

Linkage disequilibrium between single nucleotide polymorphisms and hypermutable loci

Sterling Sawaya^{*,1}, Matt Jones[§] and Matt Keller^{*}

^{*}Institute for Behavioral Genetics, University of Colorado, Boulder, CO 80302, USA, [§]Department of Psychology and Neuroscience, University of Colorado, Boulder, CO 80302, USA

ABSTRACT Some diseases are caused by genetic loci with a high rate of change, and heritability in complex traits is likely to be partially caused by variation at these loci. These hypermutable elements, such as tandem repeats, change at rates that are orders of magnitude higher than the rates at which most single nucleotides mutate. However, single nucleotide polymorphisms, or SNPs, are currently the primary focus of genetic studies of human disease. Here we quantify the degree to which SNPs are correlated with hypermutable loci, examining a range of mutation rates that correspond to mutation rates at tandem repeat loci. We use established population genetics theory to relate mutation rates to recombination rates and compare the theoretical predictions to simulations. Both simulations and theory agree that, at the highest mutation rates, almost all correlation is lost between a hypermutable locus and surrounding SNPs. The theoretical predictions break down for middle to low mutation rates, differing widely from the simulated results. The simulation results suggest that some correlation remains between SNPs and hypermutable loci when mutation rates are on the lower end of the mutation spectrum. Consequently, in some cases SNPs can tag variation caused by tandem repeat loci. We also examine the linkage between SNPs and other SNPs and uncover ways in which the linkage disequilibrium of rare SNPs differs from that of hypermutable loci.

KEYWORDS Linkage Disequilibrium, Hypermutability, Tandem Repeats, Missing Heritability, Population Genetics

Introduction

Missing heritability and hypermutable loci

Mutation can take many forms, and can occur at vastly different rates across the human genome (Rando and Verstrepen 2007). Hypermutable regions composed of tandem repeats are of particular interest because of the way in which they mutate. Tandem repeats expand and contract in repeat number at a rate that is orders of magnitude higher than the rate of single nucleotide point mutations (Ellegren 2004; Kelkar *et al.* 2008; Sun *et al.* 2012; Whittaker *et al.* 2003). These regions are able to mutate new alleles and then revert to their original form, all while maintaining their ability to expand and contract. Therefore, not only are many of these loci highly polymor-

phic, but their alleles can often be identical-by-state and not identical-by-descent. Furthermore, tandem repeats are the most common hypermutable loci in the human genome (Ellegren 2004; Rando and Verstrepen 2007), and are often found in regions of functional significance (Sawaya *et al.* 2013).

The rates of expansion and contraction at tandem repeats are known to depend on the length of the tandem repeats, the size of the repeated subunit and the sequence composition. The most mutable are tandem repeats composed of short subunits, called microsatellites (also known as short tandem repeats, or simple sequence repeats). These repeats can have mutation rates up to 10^{-2} (Ellegren 2004), but most have rates between 10^{-3} and 10^{-5} (Whittaker *et al.* 2003; Kelkar *et al.* 2008; Sun *et al.* 2012). The most hypermutable microsatellites tend to have a high A/T content and have a large number of repeated subunits. Because long microsatellites have a tendency to contract more often than they expand (Xu *et al.* 2000), microsatellites undergo a lifecycle in which they are “born” and “die” in the genome over evolutionary time (Kelkar *et al.* 2008; Buschiazio and Gemmell

2010).

Tandem repeats composed of subunits greater than nine base pairs are called minisatellites. Unlike microsatellites, these tandem repeats are not known for their extreme mutability. Their mutation rates are not as well documented (Gemayel *et al.* 2010), but a method to estimate their relative mutation rates is available (Legendre *et al.* 2007). Minisatellites are thought to expand and contract in repeat number through recombination (Jeffreys *et al.* 1998), in contrast to microsatellites which mutate primarily through polymerase slippage and subsequent mismatch repair (Ellegren 2004; Baptiste *et al.* 2013).

Tandem repeat alleles are associated with a range of human diseases (Hannan 2010; Gemayel *et al.* 2010). Of these diseases, perhaps the most well known are caused by expanded microsatellites: Fragile-X disease caused by an expanded CGG repeat (Verkerk *et al.* 1991), and Huntington's disease caused by an expanded CAG repeat (MacDonald *et al.* 1993). Both of these repeats are found in promoters, functional regions near the start of a gene. Promoters have a relatively high density of tandem repeats, suggesting that these hypermutable sequences may play a role in regulating gene expression (Vinces *et al.* 2009; Sawaya *et al.* 2013).

Although tandem repeats are potential sources of heritable disease, recent attention has focused on SNPs for genetic association studies due to technology that allows them to be inexpensively and rapidly genotyped genome-wide. Common SNP variants can be used to measure genome-wide relatedness, and this relatedness can explain a moderate portion of the heritability for complex traits (Yang *et al.* 2011). However, many SNP studies have failed to uncover variants with significant associations (Maher 2008). Furthermore, even SNPs with the strongest associations can only explain a small fraction of heritable genetic variation (Manolio *et al.* 2009).

This lack of significant GWAS hits has been referred to as "missing heritability" (Maher 2008; Manolio *et al.* 2009), and the heritability still not explained by modeling all genome-wide SNPs simultaneously has been termed the "still-missing heritability" (Witte *et al.* 2014; Wray *et al.* 2014). Tandem repeats have been hypothesized to be partially responsible for missing heritability (Hannan 2010; Press *et al.* 2014), and may also be partially responsible for some of the still-missing heritability. Due to their high mutability, tandem repeats can mutate away from linkage with surrounding SNPs, and therefore SNP association studies are not expected to pick up all of the heritability caused by hypermutable variants. Studies using large numbers of tandem repeat loci have shown that tandem repeat variants are usually very weakly linked with surrounding SNPs (Willems *et al.* 2014; Payseur *et al.* 2008; Brahmachary *et al.* 2014). These studies highlight how SNP data can be uninformative about tandem repeat variation, providing further support for the hypothesis that missing heritability might be caused by these hypermutable loci (Willems *et al.* 2014).

However, not all tandem repeat variants are weakly tagged by SNPs. A recent genome wide association study of amyotrophic lateral sclerosis (ALS) in the Finnish population (Laaksovirta *et al.* 2010) uncovered a microsatellite tandem repeat as the most prevalent cause of familial ALS found to

date (DeJesus-Hernandez *et al.* 2011). In the C9ORF72 gene, expansion of a CCGGGG repeat in the first intron results in a dominant allele that causes ALS and can also cause frontal-temporal dementia (DeJesus-Hernandez *et al.* 2011). The expanded repeat allele is in strong linkage disequilibrium with surrounding SNPs (Laaksovirta *et al.* 2010; Mok *et al.* 2012; Majounie *et al.* 2012). Studies of the associated haplotype reveal that the expanded repeat likely arose only once (Mok *et al.* 2012; Majounie *et al.* 2012) and then spread around the globe, possibly along with Viking conquests (Pliner *et al.* 2014). This discovery demonstrates that tandem repeat diseases can be uncovered from SNP association studies.

The 5HTTLPR gene provides another example of how SNPs can be associated with functional tandem repeat variants. Variation in a minisatellite within the 5HTTLPR promoter may be associated with a range of personality phenotypes and neurological diseases (Lesch *et al.* 1996; Wray *et al.* 2009). Two SNPs adjacent to the promoter repeat are in strong linkage disequilibrium with the repeat alleles that have been associated with disease ($r^2=0.72$; Wray *et al.* (2009)).

Together, these studies raise the possibility that more tandem repeat alleles can be uncovered as sources of disease using SNP data. But how quickly do tandem repeats need to mutate to lose their linkage with SNPs and therefore be hidden in SNP association studies? Due to the size of their repeated subunit and their C/G content, the C9ORF72 repeat and the 5HTTLPR repeat are both predicted to have a mutation rate that is lower than most tandem repeats. This suggests that low-mutating tandem repeats have the potential to be tagged by SNPs. To explore this possibility we utilize established population genetics theory and simulations to investigate how mutation rate is related to linkage disequilibrium between a hypermutable locus and surrounding SNPs.

Materials and Methods

Theory relating linkage disequilibrium with mutation rates

We examine the linkage disequilibrium between a hypermutable locus, A/a , and an adjacent SNP marker, B/b , defined by the following mutation dynamics:

$$A \xrightleftharpoons[\mu_a]{\mu_A} a$$

$$B \xrightleftharpoons[\mu_b]{\mu_B} b.$$

We model the hypermutable locus (A/a) as having only two alleles, with equal forward and backward mutation rates (so that $\mu_A = \mu_a$), although it does not perfectly correspond to hypermutable tandem repeat loci. This allows for a simple measure of correlation between the two loci, fitting the population genetics theory outlined below.

We assume the SNP locus (B/b) has a standard low mutation rate and the hypermutable locus has a high mutation rate, such that $\mu_A + \mu_a \gg \mu_B + \mu_b$. The allele frequencies at locus B will be primarily influenced by drift, while the allele frequencies

at A will be influenced by both drift and mutation (we ignore the possibility of selection). Denote the allele frequency of A (B) as p_A (p_B). The allele frequency at locus A is influenced by mutational equilibrium, in which:

$$p_A \approx \frac{\mu_a}{\mu_A + \mu_a}. \quad (1)$$

In a large population with limited drift, the frequency of allele A primarily depends on its forward and backward mutation rates. As population sizes get smaller, and/or the mutation rate gets lower, the allele frequencies are increasingly influenced by population dynamics (as shown in the results).

The allele frequencies at each locus are important because there is an important relationship between the standardized measure of linkage disequilibrium (LD), r^2 , and relative allele frequencies (VanLiere and Rosenberg 2008; Wray 2005; Eberle *et al.* 2007; Hedrick 1987; Hill and Robertson 1968). The maximum possible value of r^2 between two loci is inversely related to the difference between the minor allele frequencies, so if there is a large difference in frequency between the two loci, r^2 cannot be large (Hedrick 1987; Wray 2005; VanLiere and Rosenberg 2008).

Our primary interest is the expected correlation between two loci when one locus has a high mutation rate. For this, the frequency of haplotype AB will be defined as p_{AB} . Linkage disequilibrium, D , is defined as:

$$D = p_{AB} - p_A p_B. \quad (2)$$

The square of the correlation between allele frequencies, r^2 , provides the proportion of variance at one locus that can be explained by another locus, and acts as a standardized measure of LD (Hill and Robertson 1968):

$$r^2 = \frac{D^2}{p_A(1-p_A)p_B(1-p_B)}. \quad (3)$$

How much correlation is expected between loci? To examine this, Ohta and Kimura (1969) define a new variable, ρ^2 , as an approximation for $E(r^2)$. They use the approximation $E(x/y) \approx E(x)/E(y)$ to find an approximation for $E(r^2)$,

$$E(r^2) \approx \rho^2 \equiv \frac{E(D^2)}{E[p_A(1-p_A)p_B(1-p_B)]}. \quad (4)$$

Ohta and Kimura (1969) then solve for the expected values of the numerator and denominator for a diffusion model, obtaining:

$$\rho^2 = \frac{1}{3 + 4N(c+k) - 4/(5 + 2N(c+k) + 2Nk)}, \quad (5)$$

where N is effective population size, and c is the recombination rate between these two loci (here measured in Morgans, M). The variable k is the sum of the mutation rates across both loci, $k \equiv \mu_A + \mu_a + \mu_B + \mu_b$, which is dominated by the mutation rates at the hypermutable locus ($k \approx \mu_A + \mu_a$). To simplify notation, the forward/backward mutation rates at the hypermutable loci will be referred to as simply μ , such that $k \approx 2\mu$.

Somewhat counterintuitively, allele frequency is not present in the approximation for ρ^2 (5). Although allele frequencies are present in the numerator, $E(D^2)$, and denominator, $E[p_A(1-p_A)p_B(1-p_B)]$, their terms cancel resulting in an expression that only involves population size, N , recombination rate, c and the sum of mutation rates, k (Ohta and Kimura 1969). As discussed above, the maximum r^2 value is determined by relative allele frequencies, but these results suggest that, on average at equilibrium, r^2 is a function of only N , c and k . This prediction is examined here using simulated data (see next section). The simulations also use the diffusion model, so the equivalence of (4) and (5), as well as all of our results, rely on the assumptions of the model.

Furthermore, Ohta and Kimura (1969) showed that ρ^2 is only an accurate approximation of $E(r^2)$ when $N(c+k)$ is sufficiently larger than one. In this case ρ^2 is approximated as:

$$\rho^2 \approx \frac{1}{4N(c+k)}. \quad (6)$$

This approximation suggests that mutation and recombination act similarly to reduce linkage disequilibrium. Mutation is slightly different than recombination, however, because it changes allele frequencies, but this effect is reduced if the locus is in mutational equilibrium. More importantly, (6) also suggests that the expected correlation between allele frequencies is very small when $N(c+k)$ is large. Therefore, if the mutation rate is large one would expect a weak correlation between a hypermutable locus and an adjacent SNP marker, unless the effective population size is small.

A. Simulations

Using the coalescent simulation program FastSimCoal (Excoffier and Foll 2011), we simulated a population of 10,000 individuals for a region of 100,000 base pairs (100kb). At the center of the 100kb region we placed hypermutable locus (referred to as a “microsatellite” in FastSimCoal documentation) limited to only two alleles (A and a), with equal forward and backward mutation rates (μ) set to 10^{-3} , 10^{-4} , and 10^{-5} for different simulations. Two-thousand simulation results were obtained for each mutation rate. The recombination rate between adjacent base pairs was set to 10^{-8} , and the mutation rates at surrounding DNA loci were set to $5 \cdot 10^{-8}$. The positions of the polymorphic locus, i.e. loci with a non-zero minor allele frequency, their variants, and the variants at the central hypermutable locus were retrieved from FastSimCoal. These results were converted to necessary file types using custom python scripts, and analyzed in python and R. There were 46 simulations for $\mu = 10^{-5}$ that were excluded because hypermutable loci were not polymorphic.

For each simulation, four statistics were calculated. First, the r^2 values between the central hypermutable locus and surrounding SNPs were calculated. The mean of this value across simulations is referred to as “mean r^2 ”. We expect this simulated measure of LD to be the most accurate estimate of the true degree of association because it does not rely on as many assumptions as the analytical approximation. Second, the average empirical values for D^2 and $p_A(1-p_A)p_B(1-p_B)$ were calculated from the simulations. We refer to the ratio of these two measures as “empirical ρ^2 ”. Next, the values of ρ^2

from (5) were calculated using the three parameters, N , c , and k , that were used in the simulation. We expect the analytical approximation ρ^2 from (5) and empirical ρ^2 to closely match because both the simulations and the statistical approach of Ohta and Kimura (1969) rely on the diffusion approximation. Finally, the position and r^2 for the individual SNP with the highest r^2 value were recorded from each individual simulation.

The simulation results were binned into regions of 100 base-pairs, corresponding to regions along the simulated chromosome relative to the position of the hypermutable locus. The values for r^2 and empirical ρ^2 were calculated and then averaged across SNPs for each 100 base pair bin. The resulting plots were smoothed with LOESS smoothing.

To compare the hypermutable results with SNP-SNP correlations, we simulated a 150-kb region 50 times, with the same parameters as above (10,000 population size, recombination rates of 10^{-8} , and mutation rate of $5 \cdot 10^{-8}$). For each simulation, we used SNPs that were at least 50-kb from the end of the region. Each SNP in this central region was examined separately for its correlations with surrounding SNPs at most 50kb away. This is equivalent to a central SNP in a 100kb region, thus making the LD between two SNPs comparable to the LD between SNPs and hypermutable loci.

1. Results

A. Allele frequencies from simulations

Figures 1 (a)-(c) display the minor allele frequencies (MAFs) for the hypermutable loci, for each mutation rate. At mutation rates of 10^{-3} or 10^{-4} most of the hypermutable alleles have a high MAF. These high mutation rates drive the allele frequencies toward their mutational equilibria of 0.5. In contrast, the allele frequencies for loci with the mutation rate of 10^{-5} are strongly right skewed, with mostly rare alleles. At this lower mutation rate, the allele frequencies appear to be strongly influenced by population dynamics.

The simulated SNP allele frequencies are also strongly influenced by population dynamics, and the MAFs for most of these loci are very low (Figure 1 (d)). As discussed previously, the difference in allele frequencies between two loci influences their maximum possible r^2 . Hypermutable loci with a mutation rate of 10^{-3} have, on average, a high MAF, whereas the average SNP MAF is very low. Therefore, a large difference in allele frequencies exists between rare SNPs and most hypermutable loci, limiting their maximum r^2 .

B. Comparing r^2 estimates with simulated results

For each mutation rate we plot the mean r^2 between a central hypermutable locus and SNPs with any MAF across the entire simulated region (Figure 2, green line). These mean r^2 values are primarily influenced by associations between hypermutable loci and rare SNPs. The mean r^2 values for simulations with a mutation rate of 10^{-3} are very low (Figure 2 (c)), increasing slightly for 10^{-4} (Figure 2 (b)), and more so for 10^{-5} (Figure 2 (a)). We also plot the estimate of ρ^2 made by Ohta and Kimura (1969), equation (5), in red. This approximation is greater than the mean r^2 value for each scenario examined here, and much greater when the mutation rate is low or the inter-locus distance

is short. Importantly, when mutation rates are low or loci are in close proximity, the value of $N(c+k)$ is much less than 1. Consequently, as predicted by Ohta and Kimura (1969), this causes the estimate of ρ^2 to differ from the mean r^2 .

Because the simulations use the same diffusion approximation assumptions as the analytical approach of Ohta and Kimura (1969), we expect the empirical ρ^2 to match the approximation ρ^2 from (5). Empirical ρ^2 and the approximation (5) are nearly identical for the simulations using a hypermutable mutation rate of 10^{-5} or 10^{-4} , but not for 10^{-3} (Figure 2, blue and red lines). We cannot explain this discrepancy. Nevertheless, for a mutation rate of 10^{-3} all three measures of r^2 are very small.

Importantly, the mean r^2 measured here uses hypermutable loci and SNPs with any allele frequencies above 0 (following the assumptions of Ohta and Kimura (1969)). This corresponds to a study in which all, or most, SNPs are genotyped, such as a sequencing study. If a study only uses common alleles, such as on a SNP chip with only common SNPs (MAF > 0.05), then the mean r^2 values found between these common SNPs and a hypermutable site should be different.

To address how SNP minor allele frequencies influence the r^2 between the SNPs and hypermutable loci, we examine the r^2 values for SNPs with different MAFs, averaged across all regions. The horizontal black line in Figure 3 shows the mean empirical r^2 for SNPs binned by MAF value, for each mutation rate. The outer ends of the red vertical lines in this figure indicate the range between the 25th and 75th percentiles (5th and 95th for the ends of the blue lines).

In general, the SNP MAF only has a weak effect on the mean r^2 ; the range of r^2 values is similar for most SNP MAFs. However, for the lowest-MAF SNPs, the maximum possible r^2 values are very small and the distribution of r^2 shows that almost all low-MAF SNPs have very weak associations with the hypermutable locus. More importantly, Figures 3 (b) and (c) show that common SNPs (MAF > 0.1) can sometimes be in relatively high LD ($r^2 > 0.2$) with hypermutable loci at the lower range of mutation rates ($\mu = 10^{-4}$ to 10^{-5}).

C. SNP-SNP correlations

To put all of the above results in context, we examine how SNPs are correlated with each other. We find that, on average, SNPs have an extremely low mean r^2 value with other SNPs (Figure 4 (a)). The maximum mean r^2 value, provided by SNPs in close proximity to the central SNP, is less than 0.05. Importantly, most SNPs have extremely low MAF (Figure 1 (d)), and the mean r^2 value is strongly influenced by weak associations with rare SNPs (not shown). The correlation between common SNPs and rare SNPs is known to be weak (Sun et al. 2011), so the lack of a regional association between a single rare SNP and surrounding SNPs is expected. Furthermore, this scenario represents a breakdown of the approximation; the value of $N(c+k)$ is too small for the approximation to be accurate. Therefore the predicted and empirical ρ^2 of almost 0.45 for the SNPs that are in close proximity are clearly not a good approximation for the mean r^2 .

Because hypermutable elements tend to have higher

MAFs, perhaps a more appropriate comparison is to examine a central SNP only if its MAF is above 0.05. When these common central SNPs are examined for their correlations with surrounding SNPs with any MAF, the mean r^2 values increase, but again the approximation (5) is not a good approximation for $E(r^2)$ because again $N(c+k)$ is too small (Figure 4(b)). To explore how the MAF of surrounding SNPs affects these values, we plot the r^2 values for correlations between a central common SNP and surrounding SNPs with binned MAF (Figure 4 (c)). Again the rare SNPs (MAF < 0.05) show a very weak association, and common SNPs show a higher correlation. Intriguingly, common SNPs tag rare SNPs worse than they tag (the often common) hypermutable elements.

The correlations found using common central SNPs are similar to those found with hypermutable elements with a mutation rate of 10^{-5} (Figure 2). However, the distribution of the r^2 values for common central SNPs (Figure 4 (c)) indicates that the upper 95th percentile of r^2 values for common SNP associations are higher than those of any hypermutable element (Figure 3 (c)). Therefore, large r^2 values (e.g. $r^2 > 0.5$) will be more frequent between common SNPs than between any hypermutable element and surrounding SNPs.

D. Relating hypermutable locus-SNP correlations with SNP-SNP correlations

To compare the mean r^2 values for each scenario used, we plot all of the mean r^2 values for all simulations together (Figure 5). This plot demonstrates the relatively high mean r^2 values for common SNPs (peaking just below 0.15), and a lower mean r^2 values for loci with a mutation rate of 10^{-5} . Additionally, loci with a mutation rate of 10^{-4} provide an interesting comparison to the analysis using all SNPs. In close proximity, the mean r^2 measured on all SNPs is higher than that for loci with a mutation rate of 10^{-4} , but the correlation decays with distance much more rapidly for the SNPs. At a distance of 4000 bp the mean r^2 is nearly zero for all SNPs, but it remains above 0.1 at 4000 bp for hypermutable loci with mutation rates of 10^{-4} and 10^{-5} .

To further investigate these simulation results, we examine the locus with the largest r^2 found in each simulation, 2000 simulations per scenario. The maximum r^2 that occurs in an individual population is of interest because GWAS associations typically focus on SNPs with the lowest p-values. The scatter plot of the maximum per-simulation r^2 for a central hypermutable locus (Figure 6 (a)) demonstrates that SNPs with the strongest associations are more centralized in the simulations using lower mutation rates than in those using higher mutation rates. There is almost no localization in the simulations with $\mu = 10^{-3}$ (Figure 6 (c)). Furthermore, the maximum r^2 values under the mutation rate of 10^{-3} are always small; the largest maximum r^2 was only 0.202.

When the central locus is a common SNP, the maximum r^2 values are often near one (Figure 6 (b)). When the central SNP is rare, the maximum r^2 for the simulation is usually either very low or near one. Rare SNPs often have no association with surrounding loci, but occasionally a rare central SNP will be in perfect LD with another rare SNP, and this surrounding SNP in perfect LD is sometimes at a great distance. The maximum r^2

for common central SNPs is often relatively large and localized to the central region (Figure 6 (d)).

2. Discussion

A. Comparing results from the approximation with simulations

The approximation made by Ohta and Kimura (1969), $E(r^2) \approx 1/[4N(c+k)]$, provides a useful way to think about how mutation rates are related to linkage: the effects of mutation are similar to the effects of recombination, breaking linkage disequilibrium between loci. Although this approximation is only accurate when $N(c+k)$ is large, one can nevertheless use it to build intuition about how mutation reduces correlations between loci. A forward-backward mutation rate of 10^{-3} acts like a genetic distance of 0.002 M, about 200kb ($k \approx 2\mu = 0.002$, corresponding to $c=0.002$). Loci at a distance of 200kb are essentially unlinked. Therefore, even SNPs in close proximity to a hypermutable element with such a high mutation rate will be unlinked. This simple approximation makes it clear that SNPs do not tag variation caused by the most hypermutable loci in the genome. Furthermore, the simulations demonstrate that the Ohta and Kimura (1969) approximation over-estimates $E(r^2)$. When a site mutates rapidly, almost all of its correlation with surrounding loci is lost.

The approximation breaks down when $N(c+k)$ is smaller than one (Ohta and Kimura 1969), which is the case for most of the scenarios examined here. In these scenarios, the ratio of expectations in (4), ρ^2 is a poor approximation for the expectation of the ratio given in (3). The only scenario in which $N(c+k)$ is larger than one is when the mutation rate is 10^{-3} (Figure 2 (c)). Oddly, this is also the only scenario in which empirical ρ^2 does not appear to match the analytical approximation ρ^2 of equation (5).

Therefore, although the approximation made by Ohta and Kimura (1969) can be helpful for understanding how mutation rates relate to recombination distance, simulations are required to estimate the mean r^2 values for hypermutable elements with mutation rates larger than 10^{-3} . For investigating these mutation rates, neither decreasing the population size nor increasing genetic distance would increase the accuracy or utility of the approximation. The diffusion approximation breaks down as population sizes decrease. Furthermore, our interest here is to understand how SNPs can tag nearby hypermutable elements, and examining SNPs that are a great distance to a hypermutable element provides limited utility because a tiny r^2 is expected across large genetic distances. Thus the approximation ρ^2 has many limitations when studying hypermutable elements.

The simulation results provide useful insight into how SNPs correlate with hypermutable elements. For most hypermutable elements, the mean r^2 values with nearby SNPs are small, especially in comparison to common SNP-SNP associations (Figure 5). However, for hypermutable elements with mutation rates of 10^{-5} not all of the correlation is lost. The mean r^2 value for mutation rates of 10^{-5} is approximately half that of common SNP-SNP associations (Figure 5). Furthermore, for a mutation rate of 10^{-5} the top 5th percentile of r^2 values are all above 0.3 when the surrounding SNPs have an MAF above

0.2 (Figure 3 (c)). Stronger associations exist between common SNPs and other common SNPs (Figure 4 (c)), but the scenario with mutation rates of 10^{-5} is somewhat comparable.

Rare SNPs are known to have a small r^2 value with other SNPs (Sun *et al.* 2011), and rare SNPs are a potential explanation for missing heritability (Manolio *et al.* 2009) and still-missing heritability (Wray *et al.* 2014; Witte *et al.* 2014). The simulations indicate that rare SNPs have a low mean r^2 with other SNPs, comparable to hypermutable elements with mutation rates of 10^{-4} or smaller. However, the mean r^2 diminishes across genetic distance faster for SNPs than for hypermutable loci (Figure 5). This suggests that although hypermutable elements may behave similarly to rare SNPs, associations with hypermutable elements may show weaker localization. This delocalization spreads associations with hypermutable loci around the genome. Therefore, methods that use all SNPs together to measure overall genetic effects, such as GCTA (Yang *et al.* 2011), may be able to recover information about causal hypermutable loci.

B. Implications for GWAS

Hypermutable tandem repeat loci may be partially responsible for missing heritability (Hannan 2010; Press *et al.* 2014) and also still-missing heritability. The results presented here suggest that loci with high mutation rates are not well tagged by SNPs, and therefore much of the heritable variation caused by such loci will not have been captured in modern GWAS analyses. Scientists have just recently begun to estimate the mutation rates of hypermutable elements in the human genome (Whittaker *et al.* 2003; Kelkar *et al.* 2008; Sun *et al.* 2012; Legendre *et al.* 2007), and a database of known tandem repeat variants has recently been developed (Willems *et al.* 2014). As more tandem repeat variants are cataloged, understanding how these variants can be tagged by SNPs will allow researchers to measure their relative contributions to phenotypes.

When a tandem repeat has a lower mutation rate (such as G/C rich microsatellites or minisatellites), studies have shown that SNPs can be linked to disease repeat alleles (Wray *et al.* 2009; Laaksovirta *et al.* 2010; Mok *et al.* 2012; Majounie *et al.* 2012; Pliner *et al.* 2014), and our results corroborate this. For the C9orf72 repeat expansion there appears to have been a single repeat expansion in the European population, with nearby SNPs in strong linkage disequilibrium (Mok *et al.* 2012). This finding, along with the analyses here, suggests that other GWAS could be picking up phenotypic variation caused by tandem repeats. However, because SNPs are the main focus of contemporary genetics research, tandem repeats are often overlooked as potential causal variants. Furthermore, due to the limitations of PCR and next-gen sequencing technologies, tandem repeats are often difficult to genotype or sequence (Treangen and Salzberg 2012; Loomis *et al.* 2013; Press *et al.* 2014; Gymrek *et al.* 2012; Brahmachary *et al.* 2014; Krsticevic *et al.* 2015; Ummat and Bashir 2014; Doi *et al.* 2014). Consequently, researchers could easily miss a causal tandem repeat variant while investigating a GWAS signal using DNA sequencing.

An important consideration when investigating a GWAS signal is the distance between the SNP with the lowest p-value and the variant(s) driving the association. The position of the lowest p-value SNP is often used to link a gene with a phenotype. Our results suggest that the top SNP associations

are far less localized for hypermutable elements, with almost no localization for elements with a mutation rate of 10^{-3} (Figure 6 (c)). Therefore, if a hypermutable element is causing a SNP association, the strongest SNP association may occur at a great distance from the causal element. Associations with hypermutable elements are also spread across a larger region (Figure 5), providing an association signature that may be noticeably distinct from other types of associations.

Finally, because traits can be influenced by hypermutable elements and/or low frequency variants, SNP data alone cannot be used to exclude a gene or region of the genome as causal. If a gene is affected by hypermutable elements and/or rare variants, then SNPs will often fail to find an association. Regions or genes that contain potentially functional hypermutable elements require further genotyping of these elements before they can be totally excluded as potentially impacting a trait. Furthermore, many sequencing technologies have a limited ability to genotype some tandem repeat variants (Treangen and Salzberg 2012; Loomis *et al.* 2013; Press *et al.* 2014; Krsticevic *et al.* 2015; Doi *et al.* 2014), so our results apply to any data that is limited to SNPs. Recent advances in sequencing technology (Loomis *et al.* 2013; Ummat and Bashir 2014; Krsticevic *et al.* 2015) and tandem repeat genotyping (Gymrek *et al.* 2012; Carlson *et al.* 2015; Brahmachary *et al.* 2014; Doi *et al.* 2014) provide hope that some hypermutable elements will be included in future studies of genetic heritability and genetic disease. Nevertheless, some of the missing heritability cause by hypermutable elements may remain missing, at least for the near future.

C. Limitations and potential extensions

This study only use two possible states at each locus, and the forward and backward mutations are equal. This simplifies both the analytical approach as well as the simulations, and can be used as a simple model of tandem repeat evolution. Tandem repeats often have more than two states, but diseases caused by tandem repeats are often caused by expansion (Hannan 2010; Gemayel *et al.* 2010). Therefore, tandem repeats can sometimes fit into a two-allele model as was done here (short versus long). However, transitions between a short allele and a long allele depend on the repeat length (Kelkar *et al.* 2008), and thus forward and backward mutation rates are not necessarily equivalent. A step-wise mutation model, allowing multiple allele sizes at the hypermutable locus and binning them as short or long, may provide a more accurate model of tandem repeat diseases. These more complicated models are likely to return similar results because empirical data indicates that small r^2 values are found between SNPs and tandem repeat loci, whether they are bi-allelic or multi-allelic (Willems *et al.* 2014).

In addition, the use of a stable population with an effective size of 10,000 without population history may further limit the direct application of these results. The results from smaller population sizes might drastically change because the diffusion approximation does not work well for small effective population sizes. In addition, complicated population histories may change these results in unexpected ways, especially because tandem repeats and SNPs provide different information about population histories (Payseur and Jing 2009). Future simulations could address these possibilities.

Equation (6) suggests that increasing the population size will result in an approximately harmonic decrease in the mean r^2 . Therefore, one can expect the mean r^2 from an effective population size of 20,000 to be approximately half of the mean r^2 found here with an effective population size of 10,000. Extrapolating the results presented here to smaller population sizes would not be as straightforward. Due to the aforementioned effect that small populations have on the accuracy of the diffusion approximation, estimating how these results would change if one used a smaller population size is not as simple as applying a linear transformation.

D. Summary and conclusion

As shown by [Ohta and Kimura \(1969\)](#), mutation and recombination act in a similar fashion to break up linkage between loci. The magnitude of the mutation rate can be approximately equated to recombination distance in Morgans. However, this approximation only holds when the mutation rates are high and/or population sizes are large. With lower mutation rates the approximation breaks down and simulations must be used to estimate the expected linkage between loci.

The simulations reported here suggest that the variation caused by some hypermutable elements can be captured using SNPs. At mutation rates of 10^{-5} or smaller the associations between hypermutable loci and SNPs is comparable to, although lower than, common SNP - common SNP associations. On the other hand, the correlations between SNPs and loci with mutation rates of 10^{-4} and 10^{-3} are relatively low, and therefore variation caused by loci with these mutation rates are likely to show only weak association with SNPs of any MAF.

Heritable variation can be caused by genetic loci with a range of mutation rates ([Rando and Verstrepen 2007](#)). Hypermutable loci can remain highly polymorphic in a population, and they may be important causes of human disease and heritability of complex traits. Common SNP variants are currently inexpensive and widely used to search for genes that contribute to heritable variation. Unfortunately, many hypermutable loci will have poor linkage with SNPs, and therefore these loci will be unlikely to be uncovered using SNP GWAS methods. Direct genotyping will be necessary to uncover the effects that many hypermutable loci have on genetic variation. We hope that this work will help researchers investigating the sources of human diseases and heritable traits.

Figures

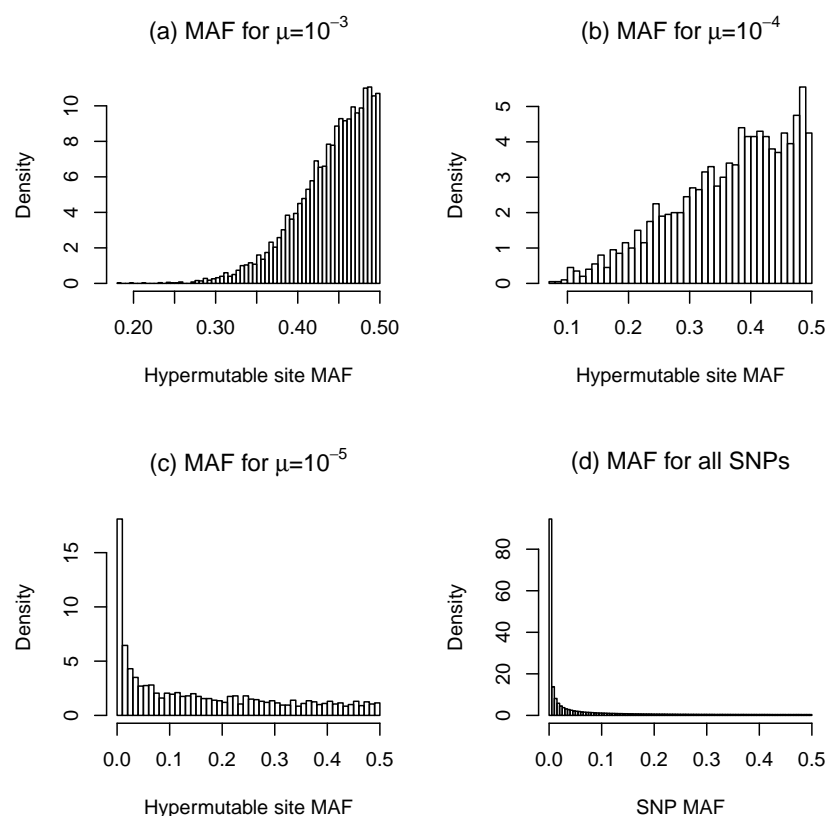


Figure 1 Histograms of allele frequencies from the simulations. The minor allele frequencies for bi-allelic hypermutable sites with mutation rates of 10^{-3} (a) 10^{-4} (b) and 10^{-5} (c) are shown. Only simulations with non-zero allele frequencies were used. Plot (d) shows a histogram of minor allele frequencies for SNPs in the simulation.

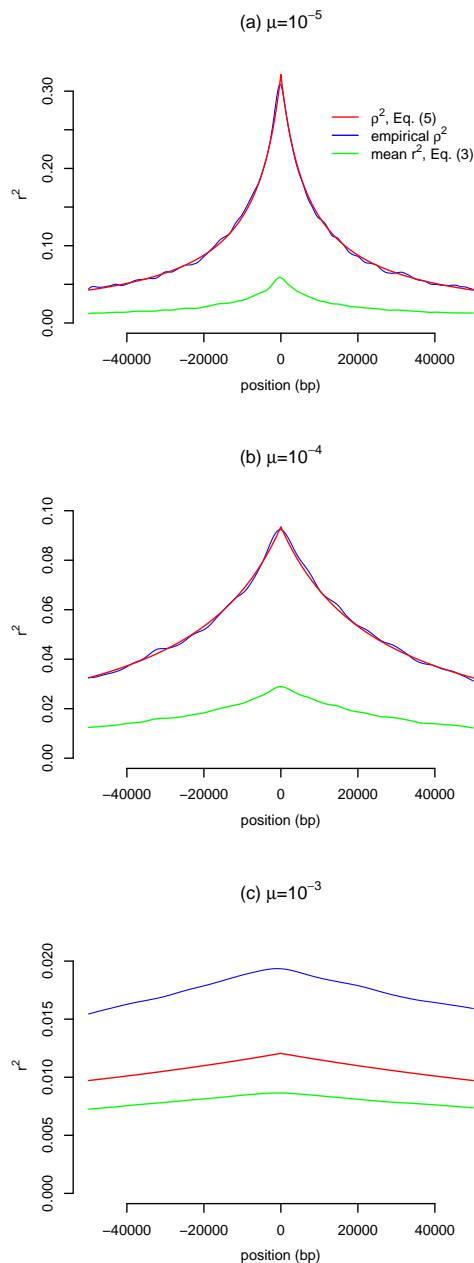


Figure 2 Plots comparing mean r^2 from simulations (green), its approximation, ρ^2 (red), and the empirical value of ρ^2 (blue). The hypermutable locus is central (position 0), and r^2 values were calculated between the central hypermutable element and surrounding SNPs. Results for simulations using hypermutable mutation rates of 10^{-3} (a), 10^{-4} (b), and 10^{-5} (c) are shown. The values of ρ^2 are far greater than the mean r^2 , with the greatest difference found for low mutation rates. The values were calculated for bins of 100 base-pairs, and a line was drawn between these binned values using LOESS smoothing. Note the change in scale on the vertical axes between plots of different mutation rates.

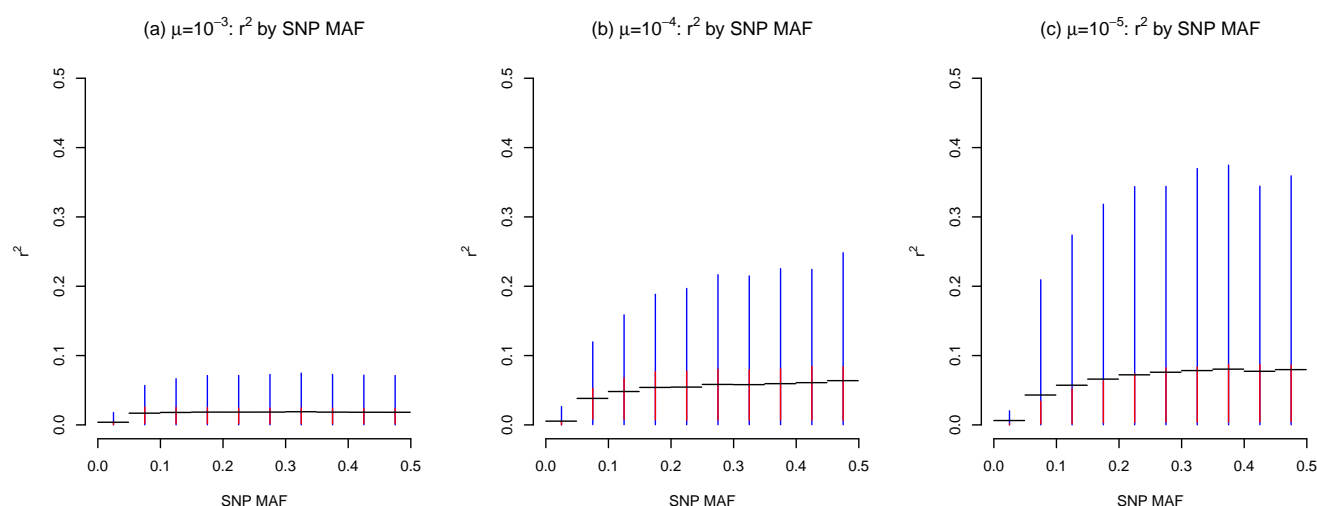


Figure 3 Mean r^2 values between the hypermutable locus and SNPs with varying MAF. The mean r^2 values are represented by the horizontal black line. The top (bottom) of the vertical red line represents the 75th (25th) percentile, and the top (bottom) of the blue lines represents the 95th (5th) percentile. Results for simulations using a hypermutable mutation rate of 10^{-3} (a), 10^{-4} (b), and 10^{-5} (c) are shown.

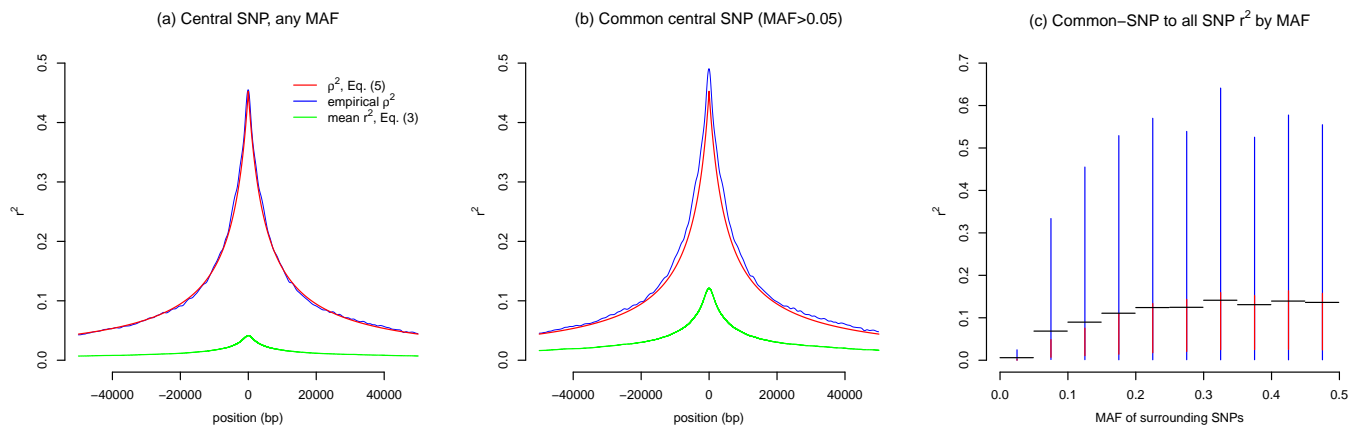


Figure 4 The r^2 between a central SNP and surrounding SNPs. (a) Mean r^2 values for SNP-SNP pairs, using a central SNP with any MAF (green). Also the analytical approximation for r^2 (ρ^2 , red), and empirical ρ^2 (blue). (b) Same as in (a) but for central SNPs with an MAF above 0.05, i.e. common central SNPs. (c) Distribution of r^2 for comparisons between a central SNP with MAF above 0.05 and surrounding SNPs binned by their MAF. The mean r^2 values are represented by the horizontal black line. The top (bottom) of the vertical red line represents the 75th (25th) percentiles, and the top (bottom) of the blue lines represent the 95th (5th) percentiles.

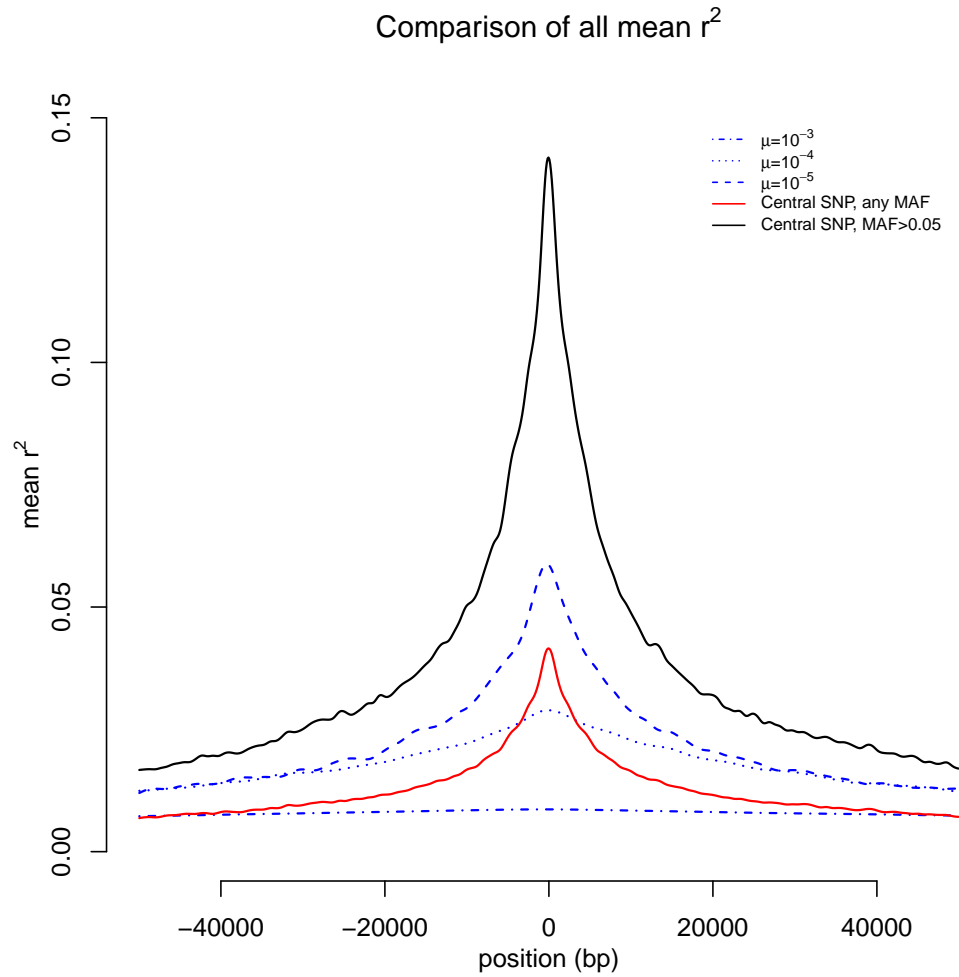


Figure 5 The mean r^2 between surrounding SNPs, with any MAF, and a central variant with these classes: mutation rates of 10^{-3} , 10^{-4} , and 10^{-5} , as well as a central SNP with any MAF and also a central common SNP (MAF>0.05).

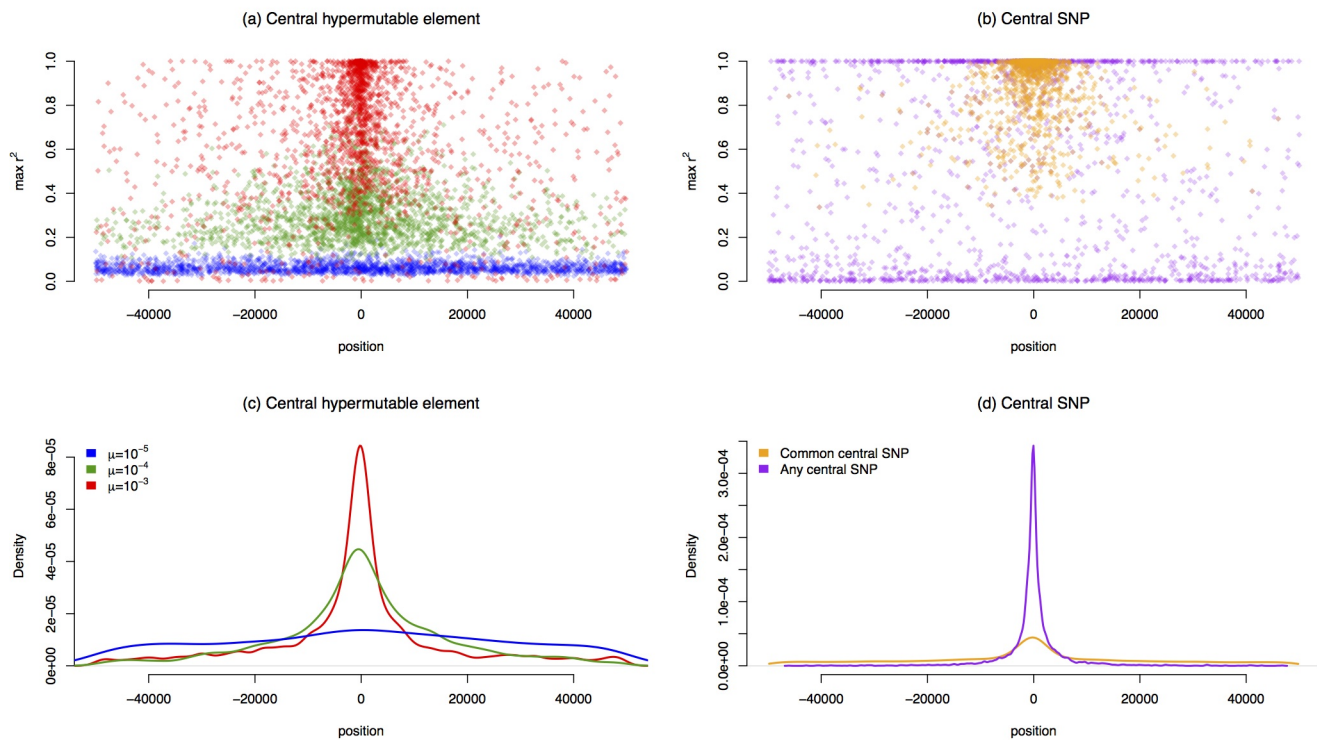


Figure 6 Characteristics of the maximum r^2 between a central element and surrounding SNPs from each individual simulation, 2000 in total. (a) Scatterplot of the maximum r^2 against the position relative to a central hypermutable element. Colors indicate mutation rate of the hypermutable element. (b) Scatterplot of the maximum r^2 against the position relative to a central SNP (i.e., a central locus with normal mutation rate). Colors indicate MAF of the central SNP (common or unconstrained). (c) Density of the position of the locus with maximum r^2 , relative to a central hypermutable element. (d) Density of the position of the locus with maximum r^2 , relative to a central SNP.

Literature Cited

- Baptiste, B. A., G. Ananda, N. Strubczewski, A. Lutzkanin, S. J. Khoo, A. Srikanth, N. Kim, K. D. Makova, M. M. Krasilnikova, and K. A. Eckert, 2013 Mature microsatellites: mechanisms underlying dinucleotide microsatellite mutational biases in human cells. *G3 (Bethesda)* 3: 451–463.
- Brahmachary, M., A. Guilmatre, J. Quilez, D. Hasson, C. Borel, P. Warburton, and A. J. Sharp, 2014 Digital genotyping of macrosatellites and multicopy genes reveals novel biological functions associated with copy number variation of large tandem repeats. *PLoS Genet.* 10: e1004418.
- Buschiazio, E. and N. J. Gemmell, 2010 Conservation of human microsatellites across 450 million years of evolution. *Genome Biol Evol* 2: 153–165.
- Carlson, K. D., P. H. Sudmant, M. O. Press, E. E. Eichler, J. Shendure, and C. Queitsch, 2015 Mipstr: a method for multiplex genotyping of germline and somatic str variation across many individuals. *Genome Research* .
- DeJesus-Hernandez, M., I. R. Mackenzie, B. F. Boeve, A. L. Boxer, M. Baker, N. J. Rutherford, A. M. Nicholson, N. A. Finch, H. Flynn, J. Adamson, N. Kouri, A. Wojtas, P. Sengdy, G. Y. Hsiung, A. Karydas, W. W. Seeley, K. A. Josephs, G. Coppola, D. H. Geschwind, Z. K. Wszolek, H. Feldman, D. S. Knopman, R. C. Petersen, B. L. Miller, D. W. Dickson, K. B. Boylan, N. R. Graff-Radford, and R. Rademakers, 2011 Expanded GGGGCC hexanucleotide repeat in noncoding region of C9ORF72 causes chromosome 9p-linked FTD and ALS. *Neuron* 72: 245–256.
- Doi, K., T. Monjo, P. H. Hoang, J. Yoshimura, H. Yurino, J. Mitsui, H. Ishiura, Y. Takahashi, Y. Ichikawa, J. Goto, S. Tsuji, and S. Morishita, 2014 Rapid detection of expanded short tandem repeats in personal genomics using hybrid sequencing. *Bioinformatics* 30: 815–822.
- Eberle, M. A., P. C. Ng, K. Kuhn, L. Zhou, D. A. Peiffer, L. Galver, K. A. Viaud-Martinez, C. T. Lawley, K. L. Gunderson, R. Shen, and S. S. Murray, 2007 Power to detect risk alleles using genome-wide tag SNP panels. *PLoS Genet.* 3: 1827–1837.
- Ellegren, H., 2004 Microsatellites: simple sequences with complex evolution. *Nature Reviews Genetics* 5: 435–445.
- Excoffier, L. and M. Foll, 2011 fastsimcoal: a continuous-time coalescent simulator of genomic diversity under arbitrarily complex evolutionary scenarios. *Bioinformatics* 27: 1332–1334.
- Gemayel, R., M. D. Vences, M. Legendre, and K. J. Verstrepen, 2010 Variable tandem repeats accelerate evolution of coding and regulatory sequences. *Annu. Rev. Genet.* 44: 445–477.
- Gymrek, M., D. Golan, S. Rosset, and Y. Erlich, 2012 lobstr: A short tandem repeat profiler for personal genomes. *Genome Research* 22: 1154–1162.
- Hannan, A., 2010 Tandem repeat polymorphisms: modulators of disease susceptibility and candidates for ‘missing heritability’. *Trends in Genetics* 26: 59–65.
- Hedrick, P. W., 1987 Gametic disequilibrium measures: proceed with caution. *Genetics* 117: 331–341.
- Hill, W. G. and A. Robertson, 1968 Linkage disequilibrium in finite populations. *Theor. Appl. Genet.* 38: 226–231.
- Jeffreys, A. J., D. L. Neil, and R. Neumann, 1998 Repeat instability at human minisatellites arising from meiotic recombination. *EMBO J.* 17: 4147–4157.
- Kelkar, Y. D., S. Tyekucheva, F. Chiaromonte, and K. D. Makova, 2008 The genome-wide determinants of human and chimpanzee microsatellite evolution. *Genome Res.* 18: 30–38.
- Krsticevic, F. J., C. G. Schrago, and A. B. Carvalho, 2015 Long-Read Single Molecule Sequencing To Resolve Tandem Gene Copies: The Mst77Y Region on the Drosophila melanogaster Y Chromosome. *G3 (Bethesda)* .
- Laaksovirta, H., T. Peuralinna, J. C. Schymick, S. W. Scholz, S. L. Lai, L. Myllykangas, R. Sulkava, L. Jansson, D. G. Hernandez, J. R. Gibbs, M. A. Nalls, D. Heckerman, P. J. Tienari, and B. J. Traynor, 2010 Chromosome 9p21 in amyotrophic lateral sclerosis in Finland: a genome-wide association study. *Lancet Neurol* 9: 978–985.
- Legendre, M., N. Pochet, T. Pak, and K. J. Verstrepen, 2007 Sequence-based estimation of minisatellite and microsatellite repeat variability. *Genome Res.* 17: 1787–1796.
- Lesch, K. P., D. Bengel, A. Heils, S. Z. Sabol, B. D. Greenberg, S. Petri, J. Benjamin, C. R. Muller, D. H. Hamer, and D. L. Murphy, 1996 Association of anxiety-related traits with a polymorphism in the serotonin transporter gene regulatory region. *Science* 274: 1527–1531.
- Loomis, E. W., J. S. Eid, P. Peluso, J. Yin, L. Hickey, D. Rank, S. McCalmon, R. J. Hagerman, F. Tassone, and P. J. Hagerman, 2013 Sequencing the unsequenceable: expanded CGG-repeat alleles of the fragile X gene. *Genome Res.* 23: 121–128.
- MacDonald, M. E., C. M. Ambrose, M. P. Duyao, R. H. Myers, C. Lin, L. Srinidhi, G. Barnes, S. A. Taylor, M. James, N. Groot, H. MacFarlane, B. Jenkins, M. A. Anderson, N. S. Wexler, J. F. Gusella, G. P. Bates, S. Baxendale, H. Hummerich, S. Kirby, M. North, S. Youngman, R. Mott, G. Zehetner, Z. Sedlacek, A. Poustka, A.-M. Frischauf, H. Lehrach, A. J. Buckler, D. Church, L. Doucette-Stamm, M. C. O’Donovan, L. Ribar-Ramirez, M. Shah, V. P. Stanton, S. A. Strobel, K. M. Draths, J. L. Wales, P. Dervan, D. E. Housman, M. Altherr, R. Shi-ang, L. Thompson, T. Fielder, J. J. Wasmuth, D. Tagle, J. Valdes, L. Elmer, M. Allard, L. Castilla, M. Swaroop, K. Blanchard, F. S. Collins, R. Snell, T. Holloway, K. Gillespie, N. Datson, D. Shaw, and P. S. Harper, 1993 A novel gene containing a trinucleotide repeat that is expanded and unstable on Huntington’s disease chromosomes. The Huntington’s Disease Collaborative Research Group. *Cell* 72: 971–983.
- Maher, B., 2008 Personal genomes: The case of the missing heritability. *Nature* 456: 18–21.
- Majounie, E., A. E. Renton, K. Mok, E. G. Dopper, A. Waite, S. Rollinson, A. Chio, G. Restagno, N. Nicolaou, J. Simon-Sanchez, J. C. van Swieten, Y. Abramzon, J. O. Johnson, M. Sendtner, R. Pamphlett, R. W. Orrell, S. Mead, K. C. Siddle, H. Houlden, J. D. Rohrer, K. E. Morrison, H. Pall, K. Talbot, O. Ansorge, D. G. Hernandez, S. Arepalli, M. Sabatelli, G. Mora, M. Corbo, F. Giannini, A. Calvo, E. Englund, G. Borghero, G. L. Floris, A. M. Remes, H. Laaksovirta, L. McCluskey, J. Q. Trojanowski, V. M. V. Deerlin, G. D. Schellenberg, M. A. Nalls, V. E. Drory, C.-S. Lu, T.-H. Yeh, H. Ishiura, Y. Takahashi, S. Tsuji, I. L. Ber, A. Brice, C. Drepper, N. Williams, J. Kirby, P. Shaw, J. Hardy, P. J. Tienari, P. Heutink, H. R. Morris, S. Pickering-Brown, and B. J. Traynor, 2012 Frequency of the c9orf72 hexanucleotide repeat expansion in patients with amyotrophic lateral sclerosis and frontotemporal dementia: a cross-sectional study. *The Lancet Neurology* 11: 323–330.
- Manolio, T. A., F. S. Collins, N. J. Cox, D. B. Goldstein, L. A. Hindorf, D. J. Hunter, M. I. McCarthy, E. M. Ramos, L. R. Cardon, A. Chakravarti, J. H. Cho, A. E. Guttacher, A. Kong, L. Kruglyak, E. Mardis, C. N. Rotimi, M. Slatkin, D. Valle, A. S. Whittemore, M. Boehnke, A. G. Clark, E. E. Eichler, G. Gibson, J. L. Haines, T. F. Mackay, S. A. McCarroll, and P. M. Visscher, 2009 Finding the missing heritability of complex diseases. *Nature* 461: 747–753.

- Mok, K., B. J. Traynor, J. Schymick, P. J. Tienari, H. Laaksovirta, T. Peuralinna, L. Myllykangas, A. Chio, A. Shatunov, B. F. Boeve, A. L. Boxer, M. DeJesus-Hernandez, I. R. Mackenzie, A. Waite, N. Williams, H. R. Morris, J. Simon-Sanchez, J. C. van Swieten, P. Heutink, G. Restagno, G. Mora, K. E. Morrison, P. J. Shaw, P. S. Rollinson, A. Al-Chalabi, R. Rademakers, S. Pickering-Brown, R. W. Orrell, M. A. Nalls, and J. Hardy, 2012 Chromosome 9 ALS and FTD locus is probably derived from a single founder. *Neurobiol. Aging* **33**: 3–8.
- Ohta, T. and M. Kimura, 1969 Linkage disequilibrium at steady state determined by random genetic drift and recurrent mutation. *Genetics* **63**: 229–238.
- Payseur, B. A. and P. Jing, 2009 A genomewide comparison of population structure at STRPs and nearby SNPs in humans. *Mol. Biol. Evol.* **26**: 1369–1377.
- Payseur, B. A., M. Place, and J. L. Weber, 2008 Linkage disequilibrium between STRPs and SNPs across the human genome. *Am. J. Hum. Genet.* **82**: 1039–1050.
- Pliner, H. A., D. M. Mann, and B. J. Traynor, 2014 Searching for Grendel: origin and global spread of the C9ORF72 repeat expansion. *Acta Neuropathol.* **127**: 391–396.
- Press, M. O., K. D. Carlson, and C. Queitsch, 2014 The overdue promise of short tandem repeat variation for heritability. *Trends Genet.* **30**: 504–512.
- Rando, O. J. and K. J. Verstrepen, 2007 Timescales of genetic and epigenetic inheritance. *Cell* **128**: 655–668.
- Sawaya, S., A. Bagshaw, E. Buschiazzi, P. Kumar, S. Chowdhury, M. A. Black, and N. Gemmell, 2013 Microsatellite tandem repeats are abundant in human promoters and are associated with regulatory elements. *PLoS ONE* **8**: e54710.
- Sun, J. X., A. Helgason, G. Masson, S. S. Ebenesersdottir, H. Li, S. Mallick, S. Gnerre, N. Patterson, A. Kong, D. Reich, and K. Stefansson, 2012 A direct characterization of human mutation based on microsatellites. *Nat. Genet.* **44**: 1161–1165.
- Sun, X., J. Namkung, X. Zhu, and R. C. Elston, 2011 Capability of common SNPs to tag rare variants. *BMC Proc* **5 Suppl 9**: S88.
- Treangen, T. J. and S. L. Salzberg, 2012 Repetitive DNA and next-generation sequencing: computational challenges and solutions. *Nat. Rev. Genet.* **13**: 36–46.
- Ummat, A. and A. Bashir, 2014 Resolving complex tandem repeats with long reads. *Bioinformatics* **30**: 3491–3498.
- VanLiere, J. M. and N. A. Rosenberg, 2008 Mathematical properties of the r^2 measure of linkage disequilibrium. *Theor Popul Biol* **74**: 130–137.
- Verkerk, A. J., M. Pieretti, J. S. Sutcliffe, Y. H. Fu, D. P. Kuhl, A. Pizzuti, O. Reiner, S. Richards, M. F. Victoria, and F. P. Zhang, 1991 Identification of a gene (FMR-1) containing a CGG repeat coincident with a breakpoint cluster region exhibiting length variation in fragile X syndrome. *Cell* **65**: 905–914.
- Vinces, M. D., M. Legendre, M. Caldara, M. Hagihara, and K. J. Verstrepen, 2009 Unstable tandem repeats in promoters confer transcriptional evolvability. *Science* **324**: 1213–1216.
- Whittaker, J. C., R. M. Harbord, N. Boxall, I. Mackay, G. Dawson, and R. M. Sibly, 2003 Likelihood-based estimation of microsatellite mutation rates. *Genetics* **164**: 781–787.
- Willems, T. F., M. Gymrek, G. Highnam, T. G. Project, D. Mittelman, and Y. Erlich, 2014 The landscape of human str variation. *Genome Research*.
- Witte, J. S., P. M. Visscher, and N. R. Wray, 2014 The contribution of genetic variants to disease depends on the ruler. *Nat. Rev. Genet.* **15**: 765–776.
- Wray, N. R., 2005 Allele frequencies and the r^2 measure of linkage disequilibrium: impact on design and interpretation of association studies. *Twin Res Hum Genet* **8**: 87–94.
- Wray, N. R., M. R. James, S. D. Gordon, T. Dumenil, L. Ryan, W. L. Coventry, D. J. Statham, M. L. Pergadia, P. A. Madden, A. C. Heath, G. W. Montgomery, and N. G. Martin, 2009 Accurate, large-scale genotyping of 5httlpr and flanking single nucleotide polymorphisms in an association study of depression, anxiety, and personality measures. *Biological Psychiatry* **66**: 468 – 476, Medical Consequences and Contributions to Depression.
- Wray, N. R., S. H. Lee, D. Mehta, A. A. Vinkhuyzen, F. Dudbridge, and C. M. Middeldorp, 2014 Research review: Polygenic methods and their application to psychiatric traits. *J Child Psychol Psychiatry* **55**: 1068–1087.
- Xu, X., M. Peng, and Z. Fang, 2000 The direction of microsatellite mutations is dependent upon allele length. *Nat. Genet.* **24**: 396–399.
- Yang, J., S. H. Lee, M. E. Goddard, and P. M. Visscher, 2011 GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**: 76–82.