

# Pervasive genetic hitchhiking and clonal interference in forty evolving yeast populations

Gregory I. Lang<sup>1\*†</sup>, Daniel P. Rice<sup>2\*</sup>, Mark J. Hickman<sup>3</sup>, Erica Sodergren<sup>4</sup>, George M. Weinstock<sup>4</sup>, David Botstein<sup>1</sup> & Michael M. Desai<sup>2</sup>

The dynamics of adaptation determine which mutations fix in a population, and hence how reproducible evolution will be. This is central to understanding the spectra of mutations recovered in the evolution of antibiotic resistance<sup>1</sup>, the response of pathogens to immune selection<sup>2,3</sup>, and the dynamics of cancer progression<sup>4,5</sup>. In laboratory evolution experiments, demonstrably beneficial mutations are found repeatedly<sup>6–8</sup>, but are often accompanied by other mutations with no obvious benefit. Here we use whole-genome whole-population sequencing to examine the dynamics of genome sequence evolution at high temporal resolution in 40 replicate *Saccharomyces cerevisiae* populations growing in rich medium for 1,000 generations. We find pervasive genetic hitchhiking: multiple mutations arise and move synchronously through the population as mutational ‘cohorts’. Multiple clonal cohorts are often present simultaneously, competing with each other in the same population. Our results show that patterns of sequence evolution are driven by a balance between these chance effects of hitchhiking and interference, which increase stochastic variation in evolutionary outcomes, and the deterministic action of selection on individual mutations, which favours parallel evolutionary solutions in replicate populations.

Evolutionary adaptation is driven by the accumulation of beneficial mutations. The traditional view is that these dynamics are dominated by rare beneficial ‘driver’ mutations that occasionally survive drift and increase in frequency until they fix (a ‘selective sweep’)<sup>9,10</sup>. This implicitly assumes that at most a single beneficial mutation is present in the population at once. However, recent experiments have shown that even for modestly sized populations of microbes and viruses, beneficial mutation rates are large enough<sup>11,12</sup> that multiple driver mutations spread simultaneously, an effect known as ‘clonal interference’. This means that the fate of each mutation depends not only on its own effect on fitness but also on the rest of the variation in the population: neutral or deleterious mutations can fix if they occur in very fit genetic backgrounds, and beneficial mutations occurring in unfit lineages cannot succeed<sup>13–17</sup>.

Recent work has uncovered important consequences of these clonal interference effects. For example, interference alters the rate of adaptation<sup>18,19</sup>, the fate of marked lineages<sup>20,21</sup>, and the distribution of fitness effects of fixed mutations<sup>12,16,17</sup>. However, the underlying basis of these effects at the genomic sequence level has not been observed directly. A number of questions arise regarding the fate of those mutations that occur, the changes in frequency of each mutation over time, and the way in which these sequence-level dynamics determine the rate and repeatability of adaptation. Recent studies have sequenced clones or whole-population samples from microbial evolution experiments<sup>6,22–25</sup>, but apart from studies in viral systems<sup>26–28</sup>, this work has been limited to individual clones or populations or to widely separated time points that lack the temporal resolution to address these questions.

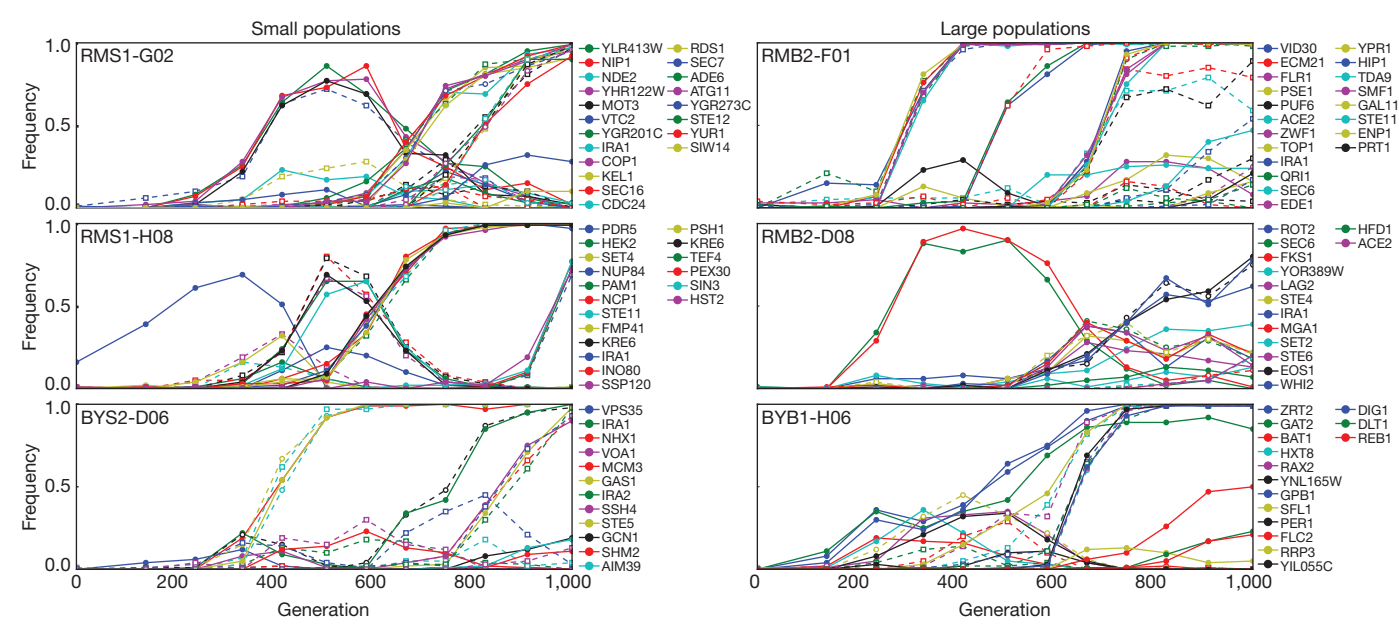
Here we describe the first direct and detailed view of the dynamics of genomic sequence evolution across many replicate microbial populations. In previous work<sup>21</sup>, we adapted approximately 600 replicate haploid yeast populations to growth in rich medium for 1,000 generations, half at ‘large’ (10<sup>6</sup>) and half at ‘small’ (10<sup>5</sup>) population sizes. Here we report the sequencing of whole-population samples from 40 of these populations (14 large and 26 small), chosen because we previously followed a single marker above a frequency of 0.1 (ref. 21). Each population was sequenced to 100-fold depth at 12 time points (approximately every 80 generations) for a total of 480 sequenced time points. Distinguishing mutations from sequencing errors in this whole-population sequence data is challenging. However, the high temporal resolution of our data permits the identification of mutations even at relatively low frequency by leveraging multiple time points. We developed two independent pipelines for this purpose, which rely on the fact that real mutations (but not sequencing or alignment errors) have frequencies that are correlated through time (Methods). This strategy allowed us to identify mutations that rose to a frequency of at least approximately 0.1 and to track these mutations through the rest of the timecourse. Across the 40 populations, we identified a total of 1,020 mutations, 253 of which fix; we annotated each to a gene or intergenic region and classified coding mutations as synonymous or nonsynonymous (Supplementary Table 1). Figure 1 shows six representative populations; the remaining populations exhibit similar patterns (Supplementary Fig. 1).

Averaged across all 40 populations, the rate at which mutations appeared and subsequently went extinct or fixed was constant through 1,000 generations (Fig. 2a). The average within-population polymorphism increased steadily through the first 600 generations, before saturating (Fig. 2a). In individual populations, however, the appearance of mutations is highly punctuated. This leads to the most striking feature of our results: selective sweeps are rarely single mutation or single phase events. Instead, mutations often move through the populations as temporal clusters (‘cohorts’) of functionally unrelated mutations, synchronously escaping drift and tracking tightly with one another through time. We quantify this temporal clustering of mutations in Fig. 2b, showing that it leads to a significant overrepresentation of time points at which either many or no mutations appeared, compared to the null expectation of mutations reaching detectable frequency at a constant rate ( $P < 10^{-6}$ ).

Multiple mutations or cohorts of mutations are often present simultaneously, as is apparent in Fig. 1, and selective sweeps are often ‘nested’; that is, one sweep initiates before the preceding sweep has completed. Cohorts and nesting of mutations are forms of genetic hitchhiking, in which individual mutations are helped (or hindered) by the genetic background in which they happen to arise. This includes both hitchhiking of likely neutral synonymous mutations, as well as ‘quasi-hitchhiking’ of multiple beneficial mutations that act together as co-drivers. In addition, frequent interference between competing cohorts often leads to

<sup>1</sup>Lewis-Sigler Institute for Integrative Genomics and Department of Molecular Biology, Princeton University, Princeton, New Jersey 08544, USA. <sup>2</sup>Departments of Organismic and Evolutionary Biology and of Physics, and FAS Center for Systems Biology, Harvard University, Cambridge, Massachusetts 02138, USA. <sup>3</sup>Departments of Biological Sciences and Chemistry & Biochemistry, Rowan University, Glassboro, New Jersey 08028, USA. <sup>4</sup>The Genome Institute, Washington University School of Medicine, St. Louis, Missouri 63108, USA. <sup>†</sup>Present address: Department of Biological Sciences, Lehigh University, Bethlehem, Pennsylvania 18015, USA.

\*These authors contributed equally to this work.



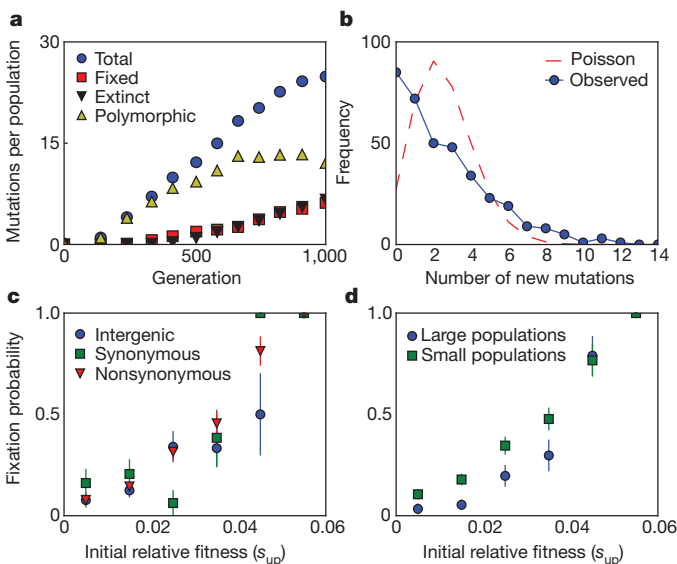
**Figure 1 | The fates of individual spontaneously arising mutations.** We show the frequency of all identified mutations through 1,000 generations in 6 of the 40 sequenced populations. Non-synonymous mutations are solid lines with solid circles, and synonymous and intergenic mutations are dotted lines with

open circles and squares, respectively. Populations in the left and right columns were evolved at small ( $10^5$ ) and large ( $10^6$ ) population sizes, respectively. We observe qualitatively similar patterns in the other populations (Supplementary Fig. 1).

the extinction of beneficial mutations even after they reach substantial frequency. Drawing from the full aggregate data set as well as individual ‘case study’ populations, we now show how this pervasive hitchhiking and interference strikes a balance between chance and determinism in governing evolutionary outcomes.

To investigate the repeatability of adaptation and identify those mutations that are driving adaptation, we looked for genes in which we observed mutations more often than expected by chance. Of the 995 nuclear mutations we identified, 723 fall within coding regions. If these mutations were distributed randomly over the 5,799 yeast genes, we expect only two genes with three or more mutations. Instead, we find

**24 genes hit three or more times** (Table 1, Supplementary Table 2 and Supplementary Fig. 2). This parallelism is at the gene level; mutations in different populations are different at the nucleotide level, with four exceptions (Methods). These 24 putative drivers represent approximately 0.6% of the yeast genome by size but account for 14% of the observed mutations, and are more likely to fix in the population (52/140, 37%) compared to all other nonsynonymous mutations (110/476, 23%,  $P < 0.005$ ). Only 1 of the 141 mutations in these putative drivers is synonymous (<1%), compared to 19% for the 472 mutations that fall in genes that are hit only once (Supplementary Table 3). Putative drivers are similarly depleted for missense mutations, and are enriched for nonsense and frameshift mutations (Supplementary Table 3). This mutational spectrum differs between functional categories of putative driver mutations. For genes in the mating pathway and



**Figure 2 | Statistical analysis across 40 replicate populations.** **a**, The per-population number of total mutations, fixed mutations, extinct mutations and mutations that are currently polymorphic over the course of the 1,000 generations. **b**, The distribution of the number of new mutations detected at each time point (solid blue line; see Methods for details) and a Poisson distribution with the same mean (dashed red line). **c**, **d**, Mutation fixation probability as a function of initial relative fitness. Data are mean  $\pm$  s.e.m.

Table 1   Repeatedly hit genes are putative drivers of adaptation			
Gene	Hits	Fixed	Biological process*
<i>IRA1</i>	21	10	Negative regulator of Ras
<i>ROT2</i>	11	2	Cell wall biogenesis
<i>YUR1</i>	11	5	Cell wall biogenesis
<i>ACE2</i>	9	4	Cytokinesis
<i>STE11</i>	9	1	Mating
<i>STE12</i>	9	2	Mating
<i>PDR5</i>	8	5	Multidrug transport
<i>WHI2</i>	7	2	General stress response
<i>STE4</i>	6	1	Mating
<i>IRA2</i>	5	3	Negative regulator of Ras
<i>KRE6</i>	4	1	Cell wall assembly
<i>SFL1</i>	4	1	Regulation of flocculation genes
<i>STE5</i>	4	3	Mating
<i>ANP1</i>	3	1	Protein glycosylation
<i>CNE1</i>	3	2	Protein folding
<i>GAS1</i>	3	3	Cell wall assembly
<i>GCN1</i>	3	1	Regulation of translation
<i>GPB1</i>	3	1	Negative regulator of Ras
<i>GPB2</i>	3	1	Negative regulator of Ras
<i>KEG1</i>	3	0	Cell wall assembly
<i>KRE5</i>	3	1	Cell wall assembly
<i>RPO31</i>	3	0	RNA polymerase III transcription
<i>SET4</i>	3	2	Unknown
<i>YJL171C</i>	3	0	Unknown

\* Biological process was manually curated from the *Saccharomyces* Genome Database (<http://www.yeastgenome.org>).

negative regulators of *Ras*, we observe 14 missense, 8 nonsense, and 10 frameshift mutations, suggesting that selection at these loci is for loss of function (Supplementary Fig. 2). In contrast, all 13 mutations observed in cell-wall assembly genes are missense, suggesting that at these loci selection is for alteration or attenuation, not loss, of function (Supplementary Fig. 2).

This evidence argues that mutations in multi-hit genes provided strong fitness advantages that made them parallel adaptive solutions in multiple replicate populations (related arguments have been made in bacterial<sup>6,7</sup> and viral<sup>29</sup> systems). However, the fate of each mutation also depends on random hitchhiking and interference effects, which increase variation in evolutionary outcomes. Even beneficial driver mutations must often quasi-hitchhike as co-drivers with others in a larger cohort if they are to succeed. For example, in population BYB1-G07, a mutation in *SPC3* began to sweep within the first 300 generations (Fig. 3), before a competing cohort appeared containing mutations in the multi-hit genes *WHI2* and *ROT2*. The *WHI2*–*ROT2* cohort rose in frequency at the expense of *SPC3*, until the *SPC3* genotype was partially rescued by a mutation in the multi-hit gene *YURI*. Finally, a second and distinct mutation in *WHI2* appeared in the *SPC3*–*YURI* background. This genotype fixed, forcing the *WHI2*–*ROT2* cohort to extinction. These dynamics illustrate how a balance between the fitness advantages of individual driver mutations and random hitchhiking and interference effects determines evolutionary outcomes.

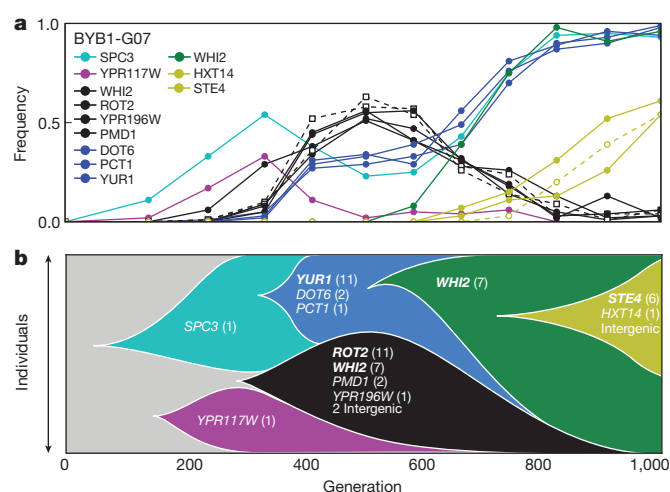
Although the dynamics of any individual population are highly stochastic, a statistical analysis across replicate populations sheds light on the factors that determine the fate of each mutation. To this end, we measured the initial rate of increase in frequency of each mutation (Methods). We have previously<sup>21</sup> referred to this ‘initial relative fitness’ as  $s_{up}$ . It measures the combined fitness effect of a mutation together with the genetic background in which it arose, relative to the average of all other genetic backgrounds currently in the population. The probability that a mutation fixes increases with  $s_{up}$  (Fig. 2c). Non-synonymous mutations tended to have higher  $s_{up}$  than synonymous mutations ( $P < 0.05$ ) and non-synonymous mutations in multi-hit genes tended to have higher  $s_{up}$  than those in single-hit genes ( $P < 0.02$ ), as we would expect if the former classes tend to confer a larger fitness advantage. However, given a particular value of  $s_{up}$ , all types of mutations were equally likely to succeed. In other words, a

weak or neutral mutation on a good background is just as likely to fix as a strongly beneficial mutation on a poor background; all that matters is the initial relative fitness of the mutation combined with the background in which it occurred.

In theory, population size could be predicted to either increase or decrease the patterns of reproducibility between replicate populations. Larger populations will sample more possible mutations, and thus favour the best genotypes in replicate populations<sup>13,16</sup>. But larger populations also maintain more genetic variation, making each mutation more likely to be influenced by chance associations<sup>16</sup>. Our data make it possible to determine experimentally the influence of population size on the reproducibility of evolutionary outcomes. Of the 40 sequenced populations, 14 were evolved at a large ( $10^6$ ) and 26 at a small ( $10^5$ ) population size. We find that putative driver mutations are more commonly observed in large populations (Table 2,  $P < 0.025$ ). However, given a particular value of  $s_{up}$ , a mutation is less likely to fix in a large population—that is, subsequent chance associations are more likely to interfere (Fig. 2d,  $P < 10^{-5}$ ). Together, these results show that beneficial mutations occur more consistently in larger populations, but that each mutation has a more random fate once it has occurred.

To demonstrate how our system can be used to dissect the fitness effects of the individual mutations that underlie these dynamics, we chose for genetic dissection a population that displayed simple sequence-level dynamics. In BYB1-A08, two mutations (*ELO1* and *GAS1*) have fixed and a third (*STE12*) is on its way to fixation by generation 545 (Fig. 4a). Clones from this time point had gained, on average, a 4.3% fitness advantage relative to the ancestor. To determine how these three mutations contribute to fitness, we crossed three clones from generation 545 to the ancestor and isolated 80 haploid progeny. Each haploid was genotyped at these three loci and assayed for fitness, enabling us to quantify the fitness effect of each mutation individually and in combination. We find that mutations in both *GAS1* and *STE12* provide a selective advantage, while the mutation in *ELO1* is a neutral hitchhiker (Fig. 4b). Consistent with this, mutations in *GAS1* and *STE12* are observed in three and nine replicate populations, respectively (Table 1), but *ELO1* only once.

Our analysis has shown that the combination of experimental evolution and whole-genome whole-population sequencing over a dense timecourse is a powerful tool. Our data demonstrate the importance of pervasive hitchhiking and clonal interference among cohorts of mutations in determining the molecular dynamics of adaptation. Further work is needed to determine the mechanism underlying the formation of these cohorts. Interestingly, cohorts and genetic hitchhiking have been described in other systems, such as influenza evolution<sup>2</sup> and the somatic evolution of cancers<sup>30</sup>, suggesting that these dynamics represent a general mode of adaptation. Our data also highlight the relatively small subset of genes that repeatedly provide driver mutations, suggesting a limited number of open pathways to substantially increased fitness. This work is a first step towards a complete understanding of the dynamics of adaptation under conditions where multiple beneficial mutations spread simultaneously, and illustrates the importance of both chance and selection in determining evolutionary outcomes.

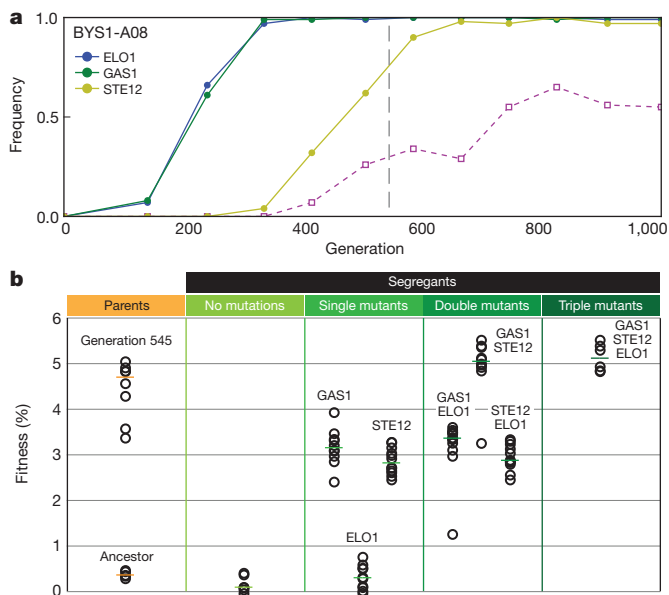


**Figure 3 | The dynamics of sequence evolution in BYB1-G07.** **a**, The trajectories of the 15 mutations that attain a frequency of at least 30%, hierarchically clustered into several distinct mutation ‘cohorts’, each of which is represented by a different colour (Methods). **b**, Muller diagram showing the dynamics of the six main cohorts in the population. The number of times a mutation was observed in a given gene across all 40 populations is indicated in parentheses. Mutations in genes observed in more than three replicate populations (Table 1) are indicated in bold.

**Table 2 | Summary of the fates of nuclear mutations observed throughout the experiment**

Class of mutation	All populations (40)			Small populations (26)			Large populations (14)		
	Total	No. fixed	Fixed (%)	Total	No. fixed	Fixed (%)	Total	No. fixed	Fixed (%)
All	995	246	25	703	191	27	292	55	19
Intergenic	272	58	21	207	48	23	65	10	15
Synonymous	107	26	24	77	21	27	30	5	17
Nonsynonymous	616	162	26	419	122	29	197	40	20
Hit 1 ×	381	86	23	273	65	24	108	21	19
Hit 2 ×	95	24	25	63	17	27	32	7	22
Hit ≥3 ×	140	52	37	83	40	48	57	12	21





**Figure 4 | Genetic dissection of BYSI-A08.** **a**, The trajectories of observed mutations. **b**, We crossed evolved clones from generation 545 (dotted grey line in **a**) to the ancestor; shown here are the fitnesses and genotypes of parental clones and 80 haploid progeny.

## METHODS SUMMARY

Whole-population DNA samples were sequenced using the Illumina HiSeq platform. Mutations were identified using Breqseq (<http://www.barricklab.org/twiki/bin/view/Lab/ToolsBacterialGenomeResequencing>; Pipeline 1), or BWA and FreeBayes (<http://www.bio-bwa.sourceforge.net> and Marth Laboratory, Boston College; Pipeline 2). Time-course information was used to distinguish real mutations from spurious calls due to sequencing or alignment error. Mutations were annotated to the yeast genome using NCBI-BLAST. Fitness of evolved or reconstructed clones was determined by competition against a fluorescently-labelled reference strain as described previously<sup>21</sup>.

**Full Methods** and any associated references are available in the online version of the paper.

Received 11 January; accepted 3 June 2013.

Published online 21 July 2013.

1. Weinreich, D. M., Delaney, N. F., DePristo, M. A. & Hartl, D. L. Darwinian evolution can follow only very few mutational paths to fitter proteins. *Science* **312**, 111–114 (2006).
2. Strelkowa, N. & Lässig, M. Clonal interference in the evolution of influenza. *Genetics* **192**, 671–682 (2012).
3. Levin, B. R. & Bull, J. J. Short-sighted evolution and the virulence of pathogenic microorganisms. *Trends Microbiol.* **2**, 76–81 (1994).
4. Greaves, M. & Maley, C. C. Clonal evolution in cancer. *Nature* **481**, 306–313 (2012).
5. Sprouffske, K., Merlo, L. M. F., Gerrish, P. J., Maley, C. C. & Sniegowski, P. D. Cancer in light of experimental evolution. *Curr. Biol.* **22**, R762–R771 (2012).
6. Tenaillon, O. *et al.* The molecular diversity of adaptive convergence. *Science* **335**, 457–461 (2012).
7. Woods, R., Schneider, D., Winkworth, C. L., Riley, M. A. & Lenski, R. E. Tests of parallel molecular evolution in a long-term experiment with *Escherichia coli*. *Proc. Natl Acad. Sci. USA* **103**, 9107–9112 (2006).
8. Saxer, G., Doebeli, M. & Travisano, M. The repeatability of adaptive radiation during long-term experimental evolution of *Escherichia coli* in a multiple nutrient environment. *PLoS ONE* **5**, e14184 (2010).
9. Atwood, K. C., Schneider, L. K. & Ryan, F. J. Periodic selection in *Escherichia coli*. *Proc. Natl Acad. Sci. USA* **37**, 146–155 (1951).

1. Paquin, C. & Adams, J. Frequency of fixation of adaptive mutations is higher in evolving diploid than haploid yeast populations. *Nature* **302**, 495–500 (1983).
11. Joseph, S. B. & Hall, D. W. Spontaneous mutations in diploid *Saccharomyces cerevisiae*: more beneficial than expected. *Genetics* **168**, 1817–1825 (2004).
12. Perfeito, L., Fernandes, L., Mota, C. & Gordo, I. Adaptive mutations in bacteria: high rate and small effects. *Science* **317**, 813–815 (2007).
13. Gerrish, P. J. & Lenski, R. The fate of competing beneficial mutations in an asexual population. *Genetica* **102–103**, 127–144 (1998).
14. Desai, M. M. & Fisher, D. S. Beneficial mutation-selection balance and the effect of linkage on positive selection. *Genetics* **176**, 1759–1798 (2007).
15. Rouzine, I. M., Wakeley, J. & Coffin, J. The solitary wave of asexual evolution. *Proc. Natl Acad. Sci. USA* **100**, 587–592 (2003).
16. Good, B. H., Rouzine, I. M., Balick, D. J., Hallatschek, O. & Desai, M. M. The rate of adaptation and the distribution of fixed beneficial mutations in asexual populations. *Proc. Natl Acad. Sci. USA* **109**, 4950–4955 (2012).
17. Schifffels, S., Szöllösi, G. J., Mustonen, V. & Lässig, M. Emergent neutrality in adaptive asexual evolution. *Genetics* **189**, 1361–1375 (2011).
18. Desai, M. M., Fisher, D. S. & Murray, A. W. The speed of evolution and maintenance of variation in asexual populations. *Curr. Biol.* **17**, 385–394 (2007).
19. de Visser, J. A. G. M., Zeyl, C. W., Gerrish, P. J., Blanchard, J. L. & Lenski, R. E. Diminishing returns from mutation supply rate in asexual populations. *Science* **283**, 404–406 (1999).
20. Kao, K. C. & Sherlock, G. Molecular Characterization of clonal interference during adaptive evolution in asexual populations of *Saccharomyces cerevisiae*. *Nature Genet.* **40**, 1499–1504 (2008).
21. Lang, G. I., Botstein, D. & Desai, M. M. Genetic variation and the fate of beneficial mutations in asexual populations. *Genetics* **188**, 647–661 (2011).
22. Barrick, J. E. *et al.* Genome evolution and adaptation in a long-term experiment with *Escherichia coli*. *Nature* **461**, 1243–1247 (2009).
23. Barrick, J. E. & Lenski, R. E. Genome-wide mutational diversity in an evolving population of *Escherichia coli*. *Cold Spring Harb. Symp. Quant. Biol.* **74**, 119–129 (2009).
24. Dettman, J. R. *et al.* Evolutionary insight from whole-genome sequencing of experimentally evolved microbes. *Mol. Ecol.* **21**, 2058–2077 (2012).
25. Gresham, D. *et al.* The repertoire and dynamics of evolutionary adaptations to controlled nutrient-limited environments in yeast. *PLoS Genet.* **4**, e1000303 (2008).
26. Bollback, J. P. & Huelsenbeck, J. P. Clonal interference is alleviated by high mutation rates in large populations. *Mol. Biol. Evol.* **24**, 1397–1406 (2007).
27. Betancourt, A. J. Genomewide patterns of substitution in adaptively evolving populations of the RNA bacteriophage MS2. *Genetics* **181**, 1535–1544 (2009).
28. Miller, C. R., Joyce, P. & Wichman, H. A. Mutational effects and population dynamics during viral adaptation challenge current models. *Genetics* **187**, 185–202 (2011).
29. Wichman, H. A., Badgett, M. R., Scott, L. A., Boulianne, C. M. & Bull, J. J. Different trajectories of parallel evolution during viral adaptation. *Science* **285**, 422–424 (1999).
30. Nik-Zainal, S. *et al.* The life history of 21 breast cancers. *Cell* **149**, 994–1007 (2012).

**Supplementary Information** is available in the online version of the paper.

**Acknowledgements** We thank the production team led by L. Fulton and R. Fulton at the Genome Institute at Washington University for sample management and data production, and E. Lobos for coordinating the project. We thank L. Parsons and J. Wiggins for assistance with data management, P. Gibney for assistance with sample preparation, and T. DeCoste for assistance with flow cytometry. We thank K. Kosheleva for discussions, and A. Murray, C. Marx, M. McDonald, G. Sherlock and D. Kvitsek for comments on the manuscript. D.P.R. acknowledges support from an NSF Graduate Research Fellowship. D.B. acknowledges support from NIGMS Centers of Excellence grant GM071508 and NIH grant GM046406. M.M.D. acknowledges support from the James S. McDonnell Foundation, the Alfred P. Sloan Foundation, and the Harvard Milton Fund.

**Author Contributions** G.I.L., D.B. and M.M.D. designed the project; E.S. and G.M.W. generated the sequencing data; G.I.L., D.P.R., M.J.H. and M.M.D. analysed the sequencing data; G.I.L. performed the experiments; G.I.L., D.P.R., D.B. and M.M.D. wrote the paper. Co-senior authors, D.B. and M.M.D.

**Author Information** Genome sequence data have been deposited to GenBank under the BioProject identifier PRJNA205542. Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to G.I.L. (glang@lehigh.edu) or M.M.D. (mmdesai@fas.harvard.edu).

## METHODS

**DNA sequencing.** Cells were grown by inoculating 8 µl of each frozen population from our earlier experiment<sup>21</sup> into 20 ml YPD (yeast extract, peptone, dextrose) + ampicillin (100 µg ml<sup>-1</sup>) and tetracycline (25 µg ml<sup>-1</sup>) and grown overnight to saturation. Cells were pelleted and washed once with water. Genomic DNA was prepared using a modified glass bead lysis method. Cells were resuspended in 400 µl of DNA extraction buffer (2% Triton X-100, 1% SDS, 100 mM NaCl, 10 mM Tris, pH 8.0, and 1 mM EDTA). To the resuspended cells, 600 µl of acid washed glass beads (425–600 µm, acid-washed; Sigma) and 400 µl of phenol:chloroform:isoamyl alcohol (25:24:1, Tris saturated) was added and the cells were mechanically lysed for 2.5 min using a bead beater. After centrifugation, the supernatant was removed and incubated at 37 °C with RNaseA for 1 h, followed by a second phenol:chloroform:isoamyl alcohol extraction. The aqueous supernatant was removed and genomic DNA was precipitated with ethanol and resuspended in water. Paired-end Illumina sequencing libraries of 500-bp fragments were prepared at The Genome Institute, Washington University School of Medicine, and the libraries were run on the Illumina HiSeq with average of 100-fold coverage.

**Identifying mutations from raw sequencing data.** We developed two independent methods for identifying mutations from the raw sequencing data and for distinguishing bona fide mutations from spurious calls that resulted from either sequencing or alignment errors by leveraging time course information. Both pipelines identified base-pair substitutions (BPS), small insertion and deletion mutations (InDel) and complex mutations involving both BPS and InDels. We note, however, that neither pipeline is well suited to identify certain types of mutations, such as copy number variation, inversions, or large insertions or deletions. Both pipelines produced similar results. The data presented in the paper were produced using Pipeline 1. Supplementary Table 1 reports the results of both pipelines.

In Pipeline 1, we used the software package Breseq (<http://www.barricklab.org/twiki/bin/view/Lab/ToolsBacterialGenomeResequencing>) to align Illumina reads and make initial polymorphism calls. We ran Breseq on each time point of each population independently and constructed a list of all mutations called in any time point. For each mutation, we used SAMTOOLS<sup>31</sup> to calculate the frequency of reads supporting the mutation in all time points of the population where the mutation was called. We then applied a series of filters based on the frequency trajectories to eliminate false positives. Mutations that did not change frequency over the course of the entire experiment are likely to be sequencing or alignment errors. Therefore, we required the maximum frequency to be at least 0.1 greater than the minimum frequency. We also required the absolute difference between the maximum or minimum frequency and the frequency at generation zero to be at least 0.1. The frequency trajectories of real mutations are expected to be autocorrelated, whereas those of false positives should be uncorrelated from time point to time point. We rejected any mutation with an autocorrelation coefficient less than 0.2. Generation zero was not expected to contain any mutations. Therefore, we rejected any mutation detected by Breseq in generation zero of more than five populations. Also, for any mutation detected by Breseq in generation zero of more than two populations, we required the autocorrelation coefficient to be at least 0.5. Finally, for any mutation with a frequency greater than 0.01 in generation zero, we required the autocorrelation coefficient to be at least 0.35.

In Pipeline 2, for each population and for each time point, we aligned the raw reads to a SNP/Indel corrected W303 reference genome (reference available upon request) using BWA for Illumina version 1.2.2 (ref. 32) using default parameters (except 'Disallow insertion/deletion within [value] bp towards the end' set to 0 and 'Gap open penalty' set to 5). Mutations were called relative to the SNP/Indel corrected W303 reference genome using Freebayes version 0.8.9.a. (Marth Laboratory, Boston College) using default parameters (except 'Pooled' set to 'True', and 'Base alignment quality (BAQ) adjustment' set to 'True'). For each population we merged the 12 resulting .vcf files (one for each time point) using the 'vcf-merge' included in the VCFtools package ([http://vcftools.sourceforge.net/perl\\_module.html](http://vcftools.sourceforge.net/perl_module.html)). We wrote two perl scripts to analyse the resulting merged .vcf file (programs available upon request). The script 'allele\_counts.pl' calculated the frequencies of mutant alleles for each time point in the series and 'composite\_scores.pl' scored the trajectories of each mutation across the twelve time points based on six attributes: autocorrelation, area under the curve relative to time zero, minimum frequency, maximum frequency, max step (the largest difference in frequency in adjacent time points), and the number of called alternate alleles. We developed a heuristic composite score with which to rank the trajectories by their likelihood of being a bona fide mutation.

Any mutation called in either Pipeline was validated manually using the Integrative Genome Viewer<sup>33,34</sup>.

**Annotating mutations.** For each mutation, we aligned the surrounding 2-kb region to the annotated s288c genome using NCBI-BLAST<sup>35</sup>. We then used NCBI-BLAST's CDS feature option to identify the gene or intergenic region containing the mutation and the identity of any amino-acid changes.

All of the observed nuclear mutations represent unique alterations to the yeast genome with four exceptions: two cases of recurrent mutation at the same position and two instances of pre-existing mutations in the seed culture that reached detectable frequency during the evolution experiment. In *ROT2* and *STE12*, recurrent frameshift mutations were observed within homopolymeric runs of seven T's and eight G's, respectively. For *ROT2* all four occurrences of mutations in this homopolymeric run were T insertions. For *STE12*, two mutations were G insertions and two were G deletions. In addition to recurrent mutations, we observed two pre-existing mutations. In the initial evolution experiment, two nearly isogenic haploid ancestral strains (B and R) were used to seed approximately 300 populations each. Of the sequenced populations reported here, 30 are derived from the B progenitor and 10 from the R progenitor. We observed several occurrences where the same mutation was observed in multiple populations. The same single base-pair deletion in *IRA1* was observed in four populations derived from the B ancestor. In each case this allele was observed early and before the first selective sweep suggesting that this mutation was present at low frequency in the starting B population. In all 10 R populations, the same T to C substitution in *PDR5* was initially at 15% at Generation 0. This mutation quickly fixed in two populations, slowly fixed in another, rose to above 50% before going extinct in two and quickly went extinct in the other five (Supplementary Fig. 3).

**Analysis of trajectories.** To assess the relationship between fitness and fixation probability, we estimated  $s_{up}$ , the fitness of clones containing a given mutation relative to the mean fitness of the population when we first detected the mutation. For each mutation, we identified  $t_1$  and  $t_2$ , the first consecutive time points such that the frequency of the mutation at  $t_1$  was greater than zero and the frequency at  $t_2$  was greater than 0.1. We then calculated

$$s_{up} = \frac{1}{t_2 - t_1} \left( \ln \frac{f(t_2)}{1 - f(t_2)} - \ln \frac{f(t_1)}{1 - f(t_1)} \right)$$

where  $f(t)$  is the frequency of the mutation at time  $t$ . This quantity estimates the combined effects of the focal mutation and the background it occurred on. For instance, a mutation conferring a 3% fitness advantage on a neutral background will have the same value of  $s_{up}$  as a neutral mutation occurring on a background that is 3% fitter than the population average.

**Identifying mutation cohorts.** The most notable feature of our results is that mutations often move through populations as temporal clusters of functionally unrelated mutations, tracking tightly with one another through time. We have termed these 'cohorts'. To empirically assign mutations to cohorts, we treated each frequency trajectory as a vector in twelve dimensions. We used the hierarchical clustering package in SciPy (<http://www.scipy.org>) to cluster the mutations in each population based on the Euclidean distance between frequency vectors. Because low-frequency mutations contain too little information for reliable clustering, we excluded mutations with maximum frequencies less than 0.3. We then flattened the hierarchies using a cutoff distance of 0.275.

**Fitness assays and genetic dissection.** Fitness assays were performed as described previously<sup>21</sup>. To measure the fitness of evolved clones from frozen stock, we struck to singles from population BYS1-A08 from generation 545. We selected seven single colonies at random and measured their fitness relative to an mCherry-expressing reference strain. The experimental and reference strains were grown separately in 96-well plates, then mixed 50:50 and propagated by diluting 1:1,024 every 24 h. At generations 10, 20, 30 and 40, we transferred 4 µl of saturated culture into 100 µl of cold PBST and the ratio of nonfluorescent (experimental) and mCherry-positive (reference) cells was determined by flow cytometry using an LSRII flow cytometer (BD Biosciences) counting 50,000 total cells for each sample. The fitness difference between the experimental and reference strain was calculated as the rate of the change in the  $\ln$  ratio of experimental to reference versus generations<sup>36</sup>. To determine fitness effects of the three evolved mutations in BYS1-A08 (*GAS1*, *ELO1*, *STE12*), we chose three of the seven clones and backcrossed them to a MAT $\alpha$  version of the ancestral strain. From these three diploids we sporulated and selected 80 haploid MAT $\alpha$  segregants. Each segregant was genotyped by SNP-specific PCR using the following primers: *GAS1*\_Forward (5'-TTTTCGTGCGCAAACGTGG-3'), *GAS1*\_WT\_Reverse (5'-ATTGGAAGAGTAGCCAACCTG-3'), *GAS1*\_Mutant\_Reverse (5'-ATTGGAAGAGTAGCCAACCTA-3'), *ELO1*\_Forward (5'-AACACAACAAATCGCAAGCC-3'), *ELO1*\_WT\_Reverse (5'-TAACCAACCAATGATTATA-3'), *ELO1*\_Mutant\_Reverse (5'-TAACCAACCAATGATTATG-3'), *STE12*\_Reverse (5'-TGAGCAGAACTCTTCGTCACC-3'), *STE12*\_WT\_Forward (5'-AATCTCACAACCTCTGGCCAG-3'), and *STE12*\_Mutant\_Forward (5'-AAATCTCACAACCTCTGCCAA-3'). The fitness of each of the haploid segregants was measured relative to the mCherry-expressing reference strain as described above.

31. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
32. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
33. Robinson, J. T. *et al.* Integrative genomics viewer. *Nature Biotechnol.* **29**, 24–26 (2011).
34. Thorvaldsdóttir, H., Robinson, J. T. & Mesirov, J. P. Integrative Genomics Viewer (IGV): high-performance genomics data visualization and exploration. *Brief. Bioinform.* **14**, 178–192 (2013).
35. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
36. Hartl, D. *A Primer of Population Genetics*. (Sinauer Associates, 2000).