var}(f_{ij,s}) = h_s(1 - h_s)$$

220 (note that  $h_s$  is the heterozygosity at locus  $s$ ).

221 The expectation of the indicator function for ordered pairs on pairs in-

222 cludes a covariance term, namely,

$$\begin{aligned} \kappa_s &= E(f_{ij,s} f_{jk,s}) = p(z_{i_s,s} \neq z_{j_s,s}, z_{j_s,s} \neq z_{k_s,s}) = p(z_{i_s,s} = z_{k_s,s} \neq z_{j_s,s}) \\ &= p(z_{i_s,s} = z_{k_s,s} = 1, z_{j_s,s} = 0) + p(z_{i_s,s} = z_{k_s,s} = 0, z_{j_s,s} = 1) \\ &= p_s^2(1 - p_s) + (1 - p_s)^2 p_s = h_s/2. \end{aligned}$$

The sampling estimator for  $\hat{\pi}$  under ILS is given by

$$\hat{\pi}_{ILS} = \sum_s \left\{ \frac{1}{N_2} \sum_{i < j} I(z_{i_s,s} \neq z_{j_s,s}) \right\} = \sum_s \left\{ \frac{1}{N_2} \sum_{i < j} f_{ij,s} \right\}.$$

223 From the assumption of statistical independence among sites  $s$  located on

224 different reads under ILS, it follows (Appendix A1) that for a sample of  $n$ ,

$$\text{var}(\hat{\pi}_{ILS}) = \sum_s \text{var}(\hat{f}_{n,s}) = \sum_s \frac{1}{N_2} h_s \left\{ (1 - h_s) + \frac{N_3}{N_2} (1 - 2h_s) \right\} \quad (4)$$

225 We remark that in practice, the assumption of independence requires that

226 the number of possible samples of size  $n$  is much larger than the number of

227 segregating sites (i.e.  $N \gg n$  so that  $\binom{N}{n} \gg S_N$ ).

228 **2.2 Case 2: Whole Haplotype Sampling (WHS)**

229 Computing pairwise differences for independent samples of  $z = z_i$  under  
 230 WHS involves computing moments of sums rather than sums of moments,  
 231 i.e.

$$g_{ij} = \sum_s I(z_{is} \neq z_{js}) = \sum_s f_{ij,s}$$

232 For samples of individual haplotypes  $i = 1 \dots n$ , consider  $z_i \sim p(z)$  with  
 233  $p(z_{is} = 1) = p_s$  as before, but with correlated  $z_{is}, z_{ir}$  due to linkage disequi-  
 234 librium (LD) between sites, i.e. for (arbitrarily labeled) alleles  $R, r$  and  $S, s$   
 235 at the two sites, and defining  $q_s, q_r = 1 - p_s, 1 - p_r$  (Lewontin and Kojima,  
 236 1960),

$$\begin{aligned} p(RS) &= p(R)p(S) + D_{sr} = p_r p_s + D_{sr} \\ p(rs) &= p(r)p(s) + D_{sr} = q_r q_s + D_{sr} \\ p(Rs) &= p(R)p(s) - D_{sr} = p_r q_s - D_{sr} \\ p(rS) &= p(r)p(S) - D_{sr} = q_r p_s - D_{sr}. \end{aligned}$$

237 As with ILS, we have, for  $h_s = 2p_s q_s$ ,

$$\mu_f = E(f_{ij,s}) = h_s \text{ and } \sigma_f^2 = \text{var}(f_{ij,s}) = h_s(1 - h_s)$$

238 With non-zero LD, the probability of different identity among sites  $s, r$  in  
 239 a sample pair  $i, j$  is  $p(f_{ijs} f_{ijr} = 1) = p(RS, rs) + p(rs, RS) + p(Rs, rS) +$   
 240  $p(rS, Rs)$ , where  $(RS, rs) = (z_{i,rs} = RS, z_{j,rs} = rs)$  etc. Therefore

$$\gamma_{sr} = E(f_{ij,s} \cdot f_{ij,r}) = 2(p_s p_r + D_{sr})(q_s q_r + D_{sr}) + 2(p_s q_r - D_{sr})(q_s p_r - D_{sr})$$

241 and similarly, considering triplet samples with shared element  $j$  paired with  
 242  $i$  and  $k$ , the probability of different identity between  $j$  and  $j$  at site  $s$  and  $j$   
 243 vs.  $k$  at site  $r$  is  $p(f_{ijs}f_{jkr} = 1) = p(R, rS, s) + p(R, rs, S) + p(r, RS, s) + \dots$   
 244 Using these terms, we compute the expectation:

$$\delta_{sr} = E(f_{ij,s} \cdot f_{jk,r}) = 2p_s(q_s p_r - D_{sr})(q_s q_r + D_{sr}) + 2(p_s q_r - D_{sr})(q_s p_r - D_{sr})$$

245 Assuming independence (linkage equilibrium,  $D_{sr} = 0$  for all  $s, r$ ) gives re-  
 246 sults equivalent to ILS, i.e. both equations simplify to  $\gamma_{sr} = \delta_{sr} = 4p_s q_s p_r q_r$ .  
 247 The mean and sample variance terms for the expected pairwise distances are,  
 248 respectively,

$$\begin{aligned} \mu &= E(g_{ij}) = \sum_s h_s, \\ \sigma^2 &= \text{var}(g_{ij}) = \sum_s \text{var}(f_{ij,s}) + 2 \sum_{r < s} \text{cov}(f_{ij,r}, f_{ij,s}) \\ &= \sum_s h_s(1 - h_s) + 2 \sum_{r < s} (\gamma_{sr} - h_s h_r), \end{aligned}$$

249 while the covariance  $\kappa$  for the ordered pair of pairs with a shared  $j$  element  
 250 is:

$$\begin{aligned} \kappa &= E(g_{ij} g_{jk}) = E \left\{ \sum_s f_{ij,s} \cdot \sum_s f_{jk,s} \right\} = \\ &= E \left\{ \sum_s I(z_{is} = z_{ks} \neq z_{js}) + 2 \sum_{r < s} (I(z_{is} \neq z_{js}) I(z_{jr} \neq z_{kr})) \right\} \\ &= \sum_s h_s/2 + 2 \sum_{r < s} \delta_{sr}. \end{aligned}$$

By incorporating  $\kappa$ , we can construct the sample estimate and variances for  $g_{ij}$ . For the WHS model,  $\mathbf{z}_i \sim p(\mathbf{z})$ , independently, from which we construct

the sample estimate for  $g_n$  as:

$$\hat{\pi}_{WHS} \equiv \hat{g}_n = \frac{1}{N_2} \sum_{i < j} g_{ij},$$

now averaging over haplotypes  $\mathbf{z}_i$  (rather than independent counts for each site).

Note that  $\hat{g}_n$  is again an average across pairs, like  $\hat{f}_n$  in the ILS case.

We again apply the result in Eqn. (3) to find

$$\begin{aligned} \text{var}(\hat{\pi}_{WHS}) &= \frac{\sigma^2}{N_2} + 2 \frac{N_3}{N_2^2} (\kappa - \mu^2) = \\ &= \frac{1}{N_2} \underbrace{\left( \sum_s h_s(1-h_s) + 2 \sum_{r < s} (\gamma_{sr} - h_s h_r) \right)}_{\sigma^2} + \frac{2N_3}{N_2^2} \left[ \underbrace{\sum_s h_s/2 + 2 \sum_{r < s} \delta_{sr}}_{\kappa} - \left( \sum_s h_s \right)^2 \right] \end{aligned} \quad (5)$$

### 251 **2.3 Difference and independence**

252 Using the results in Eqns. (4) and (5), we derive the difference between the  
 253 sample variances in pairwise differences under WHS vs. ILS as

$$\Delta = \text{var}(\hat{\pi}_{WHS}) - \text{var}(\hat{\pi}_{ILS}) = \frac{2}{N_2} \sum_{r < s} (\gamma_{sr} - h_s h_r) + \frac{4N_3}{N_2^2} \sum_{r < s} (\delta_{sr} - h_s h_r) \quad (6)$$

254 By collecting terms, we can rewrite the above as

$$\Delta = \frac{2}{N_2} \sum_{r < s} B_{sr} + \frac{4N_3}{N_2^2} \sum_{r < s} A_{sr},$$

255 where

$$\begin{aligned}
 A_{sr} &= \delta_{sr} - h_s h_r = (p_s p_r + q_s q_r - p_s q_r - p_r q_s) D_{sr} + 4p_s q_s p_r q_r - 4p_s q_s p_r q_r \\
 &= (p_s - q_s)(p_r - q_r) D_{sr} = (2p_s - 1)(2p_r - 1) D_{sr} \\
 B_{sr} &= \gamma_{sr} - h_s h_r = 4D_{sr}^2 + 2(p_s p_r + q_s q_r - p_s q_r - p_r q_s) D_{sr} + 4p_s q_s p_r q_r \\
 &\quad - 4p_s q_s p_r q_r \\
 &= 4D_{sr}^2 + 2A_{sr}
 \end{aligned}$$

256 For notational convenience, we define:

$$E[A_{sr}] = \frac{1}{N_2} \sum_{r < s} A_{rs}$$

257 In the absence of linkage disequilibria among pairs ( $D_{sr} = 0$  and therefore  
 258  $A_{sr}, B_{sr} = 0$  for all  $s, r$  pairs),  $\gamma_{sr} = \delta_{sr} = h_s h_r$  and  $\Delta = 0$ , i.e. the sample  
 259 variances under WHS and ILS are equal. Otherwise, because  $B_{sr} \geq A_{sr}$   
 260 for  $A_{sr} > 0$ ,  $E[A_{sr}] > 0$  is a sufficient condition for  $\Delta > 0$ . This condition  
 261 is satisfied provided that the sum of weighted linkage disequilibria  $A_{sr}$  is  
 262 positive, i.e.

$$\sum_{sr} A_{sr} = \sum_{sr} (2p_s - 1)(2p_r - 1) D_{sr} > 0. \quad (7)$$

263 While  $E[A_{sr}] > 0$  is a sufficient condition for  $\Delta > 0$ , it is not a necessary  
 264 condition. In fact, the variance in mean pairwise distance under ILS may in  
 265 some cases still be lower than under WHS even for  $E[A_{sr}] < 0$ . This follows  
 266 because negative  $A_{sr}$  may be offset by the positive contributions of  $D_{sr}^2$  to  
 267 the  $B_{sr}$  term when pairwise LD values in the population are sufficiently  
 268 high. However, for large sample sizes, the  $A_{sr}$  term dominates because it  
 269 scales as  $\sim 1/n$  while the  $B_{sr}$  term scales as  $\sim 1/n^2$ , which means that for  
 270 many practical cases the sign of  $E[A_{sr}]$  predicts that of  $\Delta$ .

271 In order to have  $E[A_{sr}] > 0$ , it is required that on average  $A_{sr}$  is positive,  
272 i.e. that for most pairs of loci  $s, r$ , the "major" alleles (those with  $p_s, p_r > 0.5$ )  
273 are in positive LD, while major and minor allele pairs ( $p_s > 0.5, p_r < 0.5$  or  
274 vice-versa) are in negative LD. The weighted LD  $A_{sr}$  provides a measure of  
275 the extent to which major alleles are in positive LD, regardless of whether the  
276 more common allele is a reference/wildtype or variant/mutant at a particular  
277 site. Our results predict that when the mean weighted LD is positive, the  
278 sample variance (error) in estimated pairwise genetic distance will be lower  
279 under ILS than under WHS.

## 280 **2.4 Implications**

281 To understand the conditions under which  $\Delta > 0$  holds, we consider the dis-  
282 tribution of allele frequencies and pairwise LD under different evolutionary  
283 scenarios. Specifically, we ask whether positive weighted linkage disequi-  
284 libria (the conditions in Eqn. (7)) are general enough to assume that ILS  
285 generally leads to a reduced error in estimated genetic distance relative to  
286 WHS.

287 Consider a population undergoing random mutation under an infinite  
288 sites model and Fisher-Wright genetic drift in a finite population. At an  
289 equilibrium of new alleles acquired via mutations and those lost by genetic  
290 drift, the expected number of sites  $\eta_k$  that have  $k$  copies of a mutant allele  
291 is

$$E[\eta_k] = \theta/k,$$

292 (Watterson 1975, see also e.g. Ewens 2004 Ch 9, Ch. 2 in Durrett 2008),  
293 so that the expected frequency of alleles occurring as  $k$ -tuples is  $\theta/(S_N k)$ .

294 Because of this harmonic relationship, the majority of mutant alleles in a  
295 population are represented as singletons and as other small  $k$ -tuples (e.g.  
296  $k = 2, 3$ , etc). This is consistent with a majority of alleles in the population  
297 being rare and of recent origin, with variant allele frequencies close to  $p \sim$   
298  $1/N$ . These rare alleles of recent origin are usually lost from the population,  
299 while a much smaller subset of alleles in the sample have frequencies  $p > 0.5$   
300 and consequently a high probability  $p$  of eventual fixation in the population.  
301 As a result, for the majority of variant allele pairs in a sample, we have  
302  $p_s, p_r \ll 0.5$ .

303 In the absence of recombination, multilocus haplotypes behave as alle-  
304 les at a single locus, so that the infinite sites model becomes effectively an  
305 infinite alleles model (Tajima, 1996). Therefore, every new mutation is in  
306 positive LD with the other variant alleles with which it co-occurs and in  
307 negative LD with non-co-occurring mutations on other haplotypes. We con-  
308 sider the following scenarios: A) LD among rare, typically non-co-occurring  
309 alleles on different haplotypes, B) LD between rare and common alleles when  
310 a recent mutation appears on a common haplotype as a genetic background  
311 and C) co-occurrence of common alleles on dominant haplotypes.

312 A) Following the Watterson distribution of  $k$ -tuples at equilibrium, there  
313 are a large numbers of rare alleles  $p_s, p_r \sim 1/N$ . However, most of these  
314 rare alleles do not co-occur with one another, consequently  $P(sr) \sim 0$  and  
315  $D_{sr} \sim -1/N^2$ . B) Rare alleles typically appear against a background of  
316 common haplotypes or subclones defined by high-frequency variant alle-  
317 les. If we have a recent mutation with frequency  $p_s \sim 1/N$  appearing  
318 on a background of common alleles at other loci  $p_r \sim 0.1$ , then the LD  
319  $D_{sr} \sim p_s - p_s p_r > 0 \sim 1/N$ , because the frequency of the  $P(sr)$  haplotype  
320 is  $p_s$ . By symmetry, new alleles that happen to co-occur on rare haplotypes

321 will have  $D_{sr} \sim -1/N$  with respect to the sites on their genetic background,  
322 but there will be an order of magnitude fewer such associations because the  
323 majority of new mutations will appear against a genetic background of com-  
324 mon haplotypes. C) Similarly, common alleles that co-occur on dominant  
325 subclones have  $D_{sr} = p_s - p_s p_r \sim 0.1$  (where dominant haplotypes, and  
326 therefore allele frequencies can potentially be  $p_s, p_r > 0.5$ ).

327 These heuristic considerations of scale suggest that the distribution of  
328  $D_{sr}$  in the clonal population will be highly skewed, consisting of large num-  
329 bers of negative but near-zero LD values for the many  $p_s \sim 1/N$  rare alleles,  
330 and a smaller number of large positive associations associated with muta-  
331 tions defining the common haplotypes. This conclusion is consistent with  
332 the highly skewed sampling distributions of  $D_{sr}$  for non-recombining loci  
333 computed numerically in Golding (1984), i.e. large numbers of weakly neg-  
334 ative associations and a small number of high positive LD.

335 Because of this skew, we hypothesize that populations where the al-  
336 lele frequency distributions are in approximate equilibrium under mutation  
337 and drift will have positive  $E[A_{sr}]$  and therefore higher sample variance in  
338 pairwise genetic distance under WHS than under ILS. In contrast, among  
339 populations where all allelic variation is of recent origin and characterized by  
340 low frequencies (such as in newly emergent tumors, or in populations that  
341 have experienced recent bottlenecks), the negative associations  $D_{sr} < 0$  will  
342 dominate the distribution due to the fact that recent mutations will ini-  
343 tially occur on different reference haplotypes. Because most mutations will  
344 occur on disjoint branches of the genealogy, very few haplotypes with sig-  
345 nificant numbers of co-occurring mutations will have attained high enough  
346 frequencies to offset the small magnitude but negative LD values. There-  
347 fore, it is possible to observe  $E[A_{sr}] < 0$  (albeit with  $|E[A_{sr}]|$  and  $\Delta \sim 0$ )

348 in populations where most variant alleles occur at near zero frequencies.

349 We will assess these heuristic predictions about the sign and magni-  
350 tude of  $E[A_{sr}]$  and  $\Delta$  under different frequency distributions of  $p$  and  $D_{sr}$   
351 through simulations of mutation and genetic drift for a range of population  
352 parameters.

### 353 **3 Comparison to individual-based simulations**

354 To simulate Fisher-Wright genetic drift in an infinite sites model, we initial-  
355 ized a population of  $N$  haploid genotypes characterized by  $K = 10^8$  sites  
356 with reference genotypes (all alleles set to 0 value, to distinguish them from  
357 variant mutations set to 1). In every generation,  $N$  individuals were sampled  
358 with replacement from the existing pool, with each individual sampled pro-  
359 ducing a single progeny. The number of mutations  $m$  for each progeny was  
360  $m \sim Poiss(Ku)$ , with the mutations randomly distributed among the  $K$   
361 sites. This process was iterated over  $T$  generations; in order to approximate  
362 a distribution of mutation frequencies near equilibrium, we chose  $T \sim 4N$   
363 (because expected coalescent time for all  $N$  haplotypes in a population is  
364  $E[T_C] = 2N$ ). In addition, simulations were run for a range of values  $T < N$   
365 for comparison to non-equilibrium distributions of allele frequencies and  
366 pairwise LD. For each combination of parameters, the simulation cycle was  
367 run over 100 replicates.

368 In order to simulate WHS sampling,  $n$  haplotypes were randomly selected  
369 without replacement from the model population. The Hamming distances  
370 were calculated for all pairs in a sample, while variant allele frequencies and  
371 linkage disequilibria were calculated for all individuals and all pairs in the  
372 model population. ILS sampling was simulated by selecting  $n$  alleles without  
373 replacement at every segregating site, summing pairwise distances over all

374 sites (this can be thought of as sampling with replacement with respect  
375 to genomes, but without replacement with respect to each locus).  $\Delta$  was  
376 estimated as the difference in the sample variances between the WHS and  
377 ILS pairwise distances. For each simulation replicate,  $A_{sr}$  was calculated  
378 from the mutation frequencies  $p_s, p_r$  and from  $D_{sr}$  using Eqn. (6). All  
379 simulations were implemented using Python 2.7.3, the code is available from  
380 the corresponding author upon request.

381 Simulation output for population sizes  $N = 200, 500$ , a sample size of  
382  $n = 20$  and a range of generation times  $T$  are summarized in Tables 1 and  
383 2. The first table shows the estimated parameter values from which  $\Delta$  is  
384 calculated - including the number of polymorphic sites  $S_N$  in the population  
385 (as opposed to the sample number of segregating sites  $S_n$ ), the population  
386 mean allele frequency across polymorphic sites, the sample mean pairwise  
387 genetic distances under WHS and ILS (for  $n = 20$ ), as well as their respec-  
388 tive sample variances over 100 replicates.

389

390 TABLES 1 and 2 HERE

391

392 For small time intervals  $T < N$ , there are few ( $\sim 100$ ) polymorphic sites,  
393 all of which are characterized by low variant allele frequencies. Consequently,  
394 the mean and variance of genetic distances are of the order  $\sim 1, \sim 0.1$ , re-  
395 spectively. For  $T \sim 4N$ , allele frequencies and genetic distances tend towards  
396 the equilibrium values predicted under the neutral infinite sites model, e.g.  
397 the estimated pairwise genetic distance  $\hat{\pi}$  converges to the Tajima estimator  
398 for haploids  $\theta = 2Nu$ , which is  $\hat{\pi} = 150, 300$  for  $N = 200, 500$ , respectively.  
399 Table 2 shows the population mean LDs  $\bar{D}_{sr}$  and the sum of weighted LD  
400 values  $\sum A_{sr} = \bar{A}_{sr}N_2$ . We remark that while the mean values of LD are

401 effectively 0 even for large values of  $T$  and  $S_N$ , this is not due to individual  
402 LD values being near 0. Rather,  $\bar{D}_{sr} \sim 0$  is the result of large numbers of  
403 positive and negative LD values with high absolute value, as can be seen  
404 from the large magnitudes of the summed weighted LD. Figures 2a and 2b  
405 show frequency distributions of pairwise LD and weighted linkage LD for a  
406 representative model population.

407

408       FIGURE 2a-b HERE

409

410       Using  $\sum A_{sr}$ , we compute the predicted difference between WHS and  
411 ILS variances  $\Delta_P$  from Eqn. (6). This predicted value is compared to  
412 the simulation estimate  $\Delta_S = var_{WHS} - var_{ILS}$ . The close correspon-  
413 dence between observed and predicted values of  $\Delta$  is confirmed by the fact  
414 that even the largest deviations are within less than two standard error  
415  $SE_{\Delta_S} = \sqrt{var(\Delta_S)/n}$  units with respect to the point estimate  $\Delta_S$ . The fit  
416 between analytical predictions and observed values improves for longer time  
417 intervals (i.e. as the population distribution of allele frequencies and pair-  
418 wise LD approach equilibrium), in part because of the much larger number  
419 of polymorphic sites and the higher frequency of variant alleles at those sites.

420

421       With the exception of populations where there are very few mutations  
422 and where weighted LD values are very close to 0, we have  $\Delta > 0$  for most of  
423 the simulated populations. These results conform to our hypothesis that the  
424 error in genetic distance estimates based on WHS will be greater than those  
425 for ILS for the majority of natural and model populations. The reduction  
426 of error through ILS is strongest for near-equilibrium distributions of allele  
427 frequencies, for large numbers of segregating sites, and for small sample sizes

428 (corresponding to low coverage depth with NGS).  $\Delta$  scales approximately  
429 as  $\sim 1/n$  for sufficiently large  $n$ ; consequently, for sample numbers and  
430 coverage depths of the order  $\sim 100$ ,  $\Delta$  will be smaller by nearly an order of  
431 magnitude relative to the values shown in Table 2 for  $n = 20$  (simulations  
432 were performed for  $n = 10, 50$ , the results are not shown due to qualitative  
433 similarity to the data in Tables 1-2).

434 The two observed cases with  $\bar{A}_{sr} < 0$  are for  $T = 10$  at both simulated  
435 population sizes, with a negative predicted value  $\Delta_P$  for  $N = 500$  (though  
436 not for  $N = 200$ ). In these cases, the  $\Delta$  values are effectively zero within  
437 a standard error unit, so whether positive or negative values are observed  
438 is of purely formal interest (note that for even smaller time intervals  $T = 5$   
439 and even fewer polymorphic sites, both  $\bar{A}_{sr}$  and  $\Delta > 0$ , albeit very small).  
440 This suggests that at least under neutral evolution,  $E[A_{sr}] < 0$  occurs under  
441 rather restricted conditions corresponding to very small absolute values of  $\Delta$   
442 and negligible reduction of error in estimating  $\hat{\pi}$  through either WHS or ILS,  
443 while for large numbers of segregating sites and increasing allele frequencies,  
444 there can be considerable increases in error when  $\hat{\pi}$  is estimated via WHS  
445 rather than ILS.

## 446 **4 Analysis of cancer sequence data**

447 We apply the results of our derivations and numerical analyses to genomic  
448 data by estimating  $\sum A_{rs}$  and  $\Delta$  to haplotype frequencies estimated from  
449 a lung adenocarcinoma tumor sequence data. The data was obtained from  
450 whole-exome sequencing of 4 sections of a primary solid tumor taken from  
451 a lung cancer patient. DNA from the samples was extracted using Agilent  
452 SureSelect capture probes. The exome library was sequenced using paired-  
453 end 100 bp reads on the Illumina HiSeq 2000 platform. Reads were mapped

454 onto the human genome HG19 using BWA (Li and Durbin, 2009), giving  
455 a post-mapping mean coverage (depth) of 60-70 fold across sites. Variant  
456 calls were performed using GATK (McKenna et al., 2010). The unpublished  
457 data were provided to the authors as summaries of variant frequencies and  
458 haplotypes by K. Gulukota and Y. Ji.

459 Through the matching of read ends, somatic mutations co-occurring  
460 within  $\sim 100$  bp in single genomes were identified (Sengupta et al. 2015,  
461 unpublished). These mutation pairs define two locus haplotypes that can  
462 be tallied without the need of phasing. This allows us to estimate the fre-  
463 quencies of haplotypes defined at two adjacent loci directly from the read  
464 counts, along with individual allele frequencies. Following the terminology of  
465 this paper, while non-adjacent polymorphic sites are sampled as (effectively)  
466 ILS, adjacent sites are effectively sampled as whole haplotypes. Because re-  
467 production in tumor cells is asexual and ameiotic, estimates of  $D_{sr}$  and  $A_{sr}$   
468 using a subset of nearly adjacent sites is as representative of other haplotype  
469 pairs as if they came from more distant sites or on different chromosomes.  
470 The adenocarcinoma data contain estimated frequencies of 69 haplotypes  
471 defined by variant alleles at two sites on a single read, and allele frequen-  
472 cies for a total of 138 sites (comparable to the number of somatic mutations  
473 identified in the exomes of lung adenocarcinoma and other cancer types, e.g.  
474 TCGA 2014, Hoadley et al. 2014). The provided haplotype data is used to  
475 determine how the LD values and allele frequencies would effect the error in  
476 estimation of  $\hat{\pi}$  for this data set under WHS vs. ILS sampling.

477 A naive application of Eqn. (6) to the distribution of mutation frequen-  
478 cies and LD values gives  $\Delta \sim 0.1$  for  $n = 65$ , suggesting lower error in  $\hat{\pi}$   
479 estimates from ILS for this data. However, several aspects of cancer ge-  
480 netics complicate this estimate. First, because cancer cells reproduction is

481 clonal, somatic mutations appear in heterozygous genotypes in the absence  
482 of mitotic recombination and gene conversion. A SNV frequency of  $p = 0.5$   
483 corresponds to "fixation" of a somatic mutation in a population of asexual  
484 diploids. Therefore, if we have heterozygous fixation at a single SNV site,  
485 a population consisting of 0/1 (reference and variant) genotypes, a mean  
486 genetic distance measure of  $\hat{\pi} = 1/2$  is meaningless because the population  
487 is homogeneous with respect to the 0/1 genotype. Variant allele frequencies  
488 must be rescaled to reflect these considerations.

489 Figure 3 shows the distribution of mutant allele frequencies in Sample 1,  
490 note the high frequency of values near  $p = 0.5$ , and the fact that this distri-  
491 bution is not consistent with an equilibrium neutral distribution of  $\sim \theta/k$   
492 k-tuples, due to the scarcity of detected rare variants.

493

494 FIGURE 3 HERE

495

496 Williams et al. (2016a,b) (see also Ling et al. 2015) address the issue  
497 of the fixation of heterozygous genotypes by only considering polymorphic,  
498 segregating sites when comparing allele frequencies in tumors to those pre-  
499 dicted from the neutral model, to the exclusion of sites that are  $\geq 0.5$  within  
500 a margin of sampling error. This also excludes those sites with frequencies  
501  $p > 0.5$  due to loss of heterozygosity. In addition, with a range of allele  
502 frequencies  $p = [0, 0.5]$ , the frequencies are rescaled to reflect the frequency  
503 of the heterozygous genotype, which for diploids means mapping  $p' = 2p$ , or  
504 more generally,  $p' = p/f_c$  where  $f_c$  is the cutoff for the inference of fixation.  
505 With this mapping, the genetic distance for a sample where all genotypes  
506 at a variant site are 0/1 is 0.

507 With the assumption of diploidy at all of the genotyped SNV sites and

508 defining fixation as  $p = 0.5$ , we find that for  $n = 65$ , the binomial prob-  
509 ability of observing fewer than  $x = 26$  mutant alleles is  $Bin(x \leq 25|n =$   
510  $65, p = 0.5) = 0.041$ . Thus, we use  $f_c = 0.4$  as a cutoff defining poly-  
511 morphic sites. Using this criterion, and the rescaling  $p' = p/f_c$ , there are  
512 only between 6 (sample 4) and 10 (sample 3) adjacent segregating sites, and  
513 consequently between 3 and 5 haplotypes defined by such a pair out of the  
514 original 69. The LD and  $\Delta$  values for this subset of haplotypes are summa-  
515 rized in Table 3. The differences in variances  $\Delta$  remain positive, consistent  
516 with sample variance under WHS being greater than under ILS as before.  
517 However  $\Delta$  is small ( $0.034 \leq \Delta \leq 0.070$ ), suggesting that in practice the  
518 estimation errors for  $\hat{\pi}$  are negligibly different for this data set. The small  $\Delta$   
519 are partly a consequence of the small number of segregating sites (because  
520  $\hat{\pi}_{max} = S_n/2$ ). Therefore, the variance in  $\hat{\pi}$  estimation under WHS may  
521 be expected to increase for greater numbers of segregating sites, as was the  
522 case in the simulation data for larger time intervals and  $S$ .

523

524 TABLES 3a-b HERE

525

526 The values of  $\Delta$  are also sensitive to the choice of truncation, as many  
527 of the SNVs occur in genotypes that are close to fixation in the tumor.  
528 For example, if we use  $f_c = 0.49, x = 32$  as a cutoff to define segregat-  
529 ing sites rather than  $f_c = 0.40$ , we obtain  $\bar{A}_{sr} < 0$  and  $\Delta < 0$  (of the  
530 order  $\sim 0.1$ ). The sign reversal results from some lower frequency SNVs  
531 uniquely co-occurring in genomes with other SNVs that are close to fixation.  
532 The remaining allele and haplotype distributions contribute negative link-  
533 age disequilibria between the high frequency SNVs at one locus and high  
534 frequency reference alleles at the other site. The greater absolute value of

535  $\Delta$  is a consequence of the fact that with a cutoff of  $f_c = 0.49$ , there are  
536 now 21-28 haplotypes (and 42-56 segregating sites) rather than the 6-10 for  
537 the  $f_c = 0.40$  cutoff. The negative weighted LDs and  $\Delta$  with this cutoff  
538 are shown in the second panel Table 3b, as an illustration of how for some  
539 samples, the variances in  $\hat{\pi}$  may actually be lower under WHS than under  
540 ILS.

## 541 **5 Discussion**

542 Heuristically, the higher error in estimated genetic distance under WHS  
543 when weighted LD are positive on average reflects the loss of information  
544 due to non-independence across sites. If for most pairs of sites, the most  
545 frequent (major) alleles are in positive LD, then any error in estimating  
546 frequency and heterozygosity at one site covaries with the error at the other  
547 sites. In contrast, with ILS, each site provides independent information and  
548 the error across sites is uncorrelated. If there are  $S_n$  segregating sites in a  
549 sample of  $n$  and the variance in estimated genetic distance per site is  $\sigma^2$ ,  
550 then with independent sampling the error across sites will approach  $\sigma^2/S_n$ .  
551 In contrast, in the extreme case where allele frequencies across sites are  
552 nearly identical (complete linkage), the sample variance is  $\sigma^2$  independent  
553 of the number of sites. In the case of negative LD (i.e. negative association  
554 among common alleles), there is an information gain across sites.

555 On the other hand, a negative association of allele frequencies across  
556 pairs of sites means that an error in estimated distance at one site will on  
557 average be compensated by an error in the opposite direction at another  
558 site, leading to reduction in variance under WHS (analogous to improved  
559 estimation of the mean by sampling positive and negative extremes of a  
560 distribution). Both heuristic considerations and simulation results suggest

561 that such a scenario is unlikely except for distributions of allele frequencies  
562 that give very small error values regardless.

563 Because  $\Delta$  will either be positive or close to 0 for most distributions  
564 of allele frequencies, our results suggest that ILS should be used to mini-  
565 mize error in genetic distance estimation for most natural and experimental  
566 populations. However, there are several caveats to this conclusion, some the-  
567 oretical, others practical. For example, we know that when most pairwise  
568 LD are approximately 0, the difference  $\Delta$  between WHS and ILS estimates  
569 will be very small. A number of recent studies have shown that LD are  
570 generally among sites that are not physically linked in the genomes of sex-  
571 ually reproducing model organisms, including *Drosophila* (Andolfatto and  
572 Przeworski, 2000) and humans (Peterson et al., 1995; Reich et al., 2001).  
573 This suggests that any error introduced by sampling alleles from genomes  
574 (WHS) rather than individually via ILS will be negligible.

575 In contrast, for the genomes of clonal, ameiotic organisms or for regions  
576 of genome under very low recombination in sexually reproducing organisms,  
577 LD values will be high. Depending on the distribution of allele frequencies,  
578  $\Delta$  will be large when evaluated over many polymorphic sites. In the cases of  
579 cancer and microbial genomics, the standard NGS approach to sequencing  
580 reads from large numbers of cells (approximating ILS) suggests an improved  
581 estimation of  $\hat{\pi}$  (and consequently,  $\theta$  and  $N_e$ ) relative to what would be  
582 obtained from more expensive single cell sequencing approaches. Moreover,  
583 single-cell sequencing usually entails a much smaller sample size  $n$  than the  
584 coverage depths of 100-1000 that are standard for NGS. Even in cases where  
585  $\Delta < 0$  (such as for some of the simulated data with small numbers of rare  
586 mutations, or for some truncations of the lung cancer data), the magnitude  
587 of the effect is going to be small and outweighed by the reduction of error

588 through high coverage. Moreover,  $\Delta$  is defined on the assumption of the  
589 same effective sample size  $n$  for both WHS and ILS, if ILS allows for much  
590 larger  $n$ , as is often the case, then this is often sufficient to reverse the sign  
591 of  $var(\hat{\pi}_{WHS}) - var(\hat{\pi}_{ILS})$ .

592 In addition to providing a summary statistic of genetic variation in a  
593 population,  $\hat{\pi}$  is an estimator of population mutation rate  $\theta$  (and, with a  
594 known mutation rate, effective population size  $N_e$ ) under a neutral model  
595 of sequence evolution. As noted in the introduction, these parameter es-  
596 timates can be used to detect the population genetic signatures of natural  
597 selection and/or demographic histories when compared to  $\theta$  estimates from  
598 the sample number of segregating sites  $S_n$ . Consequently, our derivation  
599 of the expectation and sample variance in  $\hat{\pi}$  under WHS and ILS are key  
600 to calculating the error in estimates of  $\theta$  and  $N_e$ . Sampling error in the  
601 Tajima D statistic can be estimated using our derivation of  $\Delta$  together with  
602 an analogous estimate for the sampling error of  $S_n$ .

603 Another future research direction suggested by our results is deriving  
604 analytically the conditions under which  $E[A_{sr}], \Delta > 0$ . Eqn. (7) provides  
605 the conditions in terms of allele frequencies and LD under which  $\Delta > 0$ ,  
606 but does not specify the population genetic conditions under which these  
607 distributions hold. For example, showing that an equilibrium distribution  
608 of allele frequencies under Fisher-Wright drift both without recombination  
609 and for a range of recombination rates leads to  $\Delta > 0$  requires deriving a  
610 population distribution (as opposed to the distribution within the sample)  
611 of pairwise LD values  $D_{sr}$ . Computing  $E[A_{sr}]$  over a distribution of al-  
612 lele frequencies and pairwise LD would essentially formalizing the heuristic  
613 argument presented in subsection 2.4

614 Finally, we remark that this study was to a large part motivated by

615 efforts to apply the methods and theory of population genetics to cancer  
616 biology, where whole haplotype versus individual locus sampling appear as  
617 options under single cell sequencing versus WGS of multicell samples, re-  
618 spectively. The case study from lung cancer data in the previous section was  
619 used as proof of principle. A more accurate and refined analysis would have  
620 to take into consideration a number of potentially confounding variables.  
621 These include polyploidy and aneuploidy (so that with ploidy  $X$ , fixation  
622 corresponds to  $p = 1/X$ ), as well as accounting for the loss of heterozygosity  
623 through mitotic recombination, reflected in frequencies  $p > 0.5$ . The sensi-  
624 tivity of  $\Delta$  to the choice of cutoff  $f_c$  defining fixation, even in the diploid  
625 cases, bears further investigation as well.

## 626 **6 Acknowledgments**

627 MS and JL were supported by the St. David's Foundation impact fund. YN  
628 and PM were supported by NIH grant 2R01CA132897-06A1 . We thank  
629 Kalamakar Gulukota and Yuan Ji at NorthShore University HealthSystem  
630 for providing the lung cancer data used in this paper. We also thank the  
631 following individuals for their helpful comments: Matthew Cowperthwaite,  
632 Habil Zare, Mark Kirkpatrick, Jeffrey Townsend, Vincent Cannataro, and  
633 Andrea Sottoriva.

## 634 **7 Statement of Effort**

635 MS proposed the study, wrote the manuscript, ran the simulations, and  
636 analyzed the data. YN and PM derived most of the equations in section 2  
637 and in the Appendix. JL wrote the python code used for the simulations.

## 638 8 Figures and Tables

639 **Figure 1.** Illustration of whole haplotype sampling (WHS) versus indi-  
640 vidual locus sampling (ILS). In this example, the population consists of 8  
641 haploid organisms  $G1\dots G8$  characterized by 4 segregating sites  $S1\dots S4$ . We  
642 assume a sampling depth of  $n = 3$  and sufficiently many reads to capture  
643 all segregating sites. In the left panel, we have a random instance of WHS  
644 via the sampling of  $G2, G4, G5$  (gray ovals representing sampling), giving a  
645 mean pairwise distance of  $\hat{\pi} = 2$ . In the right panel, we have a random ILS  
646 such that  $G1, G3, G8$  are sampled at  $S1, G4, G5$  and  $G8$  at  $S2$ , etc, giving  
647 a mean genetic distance  $\hat{\pi} = 8/3$ .

648

649 **Figures 2a-b.** Population distributions of pairwise linkage disequilibria  $D_{sr}$   
650 (2a) and weighted linkage disequilibria  $A_{sr}$  (2b) for a simulated population  
651 with  $N = 500$  haploid genotypes after  $T = 2500$  generations of mutation  
652 and Fisher-Wright genetic drift, corresponding an approximate equilibrium  
653 allele frequency distribution.

654

655 **Figure 3.** Distribution of allele frequencies  $p$  in the first lung adenocarci-  
656 noma sample, for  $S_n = 138$  polymorphic sites. Values of  $p$  near 0.5 indicate  
657 heterozygous variant genotypes near fixation. Values  $p > 0.5$  are a conse-  
658 quence of loss of heterozygosity via gene conversion during mitotic recom-  
659 bination, these are excluded from our analyses.

660

661 **Table 1.** A summary of results for a Fisher-Wright model of genetic drift  
662 with infinite sites. The table shows a comparison of  $\Delta_P$  values predicted  
663 from Eqn. (6) with simulation the values  $\Delta_S$  for  $N = 200, 500$  and sample  
664 size/coverage depth  $n$  for a range of time intervals (the last pair of time

665 values for each population size is of the order  $4N$ , corresponding to an ap-  
666 proximate equilibrium in allele frequencies). The standard error of  $\Delta_S$  is  
667 also shown, where  $\Delta P$  lies within less than two SE units from  $\Delta_S$  even for  
668 small time intervals where there are few mutations. Mean population pair-  
669 wise linkage disequilibrium values are all essentially zero for all simulations,  
670 while the magnitudes of  $A_{sr}$  increase with  $T$  as predicted.  $p$  is the mean  
671 variant allele frequency across all segregating sites.

672

673 **Table 2.** This table shows the number of segregating sites  $S_n$  in a sample  
674 of  $n = 20$ , the mean pairwise genetic distances  $\hat{\pi}_W, \hat{\pi}_I$  (for WHS and ILS,  
675 respectively), and the variances in pairwise genetic distance for WHS and  
676 ILS. The latter are used to compute  $\Delta_S$  in Table 1.

677

678 **Table 3.** Calculation of  $\Delta$  from haplotype and allele frequencies in the  
679 lung adenocarcinoma sequence data, where haplotype frequencies for sites  
680 on individual long reads are known. Note that  $\hat{A} > 0$  and  $\Delta > 0$  for all  
681 4 samples, indicating that the error in pairwise genetic distance estimates  
682 for this data set are greater under WHS than under ILS, albeit weakly  
683 given the small number of unique haplotypes.  $\Delta$  is computed using the  
684 actual mean coverage depth  $n = 65$  for two different cutoffs used to define  
685 polymorphic sites. The upper panel shows the values for a cutoff of  $f_c = 0.40$ ,  
686 selected based on a binomial probability. The lower panel shows the same  
687 for  $f_c = 0.49$ , selected arbitrarily close to  $p = 0.5$  to show the sensitivity  
688 of  $\Delta$  to the cutoff. The  $f_c = 0.40$  calculations are based on 6-10 remaining  
689 polymorphic sites, the  $f_c = 0.49$  on 42-56 sites, depending on the sample.  
690 Note that  $\bar{p}'$  is based on  $p' = p/f_c$ , rescaled with respect to the diploid cutoff  
691 value.

## 692 **References**

- 693 Andolfatto, P. and Przeworski, M. (2000). A genome-wide departure from  
694 the standard neutral model in natural populations of drosophila. *Genetics*,  
695 156(1):257–268.
- 696 Burrell, R. A., McGranahan, N., Bartek, J., and Swanton, C. (2013). The  
697 causes and consequences of genetic heterogeneity in cancer evolution. *Nature*,  
698 501(7467):338–345.
- 699 Dexter, D. L. and Leith, J. T. (1986). Tumor heterogeneity and drug resis-  
700 tance. *Journal of Clinical Oncology*, 4(2):244–257.
- 701 Durrett, R. (2008). *Probability models for DNA sequence evolution*. Springer  
702 Science & Business Media.
- 703 Ellegren, H. and Galtier, N. (2016). Determinants of genetic diversity. *Nature Reviews Genetics*.  
704
- 705 Ewens, W. J. (2004). Mathematical population genetics. i. theoretical in-  
706 troduction. *interdisciplinary applied mathematics*, 27.
- 707 Fu, Y.-B. (2015). Understanding crop genetic diversity under modern plant  
708 breeding. *Theoretical and Applied Genetics*, 128(11):2131–2142.
- 709 Gawad, C., Koh, W., and Quake, S. R. (2016). Single-cell genome sequenc-  
710 ing: current state of the science. *Nature Reviews Genetics*, 17(3):175–188.
- 711 Golding, G. (1984). The sampling distribution of linkage disequilibrium.  
712 *Genetics*, 108(1):257–274.
- 713 Goodwin, S., McPherson, J. D., and McCombie, W. R. (2016). Coming of  
714 age: ten years of next-generation sequencing technologies. *Nature Reviews*  
715 *Genetics*, 17(6):333–351.

- 716 Hansson, B. and Westerberg, L. (2002). On the correlation between  
717 heterozygosity and fitness in natural populations. *Molecular Ecology*,  
718 11(12):2467–2474.
- 719 Hoadley, K. A., Yau, C., Wolf, D. M., Cherniack, A. D., Tamborero, D.,  
720 Ng, S., Leiserson, M. D., Niu, B., McLellan, M. D., Uzunangelov, V.,  
721 et al. (2014). Multiplatform analysis of 12 cancer types reveals molecular  
722 classification within and across tissues of origin. *Cell*, 158(4):929–944.
- 723 Kimura, M. (1969). The number of heterozygous nucleotide sites maintained  
724 in a finite population due to steady flux of mutations. *Genetics*, 61(4):893.
- 725 Lewontin, R. and Kojima, K.-i. (1960). The evolutionary dynamics of com-  
726 plex polymorphisms. *Evolution*, pages 458–472.
- 727 Lewontin, R. C. et al. (1974). *The genetic basis of evolutionary change*,  
728 volume 560. Columbia University Press New York.
- 729 Li, H. and Durbin, R. (2009). Fast and accurate short read alignment with  
730 burrows–wheeler transform. *Bioinformatics*, 25(14):1754–1760.
- 731 Ling, S., Hu, Z., Yang, Z., Yang, F., Li, Y., Lin, P., Chen, K., Dong, L.,  
732 Cao, L., Tao, Y., et al. (2015). Extremely high genetic diversity in a single  
733 tumor points to prevalence of non-darwinian cell evolution. *Proceedings*  
734 *of the National Academy of Sciences*, 112(47):E6496–E6505.
- 735 MacLean, R. C., Hall, A. R., Perron, G. G., and Buckling, A. (2010). The  
736 population genetics of antibiotic resistance: integrating molecular mech-  
737 anisms and treatment contexts. *Nature Reviews Genetics*, 11(6):405–414.
- 738 Martinez, J. and Baquero, F. (2000). Mutation frequencies and antibiotic  
739 resistance. *Antimicrobial agents and chemotherapy*, 44(7):1771–1777.

- 740 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernyt-  
741 sky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M., et al. (2010).  
742 The genome analysis toolkit: a mapreduce framework for analyzing next-  
743 generation dna sequencing data. *Genome research*, 20(9):1297–1303.
- 744 National Research Council (1993). *Livestock*. The National Academies Press,  
745 Washington, DC.
- 746 Navin, N. E. (2015). The first five years of single-cell cancer genomics and  
747 beyond. *Genome research*, 25(10):1499–1507.
- 748 Peterson, A. C., Di Rienzo, A., Lehesjoki, A.-E., de la Chapelle, A., Slatkin,  
749 M., and Frelmer, N. B. (1995). The distribution of linkage disequilibrium  
750 over anonymous genome regions. *Human molecular genetics*, 4(5):887–  
751 894.
- 752 Reich, D. E., Cargill, M., Bolk, S., Ireland, J., Sabeti, P. C., Richter, D. J.,  
753 Lavery, T., Kouyoumjian, R., Farhadian, S. F., Ward, R., et al. (2001).  
754 Linkage disequilibrium in the human genome. *Nature*, 411(6834):199–204.
- 755 Sengupta, S., Wang, J., Lee, J., Müller, P., Gulukota, K., Banerjee, A.,  
756 and Ji, Y. (2015). Bayclone: Bayesian nonparametric inference of tumor  
757 subclones using ngs data. In *Proceedings of The Pacific Symposium on*  
758 *Biocomputing (PSB)*, volume 20.
- 759 Sun, X.-x. and Yu, Q. (2015). Intra-tumor heterogeneity of cancer cells and  
760 its implications for cancer treatment. *Acta Pharmacol Sin*, 36(10):1219–  
761 1227.
- 762 Tajima, F. (1989). Statistical method for testing the neutral mutation hy-  
763 pothesis by dna polymorphism. *Genetics*, 123(3):585–595.

- 764 Tajima, F. (1996). Infinite-allele model and infinite-site model in population  
765 genetics. *Journal of Genetics*, 75(1):27–31.
- 766 TCGA (2014). Comprehensive molecular profiling of lung adenocarcinoma.  
767 *Nature*, 511(7511):543–550.
- 768 Van Dyke, F. (2008). *Conservation biology: foundations, concepts, applica-*  
769 *tions*. Springer Science & Business Media.
- 770 Watterson, G. (1975). On the number of segregating sites in genetical models  
771 without recombination. *Theoretical population biology*, 7(2):256–276.
- 772 Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A., and Sottoriva,  
773 A. (2016a). Identification of neutral tumor evolution across cancer types.  
774 *Nature genetics*.
- 775 Williams, M. J., Werner, B., Curtis, C., Barnes, C., Sottoriva, A., and  
776 Graham, T. A. (2016b). Quantification of subclonal selection in cancer  
777 from bulk sequencing data. *bioRxiv*, page 096305.

## 778 9 Appendix A1: Ordered Pairs of Pairs

779 Recall the definitions  $\mu = E(\phi_{ij})$ ,  $\sigma^2 = \text{var}(\phi_{ij})$  and  $\kappa = E(\phi_{ij}, \phi_{jk})$ .

**Lemma 1.** *Let  $\mu = E(\phi_{ij})$  where the expectation is over pairs  $x_i \sim p(x)$  and  $x_j \sim p(x)$ , independently. Let  $\hat{\phi}_n = \frac{1}{N_2} \sum_{i < j} \phi_{ij}$ , denote a sample estimate for  $\mu$ , averaging over all pairs  $(i, j)$  of samples. Then  $\hat{\phi}_n$  is unbiased,  $E(\hat{\phi}_n) = \mu$ , and*

$$\text{var}(\hat{\phi}_n) = \frac{\sigma^2}{N_2} + 2 \frac{N_3}{N_2^2} (\kappa - \mu^2).$$

780 *Proof.* Unbiasedness is straightforward:

$$E(\hat{\phi}_n) = E\left(\frac{1}{N_2} \sum_{i < j} \phi_{ij}\right) = \frac{1}{N_2} \sum_{i < j} E(\phi_{ij}) = \mu.$$

781 For the variance, note that

$$\text{cov}(\phi_{ij}, \phi_{kl}) = E(\phi_{ij}\phi_{kl}) - E(\phi_{ij})E(\phi_{kl}) = \begin{cases} 0 & \text{when } \{i, j\} \cap \{k, \ell\} = \emptyset \\ \kappa - \mu^2 & \text{when } |\{i, j, k, \ell\}| = 3 \end{cases}$$

782 Then

$$\text{var}(\hat{\phi}_n) = \frac{\sigma^2}{N_2} + \frac{1}{N_2^2} \sum_P \text{cov}(\phi_{ij}, \phi_{kl}) = \frac{\sigma^2}{N_2} + \frac{2}{N_2^2} N_3(\kappa - \mu^2).$$

783

□

784 *Proof of Eqn. (4).* Let  $\hat{f}_{ns} = \frac{1}{N_2} \sum_{i < j} f_{ij,s}$ . From the statistical indepen-  
 785 dence among sites  $s$  located on different reads under ILS, it follows that for  
 786 a sample of  $n$ ,

$$\text{var}(\hat{\pi}_1) = \sum_s \text{var}(\hat{f}_{ns})$$

787 with

$$\begin{aligned} \text{var}(\hat{f}_{ns}) &= \frac{\sigma_s^2}{N_2} + 2\frac{N_3}{N_2^2}(\kappa_s - \mu_s^2) = \frac{1}{N_2}h_s(1 - h_s) + 2\frac{N_3}{N_2^2}(h_s/2 - h_s^2) \\ &= \frac{1}{N_2}h_s \left\{ 1 - h_s + \frac{N_3}{N_2}(1 - 2h_s) \right\} \end{aligned}$$

788 where the first equality is due to Eqn. (3).

**Table 1**

$N$	$T$	$S_N$	$\bar{p}'$	$\hat{\pi}_w$	$\hat{\pi}_I$	$var_w$	$var_I$
200	5	112.7	0.012	2.94	2.94	0.249	0.241
200	10	181.5	0.016	5.82	5.81	0.468	0.476
200	20	250.1	0.024	11.47	11.47	0.878	0.819
200	50	351.2	0.043	26.79	26.80	2.27	1.58
200	800	770.6	0.315	115.59	114.58	272.0	3.70
200	1000	847.7	0.361	122.90	122.77	415.7	3.97
500	5	308.4	0.0049	2.96	2.96	0.279	0.271
500	10	455.9	0.0066	5.97	5.97	0.537	0.543
500	20	616.9	0.0096	11.61	11.60	1.04	1.04
500	50	875.5	0.0172	28.61	28.68	2.53	2.27
500	100	1078.3	0.0281	54.67	54.67	6.23	3.71
500	2000	2202.4	0.296	301.16	301.22	1532.1	9.09
500	2500	2395.1	0.153	316.06	315.64	2089.8	10.53

**Table 2**

$N$	$T$	$\bar{D}_{sr}$	$\sum A_{sr}$	$\Delta P$	$\Delta S$	$SE(\Delta_S)$
200	5	$-4.57 \times 10^{-7}$	$4.18 \times 10^{-3}$	$-9.57 \times 10^{-4}$	$8.01 \times 10^{-3}$	$4.58 \times 10^{-3}$
200	10	$1.16 \times 10^{-6}$	0.0997	0.0129	$-7.85 \times 10^{-3}$	0.012
200	20	$-3.16 \times 10^{-7}$	0.0529	0.0482	0.0587	0.0226
200	50	$-9.30 \times 10^{-8}$	1.14	0.766	0.687	0.0927
200	800	$-8.89 \times 10^{-6}$	660.5	297.96	268.34	44.56
200	1000	$1.49 \times 10^{-5}$	1009.0	444.94	411.68	51.51
500	5	$1.16 \times 10^{-7}$	$4.58 \times 10^{-3}$	$2.13 \times 10^{-3}$	$8.08 \times 10^{-3}$	$3.00 \times 10^{-3}$
500	10	$-2.83 \times 10^{-7}$	-0.0242	$-7.92 \times 10^{-3}$	$-6.23 \times 10^{-3}$	$7.53 \times 10^{-3}$
500	20	$7.70 \times 10^{-8}$	0.393	0.00	0.0269	0.0213
500	50	$-1.43 \times 10^{-7}$	0.269	0.256	0.259	0.0546
500	100	$-9.80 \times 10^{-8}$	4.35	2.74	2.52	0.182
500	2000	$-4.46 \times 10^{-6}$	3362.8	1606.1	1523.0	213.90
500	2500	$5.31 \times 10^{-6}$	4871.9	2241.3	2079.3	273.37

**Table 3a**

<b>Pr=0.40</b>	<b>S</b>	$\bar{p}$	$\bar{D}_{sr}$	$\sum A_{sr}$	$\Delta$
Sample 1	8	0.492	0.223	0.321	0.045
Sample 2	8	0.423	0.555	0.225	0.034
Sample 3	10	0.457	0.408	0.380	0.054
Sample 4	6	0.328	0.500	0.510	0.070

**Table 3b**

<b>Pr=0.49</b>	<b>S</b>	$\bar{p}$	$\bar{D}_{sr}$	$\sum A_{sr}$	$\Delta$
Sample 1	42	0.753	-0.713	-1.951	-0.653
Sample 2	56	0.760	0.0352	-3.077	-1.040
Sample 3	46	0.754	-0.0907	-1.998	-0.653
Sample 4	56	0.759	-0.474	-2.422	-0.778

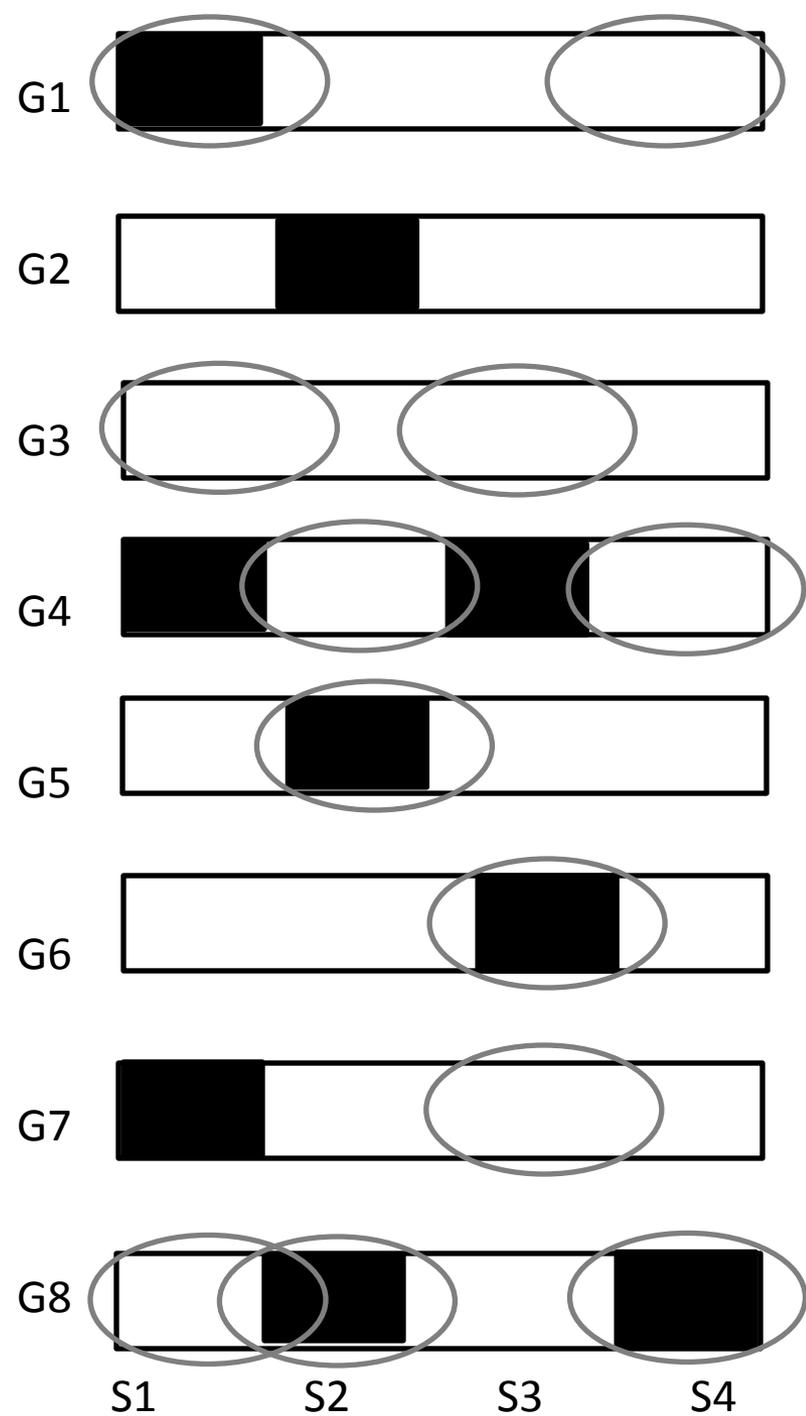
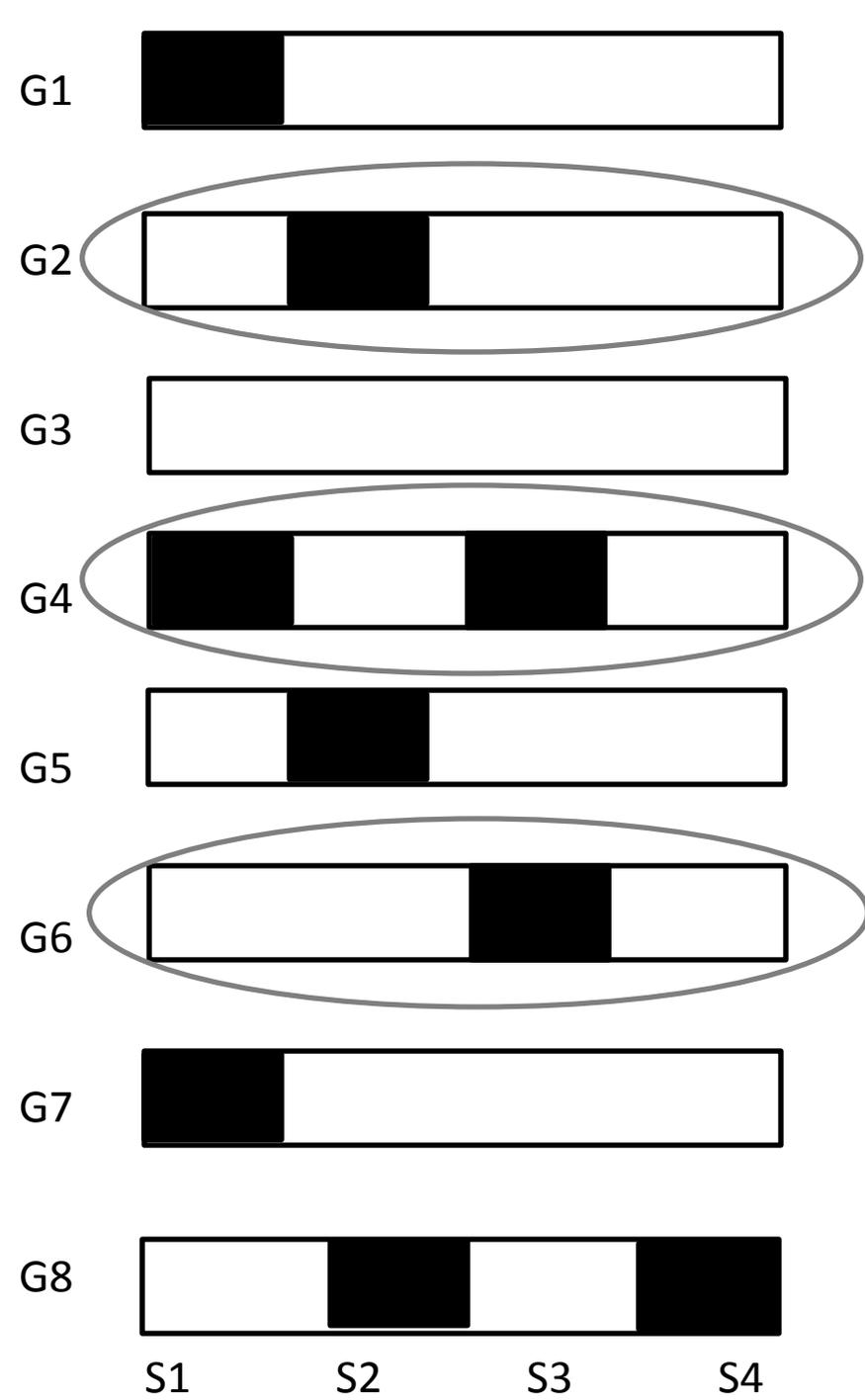


Figure 2a

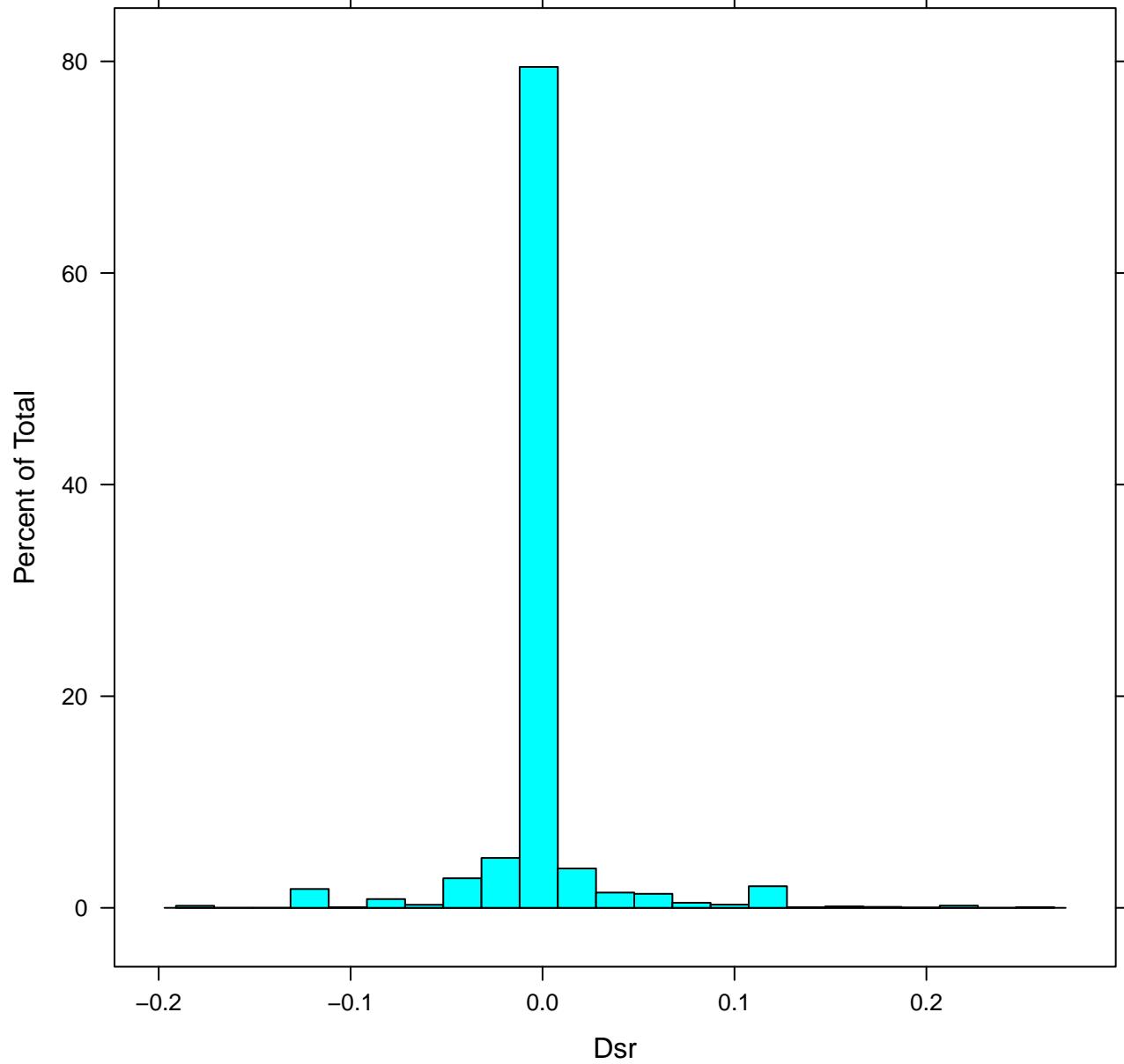


Figure 2b

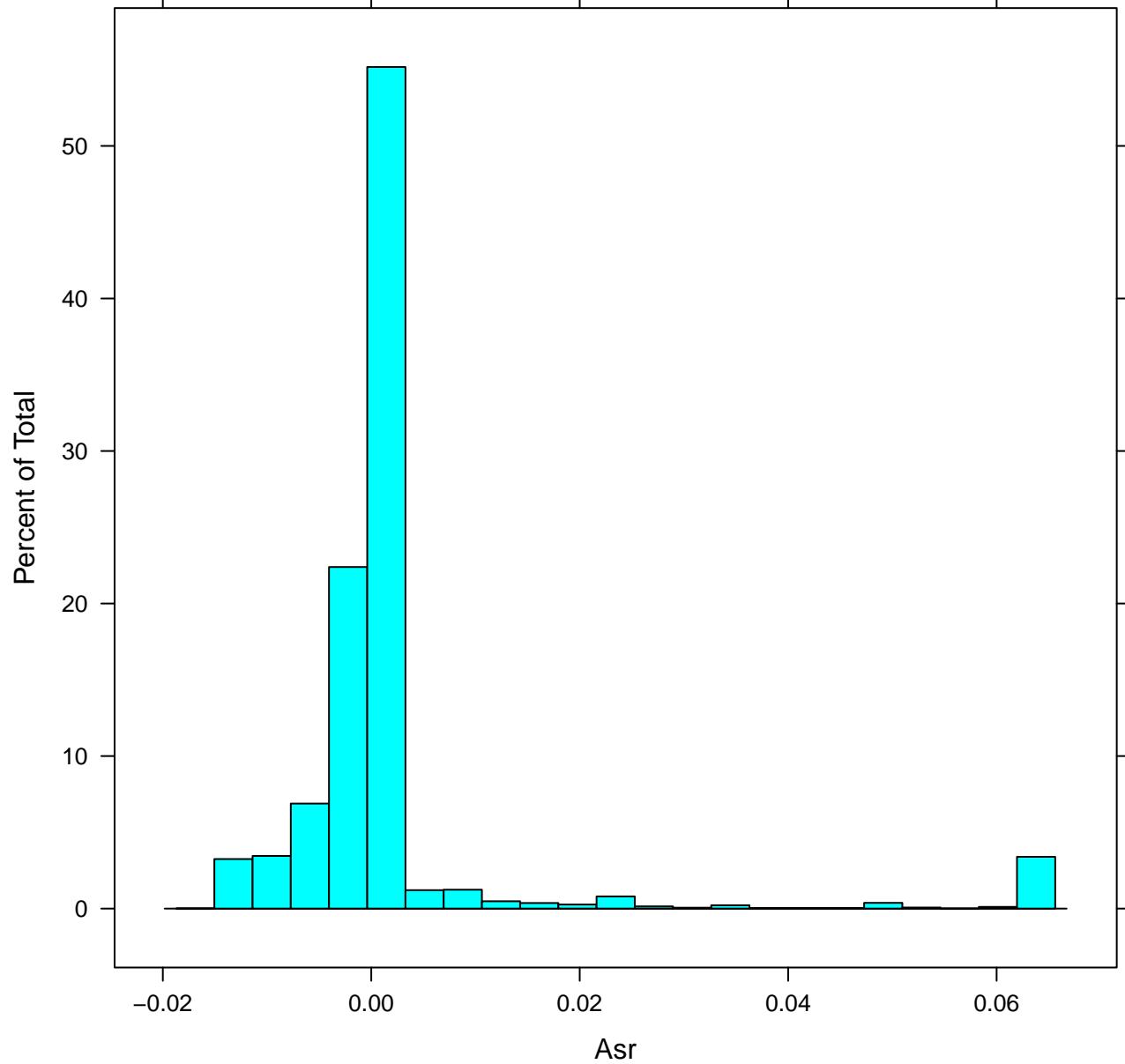


Figure 3

