

# Extensive *de novo* mutation rate variation between individuals and across the genome of *Chlamydomonas reinhardtii*

Rob W. Ness<sup>\*</sup>, Andrew D. Morgan, Radhakrishnan B. Vasanthakrishnan, Nick Colegrave<sup>1</sup>, and Peter D. Keightley<sup>1</sup>

Institute of Evolutionary Biology, University of Edinburgh  
Ashworth Labs, King's Buildings, West Mains Road  
Edinburgh EH9 3JT

<sup>\*</sup>Corresponding author: email: [rob.ness@ed.ac.uk](mailto:rob.ness@ed.ac.uk), Telephone +44(0)131 650 7334

<sup>1</sup>These authors contributed equally to this work.

**Short title:** Mutation rate variation in *Chlamydomonas*

**Key words:** Mutation rate, Spontaneous mutation, *Chlamydomonas reinhardtii*

## 1 **Abstract**

2 Describing the process of spontaneous mutation is fundamental for understanding the genetic basis of  
3 disease, the threat posed by declining population size in conservation biology, and in much  
4 evolutionary biology. However, directly studying spontaneous mutation is difficult because of the rarity  
5 of *de novo* mutations. Mutation accumulation (MA) experiments overcome this by allowing mutations  
6 to build up over many generations in the near absence of natural selection. In this study, we  
7 sequenced the genomes of 85 MA lines derived from six genetically diverse wild strains of the green  
8 alga *Chlamydomonas reinhardtii*. We identified 6,843 spontaneous mutations, more than any other  
9 study of spontaneous mutation. We observed seven-fold variation in the mutation rate among strains  
10 and that mutator genotypes arose, increasing the mutation rate dramatically in some replicates. We  
11 also found evidence for fine-scale heterogeneity in the mutation rate, driven largely by the sequence  
12 flanking mutated sites, and by clusters of multiple mutations at closely linked sites. There was little  
13 evidence, however, for mutation rate heterogeneity between chromosomes or over large genomic  
14 regions of 200Kbp. Using logistic regression, we generated a predictive model of the mutability of sites  
15 based on their genomic properties, including local GC content, gene expression level and local  
16 sequence context. Our model accurately predicted the average mutation rate and natural levels of  
17 genetic diversity of sites across the genome. Notably, trinucleotides vary 17-fold in rate between the  
18 most mutable and least mutable sites. Our results uncover a rich heterogeneity in the process of  
19 spontaneous mutation both among individuals and across the genome.

20

## 1 Introduction

2 Understanding the processes that generate new genetic variation from mutation is a key goal of  
3 genetics research. It is widely believed that the majority of new mutations that affect functional  
4 elements of the genome are deleterious. In humans, new mutations cause Mendelian genetic  
5 disorders, play a direct role in polygenic disease (e.g. Veltman and Brunner 2012), and are a major  
6 factor in cancers (e.g. Alexandrov et al. 2013a). New mutations also play a central role in evolutionary  
7 biology, since the variation that fuels adaptive evolution is ultimately derived from advantageous  
8 mutations. For example, the input of new variation from mutation is pivotal for theory to explain the  
9 evolution of recombination and sex (reviewed in Otto 2009).

10 If new mutations are harmful, theory predicts that the mutation rate should evolve towards zero,  
11 because individuals with higher mutations rates will suffer a greater mutational load. However, the  
12 mutation rate is always greater than zero in nature, ranging over seven orders of magnitude (reviewed  
13 by Drake 2006), and two main explanations have been proposed for this. One explanation is that there  
14 is a limit to the fidelity of DNA repair, due to a trade-off between the benefit of further reducing the  
15 mutation rate and the costs of increased fidelity (Kimura 1967). Alternatively, a 'selection-drift' barrier  
16 may constrain progress toward lower mutation rate when the selective advantage of further  
17 improvement becomes so small that new mutations decreasing the mutation rate are effectively  
18 neutral (Lynch 2010). Evidence for a selection-drift barrier comes from the negative correlation  
19 between the mutation rate per generation and effective population size ( $N_e$ ) (Sung et al. 2012).  
20 However, when mutation rate is expressed per cell division, there is much less variation between  
21 species and little relationship with  $N_e$ , consistent with the constraint on the fidelity of replication  
22 hypothesis. It is currently difficult to fully evaluate the support for these hypotheses, however, because  
23 studies of mutation are restricted to a small number of taxa, few genotypes per species and a limited  
24 number of mutation events.

25 Although there is clear evidence for variation between species, we know relatively little about the  
26 extent of mutation rate variation within species. Individuals with unusually high mutation rate have  
27 been isolated from natural populations of prokaryotes (Matic et al. 1997; Sundin and Weigand 2007),  
28 but no natural mutators have been found in eukaryotes. This discrepancy likely stems from the fact  
29 that prokaryotes are asexual whereas eukaryotes are predominantly sexual. Theory predicts that in an  
30 asexual population, a mutator allele can hitchhike to high frequency if it causes a beneficial allele on  
31 the same genetic background (Johnson 1999). In contrast, recombination in sexual populations  
32 uncouples a mutator from a linked beneficial allele, so the mutator allele is then expected to be  
33 selected against because of its association with linked deleterious mutations (reviewed by Drake et al.  
34 1998). Although a smaller amount of mutation rate variation is expected in sexual than asexual

1 species, mutations that alter the mutation rate are nevertheless expected to occur, and potentially  
2 provide the basis for mutation rate evolution. Mutation rate variation within a species may also reflect  
3 mutation-selection balance, whereby new deleterious alleles that alter the mutation rate continually  
4 arise and are purged by selection. In this scenario, intraspecific mutation rate variation will reflect the  
5 distribution of phenotypic effects of mutations that alter DNA repair and stability and the effectiveness  
6 of selection against them. In the largest study of spontaneous mutation in humans, there was little  
7 evidence for mutation rate variation among individuals after accounting for parental age (Kong et al.  
8 2012). Father's age was also an important factor explaining mutation rate variation in chimpanzees  
9 (Venn et al. 2014). Similarly, there was no evidence of mutation rate variation between two strains in  
10 both *Caenorhabditis elegans* and *C. briggsae* (Denver et al. 2012). There is evidence from *Drosophila*  
11 that individuals in poor condition have elevated mutation rates (Sharp and Agrawal 2012) and a  
12 separate study comparing two inbred lines revealed a 2.4-fold difference in the rate of mutation  
13 (Schridder et al. 2013). Moreover, two independent experiments in *Chlamydomonas reinhardtii*  
14 suggested that there is a 5-fold difference in the mutation rate between two natural strains (Ness et al.  
15 2012; Sung et al. 2012).

16 In addition to mutation rate variation within and between species, there is also evidence that mutation  
17 rate varies across the genome. Such heterogeneity is expected to alter the rate of evolution across the  
18 genome and to create variation in the susceptibility of genes or sites to deleterious or beneficial  
19 mutations. There is clear evidence for fine-scale variation in the rate of mutation. At the scale of  
20 individual sites, G:C positions tend to mutate at higher rates than A:T positions, and transitions from  
21 G:C→A:T are the most common change in a broad range of species (for example bacteria (Hershberg  
22 and Petrov 2010), animals (Kong et al. 2012; Schridder et al. 2013), fungi (Zhu et al. 2014) and plants  
23 (Ness et al. 2012)). Similarly, the bases surrounding a mutated site have a strong effect on mutability.  
24 For example, the high frequency of G:C→A:T transitions in mammals is driven by the deamination of  
25 methylated C<sub>p</sub>G sites (Ehrlich and Wang 1981). In general, the bases flanking a particular site,  
26 referred to as the 'sequence context', are one of the best predictors of mutation rate (Michaelson et al.  
27 2012; Neale et al. 2012; Samocha et al. 2014; Zhu et al. 2014). However, the underlying mechanisms  
28 and the consistency of the effect of sequence context on mutability across species is unknown.

29 At a broader genomic scale, evidence for mutation rate heterogeneity is weaker. Sequencing of MA  
30 lines in *S. cerevisiae* (Zhu et al. 2014) and *D. melanogaster* (Schridder et al. 2013) found no evidence  
31 of mutation rate variation between chromosomes. Although there is evidence that mutation rate  
32 increases as a function of replication timing (Stamatoyannopoulos et al. 2009; Lang and Murray 2011),  
33 this finding has not been supported by direct estimates of mutation rate (Samocha et al. 2014; Zhu et  
34 al. 2014). A variety of other genomic properties have been linked to increased susceptibility to

1 mutation, including transcription level, nucleosome occupancy, DNase hypersensitivity and  
2 recombination rate (e.g. Michaelson et al. 2012). If these factors strongly influence mutation and  
3 generate variation between sites or large scale patterns of mutation rate variation, it is important to  
4 quantify their effects, in order to facilitate better predictive models of DNA sequence evolution.

5 Detailed investigations of the process of spontaneous mutation and the extent of mutation rate  
6 variation are limited. This is because spontaneous mutations are very rare, constraining direct  
7 observation of sufficient numbers of mutations to infer the underlying biology. Sequencing of parents  
8 and their offspring is an increasingly common method for directly identifying de novo mutations (e.g.  
9 Keightley et al. 2014a; Keightley et al. 2014b). Although this approach has advantages, it is currently  
10 very expensive to sequence enough families to observe large numbers of mutations and has therefore  
11 only been applied on a large scale in humans (Kong et al. 2012). Another approach is to maintain  
12 experimental populations for many generations under minimal natural selection to allow mutations to  
13 accumulate regardless of their fitness consequences. Increasing the strength of genetic drift by  
14 bottlenecking the population each generation allows random, unbiased accumulation of all but the  
15 strongest deleterious mutations. These ‘mutation accumulation’ (MA) experiments have been used in  
16 a variety of species to investigate the phenotypic effects of new mutations (reviewed in Halligan and  
17 Keightley 2009) and are now being paired with whole genome sequencing to identify individual  
18 mutations. MA studies have generally been limited to sequencing a small number of genomes, and  
19 only two studies (Schridder et al. 2013; Zhu et al. 2014) have tested for heterogeneity in mutation rate  
20 across the genome, and no study has included more than two ancestral genotypes from a single  
21 species. In this study, we sequenced the genomes of 85 MA lines derived from six genetically diverse  
22 wild strains of the model green alga *C. reinhardtii*. We identified 6,843 mutations, seven-fold more  
23 than any previous MA study, and integrate this data with detailed annotation of genomic properties to  
24 investigate the process of spontaneous mutation with unprecedented detail. Specifically, we address  
25 the following fundamental questions (1) What is the relative frequency of different kinds of mutation,  
26 including the base spectrum and rate of insertion and deletion mutations? (2) What is the extent of  
27 mutation rate variation between individuals within a species? (3) Is there evidence of mutation rate  
28 heterogeneity across the genome and what genomic properties predict the rate of mutation at  
29 individual sites?

## 30 **Results**

31 We conducted a mutation accumulation experiment in six genetically diverse wild strains of *C.*  
32 *reinhardtii* that were chosen to broadly cover the geographic range of known *C. reinhardtii* samples in  
33 North America (Table 1). 15 replicate MA lines from each of the six ancestral strains were initiated for  
34 a total of 90 MA lines. 85 of the initial 90 MA lines survived to the end of the experiment. The mean

1 number of generations per MA line was 940 (range 403 to 1,130). DNA was extracted from each line  
2 and sequenced using Illumina whole genome sequencing allowing us to identify mutations in an  
3 average of 75.4Mbp per line (72.5% of genome, range 58.5-84.9Mbp; See Materials and Methods for  
4 details on mutation calling). In total, we identified 6,843 mutations, including 5,716 single nucleotide  
5 mutations (SNMs) and 1,127 short indels. To confirm these mutation calls, we Sanger sequenced a  
6 random sample of 192 mutations. 138 were successfully amplified and sequenced, 115 of 117 SNMs  
7 were confirmed and 19 of 21 indels were confirmed, resulting in an accuracy of 98.3% and 90.5% for  
8 SNMs and indels respectively.

### 9 **Mutation rate variation among genotypes.**

10 Including all MA lines, the total mutation rate was,  $\mu = 1.15 \times 10^{-9}$  muts/site\*generation, with SNM and  
11 indel mutation rates of  $\mu_{\text{SNM}} = 9.63 \times 10^{-10}$  and  $\mu_{\text{INDEL}} = 1.90 \times 10^{-10}$ , respectively. The mutation rate  
12 varied substantially among the MA replicates and between ancestral strains. Mutation rates of the  
13 individual MA lines ranged over nearly two orders of magnitude from  $\mu_{\text{CC-1952-MA4}} = 5.65 \times 10^{-11}$  to  $\mu_{\text{CC-}}$   
14  $_{2344\text{-MA1}} = 4.94 \times 10^{-9}$ . There was significant variation in the mean mutation rate among the strains ( $F_{1,5}$   
15  $= 30.96$ ,  $P < 0.0001$ , see Fig. 1). Post hoc Tukey tests showed strain CC-1373 had an average mutation  
16 rate significantly higher than all other strains ( $\mu = 28.1 \times 10^{-10}$ ,  $P$  0.01 to  $< 0.001$ ). This rate was nearly 7-  
17 fold higher than CC-1952 ( $\mu = 4.05 \times 10^{-10}$ ), which had the lowest mutation rate, and was significantly  
18 lower than CC-1373 ( $P < 0.001$ ), CC-2931 ( $\mu = 15.6 \times 10^{-10}$ ,  $P < 0.001$ ) and CC-2342 ( $\mu = 11.1 \times 10^{-10}$ ,  
19  $P < 0.01$ ). Within strains CC-2344 and CC-2931, there were individual MA lines with significantly higher  
20 mutation rates, 3.5 $\times$  and 8.0 $\times$  above their respective strain means ( $\mu$  estimates are outside the  
21 99.99% CI of their ancestral strain mutation rates,  $\mu_{\text{CC-2344-MA1}} = 56.9 \times 10^{-10}$ , CC-2344 CI = 2.6 -  $12.0 \times 10^{-}$   
22  $10$ ;  $\mu_{\text{CC-2931-MA5}} = 36.2 \times 10^{-10}$ , CC-2931 CI = 7.2 -  $20.0 \times 10^{-10}$ ).

### 23 **Indel mutations.**

24 There were significantly more short deletions (613) than insertions (514) ( $\chi^2 = 8.7$ ,  $P < 0.005$ ) and these  
25 tended to be larger (mean length -7.9 and +5.9, respectively, Mann-Whitney U test,  $W = 112604.5$ ,  $P <$   
26  $2.2 \times 10^{-16}$ ), but the difference was not significant. MA lines of strain CC-2931 had an unusually high  
27 number of indels (408) due to an abundance of 9bp deletions. 120 of 408 indels in CC-2931 were 9bp  
28 deletions compared to a mean of five 9bp deletions in each of the other strains. These deletions did  
29 not appear to have any shared sequence motif nor were they associated with coding exons, repetitive  
30 sequence or any genomic property that we could identify. After adjusting for the excess of 9bp  
31 deletions in CC-2931 by substituting the mean number of 9bp deletions found in the other strains,  
32 there were similar numbers of insertions and deletions, but deletions were still significantly longer  
33 ( $W = 100759.5$ ,  $P = 3.3 \times 10^{-9}$ ).

## 1 **Spatial heterogeneity.**

2 When mutation rate was measured in 200kbp sliding windows,  $\mu$  ranged from 0.0 to  $23.5 \times 10^{-10}$ . By  
3 comparing the distribution of mutation rates for each window with a simulation distribution, much of  
4 this variation could be accounted for as noise around the genome average mutation rate (KS test  $D =$   
5  $0.038$ ,  $P = 0.43$ ). In 1,000 simulations where mutation positions were randomized, the 95% confidence  
6 interval (CI) was  $\mu = 5.3 - 18.3 \times 10^{-10}$  compared to a 95% CI of  $\mu = 4.8 - 19.4 \times 10^{-10}$  in the observed  
7 data. 8% of 200kbp windows were above the 95th percentile of simulated mutation rates, suggesting a  
8 very slight excess of windows with a high mutation rate. Notably, the chloroplast genome had a  
9 mutation rate of  $\mu_{\text{cpDNA}} = 5.17 \times 10^{-9}$ , nearly  $4.5 \times$  the genome average.

10 We detected a significant deviation in the distribution of minimum intermutation distance compared to  
11 those expected under simulation (Fig. 2, K-S test:  $D = 0.048$ ,  $P = 4.5 \times 10^{-14}$ ). There was a large excess  
12 of mutations clustered very near to one another (<100bp apart) and most of this excess was caused  
13 by mutations at adjacent sites. Specifically, we expected zero adjacent mutations, but identified 55  
14 mutations where two adjacent sites were mutated, each of which was visually inspected in the  
15 Integrated Genomics Viewer, IGV (Thorvaldsdóttir et al. 2012). 27 of these clustered mutations  
16 occurred at CC (or GG) sites, and 25 of 27 mutated to AA/AT/TA/TT. We also found a number of  
17 indels where a short amount of sequence was replaced by an unrelated stretch of sequence. These  
18 complex indels are often reported by GATK HaplotypeCaller as a separate deletion and insertion  
19 rather than a single event. The excessive clustering occurred only within MA lines and when we  
20 limited our analysis to test for the presence of clustering of mutations found in different lines there was  
21 no evidence for this effect (K-S test  $D = 0.02$ ,  $P = 0.13$ ).

## 22 **Base composition.**

23 Treating the strand symmetrically we found a significantly non-random distribution of the six possible  
24 SNMs ( $\chi^2 = 1630.3$ ,  $P < 0.0001$ ; Fig. 3). Mutations occurring at C:G sites were  $4.2 \times$  more frequent than  
25 mutations at A:T sites, after correcting for genomic base composition, and this pattern was consistent  
26 across all MA lines and ancestral strains. Transitions from C:G  $\rightarrow$  T:A were over-represented nearly  
27 two-fold compared to random expectation. While transitions from A:T  $\rightarrow$  G:C were more common than  
28 the other mutations possible at A:T sites, they were still less common than any mutation at C:G sites.  
29 Transversions from A:T  $\rightarrow$  C:G or T:A were the least common and were found  $2.4 \times$  less frequently than  
30 expected by chance.

31 To assess the effect of the local sequence context on mutation rate, we measured the frequency of  
32 the bases surrounding random A:T and C:G sites in the genome and compared this to the base  
33 frequencies in the window surrounding SNMs (Fig. 4). We found non-random patterns surrounding all

1 six kinds of mutation, but the extent of the deviation was strongest for mutations at C:G sites. The  
2 deviation was particularly strong in the 2-4bp upstream of mutations at C:G sites and to a lesser extent  
3 1bp downstream of all mutation types. Specifically, the composition of the two nucleotides  
4 immediately upstream of mutated C:G sites was strongly biased. In the case of the CTC trinucleotide,  
5 for example, where the final C was mutated, that mutation rate was 4.5x the background rate.

## 6 **Mutability.**

7 To determine which genomic properties influenced the mutability of individual sites we used logistic  
8 regression to differentiate between the identified mutations and randomly selected not mutated sites.  
9 Using this model, we then calculated the probability of mutation, or 'mutability', for each site in the  
10 genome (See Materials and Methods for details). To assess the accuracy of the model we binned  
11 sites in the genome based on their mutability (0-1) and calculated the observed mutation rate in each  
12 bin (bin width = 0.01). The predicted mutability of sites was strongly correlated with observed mutation  
13 rate (Fig. 5,  $R^2=0.953$ , weighted by number of site-generations per bin). To ensure that the fit was not  
14 due to using the same mutations to generate the model and assess its fit, we also trained a model  
15 using a random subset of 1,000 mutations and excluded these sites when assessing the fit. As with  
16 the full data set, predicted and observed mutability were highly correlated ( $R^2=0.88$ ). The fit was  
17 slightly reduced, presumably because using fewer mutations to calculate mutation rates led to more  
18 noise. Although mutability ranged from nearly 0 to 1.0, we found that 99.9% of the genome had  
19 mutability values between 0.01 and 0.30, corresponding to mutation rates of  $0.25-55.9 \times 10^{-10}$ . The top  
20 25% of genome by mutability accounts for 57% of all mutations. Mutability was highest for sites in the  
21 3' and 5' UTRs (predicted  $\mu = 1.37 \times 10^{-9}$ ) and lowest for 0-fold and 4-fold degenerate sites (predicted  $\mu$   
22  $= 7.92 \times 10^{-10}$ ).

23 In neutrally evolving haploid DNA the level of nucleotide diversity ( $\theta_\pi$ ) is expected to be twice the  
24 product of mutation rate and the effective population size ( $2N_e\mu$ ). We binned silent sites (intergenic,  
25 intronic and 4-fold degenerate sites) into 100 uniformly spaced mutability categories from 0.0-1.0 and  
26 calculated  $\theta_\pi$  for each bin using natural variation in the six ancestral strains used to initiated the MA  
27 lines. We found that, as predicted, sites with higher mutability have higher neutral genetic diversity  
28 (Fig. 6).

## 29 **Factors influencing mutability.**

30 From the model of mutation rate, we extracted the relative contribution of different genomic properties  
31 to mutability. To allow comparison among the genomic properties, we scaled continuous predictors so  
32 that a change from 0 to 1 was a change of one standard deviation. We found that GC-content of the  
33 surrounding genome strongly influenced the mutability at a site. Increasing the GC content of the 10bp



1 surrounding a site increased its mutability (GC% 10bp, odds ratio = 1.38), but at larger scales GC  
2 content was negatively related to mutability (GC% 1000bp, odds ratio = 0.12). The negative  
3 relationship between GC-content and mutation rate was supported by a highly significant correlation  
4 between the observed mutation rate and GC content across the genome (see Supplementary Fig. S1,  
5  $R^2=0.831$ ,  $P<0.001$ ). Reflecting similar patterns of sequence context described above, the trinucleotide  
6 sequence in which a mutation occurred also had a strong effect on mutability. The most mutable  
7 trinucleotides were  $CTC$  and  $CAC$ , where the final C was the mutant position (odds ratio = 3.54 and  
8 2.02 respectively), and the least mutable were  $GTT$  and  $AGA$  (odds ratio = 0.57 and 0.58  
9 respectively). It was not possible to combine the triplets into a single predictor, but the maximum  
10 difference in mutability between triplets indicated a strong effect of sequence context on mutability. A  
11 number of other genomic properties increased mutability, such as gene density (odds ratio =1.17) and  
12 being upstream of a transcription start site (odds ratio =1.13). Interestingly, although a change of one  
13 standard deviation in transcription level had little effect on mutability (odds ratio =1.02), the most  
14 highly transcribed sites in the genome were 3.7× more mutable than untranscribed sites.

## 15 Discussion

16 In total we detected 6,843 mutations, the largest set of characterized spontaneous mutations to date.  
17 The overall rate of mutation across all lines was  $\mu = 11.5 \times 10^{-10}$ /site/generation, and the mutation rate  
18 for SNMs was  $9.63 \times 10^{-10}$  and  $1.90 \times 10^{-10}$  for small indels. There are therefore five SNMs for each small  
19 indel, consistent with previous results in *C. reinhardtii*, and similar to *Arabidopsis thaliana* (~5:1), but  
20 substantially lower than the ratios recently reported from MA studies in *S. cerevisiae* (33:1) and *D.*  
21 *melanogaster* (12:1). This large set of mutations, and the inclusion of multiple natural genotypes,  
22 allowed detailed examination of mutation rate variation between individuals within a species and  
23 mutation rate heterogeneity across the genome. In what follows we discuss the key results as they  
24 relate to the extent of mutation rate variation between natural strains and across the genome.

### 25 Within species mutation rate variation.

26 Our estimate of total mutation rate in *C. reinhardtii* is 14.2-fold and 4.6-fold higher than two previous  
27 estimates (Ness et al. 2012; Sung et al. 2012). The current estimate of mutation rate was partly driven  
28 by the higher rate in MA lines derived from ancestor CC-1373, but even after excluding this line the  
29 mutation rate is still substantially higher than previous estimates. The two MA lines (CC-2937-MA1,  
30 CC-2937-MA2) that were used to estimate mutation rate by Ness et al. (2012) continued to  
31 accumulate mutations for an average of ~611 generations additional generations, and the final  
32 mutation rate estimate for each of these two lines is within the confidence interval of the earlier  
33 estimate. Unfortunately, our experiment did not include strain CC-124 used in Sung et al. (2012), and  
34 so we can not directly compare mutation rates to this study. Only a single MA line (CC-1952-MA4) had

1 a mutation rate as low as the Sung et al. (2012) estimate and the mean of all MA lines derived from  
2 that ancestor was 9 times higher. Whether this variation is the result of methodological differences or  
3 biological variation between strain CC-124 and the six strains included in our study remains to be  
4 determined.

5 We observed a large degree of within-species variation for the mutation rate (Fig. 1). MA lines derived  
6 from strain CC-1373 had an average mutation rate more than three times that of the other strains. MA  
7 experiments in diploid species generally start with inbred lines, and it has been argued that mutation  
8 could be affected by recessive mutation rate modifiers that are not expressed in nature. However, *C.*  
9 *reinhardtii* is haploid, so the elevated rate in CC-1373 must be caused by a mutation modifier that  
10 arose since collection or by natural variation expressed in nature. CC-1373 is the slowest growing of  
11 the ancestral strains, indicating that it is not well adapted to laboratory conditions. A MA experiment in  
12 *Drosophila* provided evidence that individuals in poor condition have a higher mutation rate (Sharp  
13 and Agrawal 2012), so it is possible that the higher mutation rate in CC-1373 reflects its poor  
14 condition. At the other end of the spectrum, CC-1952 had the lowest mutation rate, nearly seven-fold  
15 lower than that of CC-1373. The extent of intraspecific mutation rate variation we found implies that  
16 measuring the mutation rate for a species from a single genotype may not adequately reflect the  
17 species as a whole, and interspecific differences in mutation rate may actually reflect poor sampling  
18 within species.

19 In general, theory predicts that selection are expected to drive the mutation rate towards zero,  
20 because alleles that increase the mutation rate will generate deleterious alleles and thereby reduce  
21 fitness (reviewed by Sniegowski and Raynes 2013). However, mutation rates are always above zero  
22 in nature, which is usually explained by the cost of increased fidelity or by the 'selection-drift barrier'  
23 imposed when selection for increasingly small improvements becomes too weak to counteract genetic  
24 drift. Under both hypotheses, the extent of intraspecific mutation rate variation may reflect mutation-  
25 selection balance in genes that affect DNA-repair, replication fidelity or the susceptibility to DNA  
26 damage. In our experiment, we detected at least two MA lines with mutation rates significantly higher  
27 than their strain means ( i.e., CC-2344-MA1 and CC-2931-MA5 had mutation rates 8.0× and 3.5×  
28 above their respective strain means, Fig. 1). It is likely that these two lines acquired mutations that  
29 damaged DNA repair or stability, concordant with the presence of two mutations in DNA repair  
30 proteins in CC-2344-MA1 and five such mutations in CC-2931-MA5. However, 26 of 85 MA lines also  
31 acquired one or more mutations that affect DNA repair associated proteins, but did not have elevated  
32 mutation rates. It is possible that many of these mutations did not significantly alter the mutation rate,  
33 or that the mutations arose too late in the experiment to cause a detectable elevation of mutation rate.  
34 The increase in mutation rate in line CC-2344-MA1 was greater than the extent of natural variation

1 among ancestral strains, suggesting that mutations that strongly alter mutation rate are common, and  
2 may segregate in natural populations until purged by selection. Therefore the high mutation rate of  
3 CC-1373 may be caused by a naturally occurring mutator allele. Alternatively, if *C. reinhardtii* is  
4 primarily asexual in nature, theory predicts that if a mutator allele results in a linked beneficial allele,  
5 the mutator will hitchhike to high frequency. A key parameter determining whether selection will favor  
6 higher mutation rates is the rate of recombination, but the frequency of sex and recombination in  
7 natural populations of *C. reinhardtii* is unknown.

## 8 **Spatial heterogeneity in mutation rate.**

9 By examining the spectrum of mutations and the local sequence in which they occur, we found clear  
10 evidence for heterogeneity in mutation rate at fine-scales. In particular, the rate of mutation at C:G  
11 sites ( $12.2 \times 10^{-10}$ ) was 2.4x higher than at A:T sites ( $5.19 \times 10^{-10}$ ) and transitions from C:G→T:A  
12 occurred at twice the rate expected if all mutations occurred at even rates (Fig. 2). An AT-biased  
13 mutation spectrum is consistent with a growing body of evidence suggesting that it might be universal  
14 in prokaryotes (Hershberg and Petrov 2010) and eukaryotes (e.g. Zhu et al. 2014). Additionally, we  
15 found that the sequence flanking a mutated site strongly influenced the mutation rate. In mammals  
16 methylated CpG sites are frequently deaminated, causing C to T transitions, but in *C. reinhardtii* there  
17 is only weak evidence of CpG methylation, and our data reveals only a small excess of CpG motifs in  
18 C to T mutations (Fig. 3). The most mutable triplet (CTC) had a mutation rate more than 10x higher  
19 than the least mutable triplet (GCA), and after accounting for background triplet frequencies, a  
20 mutation from CTC to CTT was 17x more likely than a mutation from AAA to AAG. Interestingly, this  
21 CTC triplet appears to be highly mutable across a very wide diversity of organisms, including fungi  
22 (Zhu et al. 2014), plants and animals (Alexandrov et al. 2013b). In human tumor genomes, there is a  
23 predominance of C to T and C to G mutations in the same CTCG sequence motif, which has been  
24 linked with the APOBEC family of cytidine deaminases (Alexandrov et al. 2013b). Given that this motif  
25 has been found repeatedly, it seems probable that the mutability of other sequence motifs may be  
26 shared across species, however the mechanisms underlying this phenomenon are unknown. The fact  
27 that the mutation rate can vary to this extent over very short scales has consequences for the  
28 evolution of DNA and protein sequence. In the future, incorporation of direct measurements of  
29 mutability into models of sequence change will facilitate better predictions of disease susceptibility and  
30 molecular evolution (see Michaelson et al. 2012; Neale et al. 2012; Samocha et al. 2014).

31 By comparing the distribution of intermutation distances to a random expectation, we found that there  
32 is an excess of mutations clustered within 1-10bp of one another (Fig. 4). The fact that these clusters  
33 all occur within MA-lines suggests that each represents a single multinucleotide mutation (MNM)  
34 event. In total there were 80 pairs and two trios of MNMs within 10bp of one another, implying that

1 2.8% of SNMs arise through clustered mutations. The average proportion of MNMs was similar in MA  
2 studies of *S. cerevisiae*, *D. melanogaster*, *C. elegans* and *A. thaliana* (3.4%), and genome sequencing  
3 of humans (1-4% Schrider et al. 2011; Harris and Nielsen 2014). The generation of these clusters has  
4 been linked to error prone polymerases such as Pol  $\zeta$  in *S. cerevisiae* (Stone et al. 2012; Northam et  
5 al. 2013). In human and *S. cerevisiae* the Pol  $\zeta$  enzyme creates an excess of GC to AA or TT MNMs  
6 (Northam et al. 2013; Harris and Nielsen 2014). Although we did not observe a similar excess of  
7 mutations at GC sites, we found that 27 of 55 dinucleotide MNMs occur at CC sites and that 25 of  
8 these resulted in AA/AT/TA/TT dinucleotides. The consistency of these results across taxa suggests  
9 that we cannot consider MNMs as an oddity of the mutational process. MNMs violate the assumption  
10 of independence between SNP sites and could potentially lead to mis-inferences about the nature of  
11 selection in the genome. Additionally, by altering two or more nearby sites, MNMs have the potential  
12 to move between fitness peaks that would otherwise require maladaptive single mutations as  
13 intermediates.

14 At large genomic scales, we found little evidence for heterogeneity of the mutation rate. For example,  
15 the mutation rate variation among 200Kbp windows could be largely accounted for by random  
16 fluctuations. Although we found clear evidence of fine-scale variation in mutation rate, the variation  
17 appears to be evenly spread along the chromosome. This effect can be seen in our predictive model  
18 of mutation, where the mutability of sites in 200Kbp windows averages out, so that the standard  
19 deviation among windows equates to  $\sim 7.5\%$  of the mean (i.e., mean mutability = 0.069, SD = 0.005).  
20 Our findings are consistent with direct measurements of mutation rate in *D. melanogaster* (Schrider et  
21 al. 2013), *S. cerevisiae* (Zhu et al. 2014) and humans (Kong et al. 2012), where no evidence of large  
22 scale variation in the mutation rate was detected. Although, comparative evidence suggests that  
23 substitution rate varies at the scale of megabases in mammals, this may be driven by selection or  
24 biased gene conversion during recombination. From our observations and direct estimates of mutation  
25 rate variation in other species, we conclude that the causes of mutational heterogeneity do not appear  
26 to operate at the scale of kilobases, and if heterogeneity exists at this scale it will require even more  
27 precise measurements of the mutation rate.

## 28 **Factors that predict mutability.**

29 Our model of mutability identified a number of other genomic properties that predict the rate of  
30 spontaneous mutation and create heterogeneity between sites. For example, the %GC of the 10bp  
31 around a mutated site was positively correlated with mutability (Odds ratio = 1.38, 1-SD=16.3%),  
32 probably because G:C bases and GC-rich triplets were more mutable. However, the GC-content of the  
33 1,000bp surrounding a site was negatively associated with its mutability (e.g., %GC of 1,000bp  
34 window, Odds ratio = 0.12, 1-SD=5.4%). A negative correlation between mutability and GC content in

1 humans has been attributed to higher melting temperatures of GC-rich DNA (Fryxell and Moon 2005).  
2 Because cytosine deamination is one of the most common sources of mutation and only occurs while  
3 DNA is single stranded, mutation is less common in regions with high melting temperature (Frederico  
4 et al. 1993). An alternate explanation for our observations is that sites with a high mutation rate, for an  
5 unknown reason, evolve low GC-content because mutation is AT-biased.

6 Our model of mutability also revealed an effect of gene expression when comparing untranscribed  
7 DNA to the most highly transcribed genes (odds ratio = 3.71). However, because most regions are  
8 untranscribed and the variance of transcription in expressed genes is relatively low, transcription level  
9 overall had little effect on mutability (odds ratio 1.02, 1-SD=108.3 FPKM). It is commonly reported that  
10 highly expressed genes are the most evolutionarily conserved, therefore an elevated mutation rate  
11 would predict that more deleterious mutations should occur in high expression genes and therefore  
12 more purifying would be required to conserve these sequences. The mean mutability score varied  
13 across sites with different annotations. The 5' and 3' UTRs had the highest mutability (predicted  $\mu =$   
14  $1.5 \times 10^{-9}$ ), which is consistent with the observation in other species that these regulatory regions are  
15 often found in open chromatin, allowing binding of transcription factors, and potentially leading to more  
16 damage to the DNA. Consistent with an increased mutation rate and AT-biased mutation, UTRs have  
17 the lowest GC content of any broad category of sites (56.7%). Although the model predicted a higher  
18 rate in UTRs, we did not observe an elevation in observed mutation rate, possibly because even with  
19 nearly 7,000 mutations there was still insufficient power to detect such subtle variation. Overall, the  
20 model accurately predicted the observed mutation rate, demonstrating that average mutation rate can  
21 be predicted from key genomic properties (Fig. 5). However, variation in mutability may not be fully  
22 captured with this approach Eyre-Walker and Eyre-Walker (2014). For a close fit between observed  
23 and predicted mutability, only the average mutability of each bin needs to be accurately predicted.  
24 There may still be unexplained variation around the mean within each bin and we should be cautious  
25 about predictions of mutability for very small numbers of sites. However, for large groups of sites the  
26 model accurately predicts the average mutation rate and we can be confident in the genomic  
27 properties that best predict mutation rate. Mutability also revealed that mutation rate variation affects  
28 patterns of neutral genetic variation. We found a clear positive relationship between mutability and  
29 nucleotide diversity at silent sites (Fig. 6). The model identifies the genomic properties of sites that  
30 mutated in our experiment and we show that using these genomic properties we are able to predict  
31 natural levels of genetic diversity. This implies that our model captures the variation in mutation rate  
32 that exists under natural conditions.

33 This study characterized the largest set of spontaneous mutations to date and demonstrated the  
34 insights that can be gained by combining MA with whole genome sequencing. We found 7-fold

1 variation in mutation rate among natural strains of *C. reinhardtii*. Although the mutation rate did not  
2 vary across large genomic windows, the mutation rate of individual sites was strongly affected by their  
3 flanking sequence, resulting in fine-scale heterogeneity of mutation rate. Other genomic properties,  
4 such as GC content, gene density and expression level, also influenced mutability. Similar results  
5 across a wide diversity of species suggests that general properties of mutation exist and that models  
6 of sequence evolution could be improved to reflect these properties and better detect selection in the  
7 genome or estimate phylogenetic relationships. In the near future rapidly evolving sequencing  
8 technologies will facilitate even more detailed investigation into the process of mutation from both MA  
9 and parent-offspring sequencing. One important avenue of future research will be a synthesis of  
10 findings from studies like ours with the underlying DNA repair and damage mechanisms to provide  
11 explanations for patterns mutational heterogeneity between individuals and across the genome.

## 12 **Methods**

### 13 **Mutation accumulation experiment.**

14 We conducted a mutation accumulation experiment in six genetically diverse wild strains of *C.*  
15 *reinhardtii* obtained from the Chlamydomonas Resource Center ([chlamycollection.org](http://chlamycollection.org)). The strains  
16 were chosen to broadly cover the geographic range of known *C. reinhardtii* samples in North America  
17 (Table 1). To initiate the MA lines, a single colony from each of the six ancestral strains was streaked  
18 out, and we randomly selected 15 individual colonies to start the replicated MA lines (for a total of 90  
19 MA lines). We bottlenecked the MA lines at regular intervals by selecting a random colony which was  
20 streaked onto a fresh agar plate. We estimated the number of generations undergone by each MA line  
21 over the course of the experiment by measuring the number of cells in colonies grown on agar plates  
22 after a period of growth equivalent to the times between transfers in the experiment. The details of the  
23 MA line creation and generation time estimation can be found in Morgan et al. (2014).

### 24 **Sequencing and alignment.**

25 To extract DNA, we grew cells on 1.5% Bold's agar for 4 days until there was a high density of cells, at  
26 which point the cells were collected and frozen at -80°C. We disrupted the frozen cells using glass  
27 beads, and extracted DNA using a standard phenol-chloroform extraction. Whole-genome re-  
28 sequencing was conducted using the Illumina GAII platform at the Beijing Genomics Institute (BGI-  
29 HongKong Co., Ltd, Hong Kong). The sequencing protocol was modified to accommodate the  
30 unusually high GC content of the *C. reinhardtii* genome (mean GC= 63.9%). Variation in GC-content is  
31 known to cause uneven representation of sequenced fragments, especially when GC > 55% (Aird et  
32 al. 2011). We therefore used a modified PCR step in sequencing library preparation, following Aird et  
33 al. (2011) (3 min at 98°C; 10 × [80 sec at 98°C, 30 sec at 65°C, 30 sec at 72°C]; 10 min at 72°C, with

1 2M betaine and slow temperature ramping 2.2°C/sec). We obtained ~30× coverage of the genome  
2 (3Gbp of 100bp paired-end sequence) for each of the MA lines.

3 We aligned reads to the *C. reinhardtii* reference genome (version 5.3) using BWA 0.7.4-r385 (Li and  
4 Durbin 2009). We included the plastid genome (NCBI accession NC\_005353), mitochondrial genome  
5 (NCBI accession NC\_001638) and the MT- locus (NCBI accession GU814015) to avoid misalignment  
6 of reads derived from these loci onto other parts of the nuclear genome. We tested a variety of values  
7 for the fraction of mismatching bases allowed in alignments, but variation about the default (n=0.04)  
8 did not improve the number of high quality reads mapped or genome coverage (results not shown).  
9 After alignment, we removed duplicate reads with the Picard tool MarkDuplicates (v1.90). To avoid  
10 calling false variants due to alignment errors, we used the GATK (v2.8-1) tools  
11 RealignerTargetCreator and IndelRealigner (McKenna et al. 2010; DePristo et al. 2011) to realign  
12 reads flanking potential insertions and deletions. We realigned all replicate MA lines from each starting  
13 strain together to ensure that the same alignment solutions were chosen in all lines derived from that  
14 strain. The realigned BAM files included all MA lines from given ancestral strain and were then used to  
15 jointly call genotypes using the UnifiedGenotyper from GATK. We used the “--output\_mode  
16 EMIT\_ALL\_SITES” option to output all genomic positions so that we could identify both high quality  
17 sites regardless of whether they had mutated. We used a “heterozygosity” parameter of 0.01, but  
18 previous testing in *C. reinhardtii* showed that our genotyping is not sensitive to this prior as long as  
19 read depth is high, as it is in the present experiment (Ness et al. 2012). To identify short insertions and  
20 deletions (indels) we used the GATK v(2.8-1) tool ‘HaplotypeCaller’, which performs local re-assembly  
21 of reads (i.e., indels called with UnifiedGenotyper were ignored). The six resulting Variant Call Format  
22 files (VCFs) (one per ancestral strain) were converted to wormtable databases using the python  
23 package WormTable v0.1.0 (Kelleher et al. 2013) which enabled efficient exploration of quality filters  
24 for mutation identification.

## 25 **Mutation identification.**

26 MA lines within an ancestral strain were genetically identical at the start of the experiment, so any  
27 unique allele carried by a replicate within a strain was a candidate mutation. We applied a number of  
28 filters to genotype calls to identify mutations, while minimizing false positive and false negative calls. A  
29 site was called as a mutation if within that ancestral strain:

- 30 (1) The mapping quality (MQ)  $\geq 90$  and the PHRED called site quality (QUAL)  $\geq 100$
- 31 (2) All MA lines were ‘homozygous’; *C. reinhardtii* is haploid therefore this filter avoided  
32 mapping errors due to paralogous loci.
- 33 (3) The genotype of exactly one MA line differed from the rest of the lines
- 34 (4) All non-mutated lines shared the same genotype

1 (5) At least two sequences have confident genotype calls

## 2 **Callable sites.**

3 To calculate mutation rates and define null expectations, we needed to know the total number of sites  
4 with equivalent quality to the new mutations, hereafter referred to as “callable” sites. However, the  
5 definitions and distributions of quality scores are often different for variant and invariant sites. We  
6 therefore inferred a second measures of quality for invariant sites that was comparable to that used for  
7 mutant sites. For each mutant site we extracted the QUAL and MQ for the mutation and the nearest  
8 invariant site, under the assumption that because most reads are shared between adjacent sites the  
9 quality characteristics of the sites will be similar. We then estimated the correlation and relationship  
10 between quality scores at neighboring mutant and invariant sites using a linear model (MQ:  
11  $R^2=0.9996$ ,  $P < 0.001$ , QUAL:  $R^2=0.38$ ,  $P < 0.001$ ). The linear relationships between invariant and  
12 variant quality scores were used to predict appropriate MQ and QUAL thresholds for invariant sites  
13 (invariant MQ threshold = 90, invariant QUAL threshold =36.4). Analogous to the mutation calling, a  
14 site was callable within an ancestral strain if no line was called as a heterozygote, all lines with  
15 mapped reads had the same genotype call and at least two MA lines had genotype calls.

## 16 **Sanger confirmation.**

17 We estimated the accuracy of our mutation calls using Sanger sequencing. We randomly selected 192  
18 mutation calls (32 per ancestral strain) including both short indels and SNMs. We amplified each locus  
19 in the putative mutant MA line and a non-mutated MA line from the same ancestral strain. Sequences  
20 were then visually inspected in SeqTrace v0.9.0 to confirm the presence of the mutated site.

## 21 **Mutation rate calculations.**

22 We calculated the mutation rate ( $\mu$ ) in each replicate as,  $\mu = \text{mutations} / (\text{callable sites} \times \text{MA}$   
23  $\text{generations})$ . Whenever multiple MA lines were combined for mutation rate calculations, the number  
24 of callable sites and MA generations (site-generations) for each MA line was included to accurately  
25 account for differences amongst replicate lines. Similarly, all null expectations and mutation rate  
26 estimates for particular classes of sites take into account the number of site-generations for the  
27 specific positions included. To compare the average mutation rate of the six ancestral strains, we used  
28 the GLS function in R to fit a linear model to the individual mutation rate estimates of the MA lines. The  
29 model included mutation rate as the response variable and ancestral strain as a fixed factor. We  
30 allowed the variance to differ among ancestral lines using the varIdent function (Zuur et al. 2009). We  
31 then used the gHt function to generate linear contrasts, allowing us to further explore differences  
32 among the ancestors.



## 1 **Base composition and sequence context.**

2 Throughout our analyses of the mutation spectrum, we treated complementary mutations (C:G and  
3 A:T) symmetrically, such that there were six distinct SNMs (A:T→C:G, A:T→G:C, A:T→T:A,  
4 C:G→A:T, C:G→G:C, C:G→T:A). To assess the base spectrum of mutations, we calculated the  
5 frequency of each of the six mutation types relative to the expected frequency calculated from the  
6 base composition of the callable sites. To analyze the local sequence context in which mutations  
7 occurred, we measured base composition at each of the positions 5bp upstream and downstream of  
8 the mutated site. To calculate the null expectation for sequence context we estimated base  
9 composition in analogous windows surrounding  $10^6$  randomly selected callable sites. Separate  
10 expectations were generated for sites centered on A:T and C:G.

## 11 **Spatial heterogeneity of mutation.**

12 To assess whether there was spatial heterogeneity in mutation rate we calculated the mutation rate  
13 across the genome in sliding windows. We conducted the analysis with windows of 100Kbp, 200Kbp,  
14 500Kbp and 1Mbp but because the results were qualitatively similar and we report only the 200Kbp  
15 analysis. The mutation rate of each window was calculated as the number of mutations in that window  
16 divided by the total number of callable site\*generations. To assess how the mutation rate in these  
17 windows varied relative to null expectations, we simulated a random distribution of mutations. For  
18 each MA line we generated a corresponding simulated line where the number of mutations carried by  
19 that line was distributed amongst the 200Kbp windows in proportion to the number of callable site-  
20 generations in each window. This procedure was repeated 1,000 times to generate an expected  
21 distribution of mutation rates across the 200Kbp windows.

22 We also tested for the presence of a non-random spatial distribution of mutations by comparing the  
23 observed distribution of intermutation distances to a simulated distribution. This approach differs from  
24 the analysis above because it can detect fine scale clusters of mutations. We simulated data under a  
25 model where mutations occur randomly across the genome, while retaining the same number of  
26 mutations per MA line and accounting for differences in the callable genome positions. For each MA  
27 line we generated a corresponding simulated sample by randomly assigning the number of mutations  
28 that occurred in that MA line to individual callable positions. This allowed us to assess whether there  
29 was significantly more clustering within and between lines while accounting for line-specific differences  
30 in callable sites. The observed and simulated distributions of intermutation distances were compared  
31 using the Kolmogorov–Smirnov (KS) test in R.

## 1 **Mutability.**

2 To determine which genomic properties influenced the mutability of individual sites we used  
3 regularized logistic regression to differentiate between the identified mutations and randomly selected  
4 callable sites. Our analysis was loosely based on the approach of Michaelson et al. (2012) . For all  
5 6,843 mutations and  $10^5$  non-mutated sites, we collated a table of genomic properties and annotations  
6 to use as predictors in the logistic regression. Genomic properties included %GC, gene density,  
7 transcription level, recombination rate, nucleosome occupancy and the trinucleotide sequence in  
8 which the site occurs (see Supplementary table S1 for details). A number of genomic properties were  
9 calculated for each site in windows of varying size from 10bp up to 1Mbp. Categorical predictors were  
10 converted to multiple binary predictors (0/1 for each category level) to be fitted in the same model with  
11 numeric predictors.

12 With these predictors we used the R package GLMnet (v1.9-8) (Friedman et al. 2010) to fit a logistic  
13 regression, where mutation class (mutant (1) or background (0)) was the binary response variable.  
14 From the model, we estimated mutability at each site in the genome as its probability of belonging to  
15 class 'mutation' given the genomic predictors at a given site. We assessed the accuracy of the  
16 predicted mutability by binning sites into 100 uniformly spaced mutability categories from 0.0-1.0. The  
17 exact value of the mutability score was not relevant, since it only reflected the ratio of mutations to  
18 background sites in the original training set. Within each mutability category the number of observed  
19 mutations divided by the total number of site-generations in that category was used to estimate the  
20 mutation rate. The observed mutation rate was predicted to be positively correlated with the mid-point  
21 mutability of the category. To test whether mutability predicted long term effects of mutation rate  
22 variation, we also calculated the relationship between mutability and natural levels of nucleotide  
23 diversity in the six ancestral strains used to start the MA lines. In neutrally evolving haploid DNA the  
24 level of nucleotide diversity ( $\theta_{\pi}$ ) is expected to be twice the product of mutation rate and the effective  
25 population size ( $2N_e\mu$ ), we therefore predict that the mutation rate should correlate positively with  
26 mutability. For this analysis whether a site was variant was omitted from the model in order to avoid  
27 circularity in the relationship between diversity and mutability. We binned silent sites (intergenic,  
28 intronic and 4-fold degenerate sites) into 100 uniformly spaced mutability categories from 0.0-1.0 and  
29 calculated  $\theta_{\pi}$  for all sites in each bin.

30 To assess the relative contributions of each genomic property to mutability, we extracted the  
31 coefficients of each predictor from the model. GLMnet can handle highly correlated predictors using an  
32 elastic-net penalty ( $\alpha$ ) that will either shrink coefficients toward each other (ridge penalty,  $\alpha=0$ ) or  
33 keeps one and discard the others (lasso,  $\alpha=1$ ). The fit of the model was unchanged by the selection of  
34  $\alpha$  and all results presented used  $\alpha=0.01$ . To compare the log(odds ratio) of each genomic property on

1 mutability, we scaled each predictor so that a change from 0.0 to 1.0 was a change of one standard  
2 deviation. As alternate scaling we also normalized the predictors such that each ranged from exactly  
3 zero to one.

#### 4 **Data Access**

5 All genomic data generated as part of this project is publicly available through the NCBI Sequence  
6 Read Archive (SRA BioProject Accession: SRP052900)

#### 7 **Figure Legends**

##### 8 **Figure 1. Variation in mutation rate between strains.**

9 Total mutation rate ( $\mu = \text{total mutations} / (\text{site} \times \text{generation})$ ) for each of the MA lines, categorized  
10 based on their ancestral strain. The boxes outline the 1st to 3rd quartile of the mutation rate in lines  
11 from a given ancestral strain, the thick horizontal line indicates the median mutation rate and the  
12 whiskers extend to the last data point that is within 1.5 $\times$  the interquartile range, points outside the  
13 whiskers are filled black.

##### 14 **Figure 2. Expected and observed distributions of intermutation distance.**

15 Comparison of observed (red) and expected (blue) distributions of the distance between mutations. In  
16 this plot, intermutation distance was measured as the nearest mutation irrespective of the MA line or  
17 strain it occurred in. The expected distribution was generated by randomizing the location of mutations  
18 in each MA line and recalculating the intermutation distances. The simulation was repeated 1000  
19 times and the average of those iterations is shown here.

##### 20 **Figure 3. Mutation base spectrum of single nucleotide mutations.**

21 Base mutation spectrum of 5716 single nucleotide mutations (SNMs). The deviation of the mutation  
22 rate for each of the six possible SNMs relative to its expectation based on equal mutation rates was  
23 calculated as the observed number of mutations of each kind divided the number of mutations  
24 expected if mutations occurred randomly with respect to base. Background base composition was  
25 calculated only from sites that have high quality genotype calls (callable sites).

##### 26 **Figure 4. Sequence context of spontaneous mutations.**

27 Deviations in the local sequence context of the 2bp flanking mutated sites. Deviations were calculated  
28 from the observed frequency of each base (A, T, C, G) in the flanks of mutated sites and the expected  
29 background composition based on flanking sequences of 10<sup>6</sup> random A:T or C:G sites. Each horizontal  
30 panel represents one of the six possible mutations indicated in the centre. Significant deviations from  
31 the background base composition at each position were detected with tests and indicated as \*P < 0.05,  
32 \*\*P < 0.01, \*\*\*P < 0.001 (alpha-values were adjusted for multiple tests using a Bonferroni correction).

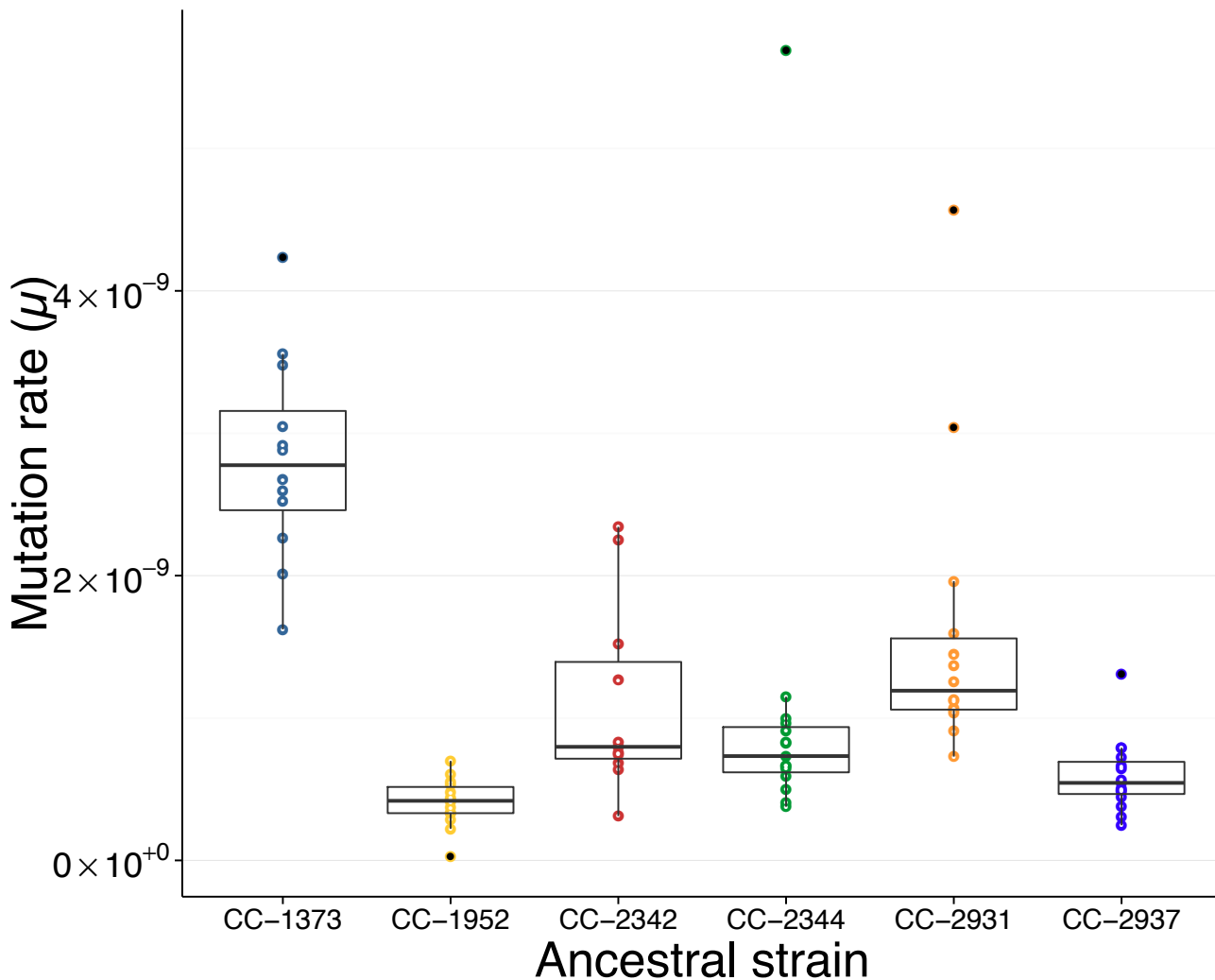
1 **Figure 5. Linear fit between observed mutation rate and predicted mutability.**

2 Mutability was estimated using a logistic regression, where the presence or absence of a mutation  
3 was the response variable, and a variety of genomic properties were used as predictors (see  
4 Supplementary table S1). Each point represents multiple genomic sites placed in discrete bins (width  
5 = 0.01) based on each site's mutability score. The size of each point is proportional to the number of  
6 sites in the genome with a given mutability. Observed mutation rates for each point were calculated as  
7 the number of observed mutations divided by the total number of callable sites-generations in that bin.  
8 The linear regression was weighted by the number of sites in each bin and the shaded grey area  
9 around the line represents the 95% confidence region.

10 **Figure 6. Relationship between natural genetic diversity and predicted mutability.**

11 Each point represents multiple genomic sites placed in discrete bins (width = 0.01) based on the  
12 predicted mutability of each site. Only putatively neutral sites (intronic, intergenic and 4-fold  
13 degenerate sites) were included in this figure. Nucleotide diversity ( $\theta_n$ ) was calculated in each bin from  
14 the six ancestral strains used to start the mutation accumulation lines. The size of each point is  
15 proportional to the number of sites in the genome with a given mutability.

16

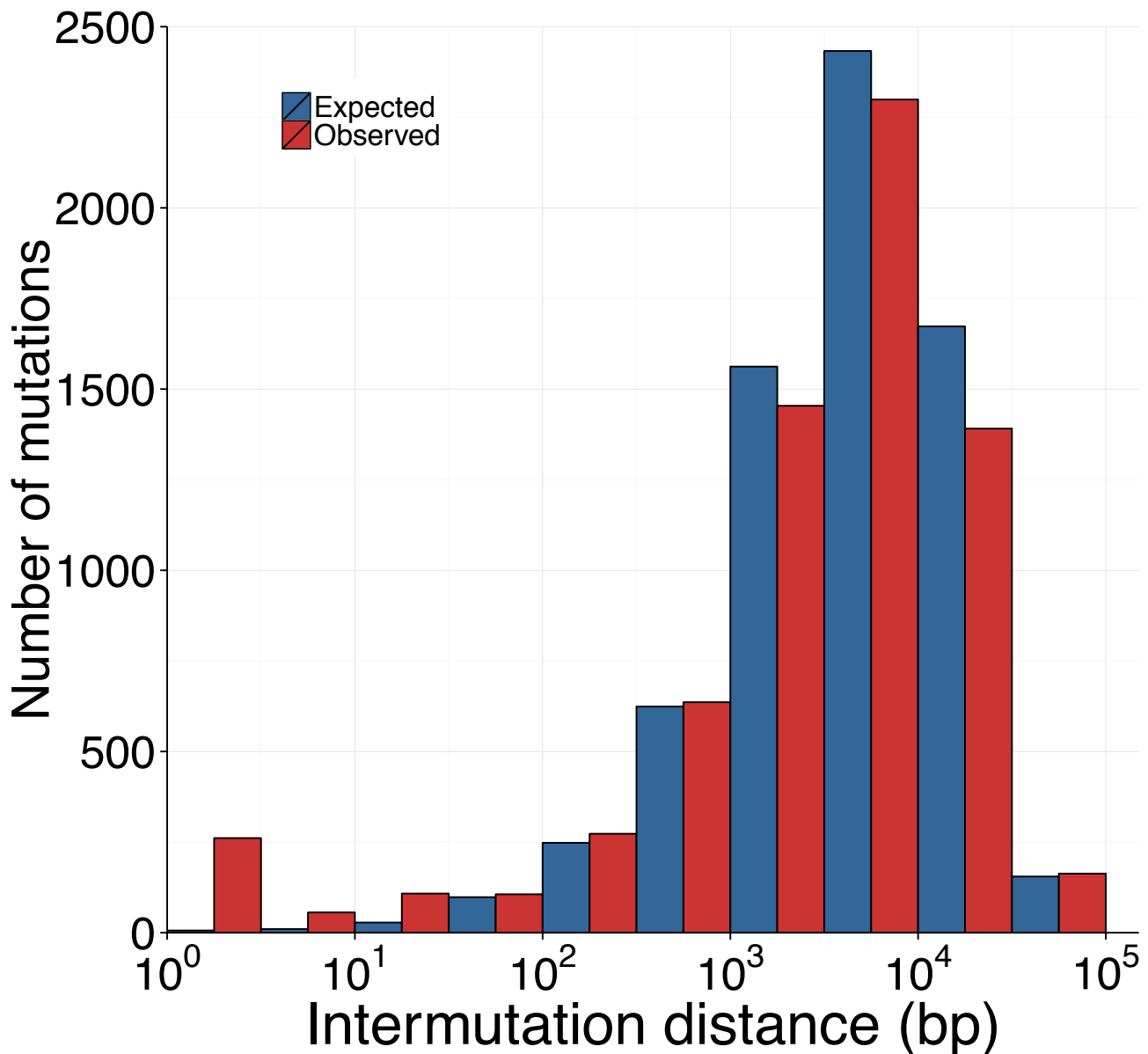


1

2 **Figure 1. Variation in mutation rate between strains.**

3 Total mutation rate ( $\mu = \text{total mutations} / (\text{site} \times \text{generation})$ ) for each of the MA lines, categorized  
4 based on their ancestral strain. The boxes outline the 1st to 3rd quartile of the mutation rate in lines  
5 from a given ancestral strain, the thick horizontal line indicates the median mutation rate and the  
6 whiskers extend to the last data point that is within  $1.5 \times$  the interquartile range, points outside the  
7 whiskers are filled black.

8

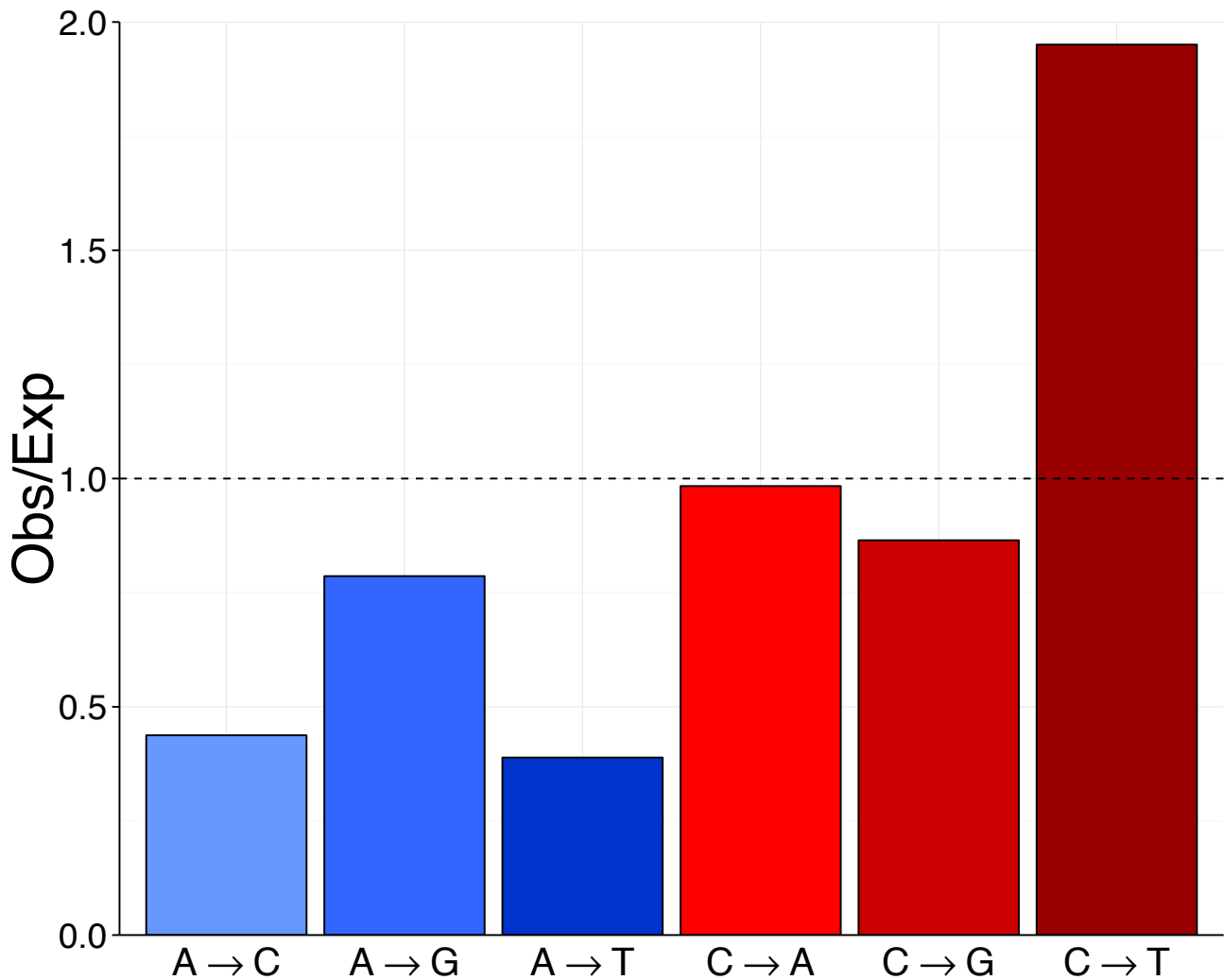


1

2 **Figure 2. Expected and observed distributions of intermutation distance.**

3 Comparison of observed (red) and expected (blue) distributions of the distance between mutations. In  
4 this plot, intermutation distance was measured as the nearest mutation irrespective of the MA line or  
5 strain it occurred in. The expected distribution was generated by randomizing the location of mutations  
6 in each MA line and recalculating the intermutation distances. The simulation was repeated 1000  
7 times and the average of those iterations is shown here.

8

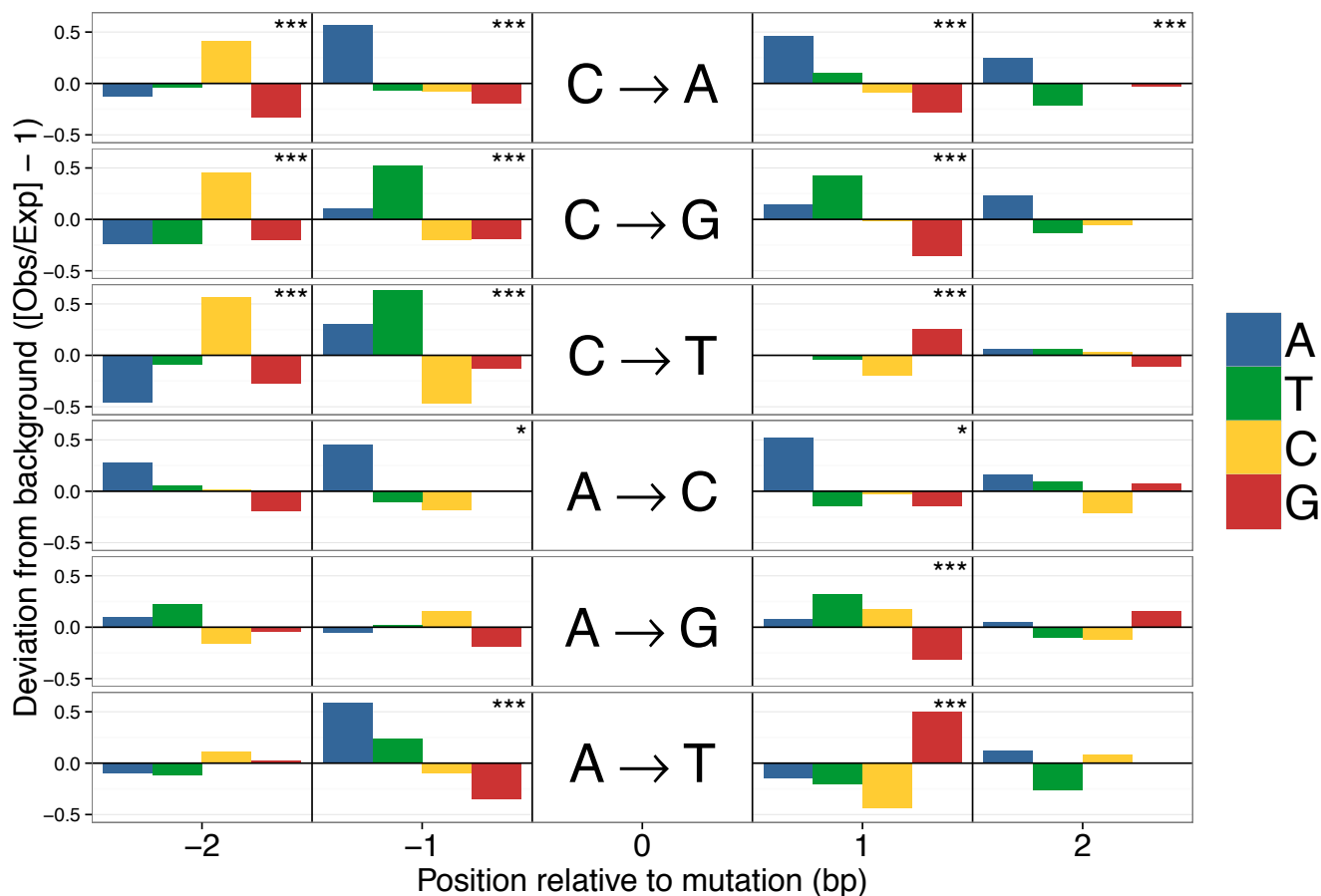


1

2 **Figure 3. Mutation base spectrum of single nucleotide mutations.**

3 Base mutation spectrum of 5716 single nucleotide mutations (SNMs). The deviation of the mutation  
4 rate for each of the six possible SNMs relative to its expectation based on equal mutation rates was  
5 calculated as the observed number of mutations of each kind divided the number of mutations  
6 expected if mutations occurred randomly with respect to base. Background base composition was  
7 calculated only from sites that have high quality genotype calls (callable sites).

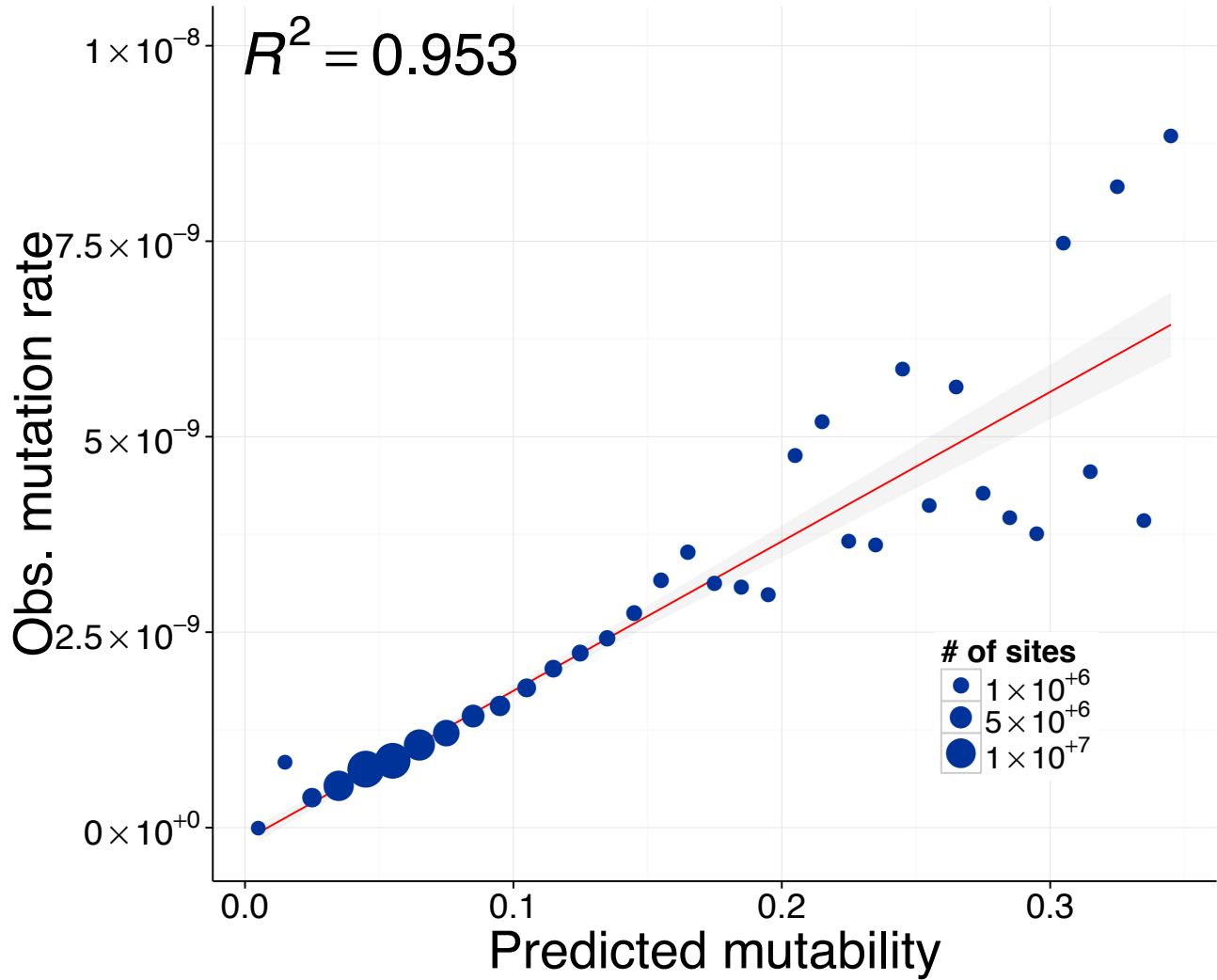
8



**Figure 4. Sequence context of spontaneous mutations.**

Deviations in the local sequence context of the 2bp flanking mutated sites. Deviations were calculated from the observed frequency of each base (A, T, C, G) in the flanks of mutated sites and the expected background composition based on flanking sequences of  $10^6$  random A:T or C:G sites. Each horizontal panel represents one of the six possible mutations indicated in the centre. Significant deviations from the background base composition at each position were detected with tests and indicated as \* $P < 0.05$ , \*\* $P < 0.01$ , \*\*\* $P < 0.001$  (alpha-values were adjusted for multiple tests using a Bonferroni correction).





1

2 **Figure 5. Linear fit between observed mutation rate and predicted mutability.**

3 Mutability was estimated using a logistic regression, where the presence or absence of a mutation

4 was the response variable, and a variety of genomic properties were used as predictors (see

5 Supplementary table S1). Each point represents multiple genomic sites placed in discrete bins (width

6 = 0.01) based on each site's mutability score. The size of each point is proportional to the number of

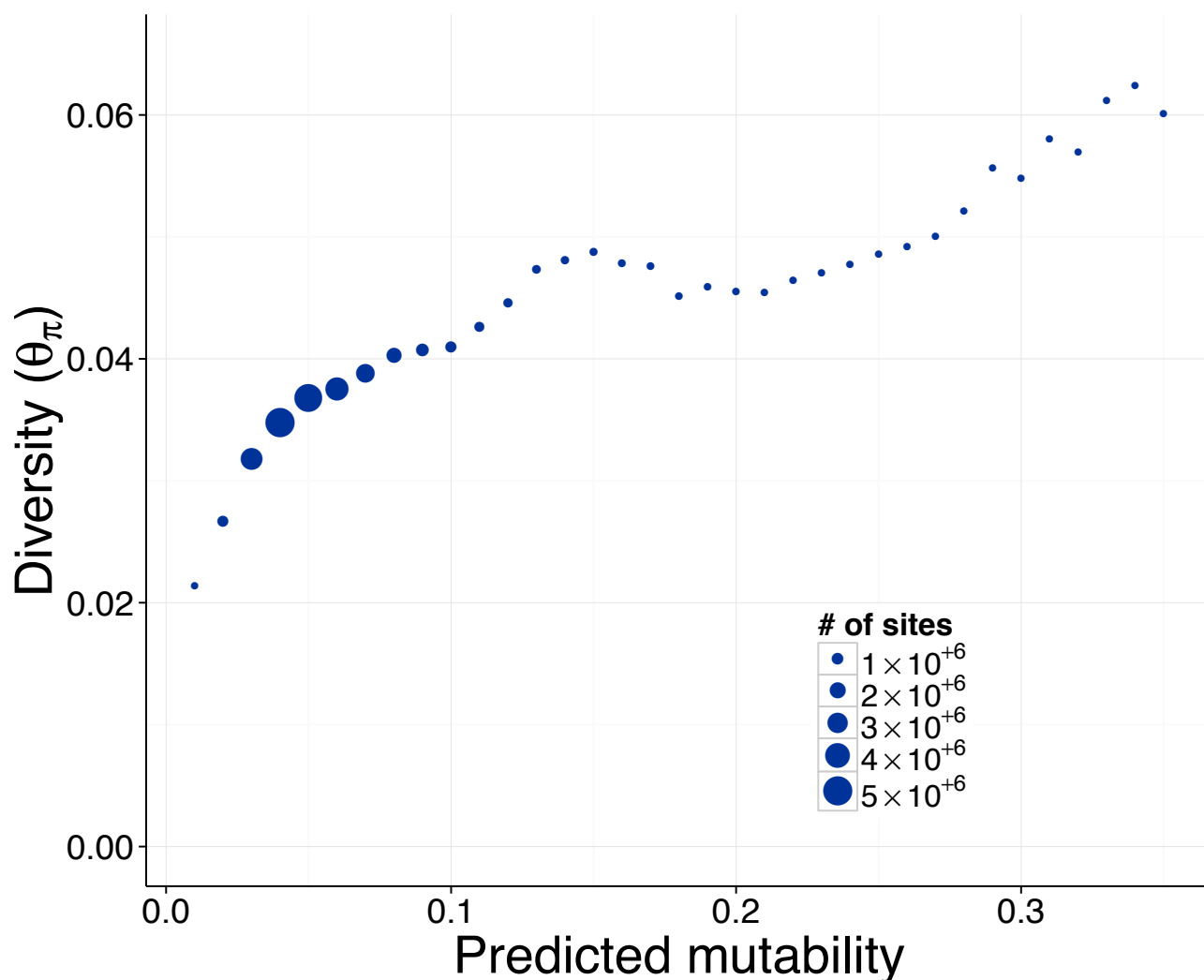
7 sites in the genome with a given mutability. Observed mutation rates for each point were calculated as

8 the number of observed mutations divided by the total number of callable sites-generations in that bin.

9 The linear regression was weighted by the number of sites in each bin and the shaded grey area

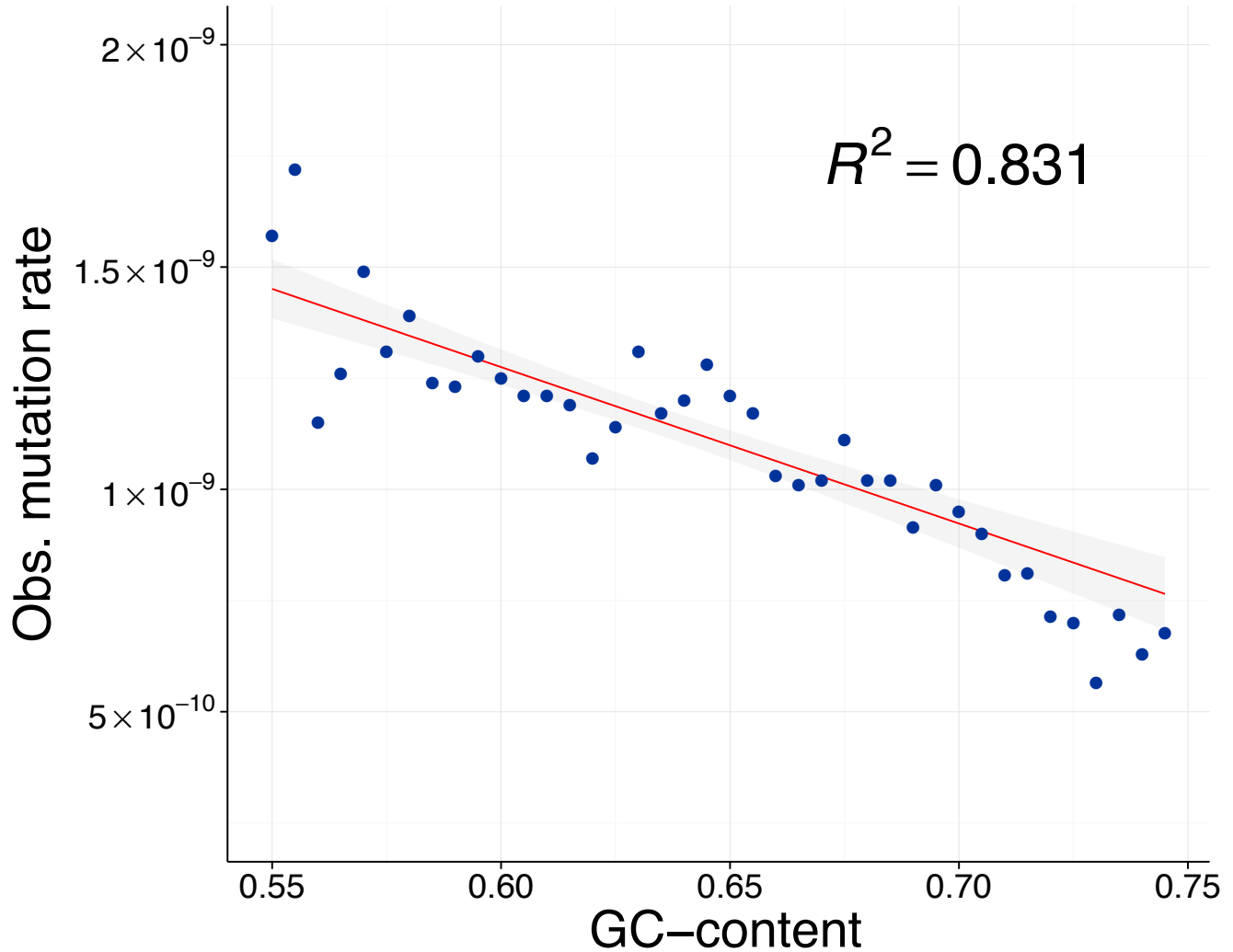
10 around the line represents the 95% confidence region.

11



**Figure 6. Relationship between natural genetic diversity and predicted mutability.**

Each point represents multiple genomic sites placed in discrete bins (width = 0.01) based on the predicted mutability of each site. Only putatively neutral sites (intronic, intergenic and 4-fold degenerate sites) were included in this figure. Nucleotide diversity ( $\theta_\pi$ ) was calculated in each bin from the six ancestral strains used to start the mutation accumulation lines. The size of each point is proportional to the number of sites in the genome with a given mutability.



1

2 **Supplementary Figure S1. Relationship between mutation rate and GC content.**

3 Linear fit between observed mutation rate and GC content of the 1000bp surrounding a site. Each  
4 point represents multiple genomic sites placed in discrete bins (width = 0.005) based on each site's  
5 GC content. Observed mutation rate for each point was calculated as the number of observed  
6 mutations divided by the total number of callable sites-generations in that bin. The shaded grey area  
7 around the line represents the 95% confidence region.

8

1 **Table 1. Ancestral strains of *Chlamydomonas reinhardtii* used for mutation accumulation (MA).**  
2 Each of the six strains was used to generate 11-15 replicate MA lines. The original sampling location,  
3 date and mating type (+/-) are indicated. The total number of single nucleotide mutations (SNMs) and  
4 short indels (<50bp) identified across all replicates of each strain are reported, along with the mean  
5 number of high quality ('callable') genomic sites sequenced in each strain.

Ancestral Strain	Collection Location/Year	Mating Type	MA lines	Mutations (SNMs / short indels)	Mean callable sites (Mbp)
CC-1373	Massachusetts/1945	+	12	1696/222	78.8
CC-1952	Minnesota / 1986	-	14	366/66	74.4
CC-2342	Pennsylvania / 1989	-	11	824/73	72.0
CC-2344	Pennsylvania / 1989	+	15	946/181	75.3
CC-2931	North Carolina / 1991	-	14	1215/405	72.5
CC-2937	Quebec / 1993	+	15	508/149	78.6

6

## References

- Aird D, Ross MG, Chen W-S, Danielsson M, Fennell T, Russ C, Jaffe DB, Nusbaum C, Gnirke A. 2011. Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries. *Genome Biol* **12**(2): R18.
- Alexandrov LB, Nik-Zainal S, Wedge DC, Aparicio SAJR, Behjati S, Biankin AV, Bignell GR, Bolli N, Borg Å, Børresen-Dale A-L et al. 2013a. Signatures of mutational processes in human cancer. *Nature* **500**(7463): 415-421.
- Alexandrov LB, Nik-Zainal S, Wedge DC, Campbell PJ, Stratton MR. 2013b. Deciphering signatures of mutational processes operative in human cancer. *Cell Reports* **3**(1): 246-259.
- Denver DR, Wilhelm LJ, Howe DK, Gafner K, Dolan PC, Baer CF. 2012. Variation in base-substitution mutation in experimental and natural lineages of *Caenorhabditis* nematodes. *Genome Biol Evol* **4**(4): 513-522.
- Depristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del Angel G, Rivas MA, Hanna M et al. 2011. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet* **43**(5): 491-498.
- Drake JW. 2006. Chaos and order in spontaneous mutation. *Genetics* **173**(1): 1-8.
- Drake JW, Charlesworth B, Charlesworth D, Crow JF. 1998. Rates of spontaneous mutation. *Genetics* **148**(4): 1667-1686.
- Ehrlich M, Wang RY. 1981. 5-Methylcytosine in eukaryotic DNA. *Science* **212**(4501): 1350-1357.
- Eyre-Walker A, Eyre-Walker YC. 2014. How much of the variation in the mutation rate along the human genome can be explained? *G3:Genes|Genomes|Genetics* **4**(9): 1667-1670.
- Frederico LA, Kunkel TA, Shaw BR. 1993. Cytosine deamination in mismatched base pairs. *Biochemistry* **32**(26): 6523-6530.
- Friedman J, Hastie T, Tibshirani R. 2010. Regularization Paths for Generalized Linear Models via Coordinate Descent. *J Stat Soft* **33**(1): 1-22.
- Fryxell KJ, Moon W-J. 2005. CpG mutation rates in the human genome are highly dependent on local GC content. *Mol Biol Evol* **22**(3): 650-658.
- Halligan DL, Keightley PD. 2009. Spontaneous mutation accumulation studies in evolutionary genetics. *Annu Rev Ecol Evol Syst* **40**: 151-172.
- Harris K, Nielsen R. 2014. Error-prone polymerase activity causes multinucleotide mutations in humans. *Genome Res* **24**(9): 1445-1454.
- Hershberg R, Petrov DA. 2010. Evidence that mutation is universally biased towards AT in bacteria. *PLoS Genet* **6**(9).
- Johnson T. 1999. Beneficial mutations, hitchhiking and the evolution of mutation rates in sexual populations. *Genetics* **151**(4): 1621-1631.
- Keightley PD, Ness RW, Halligan DL, Haddrill PR. 2014a. Estimation of the spontaneous mutation rate per nucleotide site in a *Drosophila melanogaster* full-sib family. *Genetics* **196**(1): 313-320.
- Keightley PD, Pinharanda A, Ness RW, Simpson F, Dasmahapatra KK, Mallet J, Davey JW, Jiggins CD. 2014b. Estimation of the spontaneous mutation rate in *Heliconius melpomene*. *Mol Biol Evol*.
- Kelleher J, Ness RW, Halligan DL. 2013. Processing genome scale tabular data with Wormtable. *BMC Bioinformatics* **14**: 356.
- Kimura M. 1967. On evolutionary adjustment of spontaneous mutation rates. *Genet Res* **9**(1): 23-34.
- Kong A, Frigge ML, Masson G, Besenbacher S, Sulem P, Magnusson G, Gudjonsson SA, Sigurdsson A, Jonasdottir A, Jonasdottir A et al. 2012. Rate of *de novo* mutations and the importance of father's age to disease risk. *Nature* **488**(7412): 471-475.
- Lang GI, Murray AW. 2011. Mutation rates across budding yeast chromosome VI are correlated with replication timing. *Genome Biol Evol* **3**: 799-811.

- Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**(14): 1754-1760.
- Lynch M. 2010. Evolution of the mutation rate. *Trends Genet* **26**(8): 345-352.
- Matic I, Radman M, Taddei F, Picard B, Doit C, Bingen E, Denamur E, Elion J. 1997. Highly variable mutation rates in commensal and pathogenic *Escherichia coli*. *Science* **277**(5333): 1833-1834.
- Mckenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M et al. 2010. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**(9): 1297-1303.
- Michaelson JJ, Shi Y, Gujral M, Zheng H, Malhotra D, Jin X, Jian M, Liu G, Greer D, Bhandari A et al. 2012. Whole-genome sequencing in autism identifies hot spots for *de novo* germline mutation. *Cell* **151**(7): 1431-1442.
- Morgan AD, Ness RW, Keightley PD, Colegrave N. 2014. Spontaneous mutation accumulation in multiple strains of the green alga, *Chlamydomonas reinhardtii*. *Evolution* **68**(9): 2589-2602.
- Neale BM, Kou Y, Liu L, Ma'ayan A, Samocha KE, Sabo A, Lin C-F, Stevens C, Wang L-S, Makarov V et al. 2012. Patterns and rates of exonic *de novo* mutations in autism spectrum disorders. *Nature* **485**(7397): 242-245.
- Ness RW, Morgan AD, Colegrave N, Keightley PD. 2012. An estimate of the spontaneous mutation rate in *Chlamydomonas reinhardtii*. *Genetics* **192**: 1447-1454.
- Northam MR, Moore EA, Mertz TM, Binz SK, Stith CM, Stepchenkova EI, Wendt KL, Burgers PMJ, Shcherbakova PV. 2013. DNA polymerases and Rev1 mediate error-prone bypass of non-B DNA structures. *Nucleic Acids Res* **42**(1): 290-306.
- Otto SP. 2009. The evolutionary enigma of sex. *Am Nat* **174** **Suppl 1**: S1-S14.
- Samocha KE, Robinson EB, Sanders SJ, Stevens C, Sabo A, McGrath LM, Kosmicki JA, Rehnström K, Mallick S, Kirby A et al. 2014. A framework for the interpretation of *de novo* mutation in human disease. *Nat Genet* **46**(9): 944-950.
- Schrider DR, Houle D, Lynch M, Hahn MW. 2013. Rates and genomic consequences of spontaneous mutational events in *Drosophila melanogaster*. *Genetics* **194**(4): 937-954.
- Schrider DR, Hourmozdi JN, Hahn MW. 2011. Pervasive multinucleotide mutational events in eukaryotes. *Curr Biol* **21**(12): 1051-1054.
- Sharp NP, Agrawal AF. 2012. Evidence for elevated mutation rates in low-quality genotypes. *Proc Natl Acad Sci USA* **109**(16): 6142-6146.
- Sniegowski P, Raynes Y. 2013. Mutation Rates: How Low Can You Go? *Curr Biol* **23**(4): R147-R149.
- Stamatoyannopoulos JA, Adzhubei I, Thurman RE, Kryukov GV, Mirkin SM, Sunyaev SR. 2009. Human mutation rate associated with DNA replication timing. *Nat Genet* **41**(4): 393-395.
- Stone JE, Lujan SA, Kunkel TA. 2012. DNA polymerase zeta generates clustered mutations during bypass of endogenous DNA lesions in *Saccharomyces cerevisiae*. *Environ Mol Mutagen* **53**(9): 777-786.
- Sundin GW, Weigand MR. 2007. The microbiology of mutability. *FEMS Microbiol Lett* **277**(1): 11-20.
- Sung W, Ackerman MS, Miller SF, Doak TG, Lynch M. 2012. Drift-barrier hypothesis and mutation-rate evolution. *Proc Natl Acad Sci USA* **109**(45): 18488-18492.
- Thorvaldsdóttir H, Robinson JT, Mesirov JP. 2012. Integrative Genomics Viewer (IGV): high-performance genomics data visualization and exploration. *Briefings in Bioinformatics*.
- Veltman JA, Brunner HG. 2012. *De novo* mutations in human genetic disease. *Nat Rev Genet* **13**(8): 565-575.
- Venn O, Turner I, Mathieson I, de Groot N, Bontrop R, Mcvean G. 2014. Nonhuman genetics. Strong male bias drives germline mutation in chimpanzees. *Science* **344**(6189): 1272-1275.
- Zhu YO, Siegal ML, Hall DW, Petrov DA. 2014. Precise estimates of mutation rate and spectrum in yeast. *Proc Natl Acad Sci USA* **111**(22): E2310-2318.
- Zuur A, Ieno EN, Walker N, Saveliev AA, Smith GM. 2009. *Mixed Effects Models and Extensions in Ecology with R*. Springer.