

Detection of Heterozygous Mutations in the Genome of Mismatch Repair Defective Diploid Yeast Using a Bayesian Approach

Sarah Zanders,^{*,1} Xin Ma,^{†,1} Arindam RoyChoudhury,^{†,1,2} Ryan D. Hernandez,^{†,3}
Ann Demogines,^{*,4} Brandon Barker,[§] Zhenglong Gu,[§]
Carlos D. Bustamante^{†,5,6} and Eric Alani^{*,6}

^{*}Department of Molecular Biology and Genetics, Cornell University, Ithaca, New York 14853, [†]Department of Biological Statistics and Computational Biology, Cornell University, Ithaca, New York 14853, [‡]Department of Human Genetics, University of Chicago, Chicago, Illinois 60637 and [§]Division of Nutritional Sciences, Cornell University, Ithaca, New York 14853

Manuscript received June 21, 2010
Accepted for publication July 19, 2010

ABSTRACT

DNA replication errors that escape polymerase proofreading and mismatch repair (MMR) can lead to base substitution and frameshift mutations. Such mutations can disrupt gene function, reduce fitness, and promote diseases such as cancer and are also the raw material of molecular evolution. To analyze with limited bias genomic features associated with DNA polymerase errors, we performed a genome-wide analysis of mutations that accumulate in MMR-deficient diploid lines of *Saccharomyces cerevisiae*. These lines were derived from a common ancestor and were grown for 160 generations, with bottlenecks reducing the population to one cell every 20 generations. We sequenced to between 8- and 20-fold coverage one wild-type and three mutator lines using Illumina Solexa 36-bp reads. Using an experimentally aware Bayesian genotype caller developed to pool experimental data across sequencing runs for all strains, we detected 28 heterozygous single-nucleotide polymorphisms (SNPs) and 48 single-nt insertion/deletions (indels) from the data set. This method was evaluated on simulated data sets and found to have a very low false-positive rate ($\sim 6 \times 10^{-5}$) and a false-negative rate of 0.08 within the unique mapping regions of the genome that contained at least sevenfold coverage. The heterozygous mutations identified by the Bayesian genotype caller were confirmed by Sanger sequencing. All of the mutations were unique to a given line, except for a single-nt deletion mutation which occurred independently in two lines. All 48 indels, composed of 46 deletions and two insertions, occurred in homopolymer (HP) tracts [*i.e.*, 47 poly(A) or (T) tracts, 1 poly(G) or (C) tract] between 5 and 13 bp long. Our findings are of interest because HP tracts are present at high levels in the yeast genome ($>77,400$ for 5- to 20-nt HP tracts), and frameshift mutations in these regions are likely to disrupt gene function. In addition, they demonstrate that the mutation pattern seen previously in mismatch repair defective strains using a limited number of reporters holds true for the entire genome.

MUTATION rates in prokaryotic and eukaryotic organisms are typically determined by measuring reversion or forward mutation for specific marker alleles. These values are then extrapolated to obtain genome-wide estimates. Mutation rates in higher eukaryotes are also estimated by analyzing sequence divergence between

different strains or species, followed by reconstructing the accumulation of mutations since divergence (reviewed in NISHANT *et al.* 2009). These approaches suffer from two main limitations. First, recent studies have shown that mutation rate and repair efficiency vary across the genome and are affected by parameters that include base composition, local recombination rate, gene density, transcriptional activity, repair efficiency, chromatin structure, nucleosome position, and replication timing (WOLFE *et al.* 1989; DATTA and JINKS-ROBERTSON 1995; MATASSI *et al.* 1999; HARDISON *et al.* 2003; ARNDT *et al.* 2005; HAWK *et al.* 2005; TEYTELMAN *et al.* 2008; WASHIETL *et al.* 2008; STAMATOYANNOPOULOS *et al.* 2009). Second, genomic comparisons can yield inaccurate rate measurements because DNA repair and subsequent purifying natural selection can bias the number and type of mutations that remain in the population, especially for mutations that occur in coding regions (reviewed in NISHANT *et al.* 2009).

¹These authors contributed equally to this work.

²Present address: Department of Biostatistics, Columbia University, New York, NY 10032.

³Present address: Department of Bioengineering and Therapeutic Sciences, University of California, San Francisco, CA 94143.

⁴Present address: Section of Molecular Genetics and Microbiology, University of Texas, Austin, TX 78712.

⁵Present address: Department of Genetics, Stanford University, Stanford, CA 94305.

⁶Corresponding authors: Cornell University, Department of Molecular Biology and Genetics, 459 Biotechnology Bldg., Ithaca, NY 14853-2703. E-mail: eea3@cornell.edu; and Stanford University School of Medicine, Department of Genetics, 300 Pasteur Dr., Stanford, CA 94305. E-mail: cdbustam@stanford.edu

The DNA mismatch repair system improves the fidelity of DNA replication by about 1000-fold by excising DNA mismatches in the newly replicated strand that arise from polymerase misincorporation and slippage (reviewed in MODRICH and LAHUE 1996; KUNKEL and ERIE 2005; McCULLOCH and KUNKEL 2008). Eukaryotes contain multiple MutS (MSH) and MutL (MLH) homologs (reviewed in KUNKEL and ERIE 2005). In *Saccharomyces cerevisiae*, two heterodimeric MutS homolog complexes, MSH2–MSH3 and MSH2–MSH6, act in mismatch recognition. MSH2–MSH6 is primarily involved in repairing base–base and small insertion/deletion loop mismatches. MSH2–MSH3 acts primarily on insertion/deletion loop mismatches up to 17 nt in length. In the presence of ATP, both MSH complexes interact primarily with MLH1–PMS1 to form a mismatch-MSH–MLH complex that interacts with downstream repair components. Recent work in humans and yeast suggests that MLH1–PMS1 contains an ATP–Mn²⁺-dependent latent endonuclease activity that acts near the mismatch and is essential for MMR, most likely in excision steps (KADYROV *et al.* 2006, 2007). Null mutations in MSH2 and MLH1, the key partners in the MSH and MLH complexes, confer severe defects in MMR; reporter assays have shown that strains bearing these mutations display high rates of base substitutions and DNA slippages. For example, in an assay that measures frameshift mutations in homopolymeric runs, *msh2Δ* and *mlh1Δ* mutations confer mutation rates that are ~10,000-fold higher than wild type (MARSISCHKY *et al.* 1996; TRAN *et al.* 1997, 2001; GRAGG *et al.* 2002).

Our goal in this study was to analyze with limited bias the rate at which mutations occur in MMR-defective lines due to DNA polymerase errors during DNA replication, and to identify novel genomic features associated with these errors. The baker's yeast *S. cerevisiae* is an ideal model system in which to perform these studies because genetic analysis of many of the key MMR factors has been performed; more importantly the effect of null mutations in these factors has been extensively characterized using a variety of mutator assays (KUNKEL and ERIE 2005). Previously, one of our groups (HECK *et al.* 2006) grew wild-type and conditional *mlh1* (*mlh1-7^{ts}*) diploid strains of *S. cerevisiae* for 160 generations with bottlenecks that reduced the population size to one cell every 20 generations. These lines were grown at 35°, the nonpermissive temperature for *mlh1-7^{ts}*. A conditional *mlh1* allele was chosen instead of a null so that mutation accumulation in the absence of MMR could be limited to 160 generations by shifting cells at generation 160 to the permissive temperature for MMR function. The *mlh1-7^{ts}* mutation contains two mutations within the ATP-binding domain of MLH1 (K67A, D69A). Unlike *mlh1Δ* strains that display poor spore viability due to defects in meiotic crossing over, *mlh1-7^{ts}* lines display wild-type spore viability at the permissive temperature. Such a phenotype allowed us to easily identify recessive

lethal mutations (HECK *et al.* 2006). At the nonpermissive temperature, the *mlh1-7^{ts}* mutation conferred a phenotype similar to the null in the canavanine resistance mutation assay and a mutator phenotype in the *lys2-Δ14* reversion assay that was 1000-fold higher than MLH1 but 4-fold lower than the null (HECK *et al.* 2006); J. HECK and E. ALANI, unpublished observations).

Tetrad analysis showed that the *mlh1-7^{ts}* bottleneck lines would be ideally suited for a high-throughput DNA sequencing approach that would identify mutagenesis patterns. First, the wild-type lines maintained high spore viability (~94%) at generation 160. In contrast, *mlh1-7^{ts}* lines displayed spore viabilities that ranged from 1.1 to 77%, demonstrating that the lines had accumulated recessive lethal mutations. Second, comparative genome hybridization (CGH) and pulse-field (PFGE) analyses of the *mlh1-7^{ts}* strains indicated that they did not undergo major genome rearrangements (HECK *et al.* 2006). Third, because the lines were grown as diploids for a limited number of generations, secondary mutations, dominant or recessive, that alter the rate or type of mutagenesis should rarely occur. Also, because there is no sexual reproduction and mutations should clonally propagate after escaping the initial bottleneck, newly arising mutations should appear as heterozygous sites. Finally, the above strategy should limit biases in mutation accumulation because the diploid cells were grown in rich media under minimal selection pressure where deleterious mutations could accumulate (HECK *et al.* 2006).

As described below, a Bayesian method was developed to detect heterozygous mutations in one wild-type and three *mlh1-7^{ts}* lines using whole-genome sequencing. We detected 28 heterozygous single-nucleotide polymorphisms (SNPs) and 48 single-nt insertion/deletion (indels) in the mutator lines, all of which mapped to homopolymeric runs of nucleotides (HP tracts). The mutation spectra match closely with that seen in MMR defective strains using different reporter constructs (MARSISCHKY *et al.* 1996; TRAN *et al.* 1997, 2001). This demonstrates that the mutation pattern seen previously using a limited number of reporters holds true for the entire genome. In addition, we were able to correlate genotype to phenotype for one locus in one mutator line. Together this work provides new insights into how mismatch repair can shape genome stability and dynamics, mutation mechanisms, and evolution.

MATERIALS AND METHODS

Whole-genome sequencing analysis of Mut lines: Bottleneck experiments involving 10 independent wild-type (*MATa/MATα*, *his3/HIS3*, *LEU2/leu2*, *cyl1/cyl1*, *ade2/ADE2*, *ura3/ura3*, *trp1/trp1*) and *mlh1-7^{ts}* (*MATa/MATα*, *mlh1-7::KanMX4/mlh1-7::KanMX4*, *his3/HIS3*, *LEU2/leu2*, *cyl1/cyl1*, *ade2/ADE2*, *ura3/ura3*, *trp1/trp1*) lines were performed previously (HECK *et al.* 2006). Three of

the 10 *mlh1-7^s* lines at generation 160 were analyzed by whole-genome sequencing. These lines were chosen to ensure a reasonable sample set of mutations and displayed a lower range of spore viabilities (2.5–15.6%) following tetrad dissection compared to the entire set (1.1–77%).

Whole-genome sequencing was performed at the Cornell University Life Sciences Core Laboratory Center (CLC) using an Illumina Genome Analyzer (<http://www.illumina.com>). Yeast genomic DNA for whole genome sequencing was prepared using a Qiagen genomic DNA preparation kit (<http://www.qiagen.com>). Sequencing was performed using the Illumina pipeline for 36-bp single-end reads. Reads were aligned onto the S288c genome (<http://genome.ucsc.edu/cgi-bin/hgGateway>) using Novoalign (<http://www.novocraft.com>), a program that performs a gapped alignment with high specificity and sensitivity.

Detection of DNA sequence heterozygosity using a Bayesian approach: We analyzed five diploid strains in this study: a wild-type strain at generations 0 and 160 (Wt0, Wt160) and three derived *mlh1-7^s* mutator lines grown vegetatively (*i.e.*, no meiosis) and bottlenecked to one cell every 20 generations until generation 160 (Mut2, Mut3, Mut4). Several aspects of the experiment required us to develop a novel approach for calling genotypes from the sequencing data. First, the initial wild-type strain (Wt0) likely contained SNPs and indels that distinguish it from the reference yeast genome. Because all lines were grown vegetatively, they were all expected to have these “propagated” SNPs and indels. Thus reads from the five sequenced lines were used to identify these variants. Furthermore, we expect new mutations (as this those occurring in Wt160, Mut2, Mut3, or Mut4 during generations 1–160) to be heterozygous at the end of the experiment and few, if any, variants are expected to be shared (as this would require independent hits in replicate lines). Finally, the sequencing depth ($\sim 8\text{--}20\times$) suggests moderate but not exceptional power to detect heterozygous mutations from the sequence of a single line on its own. Therefore, we developed a Bayesian SNP caller that (1) aligns all reads to the genome and (2) uses read depth and quality scores at a given position to call genotypes for all five lines simultaneously.

Importantly, our Bayesian model allows us to distinguish between a propagated mutation, (defined as a variant seen in all five strains in either heterozygous or homozygous state from Wt0) and a derived mutation, defined as a DNA sequence variant that arose in only a single line. First, we indexed the five diploid strains as $s = 1, 2, 3, 4, 5$ for Wt0, Wt160, Mut2, Mut3, and Mut4, respectively. We set the prior probability of strain s being heterozygous as $\text{Prior}_s = 10^{-7}, 10^{-8}, 10^{-5}, 10^{-5}, 10^{-5}$ for $s = 1, 2, 3, 4,$ and 5 , respectively, according to mutation rates previously determined in wild-type and mismatch-repair defective organisms (DENVER *et al.* 2005; IYER *et al.* 2006; NISHANT *et al.* 2009). It is important to note that Wt160 was assigned a lower prior probability of being heterozygous relative to Wt0. This is because a heterozygosity in Wt0 is defined as the difference between the Wt0 strain (HECK *et al.* 2006) and the S288c reference genome (<http://genome.ucsc.edu/cgi-bin/hgGateway>). There were a significant number of differences between the two strains. On the other hand, a heterozygosity in Wt160 was defined as one that occurred during the bottleneck experiment (propagated). Because there were only 160 generations between Wt0 and Wt160, we expected the number of differences between the lines to be small; in fact, none were detected.

At a given locus, let A and a be the major and minor allele types, respectively, based on the allele counts from all the strains. Let N_s be the total number of alleles observed for strain s ; let $A_{j,s}$ be the type of the j th allele copy among these N_s alleles, $j = 0, 1, \dots, N_s$. Let e_j be the probability that the j th allele has been assigned the wrong allele type. We estimated e_j

from the error rates given by DOHM *et al.* (2008) for 36-bp Solexa reads as a function of read position.

To call SNPs and indels in Wt0, we used the allele count data from Wt0 along with that from the other four strains. The posterior probabilities of a given genomic position being homozygous or heterozygous in Wt0 are

$$\begin{aligned} P_1(\text{Heter.} | \text{Data}) &= \frac{P_1(\text{Data} | \text{Heter.}) \times P_1(\text{Heter.})}{P(\text{Data})} \\ &\propto \text{Prior}_1 \times (0.5)^{\sum_{j=1}^{N_1} N_j} \\ P_1(\text{Homo.} | \text{Data}) &= \frac{P_1(\text{Data} | \text{Homo.}) \times P_1(\text{Homo.})}{P(\text{Data})} \\ &\propto (1 - \text{Prior}_1) \times \prod_{j=1}^{N_1} (1 - e_j)^{I_{(A_j=A)}} \times e_j^{I_{(A_j=a)}} \\ &\quad \times \prod_{s=2}^5 \left((1 - \text{Prior}_s) \prod_{j=1}^{N_s} (1 - e_j)^{I_{(A_j=A)}} (e_j)^{I_{(A_j=a)}} + \text{Prior}_s \times (0.5)^{N_s} \right), \end{aligned}$$

where $P_s(\cdot)$ denotes the probability in the context of strain s . On the basis of the posterior probabilities above, we classified each locus as homozygous or heterozygous for Wt0. If a locus was classified as heterozygous for Wt0, then it was assumed to have a propagated mutation in the rest of the strains. To call derived mutation in strains $s = 2, 3, 4, 5$, we use similar logic:

$$\begin{aligned} P_s(\text{Heter.} | \text{Data}) &= \frac{P_s(\text{Data} | \text{Heter.}) \times P_s(\text{Heter.})}{P_s(\text{Data})} \\ &\propto \text{Prior}_s \times (0.5)^{N_s} \\ P_s(\text{Homo.} | \text{Data}) &= \frac{P_s(\text{Data} | \text{Homo.}) \times P_s(\text{Homo.})}{P_s(\text{Data})} \\ &\propto (1 - \text{Prior}_s) \times \prod_{j=1}^{N_s} (1 - e_j)^{I_{(A_j=A)}} \times e_j^{I_{(A_j=a)}}. \end{aligned}$$

We use the posterior probabilities calculated above, to make a decision as to whether a site is called as heterozygous for a new mutation, heterozygous for a propagated mutation, or invariant for the four evolved strains: $s = 2, 3, 4, 5$. Specifically, if the posterior probability of heterozygosity was greater than 50% at a given position, then we classified the site as containing a SNP or indel. Visual inspection of the alignments for some of the inferred indel positions revealed that pairwise alignment of reads could induce false positives across multiple lines due to variations on how the alignment software interprets the alignment of different reads around a given position. These are characterized by one allele count being much smaller (but nonzero) compared to the other, across multiple strains. To bioinformatically cull such sites from our data set, we carried out an additional likelihood-ratio test for the allele frequencies to be equal (*i.e.*, a propagated SNP had to have statistical support for the model of 50% frequency across Wt0, Wt160, Mut2, Mut3, and Mut4; a derived SNP had to have statistical support for 50% in one of the evolved lines, and 0% in all the others). If the hypothesis of equality was rejected for an indel, we flagged it as low confidence (Figure 1).

We expected, on the basis of previous estimates of mutation rate in MMR defective strains, to find ~ 125 mutations for each of the MMR deficient strains (approximately one mutation per line generation). This corresponds to a prior mutation rate of 10^{-5} mutations/site/generation. However, we detected 12, 24, 40 mutations for each of the MMR-deficient strains, which yield mutation rates of 1×10^{-6} , 2×10^{-6} , and 3×10^{-6} in each line, respectively. Although our estimated prior values differ somewhat from the real data, the alignment analysis allowed us to calculate very accurate posterior subjective probabilities. This accuracy is due to the large number of observations and

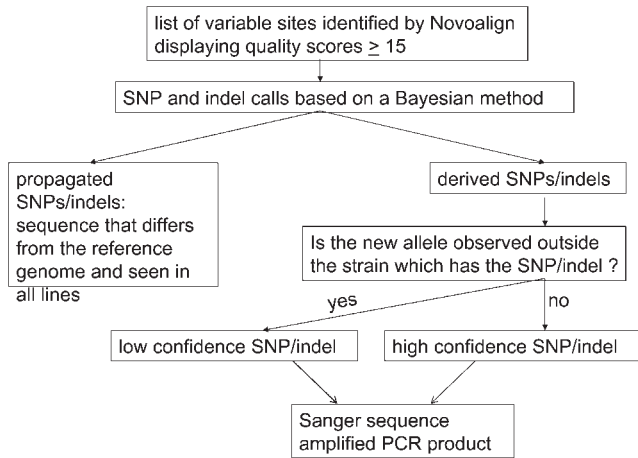


FIGURE 1.—Flow chart describing bioinformatic methods used to identify heterozygous mutations from Illumina GA whole-genome sequencing. See text for details.

has in practice made the influence of the prior negligible. Thus given the high coverage for the Mut lines, the difference in our prior estimates does not influence our analysis. Even with low coverage data where accurate estimates of prior are critical, a higher prior value would yield a larger number of false positives. The majority of mutations (and all low confidence mutations) were verified by Sanger sequencing, suggesting that false positives were rare, but we may have false negatives (*i.e.*, missed variants) due to the medium coverage ($\sim 8\text{--}20\times$) of the lines.

Simulation study: To estimate the false-positive and false-negative rates, as well as to check our bioinformatics and SNP/indel calling pipelines, we set up a simulation to test the accuracy of our Bayesian approach. We started with a complete genome of a yeast S288c strain (<http://genome.ucsc.edu/cgi-bin/hgGateway>; June 2008 assembly from the Saccharomyces Genome Database (SGD, <http://www.yeastgenome.org/>) and introduced SNPs and indels to simulate five strains: Wt0, Wt160, Mut2, Mut3, and Mut4. To simulate Wt0, we duplicated the S288c genome to create a diploid. We then randomly selected n_μ and n_d positions for SNPs and indels respectively. ($n_\mu = 2$, $n_d = 8$; the values of n_μ and n_d were chosen to mimic changes between S288c and the Wt0 strain used in the bottleneck experiment). One of the two copies of S288c was randomly selected to incur each SNP or indel. For an indel mutation, the nt in that copy was deleted, or a new randomly chosen allele was inserted after it. For a SNP position, the nt was randomly changed to another nt. The resulting two copies of the genome were defined as the Wt0 diploid. The other four strains were all simulated directly from Wt0 by introducing SNPs and indels in the two copies of Wt0. The mechanism of adding SNPs and indels was exactly as described above. The values of n_μ and n_d for each of the simulations are given below. These values mimic the number of mutations that were expected in the bottleneck experiments. One distinction between the simulations and the real data is that the SNPs and indels in the simulations were not introduced into HP tracts. As described below, we believe that our ability to detect indels in HP tracts is lower because indels in HP tracts can be identified only if the entire tract and sequence flanking both sides are present in a 36-nt read.

Next, we simulated 32-nt Illumina GA reads from each of the five strains by randomly choosing read-start positions and copying 32 nt of strain s starting from that position. For each strain, the number of reads simulated matches the coverage

	0 → Wt0	Wt0 → Wt160	Wt0 → Mut2	Wt0 → Mut3	Wt0 → Mut4
n_μ	2	1	25	25	25
n_d	1	1	100	100	100

achieved in the real sequencing experiment. We also simulated a quality score for each position of each read, following the error rate distribution given in DOHM *et al.* (2008). The reads were aligned with S288c using Novoalign (<http://www.novocraft.com>). Based on the alignment, we listed the allele counts and associated quality scores in each of the variable, potentially heterozygous, positions. We used this list as the input to a computer program created on the basis of our method of heterozygosity detection, which went through all the steps described in the last section. The rates of false positives and negatives (based on the output of the program) are given in Table 1. We believe that these rates are similar to those seen in the bottleneck experiment.

Verifying mutations identified using Bayesian method: Our method for heterozygous mutation calling from the whole-genome sequencing data yielded both low- and high-confidence predictions (see above). All low-confidence predictions (10 in total) were verified and either validated ($n = 4$) or disproved ($n = 6$) using Sanger sequencing. Briefly, to assay heterozygous mutations predicted from the whole-genome sequence data, genomic DNA was prepared from wild-type generation zero, and mutation accumulation lines Mut2, Mut3, and Mut4 using standard techniques. Approximately 400 bp of DNA flanking the predicted mutated site was amplified in all lines using PCR and Sanger sequenced at the Cornell CLC using an Applied Biosystems Automated 3730 DNA analyzer. The sequencing traces were all analyzed visually. A heterozygous base change mutation was confirmed if a doublet representing both alleles was observed only in the sequencing trace of the predicted Mut line, but all other lines showed only a singlet representing the parental allele. A heterozygous indel mutation was confirmed if the sequencing reaction failed (*i.e.*, tall singlet peaks fall to small doublet peaks or random noise) at the predicted location only in the predicted Mut line, but the sequencing reactions in all other lines were able to successfully sequence past the site.

For the high confidence predictions, 31 (of 65) were sequenced and verified using the methods described above. Of those 31 mutations, 10 were further verified by genotyping the haploid progeny of the diploid containing the heterozygous mutation via Sanger sequencing. Both alleles comprising the heterozygote were observed in the haploid progeny with the exception of the frameshift mutation in the essential *MDN1* gene. Six additional high-confidence predictions were also verified by genotyping the haploid progeny of the heterozygous diploid.

By Sanger sequencing of the diploid lines (see above) we also found and verified four heterozygous mutations that were detected in earlier, less accurate prediction protocols that were not found using the final more stringent prediction method.

RESULTS AND DISCUSSION

Identification of mutations in diploid bottleneck lines using maximum-likelihood and Bayesian methods:

One wild-type and three *mlh1-7^{ts}* lines (Mut2, Mut3, and Mut4) allowed to accumulate mutations for 160 generations were sequenced using the Illumina genome analyzer technology (MATERIALS AND METHODS; <http://www.illumina.com>). The wild-type progenitor of all the

TABLE 1
False-positive and -negative rates based on the simulation analysis

	False-positive rate (in units of no. of SNP calls)	False-negative rate (in units of no. of SNP calls)
Mutant	6×10^{-5}	0.030
Indel	0	0.089
Total	6×10^{-5}	0.078
Propagated	0	0
Derived	6×10^{-5}	0.091
Total	6×10^{-5}	0.078

See MATERIALS AND METHODS for details.

strains was also sequenced. The analysis was performed with three independent *mlh1-7^{ts}* lines to control for chance associations within an individual line and for mutations that could alter the mutation rate of a given line. The Mut2, Mut3, and Mut4 lines at generation 160 displayed 15.6, 7.1, and 2.5% spore viability, respectively (HECK *et al.* 2006). As shown below and in Tables 2 and 3, our data analysis indicated that the mutation spectra and rates in the three *mlh1-7^{ts}* lines were indistinguishable. In total, 25 million, out of 35 million sequenced, 36 nt sequence reads were uniquely mapped to the yeast genome, allowing up to two mismatches per read (MATERIALS AND METHODS). The wild-type and Mut2 generation 160 strains were sequenced to 9 \times and 8 \times average genome coverage depth, respectively. Mut3 (160) and Mut4 (160) were sequenced to average depths of 18 \times and 22 \times , respectively. We then developed and employed an “experiment aware” probabilistic framework using maximum-likelihood and Bayesian methods that utilized sequence coverage of the entire data set (\sim 70-fold; Figure 1; MATERIALS AND METHODS; DOHM *et al.* 2008). Briefly, the approach classifies each site in the yeast genome with uniquely mapping reads into one of three categories: (1) invariant across all strains, (2) heterozygous in the wild-type (and all derived strains), which we term “propagated” SNPs or indels, or (3) heterozygous in one of the mutant strains, which we term “derived” SNPs or indels. As described below, this method allowed us to pool experimental data across sequencing runs for all strains and detect with high reliability heterozygous SNPs (28 identified) and single-nt indels (48 identified) from the 36-nt read data set. This method was evaluated on simulated data sets and found to have a very low false-positive rate ($\sim 6 \times 10^{-5}$) and a false-negative rate of 0.08 within the unique mapping regions of the genome that contained at least sevenfold coverage (Table 1). The low false-positive rate was verified by PCR amplifying genomic fragments covering a specific mutation site and then confirming the presence of a heterozygous mutation by Sanger sequencing the fragment (MATERIALS AND METHODS). On the basis of

simulations, we estimated that the method, as applied to regions with at least sevenfold sequencing coverage, allowed us to detect heterozygous mutations in 60, 41, 69, and 84% of the total genome for the generation 160 wild-type, Mut2, Mut3, and Mut4 lines, respectively.

We did not detect any mutations in the wild-type generation 160 line, which was predicted on the basis of the previously calculated mutation rate of 3.3×10^{-10} mutations/base/generation (<1 expected; LYNCH *et al.* 2008). As shown in Tables 2 and 3 and Figure 2, only heterozygous mutations, composed of 28 base substitution and 48 single-nt indel mutations, were detected in the three MMR-defective lines. All of the mutations were unique between lines except for a single-nt deletion mutation between SGD (<http://www.yeastgenome.org>) coordinates 92,271–92,279 on chromosome 2, which occurred independently in both Mut2 and Mut3 (Table 2). All 48 indels, composed of 46 deletions and 2 insertions, occurred in HP tracts [47 poly(A) or (T) tracts, 1 poly(G) or (C) tract] between 5 and 13 bp long (Table 2). Due to the constraints of using 36-nt Illumina GA reads, we do not have the power to detect mutations in HP tracts larger than 13 nt, but <400 such tracts are present in the yeast genome. Visual inspection of the DNA sequences surrounding the indel mutations (\sim 400 bp; Figure 2) suggested that they were enriched for HP runs. These are primarily poly(dA:dT) tracts that are present in the yeast genome at a 20-fold higher frequency than poly(dG:dC) tracts. Consistent with this, the AT content of the genomic regions surrounding the indel mutations was significantly higher than that for unmutated HP regions (windows up to 500 bp; data not shown). Detailed bioinformatic and genetic analyses will be required to determine if this pattern is significant; however, a previous study (HARFE and JINKS-ROBERTSON 2000) showed that DNA polymerase slippage was not greatly influenced by sequence context, including nearby HP tracts.

Our analysis permitted the detection of up to two single-nt indels in a 36-nt read; these indels can be right next to each other to create a 2-nt indel or separated from each other. We assigned this limit because creating high-quality and unique alignments became very difficult when allowing indels larger than 2 nt. We were unable to detect indels of 2 nt in any of the lines. Such a result is not surprising due to previous studies of wild-type and MMR mutants analyzed for reversion of frameshift mutations in HP runs. In these studies the overwhelming majority of mutations involved single-nt deletions. For example TRAN *et al.* (1997) found that 225 of 227 reversions in +1 HP tracts in wild-type, polymerase proofreading, and mismatch repair mutants were due to deletions of a single nt. For –1 HP tracts, they found that 206 of 218 reversions were due to additions of a single nt. The remaining revertants in both HP tracts involved expansions or contractions of no greater than 2 nt in size.

TABLE 2
Genome location of mutations detected in the Mut2, 3, and 4 lines

Chromosome	SGD Position	HP tract	Strain	Gene	Amino acid change	Mutation	Distribution of sequence reads					
							A	G	T	C	Indel	
1	139,349–139,358	10	4			del	1		11			8
2	275,549–275,557	9	2			del			3			3
2	92,271–92,279	9	3			del	1		10			8
2	92,271–92,279	9	2			del			3			0
2	423,462–423,469	8	3			ins	8					16
2	662,560–662,569	10	4		fs; 103–106	del	6					5
2	653,035–653,045	11	4		fs; 176–178	del	13				15	9
3	212,451–212,457	7C	3		297; V to I	del	12	15				11
3	275,289	NA	4			G to A						
4	512,796	NA	2			A to T	8		7			
4	814,336	NA	2			G to A	10	6				
4	929,182–929,193	12	3			del	1		5			4
4	963,768	NA	3			T to G		12	18			
4	231,908–231,914	7	4		120; G to V	del	12					12
4	1,386,657–1,386,664	8	4			del			17			11
4	470,576–470,584	9	4			del	19					12
4	832,716–832,726	11	4			del	15					7
4	50,592–50,603	12	4			del			10			9
5	1,054,759	NA	4			A to T	9		11			
5	305,972	NA	2			C to A	5			5		
5	479,369	NA	2		1159; S to C	T to A	8		3			
5	225,319–225,327	9	3			Del			9			7
5	403,576	NA	3		258; A to E	G to T		8	14			
5	34,325–34,333	9	4			del						7
5	402,832–402,843	12	4			del	13					7
6	223,108–223,118	11	3			del	10					7
6	114,200–114,210	11	4			del	6					7
6	88,832	NA	4			G to T		15	13			
6	225,229	NA	4		504; D to E	G to C		9		8		
7	194,092–194,098	7	3		240; R to G	del			13			13
7	878,690–878,701	12	3		fs; 771–773	del	6					4
7	653,363–653,369	7	4			del	14					10
7	882,549–882,558	10	4			del	10					6
7	20,017–20,027	11	4			del	10					5
7	678,172–678,182	11	4			del	11					9
8	150,380–150,386	7	3			del	11					12
8	472,612–472,624	13	3			del			1			8
8	288,299	NA	3		172; K to K	C to T			10			1
8	370,253	NA	3		46; Y to F	A to T	15		5	6		1
9	270,327	NA	2		560; A to S	G to T			12	1		
9	375,856	NA	3		143; M to I	G to A	7	10	4			
9	199,995	NA	4		41; T to I	G to A	12	18				

(continued)

TABLE 2
(Continued)

Chromosome	SGD Position	HP tract	Strain	Gene	Amino acid change	Mutation	Distribution of sequence reads				
							A	G	T	C	Indel
10	445,012–445,020	9	3			del			10		6
10	131,051–131,059	9	4			del			11		15
10	469,684–469,694	11	4			del		19			6
11	162,688–162,695	8	4			del			9	1	7
11	403,466	NA	4	<i>YKL018C-A</i>	19; S to S	C to T			13	13	3
12	405,712–405,719	8	2	<i>YLR131C/ACE2</i>	fs; 369–371	del			4		5
12	32,320–32,330	11	3			del		9			5
12	964,065	NA	3	<i>YLR420W/URA4</i>	95; R to H	G to A	12				
12	1,009,007	NA	3	<i>YLR436C/ECM30</i>	746; I to V	T to C					
12	363,531–363,537	7	4	<i>YLR106C/MDN1</i>	fs; 68–70	del		11			13
12	201,846–201,856	11	4			del		14			9
12	1,047,741	NA	4	<i>YLR454W/FMP27</i>	1249; D to G	A to G	12				1
13	763,010–763,016	7	2	<i>SNR86</i> (small nucleolar RNA)		ins		5			5
13	241,855–241,867	13	3			del					6
13	311,843	NA	3			C to T			12	25	6
13	139,705–139,709	5	4	<i>YML067C/ERV41</i>	fs; 138–139	del			15		13
13	816,457–816,463	7	4	<i>YMR275C/BUL1</i>	fs; 706–708	del			11		9
14	761,792	NA	2	<i>YNR069C/BSC5</i>	267; V to V	C to T			8	6	1
14	222,733	NA	3	<i>YNL225C/CMN67</i>	580; V to M	C to T		1	17	10	10
14	435,595–435,601	7	4	<i>YNL101W/AVT4</i>	fs; 199–201	del		9			8
14	481,123–481,129	7	4			del		15			9
14	685,574–685,582	9	4			del			12	1	8
14	575,616–575,626	11	4			del					8
14	400,002	NA	4	<i>YNL121C/TOM70</i>	180; G to STOP	C to A		12		15	
14	734,521	NA	4	<i>YNR058W/BIO3</i>	77; L to L	A to G	10				5
15	854,146–854,153	8	2			del		4			7
15	874,052–874,057	6	3	<i>YOR296W</i>	fs; 1284–1286	del		10			7
15	767,667–767,673	7	3	<i>YOR228C</i>	fs; 36–38	del				10	7
15	822,829–822,835	7	3	<i>YOR267C/HRK1</i>	fs; 678–680	del				11	8
16	146,421–146,427	7	2	<i>YPL216W</i>	fs; 868–870	del		7			5
16	22,677	NA	4			C to T			20	14	
16	131,583	NA	4	<i>YPL222W/FMP40</i>	475; A to T	G to A		11	1		
16	509,632	NA	4	<i>YPL022W/RAD1</i>	980; A to S	G to T		19	18	1	
16	570,131	NA	4	<i>YPR007C/REC8</i>	415; S to M	G to A		11	11		

The type of mutation [base substitution, single-nt insertion (ins), single-nt deletion (del)] is shown, as well as the length of the HP tract that contains an indel. The specific Mut line (2, 3, or 4) is indicated under “strain.” All HP tracts were poly(A) or poly(T) except for the mutation in chromosome 3 at 212,451–212,457, which involved a poly(C) tract. For mutations that occurred within an open reading frame, both the gene name and predicted amino acid changes (fs, frameshift) are provided. NA, not applicable. Coordinates are presented as shown in the SGD (<http://www.yeastgenome.org/>). The number and distribution of the sequence reads are presented for each mutation. The frameshift mutation in *YLR106C/MDN1* conferred a recessive lethal phenotype (data not shown).

TABLE 3
Mutation rates for Mut2, Mut3, and Mut4 lines grown in bottlenecks for 160 generations

Base substitution mutations				
Strain	No. mutations	% genome $\geq 7\times$ coverage	Genome Size (bp) adjusted	Mutation rate (per base per gen $\times 10^{-9}$)
Mut2	6	41	9,898,136	3.8
Mut3	9	69	16,657,838	3.4
Mut4	13	84	20,279,107	4.0
Average				3.7

Single-nucleotide indel mutations in 5- to 13-nt HP tracts				
Strain	No. mutations	No. HP tracts $>7\times$ coverage	Mutation rate (per HP tract/generation $\times 10^{-7}$)	
Mut2	6	57,502	6.5	
Mut3	15	99,714	9.4	
Mut4	27	122,816	14	
Average			10	

Single-nucleotide indel mutations in 8- to 13-nt HP tracts				
Strain	No. mutations	No. HP tracts $>7\times$ coverage	Mutation rate (per HP tract/generation $\times 10^{-7}$)	
Mut2	4	2,820	89	
Mut3	10	7,054	89	
Mut4	19	8,696	140	
Average			110	

The base substitution mutation rate was determined by calculating the percentage of the genome in which at least sevenfold DNA sequencing coverage to unique regions was obtained. This was done because our statistical analysis did not have sufficient power to reliably detect heterozygous mutations in regions with lower coverage. This information was used to calculate the mutation rate on the basis of the following formula: (number of mutations)/(160 generations)/(adjusted genome size), with the diploid *S. cerevisiae* genome size determined as 24,141,794 bp (<http://www.yeastgenome.org/>). To obtain indel mutation rates, we first determined the number of HP tracts of a given length in unique regions of the genome which had \geq sevenfold sequence coverage. We then used the following equation to calculate mutation rate: (number of indels)/(160 generations)/(number of HP tracts with \geq sevenfold coverage).

The predominance of single-nt deletions over single-nt insertions and base substitutions was similar to previous reports for the mutational spectra in reporter genes in MMR null mutants (MARSISCHKY *et al.* 1996; TRAN *et al.* 1997, 2001; DENVER *et al.* 2005). The average mutation rate in the 5- to 13-bp HP tracts was 1.0×10^{-6} /HP tract/generation (Table 3). The rate was an order of magnitude greater (1.1×10^{-5}) if only runs between 8 and 13 bp long were considered (Table 3). These values approach the rates seen in MMR-defective yeast (*mlh1*, *msh2*) containing reporters bearing 10-bp poly(T) (2.8×10^{-4} ; TRAN *et al.* 1997) and 10-bp poly(A) (7.3×10^{-5} ; GRAGG *et al.* 2002) tracts. Low-sequence coverage provides one explanation for why the rate is lower than those seen previously in reporter assays. In our analysis, indels in HP tracts can be identified only if the entire tract and sequence flanking both sides are present in a 36-nt read; the longer the HP tract, the less likely it is to obtain reads that cover the entire tract. Thus higher sequence coverages are required to identify indels in HP tracts. Consistent with this, a higher indel mutation rate was seen in lines that had higher sequencing coverage (Table 3). In contrast, SNPs that occur outside of an HP tract should

not be as affected by sequence coverage (aside from the relationship between coverage and probability of detecting sufficient copies of the alternate base to reliably make a call). This was seen for the analysis of base substitutions (Table 3).

The average rate of base substitution mutations in *mlh1-7^s* was 3.7×10^{-9} mutations/base/generation (Table 3), which is 11-fold higher than the base substitution rate observed in wild-type haploid strains (LYNCH *et al.* 2008). Of the 28 base substitution mutations detected in the Mut2–4 lines, 16 were transitions and 12 were transversions (Table 2). Nineteen of these mutations resulted in a change from a G–C to an A–T base pair, whereas only 4 were in the opposite direction. This overall mutational bias toward A–T base pairs was seen and discussed previously (*e.g.*, LYNCH *et al.* 2008; DENVER *et al.* 2009; KEIGHTLEY *et al.* 2009). The modest increase that we observed in the base substitution rate in MMR defective strains is significantly lower than predicted (~ 100 -fold increase for base substitutions and frameshifts; DENVER *et al.* 2005; IYER *et al.* 2006). We suggest two reasons for these differences. First, our measurements were determined from a genome-wide measurement rather than by

Indel site

aattggatttgttagcaaatcagccttgcctgctcgcattctctctg**TTTTTTTTT**acttctcggctcattg**AAAAA**tcctgacgaaaatatttcaaggtccta
 ttctcgtccctagaagg**AAAAAAAAA**gagaagtttctcgcagagtg**AAAAAAAAA**gctcaagaaaagatcctacaagaacaaattatgcccgaataatgct
 tgaagactacggtgaagactacggagaaaacggacaaggcagatg**CCCCCCC**aggag**AAAAA**ctggagtcgaacttttcagggatctacgtgttcgcatg
 ttgtaaatcagt**AAAAAAAAAAAAA**ttaacag**TTTTTTTTTTT**ca**TTTTTTTTTTTTT**attcttatttatgtagtatactttattatatttctcttaattat
 agaag**AAAAA**ctacacgggcactaacatgttaaatatgataat**TTTTTTTTTTT**ataagagaatcactaccaagttacctgaactacgtcaaggaaaagcc
 cttccatttttggggcacataaggagg**AAAAA**gaaaatataaag**AAAAAAAAAAAAA**gaaat**AAAAAAA**gaaaacgggtactggatattgagataaatttctc
 tctctatggtaccagtcgatataatctgtaaatgaaacaatcc**TTTTTTTTT**gaccgtctctccaaacacgtgccaaagcctgtgtgacgagcaggat
 ccacttcacgcggttcgcatatttgtccagcctc**TTTTT**cc**AAAAAAAAAAAAA**taataataaaatgaaacggacaggaattgaacctgcaaccctt
 tcatagcgtttgacctaatcaggtataatataaag**AAAAAAA**gac**AAAAAAA**taggtatgccaacacagagcagaaccttaacaatttctctagatacaa
 attaaaggtttataaacggggtcccgaactatctccattc**TTTTTTTTTTTTT**actgaagaaaggcaaacggcgacacaaaatcatgaagtcagaa
 tatttatagtagtacaaggaatc**TTTTTTTTT**gattccacatggc**TTTTTTTTTTT**agaagattgcatctcgttacatgatcaacaatccattgtaacggt
 taaccatagttcttggaaatgtcaactgagggatattgcactc**AAAAAAAAAAAAA**ttattaaatgagactatatacagtgagcacaacctgtctaataca
 aagctgattaacagaattagaagatctgcagtaactcgttttc**TTTTTTTTTTTTT**cattaatttatatgctatccttt**AAAAA**tagacatgtcatttca
 atctccgacatgaactcctgg**AAAAA**ggatcaggaatgccaaat**TTTTTTTTT**ctcttattcaggattaccatctccttcggggctatgtgcattaacatc
 aagcatgtggtagtagacgtaataagacggtatgggccatgaataa**TTTTTTTTT**atggtagggctcctctgtgctgttgactcgatttattatcggaagcgt
 ctatg**AAAAA**tagaagagcggtagtagagcggctactcctacagg**AAAAAAA**gggcagaactgaagtccttttatataatatacatgggggtgctgacta
 accatagtgtagcagataatcaagaagtgaactccttctctat**TTTTTTTTTTT**aattg**AAAAA**ttccttctctatagcgtatagaatataatgttaca
 gaggcttgacggtgtctttcccgctgttccctctccctc**AAAAAAAAAAAAA**cttttacaatgactaaaagaaaatattgtcatttatatggtttt
 tgactcttataataaatagctcagacgtaactgccaaaatgt**AAAAAAAAAAAAA**g**AAAAA**ctattagaaaatattgtgc**AAAAAAAAAAAAA**ggtagtca
 cgcttactgtgagaataatcaaccttatagcatatgtttctc**TTTTTTTTTTTTT**gcttatggagataatgaacattgtcacacatgaacaagtggtagt
 ttgtag**AAAAA**gaacctaaatcacg**AAAAAAA**tagtatgctag**AAAAAAA**tattatgaacgtcaata**TTTTT**cttctcatgggaccaaacacatc
 tgacattcatataatatacaataagtatagtttataatagatc**AAAAAAAAA**gctctgtattaccatgtaaatagttacattgtgtataactctcgtag
 aaagcctgaaaattgaagattacgcagtatcagcgtaaacc**AAAAAAAAAAAAA**tacaatcgggtaaagtgctgttttaatttagaataatattgtttc
 aaagggcattggtgaattatgacatgaactgtactcactatgacg**TTTTTTTTT**actgc**TTTTT**cccatcctgcaagc**AAAAAAA**gcaagtcgacta
 ctactcattgttgcctgtatctccgaataggactgataagtgat**AAAAAAAAA**gaaaactggggttagcttgagcagcttcttctactaaaggaa
 atcatataatg**AAAAA**tgagaat**AAAAAAA**caatgataacaatg**AAAAAAAAAAAAA**tcaataaagagaaaataaacttatttcttaagtagtccct
 atataagtagtttataagctagctactactcagaaaataattc**AAAAAAAAAAAAA**gaagtgggcactttaaataagagattagctatagtaaatctgta
 ggttagacttgagatagatgtacctggagagaaaatctataat**AAAAAAAAAAAAA**acttctcaggtcagatattc**TTTTTTTTT**gtccctaat**AAAAA**gaaa
 aaaagtattcgtctctttgtgtttatgaaaagggaacgtgatat**AAAAAAA**catccttgggtgggacatggcctttgtttagagaatggttatcac
 aaatgagaaacgtgggagatgtcaacctcgctgactcacc**AAAAAAAAAAAAA**tcgc**AAAAA**taagcggg**AAAAA**tgctcagattctaaagttcacaacc
 ttcagaatcgttatcctggcgg**AAAAA**tcattgttaacttt**AAAAAAAAAAAAA**gccaatattcccacaaatattaagagcgcctccattattaactaaa
 atgactctctttaggtaaa**TTTTT**attagctttatttggtag**TTTTTTTTT**atcttgggcatgtacgaagagcaagtagcttttaactatacatat
 atttccctttagctcctgaaaatactatcatggcgtaaggg**AAAAAAAAAAAAA**ctaactctacgtggtctctataagatgaaacagcaagctgcataca
 ttattatagtagttatatttctgatgttgggtatagcaagcagc**TTTTTTTTT**ctctttaccataatcattgtaccaggaattttgttcataattatt
 ggggttgcctgatagcagaataaaagtacagctctaggtctat**AAAAAAAAAAAAA**tggttaaag**AAAAA**tatacaggttggtatagaaatcatttaatt
 taatc**TTTTTTTTTTT**ggagtcctttatccaacgtgccaactctc**AAAAAAA**atggttgaagaaacatgaccagtgatggaacggctatgtagta
 gtgccttgttagactccttcatc**AAAAA**tgacctagtatctag**TTTTT**ctctaaattctgctgggatagcttcccctaaagatttcatccaattccgg
 atcggattctatgcaagtcattgttaaatcaag**TTTTT**ccaatgg**TTTTTTTTT**actgctcoggaatcaactcagacaatctcaaccaatttcttggctg
 ggtttaaagtccccgggtggctcagaagggaaatttattgtcaac**AAAAAAA**ggcaagaacatcaataaattgattccgctagctcagacttcaactcgc
 ttt**CCCCC**gctcagattctcttagaacattacggaataaagg**AAAAAAA**gactggagcatcgaatctgtagact**AAAAA**ggtaatgacgcgttctt
 aaaggagatgtttgtatgat**CCCCC**agctcaaatgcatag**AAAAAAAAAAAAA**tcccgccttatattcatgatcttccacctcttagtctctggcca
 ggacagacggagctcacacaaaagttacaataatgctggaacc**TTTTTTTTTTT**cctttgaaattcttgtaaaccagtaattcgttttgaaacaccggtt

FIGURE 2.—The 100-bp region surrounding indel mutations in the Mut3 and Mut4 lines. The locations of the indel mutations are indicated in black boldface type. HP runs of ≥ 5 in this window are color coded as shown: red, A_n ; blue, T_n ; green, C_n .

extrapolation from a few marker loci. Second, the *mlh1-7^s* allele is not a complete null mutation. It phenocopies the *mlh1Δ* phenotype in the *CAN1* mutational assay, but has a fourfold lower mutation rate than *mlh1Δ* in the *lys2-A₁₄* reversion assay (HECK *et al.* 2006; data not shown). Because *mlh1-7^s* strains display residual DNA repair, it is possible that there is a bias toward the repair of specific mismatches in these strains. While we cannot rule this out, the fact that the mutation signature seen in *mlh1-7^s* appeared indistinguishable from *mlh1* null strains argues against such a possibility (MARSISCHKY *et al.* 1996; TRAN *et al.* 1997, 2001). Finally, we cannot rule out the possibility that mutation rates in MMR-defective strains are different in haploid *vs.* diploid yeast, although a recent analysis of mutation rates in diploid bottleneck lines showed that wild-type diploid yeast displayed an estimated base substitution rate that was very similar to that reported previously for haploid yeast (LYNCH *et al.* 2008; NISHANT *et al.* 2010).

Because the three lines showed viability that ranged from 2.5 to 15.6%, we expected to identify mutations that conferred a lethal phenotype. We examined whether any of the mutations that mapped to open reading frames in the Mut4 line (2.5% viability) were not detected in haploid progeny. This was done by sequencing DNA surrounding a particular mutation in 20 viable spore clones obtained by sporulating the Mut4 generation 160 line. Of these 14 mutations, only the frameshift mutation in *MDN1* was not detected, consistent with previous work showing that *mdn1Δ* mutants are inviable (GIAEVER *et al.* 2002). While it is unclear how many mutations would confer lethality in the absence of other mutations, the assortment of 5 independent lethal mutations would result in 3% spore viability, similar to that seen in the Mut4 line. We hypothesize that other lethal mutations were not identified in Mut4 and other lines because:

1. A large number of frameshift mutations in HP tracts may not have been detected because indels can be identified only if the entire tract and sequence flanking both sides are present in a 36-nt read. Identifying indels in HP tracts is very challenging using short-read sequencing. However, increasing sequence coverage and using paired-end reads of a larger size (~180 bp) should provide a good test of this idea.
2. Our sequence analysis did not cover the entire genome (84% for Mut4).
3. While previous CGH and PFGE analyses (~1-kb resolution; HECK *et al.* 2006) did not reveal rearrangements, it is possible that mutations that involved indels larger than two nt and smaller than 1 kb occurred. However, we find this to be less likely because a previous analysis of mutation spectra in MMR mutants indicated that indels greater than two nt are extremely rare (TRAN *et al.* 1997).

Closing thoughts: In the *S. cerevisiae* S288c haploid genome there are over 77,425 HP tracts five nt or greater. Frameshift mutations in coding regions that disrupt protein function are likely to have significant effects on organism fitness. In wild-type yeast, insertion/deletion mutations appear to be relatively rare compared to base substitutions; comparative analyses of multiple domestic and wild yeast strains identified ~14,000 indels compared to ~235,000 SNPs (WEI *et al.* 2007; LITI *et al.* 2009). In contrast, MMR mutants display a strong bias toward frameshifts over base substitutions in the genome. Thus our data, together with previous work, illustrate the critical role that MMR plays in preventing frameshifts in HP tracts across the genome.

We thank Amit Indap and the Cornell Core Laboratory Center (CLC), especially Peter Schweitzer, James VanEe, and Tom Stelick, for preparing samples for Illumina GA sequencing and bioinformatic analyses; Julie Heck and K. T. Nishant for technical advice and providing unpublished data, and the Alani, Bustamante and Aquadro laboratories and Nadia Singh and Dan Barbash for comments on the manuscript. E.A., S.Z., and A.D. were supported by National Institutes of Health (NIH) GM53085. S.Z. was also supported by a Cornell Presidential Fellowship and an NIH training grant in Genetics and Development. A.D. was also supported by an NIH training grant in Biochemistry, Molecular and Cell Biology. X.M. and C.D.B. were supported by NSF 0606461 and NSF 0701382. A.R. was supported by NIH grant RO1 HG003229. R.H. was supported by a National Science Foundation (NSF) Minority Postdoctoral Fellowship. Z.G. was supported by a startup fund from Cornell and NSF DEB-0949556. B.B. was supported by an NIH training grant to the Tri-Institutional Training Program in Computational Biology and Medicine.

LITERATURE CITED

- ARNDT, P. F., T. HWA and D. A. PETROV, 2005 Substantial regional variation in substitution rates in the human genome: importance of GC content, gene density, and telomere-specific effects. *J. Mol. Evol.* **60**: 748–763.
- DATTA, A., and S. JINKS-ROBERTSON, 1995 Association of increased spontaneous mutation rates with high levels of transcription in yeast. *Science* **268**: 1616–1619.
- DENVER, D. R., S. FEINBERG, S. ESTES, W. K. THOMAS and M. LYNCH, 2005 Mutation rates, spectra and hotspots in mismatch repair-deficient *Caenorhabditis elegans*. *Genetics* **170**: 107–113.
- DENVER, D. R., P. C. DOLAN, L. J. WILHELM, W. SUNG, J. I. LUCAS-LLEDÓ *et al.*, 2009 A genome-wide view of *Caenorhabditis elegans* base-substitution mutation processes. *Proc. Natl. Acad. Sci. USA* **106**: 16310–16314.
- DOHM, J. C., C. LOTTAZ, T. BORODINA and H. HIMMELBAUER, 2008 Substantial biases in ultra-short read data sets from high-throughput DNA sequencing. *Nucleic Acids Res.* **36**: e105.
- GIAEVER, G., A. M. CHU, L. NI, C. CONNELLY, L. RILES *et al.*, 2002 Functional profiling of the *Saccharomyces cerevisiae* genome. *Nature* **418**: 387–391.
- GRAGG, H., B. D. HARFE and S. JINKS-ROBERTSON, 2002 Base composition of mononucleotide runs affects DNA polymerase slippage and removal of frameshift intermediates by mismatch repair in *Saccharomyces cerevisiae*. *Mol. Cell. Biol.* **22**: 8756–8762.
- HARDISON, R. C., K. M. ROSKIN, S. YANG, M. DIEKHANS, W. J. KENT *et al.*, 2003 Covariation in frequencies of substitution, deletion, transposition, and recombination during eutherian evolution. *Genome Res.* **13**: 13–26.
- HARFE, B. D., and S. JINKS-ROBERTSON, 2000 Sequence composition and context effects on the generation and repair of frameshift intermediates in mononucleotide runs in *Saccharomyces cerevisiae*. *Genetics* **156**: 571–578.
- HAWK, J. D., L. STEFANOVIC, J. C. BOYER, T. D. PETES and R. A. FARBER, 2005 Variation in efficiency of DNA mismatch repair at different sites in the yeast genome. *Proc. Natl. Acad. Sci. USA* **102**: 8639–8643.
- HECK, J. A., D. GRESHAM, D. BOTSTEIN and E. ALANI, 2006 Accumulation of recessive lethal mutations in *Saccharomyces cerevisiae* *mh1* mismatch repair mutants is not associated with gross chromosomal rearrangements. *Genetics* **174**: 519–523.
- IYER, R. R., A. PLUCIENNIK, V. BURDETT and P. L. MODRICH, 2006 DNA mismatch repair: functions and mechanisms. *Chem. Rev.* **106**: 302–323.
- KADYROV, F. A., L. DZANTIEV, N. CONSTANTIN and P. MODRICH, 2006 Endonucleolytic function of MutLalpha in human mismatch repair. *Cell* **126**: 297–308.
- KADYROV, F. A., S. F. HOLMES, M. E. ARANA, O. A. LUKIANOVA, M. O'DONNELL *et al.*, 2007 *Saccharomyces cerevisiae* MutLalpha is a mismatch repair endonuclease. *J. Biol. Chem.* **282**: 37181–37190.
- KEIGHTLEY, P. D., U. TRIVEDI, M. THOMSON, F. OLIVER, S. KUMAR *et al.*, 2009 Analysis of the genome sequences of three *Drosophila melanogaster* spontaneous mutation accumulation lines. *Genome Res.* **19**: 1195–1201.
- KUNKEL, T. A., and D. A. ERIE, 2005 DNA mismatch repair. *Annu. Rev. Biochem.* **74**: 681–710.
- LITI, G., D. M. CARTER, A. M. MOSES, J. WARRINGER, L. PARTS, *et al.*, 2009 Population genomics of domestic and wild yeasts. *Nature* **458**: 337–341.
- LYNCH, M. W., SUNG, K. MORRIS, N. COFFEY, C. R. LANDRY *et al.*, 2008 A genome-wide view of the spectrum of spontaneous mutations in yeast. *Proc. Natl. Acad. Sci. USA* **105**: 9272–9277.
- MARSISCHKY, G. T., N. FILOSI, M. F. KANE and R. KOLODNER, 1996 Redundancy of *Saccharomyces cerevisiae* MSH3 and MSH6 in MSH2-dependent mismatch repair. *Genes Dev.* **10**: 407–420.
- MATASSI, G., P. M. SHARP and C. GAUTIER, 1999 Chromosomal location effects on gene sequence evolution in mammals. *Curr. Biol.* **9**: 786–791.
- MCCULLOCH, S. D., and T. A. KUNKEL, 2008 The fidelity of DNA synthesis by eukaryotic replicative and translesion synthesis polymerases. *Cell Res.* **18**: 148–161.
- MODRICH, P., and R. S. LAHUE, 1996 Mismatch repair in replication fidelity, genetic recombination, and cancer biology. *Ann. Rev. Biochem.* **65**: 101–133.
- NISHANT, K. T., N. D. SINGH and E. ALANI, 2009 Genomic mutation rates: what high-throughput methods can tell us. *BioEssays* **31**: 912–920.
- NISHANT, K. T., W. WEI, E. MANCERA, J. L. ARGUESO, A. SCHLATTL *et al.*, 2010 The baker's yeast diploid genome is remarkably stable in vegetative growth and meiosis. *PLoS Genetics* (in press).
- STAMATOYANNOPOULOS, J. A., I. ADZHUBEI, R. E. THURMAN, G. V. KRYUKOV, S. M. MIRKIN *et al.*, 2009 Human mutation rate

- associated with DNA replication timing. *Nat. Genet.* **41**: 393–395.
- TEYTELMAN, L., M. B. EISEN and J. RINE, 2008 Silent but not static: accelerated base-pair substitution in silenced chromatin of budding yeasts. *PLoS Genet.* **4**: e1000247.
- TRAN, H. T., J. D. KEEN, M. KRICKER, M. A. RESNICK and D. A. GORDENIN, 1997 Hypermutability of homonucleotide runs in mismatch repair and DNA polymerase proofreading yeast mutants. *Mol. Cell. Biol.* **17**: 2859–2865.
- TRAN, P. T., J. A. SIMON and R. M. LISKAY, 2001 Interactions of Exo1p with components of MutLa in *Saccharomyces cerevisiae*. *Proc. Natl. Acad. Sci. USA* **98**: 9760–9765.
- WASHIETL, S., R. MACHNÉ and N. GOLDMAN, 2008 Evolutionary footprints of nucleosome positions in yeast. *Trends Genet.* **24**: 583–587.
- WEI, W., J. H. MCCUSKER, R. W. HYMAN, T. JONES, Y. NING *et al.*, 2007 Genome sequencing and comparative analysis of *Saccharomyces cerevisiae* strain YJM789. *Proc. Natl. Acad. Sci. USA* **104**: 12825–12830.
- WOLFE, K. H., P. M. SHARP and W. H. LI, 1989 Mutation rates differ among regions of the mammalian genome. *Nature* **337**: 283–285.

Communicating editor: S. KEENEY