

1 Systematic perturbation of yeast essential genes using base 2 editing

3

4 Philippe C Després^{1,2,3}, Alexandre K Dubé^{1,2,3,4}, Motoaki Seki⁵, Nozomu Yachie^{*5,6,7} and
5 Christian R Landry^{*1,2,3,4}

6

7

8 1. Département de Biochimie, Microbiologie et Bio-informatique, Faculté de sciences et
9 génie, Université Laval, Québec, Québec, G1V 0A6, Canada

10 2. PROTEO, le regroupement québécois de recherche sur la fonction, l'ingénierie et les
11 applications des protéines, Université Laval, Québec, Québec, G1V 0A6, Canada

12 3. Centre de Recherche en Données Massives (CRDM), Université Laval, Québec,
13 Québec, G1V 0A6, Canada

14 4. Département de Biologie, Faculté de sciences et Génie, Université Laval, Québec,
15 Québec, G1V 0A6, Canada

16 5. Research Center for Advanced Science and Technology, Synthetic Biology Division,
17 University of Tokyo, Tokyo, 4-6-1 Komaba, Meguro-ku, 153-8904, Japan

18 6. Department of Biological Sciences, Graduate School of Science, the University of
19 Tokyo, Tokyo, Japan

20 7. Institute for Advanced Biosciences, Keio University, Tsuruoka, Japan

21 *To whom correspondence should be addressed. CRL: Tel: 1-418-656-3954, Fax 1-418-
22 656-7176, christian.landry@bio.ulaval.ca NY: Tel +81-3-5452-5242 (x55242), Fax +81-
23 3-5452-5241 (x55241), yachie@synbiol.rcast.u-tokyo.ac.jp

24 Abstract

25 Base editors derived from CRISPR-Cas9 systems and DNA editing enzymes offer an
26 unprecedented opportunity for the precise modification of genes, but have yet to be used at a
27 genome-scale throughput. Here, we test the ability of an editor based on a cytidine deaminase,
28 the Target-AID base editor, to systematically modify genes genome-wide using the set of yeast
29 essential genes. We tested the effect of mutating around 17,000 individual sites in parallel across
30 more than 1,500 genes in a single experiment. We identified over 1,100 sites at which mutations
31 have a significant impact on fitness. Using previously determined and preferred Target-AID
32 mutational outcomes, we predicted the protein variants caused by each of these gRNAs. We
33 found that gRNAs with significant effects on fitness are enriched in variants predicted to be
34 deleterious by independent methods based on site conservation and predicted protein
35 destabilization. Finally, we identify key features to design effective gRNAs in the context of base
36 editing. Our results show that base editing is a powerful tool to identify key amino acid residues
37 at the scale of proteomes.

38 **Introduction**

39 Recent technical advances have allowed the investigation of the genotype-phenotype map at high
40 resolution by experimentally measuring the effect of all possible nucleotide substitutions in a short
41 DNA sequence. While saturated mutagenesis informs us on the effect of many mutations, it
42 usually covers a single locus or a fraction of it (Fowler and Fields 2014; Gray *et al.* 2018). Because
43 such data is only available at sufficient coverage for a very small number of proteins, general
44 rules on substitution effects must be extrapolated to other, often unrelated proteins. At a lower
45 level of resolution, genome-scale mutations data has mostly been acquired through large-scale
46 loss-of-function strain collections, where the same genetic change (for example, complete gene
47 deletion) is applied to all genes (Winzeler *et al.* 1999; Giaever *et al.* 2002; C. elegans Deletion
48 Mutant Consortium 2012). This approach is a powerful way to isolate each gene's contribution to
49 a phenotype, including fitness, but limits our understanding of the role of specific positions within
50 a locus.

51 CRISPR-Cas9 based approaches usually cause protein loss of function through indel formation
52 (Shalem *et al.* 2014) or by modifying gene expression levels (Qi *et al.* 2013; Sander and Joung
53 2014; Smith *et al.* 2016) at many loci in parallel. Again, these approaches generally limit the
54 information gain to one perturbation per locus. There is therefore a strong tradeoff between the
55 resolution of the existing assays and the number of loci or genes investigated. Recent
56 developments in the field now allow for the exploration of the effects of many mutations per gene
57 across the genome. For instance, in yeast, methods for high throughput strain library construction
58 have allowed the measurement of thousands of variant fitness effects in parallel across the
59 genome (Sharon *et al.* 2018; Bao *et al.* 2018; Roy *et al.* 2018). These approaches rely on
60 CRISPR-Cas9 based genome modifications requiring the formation of double-strand breaks
61 followed by repair using donor DNA, which often depends on complex strain and plasmid
62 constructions. An alternative approach would be to use base editors, which allow the introduction

63 of the mutations of interest directly in the genome by direct modification of DNA bases rather than
64 DNA segment replacement.

65
66 Base editors use DNA modifying enzymes fused to modified Cas9 or Cas12 proteins to create
67 specific point mutations in a target genome (Nishida *et al.* 2016; Gaudelli *et al.* 2017; reviewed in
68 Rees and Liu 2018). Such base editors have recently been used to perform site-specific forward
69 mutagenesis in human cell lines. The two main approaches, Targeted AID-mediated mutagenesis
70 (TAM) (Ma *et al.* 2016) and CRISPR-X (Hess *et al.* 2016), target specific regions of the genome
71 where they induce mutations randomly. This generates a library of mutant genotypes that can be
72 competed to find beneficial and deleterious variants under selective pressure. As the relative
73 fitness measurements depend on targeted sequencing of the locus of interest, these approaches
74 are difficult to adapt to high throughput multiplexed screens where tens of thousands of sites can
75 be targeted within the same gRNA libraries.

76
77 Here, we present a method that bridges the flexibility of Target-AID mutagenesis and the
78 multiplexing capacities of genome editing depletion screens. By using a base editor with a narrow
79 and well-defined activity window (Nishida *et al.* 2016), we selected gRNAs generating a limited
80 number of predictable edits in yeast essential genes. This allowed us to use gRNAs as a readout
81 for the effect of the mutations, similar to commonly used barcode-sequencing approaches to
82 measure fitness effects.

83 **Results**

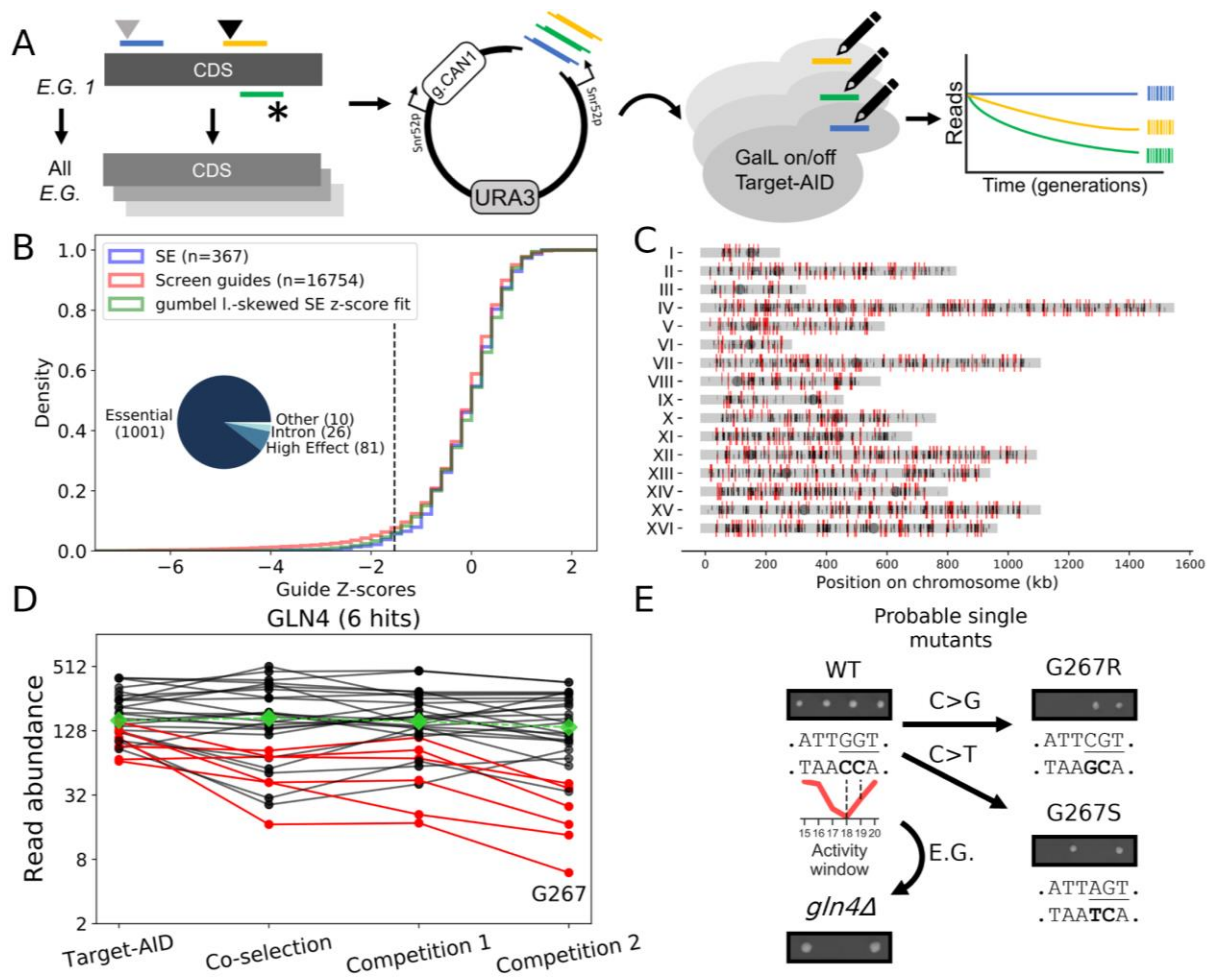
84 **Large-scale base editing screening**

85 We used Target-AID mutagenesis to simultaneously assess mutational effects at over 17,000
86 putative sites in the yeast genome. We scanned yeast essential genes for sites amenable to
87 editing by the Target-AID base editor as well as targets with other specific properties, including

88 intronic sequences (Figure 1A, Figure S1). Because all essential genes have the same fitness
89 effects when deleted (Giaever *et al.* 2002), focusing on these genes allowed to limit the variation
90 in fitness that could be due to the relative importance of individual genes for growth rather than to
91 the importance of specific positions.

92 To ensure we could predict gRNA mutational outcomes with accuracy, we included in the library
93 only gRNAs with one to two nucleotides with a high probability of being edited based on the known
94 activity window of Target-AID in yeast (Nishida *et al.* 2016). We could then predict mutagenesis
95 outcomes for gRNAs computationally. We took into account that Target-AID produces both C-
96 to-G and C-to-T mutations in yeast, with a 1.5 to 2 fold preference for C-to-G (Nishida *et al.* 2016;
97 Després *et al.* 2018). We also extended the analysis to include other point mutants at possible
98 secondary editing sites within the activity window (see methods). As such, we could associate
99 most gRNAs targeting protein-coding DNA to a primary C-to-G and C-to-T outcome (C-to-G #1
100 and C-to-T #1), as well as to possible secondary outcomes if applicable (C-to-G #2 and C-to-T
101 #2). We did not consider gRNAs that did not target between the 0.5th and 75th percentile of the
102 length of annotated genes to limit position biases that could influence the efficiency of stop-codon
103 generating guides (Doench *et al.* 2014; Michel *et al.* 2017).

104 The gRNA library was cloned into a high-throughput co-selection base editing vector (Després *et*
105 *al.* 2018). We performed pooled mutagenesis followed by bulk competition (Figure S2) to identify
106 mutations with significant fitness effects. As the relative abundance of each gRNA in the extracted
107 plasmid pool depends on the abundance of the subpopulation of cells bearing these gRNAs, any
108 fitness effect caused by the mutation they induce will influence their relative abundance. Variation
109 in plasmid abundance was measured using targeted next-generation sequencing of the variable
110 gRNA locus on the base editing vector in a manner similar to GeCKO approaches (Sanjana *et al.*
111 2014; Shalem *et al.* 2014).



112

113 **Figure 1. High-throughput forward mutagenesis by Target-AID base editing identifies sensitive sites**
 114 **across the yeast genome. A)** Experimental design. Essential genes were scanned for sites appropriate
 115 for Target-AID mutagenesis. Mutational outcomes include silent (grey triangle) and missense (black
 116 triangle) mutations, as well as stop codons (*). DNA fragments bearing the gRNA sequences were
 117 synthesized as an oligonucleotide pool and cloned into a co-selection base editing vector. Using gRNAs as
 118 molecular barcodes, the abundance of cell subpopulations bearing mutations was measured during
 119 mutagenesis and bulk competition. Mutations with fitness effects were inferred from a reduction in the
 120 relative barcode read count. **B)** Cumulative distribution of z-scores of the log₂ fold-change in gRNA
 121 abundance between mutagenesis and the end of the bulk competition experiment averaged between
 122 replicates (see Figure S2). A 5% false positive threshold was calculated by fitting a distribution of
 123 abundance variation z-score of the sequenced gRNAs with synthesis errors (SE gRNAs) and is represented
 124 by a dotted black line. The distribution of target types in the 1,118 gRNAs with Negative Effects (GNE) is
 125 shown in the inset. **C)** Positions of base editing target sites in the yeast genome. Telomeric regions are
 126 depleted in target sites because very few essential genes are located there. GNEs are shown in red,
 127 and other gRNAs are in black. The orientation of the line matches the targeted strand relative to the annotated
 128 coding sequence. **D)** Decline in barcode abundance (on a log scale) between timepoints after mutagenesis
 129 for gRNAs targeting *GLN4*, a tRNA synthetase. Median barcode abundance across the entire library
 130 through time is shown in green. The red lines represent the gRNAs categorized as having a significant
 131 effect (GNE) for this gene, while non-significant gRNAs (NSG) are shown in black. The gRNA with the most
 132 extreme z-score targets residue G267. **E)** Mutagenesis of *GLN4*-G267 confirms its essential role for protein

133 function (See methods and Figure S3A). Tetrad dissection of a heterozygous deletion mutant bearing an
134 empty vector results in only two viable spores, while the wild-type copy in the same vector restores growth.
135 Dissection of the two heterozygous mutants bearing a plasmid with the most probable single mutant based
136 on the known activity window of Target-AID shows both mutations are lethal.

137

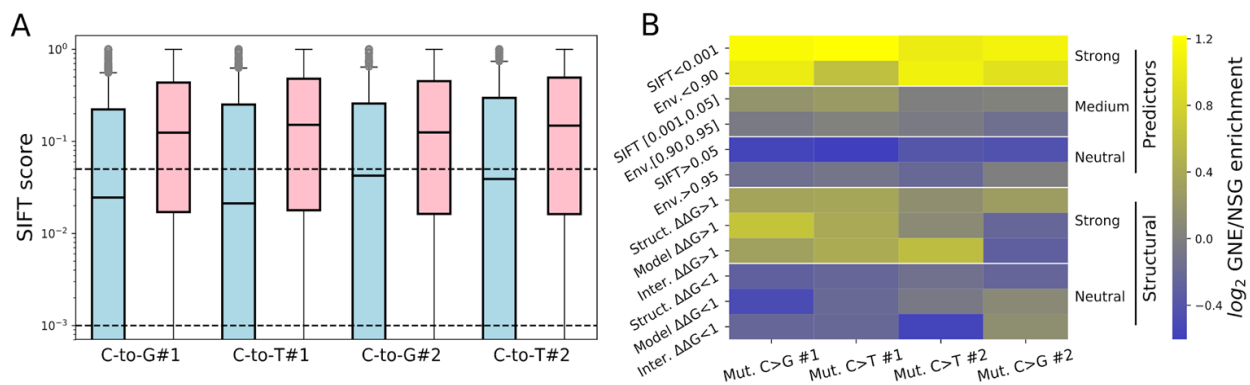
138 After applying a stringent filtering threshold based on barcode read count at the mutagenesis step
139 (Figure S2), we identified a total of ~17,000 gRNAs for which we could evaluate fitness effects.
140 Replicate data for gRNAs passing the minimal read count selection criteria show high correlation
141 across experimental time points (Figure S3) and cluster by experimental step (Figure S4),
142 showing that the approach is reproducible. Using the distribution of abundance variation of non-
143 functional gRNAs with synthesis errors as a null distribution (see methods), we identified 1,118
144 gRNAs across 605 genes or loci with significant negative effects (GNE) on cell survival or
145 proliferation using at an estimated 5% false positive rate. GNEs are distributed evenly across the
146 yeast genome (Figure 1B and 1C), suggesting no inherent bias against specific regions.

147 An example of barcode abundance variation through time for all gRNAs (both GNEs and NSGs)
148 targeting *GLN4* is shown in Figure 1D. *GLN4* is an essential gene coding for a glutamine t-RNA
149 synthetase. To confirm the deleteriousness of the predicted mutations, we transformed a
150 centromeric plasmid bearing a wild-type or mutated copy of the gene under the control of its native
151 promoter (Ho *et al.* 2009) in a heterozygous deletion background (Giaever *et al.* 1999). Following
152 dissection, spore survival was compared between wild-type and mutated copy of *GLN4* (Figure
153 S5). Using this approach, we confirmed the strong fitness effect of the best scoring GNE for *GLN4*,
154 as the most probable mutations generated are in fact lethal (Figure 1D).

155 **Comparison of GNE induced mutations with variant effect predictions**

156 If GNEs indeed induce specific deleterious mutations, these mutations should be predicted to be
157 more deleterious than those of Non-Significant gRNAs (NSG). We tested two recently published
158 resources for variant effect prediction: Envision (Gray *et al.* 2018) and Mutfunc (Wagih *et al.*
159 2018). Envision is based on a machine learning approach that leverages large-scale saturated

160 mutagenesis data of multiple proteins to perform quantitative predictions of missense mutation
 161 effects on protein function. The lower the Envision score, the higher the effect on protein function.
 162 Mutfunc aggregates multiple types of information such as residue conservation through the use
 163 of SIFT (Ng and Henikoff 2003) as well as structural constraints to provide a binary prediction of
 164 variant effect based on multiple quantitative and qualitative values. Mutations with a low SIFT
 165 score have a lower chance of being tolerated, while those with a positive $\Delta\Delta G$ are predicted to
 166 destabilize protein structure or interactions. Both Envision and the Mutfunc aggregated SIFT data
 167 cover the majority of the most probable mutations generated by the gRNA library (Figure S6A).
 168 The structural modeling information had much lower coverage, covering at best around 12% of
 169 the most probable mutations (Figure S6B).



170
 171 **Figure 2: GNE induced mutations are enriched in predicted deleterious effects** **A)** SIFT score
 172 distributions for the most likely induced mutations of both GNEs (blue) and NSGs (red). The thresholds for
 173 the categories used in the enrichment calculations in **B)** are shown as black dotted lines. SIFT scores
 174 represent the probability of a specific mutation being tolerated based on evolutionary information: the first
 175 threshold of 0.05 was set by the authors in the original manuscript (Ng and Henikoff 2003) but might be
 176 permissive considering the number of mutations tested in our experiment. All GNE vs NSG score
 177 comparisons are significant (Welch's t-test p-values: 1.19×10^{-24} , 3.01×10^{-24} , 9.00×10^{-12} , 1.55×10^{-12}). The
 178 box cutoff is due to the large fraction of mutations for which the SIFT score is 0. **B)** Enrichment folds of
 179 GNEs over NSGs for different variant effect prediction measurements. Envision score (Env.), SIFT score
 180 (SIFT), protein folding stability based on solved protein structures (Struct. $\Delta\Delta G$), protein folding based on
 181 homology models (Model $\Delta\Delta G$) and protein-protein interaction interface stability based on structure data
 182 (Inter. $\Delta\Delta G$). The raw values used to calculate ratios are shown in Supplementary table 1. The prediction
 183 based on conservation and experimental data are grouped under 'Predictors' and those based on the
 184 computational analysis of protein structures and complexes under 'Structural'.

185

186 As expected, mutations generated by GNEs showed significantly lower SIFT scores (Figure 2A)
187 and showed enrichment for strong effects predicted by SIFT, and Envision. Indeed, all four most
188 probable substitutions created by GNEs are about twice more likely to be predicted to have a
189 large deleterious effect by Envision or a very low chance of being tolerated as predicted by SIFT
190 compared to NSG gRNAs. The high homogeneity of Envision scores across the proteome makes
191 it harder to interpret. As such, the shift in score values is more subtle but supports that GNE
192 mutations are generally more likely to be deleterious as well (Figure S6C, Figure S7A).

193 Mutation with destabilizing effects as predicted by structural data also appeared to be enriched
194 for the most probable mutations but low residue coverage limits the strength of this association.
195 This is supported by the raw $\Delta\Delta G$ value distributions, which show a significant tendency (Welch's
196 t-test p-values: 0.0001, 0.0064, 0.148, 0.007) for GNE mutations to be more destabilizing (Figure
197 S7B,C,D). However, the shift in distribution only achieved significance for certain mutation
198 predictions based on solved structures and homology models. While low residue coverage limits
199 our statistical power, this weak apparent enrichment for mutations affecting protein stability may
200 reflect the marginal stability of the target proteins (DePristo *et al.* 2005), resulting in individual
201 destabilizing mutations having a limited effects on fitness. As expected from known experimental
202 data on mutagenesis outcomes (Nishida *et al.* 2016), signal was usually stronger for the most
203 probable C to G mutation.

204 **Sensitive sites provide new biological insights**

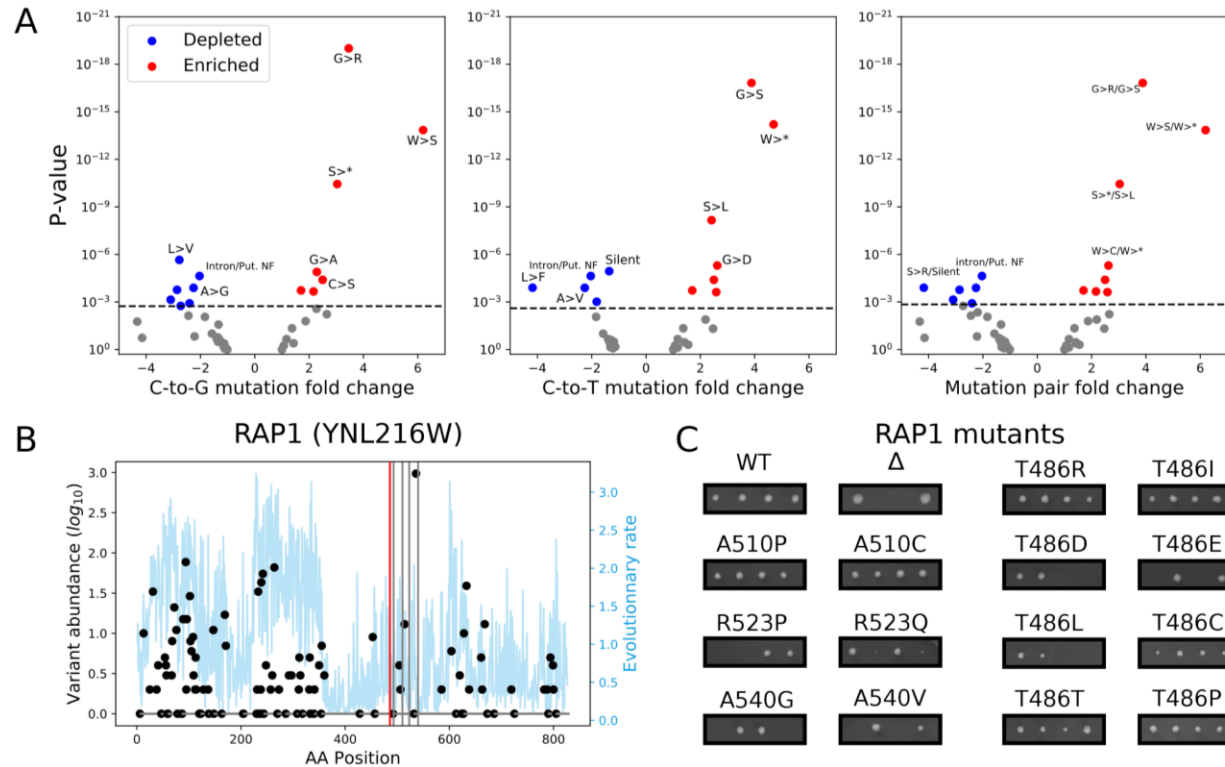
205 Because our screen specifically targeted essential genes, many gRNAs cause mutations in highly
206 conserved regions with high functional importance. To illustrate this, we focus on the highest
207 scoring GNE targeting *GLN4*, a tRNA synthetase, shown in Figure 1D. The gRNA 33725 mutates
208 a glycine at position 267 into either an arginine or a serine. Glycine 267 is part of the "HIGH" motif,
209 characteristic of class I tRNA synthetases, and is involved in ATP binding and catalysis and is
210 highly conserved through evolution (Eriani *et al.* 1990). As expected, the region around the "HIGH"

211 motif shows both a low evolutionary rate based on inter-species comparisons and a much lower
212 variant density in yeast populations compared to other domains of Gln4 (Figure S3B), showing
213 conservation both on a short and long timescales. Surprisingly, mutagenesis experiments in the
214 bacterial homolog MetRS concluded that mutating this residue from glycine to alanine did not alter
215 significantly catalysis while mutating it to proline had a strong disruptive effect (Schmitt *et al.*
216 1995). We found that mutating Gly 267 either to Arg and Ser was enough to cause protein loss of
217 function (Figure 1D). Other sensitive sites identified in *GLN4* by our screen are also clustered in
218 regions with slow evolutionary rates. Interestingly, one of these mutations affects residue R568,
219 which has been hypothesized to play a conserved role from bacteria to yeast in the anti-codon
220 and glutamine recognition process (Grant *et al.* 2013).

221 Since Target-AID can only generate a limited range of amino acid substitutions from a specific
222 coding sequence, we investigated whether any of these mutational patterns were enriched in
223 GNEs (Figure 3A, source data in Supplementary tables S2, S3, and S4). We found several
224 deviations from random expectations in both C-to-G and C-to-T mutation ratios as well as in
225 mutation combination ratios. Three out of four of the mutation pair patterns involving glycine were
226 enriched in GNEs. For example, the Glycine to Arginine or Serine substitutions (as exemplified
227 by guide 33725 targeting *GLN4*) is the second most enriched pattern, being almost four-fold
228 overrepresented in GNE outcomes. This pattern is consistent with the fact that Arginine has
229 properties highly dissimilar to those of Glycine (Sneath 1966), making these substitutions highly
230 deleterious. Furthermore, as Glycine residues are often important components of cofactor binding
231 motifs (eg.: Phosphates) (Copley and Barton 1994) this observation might reflect a tendency for
232 GNEs to alter these sites. Interestingly, genes for which more than one GNE were detected were
233 enriched for molecular function terms linked to cofactor binding (Supplementary table 5). This
234 suggests that the GNEs might indeed have a tendency to affect protein function through
235 mechanisms other than protein or interaction interface destabilization. These protein properties

236 depend on many residues, making them more robust to single amino acid substitutions, whereas
237 cofactor binding may depend specifically on a handful of residues, making these sites critical for
238 function.

239 As expected, there is a strong enrichment for patterns that result in mutation to stop codons: both
240 C-to-G patterns (Tyrosine to stop and Serine to stop) but only one C-to-T pattern (Tryptophan to
241 stop) was overrepresented significantly. Substitutions to stop codon in one outcome also drove
242 enrichment in the other: for example, the link between Serine to Stop (C-to-G) appears to be the
243 cause of the Serine to Leucine (C-to-T) overrepresentation. Both mutation pairs involving mutating
244 a Tryptophan to a stop via a C-to-T mutation: this is not surprising, as the alternative mutations
245 Tryptophan to Serine or Cysteine are also highly disruptive (Sneath 1966). Changes between
246 similar amino acids, which are expected to be tolerable, were also generally depleted in GNE (ex.:
247 the Alanine to Glycine/Valine pair). Mutations in intronic sequences and putative non-functional
248 peptides were also underrepresented, as were most patterns leading to silent mutations. These
249 results show the power of this approach to discriminate important functional sites from mre
250 mutation tolerant ones across the genome.



251

252 **Figure 3 GNE mutations are enriched for specific amino acid substitution patterns and identify**
 253 **critical sites for protein function. A)** Fold depletion and enrichment volcano plots for the most probable
 254 mutations induced by GNEs in the screen. Enrichment and depletion values were calculated by comparing
 255 the relative abundance of each mutation among GNEs and NSGs using Fisher's exact tests. Mutation
 256 patterns significantly depleted are shown in blue, while those that are enriched are in red. The significance
 257 threshold was set using the Holm-Bonferroni method at 5% FDR and is shown as a dotted grey line. **B)**
 258 Protein variant frequency among 1000 yeast isolates (black dots) and residue evolutionary rate across
 259 species (blue line) for *RAP1*. The target site for the GNEs targeting T486 is highlighted by a red line while
 260 the other detected GNEs target sites are shown by a grey line. **C)** Tetrad dissections confirm most *RAP1*
 261 GNE induced mutations indeed have strong fitness effects, as well as other substitutions targeting these
 262 sites.

263

264 The precise targeting of our method also allows us to investigate amino acid residues with known
 265 functional annotations such as post-translational modifications. We found no significant
 266 enrichment for gRNAs mutating directly annotated PTMs ($\text{ratio}^{\text{GNE PTM}} = 19/1118$, $\text{ratio}^{\text{NSG PTM}}$
 267 $243/15536$, Fisher's exact test $p=0.71$). This is consistent with the hypothesis that many PTM
 268 sites may have little functional importance (Landry *et al.* 2009) and thus their mutations may have
 269 no detectable effects for a large part. The same was also observed for gRNAs mutating residues
 270 near known PTMs that could disturb recognition sites ($\text{ratio}^{\text{GNE nearPTM}} = 130/1118$, $\text{ratio}^{\text{NSG nearPTM}}$

271 = 1698/15536, Fisher's exact test $p=0.43$). However, GNEs that do target annotated PTM sites
272 might provide additional evidence supporting the importance of these sites in particular. For
273 example, the best scoring GNE in the well-studied transcriptional regulator *RAP1* is predicted to
274 mutate residue T486. This threonine has been reported as phosphorylated in two previous studies
275 (Albuquerque *et al.* 2008; Holt *et al.* 2009), but the functional importance of this phosphorylation
276 has not been explored yet. Residue T486 is located in a disordered region in the DNA binding
277 domains (Konig *et al.* 1996), which part of the only *RAP1* fragment essential for cell growth
278 (Graham *et al.* 1999; Wu *et al.* 2018).

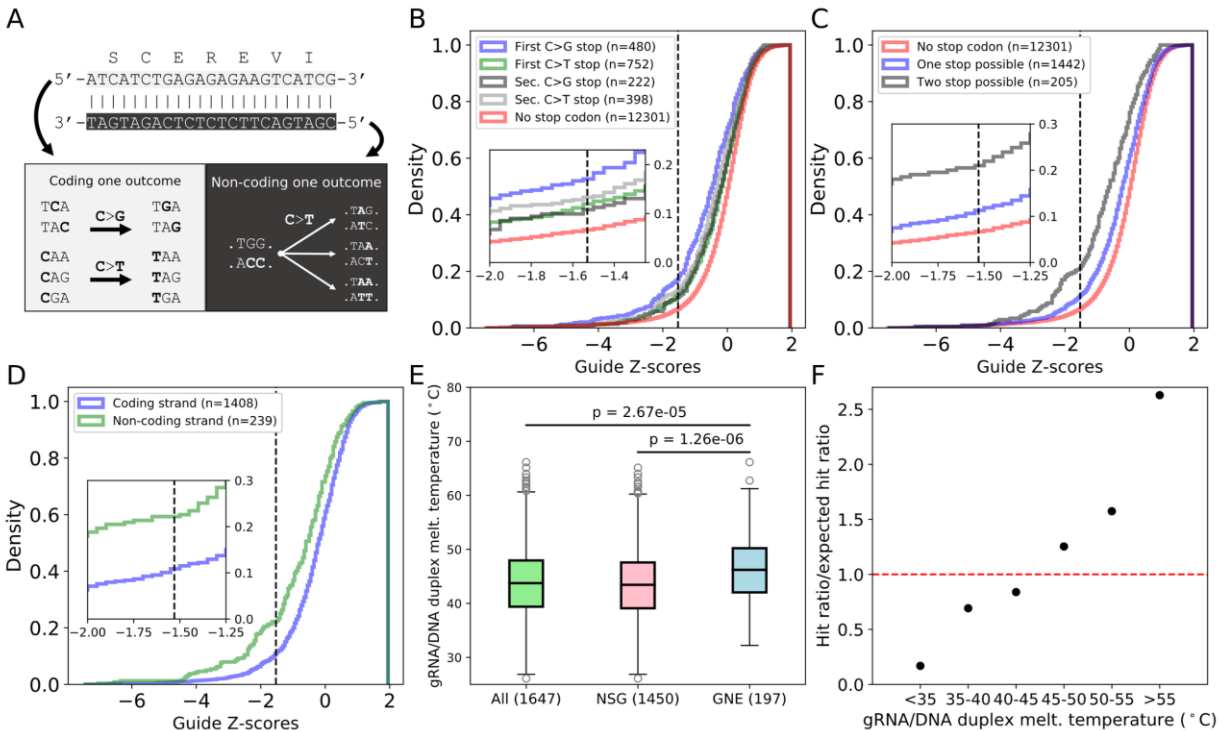
279 Because the available wild-type *RAP1* plasmid (see methods) does not complement gene
280 deletion growth phenotype, we used a different strategy for confirmation that relied on CRISPR-
281 mediated knock-in (see methods and Figure S8). While we could not confirm that the two most
282 likely mutations predicted to be caused by the GNE had a detectable fitness effect in these
283 conditions, we found that phosphomimetic mutations at this position were lethal (Figure 3C and
284 D) but most other amino acids were well tolerated. This suggests that the constitutive
285 phosphorylation of this residue would be highly deleterious. We could also confirm deleterious
286 effects for GNE induced mutations targeting residues R523 and A540, while mutations at residue
287 A510 had no detectable effect on fitness (Figure 3C and D). As we only tested progeny survival
288 on rich media and at a permissive temperature and the screen was performed in synthetic media
289 at 30°C, these mutants might still affect cell phenotype but in an environment-dependent manner.

290 **gRNA properties influence mutagenesis efficiency**

291 There are still very few high-throughput experimental datasets available that allow the
292 investigation of which gRNA properties affect editing efficiency in the context of base editing. Such
293 large-scale data was key in developing models to optimize Cas9 nuclease activity in other types
294 of genome editing experiments, which revealed that sequence specific motifs and thermodynamic
295 RNA properties can be key features (Doench *et al.* 2014, 2016; Wong *et al.* 2015). As gRNAs

296 showing high Cas9 nuclease activity might have poor base editing activity (Kim *et al.* 2017),
297 existing datasets are not easily transferable. We therefore examine what gRNA and target
298 sequence features could influence mutagenesis efficiency. To do so, we focused on the subset
299 of gRNAs with the potential to generate stop codons (stop codon generating gRNAs, SGG) in
300 essential genes (Figure 4A). Successful mutagenesis by SGGs should result in cell death or no
301 proliferation, and a sharp decrease in read abundance, also serving as a positive control for
302 fitness effects within the screen. As most gRNAs were designed to target the first 75% of the
303 coding sequences of essential genes, it is expected that stop codons in these genes would lead
304 to a loss of function.

305



306

307 **Figure 4: gRNA and target properties affect mutagenesis efficiency.** **A)** Since Target-AID can generate
 308 both C to G and C to T mutations, many codons can be targeted to create premature stop codons. **B)**
 309 Cumulative z-score density of SGGs grouped by the mutational outcome generating the stop codon. A
 310 higher rate of GNE is observed for gRNAs for which a C-to-G mutation at the highest editing activity position
 311 generates a stop codon mutation. The significance threshold is shown as a black dotted line **C)** Cumulative
 312 z-score density of SGG with a different number of mutational outcomes that could result in a stop codon.
 313 gRNAs for which more than one mutational outcome results in a stop codon show a higher mutagenesis
 314 success rate. **D)** Cumulative z-score density of SGG targeting the coding or non-coding strand. Stop codon
 315 generating gRNAs targeting the non-coding strand of essential genes show higher efficiencies compared
 316 to those targeting the coding strand. **E)** Distributions of modeled RNA/DNA duplex melting temperature for
 317 all stop codon generating gRNAs, those with no significant effects, and SGG GNE. P-values were
 318 calculated using the two-sample Kolmogorov-Smirnov test. **F)** Stop codon generating gRNAs GNE
 319 enrichment compared to the expected GNE ratio for different melting temperature ranges.

320 Data from the original Target-AID study (Nishida *et al.* 2016) suggests that the most prevalent
 321 outcome for an edited site is a C-to-G transversion. Our data support this observation, as gRNAs
 322 which would lead to a C-to-G mutation at the highest activity site of the editing window have the
 323 highest GNE detection rate (Figure 4B). It was also suggested that Target-AID could modify
 324 multiple nucleotides within the activity window that could be edited during mutagenesis. Our data
 325 support this observation, as gRNAs for which two outcomes have the potential to generate a stop
 326 codon are markedly more efficient than those with only one stop codon outcome (Figure 4C). This
 327 finding also extends to gRNAs that do not generate stop codons (Figure S9A).

328 We observed that the targeted strand relative to transcription greatly influenced editing efficiency
329 (Figure 4D). This strand effect can be explained by multiple factors. First, there are multiple
330 outcomes leading to mutation to a stop codon starting from a TGG codon (shown in Figure 4A).
331 This codon is the only one that can be targeted on the non-coding strand to generate a stop
332 codon. Second, repair efficiency has been shown to be higher for the transcribed strand in yeast
333 (Reis *et al.* 2012). Finally, as the non-coding strand is the one which is transcribed, a deamination
334 event there might lead to consequences at the protein level more rapidly because it does not
335 need DNA replication to be present on both strands. gRNAs that do not generate stop codons
336 also have a higher chance of having a fitness effect if they target the non-coding strand (Figure
337 S9B), but we did not observe any effects of the chromosomal strand on efficiency (Figure S9C).

338 One other parameter with a high impact on mutagenesis rate is the predicted melting temperature
339 of the RNA-DNA duplex formed by the gRNA sequence and its target DNA sequence (Figure 4E).
340 The distribution of the melting temperature shows a clear shift between stop codon generating
341 gRNAs that have an effect on fitness and those that do not. gRNAs with low values have a lower
342 chance of being detected as having effects, while gRNAs with higher values are enriched for GNE
343 (Figure 4F). This observation also extends to gRNAs that do not generate stop codons (Figure
344 S9D, E). This enrichment cannot be attributed to technical biases in library preparation or high-
345 throughput sequencing that would tend to lower their abundance as melting temperature shows
346 practically no correlation with read count at every time point (Figure S10). Furthermore, this effect
347 is not caused by target position bias within target genes or a strong correlation between GC
348 content and the targeted position (Figure S11). As binding energy can differ drastically even within
349 groups of gRNAs with similar GC content (Figure S9F), this could provide a useful criterion to help
350 select efficient gRNAs.

351

352 **Discussion**

353 We tested whether the Target-AID base editor is amenable for genome-wide mutagenesis. Using
354 the yeast essential genes as test cases, we identified hundreds of gRNAs targeting residues with
355 significant effects on cellular fitness when mutated. The precision and traceability of Target-AID
356 genome editing allowed us to predict the mutational outcomes of GNE and to confirm their effects
357 using orthogonal approaches. We used this data to investigate which factors influence base
358 editing efficiency and found multiple gRNAs and target properties that affect mutagenesis and
359 that could be optimized for future experiments for specific genomic space. By focusing on a few
360 highly relevant variants, we highlighted the power of our approach to generate new biological
361 insights.

362 In previously published methods such as TAM and CRISPR-X (Hess *et al.* 2016; Ma *et al.* 2016),
363 the semi-random nature of the editing forces the use of mutant allele frequencies as a readout for
364 mutational fitness effects, potentially limiting the scale of the experiments because only one
365 genomic region can be targeted at a time. To complement these approaches, we use more
366 predictable base editing to increase dramatically the number of target loci, albeit at the cost of a
367 lower mutational density. Our results demonstrate the feasibility of base editing screening at a
368 large scale with applications beyond stop codon generation, and future developments will further
369 enhance it. For instance, the use of a base editor with multiple possible mutagenesis outcomes
370 complexifies the prediction of editing outcomes, which can, in turn, make GNE confirmation
371 challenging. Using a base editor that channels mutational outcomes such as cytidine deaminase-
372 uracil glycosylase inhibitor (UGI) fusion can address this problem but decreases the number of
373 mutations explored during the experiment. However, recently published data on cytidine
374 deaminase-UGI fusion has shown they could lead to off-target editing in vivo at a much higher
375 rate compared to adenine base editors or the Cas9 nuclease (Jin *et al.* 2019; Zuo *et al.* 2019).
376 Although there is currently no high throughput data on the off-target activity of Target-AID, data

377 generated in yeast in the original publication suggests far lower rates than those recently reported
378 in mammalian cells (Nishida *et al.* 2016).

379 We provide key empirical data on parameters that can be used to optimize base editing efficiency,
380 based on gRNA dependent properties such as target strand and GC content. The results we
381 observed differ from what has been reported for Cas9-based genome editing, in which high gRNA
382 RNA/DNA duplex binding has been associated with lower mutagenesis efficiency (Wong *et al.*
383 2015). Our data thus confirms the observation that parameters associated with Cas9 editing
384 cannot readily be transferred to base editors (Kim *et al.* 2017). Furthermore, the temperature at
385 which experiments are performed might affect efficiency for certain gRNAs with low gRNA-DNA
386 duplex binding energy and should be considered when designing base editing experiments in
387 different organisms. However, it remains to be confirmed whether the enrichment for certain
388 gRNA properties we observed are specific to Target-AID or will also be transferable to other base
389 editors as this may depend on the enzymatic properties of these proteins.

390 The field of base editing is rapidly evolving, with new tools being developed constantly. One of
391 the most recent additions to this fast-growing toolkit is engineered Cas9 enzymes with broadened
392 PAM specificities (Nishimasu *et al.* 2018), which have already been shown to be compatible with
393 base editors. More flexible PAM requirements are especially useful for base editing applications,
394 as they increase the number of sites to be edited and also the number of potential gRNAs per
395 site, increasing the chances of choosing optimal properties and thus greater efficiency (Dandage
396 *et al.* 2019). Our method allows an experimental scale which bridges saturation mutagenesis
397 methods and genome-wide knock-out studies, alleviating the current trade-off between mutational
398 diversity and the number of targets genes to generate new biological insights

399

400 **Methods**

401 **Generation of a gRNA library for Target-AID mutagenesis of essential genes in yeast**

402 The Target-AID base editor has an activity window between base 15 to 20 in the gRNA sequence
403 starting from the PAM, and the efficiency at these different positions was characterized in Nishida
404 *et al.* 2016. This allowed us to predict the mutational outcomes for a specific gRNA provided the
405 number of editable bases in the window is not too high. To select gRNAs, we parsed a database
406 of gRNA targets for the *S. cerevisiae* reference genome sequences (strain S288c) (Dicarlo *et al.*
407 2013) and applied several selection criteria. Since the screen was to be performed in the BY4741
408 strain, all gRNAs (unique seed sequence, no NAG site) within the database were aligned to the
409 reference genome of that strain using Bowtie (Langmead *et al.* 2009). Only gRNAs with a single
410 perfect alignment were kept for subsequent steps. To select gRNAs amenable to Target-AID base
411 editing, we selected gRNAs with cytosines within the highest activity window of the editor
412 (positions -17 to -19 starting from the PAM). To limit the total number of possible mutational
413 outcomes, gRNAs with three cytosines within the window were removed as well as those with two
414 cytosines at the highest activity positions. Next, we filtered out any gRNA containing a BsaI
415 restriction site to prevent errors during the library cloning step.

416 The list of essential genes (n=1156) (Winzeler *et al.* 1999; Giaever *et al.* 2002) was used to
417 discriminate between gRNAs targeting essential or non-essential genes (retrieved from
418 http://www-sequence.stanford.edu/group/yeast_deletion_project/Essential_ORFs.txt). Among
419 non-essential genes, data from Qian *et al.* 2012 (Qian *et al.* 2012) was used to create categories
420 of fitness effects. If the fitness score (averaged across media and replicates) of a gene was below
421 0.75, it was categorized as “high effect” on fitness. We excluded auxotrophic marker genes as
422 well as *CAN1*, *LYP1*, and *FCY1* because those could be used as co-selection markers (Després
423 *et al.* 2018). Gene deletions with an averaged fitness score between 0.999 and 1.001 were
424 categorized as having “no detectable effect” on fitness. We selected gRNAs targeting essential

425 and high effect genes, as well as gRNAs targeting a set of 38 randomly chosen no effect genes.
426 To further limit the space of gRNAs examined, only gRNAs mapping from the 0.5th percent to the
427 75th percent of coding sequences were chosen. We also added gRNAs targeting all known yeast
428 introns (Ares lab Database 4.3) (Grate and Ares 2002) and putative non-functional peptides
429 (Smith *et al.* 2014) selected with the same strategy except for the constraints on gRNA position
430 within the sequence of interest. This resulted in a set of 39,989 gRNAs: library properties are
431 summarized in Figure S1.

432 **Library construction**

433 The plasmids, oligonucleotides, and media used in this study are presented as supplementary
434 tables S6, S7 and S8 respectively. The oligo pool was synthesized by Arbor Biosciences
435 (Michigan, USA) and was cloned into the pDYSCO vector using Golden Gate Assembly (New
436 England Biolabs, Massachusetts, USA) with the following reaction parameters:

NEB GG buffer 10X	2 μ l
pDYSCO [75ng/ μ l]	1 μ l
Oligo pool [2ng/ μ l]	1 μ l
NEB GG mix	1 μ l
Water	15 μ l

437
438 The ligation mix was transformed in *E. coli* strain MC1061 (*[araD139]_{Br} Δ (araA-leu)7697 Δ lacX74*
439 *galK16 galE15(GalS) λ - e14- mcrA0 relA1 rpsL150(strR) spoT1 mcrB1 hsdR2*, Casadaban and
440 Cohen 1980) using a standard chemical transformation protocol and plated on ampicillin selective
441 media to select for transformants. Serial dilution of cells after outgrowth were plated and then
442 used to calculate the total number of clones produced by the cloning reaction. Quality control of
443 the assembly was performed by Sanger sequencing ~10 clones per assembly reaction. Cells
444 were scraped from plates by adding ~5 ml of sterile water, incubating a few minutes at room
445 temperature, and then using a glass rake to resuspend colonies. Resuspended plates were then

446 pooled together in a single flask per reaction, which was then used to make glycerol stocks of the
447 library and cell pellets for plasmid extraction. The Qiagen Midi-Prep kit (Qiagen, Germany) was
448 used to extract plasmid DNA from cell pellets by following the manufacturer's instructions. The
449 DNA concentration of each eluate was then measured using a NanoDrop (ThermoFisher,
450 Massachusetts, USA), and a normalized master library for yeast transformation was assembled
451 by combining equal quantities of each assembly pool.

452 **Library transformation in yeast**

453 Competent BY4741 (*MATa his3Δ1 leu2Δ0 met15Δ0 ura3Δ0*) cells were first transformed with the
454 pKN1252 (p415-Gall-Target-AID) plasmid using a standard lithium acetate method (Gietz and
455 Schiestl 2007). Transformants were selected by plating cells on SC-L. After 48 h of growth,
456 multiples colonies were used to inoculate a starter liquid culture for competent cells preparation
457 using the standard lithium acetate protocol (Gietz and Schiestl 2007): a culture volume of 200 ml
458 was used to generate enough competent cells for mass transformation. The large-scale library
459 transformation was performed by combining 40 transformation reactions performed with 40 ul of
460 competent cells and 5 ul plasmid library (240 ng/ul) after the outgrowth stage and plating 100 ul
461 aliquots on SC-UL: cells were then allowed to grow at 30°C for 48 h. A 1/1000 serial dilution of
462 the cell recovery was plated in 5 replicates and used to calculate the number of transformants
463 obtained. The total number of transformants reached 3.48×10^6 CFU, corresponding to about
464 100X coverage of the plasmid pool.

465 **Target-AID mutagenesis and competition screening**

466 The mutagenesis protocol is an upscaled version of our previously published method and is
467 shown in Figure S2. Transformants were scraped by spreading 5 ml sterile water on plates and
468 then resuspending cells using a glass rake. All plates were pooled together in the same flask, and
469 the OD of the yeast resuspension was measured using a Tecan Infinite F200 plate reader (Tecan,
470 Switzerland). Pellets corresponding to about 6×10^8 cells were washed twice with SC-UL without

471 a carbon source and then used to inoculate a 100 ml SC-UL +2% glucose culture at 0.6 OD two
472 times to generate replicates A and B. Cells were allowed to grow for 8 hours before 1×10^9 cells
473 were pelleted and used to inoculate a 100 ml SC-UL + 5% glycerol culture. After 24 hours, $5 \times$
474 10^8 cells were pelleted and either put in SC-UL + 5% galactose for mutagenesis or SC-UL + 5%
475 glucose for a mock induction control. Target-AID expression (from pKN1252) was induced for 12
476 hours before 1×10^8 cells were pelleted and used to inoculate a canavanine (50 μ g/ml) co-
477 selection culture in SC-ULR. After 16 hours of incubation, 5×10^7 cells of each culture were used
478 to inoculate 100 ml SC-UR, which was grown for 12 hours before 5×10^7 cells were used to
479 inoculate a final 100 ml SC-UR culture which was grown for another 12 hours. Cell pellets were
480 washed with sterile water between each step, and all incubation occurred at 30°C with agitation.
481 $\sim 2 \times 10^7$ cells were taken for plasmid DNA extraction at the end of each mutagenesis and
482 competition screening step.

483 **Yeast plasmid DNA extraction**

484 Yeast plasmid DNA was extracted using the ChargeSwitch Plasmid Yeast Mini Kit (Invitrogen,
485 California, USA) by following the manufacturer's protocol with minor modifications: Zymolase
486 4000 U/ml (Zymo Research, California, USA) was used instead of lyticase, and cells were
487 incubated for 1 hour at room temperature, one min at -80°C, and then incubated for another 15
488 minutes at room temperature before the lysis step. Plasmid DNA was eluted in 70 μ l of E5 buffer
489 (10 mM Tris-HCl, pH 8.5) and stored at -20°C for use in library preparation.

490 **Next-generation library sequencing preparation**

491 Libraries were prepared by using two PCR amplification steps, one to amplify the gRNA region of
492 the pDSYCKO plasmid pool and the second to add sample barcodes as well as the Illumina p5
493 and p7 sequences (Yachie *et al.* 2016). Oligonucleotides for library preparation are shown in the
494 first part of the oligonucleotide table. Reaction conditions for the first PCR were as follows:

495

Phusion HF buffer (NEB) 5X	5 μ l
dNTPs 10 mM	0.5 μ l
pDYSCO_gRNA_for 10 μ M	1.25 μ l
pDYSCO_gRNA_rev 10 μ M	1.25 μ l
Phusion polymerase	0.5 μ l
Template DNA (<1 ng/ μ l)	5 μ l
PCR grade water	11.7 μ l

496

497 Thermocycler protocol:

Temperature ($^{\circ}$ C)	Time (s)	Cycles
98	30	1
98	10	16
58	15	
72	5	
72	5	1

498

499 The resulting product was verified on a 2% agarose gel colored with Midori Green Advance
500 (Nippon Genetics, Japan) and then gel-extracted and purified using the FastGene Gel/PCR
501 Extraction Kit (Nippon Genetics, Japan). The purified products were used as the template for the
502 second PCR reaction, with the following conditions:

Phusion Mastermix-HF (NEB)	10 μ l
P5-barcode-X oligo 1.333 μ M	3.75 μ l
P7-barcode-Y oligo 1.333 μ M	3.75 μ l
Template DNA (~1 ng/ μ l)	2.5 μ l

503

504 Thermocycler protocol:

Temperature (°C)	Time (s)	Cycles
98	30	1
98	10	15
60	10	
72	60	
72	300	1

505

506 PCR products were verified on a 2% agarose gel colored with Midori Green Advance (Nippon
507 Genetics, Japan) and then gel-extracted and purified using the FastGene Gel/PCR Extraction Kit
508 (Nippon Genetics, Japan). Library quality control and quantification were performed using the
509 KAPA Library Quantification Kit for Illumina platforms (Kapa Biosystems, Massachusetts, USA)
510 following the manufacturer's instructions. Libraries were then run on a single lane on HiSeq 2500
511 (Illumina, California, USA) with paired-end 150 bp in fast mode.

512 **Large-scale screen sequencing data analysis**

513 The custom Python scripts used to analyze the data will be made available on github
514 (<https://docker.pkg.github.com/Landrylab>), packages and software used are presented in
515 Supplementary table 9. Raw sequencing files have been deposited on the NCBI SRA, accession
516 number PRJNA552472. Briefly, reads were separated into three subsequences for alignment: the
517 P5 barcode, the gRNA, and the P7 barcode. Each of these was aligned using Bowtie (Langmead
518 *et al.* 2009) to an artificial reference genome containing either the barcodes or gRNA sequences
519 flanked by the common amplicon sequences. The gRNA sequences are aligned both with 0 or 1
520 mismatch allowed, and misalignment position and type were stored. Information on barcode and
521 gRNA alignment for each read was stored and combined to generate a barcode count per library
522 table, a list of mismatches in alignments for each gRNA in each library, as well as mismatch types
523 and counts for the same gRNA across all libraries.

524 Synthesis error within oligonucleotide libraries is one of the major limits of current large-scale
525 genome editing screening methods. These errors can introduce gRNA sequences that cannot
526 perform mutagenesis because the gRNA sequence does not match a site in the genome. We
527 refer to those gRNAs as SE gRNAs. In our experiment, the stringent selection criteria used to
528 select gRNAs limited the risk of off-target effects even for gRNAs with one mismatch, minimizing
529 the risk that a synthesis error gRNA could lead to editing at another site in the genome. We
530 therefore decided to use highly abundant SE gRNAs as negative controls to obtain a null
531 distribution of abundance variation for gRNAs with no fitness effects. To differentiate synthesis
532 errors from sequencing errors, we used the mismatch type and count table to assess whether a
533 particular mismatched gRNA constitutes a too large fraction of the reads associated with a gRNA
534 to be simply a repeated sequencing error. For each error, we test if:

535
$$\frac{N_{readsformismatch}}{N_{perfectalignment}} > 0.075$$

536 and discarded the reads associated with the specific mismatch alignment. This threshold was
537 obtained by iteratively testing different threshold values in an effort to maximize the gain in gRNA
538 counts while minimizing the noise added by incorrect assignments. Read counts per library for
539 abundant ($N_{readsformismatch} > 1,000$) SE gRNAs were kept to serve as negative controls when
540 measuring fitness effects, resulting in a set of 1,032 abundant SE gRNAs. gRNAs absent from
541 more than half of the libraries (4446 out of 39,989) were removed from the analysis before gRNA
542 abundance calculations.

543 **Detecting mutations with high fitness effects**

544 Barcode sequencing competition experiments use DNA barcodes to measure the relative
545 abundance of many different subpopulations of cells grown in the same pool (Robinson *et al.*
546 2014). Since each gRNA is linked to its possible mutagenesis outcomes, we can use relative
547 gRNA abundance to detect mutations with significant fitness effects. To do so, the \log_2 of the

548 relative abundance of a barcode after mutagenesis is compared with its abundance at the end of
549 the screen:

$$550 \quad \Delta \log_2_{gRNA} = \log_2 \left(\frac{N_{reads_{gRNA}t_1}}{N_{readst_1}} \right) - \log_2 \left(\frac{N_{reads_{gRNA}t_0}}{N_{readst_0}} \right)$$

551 For each gRNA, the measured fitness effect is the product of the effect of the mutational outcomes
552 on growth and of the mutation rate within the cell subpopulation bearing this particular gRNA.
553 Relative counts will also vary stochastically because of variation in sequencing coverage
554 depending on the time point and replicate. To reduce the impact of these effects, a minimal read
555 count at the end of the galactose induction step was used to filter out low abundance gRNAs. We
556 found a minimal read threshold of $n=54$ provided a good tradeoff between the number of gRNAs
557 eligible for analysis and inter-replicate correlation.

558 Using the distribution of $\Delta \log_2$ values, we calculated a z-score for each gRNA in both replicates.
559 We then averaged z-scores between replicates and compared the score distributions between
560 SE and Non-SE gRNAs. This revealed the presence of a left-skewed tail in the z-score distribution
561 of valid gRNAs, which is absent in the SE. Because the number of SE gRNAs is smaller than the
562 one of functional gRNAs by almost two orders of magnitude, a type I error (false positives)
563 empirical threshold based solely on a weighted SE z-score distribution was not practical. To
564 resolve this, we fitted a Gumbell left skewed distribution to the SE gRNAs z-score distribution and
565 used it to approximate the type I error rate as a function of the z-score. We set a significance
566 threshold such as that all gRNAs at z-scores for which the estimated false positive rate is below
567 or equal to 5% are considered GNEs.

568 **Complementation assays**

569 Experiments were performed in heterozygous deletion mutants from the YKO project
570 heterozygous deletion strain set (Dharmacon, Colorado, USA). For each gene, a single colony

571 streaked from the glycerol stock was used to prepare competent cells using the previously
572 described lithium acetate protocol. To generate mutant alleles of the genes of interest, we
573 performed site-directed mutagenesis on the appropriate MoBY collection plasmid (Ho *et al.* 2009).
574 These centromeric plasmids encode the yeast gene of interest under the control of their native
575 promoters and terminators. Mutagenesis reactions were performed with the following reaction
576 setup:

577

Kapa HiFi buffer (Kapa biosciences) 5X	5 μ l
dNTPs 10 μ M	0.75 μ l
mutation_for 10 μ M (see table 7)	0.75 μ l
mutation_rev 10 μ M (see table 7)	0.75 μ l
Kapa Hot-start polymerase	0.5 μ l
Template plasmid DNA (15ng/ μ l)	0.75 μ l
PCR grade water	16.5 μ l

578

579 Thermocycler protocol:

Temperature ($^{\circ}$ C)	Time (s)	Cycles
95	300	1
98	20	20
60	15	
72	720	
72	1080	1

580

581 After amplification, the mutagenesis product was digested with DpnI for 2 hours at 37 $^{\circ}$ C and 5 μ l
582 was transformed in *E. coli* strain BW23474 (F⁻, Δ (*argF-lac*)169, Δ *uidA4::pir-116*, *recA1*,
583 *rpoS396(Am)*, *endA9(del-ins)::FRT*, *rph-1*, *hsdR514*, *rob-1*, *creC510*, Haldimann *et al.* 1996).
584 Transformants were plated on 2YT+Kan+Chlo and grown at 37 $^{\circ}$ C overnight. Plasmid DNA was

585 then isolated from clones and sent for Sanger sequencing (CHUL sequencing platform, Université
586 Laval, Québec City, Canada) to confirm mutagenesis success.

587 Competent cells of target genes were transformed with the appropriate mutant plasmids as well
588 a the original plasmid bearing the wild-type gene and the empty vector (Zhao *et al.* 2016), and
589 transformants were selected by plating on SC-U (MSG). Multiple independent colonies per
590 transformation were then put on sporulation media until sporulation could be confirmed by
591 microscopy. For tetrad dissection, cells were resuspended in 100ul 20T zymolyase (200mg/ml
592 dilution in water) and incubated for 20 minutes at room temperature. Cells were then centrifuged
593 and resuspended in 50ul 1M sorbitol before being streaked on a level YPD plate. All dissections
594 were performed using a Singer SporePlay microscope (Singer Instruments, UK). Plate pictures
595 were taken after five days incubation at room temperature except for the RAP1 plasmid
596 complementation test for which the picture was taken after three days. Pictures are shown in
597 Supplementary image 1.

598

599 **Strain construction for confirmations in *RAP1***

600 Because the MoBY collection plasmid for RAP1 cannot fully complement the gene deletion
601 (Supplementary image file 1), we instead performed confirmations by engineering mutations a
602 diploid strain to create heterozygous mutants. *RAP1* was first tagged with a modified version of
603 fragment DHFR F[1,2] (the first half) of the mDHFR enzyme (Tarassov *et al.* 2008). The
604 mDHFR[1,2]-FLAG cassette was amplified using gene-specific primers and previously described
605 reaction parameters (Tarassov *et al.* 2008). Cells were transformed with the cassette using the
606 previously described transformation protocol and were plated on YPD+Nourseothricine (YPD+Nat
607 in Media table). Positive clones were identified by colony PCR and successful fragment fusion
608 was confirmed by Sanger sequencing (CHUL sequencing platform). We then mated the confirmed

609 clones with strain Y8205 (*Mata can1::STE2pr-his5 lyp1::STE3prLEU2 Δura3 Δhis3 Δleu2*, Kindly
610 gifted by Charlie Boone) by inoculating a 4ml YPD culture with overnight starter cultures of both
611 strains and letting the culture grow overnight. Cells were then streaked on YPD+Nat and diploid
612 cells were identified by colony PCR using mating type diagnosis primers (Huxley *et al.* 1990).

613 To create heterozygous deletion mutants of the target gene, we amplified a modified version of
614 the *URA3* cassettes that could then be targeted with the CRISPR-Cas9 system to integrate our
615 mutations of interest using homologous recombination at the target locus. The oligonucleotides
616 we used differ from those commonly used in that they amplify the cassette without the two LoxP
617 sites present at both ends. We found it necessary to remove those sites as one common
618 mutational outcome after introducing a double-stranded break in the *URA3* cassette was inter-
619 LoxP site recombination without the integration of donor DNA at the target locus. These modified
620 cassettes recombine with DNA upstream the target gene on one end and the mDHFR F[1,2]
621 fusion on the other, ensuring that the heterozygous deletion is always performed at the locus that
622 is already tagged. Cassettes were transformed using the standard lithium acetate method, and
623 cells were plated on SC-U (MSG) selective media. Heterozygous deletion mutants were then
624 confirmed by colony PCR.

625 **CRISPR-Cas9 mediated Knock-in of targeted mutations**

626 Mutant alleles of target genes were amplified in two fragments using template DNA from the
627 haploid tagged strain (See Figure S8). The two fragments bearing mutations are then fused
628 together by a second PCR round to form the final donor DNA. This DNA was then co-transformed
629 with a plasmid bearing Cas9 and a gRNA targeting the *URA3* cassette for HDR mediated editing
630 using a standard protocol (Ryan *et al.* 2016). Clones were then screened by PCR to verify donor
631 DNA and mutation integration at the target locus. The targeted region of *RAP1* was then Sanger
632 sequenced (CHUL sequencing platform, Université Laval, Québec City, Canada) to confirm the
633 presence of the mutation of interest. Heterozygous mutants were sporulated on solid media until

634 sporulation could be confirmed by microscopy using the same protocol previously described. The
635 plates were then replica plated on YPD+Nat media, and the pictures were taken after five days at
636 room temperature Supplementary image 2.

637 **Evolutionary rate measurements and protein variant abundance**

638 Evolutionary rates were calculated using the Rate4site software (Mayrose *et al.* 2004) using
639 multiple sequence alignments and phylogenies from PhylomeDB V4 (Huerta-Cepas *et al.* 2014)
640 as input and using the raw calculated rates as output. Variant data was compiled using data from
641 the 1002 Yeast Genome Project ([http://1002genomes.u-strasbg.fr/files/
642 allReferenceGenesWithSNPsAndIndelsInferred.tar.gz](http://1002genomes.u-strasbg.fr/files/allReferenceGenesWithSNPsAndIndelsInferred.tar.gz)). Strain-specific protein coding sequence
643 were aligned to the S288c sequence using Fastx36 (Pearson *et al.* 1997) with the following
644 parameters: `fastx36 -p -s -VT10 -T 6 -m 10 -n -3 querymultifasta.fasta
645 ref_orf.db 12\> fasta_out`. Alignments were then parsed with a custom Python script to
646 identify variants. Variant abundance was measured as the number of strains in the dataset in
647 which a specific variant was found. If the coding sequence contained ambiguous nucleotides (ex.:
648 R or Y), separate coding sequences were generated for each possibility and each possible variant
649 was considered as a separate occurrence.

650 **Analysis of the properties of stop codon generating gRNAs**

651 To analyse the sequence and target properties of gRNA inducing the creation of stop codons,
652 data from multiple sources was compiled. For each target gene, length and chromosomal strand
653 was obtained from the Saccharomyces Genome Database using the Yeastmine query interface
654 (Cherry *et al.* 2012). Distance to centromere was obtained by calculating the minimal distance
655 between the start of the gene and one extremity of the centromere coordinates. RNA:DNA duplex
656 melting temperature of gRNA sequence with target genomic DNA was calculated using the
657 MeltingTemp module from Biopython (Cock *et al.* 2009), which uses values taken from Sugimoto

658 et al (Sugimoto *et al.* 1995). Correlation between gRNA/DNA duplex melting temperatures was
659 assessed using Spearman's rank correlation.

660 **Variant effect prediction resources analysis and GO enrichment**

661 All prediction data except the Envision scores were extracted from the aggregated data of the
662 Mutfunc database (Wagih *et al.* 2018). Precomputed values were downloaded directly from the
663 FTP server (http://ftp.ebi.ac.uk/pub/databases/mutfunc/mutfunc_v1/yeast/). This database
664 includes precomputed SIFT scores for 5498 yeast proteins, as well as predicted variant ddG
665 values based on protein structure (n=1057), homology models (n=1703) and protein-protein
666 interaction interfaces (n=1109). Mutations with $\Delta\Delta G > 1$ considered destabilizing.

667 Precomputed values from Envision (Gray *et al.* 2018) were downloaded directly from the database
668 website (https://envision.gs.washington.edu/shiny/envision_new/, file yeast_predicted_2017-03-
669 12.csv). This file contained 34857830 mutation effect predictions spread across 4011 genes. The
670 distribution of Envision scores for the genes targeted in the experiment that are included in the
671 database are shown in (Figure S6).

672 Gene enrichments were performed using the PANTHER gene list analysis tool (Mi *et al.* 2019).
673 The list of genes for which 2 or more GNEs were detected was tested for enrichment against all
674 genes targeted by the library using Fisher's exact test and False Discovery Rate calculations. The
675 Gene Ontology datasets used were: GO molecular function complete, GO biological process
676 complete, and GO cellular component complete.

677 Supplementary dataset 1 contains all gRNAs, their z-scores values as well all the information and
678 annotations used in data analysis.

679

680

681 **Acknowledgments**

682 This work was supported by the Canadian Institutes of Health Research Foundation grant 387697
683 to C.R.L., as well as project grants 364920, 384483, a Frederick Banting and Charles Best graduate
684 scholarship and a Vanier graduate scholarship to P.C.D, by Université Laval via an André
685 Darveau Fellowship to P.C.D., the Fonds Québécois de Recherche en Santé via a Master's
686 training award to P.C.D. and the Japan Society for the Promotion of Science grant numbers
687 S15734 and S17161 to C.R.L. and N.Y. The authors thank Mathieu Hénault, Johan Hallin, and
688 Dan Yamamoto Evans for comments on the manuscript, as well as Maria Isabel Acosta Lopez
689 for assistance during the strain construction process.

690 **Author contributions**

691 PCD, AKD, NY and CRL designed research. PCD and AKD performed experiments. PCD and
692 MS generated NGS sequencing data. All data analysis was performed by PCD with input from
693 CRL. PCD and CRL wrote the manuscript with input from all authors.

694 **Conflict of interest**

695 None to declare

696 **References**

- 697 Albuquerque C. P., M. B. Smolka, S. H. Payne, V. Bafna, J. Eng, *et al.*, 2008 A multidimensional
698 chromatography technology for in-depth phosphoproteome analysis. *Mol. Cell. Proteomics*
699 7: 1389–96. <https://doi.org/10.1074/mcp.M700468-MCP200>
- 700 Bao Z., M. Hamedirad, P. Xue, H. Xiao, I. Tasan, *et al.*, 2018 Genome-scale engineering of
701 *Saccharomyces cerevisiae* with single-nucleotide precision. *Nat. Biotechnol.*
702 <https://doi.org/10.1038/nbt.4132>
- 703 *C. elegans* Deletion Mutant Consortium T. *C. elegans* D. M., 2012 Large-Scale Screening for
704 Targeted Knockouts in the *Caenorhabditis elegans* Genome. *G3;*
705 *Genes|Genomes|Genetics* 2: 1415–1425. <https://doi.org/10.1534/g3.112.003830>
- 706 Casadaban M. J., and S. N. Cohen, 1980 Analysis of gene control signals by DNA fusion and
707 cloning in *Escherichia coli*. *J. Mol. Biol.* 138: 179–207. <https://doi.org/10.1016/0022->

- 708 2836(80)90283-1
- 709 Cherry J. M., E. L. Hong, C. Amundsen, R. Balakrishnan, G. Binkley, *et al.*, 2012
710 Saccharomyces Genome Database: The genomics resource of budding yeast. *Nucleic*
711 *Acids Res.* <https://doi.org/10.1093/nar/gkr1029>
- 712 Cock P. J. A., T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox, *et al.*, 2009 Biopython: Freely
713 available Python tools for computational molecular biology and bioinformatics.
714 *Bioinformatics.* <https://doi.org/10.1093/bioinformatics/btp163>
- 715 Copley R. R., and G. J. Barton, 1994 A Structural Analysis of Phosphate and Sulphate Binding
716 Sites in Proteins. *J. Mol. Biol.* 242: 321–329. <https://doi.org/10.1006/jmbi.1994.1583>
- 717 Dandage R., P. C. Després, N. Yachie, and C. R. Landry, 2019 beditor: A Computational
718 Workflow for Designing Libraries of Guide RNAs for CRISPR-Mediated Base Editing.
719 *Genetics* 212: 377–385. <https://doi.org/10.1534/genetics.119.302089>
- 720 DePristo M. A., D. M. Weinreich, and D. L. Hartl, 2005 Missense meanderings in sequence
721 space: a biophysical view of protein evolution. *Nat. Rev. Genet.* 6: 678–687.
722 <https://doi.org/10.1038/nrg1672>
- 723 Després P. C., A. K. Dubé, L. Nielly-Thibault, N. Yachie, and C. R. Landry, 2018 Double
724 Selection Enhances the Efficiency of Target-AID and Cas9-Based Genome Editing in
725 Yeast. *G3 (Bethesda)*. g3.200461.2018. <https://doi.org/10.1534/g3.118.200461>
- 726 Dicarolo J. E., J. E. Norville, P. Mali, X. Rios, J. Aach, *et al.*, 2013 Genome engineering in
727 *Saccharomyces cerevisiae* using CRISPR-Cas systems. *Nucleic Acids Res.* 41: 4336–
728 4343. <https://doi.org/10.1093/nar/gkt135>
- 729 Doench J. G., E. Hartenian, D. B. Graham, Z. Tothova, M. Hegde, *et al.*, 2014 Rational design
730 of highly active sgRNAs for CRISPR-Cas9-mediated gene inactivation. *Nat. Biotechnol.* 32:
731 1262–1267. <https://doi.org/10.1038/nbt.3026>
- 732 Doench J. G., N. Fusi, M. Sullender, M. Hegde, E. W. Vaimberg, *et al.*, 2016 Optimized sgRNA
733 design to maximize activity and minimize off-target effects of CRISPR-Cas9. *Nat.*
734 *Biotechnol.* 34: 184–191. <https://doi.org/10.1038/nbt.3437>
- 735 Eriani G., M. Delarue, O. Poch, J. Gangloff, and D. Moras, 1990 Partition of tRNA synthetases
736 into two classes based on mutually exclusive sets of sequence motifs. *Nature* 347: 203–
737 206. <https://doi.org/10.1038/347203a0>
- 738 Fowler D. M., and S. Fields, 2014 Deep mutational scanning: a new style of protein science.
739 *Nat. Methods* 11: 801–7. <https://doi.org/10.1038/nmeth.3027>
- 740 Gaudelli N. M., A. C. Komor, H. A. Rees, M. S. Packer, A. H. Badran, *et al.*, 2017
741 Programmable base editing of A•T to G•C in genomic DNA without DNA cleavage. *Nature*
742 551: 464–471. <https://doi.org/10.1038/nature24644>
- 743 Giaever G., D. D. Shoemaker, T. W. Jones, H. Liang, E. A. Winzeler, *et al.*, 1999 Genomic
744 profiling of drug sensitivities via induced haploinsufficiency. *Nat. Genet.* 21: 278–83.
745 <https://doi.org/10.1038/6791>
- 746 Giaever G., A. M. Chu, L. Ni, C. Connelly, L. Riles, *et al.*, 2002 Functional profiling of the
747 *Saccharomyces cerevisiae* genome. *Nature* 418: 387–391.
748 <https://doi.org/10.1038/nature00935>

- 749 Gietz R. D., and R. H. Schiestl, 2007 High-efficiency yeast transformation using the LiAc/SS
750 carrier DNA/PEG method. *Nat. Protoc.* 2: 31–34. <https://doi.org/10.1038/nprot.2007.13>
- 751 Graham I. R., R. A. Haw, K. G. Spink, K. A. Halden, and A. Chambers, 1999 In vivo analysis of
752 functional regions within yeast Rap1p. *Mol. Cell. Biol.* 19: 7481–90.
753 <https://doi.org/10.1128/mcb.19.11.7481>
- 754 Grant T. D., J. R. Luft, J. R. Wolfley, M. E. Snell, H. Tsuruta, *et al.*, 2013 The structure of yeast
755 glutaminyl-tRNA synthetase and modeling of its interaction with tRNA. *J. Mol. Biol.* 425:
756 2480–2493. <https://doi.org/10.1016/j.jmb.2013.03.043>
- 757 Grate L., and M. Ares, 2002 Searching yeast intron data at Ares lab web site. *Methods*
758 *Enzymol.* [https://doi.org/10.1016/S0076-6879\(02\)50975-7](https://doi.org/10.1016/S0076-6879(02)50975-7)
- 759 Gray V. E., R. J. Hause, J. Luebeck, J. Shendure, and D. M. Fowler, 2018 Quantitative
760 Missense Variant Effect Prediction Using Large-Scale Mutagenesis Data. *Cell Syst.*
761 <https://doi.org/10.1016/j.cels.2017.11.003>
- 762 Haldimann A., M. K. Prahalad, S. L. Fisher, S. K. Kim, C. T. Walsh, *et al.*, 1996 Altered
763 recognition mutants of the response regulator PhoB: a new genetic strategy for studying
764 protein-protein interactions. *Proc. Natl. Acad. Sci. U. S. A.* 93: 14361–6.
765 <https://doi.org/10.1073/pnas.93.25.14361>
- 766 Hess G. T., L. Frésard, K. Han, C. H. Lee, A. Li, *et al.*, 2016 Directed evolution using dCas9-
767 targeted somatic hypermutation in mammalian cells. *Nat. Methods.*
768 <https://doi.org/10.1038/nmeth.4038>
- 769 Ho C. H., L. Magtanong, S. L. Barker, D. Gresham, S. Nishimura, *et al.*, 2009 A molecular
770 barcoded yeast ORF library enables mode-of-action analysis of bioactive compounds. *Nat.*
771 *Biotechnol.* 27: 369–377. <https://doi.org/10.1038/nbt.1534>
- 772 Holt L. J., B. B. Tuch, J. Villén, A. D. Johnson, S. P. Gygi, *et al.*, 2009 Global analysis of Cdk1
773 substrate phosphorylation sites provides insights into evolution. *Science* 325: 1682–6.
774 <https://doi.org/10.1126/science.1172867>
- 775 Huerta-Cepas J., S. Capella-Gutiérrez, L. P. Pryszcz, M. Marcet-Houben, and T. Gabaldón,
776 2014 PhylomeDB v4: Zooming into the plurality of evolutionary histories of a genome.
777 *Nucleic Acids Res.* <https://doi.org/10.1093/nar/gkt1177>
- 778 Huxley C., E. D. Green, and I. Dunham, 1990 Rapid assessment of *S. cerevisiae* mating type by
779 PCR. *Trends Genet.* 6: 236.
- 780 Jin S., Y. Zong, Q. Gao, Z. Zhu, Y. Wang, *et al.*, 2019 Cytosine, but not adenine, base editors
781 induce genome-wide off-target mutations in rice. *Science* eaaw7166.
782 <https://doi.org/10.1126/science.aaw7166>
- 783 Kim D., K. Lim, S. T. Kim, S. H. Yoon, K. Kim, *et al.*, 2017 Genome-wide target specificities of
784 CRISPR RNA-guided programmable deaminases. *Nat. Biotechnol.*
785 <https://doi.org/10.1038/nbt.3852>
- 786 König P., R. Giraldo, L. Chapman, and D. Rhodes, 1996 The crystal structure of the DNA-
787 binding domain of yeast RAP1 in complex with telomeric DNA. *Cell* 85: 125–36.
788 [https://doi.org/10.1016/S0092-8674\(00\)81088-0](https://doi.org/10.1016/S0092-8674(00)81088-0)
- 789 Landry C. R., E. D. Levy, and S. W. Michnick, 2009 Weak functional constraints on
790 phosphoproteomes. *Trends Genet.* 25: 193–7. <https://doi.org/10.1016/j.tig.2009.03.003>

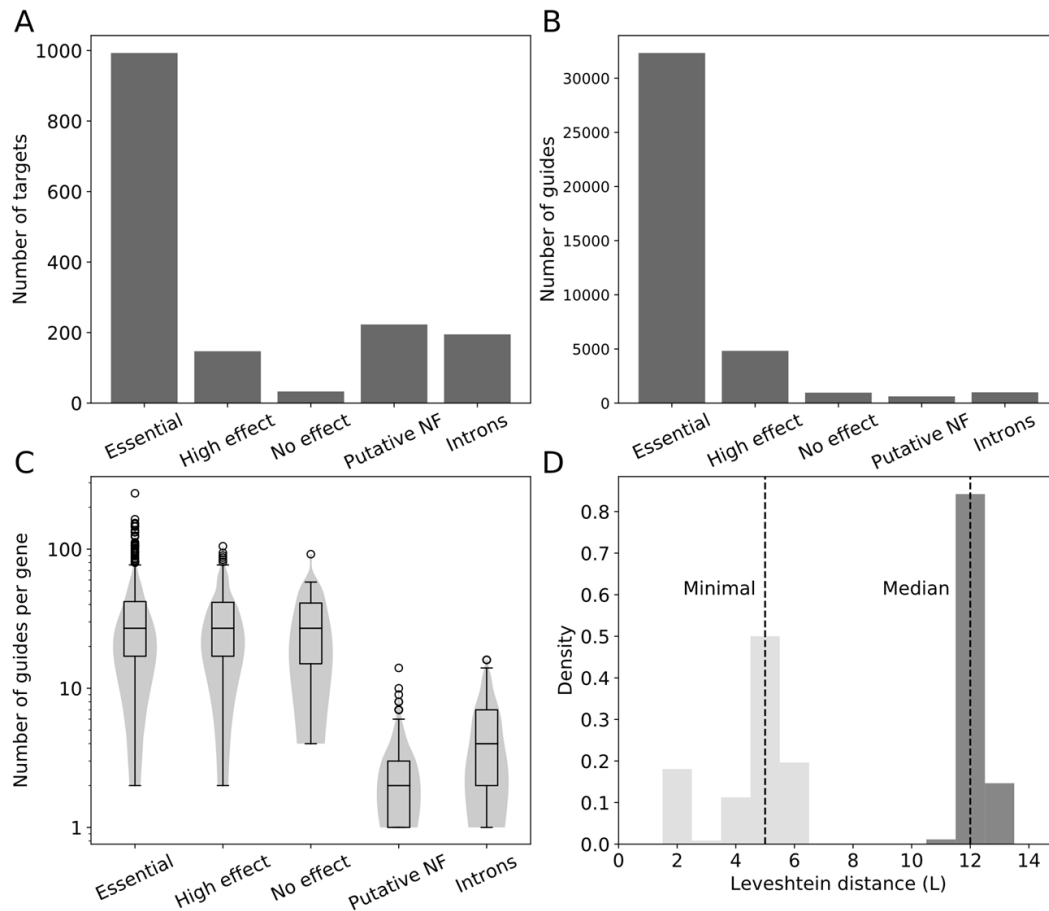
- 791 Langmead B., C. Trapnell, M. Pop, and S. L. Salzberg, 2009 Ultrafast and memory-efficient
792 alignment of short DNA sequences to the human genome. *Genome Biol.* 10: R25.
793 <https://doi.org/10.1186/gb-2009-10-3-r25>
- 794 Ma Y., J. Zhang, W. Yin, Z. Zhang, Y. Song, *et al.*, 2016 Targeted AID-mediated mutagenesis
795 (TAM) enables efficient genomic diversification in mammalian cells. *Nat. Methods* 13:
796 1029–1035. <https://doi.org/10.1038/nmeth.4027>
- 797 Mayrose I., D. Graur, N. Ben-Tal, and T. Pupko, 2004 Comparison of site-specific rate-inference
798 methods for protein sequences: Empirical Bayesian methods are superior. *Mol. Biol. Evol.*
799 <https://doi.org/10.1093/molbev/msh194>
- 800 Mi H., A. Muruganujan, X. Huang, D. Ebert, C. Mills, *et al.*, 2019 Protocol Update for large-scale
801 genome and gene function analysis with the PANTHER classification system (v.14.0). *Nat.*
802 *Protoc.* 14: 703–721. <https://doi.org/10.1038/s41596-019-0128-8>
- 803 Michel A. H., R. Hatakeyama, P. Kimmig, M. Arter, M. Peter, *et al.*, 2017 Functional mapping of
804 yeast genomes by saturated transposition. *Elife* 6. <https://doi.org/10.7554/eLife.23570>
- 805 Ng P. C., and S. Henikoff, 2003 SIFT: Predicting amino acid changes that affect protein
806 function. *Nucleic Acids Res.*
- 807 Nishida K., T. Arazoe, N. Yachie, S. Banno, M. Kakimoto, *et al.*, 2016 Targeted nucleotide
808 editing using hybrid prokaryotic and vertebrate adaptive immune systems. *Science* (80-.).
809 353: 553–563. <https://doi.org/10.1126/science.aaf8729>
- 810 Nishimasu H., X. Shi, S. Ishiguro, L. Gao, S. Hirano, *et al.*, 2018 Engineered CRISPR-Cas9
811 nuclease with expanded targeting space. *Science* eaas9129.
812 <https://doi.org/10.1126/science.aas9129>
- 813 Pearson W. R., T. Wood, Z. Zhang, and W. Miller, 1997 Comparison of DNA Sequences with
814 Protein Sequences. *Genomics* 46: 24–36. <https://doi.org/10.1006/geno.1997.4995>
- 815 Qi L. S., M. H. Larson, L. A. Gilbert, J. A. Doudna, J. S. Weissman, *et al.*, 2013 Repurposing
816 CRISPR as an RNA-guided platform for sequence-specific control of gene expression. *Cell*
817 152: 1173–83. <https://doi.org/10.1016/j.cell.2013.02.022>
- 818 Qian W., D. Ma, C. Xiao, Z. Wang, and J. Zhang, 2012 The Genomic Landscape and
819 Evolutionary Resolution of Antagonistic Pleiotropy in Yeast. *Cell Rep.* 2: 1399–1410.
820 <https://doi.org/10.1016/j.celrep.2012.09.017>
- 821 Rees H. A., and D. R. Liu, 2018 Base editing: precision chemistry on the genome and
822 transcriptome of living cells. *Nat. Rev. Genet.* 19: 770–788. <https://doi.org/10.1038/s41576-018-0059-1>
- 824 Reis A. M. C., W. K. Mills, I. Ramachandran, E. C. Friedberg, D. Thompson, *et al.*, 2012
825 Targeted detection of in vivo endogenous DNA base damage reveals preferential base
826 excision repair in the transcribed strand. *Nucleic Acids Res.* 40: 206–219.
827 <https://doi.org/10.1093/nar/gkr704>
- 828 Roy K. R., J. D. Smith, S. C. Vonesch, G. Lin, C. S. Tu, *et al.*, 2018 Multiplexed precision
829 genome editing with trackable genomic barcodes in yeast. *Nat. Biotechnol.*
830 <https://doi.org/10.1038/nbt.4137>
- 831 Ryan O. W., S. Poddar, and J. H. D. Cate, 2016 Crispr–cas9 genome engineering in
832 *Saccharomyces cerevisiae* cells. *Cold Spring Harb. Protoc.* 2016: 525–533.

- 833 <https://doi.org/10.1101/pdb.prot086827>
- 834 Sander J. D., and J. K. Joung, 2014 CRISPR-Cas systems for editing, regulating and targeting
835 genomes. *Nat. Biotechnol.* 32: 347–55. <https://doi.org/10.1038/nbt.2842>
- 836 Sanjana N. E., O. Shalem, and F. Zhang, 2014 Improved vectors and genome-wide libraries for
837 CRISPR screening. *Nat. Methods* 11: 783–784. <https://doi.org/10.1038/nmeth.3047>
- 838 Schmitt E., M. Panvert, S. Blanquet, and Y. Mechulam, 1995 Transition state stabilization by the
839 “high” motif of class I aminoacyl-tRNA synthetases: The case of *Escherichia coli* methionyl-
840 tRNA synthetase. *Nucleic Acids Res.* <https://doi.org/10.1093/nar/23.23.4793>
- 841 Schymkowitz J., J. Borg, F. Stricher, R. Nys, F. Rousseau, *et al.*, 2005 The FoldX web server:
842 an online force field. *Nucleic Acids Res.* 33: W382–W388.
843 <https://doi.org/10.1093/nar/gki387>
- 844 Shalem O., N. E. Sanjana, E. Hartenian, X. Shi, D. A. Scott, *et al.*, 2014 Genome-scale
845 CRISPR-Cas9 knockout screening in human cells. *Science* (80-.). 343: 84–87.
846 <https://doi.org/10.1126/science.1247005>
- 847 Sharon E., S. A. A. Chen, N. M. Khosla, J. D. Smith, J. K. Pritchard, *et al.*, 2018 Functional
848 Genetic Variants Revealed by Massively Parallel Precise Genome Editing. *Cell*.
849 <https://doi.org/10.1016/j.cell.2018.08.057>
- 850 Smith J. E., J. R. Alvarez-Dominguez, N. Kline, N. J. Huynh, S. Geisler, *et al.*, 2014 Translation
851 of Small Open Reading Frames within Unannotated RNA Transcripts in *Saccharomyces*
852 *cerevisiae*. *Cell Rep.* 7: 1858–1866. <https://doi.org/10.1016/j.celrep.2014.05.023>
- 853 Smith J. D., S. Suresh, U. Schlecht, M. Wu, O. Wagih, *et al.*, 2016 Quantitative CRISPR
854 interference screens in yeast identify chemical-genetic interactions and new rules for guide
855 RNA design. *Genome Biol.* 17: 45. <https://doi.org/10.1186/s13059-016-0900-9>
- 856 Sneath P. H., 1966 Relations between chemical structure and biological activity in peptides. *J.*
857 *Theor. Biol.* 12: 157–95.
- 858 Sugimoto N., S. Nakano, M. Katoh, A. Matsumura, H. Nakamuta, *et al.*, 1995 Thermodynamic
859 parameters to predict stability of RNA/DNA hybrid duplexes. *Biochemistry* 34: 11211–6.
- 860 Tarassov K., V. Messier, C. R. Landry, S. Radinovic, M. M. Serna Molina, *et al.*, 2008 An in vivo
861 map of the yeast protein interactome. *Science* 320: 1465–70.
862 <https://doi.org/10.1126/science.1153878>
- 863 Wagih O., M. Galardini, B. P. Busby, D. Memon, A. Typas, *et al.*, 2018 A resource of variant
864 effect predictions of single nucleotide variants in model organisms. *Mol. Syst. Biol.* 14:
865 e8430. <https://doi.org/10.15252/MSB.20188430>
- 866 Winzeler E. A., D. D. Shoemaker, A. Astromoff, H. Liang, K. Anderson, *et al.*, 1999 Functional
867 characterization of the *S. cerevisiae* genome by gene deletion and parallel analysis.
868 *Science* (80-.). 285: 901–906. <https://doi.org/10.1126/science.285.5429.901>
- 869 Wong N., W. Liu, and X. Wang, 2015 WU-CRISPR: characteristics of functional guide RNAs for
870 the CRISPR/Cas9 system. *Genome Biol.* 16: 218. <https://doi.org/10.1186/S13059-015-0784-0>
871
- 872 Wu A. C. K., H. Patel, M. Chia, F. Moretto, D. Frith, *et al.*, 2018 Repression of Divergent
873 Noncoding Transcription by a Sequence-Specific Transcription Factor. *Mol. Cell* 72: 942-

- 874 954.e7. <https://doi.org/10.1016/J.MOLCEL.2018.10.018>
- 875 Yachie N., E. Petsalaki, J. C. Mellor, J. Weile, Y. Jacob, *et al.*, 2016 Pooled-matrix protein
876 interaction screens using Barcode Fusion Genetics. *Mol. Syst. Biol.* 12: 863.
- 877 Zhao L., Q. Yang, J. Zheng, X. Zhu, X. Hao, *et al.*, 2016 A genome-wide imaging-based
878 screening to identify genes involved in synphilin-1 inclusion formation in *Saccharomyces*
879 *cerevisiae*. *Sci. Rep.* 6: 30134. <https://doi.org/10.1038/srep30134>
- 880 Zuo E., Y. Sun, W. Wei, T. Yuan, W. Ying, *et al.*, 2019 Cytosine base editor generates
881 substantial off-target single-nucleotide variants in mouse embryos. *Science* (80-.).
882 eaav9973. <https://doi.org/10.1126/SCIENCE.AAV9973>
- 883

884 SUPPLEMENTARY MATERIAL: Supplementary Figures 1-11

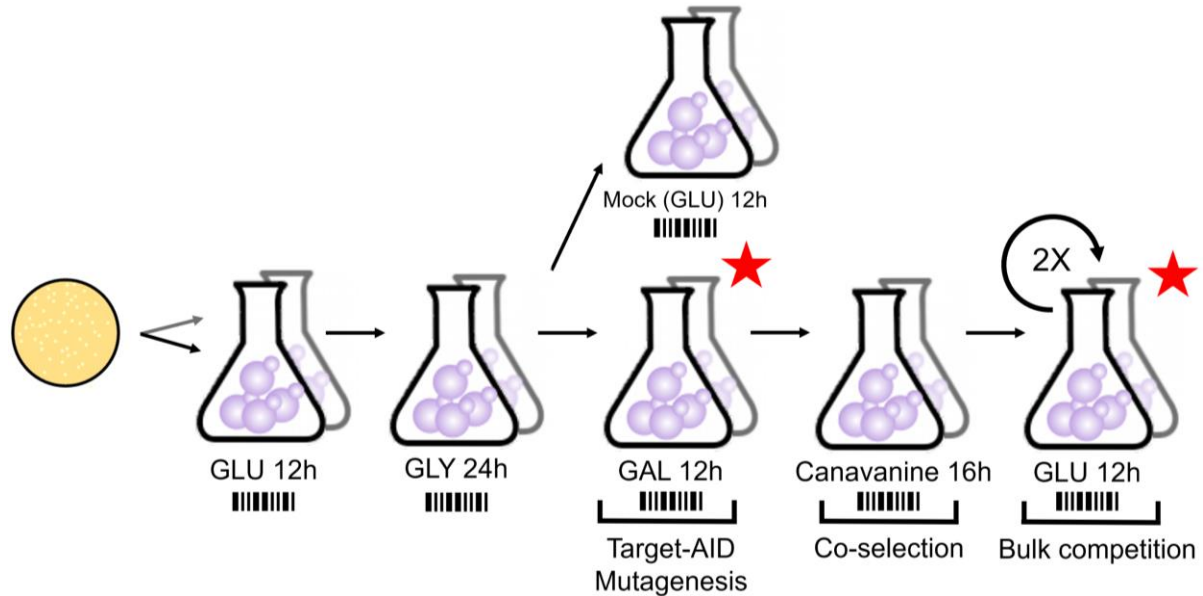
885



886

887 **Figure S1: A gRNA library for the systematic mutagenesis of yeast essential genes and**
888 **other targets of interest. A)** Number of genes targeted by the gRNA library for the different target
889 classes. **B)** Total number of gRNAs targeting genes in the different target classes. **C)** Distribution
890 of number gRNAs for each gene targeted in the different classes. **D)** Distribution of minimal (light
891 grey) and median (dark grey) pairwise sequence distance between all gRNA sequences in the
892 library.

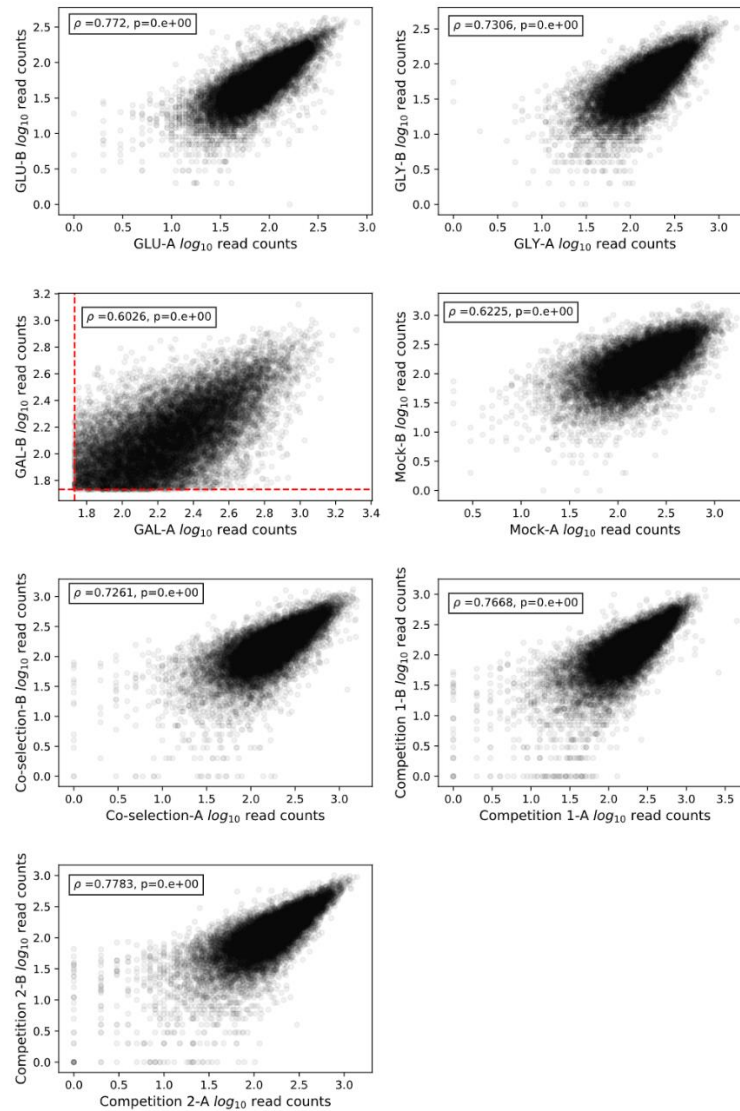
893



894

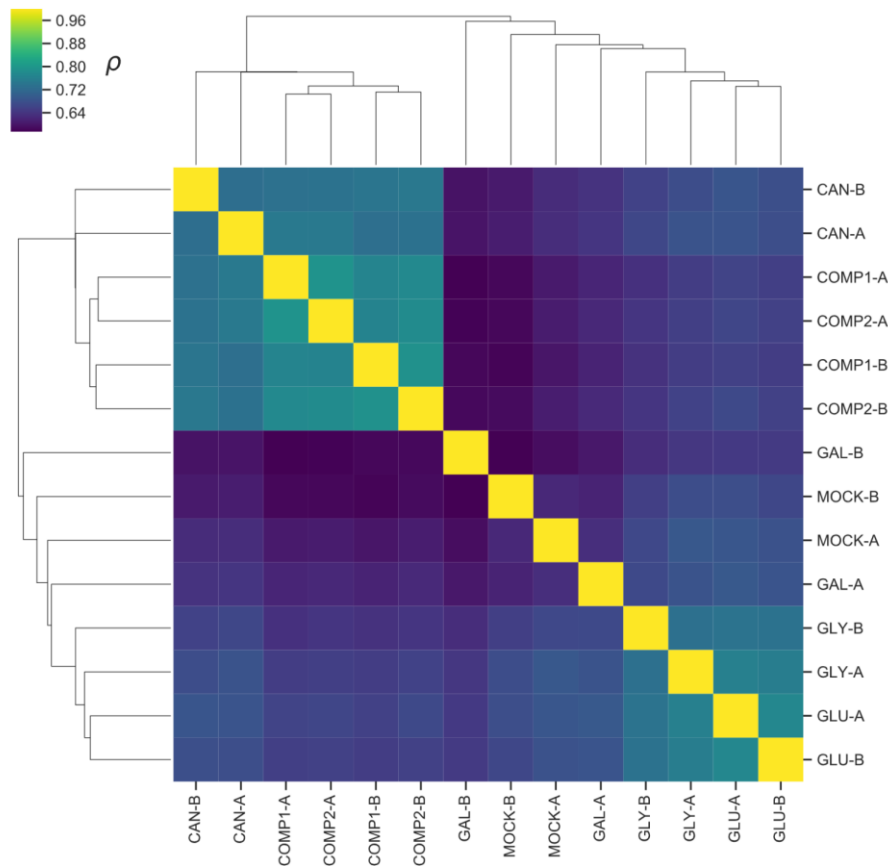
895 **Figure S2: Experimental workflow for Target-AID mutagenesis and co-selection.** The
896 mutagenesis method closely follows the base editing protocol previously described (Després *et*
897 *al.* 2018). After a pooled transformation step, cells were scraped and splitted into two replicates
898 for pre-cultures. After each step of the protocol, plasmid DNA was extracted from a cell sample
899 and used to amplify and sequence the gRNA pool. The red stars indicate time points used for
900 fitness effects analysis: read counts after galactose induction were used as T0 and were
901 compared with read counts after two rounds of competition. The mock induction steps mimics the
902 induction conditions but galactose in the media is replaced by glucose. This prevents the editing
903 enzyme from being expressed because glucose represses the GAL pathway. After canavanine
904 co-selection, cells go through two competition rounds in synthetic media where selective pressure
905 for the Target-AID bearing plasmid is lost. The entire experiment was completed within less than
906 25 generations after galactose induction, limiting the impact of compensatory and spontaneous
907 mutations.

908



909

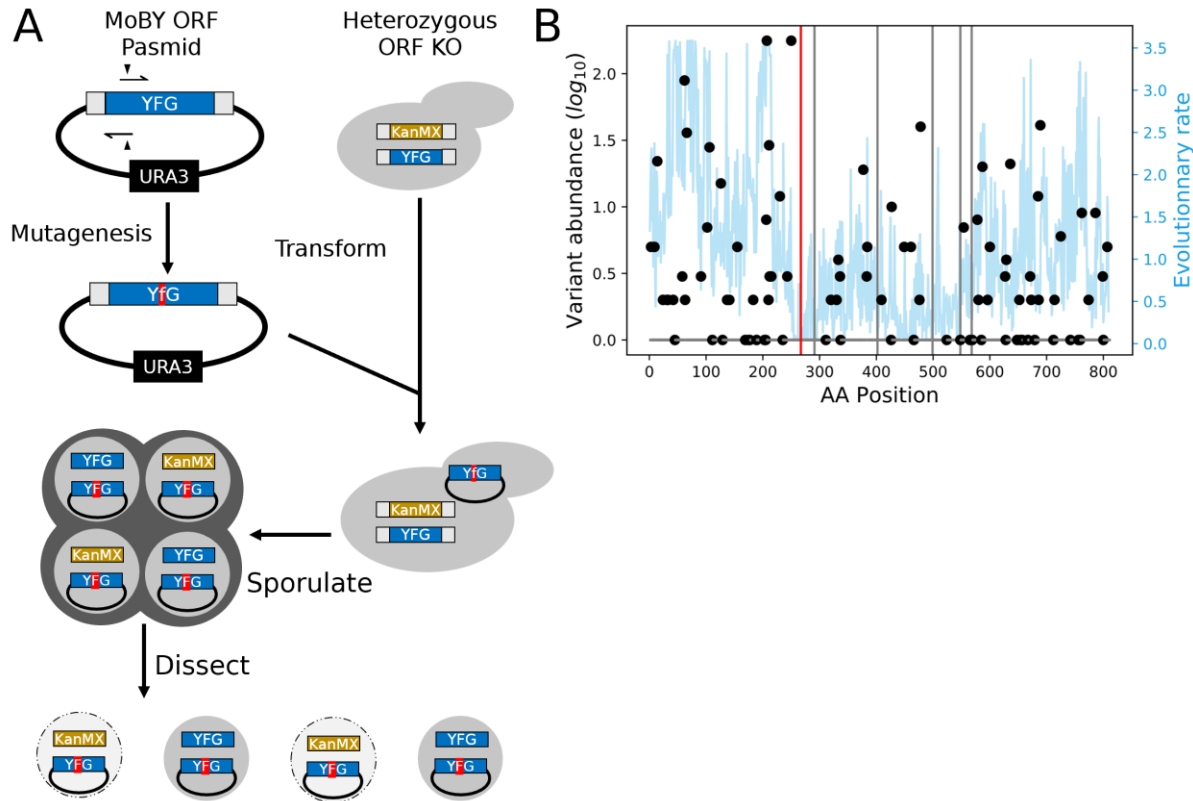
910 **Figure S3: Read abundance rank order is strongly correlated between replicates.** For
911 each time point, Spearman rank correlation of gRNA \log_{10} read abundance after basic filtering is
912 shown. The minimal read count after galactose induction, which served as the principal filtering
913 criteria, is shown on the galactose subpanel.



914

915 **Figure S4: Barcode abundance correlation clusters different experimental steps of the**
916 **screen.** Pairwise Spearman rank correlation of barcode counts was used to cluster the libraries
917 obtained at the different time points described in Figure S2. The lower level of correlation between
918 the galactose induction and mock induction timepoints compared to other associated steps could
919 reflect higher stochasticity in growth caused by cell to cell variation in the metabolic switch from
920 glycerol to sugars as the main carbon source as well as editing in the case of the galactose
921 timepoint.

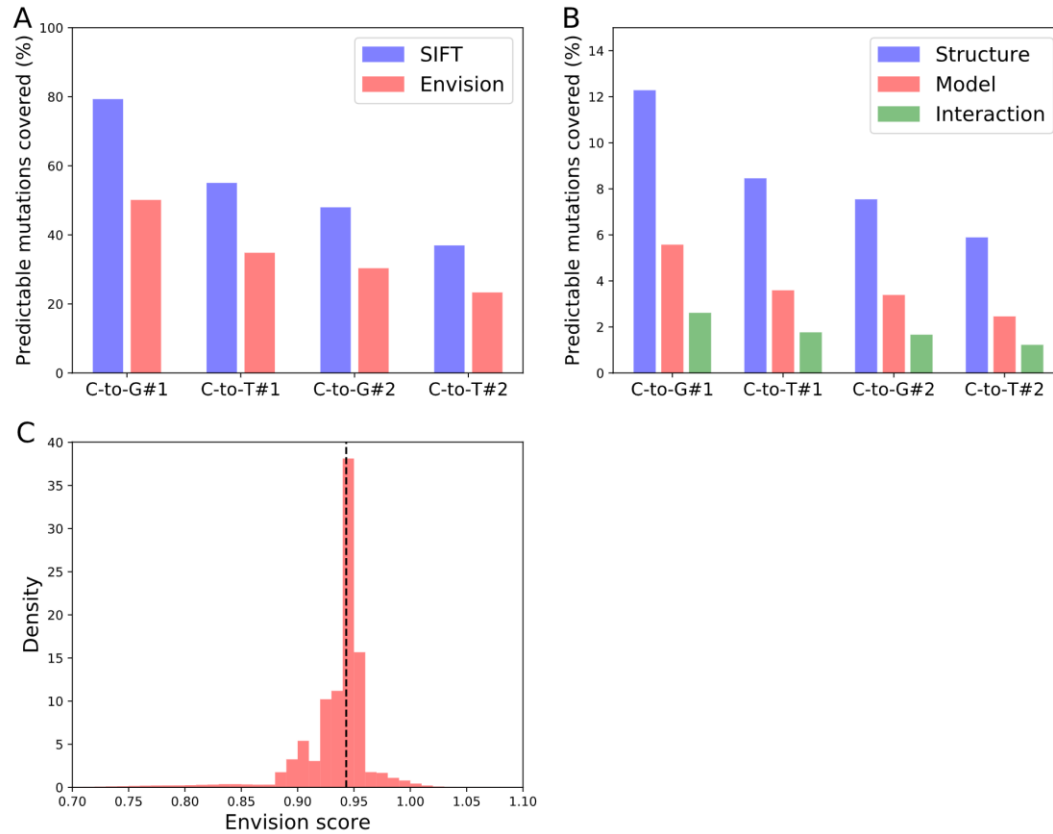
922



923

924 **Figure S5 Plasmid-based confirmation workflow by complementation test and**
 925 **evolutionary information on *GLN4*.** **A)** Detailed protocols for the different steps are presented
 926 in the methods. First, directed mutagenesis is used to introduce the mutation of interest (shown
 927 in red) in the MoBY collection plasmid of the targeted gene (YFG). This vector is then
 928 transformed into the heterozygous collection deletion strain (BY4743, *MATa* α *his3* Δ 1/*his3* Δ 1
 929 *leu2* Δ 0/*leu2* Δ 0 *LYS2*/*lys2* Δ 0 *met15* Δ 0/*MET15* *ura3* Δ 0/*ura3* Δ 0) of the gene of interest. The
 930 transformants are sporulated and their tetrads are dissected. If the mutated allele carried by the
 931 plasmid cannot complement the gene deletion, then only the two progenies bearing the wild-
 932 type copies will be viable. **B)** Protein variant frequency among 1000 yeast isolates (black dots)
 933 and residue evolutionary rate across species (blue line) for *GLN4*. The target site for the most
 934 deleterious GNE is highlighted by a red line and other GNE target sites are shown as grey lines.

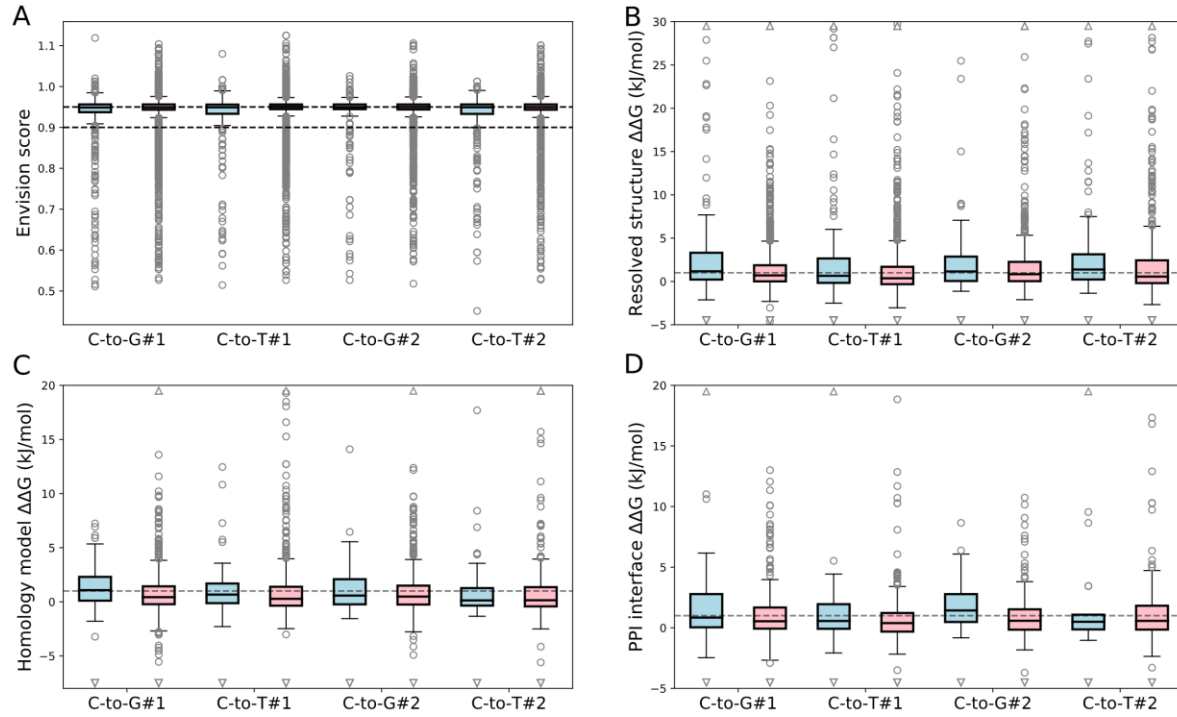
935



936

937 **Figure S6 gRNA predicted mutation coverage for Mutfunc and Envision data.** Mutfunc
938 integrates both the SIFT prediction scores and FoldX (Schymkowitz *et al.* 2005), $\Delta\Delta G$ predictions
939 for solved protein structures, homology models, and protein-protein interaction interfaces. gRNAs
940 which do not generate missense mutations were included in the calculations. **A)** Coverage for the
941 SIFT and Envision variant effect predictors for the four most probable single mutants created by
942 gRNAs detected in the experiment. **B)** Coverage for $\Delta\Delta G$ predictions for solved protein structures,
943 homology models, and protein-protein interaction interfaces for the four most probable single
944 mutants created by gRNAs detected in the experiment. **C)** Distribution of Envision scores across
945 all sites in the database for all proteins targeted by the set of gRNAs detected in the screen
946 ($n=7,556,573$). The median score is shown as a dotted black line.

947

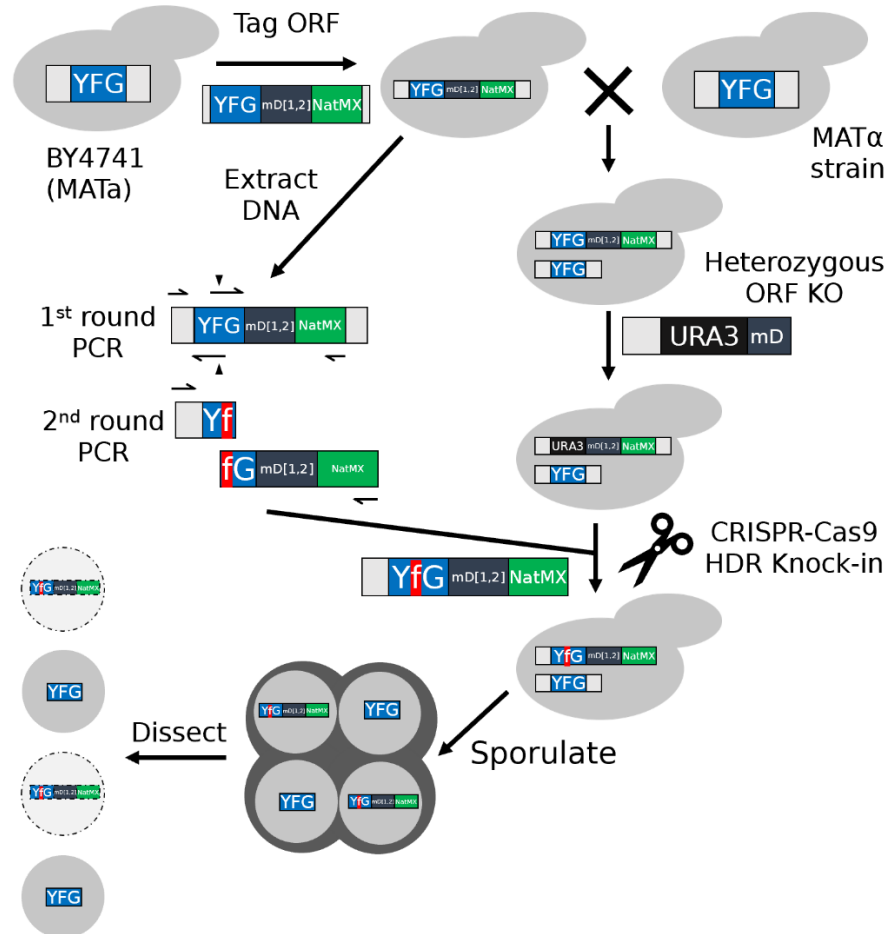


948

949 **Figure S7 GNE and non-significant gRNA effect prediction distributions.** **A)** Envision score
950 distributions for the four most probable mutations induced by GNEs (blue) and NSGs (red). Welch's
951 t-test p-values for comparisons: 5.00×10^{-6} , 0.002, 0.007, 7.75×10^{-5} . **B)** Predicted folding energy
952 variation ($\Delta\Delta G$) of GNE and NSG induced protein mutants compared to the wild-type structure
953 based on resolved protein structure. Welch's t-test p-values for comparisons: 0.0001, 0.006,
954 0.148, 0.007. **C)** Predicted folding energy variation ($\Delta\Delta G$) of GNE and NSG induced protein
955 mutants compared to the wild-type structure based on homology models of protein structure.
956 Welch's t-test p-values for comparisons: 0.016, 0.441, 0.195, 0.689. **D)** Binding energy variation
957 ($\Delta\Delta G$) of GNE and NSG induced mutant protein-protein interfaces compared to the wild-type
958 based on a resolved structure on the interface. Welch's t-test p-values for comparisons: 0.285,
959 0.303, 0.033, 0.95.

960

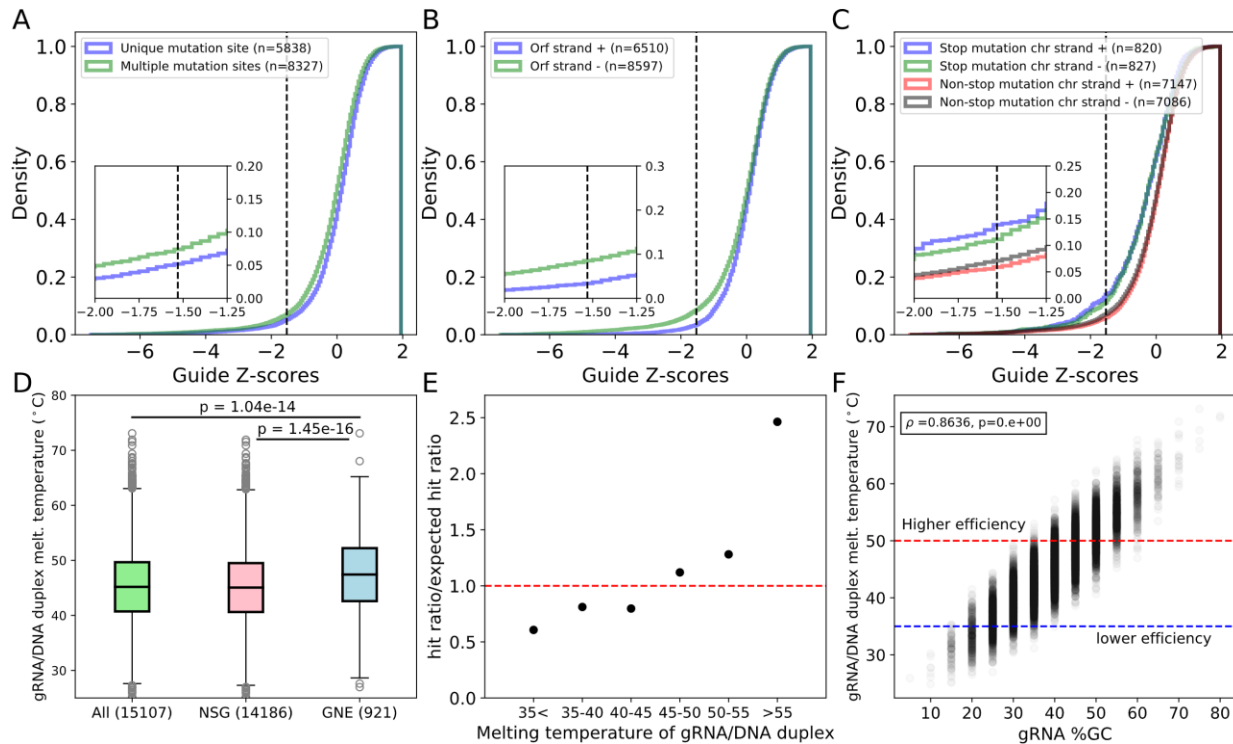
961



962

963 **Figure S8 Fitness affecting variant by CRISPR knock in confirmation workflow.** Detailed
 964 protocols for the different steps are presented in the methods. Starting from the wild-type
 965 laboratory strain BY4741, the gene of interest (*YFG*, blue) is first tagged with a modified
 966 version of the DHFR F[1,2] cassette (dark gray and green). The tagged strain is then crossed
 967 with a MAT α strain (Y8205) to create a heterozygous diploid. A *URA3* deletion cassette
 968 (black) that recombines with the *YFG* upstream sequence and the start of the mDHFR
 969 fragment is then used to generate a heterozygous KO strain. In parallel, genomic DNA is
 970 extracted from the tagged haploid strain. This DNA is then used as a template to amplify
 971 two fragments of *YFG* bearing the mutation of interest (shown in red) using a set of
 972 overhanging primers. The two fragments are then combined by fusion PCR to obtain the
 973 donor DNA used in the next step. Using a modified Cas9 vector (Ryan *et al.* 2016)
 974 that expresses a gRNA targeting the *URA3* cassette, the mutated allele is introduced
 975 at the KO locus to create a heterozygous mutant strain. The diploid cells can then be
 sporulated, and tetrad dissection allows observation of any phenotype linked with the
 mutation of interest.

976

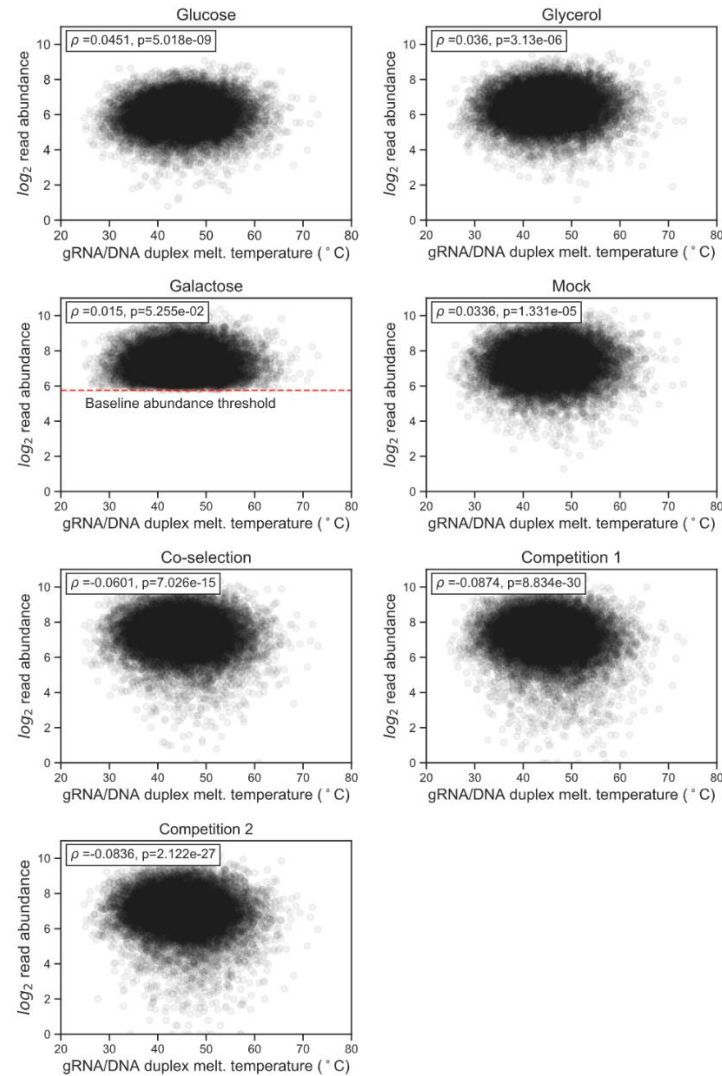


977

978 **Figure S9 Properties influencing stop codon GNEs are generalizable to non-stop codon**
 979 **generating GNEs. A)** Cumulative z-score density for gRNAs that do not generate stop codons
 980 depending on the number of mutable sites. A higher rate of GNE is observed for gRNAs which
 981 can lead to the editing of multiple nucleotides (Two-sample Kolmogorov-Smirnov test, $p=3.91 \times 10^{-29}$).
 982 The significance threshold is shown as a black dotted line. **B)** Cumulative z-score density for
 983 NSGs on orf target strand. gRNAs targeting the non-coding strand of the ORF have a higher
 984 likelihood of being GNEs (Two-sample Kolmogorov-Smirnov test, $p=1.87 \times 10^{-16}$). **C)** gRNA z-score
 985 cumulative density for both SGGs and non-SGGs grouped by the chromosomal strand they target.
 986 In SGGs, the target strand does not impact z-score distributions (Two-sample Kolmogorov-
 987 Smirnov test, $p=0.753$) and GNE proportions (Fisher's exact test, $p=0.149$). For non-SGGs, the
 988 chromosomal strand has a small influence on z-score distributions (Two-sample Kolmogorov-
 989 Smirnov test, $p=0.035$) and GNE proportions (Fisher's exact test, $p=0.002$). **D)** Distributions of
 990 modeled RNA/DNA duplex melting temperature for all non-SGGs generating gRNAs, the NSG
 991 subset, and the GNEs subset. P-values were calculated using the two-sample Kolmogorov-
 992 Smirnov test. **E)** Non-SGGs GNE enrichment compared to the expected GNE ratio for different
 993 melting temperature ranges. **F)** gRNA/DNA duplex melting temperature as a function of gRNA
 994 GC content for all gRNAs for which fitness effects were measured. The higher and lower efficiency
 995 thresholds are based on the enrichments shown in panel E and Figure 4F.

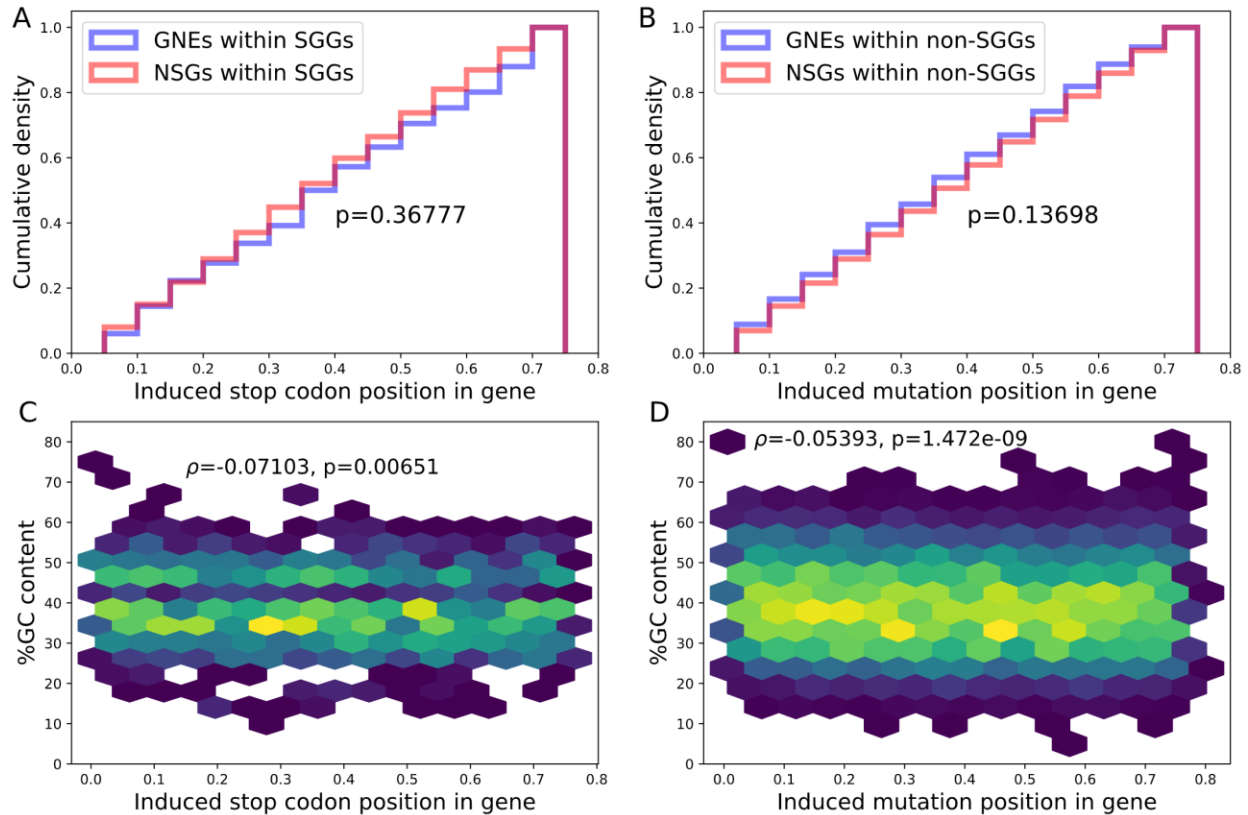
996

997



998

999 **Figure S10 gRNA/DNA duplex melting temperature is not linked to systematic sequencing**
1000 **biases.** Spearman rank correlation between replicate averaged read count and predicted
1001 gRNA/DNA duplex melting temperature is shown across timepoints. The minimal read count after
1002 galactose induction, which served as a filtering criterion, is shown on the galactose subpanels.
1003 gRNAs for which no reads were detected in one of the time points were included when computing
1004 the correlation but are not shown on the graphs because of log scaling.



1005

1006 **Figure S11 GNE density is independent of target nucleotide position bias. A)** In SGGs, GNE
1007 and NSG target sites that are evenly distributed across the target genes, and GNEs do not show
1008 any bias (Two-sample Kolmogorov-Smirnov). **B)** Non-SGG GNEs do not show any positional
1009 bias. **C)** A significant but small negative correlation is observed between gRNA target relative
1010 position and GC content of SGGs (Spearman's rank correlation). The very small observed effect
1011 coupled with the absence of position bias suggests that relative target position bias does not drive
1012 the link between GC content and gRNA efficiency. **D)** Similarly, a small but significant but small
1013 negative correlation is also observed between gRNA relative position and GC content for non-
1014 SGGs (Spearman's rank correlation).

1015