

1 **Supplementary Information**

2
3 **Predictable and precise template-free CRISPR editing of pathogenic variants**

4 Max W. Shen^{‡1,2}, Mandana Arbab^{‡3,4,5}, Jonathan Hsu^{6,7}, Daniel Worstell⁸, Olga
5 Krabbe^{8,9}, Christopher A. Cassa^{8,10}, David R. Liu^{3,4,5*}, David K. Gifford^{2,6,10,11*} and
6 Richard I. Sherwood^{8,9*}

7
8 ‡ *These authors contributed equally to this work.*

9 ¹ *Computational and Systems Biology Program, Massachusetts Institute of Technology,*
10 *Cambridge, Massachusetts, USA.*

11 ² *Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of*
12 *Technology, Cambridge, Massachusetts, USA.*

13 ³ *Merkin Institute of Transformative Technologies in Healthcare, Broad Institute of*
14 *Harvard and MIT, Cambridge, Massachusetts, USA.*

15 ⁴ *Department of Chemistry and Chemical Biology, Harvard University, Cambridge,*
16 *Massachusetts 02138, USA.*

17 ⁵ *Howard Hughes Medical Institute, Harvard University, Cambridge, Massachusetts*
18 *02138, USA.*

19 ⁶ *Department of Biological Engineering, Massachusetts Institute of Technology,*
20 *Cambridge, Massachusetts, USA.*

21 ⁷ *Molecular Pathology Unit, Center for Cancer Research, and Center for Computational*
22 *and Integrative Biology, Massachusetts General Hospital, Charlestown, Massachusetts,*
23 *USA.*

24 ⁸ *Division of Genetics, Department of Medicine, Brigham and Women's Hospital and*
25 *Harvard Medical School, Cambridge, Massachusetts, USA.*

26 ⁹ *Hubrecht Institute, Utrecht, the Netherlands.*

27 ¹⁰ *Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA.*

28 ¹¹ *Department of Electrical Engineering and Computer Science, Massachusetts Institute*
29 *of Technology, Cambridge, Massachusetts, USA.*

30
31 * Correspondence should be addressed to R.I.S. (rsherwood@rics.bwh.harvard.edu) or
32 D.K.G. (gifford@mit.edu) or D.R.L. (drliu@fas.harvard.edu).
33
34

Table of Contents

35		
36		
37	Supplementary Discussion.....	3
38		
39	Supplementary Methods.....	12
40		
41	- Plasmid and Insert Sequences.....	33
42		
43	Supplementary References.....	47
44		
45		
46		
47		

48 **Supplementary Discussion**

49

50 **Cellular repair of double-stranded DNA breaks and inDelphi**

51 DNA double-strand breaks are detrimental to genomic stability, and as such the detection
52 and faithful repair of genomic lesions is crucial to cellular integrity. A large number of
53 genes have evolved to respond to and repair DNA double-strand breaks, and these genes
54 can be broadly grouped into a set of DNA repair pathways¹, each of which differs in the
55 biochemical steps it takes to repair DNA double-strand breaks. Accordingly, these
56 pathways tend to produce characteristically distinguishable non-wildtype genotypic
57 outcomes.

58

59 The goal of our machine learning algorithm, inDelphi, is to accurately predict the identities
60 and relative frequencies of non-wildtype genotypic outcomes produced following a
61 CRISPR/Cas9-mediated DNA double-strand break. To accomplish this goal, we
62 developed parameters to classify three distinct categories of genotypic outcomes,
63 microhomology deletions, microhomology-less deletions, and insertions, informed by the
64 biochemical mechanisms underlying the DNA repair pathways that typically give rise to
65 them.

66

67 Double strand breaks are thought to be repaired via four major pathways: classical non-
68 homologous end-joining (c-NHEJ), alternative-NHEJ (alt-NHEJ), microhomology-
69 mediated end-joining (MMEJ), and homology-directed repair (HDR)¹. To create inDelphi,
70 we developed three machine learning modules to model genotypic outcomes assuming
71 characteristic of the c-NHEJ, microhomology mediated alt-NHEJ, and MMEJ pathways.
72 While template-free CRISPR/Cas9 DNA double-strand break may lead to HDR repair via
73 endogenous homology templates that exist in *trans*², we do not explicitly model HDR-
74 characteristic outcomes using our algorithm.

75

76 Before proceeding, it is important to note that while specific DNA repair pathways are
77 characteristically associated with distinct genotypic outcomes, the proteins involved in the
78 various pathways and the resulting repair products may at times overlap. This fact has
79 several implications. First, we cannot make conclusive statements about the role of
80 specific proteins or pathways in specific genotypic outcomes without perturbation
81 experiments (e.g. our comparison of wildtype and *Prkdc*^{-/-}*Lig4*^{-/-} mESCs can illuminate the
82 roles of these proteins, specifically). Second, because assigning genotypic outcomes to
83 biochemical mechanisms is likely imperfect, we use machine learning methods to identify
84 trends and patterns in genotype frequencies that refine this crude binning process.

85

86 In the first step of the inDelphi method, we separate genotypic outcomes into three
87 classes: microhomology deletions (MH deletions), microhomology-less deletions (MH-
88 less deletions), and single-base insertions (1-bp insertions) (Figure 1e). Below we outline
89 the algorithmic definitions of each genotypic outcome class, the pathways associated with
90 each class, and the DNA sequence parameters included in inDelphi training of each class.
91 For more detailed technical algorithmic definitions of the genotypic outcome classes, see
92 Supplementary Methods.

93

94 **MH deletions are predicted from MH length, MH GC content, and deletion length**

95 The majority of Cas9-mediated double-strand break repair genotypes we observe in our
96 datasets are what we classify as MH deletions (53-58% in mESC, K562, HCT116, and
97 HEK293). We hypothesize that these deletions occur through MMEJ-like processes and
98 use known features of this pathways to inform a machine learning module to predict MH
99 deletion outcomes. Following 5'-end resection as occurs in MMEJ, alt-NHEJ, and HDR¹,
100 microhomologous basepairing of single-stranded DNA (ssDNA) sequences occurs
101 across the border of the double strand breakpoint^{3,4}. To restore a contiguous double-
102 strand DNA chain, the 5'-overhangs not participating in the microhomology are removed
103 up until the paired microhomology region, and the unpaired ssDNA sequences are
104 extended by DNA polymerase using the opposing strand as a template (Figure 1d,
105 Extended Data Fig. 2).

106
107 Assuming these same processes, inDelphi calculates the set of all MH deletions available
108 given a specific sequence context and cleavage site.

109
110 As an example workflow, given the following sequence and its cleavage site:

111
112 ACGTG | CATGA
113 TGCAC | GTACT

114
115 for every possible deletion length from 1-bp to 60-bp deletions, we overlap the 3'-
116 overhang downstream of the cut site under the upstream 3'-overhang and determine if
117 there is any microhomologous basepairing. As an example, given the 4-bp deletion
118 length:

119
120 ACGTG
121 | | |
122 GTACT

123
124 we see that there are three microhomologous basepairing events.

125
126 We then choose a particular microhomology (here, the highlighted C:G):

127
128 **AC**GTG
129 | | |
130 **G**TACT

131
132 then generate its unique repair genotype by following left-to-right along the top strand and
133 jumping down to the complement of the bottom strand to simulate DNA polymerase fill-
134 in.

135
136 Here, this yields:

137
138 ACATGA
139 TGTACT

140

141 This can also be displayed as an alignment. We note that by “jumping down” after the first
142 base in the top strand, we can also describe this outcome using the delta-position 1. (See
143 section on delta-positions). A deletion at delta-position 0 yields the same genotype.

144
145 Deletion a: AC----ATGA
146 Wt: ACGTGCATGA

147
148 This same sequence context and cleavage site could produce a distinct 4-bp MH deletion
149 genotypic outcome through use of the TG:AC microhomology. This single outcome can
150 be described as using delta-positions 2, 3, or 4. inDelphi uses only the single maximum
151 delta-position (here, 4) to described a unique MH deletion.

152
153 Deletion b: ACGTG----A
154 Wt: ACGTGCATGA

155
156 Thus, there may be multiple MH deletion outcome genotypes for a given deletion length,
157 and there is always a 1:1 mapping between the microhomologous basepairing used in
158 that MH deletion and the resultant genotypic outcome. The set of MH deletions thus
159 includes all 1-bp to 60-bp deletions that can be derived from the steps above that simulate
160 the MMEJ mechanism.

161
162 MMEJ efficiency has been reported to depend on the thermodynamic favorability and
163 stability of a candidate microhomology^{3,4}. To parameterize MH deletions using the
164 biochemical sequence features that influence this form of DNA repair, inDelphi calculates
165 the MH length, MH GC content, and resulting deletion length for each possible MH
166 deletion. These features are input into a machine learning module referred to in the
167 Supplementary Methods as the microhomology neural network (MH-NN) to learn the
168 relationship between these features and the frequency of an MH deletion outcome in a
169 training CRISPR/Cas9 genotypic outcome dataset. While we predict and empirically find
170 that favored MH deletions have long MH lengths relative to total deletion length and high
171 MH GC-contents, we do not provide any explicit direction or comparative weighting to
172 these parameters at the outset. inDelphi then outputs a phi-score for any MH deletion
173 genotype (whether it was in the training data or not) that represents the favorability of that
174 outcome as predicted by MH-NN.

175
176 It is important to emphasize that the phi-score of a particular MH deletion does not itself
177 represent the likelihood of that MH deletion occurring in the context of all MH deletions at
178 a given site. Some CRISPR/Cas9 target sites may have many possible favorable MH
179 deletion outcomes while other sites have few, and thus phi-score must be normalized for
180 a given target site to generate the fractional likelihood of that genotypic outcome at that
181 site. Total unnormalized MH deletion phi-score is one factor that is further used to predict
182 the relative frequency of the different repair classes: MH deletions, MH-less deletions,
183 and insertions.

184
185

186 **MH-less deletions are predicted from their length**

187 We define MH-less deletions as all possible deletions that have not been accounted for
188 by the workflow described above for MH deletions. Mechanistically, our data analysis
189 suggests that MH deletions are associated with repair genotypes produced by c-NHEJ
190 and microhomology-mediated alt-NHEJ pathways.

191
192 Following a double-strand break, c-NHEJ-associated proteins rapidly bind the DNA
193 strands flanking the double-strand DNA breakpoint and recruit ligases, exonucleases, and
194 polymerases to process and re-anneal the breakpoint in the absence of 5'-end resection
195 (Extended Data Fig. 2)^{1,5}. Commonly, c-NHEJ repair is error-free; however, in the context
196 of Cas9-mediated cutting, faithful repair leads to repeated cutting, thereby increasing the
197 eventual likelihood of mutagenic repair. Erroneous c-NHEJ repair products are mainly
198 thought to consist of small insertions or deletions or combinations thereof that most
199 frequently occur in the direct vicinity of the DNA break point⁵⁻⁷. The resulting deletions,
200 which we refer to as medial end-joining MH-less deletions, have often lost bases both
201 upstream and downstream of the cleavage site.

202
203 Microhomology-mediated alt-NHEJ is a distinct pathway that produces MH-less deletion
204 products. In contrast to c-NHEJ, which is microhomology independent, this form of alt-
205 NHEJ repair occurs following 5'-end resection and is mediated by microhomology in the
206 sequence surrounding the double-strand break-point¹. Microhomologous basepairing
207 stabilizes the 3'-ssDNA overhangs following 5'-end resection, similarly to in MMEJ,
208 allowing DNA ligases to join the break across one of the strands of this temporarily
209 configured complex. The opposing un-annealed flap is then removed, and newly
210 synthesized DNA templated off of the remaining strand is annealed to repair the lesion
211 (Extended Data Fig. 2).

212
213 While alt-NHEJ uses microhomology, the repair products it produces do not follow the
214 predictable genotypic patterns induced by MMEJ and are thus grouped into MH-less
215 deletion genotypes. MH deletions are a direct merger of both annealed strands, in which
216 the outcome genotype switches from top to bottom strand at the exact end-point of a
217 microhomology. In contrast, while alt-NHEJ employs microhomology in its repair
218 mechanism, the deletion outcomes it generates comprise bases exclusively derived from
219 either the top or bottom strand. Mechanistically, this occurs because ligation of a 3'-
220 overhang to its downstream ligation partner results in removal of the entire opposing
221 ssDNA overhang up until the point of ligation. This process prevents any deletion from
222 occurring in the 3'-overhang strand that is first attached to the DNA backbone, while
223 inducing loss of an indeterminate length of sequence on the opposing strand. The
224 resulting deletion genotypes, which we refer to as unilateral end-joining MH-less
225 deletions, do not retain information on the exact microhomology causal to their
226 occurrence, and are thus also referred to as MH-less.

227
228 Consequently, the various mechanisms that give rise to MH-less deletions are capable of
229 generating a vast number of genotypic outcomes for any given deletion length. Having
230 less information on the biochemical mechanisms that impact the relative frequency of

231 NHEJ deletion products, inDelphi models these deletions without assuming any particular
232 mechanism.

233
234 inDelphi detects MH-less deletions from training data as the set of all deletions that are
235 not MH deletions and parameterizes them solely by the length of the resulting deletion.
236 This is based on the simple assumption that c-NHEJ and alt-NHEJ processes are most
237 likely to produce short deletions, supported by our empirical observation. As with MH
238 deletions, this assumption is not explicitly coded into the inDelphi MH-less deletion
239 prediction module, instead allowing it to be “learned” by a neural network called MHless-
240 NN.

241
242 MHless-NN optimizes a phi-score for a given MH-less deletion length, grounded in the
243 frequency of MH-less deletion outcomes of that length observed in the training data. We
244 observe that MHless-NN learns a near-exponential decaying phi-score for increasing
245 deletion length, that reflects the sum total frequency of all MH-less deletion genotypes.
246 The total unnormalized MH-less deletion phi-score for a given target and cut site is also
247 employed to inform the relative frequency of different repair classes.

248 249 **1-bp insertions are predicted from sequence context and deletion phi-scores**

250 Lastly, inDelphi predicts 1-bp insertions from both the broader sequence context and the
251 immediate vicinity of the cleavage site. We empirically find that 1-bp insertions are far
252 more common than longer insertions, so we focus on their prediction. It is classically
253 assumed that short sequence insertions are the result of c-NHEJ^{6,7}, however, little else is
254 known about their biochemical mechanism as it pertains to local sequence context to help
255 inform prediction. Nonetheless, we find powerful correlations between the identities of the
256 bases surrounding the Cas9 cleavage site and the frequency and identity of the inserted
257 base (see main text). Motivated by these empirical observations, inDelphi is fed with
258 training data on 1-bp insertion frequencies and identities at each training site
259 parameterized with the identities of the -3, -4, and -5 bases upstream of the NGG PAM-
260 sequence (when the training set is sufficiently large, and the -4 base alone when training
261 data is limited) as features. Also added as features are the precision score of the deletion
262 length distribution and the total deletion phi-score at that site. These features are
263 combined into a *k*-nearest neighbor algorithm that predicts the relative frequencies and
264 identities of 1-bp insertion products at any target site.

265 266 **The combination of the MH, MH-less, and insertion model predict genotype** 267 **fractions**

268 Altogether, informed by known paradigms of DNA repair, we build 2 neural networks and
269 a *k*-nearest neighbor model to predict genotypic outcomes following Cas9 cutting. These
270 models compete and collaborate in inDelphi to generate predictions of the relative
271 frequencies of these products. This competition within inDelphi among repair types
272 reflects empirical evidence from Lib-A and Lib-B that sequence contexts do influence
273 classes of repair outcomes. Sequence contexts with high phi scores (high
274 microhomology) have higher efficiencies of MH deletions among all editing outcomes
275 (Figure 2d, Extended Data Fig. 3), and sequence contexts with low phi scores (low
276 microhomology) have higher efficiencies of 1-bp insertions among all editing outcomes

277 (Figure 2d, Extended Data Fig. 3). While it is tempting to generalize that the competition
278 and collaboration among outcome classes modeled by inDelphi reflects interactions
279 among components of distinct DNA repair pathways, the classes of outcomes considered
280 by inDelphi do not necessarily arise from distinct DNA repair pathways as they are
281 described above. inDelphi is trained on the repair outcomes only and cannot distinguish
282 between the nature of genotypes when they may occur through MH-mediated and MH-
283 less mechanisms, and it is imaginable that some repair products result through more than
284 one repair pathway.

285
286 As an additional note, while NHEJ is generally assumed to dominate double-strand break
287 repair from environmentally induced damage⁵, we find in the context of Cas9 cutting that
288 MH deletion genotypes are more common than MH-less deletions and insertions. It is
289 possible that error-free c-NHEJ is occurring frequently in response to Cas9 cutting but
290 that its perfect repair allows for recurring Cas9 cutting that goes undetected by our
291 workflow, thus skewing the observed relative frequency profile of mutagenic outcomes
292 toward MMEJ-type repair.

293 294 **Rarer CRISPR-Cas9 outcomes**

295 Our library assay and workflow involved data processing of high-throughput sequencing
296 data using sequence alignments and a designed procedure for categorizing sequence
297 alignments into categories of CRISPR-related outcomes. Beyond simple deletions and
298 insertions, we identified other rarer outcomes that were explained as indels caused by
299 CRISPR, such as combination insertion/deletions involving and/or near the cleavage site
300 (0.5-2% of all products) and indels near but not immediately at the cleavage site (3-5% of
301 all products), which occurred more often on the PAM-distal side of the double-strand
302 break (data not shown). Our library assay is unable to observe events that occur outside
303 of our high-throughput sequencing window.

304
305 Default sequence alignment procedures can generate sequence alignments involving
306 simple CRISPR-caused deletions and insertions that do not occur immediately at the
307 cleavage site, but that can be transformed into an equal-scoring sequence alignment
308 where the indel does occur immediately at the cleavage site. This straightforward
309 processing step is not performed by the most common bioinformatic tools for sequence
310 alignment, since they were not expressly designed for CRISPR. We note here that our
311 sequence alignment procedure takes this into account (see Supplementary Methods for
312 more detailed description). This attention to detail enables us to accurately identify simple
313 indels that occur near but not immediately at the cleavage site. We observe that the
314 frequency of these indels across target sites correlates significantly with the total on-target
315 editing efficiency (measured by the frequency of non-wild-type outcomes out of all non-
316 noise outcomes) at these target sites in HEK293 and mES cells. We also observe
317 significantly higher frequencies in postCas9 treatment conditions than preCas9 control
318 conditions. Together, these observations suggest that these indels are caused by
319 CRISPR editing.

320
321 ***Prkdc*^{-/-}*Lig4*^{-/-} mutants have distinct and predictable DNA repair product**
322 **distributions**

323 While it is generally true that our work cannot establish roles for specific DNA repair
324 pathways in specific types of Cas9-mediated outcomes, we have performed an
325 experiment in which we measure Cas9-mediated genotypic outcomes from mESCs that
326 are lacking *Prkdc* and *Lig4*, two proteins known to be key in c-NHEJ⁵. We find an increase
327 in relative frequency of MH deletions as compared to MH-less deletions in *Prkdc*^{-/-}*Lig4*
328 ^{-/-} mESCs as compared to wild-type mESCs (see main text), which is suggestive of an
329 increase in MMEJ outcomes at the expense of NHEJ outcomes.

330
331 Intriguingly, we also find that *Prkdc*^{-/-}*Lig4*^{-/-} mESCs are impaired in unilateral deletions,
332 where only bases from one side of the cutsite are removed, but not medial MH-less
333 deletion outcomes that have loss of bases on both sides of the breakpoint. (Extended
334 Data Fig. 6). As discussed earlier, microhomology-mediated alt-NHEJ, which we
335 hypothesize may give rise to unilateral MH-less deletions, proceeds through a mechanism
336 in which DNA repair intermediates that mimic MMEJ-mediated repair are formed initially
337 (Extended Data Fig. 2), as microhomology base-pairing temporarily stabilizes 3'-
338 overhangs following 5'-end resection. Subsequently, ligation joins one 3' overhang with
339 the sequence on the other side of the DNA double-strand break, giving rise to a unilateral
340 deletion. If the unilateral joining products we observe in our experiments indeed arise
341 through similar mechanisms as those described by this form of alt-NHEJ, it is conceivable
342 that the MMEJ pathway may overtake 3'-end ligation at this microhomology-containing
343 intermediate step when ligation is impaired through loss of *Lig4*. Thus, cross-talk of
344 microhomology-mediated repair pathways could account for loss of unilateral end-joining
345 MH-less outcomes and concomitant increase in MH deletion outcomes. Medial joining
346 outcomes are not hypothesized to originate from intermediates that overlap with
347 microhomology-mediated deletion products (Extended Data Fig. 2). Therefore, the repair
348 genotypes generated via this orthogonal pathway may be afforded more time to be
349 completed by ligases other than *Lig4*, thus explaining why these outcomes appear
350 unaffected by NHEJ impairment.

351
352 While DNA repair products in *Prkdc*^{-/-}*Lig4*^{-/-} mESCs differ substantially from those in wild-
353 type cells, we find that these DNA repair products are also highly predictable. In particular,
354 inDelphi performed well on held-out *Prkdc*^{-/-}*Lig4*^{-/-} data when trained on *Prkdc*^{-/-}*Lig4*^{-/-} data
355 (indel genotype prediction median Pearson correlation = 0.84, indel length frequency
356 prediction Pearson correlation = 0.80), showing that our modeling approach is robustly
357 capable of learning accurate predictions for Cas9 editing data in not just wild-type
358 experimental settings but also settings with significant biochemical perturbation. As such,
359 we suggest here that inDelphi's modeling approach can be useful on additional tasks
360 unexplored here provided that inDelphi is supplied with appropriate training data.

361
362 **NU7041, DPKi3, and MLN4924 induce a distinct DNA repair product distribution**

363 We further investigated the role of DNA repair pathways by three separate experiments
364 involving HTS characterization of Lib-B in mESCs treated with three separate small
365 molecules: NU7041, a DNA dependent protein kinase (DNA-PK) inhibitor; DPKi3, another
366 DNA-PK inhibitor, and MLN4924, a NEDD8-activating enzyme (NAE) inhibitor. DNA-PK
367 and NAE are proteins involved in c-NHEJ^{5,8}.

368

369 MLN4924 is thought to inhibit the release of the Ku70/Ku80 heterodimer following proper
370 c-NHEJ repair, potentially disrupting downstream processes such as transcription and
371 replication, which may lead to decreased cell survival and a depletion of Ku70/Ku80-
372 dependent DNA repair genotypes in a population.⁸

373
374 DNA-PK is commonly recruited to DSBs during c-NHEJ and is known to phosphorylate *in*
375 *vitro* many c-NHEJ-related factors including Ku70/80, XRCC4, DNA Ligase IV, Artemis,
376 H2AX, p53, and itself. Inhibition of DNA-PK leads to DNA repair defects⁵. The catalytic
377 subunit of DNA-PK is encoded by the *Prkdc* gene, which was knocked out in *Prkdc*^{-/-}
378 *Lig4*^{-/-} cells.

379
380 From HTS data, we observed that the frequency of MH deletions among all deletions
381 clustered into three approximate groups: wild-type (median 77%) and MLN4924, then
382 DPKi3 and NU7041 (median 81%), and lastly *Prkdc*^{-/-}*Lig4*^{-/-} (median 90%) (Extended Data
383 Fig. 6). These data suggest that impairing DNA-PK (via DPKi3, NU7041 and *Prkdc*^{-/-})
384 yields a moderate 17% reduction in the frequency of MH-less deletions (23% to 19%).
385 This reduction appears to be non-redundant with knockout of *Lig4* evidenced in *Prkdc*^{-/-}
386 *Lig4*^{-/-} cells with a 57% reduction (23% to 10%) in MH-less frequency. Lastly, impairing
387 NAE did not have a significant impact on the frequency of MH-less deletions.

388
389 We observed an overall increased frequency of repair to wild-type at pathogenic
390 microduplication alleles after treatment with DPKi3, MLN4924, and NU7041 (Extended
391 Data Fig. 6). Along with *Prkdc*^{-/-}*Lig4*^{-/-} cells, the change in repair efficiency was associated
392 with deletion length ($p < 2.2 \times 10^{-3}$), with decreased efficiency compared to wild-type at
393 short deletion lengths and increased efficiency at longer deletion lengths.

394
395 The change in repair efficiency caused by separate treatments of DPKi3, MLN4924, and
396 NU7041 was highly consistent across different target sites ($r = 0.73, 0.77, \text{ and } 0.81$,
397 Extended Data Fig. 6). This is surprising since MLN4924 inhibits a different target than
398 DPKi3 and NU7041. We observed a similar but weaker relationship between the three
399 small molecules and *Prkdc*^{-/-}*Lig4*^{-/-}, with Pearson correlations of 0.09, 0.16, and 0.18.
400 Taken together, these observations suggest a relationship between DNA sequence and
401 the propensity of DNA repair outcomes through c-NHEJ.

402
403 In DPKi3, MLN4924, and NU7041 treated cells, the decrease in MH-less deletions
404 primarily occurs medial joining products (Extended Data Fig. 6), suggesting that DNA-PK
405 is a strong contributor to medial joining products. However, when both DNA-PK and *Lig4*
406 are knocked out in *Prkdc*^{-/-}*Lig4*^{-/-} cells, the average frequency of medial joining products
407 is not significantly changed, and instead the primary decrease occurs in unilateral joining
408 products.

409
410 Interestingly, MLN4924 increases the average frequency of unilateral joining events.
411 Combined with its effect of decreasing medial joining products, the overall net effect of
412 MMLN4924 is an absence of significant change to the frequency of MH-less deletions.

413

414 The frequency distribution of medial joining products in *Prkdc^{-/-}Lig4^{-/-}* reveals a decrease
415 in median frequency in combination with an inflation in high frequency outliers (target sites
416 where >80% of all deletion products are MHless medial products) which skews the
417 distribution's average to be above the median. Taken together, these data confirm that
418 both medial and unilateral products are both generally depleted in *Prkdc^{-/-}Lig4^{-/-}* cells, and
419 suggest that knocking out DNA-PK depletes medial MHless products while knocking out
420 Lig4 depletes unilateral MHless products.
421

422 **Supplementary Methods**

423

424 **Library cloning protocol**

425

426 **Synthesized oligo library sequence**

427 GATGGGTGCGACGCGTCAT [55bpTarget]AGATCGGAAGAGCACACGTCTG**AATATT**GTGGA
428 AAGGACGAAACACCG [19/20-nt PROTOSPACER depending on whether it
429 naturally starts with a G]GTTTAAGAGCTATGCTGGAAACAGC

430

431 Linker region / Oligo library amplification primer anneal region

432 Read 2 sequencing primer stub

433 **Sspl restriction site**

434 U6-promoter stub

435 sgRNA-hairpin stub

436

437

438 **1. Oligo library QPCR to determine number of amplification cycles for Oligo**
439 **Library PCR**

440 *Notes: Amplification of oligos with relatively low GC-content is less efficient than GC-rich*
441 *sequences. We found NEBNext polymerase to be the least biased in amplification of our*
442 *library. Increasing the elongation time to 1 min per cycle for all cloning and sequencing*
443 *library prep PCRs eliminates GC-skewing of library sequences and reduces the rate of*
444 *PCR-recombination.*

445

446 - Set up the following reaction:

447

0.4 ng	Synthesized Oligo Library
10 ul	NEBNext 2x Master Mix
0.5 ul	20uM OligoLib_Fw
0.5 ul	20uM OligoLib_Rv
0.2 ul	SybrGreen Dye (100x)
to 20 ul	H ₂ O

448

449 67°C annealing temperature

450

451 - Check 246bp amplicon size on 2.5% agarose gel.

452 - Determine the point that signal amplification has plateaued.

453

454 **2. Oligo Library PCR amplification**

455 - Set up the following reaction:

456

4 ng	Synthesized Oligo Library
50 ul	NEBNext 2x Master Mix
2.5 ul	20uM OligoLib_Fw
2.5 ul	20uM OligoLib_Rv
to 100 ul	H ₂ O

457
458 67°C annealing temperature, 1 minute extension time.
459 Cycle number is half the number of cycles needed to reach signal amplification plateau
460 in the QPCR in step 1, reduced by 1 cycle to scale for DNA input.

461
462 - PCR purify amplified sequence.

463
464 **3. Donor template amplification**

465 - Set up the following reaction:

466

5 ng	spCas9 sgRNA plasmid (71485)
50 ul	NEBNext 2x Master Mix
2.5 ul	20uM CircDonor_Fw
2.5 ul	20uM CircDonor_Rv
to 100 ul	H ₂ O

467

468 62°C annealing temperature

469 20 cycles

470

471 - Gel purify 167bp band from 2.5% agarose gel.

472

473 **4. Circular assembly and restriction digest linearization**

474 *Note: We use a molar ratio of donor template to amplified oligo library of 3:1. An increase*
475 *in amplified oligo library compounds cross-over within library members resulting in*
476 *mismatch of protospacer and target sequences.*

477

478 - Set up the following reaction:

479

429 ng	Donor template
239 ng	Amplified Oligo Library
30 ul	Gibson Assembly 2x Master Mix
to 60 ul	H ₂ O

480

481 50°C incubation for 1 hour.

482

483 - Exonuclease treatment

484

60 ul	Circular assembly reaction
9 ul	ATP (25mM)
9 ul	10x Plasmid Safe Buffer
3 ul	Plasmid Safe Nuclease
9 ul	H ₂ O

485

486 37°C incubation for 1 hour.

487

488 - PCR purify and elute in 50 ul.

489 - Digest to linearize library

490

50 ul	Purified assemblies
10 ul	10x CutSmart Buffer
3 ul	Sspl-HF
37 ul	H ₂ O

491

492 37°C incubation for ≥ 3 hours.

493 - Gel purify 273bp band from 2.5% agarose gel.

494

495 *Note: Band is sometimes fuzzy and poorly visible. If not clearly discernible, proceed with*
496 *gel isolation between 200-300bp.*

497

498 **5. Linearized library QPCR to determine number of amplification cycles for PCR**
499 **amplification**

500 - Set up the following reaction:

501

0.5 %	Purified linearized library
10 ul	NEBNext 2x Master Mix
0.5 ul	20uM PlasmidIns_Fw
0.5 ul	20uM PlasmidIns_Rv
0.2 ul	SybrGreen Dye (100x)
to 20 ul	H ₂ O

502

503 65°C annealing temperature

504

505 - Determine the point that signal amplification has plateaued.

506

507 **6. Linearized Library PCR amplification**

508 - Set up the following reaction:

509

50 %	Purified linearized library
50 ul	NEBNext 2x Master Mix
2.5 ul	20uM PlasmidIns_Fw
2.5 ul	20uM PlasmidIns_Rv
to 100 ul	H ₂ O

510

511 65°C annealing temperature, 1 minute extension time.

512 Cycle number is number of cycles needed to reach signal amplification plateau in the
513 QPCR in step 5, reduced by 4 cycles to scale for increased DNA input.

514

515 - Gel purify 375bp band from 2.5% agarose gel.

516

517 **7. Vector backbone digest**

518 - Set up the following reaction:

519

2 ug	spCas9 sgRNA plasmid (71485)
10 ul	10x Buffer 2.1
3 ul	BbsI
2 ul	XbaI
to 100 ul	H ₂ O

520
521
522
523
524
525
526
527
528

37°C incubation for ≥ 3 hours.

- Gel purify 5.9 kb band from 1% agarose gel.

8. Vector assembly and cleanup

Note: Include a ligation with water for insert as a control.

- Set up the following reaction:

300 ng	Digested vector backbone
42 ng	Amplified Oligo Library
30 ul	Gibson Assembly 2x Master Mix
to 60 ul	H ₂ O

529
530
531
532
533

50°C incubation for 1 hour.

- Isopropanol precipitation

40 ul	Vector assembly reaction
0.4 ul	GlycoBlue Coprecipitant
0.8 ul	50mM NaCl
38.8 ul	Isopropanol

534
535
536
537
538
539
540

- Vortex and incubate at room temperature for 15 minutes.
- Spin down at ≥15.000g for 15 minutes, and carefully remove supernatant.
- Wash pellet with 300ul 80% EtOH and repeat spin at ≥15.000g for 5 minutes.
- Carefully remove all liquid without disturbing pellet, and let air dry for 1-3 minutes.
- Dissolve dried pellet in 10 ul H₂O at 55°C for 10 minutes.

9. Transformation

Note: Electroporation competent cells give a higher transformation efficiency than chemically competent cells. We use NEB10beta electro-competent cells, however these can be substituted for other lines and transformed according to the manufacturer's instructions.

546
547
548
549
550

Note: We use DRM as recovery and culture medium to enhance yield. If substituting for a less rich medium such as LB, we recommend scaling up the culture volume to obtain similar plasmid DNA quantities.

551 *Note: Antibiotic-free recovery time should be limited to 15 minutes to prevent shedding of*
552 *transformed plasmids from replicating bacteria.*

553
554 *Note: Also transform water ligation as control.*
555

- 556 - Pre-warm 3.5mL recovery medium per electroporation reaction, at 37°C for 1 hour.
- 557 - Pre-warm LB-agar plates containing appropriate antibiotic.
- 558 - Per reaction, add 1 ul purified vector assembly to 25ul competent cells on ice.
559 Perform 8 replicate reactions.
- 560 - Electroporate according to the manufacturer's instructions.
- 561 - Immediately add 100 ul pre-warmed recovery media per cuvette and pool all
562 replicates into culture flask.
- 563 - Add 1 mL recovery media per replicate reaction to culture flask and shake at
564 200rpm 37°C for 10 – 15 minutes.
- 565 - Plate a dilution series from 1:10⁴ – 1:10⁶ on LB-agar plates containing antibiotic
566 and grow overnight at 37°C
- 567 - Add 2 mL media per replicate reaction and admix appropriate antibiotic.
- 568 - Grow overnight in shaking incubator at 200rpm 37°C
569
- 570 - Assess transformation efficiency from serial dilution LB-agar plates. Expect ~10⁶
571 clones.

572
573 *The development of this cloning protocol was guided by work described in Videgal et al.*
574 *2015.*
575

576 **Sequence alignment and data processing**

577 For library data, each sequenced pair of gRNA fragment and target was associated with
578 a set of designed sequence contexts G by finding the designed sequence contexts for all
579 gRNAs whose beginning section perfectly matches the gRNA fragment (read 1 in general
580 does not fully sequence the gRNA), and by using locality sensitive hashing (LSH) with 7-
581 mers on the sequenced target to search for similar designed targets. An LSH score on 7-
582 mers between a reference and a sequenced context reflects the number of shared 7-
583 mers between the two. If the best reference candidate scored, through LSH, greater than
584 5 higher than the best LSH score of the reference candidates obtained from the gRNA-
585 fragment, the LSH candidate is also added to G . LSH was used due to extensive (~33%
586 rate) PCR recombination between read1 and read2 which in sequenced data appears as
587 mismatched read1 and read2 pairs. The sequenced target was aligned to each candidate
588 in G and the alignment with the highest number of matches is kept. Sequence alignment
589 was performed using the Needleman-Wunsch algorithm using the parameters: +1 match,
590 -1 mismatch, -5 gap open, -0 gap extend. For library data, starting gaps cost 0. For all
591 other data, starting and ending gaps cost 0. For VO data, sequence alignments were
592 derived from SAM files from SRA.

593
594 Alignments with low-accuracy or short matching sections flanked by long (10 bp+)
595 insertions and deletions were filtered out as PCR recombination products (observed
596 frequency of ~5%). These PCR recombination products are different than that occurring
597 between read1 and read2; these occur strictly in read2. Alignments with low matching
598 rates were removed. Deletions and insertions were shifted towards the expected
599 cleavage site while preserving total alignment score. CRISPR-associated DNA repair
600 events were defined as any alignment with deletions or insertions occurring within a 4 bp
601 window centered at the expected cut site and any alignment with both deletions and
602 insertions (combination indel) occurring with a 10 bp window centered at the expected
603 cut site. All CRISPR-associated DNA repair events observed in control data had their
604 frequencies subtracted from treatment data to a minimum of 0.

605
606 We carried out replicate experiments for library data in each cell type. For each cell-type,
607 each target site not fulfilling the following data quality criteria was filtered: in each
608 replicate, data at this target site must have a total of at least 1,000 reads for all CRISPR
609 editing outcomes at that target site (see section on “Calling CRISPR editing outcomes
610 with high confidence” below for a discussion on the 1,000 reads threshold), and a Pearson
611 correlation of at least 0.85 in the frequency of microhomology-based deletion events. The
612 class of microhomology-based deletion events was used for this criterion since it is a
613 major repair class with the highest average replicability across experiments.

614
615 **Details on alignment processing**

616 All alignments with gaps were shifted as much as possible towards the cleavage site while
617 preserving the overall alignment score. Then, the following criteria were used to
618 categorize the alignments into noise, not-noise but not CRISPR-associated (for example,
619 wildtype); as well as primary and secondary CRISPR activity. All data used in modeling
620 and analysis derive solely from outcomes binned into primary CRISPR activity.

621

622 The following criteria was used to filter library alignments into “noise” categories.

623

624 Homopolymer: Entire read is homopolymer of a single nucleotide. Not considered a
625 CRISPR repair product.

626 Has N: Read contains at least one N. Discarded as noise, not considered a CRISPR
627 repair product.

628 PCR Recombination: Contains recombination alignment signature: (1) if a long indel (10
629 bp+) followed by chance overlap followed by long indel (10 bp+) of the opposite type, e.g.,
630 insertion-randommatch-deletion and deletion-randommatch-insertion. OR, if one of these
631 two indels is 30 bp+, the other can be arbitrarily short. If either criteria is true, and if the
632 chance overlap is length 5 or less, or any length with less than 80% match rate, then it
633 satisfies the recombination signature. In addition, if both indels are 30 bp+, regardless of
634 the middle match region, it satisfies the recombination signature. Finally, if randommatch
635 is length 0, then indel is allowed to be any length. Not considered a CRISPR repair
636 product.

637 Poor-Matches: 55bp designed sequence context has less than 5 bp representation (could
638 occur from 50 bp+ deletions or severe recombination) or less than 80% match rate. Not
639 considered a CRISPR repair product.

640 Cutsite-Not-Sequenced: The read does not contain the expected cleavage site.

641 Other: An alignment with multiple indels where at least one non-gap region has lower
642 than an 80% match rate. Or generally, any alignment not matching any defined category
643 above or below. In practice, can include near-homopolymers. Not considered a CRISPR
644 repair product.

645

646 The following criteria was used to filter library alignments into “main” categories.

647 Wildtype: No indels in all of alignment. Not considered a CRISPR repair product.

648

649 Deletion: An alignment with only a single deletion event. Subdivided into:

650 Deletion - Not CRISPR: Single deletion occurs outside of 4 bp window centered around
651 cleavage site. Not considered a CRISPR repair product.

652 Deletion - Not at cut: Single deletion occurring within 4 bp window centered around
653 cleavage site, but not immediately at cleavage site. Considered a CRISPR repair product.

654 Deletion: Single deletion occurring immediately at cleavage site. Considered a CRISPR
655 repair product.

656

657 Insertion: An alignment with only a single insertion event. Subdivided into:

658 Insertion - Not CRISPR: Single insertion occurs outside of 10 bp window around cleavage
659 site. Not considered a CRISPR repair product.

660 Insertion - Not at cut: Single insertion occurring within 4 bp window centered around
661 cleavage site, but not immediately at cleavage site. Considered a CRISPR repair product.

662 Insertion: Single insertion occurring immediately at cleavage site. Considered a CRISPR
663 repair product.

664

665 Combination indel: An alignment with multiple indels where all non-gap regions have at
666 least 80% match rate. Subdivided into:

667 Combination Indel: All indels are within a 10 bp window around the cleavage site.
668 Considered a primary CRISPR repair product.
669 Forgiven Combination Indel: At least two indels, but not all, are within a 10 bp window
670 around the cleavage site. Considered a rarer secondary CRISPR repair product, ignored.
671 Forgiven Single Indel: Exactly one indel is within a 10 bp window around the cleavage
672 site. Considered a rarer secondary CRISPR repair product, ignored.
673 Combination Indel - Not CRISPR: No indels are within a 10 bp window around the
674 cleavage site. Not considered a CRISPR repair product.

675
676 We note that deletion and insertion events, even those spanning many bases, are defined
677 to occur at a single location between bases. As such, events occurring up to 5 bp away
678 from the cleavage site are defined as events where there are five or fewer
679 matched/mismatched alignment positions between the event and the cleavage site,
680 irrespective of the number of gap dashes in the alignment.

681
682 **Calling CRISPR editing outcomes with high confidence**
683 Following the processing steps above, we performed the following further analysis and
684 processing steps to call high-confidence CRISPR editing outcomes. These steps largely
685 follow heuristics, and we believe that a thorough and unbiased methodological
686 standardization in counting CRISPR editing outcomes will be valuable future work.

687
688 DNA repair at Cas9-mediated double-strand breaks is known to result in a large diversity
689 of outcomes, with indels of varying length and positions around the cleavage site. The
690 frequencies of many of these editing outcomes, though enriched in Cas9-treatment data
691 over control data, are rare (<0.5% of all edited products) and can be challenging to assign
692 as a CRISPR editing outcome due to a lack of foundational biological or computational
693 models on the exact mechanisms of DNA repair. In addition, rare outcomes can
694 sometimes be attributed to sequencing errors.

695
696 In this work, we focus on CRISPR editing outcomes that are enriched in treatment data
697 over control data and agree with a relatively conservative and strict model of DNA repair,
698 in order to ensure a high degree of confidence in the editing outcomes that we call. As a
699 result, we underestimate the total number of unique CRISPR editing outcomes, though
700 we believe this underestimation is not by an order of magnitude, though it may be by a
701 factor of 2x or so.

702
703 We define high-confidence CRISPR editing outcomes as bins of alignments categorized
704 by the previously described pipeline into CRISPR-associated categories that have no
705 mismatches. Each unique deletion genotype consistent with microhomology is treated as
706 a single unique outcome, though we note that microhomology deletions may arise by
707 noise or chance though we expect this to be a rare event. Each unique insertion genotype
708 is also treated as a single unique outcome, and as with MH deletions, we note that some
709 insertions may arise by noise or chance though we anticipate this to be rare. In sum, we
710 likely overestimate by a slight amount the true number of unique microhomology deletion
711 and insertion events.

712

713 All microhomology-less deletion genotypes are binned together for a particular deletion
714 length, which almost always will bin together multiple unique MH-less deletion genotypes.
715 However, the class of MH-less deletions, in general, has lower replicate consistency and
716 higher stochasticity than MH deletions, and the space of all possible MH-less deletions is
717 orders of magnitude larger than that of MH deletions. The class of MH-less deletions is
718 also less frequent than MH deletions in all five human and mouse cell types we examined.
719 In sum, we characterize MH-less deletions as comprising a large number of rare
720 genotypes that lack high replicate consistency. As such, we conservatively count all
721 binned MH-less deletions for a particular deletion length as a single unique outcome. In
722 sum, we likely underestimate by a moderate amount the true number of unique
723 microhomology-less deletion events.

724
725 As MH-less deletions represents the larger space of possible unique genotypes, in total,
726 we are likely underestimating the total number of unique outcomes in our procedure that
727 calls unique outcomes with high confidence. We provide statistics on the total number of
728 high-confidence unique outcomes in the manuscript and in Extended Data Fig. 1 and 5.

729
730 Based on a computational simulation of subsampling the data, we empirically set 1,000
731 reads per target site as a minimum quality threshold. The diversity of editing outcomes
732 requires some minimum read count to consider the data as representative of editing
733 outcomes at that target site. In Lib-A and Lib-B data in U2OS and mESCs, we empirically
734 observe that 1,000 reads per target site lies above the “elbow” in the curve plotting the
735 number of unique high-confidence outcomes and subsampled read count. We
736 recommend this quality filtering methodology in general for future work studying CRISPR
737 editing outcomes, and based on our data, empirically suggest that 1,000 reads per target
738 site may be a useful guideline for future experimental design.

739 740 **Controlling for cell-type specific 1-bp insertion frequencies when measuring** 741 **replicability of indel frequencies across cell-types**

742 All indels not belonging to “all major repair outcomes” were filtered out. To adjust the
743 frequencies of all 1-bp insertion genotypes in a target site in two cell-types, the average
744 of the total 1-bp insertion frequency among all major repair outcomes was calculated
745 between the two cell-types, then frequencies of each 1-bp insertion genotype was
746 adjusted proportionally such that the resulting total 1-bp insertion frequency in that is
747 equal to the aforementioned average, and thereby equal to the adjusted 1-bp insertion
748 frequency in the data from the other cell-type.

749 750 **Selection of variants from disease databases**

751 Disease variants were selected from the NCBI ClinVar database (downloaded September
752 9, 2017)⁹ and the Human Gene Mutation Database (publicly available variant data from
753 before 2014.3)¹⁰ for computational screening and subsequent experimental correction.

754
755 A total of 4,935 unique variants were selected from Clinvar submissions where the
756 functional consequence is described as complete insertions, deletions, or duplications
757 where the reference or alternate allele is of length less than or equal to 30 nucleotides.
758 Variants were included where at least one submitting lab designated the clinical

759 significance as 'pathogenic' or 'likely pathogenic' and no submitting lab had designated
760 the variant as 'benign' or 'likely benign', including variants with all disease associations.
761 More complex indels and somatic variants were included. A total of 18,083 unique
762 insertion variants were selected from HGMD which were between 2 to 30 nucleotides in
763 length. Variants were included with any disease association with the HGMD classification
764 of 'DM' or disease-causing mutation.

765
766 SpCas9 gRNAs and their cleavage sites were enumerated for each disease allele. Using
767 a previous version of inDelphi, genotype frequency and indel length distributions were
768 predicted for each tuple of disease variant and unique cleavage site. Among each unique
769 disease, the single best gRNA was identified as the gRNA inducing the highest predicted
770 frequency of repair to wildtype genotype, and if this was impossible (due to, for example,
771 a disease allele with 2+ bp deletion), then the single best gRNA was identified as the
772 gRNA inducing the highest predicted frameshift repair rate. 1327 sequence contexts were
773 designed in this manner for Lib-B. An additional 265 sequence contexts were designed
774 by taking the 265 sequence contexts in any disease in decreasing order of predicted
775 wildtype repair rate, starting with Clinvar, stopping at 45% wildtype repair rate, then
776 continuing with HGMD. This yielded 1592 total sequences derived from Clinvar and
777 HGMD.

778
779

780 **Definition of Delta-Positions**

781 Using the MMEJ mechanism, deletion events can be predicted at single-base resolution.
782 For computational convenience, we use the tuple (deletion length, delta-position) to
783 construct a unique identifier for deletion genotypes. A delta-position associated with a
784 deletion length N is an integer between 0 and N inclusive (Extended Data Fig. 2). In a
785 sequence alignment, a delta-position describes the starting position of the deletion gap in
786 the read with respect to the reference sequence relative to the cleavage site. For a
787 deletion length N and a cleavage site at position C such that seq[:C] and seq[C:] yield the
788 expected DSB products where the vector slicing operation vector[index1:index2] is
789 inclusive on the first index and exclusive on the second index (python style), a delta-
790 position of 0 corresponds to a deletion gap at seq[C-N+0 : C+0], and generally with a
791 delta-position of D, the deletion gap occurs at seq[C-N+D : C+D]. Microhomologies can
792 be described with multiple delta-positions. To uniquely identify microhomology-based
793 deletion genotypes, the single maximum delta-position in the redundant set is used.
794 Microhomology-less deletion genotypes are associated with only a single delta position
795 and deletion length tuple; we use this as its unique identifier.

796
797 Another way to define delta-positions can be motivated by the example workflow in the
798 Supplementary Discussion on MH deletions describing how each microhomology is
799 associated with a deletion genotype. In that workflow, the delta-position is the number of
800 bases included on the top strand before “jumping down” to the bottom strand.

801
802 MH-less medial end-joining products correspond to all MH-less genotypes with delta-
803 position between 1 and N-1 where N is the deletion length. MH-less unilateral end-joining
804 products correspond to MH-less genotypes with delta-position 0 or N. We note that a
805 deletion genotype with delta position N does not immediately imply that it is a
806 microhomology-less unilateral end-joining product since it may contain microhomology
807 (it’s possible that delta-positions N-j, N-j+1, ..., N all correspond to the same MH deletion.)
808

809 **Definition of Precision Score**

810 For a distribution X, where |X| indicates its cardinality (or length when represented as a
811 vector):

$$812 \quad \text{PrecisionScore}(X) = 1 - \frac{-\sum_{i=1}^n P(x_i) \log(P(x_i))}{\log(|X|)}$$

813 This precision score ranges between zero (minimally precise, or highest entropy) to one
814 (maximally precise, or lowest entropy).

815 816 **inDelphi Deletion Modeling: Neural network input and output**

817 inDelphi receives as input a sequence context and a cleavage site location, and outputs
818 two objects: a frequency distribution on deletion genotypes, and a frequency distribution
819 on deletion lengths.

820
821 To model deletions, inDelphi trains two neural networks: MH-NN and MHless-NN. MH-
822 NN receives as input a microhomology that is described by two features: microhomology
823 length and GC fraction in the microhomology. Using these features, MH-NN outputs a

824 number (psi). MHless-NN receives as input the deletion length. Using this feature,
825 MHless-NN outputs a number (psi).

826

827 A phi score is obtained from a psi score using: $\text{phi}_i = \exp(\text{psi}_i - 0.25 * \text{deletion_length})$,
828 where 0.25 is a “redundant” hyperparameter that serves to increase training speed by
829 helpful scaling. This relationship between psi and phi is differentiable and encodes the
830 assumption that the frequency of an event exponentially increases with neural network
831 output psi (which empirically appears to reflect MH strength) and exponentially decreases
832 with its minimum necessary resection length (deletion length).

833

834 The architecture of the MH-NN and MHless-NN networks are *input-dimension* -> 16 -> 16
835 -> 1 for a total of two hidden layers where all nodes are fully connected. Sigmoidal
836 activations are used in all layers except the output layer. All neural network parameters
837 are initialized with Gaussian noise centered around 0. No regularization or dropout was
838 used.

839

840 **inDelphi Deletion Modeling: Making predictions**

841 Given a sequence context and cleavage site, inDelphi enumerates all unique deletion
842 genotypes as a tuple of its deletion length and its delta-position for deletion lengths from
843 1 bp to 60 bp. For each microhomology enumerated, an MH-phi score is obtained using
844 MH-NN. In addition, for each deletion length from 1 bp to 60 bp, an MHindep-phi score is
845 obtained using MHless-NN.

846

847 inDelphi combines all MH-phi and MHindep-phi scores for a particular sequence context
848 into two objects – a frequency distribution on deletion genotypes, and a frequency
849 distribution on deletion lengths – which are both compared to observations for training.
850 The model is designed to output two separate objects because both are of biological
851 interest, and separate but intertwined modeling approaches are useful for generating
852 both. By learning to generate both objects, inDelphi jointly learns about microhomology-
853 based deletion repair and microhomology-less deletion repair.

854

855 To generate a frequency distribution on deletion genotypes, inDelphi assigns a score for
856 each microhomology. Score assignment considers the concept of “full” microhomology
857 and treats full and not full MHs differently.

858

859 A microhomology is “full” if the length of the microhomology is equal to its deletion length.
860 The biological significance of full microhomologies is that there is only a single deletion
861 genotype possible for the entire deletion length, while in general, a single deletion length
862 is consistent with multiple genotypes. In addition, this single genotype can be generated
863 through not just the MH-dependent MMEJ mechanism but also through MH-less end-
864 joining, for example as mediated by Lig4. Therefore, we model full microhomologies as
865 receiving contributions from both MH-containing and MH-less mechanisms by scoring full
866 microhomologies as $\text{MH-phi}[i] + \text{MHindep-phi}[j]$ for deletion length j and microhomology
867 index i .

868

869 Microhomologies that are not “full” are assigned a score of $\text{MH-phi}[i]$ for MH index i .

870
871 Scores for all deletion genotypes assigned this way are normalized to sum to 1 to produce
872 a predicted frequency distribution on deletion genotypes.

873
874 To generate a frequency distribution on deletion lengths, inDelphi assigns a score for
875 each deletion length. Score assignment integrates contributions from both MH-dependent
876 and MH-independent mechanisms via the following procedure: For each deletion length
877 j , its score is assigned as MHindep- $\phi[j]$ plus the sum of MH- ϕ for each microhomology
878 with that deletion length. Scores for all deletion lengths are normalized to sum to 1 to
879 produce a frequency distribution.

880
881 inDelphi trains its parameters using a single sequence context by producing both a
882 predicted frequency distribution on deletion genotypes and deletion lengths and
883 minimizing the negative of the sum of two values: the mean squared Pearson correlation
884 for the deletion genotype frequency distribution at each target site in the training set plus
885 the mean squared Pearson correlation for the deletion length frequency distribution at
886 each target site in the training set. This represents a multitask learning framework.

887
888 In practice, deletion genotype frequency distributions are formed from observations for
889 deletion lengths 1-60, and deletion length frequency distributions are formed from
890 observations for deletion lengths 1-28. Both neural networks are trained jointly and
891 simultaneously on both tasks. inDelphi is trained with stochastic gradient descent with
892 batched training sets. inDelphi is implemented in Python using the autograd library. We
893 used a batch size of 200, an initial weight scaling factor of 0.10, an initial step size of 0.10,
894 and an exponential decaying factor for the step size of 0.999 per step. We observed
895 performance convergence within about 50 epochs.

896 897 **inDelphi Deletion Modeling: Multitask learning improves performance**

898 Over the course of developing our model, at an intermediate stage we considered a
899 simpler model for predicting the frequencies of MH deletions. This model featurizes all
900 sequence microhomologies at a target site using MH length and GC content and uses a
901 2x16x16x1 neural network with sigmoidal activations except at the output layer to output
902 ψ . This ψ value is adjusted using $\exp(\psi - 0.25 * \text{deletion length})$ to obtain ϕ for a
903 particular microhomology, which are normalized across all microhomologies to sum to 1
904 to achieve a predicted distribution of frequencies. Altogether, this model is identical to the
905 MH module used in inDelphi, with the notable difference of not including contributions of
906 MH-less ϕ at “full” microhomologies.

907
908 This simple model, henceforth known as the baseline model, does not recognize the
909 possibility that a MH genotype may arise from both MH and MH-independent repair
910 pathways. We compared the performance of the baseline model to inDelphi’s MH module
911 and observed a statistically significant relative improvement of 10% in model performance
912 as measured on test set data ($p \sim 0.02$). These measurements were performed using Lib-
913 A target sites in mESCs.

914

915 We note that the multitask model used in inDelphi also jointly trains MHless-NN and, in
916 addition to predicting MH deletion frequencies more accurately than the baseline, also
917 provides strong performance on deletion length frequency prediction.
918

919 Using random seed A, the baseline model mean Pearson r on held-out data was .905,
920 while the multitask model mean Pearson r on the same held-out data was 0.913, for a
921 8.5% relative improvement ($p = 0.009$, one-sided t -test). Using random seed B, the
922 baseline model mean Pearson r on held-out data was .924, while the multitask model
923 mean Pearson r on the same held-out data was 0.928, for a 5.3% relative improvement
924 ($p = 0.02$, one-sided t -test). Using random seed C, the baseline model mean Pearson r
925 on held-out data was .912, while the multitask model mean Pearson r on the same held-
926 out data was 0.917, for a 5.7% relative improvement ($p = 0.03$, one-sided t -test).
927

928 **inDelphi Deletion Modeling: Summary and Revisiting Assumptions**

929 In summary, inDelphi trains MH-NN, which uses as input (microhomology length,
930 microhomology GC content) to output a psi score which is translated into a phi score
931 using deletion length. This phi score represents the “strength” of the microhomology
932 corresponding to a particular MH deletion genotype. It also trains MHless-NN which uses
933 as input (deletion length) to directly output a phi score representing the “total strength” of
934 all MH-independent activity for a particular deletion length.
935

936 While the model assumes that microhomology and microhomology-less repair can
937 overlap in contributions to a single repair genotype, this assumption is made
938 conservatively by assuming that their contributions overlap only when there is no
939 alternative. Specifically, in the context of a deletion length with full microhomology, the
940 model assumes that they must overlap, while in the context of a deletion length without
941 full microhomology, inDelphi allows MHindep-phi to represent all MH-less repair
942 genotypes and none of the MH-dependent repair genotypes which are represented solely
943 using their MH-phi scores. This can be seen by noting that at a deletion length j without
944 full microhomology, MH genotypes are scored using their MH-phi scores, while the length
945 j is scored by MHindep-phi $[j]$ plus the sum of MH-phi for each microhomology. Therefore,
946 the subset of MH-less genotypes at this deletion length have a score MHindep-phi $[j]$.
947

948 When the subset of MH-less genotypes includes only one MH-less genotype, this single
949 genotype’s score is equal to MHindep-phi $[j]$. In general, multiple MH-less genotypes are
950 possible, in which case the total score of all of the MH-less genotypes is equal to
951 MHindep-phi $[j]$.
952

953 The relative frequency of MH deletions and MH-less deletions is learned implicitly by the
954 balancing between the sum of all MH-phi and MHindep-phi. Since MHindep-phi does not
955 vary by sequence context while MH-phi does, the model assumes that variation in the
956 fraction of deletions that use MH can at least partially be explained by varying sequence
957 microhomology as represented by MH-NN.
958

959 **inDelphi Insertion Modeling**

960 Once inDelphi is trained on both deletion tasks, it predicts insertions from a sequence
961 context and cleavage site by using the precision score of the predicted deletion length
962 distribution and total deletion phi (from all MH-phi and MHindep-phi). inDelphi also uses
963 one-hot-encoded binary vectors encoding nucleotides -4 and -3. In a training set, these
964 features are collected and normalized to zero mean and unit variance, and the fraction of
965 1-bp insertions over the sum counts of 1-bp insertions and all deletions are tabulated as
966 the prediction goal. A k -nearest neighbor model is built using the training data. inDelphi
967 uses the default parameter $k = 5$.

968
969 On test data, the above procedure is used to predict the frequency of 1-bp insertions out
970 of 1-bp insertions and all deletions for a particular sequence context. Once this frequency
971 is predicted, it is used to make frequency predictions for each of the 4 possible insertion
972 genotypes, which are predicted by deriving from the training set the average insertion
973 frequency for each base given its local sequence context. When the training set is small,
974 only the -4 nucleotide is used. When the training set is relatively large, nucleotides -5, -4,
975 and -3 are used.

976
977 To produce a frequency distribution on 1-bp insertions and 1-60 bp deletion genotypes,
978 scores for all deletion genotypes and all 1-bp insertions are normalized to sum to 1. To
979 produce a frequency distribution on indel lengths (+1 to -60), scores for all deletion lengths
980 and 1-bp insertions are normalized to sum to 1.

981 **inDelphi: Repair classes predicted at varying resolution**

982 inDelphi predicts MH-deletions and 1-bp insertions at single base resolution. Measuring
983 performance on the task of genotype frequency prediction considers this subset of repair
984 outcomes only (about 60-70% of all outcomes).

985
986 inDelphi predicts MH-less deletions to the resolution of deletion length. That is, inDelphi
987 predicts a single frequency corresponding to the sum total frequency of all unique MH-
988 less deletion genotypes possible for a particular deletion length. This modeling choice
989 was made because genotype frequency replicability among MH-less deletions is
990 substantially lower than among MH deletions.

991
992 Measuring performance on the task of indel length frequency considers MH deletions,
993 MH-less deletions, and 1-bp insertions (90% of all outcomes).

994
995 In practice, if end-users desire, they can extend inDelphi predictions to frequency
996 predictions for specific MH-less deletion genotypes by noting that MH-less deletions are
997 distributed uniformly between 0 delta-position genotypes, medial genotypes, and N delta-
998 position genotypes.

1000 **Comparison with a linear baseline model**

1001 We compared inDelphi to a baseline model with the same model structure but replacing
1002 the deep neural networks with linear models. We compared using Lib-A mESC data.
1003 While inDelphi achieves a mean held-out Pearson correlation of 0.851 on deletion
1004 genotype frequency prediction and 0.837 on deletion length frequency prediction, the
1005

1006 linear baseline model achieves a mean held-out Pearson correlation of 0.816 on deletion
1007 genotype frequency prediction and 0.796 on deletion length frequency prediction. When
1008 including the third model component for 1-bp insertion modeling and testing on genotype
1009 frequency prediction for 1-bp insertions and all deletions, inDelphi achieves a median
1010 held-out Pearson correlation of 0.937 and 0.910 on the task of indel length frequency
1011 prediction. The linear baseline model achieves a median held-out Pearson correlation of
1012 0.919 and 0.900 on the two tasks respectively.

1013
1014 From these results, we can see that much of the model's power is derived from its
1015 designed structure which is independent of the choice of linear or non-linear modeling.
1016 While the baseline does not significantly cripple the model, the use of deep nonlinear
1017 neural networks offers a substantial performance improvement (10-24%) above linear
1018 modeling. In addition, the strong performance of the linear baseline model highlights that
1019 the prediction task, given the model structure, is relatively straightforward. This suggests
1020 that our model should be able to generalize well to unseen data.

1021
1022 The deep neural network version of MH-NN learns that microhomology length is more
1023 important than % GC (Extended Data Fig. 2). The linear version learns the same
1024 concept, with a weight of 1.1585 for MH length and 0.332 for % GC.

1025
1026 **Comparison with a baseline model lacking microhomology length as a feature**
1027 Microhomology length is an important feature for MH-NN (Extended Data Fig. 2). We
1028 trained a model that uses only % GC as input to MH-NN while keeping the rest of the
1029 model structure identical. On held-out data at Lib-A target sites in mESCs, this baseline
1030 model at convergence achieves to a mean Pearson correlation of 0.59 on the task of
1031 predicting deletion genotype frequencies, and a mean Pearson correlation of 0.58 on
1032 the task of predicting deletion length frequencies. Notably, a model at iteration 0 with
1033 randomly initialized weights achieves mean Pearson correlations of 0.55 and 0.54 on
1034 the two respective tasks on held-out data. This basal Pearson correlation is relatively
1035 high due to the model structure, in particular, the exponential penalty on deletion length.
1036 In sum, removing MH length as a feature severely impacts model performance,
1037 restricting it to predictive performance not appreciably better than random chance.

1038
1039 **inDelphi training and testing on data from varying cell-types**

1040 For predicting genotype and indel length frequencies in any particular cell-type C where
1041 data D is available, we first trained inDelphi's deletion component on a subset of Lib-A
1042 mESC data. Then, we apply k -fold cross-validation on D where D is iteratively split into
1043 training and test datasets. For each cross-validation iteration, the training set is used to
1044 train the insertion frequency model (k -nearest neighbors) and insertion genotype model
1045 (matrix of observed probabilities of each inserted base given local sequence context,
1046 which is just the -4 nucleotide when the training dataset is small, and -5, -4 and -3
1047 nucleotides when the training dataset is large). For each cross-validation iteration,
1048 predictions are made at each sequence context in the test set which are compared to
1049 observations for each sequence context to yield a Pearson correlation. For any particular
1050 sequence context, the median test-time Pearson correlation across all cross-validation
1051 iterations is used as a single number summary of the overall performance of inDelphi. For

1052 all reported results, we used 100-fold cross-validation with 80%/20% training and testing
1053 splits. Empirically, we observed small variance in test-time Pearson correlation,
1054 highlighting the stability of inDelphi's modeling approach.

1055
1056 **inDelphi testing on endogenous VO data**
1057 On this task, the deletion component of inDelphi was trained on a subset of the Lib-A
1058 mESC data. For each cell type in HCT116, K562, and HEK293T, all VO sequence
1059 contexts (about 100) were randomly split into training and test datasets 100 times. During
1060 each split, the training set was used for *k*-nearest neighbor modeling of 1-bp insertion
1061 frequencies. Feature normalization to zero mean and unit variance was not performed.
1062 The average frequency of each 1-bp insertion genotype was derived from the training set
1063 as well. For each of the ~100 sequence contexts, the median test-time Pearson
1064 correlation was used for plotting in Figure 3. Due to the small size of the training set, only
1065 the -4 nucleotide was used for modeling both the insertion frequency and insertion
1066 genotype frequencies.

1067
1068 **inDelphi testing on library data**
1069 On this task, the deletion component of inDelphi was trained on a subset of the Lib-A
1070 mESC data. The remaining test set was used for measuring test-time prediction
1071 performance on Lib-A. Nucleotides -5, -4, and -3 were used for the insertion genotype
1072 model. For testing on Lib-B, Lib-B was split into training and test datasets in the same
1073 manner as with VO data. Nucleotide -4 was used for the insertion genotype model. The
1074 median test-time Pearson correlation is used as a single number summary of the overall
1075 performance of inDelphi on any particular sequence context. For reporting predictive
1076 results in Figure 4, sequence contexts with low replicability (less than 0.85 Pearson
1077 correlation) in observed editing outcome frequencies were first removed.

1078
1079 **inDelphi training and testing on *Prkdc^{-/-}Lig4^{-/-}* data**
1080 inDelphi was trained on data from 946 Lib-A sequence contexts and tested on 168 held-
1081 out Lib-A sequence contexts. Nucleotide -4 was used for insertion rate modeling, all other
1082 modeling choices were standard as described above. On held-out data, this version of
1083 inDelphi achieved a median Pearson correlation of 0.84 on predicting indel genotype
1084 frequencies, and 0.80 on predicting indel length frequencies.

1085
1086 **Training the online public version of inDelphi and its expected properties**
1087 For general-use on arbitrary cell types, we trained a version of inDelphi using additional
1088 data from diverse types of cells. Deletion modeling was trained using data from 2,464
1089 sequence contexts from high-replicability Lib-A and Lib-B data (including clinical variants
1090 and microduplications, fourbp, and longdup) in mES and data from VO sequence contexts
1091 in HEK293 and K562. Insertion frequency modeling is implemented as above. Insertion
1092 genotype modeling uses nucleotides -5, -4, and -3. The insertion frequency model and
1093 insertion genotype model are trained on VO endogenous data in K562 and HEK293T,
1094 Lib-A data in mESC, and Lib-B data (including clinical variants and microduplications,
1095 fourbp, and longdup) in mESC and U2OS.

1096

1097 Though MHless-NN, as trained on library data, never receives information on deletion
1098 lengths beyond 28, we allow it to generalize its learned function and make predictions on
1099 deletion lengths up to 60 bp to match the supported range of MH-NN.

1100
1101 inDelphi makes predictions on 1-bp insertions and 1-60-bp deletions, which we
1102 empirically show to consist of higher than 90% of all Cas9 editing outcomes in data from
1103 multiple human and mouse cell lines. Nevertheless, there is a subset of repair (about 8%
1104 on average) that inDelphi does not attempt to predict. We suggest that end-users,
1105 depending on what predictive quantities are of interest, take this into account when using
1106 inDelphi. For example, if inDelphi predicts that 60% of 1-bp insertions and 1-60-bp
1107 deletions at a disease allele correspond to repair to wildtype genotype, a quantity of
1108 interest may be the rate of wildtype repair in all Cas9 editing outcomes (including the 8%
1109 not predicted by inDelphi). In such a situation, this quantity can be calculated as
1110 $(92\% * 60\%) = 55.2\%$.

1111
1112 By the design of 1872 sequence contexts in Lib-A, our training dataset has rich and
1113 uniform representation across all quintiles of several major axes of variation including GC
1114 content, precision, and number of bases participating in microhomology as measured
1115 empirically in the human genome. This design strategy enables inDelphi to generalize
1116 well to arbitrary sequence contexts from the human genome.

1117
1118 These training data further include data in the outlier range of statistics of interest,
1119 including extremely high and low precision repair distributions, and extremely weak and
1120 strong microhomology (minimal microhomology and extensive microduplication
1121 microhomology sequences). The availability of such sequences in our training data
1122 enables inDelphi to generalize well to sequence contexts of clinical interest and sequence
1123 contexts supporting unusually high frequencies of precision repair. In particular, by
1124 training on more than 1000 examples of repair at clinical microduplications, inDelphi has
1125 received strong preparation for accurate prediction on other clinical microduplications.

1126
1127 By training on data from many cell-types, we enable inDelphi to make predictions that are
1128 generally applicable to many human cell-types. We note that the HCT116 human colon
1129 cancer cell line experiences a markedly higher frequency of single base insertions
1130 compared to all other cell lines we studied, possibly due to the MLH1 deficiency of this
1131 cell line leading to impaired DNA mismatch repair. For this reason, we excluded HCT116
1132 data from our training dataset. For best results, we suggest end-users keep in mind that
1133 repair class frequencies can be cell type-dependent, and this issue has not been well-
1134 characterized thus far.

1135
1136 We note that inDelphi's main error tendency is on the side of overestimating rather than
1137 underestimating the precision of repair (Figure 4). In general, this tendency can be
1138 explained by noting that inDelphi only considers sequence microhomology as a factor,
1139 while it's plausible and likely in biological experimental settings that even sequence
1140 contexts with very strong sequence microhomology may not yield precise results due to
1141 noise factors that are not considered by inDelphi. For best results, we recommend end-
1142 users take this tendency into account when using inDelphi predictions for further

1143 experiments. In particular, if gRNAs are designed by using a minimum precision
1144 threshold, end-users should recognize that observed repair outcomes may have empirical
1145 precision under this threshold. However, conversely, it is unlikely that a gRNA will have
1146 precision higher than what inDelphi predicts.
1147

1148 **Lib-A design**

1149 All designed sequence contexts were 55 bp in length with cutting between the 27th and
1150 28th base.
1151

1152 1872 sequence contexts were designed by empirically determining the distribution of four
1153 statistics in sequence contexts from the human genome. These four statistics are GC
1154 content, total sum of bases participating in microhomology for 3-27-bp deletions, Azimuth
1155 predicted on-target efficiency score, and the statistical entropy of the predicted 3-27-bp
1156 deletion length distribution from a previous version of inDelphi. For each of these
1157 statistics, empirical quintiles were derived by calculating these statistics in a large number
1158 of sequence contexts from the human genome. For the library, sequence contexts were
1159 designed by randomly generated DNA that categorized into each combination of quintiles
1160 across each of the four statistics. For example, a sequence context falling into the 1st
1161 quintile in GC, 2nd quintile of total MH, 1st quintile of Azimuth score, and 5th quintile of
1162 entropy, was found by random search. With four statistics and five bins each (due to
1163 quintiles), there are $5^4 = 625$ possible combinations. For each combination, we attempted
1164 to design three sequence contexts for a total of 1875; 3 sequences could not be designed
1165 (for a total of 1872) though each bin was filled. 90 sequence contexts were designed from
1166 VO sequence contexts. Other sequence contexts were also designed for a total of 2000
1167 sequence contexts in Lib-A. Lib-A sequence names, gRNAs, and sequence contexts are
1168 listed in Supplementary Table 2.
1169

1170 **Lib-B design**

1171 All designed sequence contexts were 55 bp in length with cutting between the 27th and
1172 28th base.
1173

1174 1592 sequence contexts were designed from Clinvar and HGMD (see section on
1175 Selection of variants from disease databases). Some disease sequence contexts were
1176 designed that such that the corrected wildtype or frameshift allele supports further cutting
1177 by the original gRNA; data from such sequence contexts were ignored during analysis.
1178 57 “longdup” sequence contexts were designed by repeating the following procedure
1179 three times: for $N = 7$ to 25, an N-mer was randomly generated, then duplicating and
1180 surrounded by randomly generated sequences, while ensuring that SpCas9 NGG was
1181 included and appropriately positioned for cutting between positions 27 and 28. 90
1182 sequence contexts were designed from VO sequence contexts. 228 “fourbp” sequence
1183 contexts were designed at 3 contexts with random sequences (with total phi score on
1184 average lower than VO sequence contexts) while varying positions -5 to -2; for each of
1185 the 3 “low-microhomology” contexts, 76 four bases were randomly designed while
1186 ensuring representation from all possible 2 bp microhomology patterns including no
1187 microhomology, one base of microhomology at either position, and full two bases of
1188 microhomology. Other sequence contexts were also designed for a total of 2000

1189 sequence contexts in Lib-B. Lib-B sequence names, gRNAs, and sequence contexts are
1190 listed in Supplementary Table 3.

1191

1192 **1bpInsDisLib design**

1193 12 sequence contexts were designed from Clinvar and HGMD. Pathogenic alleles were
1194 selected for a high predicted frequency of correction to the wild-type genotype via a Cas9-
1195 mediated 1-bp insertion. Sequence names, gRNAs, and sequence contexts are listed in
1196 Supplementary Table 4.

1197

1198 **PHG design**

1199 18 sequence contexts were designed using inDelphi to select SpCas9 gRNAs targeting
1200 the coding regions of genes including VEGFA, VEGFR2, PDCD1, APOB, CCR5, CD274,
1201 CXCR4, PCSK9, and APOBEC3B, such that a frameshift would be induced with higher
1202 frequency than typical SpCas9 gRNAs. Of these 18 frameshift designs, 10 were designed
1203 to induce a single deletion genotype with high precision, and 8 were designed to induce
1204 a single 1-bp insertion genotype with high precision. 6 sequence contexts were designed
1205 using inDelphi from Clinvar and HGMD where pathogenic 1-bp insertion alleles were
1206 selected based on a high predicted frequency of induction from Cas9 editing of the wild-
1207 type allele. Sequence names, gRNAs, and sequence contexts are listed in
1208 Supplementary Table 5.

1209

1210 **Generating a DNA motif for 1-bp insertion frequencies**

1211 Nucleotides from positions -7 to 0 were one-hot-encoded and used in ridge regression to
1212 predict the observed frequency of 1-bp insertions out of all Cas9 editing events in 1996
1213 sequence contexts from Lib-A mESC data. The data were split into training and testing
1214 sets (80/20 split) 10,000 times to calculate a bootstrapped estimate of linear regression
1215 weights and test-set predictive Pearson correlation. The median test-set Pearson
1216 correlation was found to be 0.62. To generate a DNA motif, any features that included 0
1217 within the bootstrapped weight range were excluded (probability that the weight is zero $>$
1218 $1e-4$). The average bootstrapped weight estimate was used as the “logo height” for all
1219 remaining features. Each feature is independent; vertical stacking of features follows the
1220 published tradition of DNA motifs.

1221

1222 **Predicting precision repair of genomic SpCas9 gRNAs**

1223 In this work, we determined the distributions of the most frequent deletion and
1224 insertion outcomes among major editing outcomes at SpCas9 gRNAs targeting human
1225 exons and introns as predicted by inDelphi trained on data from Lib-A target sites in
1226 mESCs and U2OS cells separately (Fig. 3f, Extended Data Table 1). A combination of
1227 computational constraint (the inability to make predictions at ~350 million target sites
1228 comprising all SpCas9 gRNAs in the human genome), uncertainty in the exact predictions
1229 of the model and a preference for avoiding overfitting our training data, and lack of
1230 sufficient held-out data to verify our predictions and identify potential bias, motivated us
1231 to smooth the exact predictions made by the model. We resampled each predicted value
1232 from a Gaussian centered at the predicted value with a specified standard deviation. For
1233 mESCs, we set the standard deviation as the predicted value divided by 4, up to a
1234 maximum of 3% for insertions, while for deletions we used the predicted value divided by

1235 4 with a minimum of 6%. For U2OS cells, we set the standard deviation as the predicted
1236 value divided by 4 for insertions, and the predicted value divided by 4 with a minimum of
1237 6% for deletions. The scaling of standard deviation at higher predicted values reflects the
1238 abundance of data and therefore higher relative confidence at lower predicted values.
1239 The use of symmetrical noise reflects our prior belief that our predictions are equally likely
1240 to underestimate and overestimate the true value.
1241

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282

Plasmid and insert sequences

P2T-CAG-MCS-P2A-GFP-PuroR complete plasmid sequence

CCACCTAAATTGTAAGCGTTAATATTTTGTAAAATTCGCGTTAAATTTTTGTAAAT
CAGCTCATTTTTTAACCAATAGGCCGAAATCGGCCAAAATCCCTTATAAATCAAAGA
ATAGACCGAGATAGGGTTGAGTGTGTTCCAGTTTGAACAAGAGTCCACTATTAA
AGAACGTGGACTCCAACGTCAAAGGGCGAAAAACCGTCTATCAGGGCGATGGCCC
ACTACGTGAACCATCACCTAATCAAGTTTTTTGGGGTCGAGGTGCCGTAAAGCAC
TAAATCGGAACCCTAAAGGGAGCCCCGATTTAGAGCTTGACGGGGAAAGCCGGC
GAACGTGGCGAGAAAGGAAGGGAAGAAAGCGAAAGGAGCGGGCGCTAGGGCGC
TGGCAAGTGTAGCGGTACGCTGCGCGTAACCACCACACCCGCCGCTTAATGC
GCCGCTACAGGGCGCGTCCCATTCCGCTTCCAGGCTGCGCAACTGTTGGGAAGG
GCGATCGGTGCGGGCCTCTTCGCTATTACGCCAGCTGGCGAAAGGGGGATGTGC
TGCAAGGCGATTAAGTTGGGTAACGCCAGGGTTTTCCAGTCACGACGTTGTA
CGACGGCCAGTGAGCGCGCGTAATACGACTCACTATAGGGCGAATTGGGTACCG
GCATATGGTTCTTGACAGAGGTGTAAAAAGTACTCAAAAATTTTACTCAAGTGAAAG
TACAAGTACTTAGGGAAAATTTTACTCAATTA
AAAAGTAAAAGTATCTGGCTAGAATC
TTACTTGAGTAAAAGTAAAAAAGTACTCCATTA
AAAATTGTA
CTTGAGTATTAAGGAA
GTAAAAGTAAAAGCAAGAAAGATCGATCTCGAAGGATCTGGAGGCCACCATGGTG
TCGATAACTTCGTATAGCATA
CATTATACGAAGTTATCGTGCTCGACATTGATTATT
GACTAGTTATTAATAGTAATCAATTACGGGGT
CATTAGTTCATAGCCCATATATGGA
GTTCCGCGTTACATAACTTACGGTAAATGGCCCGCCTGGCTGACCGCCCAACGAC
CCCCGCCCATTGACGTCAATAATGACGTATGTTCCCATAGTAACGCCAATAGGGAC
TTTCCATTGACGTCAATGGGTGGAGTATTTACGGTAAACTGCCCACTTGGCAGTAC
ATCAAGTGTATCATATGCCAAGTACGCCCCCTATTGACGTCAATGACGGTAAATGG
CCCGCCTGGCATTATGCCAGTACATGACCTTATGGGACTTTCCTACTTGGCAGTA
CATCTACGTATTAGTCATCGCTATTACCATGGT
CGAGGTGAGCCCCACGTTCTGCT
TCACTCTCCCCATCTCCCCCCCCCTCCCCACCCCAATTTTGTATTTATTTATTTTTTA
ATTATTTTGTGCAGCGATGGGGGCGGGGGGGGGGGGGGGGGGGCGCGCGCCAGGCG
GGGCGGGGCGGGGCGAGGGGCGGGGCGGGGCGAGGCGGAGAGGTGCGGGCGG
CAGCCAATCAGAGCGGCGCGCTCCGAAAGTTTCTTTTATGGCGAGGCGGGCGG
GGCGGCGGCCCTATAAAAAGCGAAGCGCGCGGGCGGGGAGTCGCTGCGAC
GCTGCCTTCGCCCGTGCCCCGCTCCGCGCGCCGCTCGCGCCGCCCGCCCCGG
CTCTGACTGACCGCGTTACTCCACAGGTGAGCGGGCGGGACGGCCCTTCTCCTC
CGGGCTGTAATTAGCGCTTGGTTTAATGACGGCTTGTTTCTTTTCTGTGGCTGCGT
GAAAGCCTTGAGGGGCTCCGGGAGGGCCCTTTGTGCGGGGGGAGCGGGCTCGGG
GGGTGCGTGCGTGTGTGTGCGTGGGGAGCGCCGCGTGCGGCTCCGCGCTGC
CCGGCGGCTGTGAGCGCTGCGGGCGCGGGCGGGGCTTTGTGCGCTCCGCAGT
GTGCGCGAGGGGAGCGCGGCCGGGGGGCGGTGCCCGCGGTGCGGGGGGGGCT

1283 GCGAGGGGAACAAAGGCTGCGTGCGGGGTGTGTGCGTGGGGGGGTGAGCAGGG
1284 GGTGTGGGCGCGTTCGGTTCGGGCTGCAACCCCCCTGCACCCCCCTCCCCGAGTT
1285 GCTGAGCACGGCCCCGGCTTCGGGTGCGGGGCTCCGTACGGGGCGTGGCGCGGG
1286 GCTCGCCGTGCCGGGCGGGGGGTGGCGGCAGGTGGGGGTGCCGGGCGGGGCG
1287 GGGCCGCCTCGGGCCGGGGAGGGCTCGGGGGAGGGGCGCGGGCGCCCCCGGA
1288 GCGCCGGCGGCTGTTCGAGGCGCGGGCAGCCGCAGCCATTGCCTTTTATGGTAAT
1289 CGTGCGAGAGGGCGCAGGGACTTCCTTTGTCCCAAATCTGTGCGGAGCCGAAATC
1290 TGGGAGGCGCCGCCGCACCCCCCTCTAGCGGGCGCGGGGCGAAGCGGTGCGGCG
1291 CCGGCAGGAAGGAAATGGGCGGGGAGGGCCTTCGTGCGTCCCGCGCCGCCGT
1292 CCCCTTCTCCCTCTCCAGCCTCGGGGCTGTCCGCGGGGGGACGGCTGCCTTCGG
1293 GGGGGACGGGGCAGGGCGGGGTTTCGGCTTCTGGCGTGTGACCGGGCGGCTCTAG
1294 AGCCTCTGCTAACCATGTTTCATGCCTTCTTCTTTTTCTACAGCTCCTGGGCAACGT
1295 GCTGGTTATTGTGCTGTCTCATCATTTTTGGCAAAGAATTCTCGAGCGGCCGCCAG
1296 TGTGATGGATATCGGATCCGCTAGCGCTACTAACTTCAGCCTGCTGAAGCAGGCT
1297 GGAGACGTGGAGGAGAACCCTGGACCTGGACCGGTCCGACCATGGTGAGCAAG
1298 GGCGAGGAGCTGTTACCGGGGTGGTGCCCATCCTGGTTCGAGCTGGACGGCGAC
1299 GTAAACGGCCACAAGTTCAGCGTGTCCGGCGAGGGCGAGGGCGATGCCACCTAC
1300 GGCAAGCTGACCCTGAAGTTCATCTGCACCACCGGCAAGCTGCCCGTGCCCTGG
1301 CCCACCCTCGTGACCACCCTGACCTACGGCGTGCAGTGCTTCAGCCGCTACCCCG
1302 ACCACATGAAGCAGCAGACTTCTTCAAGTCCGCCATGCCCGAAGGCTACGTCCA
1303 GGAGCGCACCATCTTCTTCAAGGACGACGGCAACTACAAGACCCGCGCCGAGGT
1304 GAAGTTCGAGGGCGACACCCTGGTGAACCGCATCGAGCTGAAGGGCATCGACTT
1305 CAAGGAGGACGGCAACATCCTGGGGCACAAGCTGGAGTACAACACTACAACAGCCAC
1306 AACGTCTATATCATGGCCGACAAGCAGAAGAACGGCATCAAGGTGAACTTCAAGAT
1307 CCGCCACAACATCGAGGACGGCAGCGTGCAGCTCGCCGACCACTACCAGCAGAA
1308 CACCCCATCGGCGACGGCCCCGTGCTGCTGCCCGACAACCACTACCTGAGCAC
1309 CCAGTCCGCCCTGAGCAAAGACCCCAACGAGAAGCGCGATCACATGGTCCTGCTG
1310 GAGTTCGTGACCGCCGCCGGGATCACTCTCGGCATGGACGAGCTGTACAAGTAAA
1311 GCGGCCGCCACCGCGGTGGAGCTCGAATTAATTCATCGATGATGATCCAGACATG
1312 ATAAGATACATTGATGAGTTTGGACAAACCACAACACTAGAATGCAGTGAAAAAATG
1313 CTTTATTTGTGAAATTTGTGATGCTATTGCTTTATTTGTAACCATTATAAGCTGCAAT
1314 AAACAAGTTAACAACAACAATTGCATTCATTTTATGTTTCAGGTTTCAGGGGGAGGTG
1315 TGGGAGGTTTTTTAAAGCAAGTAAACCTCTACAAATGTGGTATGGCTGATTATGAT
1316 CCTCTAGAGTCGGTGGGCCTCGGGGGCGGGTGCGGGGTCCGGCGGGGCGCCCC
1317 GGGTGGCTTCGGTTCGGAGCCATGGGGTTCGTGCGCTCCTTTTCGGTCCGGCGCTGC
1318 GGGTCGTGGGGCGGGCGTTCAGGCACCGGGCTTTCGGGGTCATGCACCAGGTGCG
1319 CGGTCTTCGGGCACCTCGACGTCCGGCGGTGACGGTGAAGCCGAGCCGCTCGTA
1320 GAAGGGGAGGTTGCGGGGCGCGGAGGTCTCCAGGAAGGCGGGCACCCCGGCGC
1321 GCTCGGCCGCCTCCACTCCGGGGAGCACGACGGCGCTGCCAGACCCTTGCCCT
1322 GGTGGTCCGGGCGAGACGCCGACGGTGGCCAGGAACCACGCGGGGCTCCTTGGGC

1323 CGGTGCGGCGCCAGGAGGCCTTCCATCTGTTGCTGCGCGGCCAGCCGGGAACCG
1324 CTCAACTCGGCCATGCGCGGGCCGATCTCGGCGAACACCGCCCCCGCTTCGACG
1325 CTCTCCGGCGTGGTCCAGACCGCCACCGCGGGCGCCGTCGTCCGCGACCCACACC
1326 TTGCCGATGTCGAGCCCAGCGCGCGTGAGGAAGAGTTCTTGCAGCTCGGTGACC
1327 CGCTCGATGTGGCGGTCCGGGTCGACGGTGTGGCGCGTGGCGGGGTAGTCGGC
1328 GAACGCGGGCGGGCAGGGGTGCGTACGGCCCCGGGGGACGTCGTCGCGGGGTGGCGA
1329 GGCGCACCGTGGGCTTGTACTCGGTCATGGAAGGTCGTCTCCTTGTGAGGGGTCA
1330 GGGGCGTGGGTCAGGGGATGGTGGCGGCACCGGTCGTGGCGGCCGACCTGCAG
1331 GCATGCAAGCTTTTTGCAAAAGCCTAGGCCTCCAAAAAGCCTCCTCACTACTTCT
1332 GGAATAGCTCAGAGGCCGAGGCGGCCTCGGCCTCTGCATAAATAAAAAAATTAG
1333 TCAGCCATGGGGCGGAGAATGGGCGGAACTGGGCGGAGTTAGGGGCGGGATGG
1334 GCGGAGTTAGGGGCGGGACTATGGTTGCTGACTAATTGAGATGCATGCTTTGCAT
1335 ACTTCTGCCTGCTGGGGAGCCTGGGGACTTTCCACACCTGGTTGCTGACTAATTG
1336 AGATGCATGCTTTGCATACTTCTGCCTGCTGGGGAGCCTGGGGACTTTCCACACC
1337 CTAAGTGCACACATTCCACAGAATCAAGTGATCTCCAAAAATAAGTACTTTTTG
1338 ACTGTAAATAAAATTGTAAGGAGTAAAAAGTACTTTTTTTTTCTAAAAAATGTAATTA
1339 AGTAAAAGTAAAAGTATTGATTTTTAATTGTAAGTAAAAGTAAAATCCCCAAAA
1340 ATAATACTTAAGTACAGTAATCAAGTAAAATTACTCAAGTACTTTACACCTCTGGTTC
1341 TTGACCCCCTACCTTCAGCAAGCCCAGCAGATCCGAGCTCCAGCTTTTGTTCCTT
1342 TAGTGAGGGTTAATTGCGCGCTTGCGGTAATCATGGTCATAGCTGTTTCCTGTGTG
1343 AAATTGTTATCCGCTCACAATTCCACACAACATACGAGCCGGAAGCATAAAGTGTA
1344 AAGCCTGGGGTGCCTAATGAGTGAGCTAACTCACATTAATTGCGTTGCGCTCACTG
1345 CCCGCTTTCCAGTCGGGAAACCTGTGCGTCCAGCTGCATTAATGAATCGGCCAAC
1346 GCGCGGGGAGAGGCGGTTTGCATTTGGGCGCTTTCCGCTTCTCGCTCACTGA
1347 CTGCTGCGCTCGGTGCTTCCGCTGCGGGCAGCGGTATCAGCTCACTCAAAGGC
1348 GGTAATACGGTTATCCACAGAATCAGGGGATAACGCAGGAAAGAACATGTGAGCA
1349 AAAGGCCAGCAAAAGGCCAGGAACCGTAAAAAGGCCGCGTTGCTGGCGTTTTTCC
1350 ATAGGCTCCGCCCCCTGACGAGCATCACAAAATCGACGCTCAAGTCAGAGGTG
1351 GCGAAACCCGACAGGACTATAAAGATAACCAGGCGTTTCCCCCTGGAAGCTCCCTC
1352 GTGCGCTCTCCTGTTCCGACCCTGCCGCTTACCGGATACCTGTCCGCCTTTCTCC
1353 CTTCCGGGAAGCGTGGCGCTTTCTCATAGCTCACGCTGTAGGTATCTCAGTTCGGT
1354 GTAGGTCGTTCCGCTCCAAGCTGGGCTGTGTGCACGAACCCCCCGTTCAGCCCGAC
1355 CGCTGCGCCTTATCCGGTAACTATCGTCTTGAGTCCAACCCGGTAAGACACGACTT
1356 ATCGCCACTGGCAGCAGCCACTGGTAACAGGATTAGCAGAGCGAGGTATGTAGGC
1357 GGTGCTACAGAGTTCTTGAAGTGGTGGCCTAACTACGGCTACACTAGAAGGACAG
1358 TATTTGGTATCTGCGCTCTGCTGAAGCCAGTTACCTTCGGAAAAAGAGTTGGTAGC
1359 TCTTGATCCGGCAAACAAACCACCGCTGGTAGCGGTGGTTTTTTTTGTTTGCAAGCA
1360 GCAGATTACGCGCAGAAAAAAGGATCTCAAGAAGATCCTTTGATCTTTTCTACGG
1361 GGTCTGACGCTCAGTGGAACGAAAACCTCACGTTAAGGGATTTTGGTCATGAGATTA
1362 TCAAAAAGGATCTTCACCTAGATCCTTTTAAATTAATAAATGAAGTTTTAAATCAATCT

1363 AAAGTATATATGAGTAAACTTGGTCTGACAGTTACCAATGCTTAATCAGTGAGGCAC
1364 CTATCTCAGCGATCTGTCTATTTTCGTTTCATCCATAGTTGCCTGACTCCCCGTCGTGT
1365 AGATAACTACGATACGGGAGGGCTTACCATCTGGCCCCAGTGCTGCAATGATACC
1366 GCGAGACCCACGCTCACCGGCTCCAGATTTATCAGCAATAAACCAGCCAGCCGGA
1367 AGGGCCGAGCGCAGAAGTGGTCCTGCAACTTTATCCGCCTCCATCCAGTCTATTA
1368 ATTGTTGCCGGGAAGCTAGAGTAAGTAGTTCGCCAGTTAATAGTTTGGCAACGTT
1369 GTTGCCATTGCTACAGGCATCGTGGTGTACGCTCGTCGTTTGGTATGGCTTCATT
1370 CAGCTCCGGTTCCCAACGATCAAGGCGAGTTACATGATCCCCCATGTTGTGCAAAA
1371 AAGCGGTTAGCTCCTTCGGTCCCTCCGATCGTTGTCAGAAGTAAGTTGGCCGCAGT
1372 GTTATCACTCATGGTTATGGCAGCACTGCATAATTCTCTTACTGTCATGCCATCCGT
1373 AAGATGCTTTTTCTGTGACTGGTGAGTACTCAACCAAGTCATTCTGAGAATAGTGTAT
1374 GCGGCGACCGAGTTGCTCTTGCCCGGCGTCAATACGGGATAATACCGCGCCACAT
1375 AGCAGAACTTTAAAAGTGCTCATCATTGGAAAACGTTCTTCGGGGCGAAAACCTC
1376 AAGGATCTTACCGCTGTTGAGATCCAGTTCGATGTAACCCACTCGTGCACCCAACT
1377 GATCTTCAGCATCTTTTACTTTACCAGCGTTTCTGGGTGAGCAAAAACAGGAAGG
1378 CAAAATGCCGCAAAAAAGGGAATAAGGGCGACACGGAAATGTTGAATACTCATACT
1379 CTTCTTTTTCAATATTATTGAAGCATTATCAGGGTTATTGTCTCATGAGCGGATAC
1380 ATATTTGAATGTATTTAGAAAAATAAACAAATAGGGGTTCGCGGCACATTTCCCCGA
1381 AAAGTG
1382
1383

1384 **LDLRwt**
1385 ATGGGGCCCTGGGGCTGGAAATTGCGCTGGACCGTCGCCTTGCTCCTCGCCGCG
1386 GCGGGGACTGCAGTGGGCGACAGATGCGAAAGAAACGAGTTCCAGTGCCAAGAC
1387 GGGAAATGCATCTCCTACAAGTGGGTCTGCGATGGCAGCGCTGAGTGCCAGGATG
1388 GCTCTGATGAGTCCCAGGAGACGTGCTTGTCTGTACCTGCAAATCCGGGGACTT
1389 CAGCTGTGGGGGCCGTGTCAACCGCTGCATTCCTCAGTTCTGGAGGTGCGATGGC
1390 CAAGTGGACTGCGACAACGGCTCAGACGAGCAAGGCTGTCCCCCAAGACGTGC
1391 TCCCAGGACGAGTTTCGCTGCCACGATGGGAAGTGCATCTCTCGGCAGTTCGTCT
1392 GTGACTCAGACCGGGACTGCTTGGACGGCTCAGACGAGGCCTCCTGCCCGGTGC
1393 TCACCTGTGGTCCC GCCAGCTTCCAGTGCAACAGCTCCACCTGCATCCCCAGCT
1394 GTGGGCCTGCGACAACGACCCCGACTGCGAAGATGGCTCGGATGAGTGGCCGCA
1395 GCGCTGTAGGGGTCTTTACGTGTTCCAAGGGGACAGTAGCCCCTGCTCGGCCTTC
1396 GAGTTCCACTGCCTAAGTGGCGAGTGCATCCACTCCAGCTGGCGCTGTGATGGTG
1397 GCCCGACTGCAAGGACAAATCTGACGAGGAAAACACTGCGCTGTGGCCACCTGTGC
1398 CCCTGACGAATTCCAGTGCTCTGATGGAAACTGCATCCATGGCAGCCGGCAGTGT
1399 GACCGGGAATATGACTGCAAGGACATGAGCGATGAAGTTGGCTGCGTTAATGTGA
1400 CACTCTGCGAGGGACCCAACAAGTTCAAGTGTACAGCGGGCGAATGCATCACCCT
1401 GGACAAAGTCTGCAACATGGCTAGAGACTGCCGGGACTGGTCAGATGAACCCATC
1402 AAAGAGTGCGGGACCAACGAATGCTTGGACAACAACGGCGGCTGTTCCCACGTCT
1403 GCAATGACCTTAAGATCGGCTACGAGTGCCTGTGCCCGACGGCTTCCAGCTGGT
1404 GGCCCAGCGAAGATGCGAAGATATCGATGAGTGTGAGGATCCCGACACCTGCAGC
1405 CAGCTCTGCGTGAACCTGGAGGGTGGCTACAAGTGCCAGTGTGAGGAAGGCTTC
1406 CAGCTGGACCCCCACACGAAGGCCTGCAAGGCTGTGGGCTCCATCGCCTACCTCT
1407 TCTTCACCAACCGGCACGAGGTCAGGAAGATGACGCTGGACCGGAGCGAGTACA
1408 CCAGCCTCATCCCCAACCTGAGGAACGTGGTCGCTCTGGACACGGAGGTGGCCA
1409 GCAATAGAATCTACTGGTCTGACCTGTCCCAGAGAATGATCTGCAGCACCCAGCTT
1410 GACAGAGCCCACGGCGTCTCTTCCTATGACACCGTCATCAGCAGAGACATCCAGG
1411 CCCCCGACGGGCTGGCTGTGGACTGGATCCACAGCAACATCTACTGGACCGACTC
1412 TGTCCTGGGCACTGTCTCTGTTGCGGATACCAAGGGCGTGAAGAGGAAAACGTTA
1413 TTCAGGGAGAACGGCTCCAAGCCAAGGGCCATCGTGGTGGATCCTGTTTCATGGCT
1414 TCATGTAAGTGGACTGACTGGGGAACCTCCCGCCAAGATCAAGAAAGGGGGCCTGAA
1415 TGGTGTGGACATCTACTCGCTGGTACTGAAAACATTTCAGTGGCCCAATGGCATCA
1416 CCCTAGATCTCCTCAGTGGCCGCCTCTACTGGGTTGACTCCAAACTTCACTCCATC
1417 TCAAGCATCGATGTCAATGGGGGCAACCGGAAGACCATCTTGGAGGATGAAAAGA
1418 GGCTGGCCCACCCCTTCTCCTTGGCCGTCTTTGAGGACAAAGTATTTTGGACAGAT
1419 ATCATCAACGAAGCCATTTTCAGTGCCAACCGCCTCACAGGTTCCGATGTCAACTT
1420 GTTGGCTGAAAACCTACTGTCCCCAGAGGATATGGTCCTCTTCCACAACCTCACCC
1421 AGCCAAGAGGAGTGAACCTGGTGTGAGAGGACCACCCTGAGCAATGGCGGCTGCC
1422 AGTATCTGTGCCTCCCTGCCCGCAGATCAACCCCACTCGCCCAAGTTTACCTG
1423 CGCCTGCCCGGACGGCATGCTGCTGGCCAGGGACATGAGGAGCTGCCTCACAGA
1424 GGCTGAGGCTGCAGTGGCCACCCAGGAGACATCCACCGTCAGGCTAAAGGTGAG

1425 CTCCACAGCCGTAAGGACACAGCACACAACCACCCGGCCTGTTCCCGACACCTCC
1426 CGGCTGCCTGGGGCCACCCCTGGGCTCACCACGGTGGAGATAGTGACAATGTCT
1427 CACCAAGCTCTGGGCGACGTTGCTGGCAGAGGAAATGAGAAGAAGCCCAGTAGC
1428 GTGAGGGCTCTGTCCATTGTCTCCCCATCGTGCTCCTCGTCTTCTTTGCCTGGG
1429 GGTCTTCTTCTATGGAAGAAGTGGCGGCTTAAGAACATCAACAGCATCAACTTTG
1430 ACAACCCCGTCTATCAGAAGACCACAGAGGATGAGGTCCACATTTGCCACAACCA
1431 GGACGGCTACAGCTACCCCTCGAGACAGATGGTCAGTCTGGAGGATGACGTGGC
1432 G

1433

1434 **LDLRDup252 with surrounding region**

1435 CCCCCAAGACGTGCTCCCAGGACGAGTTTCGCTGCCACGATGGGAAGTGCATCTC
1436 TCGGCAGTTCGTCTGTGACTCAGACCGGGACTGCTTGGACGGCTCAGACGAGGC
1437 CTCCTGCCCGGTGCTCACCTGTGGTCCCGCCAGCTTCCAGTGCAACAGCTCCACC
1438 TGCATCCCCCAGCTGTGGGCCTGCGACAACGACCCCGACTGCGAAGATGGCTCG
1439 GAGGCTCGGATGAGTGGCCGCAGCGCTGTAGGGGTCTTTACGTGTTCCAAGGGG
1440 ACAGTAGCCCCTGCTCGGCCTTCGAGTTCCACTGCCTAAGTGGCGAGTGCATCCA
1441 CTCCAGCTGGCGCTGTGATGGTGGCCCCGACTGCAAGGACAAATCTGACGAGGA
1442 AACTGCG

1443

1444 **LDLRDup254/255 with surrounding region**

1445 CCCCCAAGACGTGCTCCCAGGACGAGTTTCGCTGCCACGATGGGAAGTGCATCTC
1446 TCGGCAGTTCGTCTGTGACTCAGACCGGGACTGCTTGGACGGCTCAGACGAGGC
1447 CTCCTGCCCGGTGCTCACCTGTGGTCCCGCCAGCTTCCAGTGCAACAGCTCCACC
1448 TGCATCCCCCAGCTGTGGGCCTGCGACAACGACCCCGACTGCGAAGATGGCTCG
1449 GATGAGTGGCCGCAGCGCTGTAGGGGTCTTTACGTGTTCCAAGGGGACAGTAGC
1450 CCCTGCTCGGCCTTCGAGTTCCACTGCCTAAGTGGCGAGTGCATCCACTCCAGCT
1451 GGCGCTGTGATGGTGGCCCCGACTGCAAGGACAAATCTGACAGGACAAATCTGAC
1452 GAGGAAAAGTGCAGCTGTGGCCACCTGTCGCCCTGACGAATTCCAGTGCTCTGATG
1453 GAAACTGCATCCATG

1454

1455 **LDLRDup258 with surrounding region**

1456 CCCCCAAGACGTGCTCCCAGGACGAGTTTCGCTGCCACGATGGGAAGTGCATCTC
1457 TCGGCAGTTCGTCTGTGACTCAGACCGGGACTGCTTGGACGGCTCAGACGAGGC
1458 CTCCTGCCCGGTGCTCACCTGTGGTCCCGCCAGCTTCCAGTGCAACAGCTCCACC
1459 TGCATCCCCCAGCTGTGGGCCTGCGACAACGACCCCGACTGCGAAGATGGCTCG
1460 GATGAGTGGCCGCAGCGCTGTAGGGGTCTTTACGTGTTCCAAGGGGACAGTAGC
1461 CCCTGCTCGGCCTTCGAGTTCCACTGCCTAAGTGGCGAGTGCATCCACTCCAGCT
1462 GGCGCTGTGATGGTGGCCCCGACTGCAAGGACAAATCTGAGGACAAATCTGACGA
1463 GGAAAAGTGCAGCTGTGGCCACCTGTCGCCCTGACGAATTCCAGTGCTCTGATGGA
1464 AACTGCATCCATG

1465

1466 **LDLRDup261 with surrounding region**

1467 CCCCCAAGACGTGCTCCCAGGACGAGTTTCGCTGCCACGATGGGAAGTGCATCTC

1468 TCGGCAGTTCGTCTGTGACTCAGACCGGGACTGCTTGGACGGCTCAGACGAGGC

1469 CTCCTGCCCGGTGCTCACCTGTGGTCCCGCCAGCTTCCAGTGCAACAGCTCCACC

1470 TGCATCCCCCAGCTGTGGGCCTGCGACAACGACCCCGACTGCGAAGATGGCTCG

1471 GATGAGTGGCCGCAGCGCTGTAGGGGTCTTTACGTGTTCCAAGGGGACAGTAGC

1472 CCCTGCTCGGCCTTCGAGTTCCACTGCCTAAGTGGCGAGTGCATCCACTCCAGCT

1473 GGCGCTGTGATGGTGGCCCCGACTGCAAGGACAAATCTGACGACAAATCTGACGA

1474 GGAAAAGTGCCTGTGGCCACCTGTCGCCCTGACGAATTCCAGTGCTCTGATGGA

1475 AACTGCATCCATG

1476

1477 **LDLRDup264 with surrounding region**

1478 CTTTCATGTAAGTGGACTGACTGGGGAACTCCCGCCAAGATCAAGAAAGGGGGCCTG

1479 AATGGTGTGGACATCTACTCGCTGGTGGAGCTGGTGGACTGAAAACATTCAGTGGCC

1480 CAATGGCATCACCTAG

1481

1482

1483 **GAAwT**
1484 ATGGGAGTGAGGCACCCGCCCTGCTCCCACCGGCTCCTGGCCGTCTGCGCCCTC
1485 GTGTCCTTGGCAACCGCTGCACTCCTGGGGCACATCCTACTCCATGATTTCCCTGCT
1486 GGTTCCCCGAGAGCTGAGTGGCTCCTCCCCAGTCCTGGAGGAGACTCACCCAGCT
1487 CACCAGCAGGGAGCCAGCAGACCAGGGCCCCGGGATGCCCAGGCACACCCCGG
1488 CCGTCCCAGAGCAGTGCCCACACAGTGCGACGTCCCCCCCCAACAGCCGCTTCGA
1489 TTGCGCCCCTGACAAGGCCATCACCCAGGAACAGTGCGAGGCCCGCGGCTGTTG
1490 CTACATCCCTGCAAAGCAGGGGCTGCAGGGAGCCCAGATGGGGCAGCCCTGGTG
1491 CTTCTTCCCACCCAGCTACCCAGCTACAAGCTGGAGAACCTGAGCTCCTCTGAAA
1492 TGGGCTACACGGCCACCCTGACCCGTACCACCCCCACCTTCTTCCCCAAGGACAT
1493 CCTGACCCTGCGGCTGGACGTGATGATGGAGACTGAGAACCGCCTCCACTTCACG
1494 ATCAAAGATCCAGCTAACAGGCGCTACGAGGTGCCCTTGGAGACCCCGCATGTCC
1495 ACAGCCGGGCACCGTCCCCACTCTACAGCGTGGAGTTCTCCGAGGAGCCCTTCG
1496 GGGTGATCGTGCGCCGGCAGCTGGACGGCCGCGTGCTGCTGAACACGACGGTG
1497 GCGCCCCTGTTCTTTGCGGACCAGTTCCTTCAGCTGTCCACCTCGCTGCCCTCGC
1498 AGTATATCACAGGCCTCGCCGAGCACCTCAGTCCCCTGATGCTCAGCACCAGCTG
1499 GACCAGGATCACCTGTGGAACCGGGACCTTGCGCCACGCCCCGGTGCGAACCT
1500 CTACGGGTCTCACCTTTCTACCTGGCGCTGGAGGACGGCGGGTTCGGCACACGG
1501 GGTGTTCCCTGCTAAACAGCAATGCCATGGATGTGGTCCTGCAGCCGAGCCCTGCC
1502 CTTAGCTGGAGGTGACAGGTGGGATCCTGGATGTCTACATCTTCCCTGGGCCAG
1503 AGCCCAAGAGCGTGGTGCAGCAGTACCTGGACGTTGTGGGATAACCGTTCATGCC
1504 GCCATACTGGGGCCTGGGCTTCCACCTGTGCCGCTGGGGCTACTCCTCCACCGCT
1505 ATCACCCGCCAGGTGGTGGAGAACATGACCAGGGGCCACTTCCCCCTGGACGTC
1506 CAGTGGAACGACCTGGACTACATGGACTCCCGGAGGGACTTCACGTTCAACAAGG
1507 ATGGCTTCCGGGACTTCCCGGCCATGGTGCAGGAGCTGCACCAGGGCGGGCCGGC
1508 GCTACATGATGATCGTGGATCCTGCCATCAGCAGCTCGGGCCCTGCCGGGAGCTA
1509 CAGGCCCTACGACGAGGGTCTGCGGAGGGGGGTTTTTCATCACCAACGAGACCGG
1510 CCAGCCGCTGATTGGGAAGGTATGGCCCGGGTCCACTGCCTTCCCCGACTTCACC
1511 AACCCACAGCCCTGGCCTGGTGGGAGGACATGGTGGCTGAGTTCCATGACCAG
1512 GTGCCCTTCGACGGCATGTGGATTGACATGAACGAGCCTTCCAACCTTCATCAGGG
1513 GCTCTGAGGACGGCTGCCCAACAATGAGCTGGAGAACCCACCCTACGTGCCTG
1514 GGGTGGTTGGGGGGACCCTCCAGGCGGCCACCATCTGTGCCTCCAGCCACCAGT
1515 TTCTCTCCACACACTACAACCTGCACAACCTCTACGGCCTGACCGAAGCCATCGCC
1516 TCCACAGGGCGCTGGTGAAGGCTCGGGGGACACGCCCATTTGTGATCTCCCGC
1517 TCGACCTTTGCTGGCCACGGCCGATACGCCGGCCACTGGACGGGGGACGTGTGG
1518 AGCTCCTGGGAGCAGCTCGCCTCCTCCGTGCCAGAAATCCTGCAGTTTAACCTGC
1519 TGGGGGTGCCTCTGGTTCGGGGCCGACGTCTGCGGCTTCCCTGGGCAACACCTCAG
1520 AGGAGCTGTGTGTGCGCTGGACCCAGCTGGGGGCCTTCTACCCCTTCATGCGGAA
1521 CCACAACAGCCTGCTCAGTCTGCCCCAGGAGCCGTACAGCTTCAGCGAGCCGGC
1522 CCAGCAGGCCATGAGGAAGGCCCTCACCTGCGCTACGCACTCCTCCCCACCT

1523 CTACACACTGTTCCACCAGGCCACGTCGCGGGGGAGACCGTGGCCCGGCCCT
1524 CTTCTGGAGTTCCCAAGGACTCTAGCACCTGGACTGTGGACCACCAGCTCCTG
1525 TGGGGGGAGGCCCTGCTCATACCCAGTGCTCCAGGCCGGAAGGCCGAAGTG
1526 ACTGGCTACTTCCCCTTGGGCACATGGTACGACCTGCAGACGGTGCCAGTAGAGG
1527 CCCTTGGCAGCCTCCCACCCACCTGCAGCTCCCCGTGAGCCAGCCATCCACAG
1528 CGAGGGGCAGTGGGTGACGCTGCCGGCCCCCTGGACACCATCAACGTCCACCT
1529 CCGGGCTGGGTACATCATCCCCCTGCAGGGCCCTGGCCTCACAACCACAGAGTC
1530 CCGCCAGCAGCCCATGGCCCTGGCTGTGGCCCTGACCAAGGGTGGGGAGGCC
1531 GAGGGGAGCTGTTCTGGGACGATGGAGAGAGCCTGGAAGTGCTGGAGCGAGGG
1532 GCCTACACACAGGTCATCTTCTGGCCAGGAATAACACGATCGTGAATGAGCTGG
1533 TACGTGTGACCAGTGAGGGAGCTGGCCTGCAGCTGCAGAAGGTGACTGTCCTGG
1534 GCGTGGCCACGGCGCCCCAGCAGGTCCTCTCCAACGGTGTCCCTGTCTCCA
1535 CACTACAGCCCCGACACCAAGGTCCTGGACATCTGTGTCTCGCTGTTGATGGGA
1536 GAGCAGTTTCTCGTCAGCTGGTGT

1537

1538 **GAADup327/328**

1539 ATGGGAGTGAGGCACCCGCCCTGCTCCCACCGGCTCCTGGCCGTCTGCGCCCTC
1540 GTGTCCTTGGCAACCGCTGCACTCCTGGGGCACATCCTACTCCATGATTTCTGCT
1541 GGTTCCCCGAGAGCTGAGTGGCTCCTCCCCAGTCCTGGAGGAGACTCACCCAGCT
1542 CACCAGCAGGGAGCCAGCAGACCAGGGCCCCGGGATGCCCAGGCACACCCCGG
1543 CCGTCCCAGAGCAGTGCCACACAGTGCGACGTCCCCCCCAACAGCCGCTTCGA
1544 TTGCGCCCCTGACAAGGCCATCACCCAGGAACAGTGCGAGGCCCGCGGCTGTTG
1545 CTACATCCCTGCAAAGCAGGGGCTGCAGGGAGCCAGATGGGGCAGCCCTGGTG
1546 CTTCTTCCCACCCAGCTACCCAGCTACAAGCTGGAGAACCCTGAGCTCCTCTGAAA
1547 TGGGCTACACGGCCACCCTGACCCGTACCACCCACCTTCTTCCCCAAGGACAT
1548 CCTGACCCTGCGGCTGGACGTGATGATGGAGACTGAGAACCGCCTCCACTTCACG
1549 ATCAAAGATCCAGCTAACAGGCGCTACGAGGTGCCCTTGGAGACCCCGCATGTCC
1550 ACAGCCGGGCACCGTCCCCACTCTACAGCGTGGAGTTCTCCGAGGAGCCCTTCG
1551 GGGTGATCGTGCGCCGGCAGCTGGACGGCCGCGTGCTGCTGAACACGACGGTG
1552 GCGCCCCTGTTCTTTGCGGACCAGTTCCTTCAGCTGTCCACCTCGCTGCCCTCGC
1553 AGTATATCACAGGCCTCGCCGAGCACCTCAGTCCCCTGATGCTCAGCACCAGCTG
1554 GACCAGGATCACCTGTGGAACCGGGACCTTGCGCCCACGCCCGGTGCGAACCT
1555 CTACGGGTCTCACCTTTCTACCTGGCGCTGGAGGACGGCGGGTCGGCACACGG
1556 GGTGTTCTGCTAAACAGCAATGCCATGGATGTGGTCCTGCAGCCGAGCCCTGCC
1557 CTTAGCTGGAGGTCGACAGGTGGGATCCTGGATGTCTACATCTTCTGGGCCAG
1558 AGCCCAAGAGCGTGGTGCAGCAGTACCTGGACGTTGTGGGATACCCGTTTCATGCC
1559 GCCATACTGGGGCCTGGGCTTCCACCTGTGCCGCTGGGGCTACTCCTCCACCGCT
1560 ATCACCCGCCAGGTGGTGGAGAACATGACCAGGGCCCACTTCCCCCTGGACGTC
1561 CAGTGGAACGACCTGGACTACATGGACTCCCGGAGGGACTTCACGTTCAACAAGG
1562 ATGGCTTCCGGGACTTCCCGGCCATGGTGCAGGAGCTGCACCAGGGCGGGCCGGC

1563 GCTACATGATGATCGTGGATCCTGCCATCAGCAGCTCGGGCCCTGCCGGGAGCTA
1564 CAGGCCCTACGACGAGGGTCTGCGGAGGGGGGTTTTTCATCACCAACGAGACCGG
1565 CCAGCCGCTGATTGGGAAGGTATGGCCCGGGTCCACTGCCTTCCCCGACTTCACC
1566 AACCCACAGCCCTGGCCTGGTGGGAGGACATGGTGGCTGAGTTCCATGACCAG
1567 GTGCCCTTCGACGGCATGTGGATTGACATGAACGAGCCTTCCAAC TTCATCAGGG
1568 GCTCTGAGGACGGCTGCCCAACAATGAGCTGGAGAACCACCCCTACGTGCCTG
1569 GGGTGGTTGGGGGGACCCTCCAGGCGGCCACCATCTGTGCCTCCAGCCACCAGT
1570 TTCTCTCCACACACTACAACCTGCACAACCTCTACGGCCTGACCGAAGCCATCGCC
1571 TCCACAGGGCGCTGGTGAAGGCTCGGGGGACACGCCCATTTGTGATCTCCCGC
1572 TCGACCTTTGCTGGCCACGGCCGATACGCCGGCCACTGGACGGGGGACGTGTGG
1573 AGCTCCTGGGAGCAGCTCGCCTCCTCCGTGCCAGAAATCCTGCAGTTTAACCTGC
1574 TGGGGGTGCCTCTGGTCCGGGGCCGACGTCTGCGGCTTCCCTGGGCAACACCTCAG
1575 AGGAGCTGTGTGTGCGCTGGACCCAGCTGGGGGCCTTCTACCCCTTCATGCGGAA
1576 CCACAACAGCCTGCTCAGTCTGCCCCAGGAGCCGTACAGCTTCAGCGAGCCGGC
1577 CCAGCAGGCCATGAGGAAGGCCCTCACCTGCGCTACGCACTCCTCCCCACCT
1578 CTACACACTGTTCCACCAGGCCACGTCCGGGGGAGACCGTGGCCCGGCCCT
1579 CTTCTGGAGTTCCCCAAGGACTCTAGCACCTGGACTGTGGACCACCAGCTCCTG
1580 TGGGGGGAGGCCCTGCTCATCACCCAGTGCTCCAGGCCGGAAGGCCGAAGTG
1581 ACTGGCTACTTCCCCTTGGGCACATGGTACGACCTGCAGACGGTGCCAGTAGAGG
1582 CCCTTGGCAGCCTCCCACCCCCACCTGCAGCTCCCCGTGAGCCAGCCATCCACAG
1583 CGAGGGGCAGTGGGTGACGCTGCCGGCCCCCTGGACACCATCAACGTCCACCT
1584 CCGGGCTGGGTACATCATCCCCCTGCAGGGCCCTGGCCTCACAACCACAGAGTC
1585 CCGCCAGCAGCCCATGGCCCTGGCTGTGGCCCTGACCAAGGGTGGGGAGGCC
1586 GAGGGGAGCTGTTCTGGGACGATGGAGAGAGCCTGGAAGTGCTGGAGCGAGGG
1587 GCCTACACACAGGTCATCTTCTGGCCAGGAATAACACGATCGTGAATGAGCTGG
1588 TACGTGTGACCAGTGAGGGAGCTGGCCTGCAGCTGCAGAAGGTGACTGCAGAAG
1589 GTGACTGTCCTGGGCGTGGCCACGGCGCCCCAGCAGGTCCTCTCCAACGGTGTC
1590 CCTGTCTCCAACCTCACCTACAGCCCCGACACCAAGGTCCTGGACATCTGTGTCTC
1591 GCTGTTGATGGGAGAGCAGTTTCTCGTCAGCTGGTGT

1592

1593

1594 **GLB1wt**

1595 ATGCCGGGGTTCTGTTTCGCATCCTCCCTCTGTTGCTGGTTCTGCTGCTTCTGG
1596 GCCCTACGCGCGGCTTGCGCAATGCCACCCAGAGGATGTTTCAAATTGACTATAG
1597 CCGGGACTCCTTCTCAAGGATGGCCAGCCATTTTCGCTACATCTCAGGAAGCATTC
1598 ACTACTCCCGTGTGCCCCGCTTCTACTGGAAGGACCGGCTGCTGAAGATGAAGAT
1599 GGCTGGGCTGAACGCCATCCAGACGTATGTGCCCTGGAAC TTTTCATGAGCCCTGG
1600 CCAGGACAGTACCAGTTTTCTGAGGACCATGATGTGGAATATTTTCTTCGGCTGGC
1601 TCATGAGCTGGGACTGCTGGTTATCCTGAGGCCCGGGCCCTACATCTGTGCAGAG
1602 TGGGAAATGGGAGGATTACCTGCTTGGCTGCTAGAGAAAGAGTCTATTCTTCTCCG

1603 CTCCTCCGACCCAGATTACCTGGCAGCTGTGGACAAGTGGTTGGGAGTCCTTCTG
1604 CCCAAGATGAAGCCTCTCCTCTATCAGAATGGAGGGCCAGTTATAACAGTGCAGG
1605 TTGAAAATGAATATGGCAGCTACTTTGCCTGTGATTTTGACTACCTGCGCTTCCTGC
1606 AGAAGCGCTTTCGCCACCATCTGGGGGATGATGTGGTTCTGTTTACCACTGATGGA
1607 GCACATAAAACATTCTGAAATGTGGGGCCCTGCAGGGCCTCTACACCACGGTGG
1608 ACTTTGGAACAGGCAGCAACATCACAGATGCTTTCCTAAGCCAGAGGAAGTGTGA
1609 GCCCAAAGGACCCTTGATCAATTCTGAATTCTATACTGGCTGGCTAGATCACTGGG
1610 GCCAACCTCACTCCACAATCAAGACCGAAGCAGTGGCTTCCTCCCTCTATGATATA
1611 CTTGCCCGTGGGGCGAGTGTGAACTTGTACATGTTTATAGGTGGGACCAATTTTGC
1612 CTATTGGAATGGGGCCAACTCACCTATGCAGCACAGCCCACCAGCTACGACTAT
1613 GATGCCCCACTGAGTGAGGCTGGGGACCTCACTGAGAAGTATTTTGTCTCTGCGAA
1614 ACATCATCCAGAAGTTTGAAAAAGTACCAGAAGGTCCTATCCCTCCATCTACACCA
1615 AAGTTTGCATATGGAAAGGTCACCTTTGGAAAAGTTAAAGACAGTGGGAGCAGCTCT
1616 GGACATTCTGTGTCCCTCTGGGCCCATCAAAGCCTTTATCCCTTGACATTTATCCA
1617 GGTGAAACAGCATTATGGGTTTGTGCTGTACCGGACAACACTTCCTCAAGATTGCA
1618 GCAACCCAGCACCTCTCTCTTACCCCTCAATGGAGTCCACGATCGAGCATATGTT
1619 GCTGTGGATGGGATCCCCCAGGGAGTCCTTGAGCGAAACAATGTGATCACTCTGA
1620 ACATAACAGGGAAAGCTGGAGCCACTCTGGACCTTCTGGTAGAGAACATGGGACG
1621 TGTGAACTATGGTGCATATATCAACGATTTTAAGGGTTTGGTTTCTAACCTGACTCT
1622 CAGTTCCAATATCCTCACGGACTGGACGATCTTTCCTACTGGACACTGAGGATGCAG
1623 TGTGCAGCCACCTGGGGGGCTGGGGACACCGTGACAGTGGCCACCATGATGAAG
1624 CCTGGGCCCACTCACTCACTACACGCTCCCGGCCTTTTATATGGGGAACTTC
1625 TCCATTCCCAGTGGGATCCAGACTTGCCCCAGGACACCTTTATCCAGTTTCCTGG
1626 ATGGACCAAGGGCCAGGTCTGGATTAATGGCTTTAACCTTGGCCGCTATTGGCCA
1627 GCCCGGGGCCCTCAGTTGACCTTGTTTGTGCCCCAGCACATCCTGATGACCTCGG
1628 CCCCAAACACCATCACCGTGCTGGAAGTGGAGTGGGCACCCTGCAGCAGTGATGA
1629 TCCAGAACTATGTGCTGTGACGTTTCGTGGACAGGCCAGTTATTGGCTCATCTGTGA
1630 CCTACGATCATCCCTCCAAACCTGTTGAAAAAAGACTCATGCCCCCACCCTCGCAA
1631 AAAAACAAGATTCATGGCTGGACCATGTA
1632
1633

1634 **GLB1Dup84**

1635 ATGCCGGGGTTCCTGGTTCGCATCCTCCCTCTGTTGCTGGTTCTGCTGCTTCTGG
1636 GCCCTACGCGCGGCTTGGCGCAATGCCACCCAGAGGATGTTTCAAATTGACTATAG
1637 CCGGGACTCCTTCTCAAGGATGGCCAGCCATTTGCTACATCTCAGGAAGCATTG
1638 ACTACTCCCGTGTGCCCCGCTTCTACTGGAAGGACCGGCTGCTGAAGATGAAGAT
1639 GGCTGGGCTGAACGCCATCCAGACGTATGTGCCCTGGAACCTTCATGAGCCCTGG
1640 CCAGGACAGTACCAGTTTTCTGAGGACCATGATGTGGAATATTTTCTTCGGCTGGC
1641 TCATGAGCTGGGACTGCTGGTTATCCTGAGGCCCGGGCCCTACATCTGTGCAGAG
1642 TGGGAAATGGGAGGATTACCTGCTTGGCTGCTAGAGAAAGAGTCTATTCTTCTCCG
1643 CTCCTCCGACCCAGATTACCTGGCAGCTGTGGACAAGTGGTTGGGAGTCCTTCTG
1644 CCCAAGATGAAGCCTCTCCTCTATCAGAATGGAGGGCCAGTTATAACAGTGCAGG
1645 TTGAAAATGAATATGGCAGCTACTTTGCCTGTGATTTTGACTACCTGCGCTTCTGC
1646 AGAAGCGCTTTCGCCACCATCTGGGGGATGATGTGGTTCTGTTTACCACTGATGGA
1647 GCACATAAAACATTCTGAAATGTGGGGCCCTGCAGGGCCTCTACACCACGGTGG
1648 ACTTTGGAACAGGCAGCAACATCACAGATGCTTTCCTAAGCCAGAGGAAGTGTGA
1649 GCCCAAAGGACCCTTGATCAATTCTGAATTCTATACTGGCTGGCTAGATCACTGGG
1650 GCCAACCTCACTCCACAATCAAGACCGAAGCAGTGGCTTCTCCTCTATGATATA
1651 CTTGCCCGTGGGGCGAGTGTGAACTTGTACATGTTTATAGGTGGGACCAATTTTGC
1652 CTATTGGAATGGGGCCAACTCACCTATGCAGCACAGCCCACCAGCTACGACTAT
1653 GATGCCCCACTGAGTGAGGCTGGGGACCTCACTGAGAAGTATTTTCTCTGCGAA
1654 ACATCATCCAGAAGTTTAAAAAGTACCAGAAGGTCCTATCCCTCCATCTACACCA
1655 AAGTTTGCATATGGAAAGGTCACCTTTGGAAAAGTTAAAGACAGTGGGAGCAGCTCT
1656 GGACATTCTGTGTCCCTCTGGGCCCATCAAAGCCTTTATCCCTTGACATTTATCCA
1657 GGTGAAACAGCATTATGGGTTTGTGCTGTACCGGACAACACTTCTCAAGATTGCA
1658 GCAACCCAGCACCTCTCTTTCACCCCTCAATGGAGTCCACGATCGAGCATATGTT
1659 GCTGTGGATGGGATCCCCAGGGAGTCCCTTGAGCGAAACAATGTGATCACTCTGA
1660 ACATAACAGGGAAAGCTGGAGCCACTCTGGACCTTCTGGTAGAGAACATGGGACG
1661 TGTGAACTATGGTGCATATATGGTGCATATATCAACGATTTTAAAGGGTTTGGTTTCT
1662 AACCTGACTCTCAGTTCCAATATCCTCACGGACTGGACGATCTTTCCTACTGGACAC
1663 TGAGGATGCAGTGTGCAGCCACCTGGGGGGCTGGGGACACCGTGACAGTGGCCA
1664 CCATGATGAAGCCTGGGCCCAACTCATCCAACTACACGCTCCCGGCCTTTTATA
1665 TGGGGAACCTTCTCCATTCCAGTGGGATCCCAGACTTGCCCCAGGACACCTTTATC
1666 CAGTTTCTGGATGGACCAAGGGCCAGGTCTGGATTAATGGCTTTAACCTTGGCC
1667 GCTATTGGCCAGCCCGGGGCCCTCAGTTGACCTTGTGTTGTGCCCCAGCACATCCT
1668 GATGACCTCGGCCCAACACCATCACCGTGCTGGAAGTGGAGTGGGCACCCTG
1669 CAGCAGTGATGATCCAGAACTATGTGCTGTGACGTTTCGTGGACAGGCCAGTTATT
1670 GGCTCATCTGTGACCTACGATCATCCCTCCAAACCTGTTGAAAAAAGACTCATGCC
1671 CCCACCCCGCAAAAAACAAAGATTCATGGCTGGACCATGTA

1672

1673

1674 **PORCNwt**

1675 ATGGCCACCTTTAGCCGCCAGGAATTTTTCCAGCAGCTACTGCAAGGCTGTCTCCT
1676 GCCTACTGCCCAGCAGGGCCTTGACCAGATCTGGCTGCTCCTTGCCATCTGCCTC
1677 GCCTGCCGCTCCTCTGGAGGCTCGGGTTGCCATCCTACCTGAAGCATGCAAGCA
1678 CCGTGGCAGGCGGGTTCTTCAGCCTCTACCACTTCTTCCAGCTGCACATGGTTTG
1679 GGTCGTGCTGCTCAGCCTCCTGTGCTACCTCGTGCTGTTCCCTCTGCCGACATTCC
1680 CCCATCGAGGCGTCTTCCTATCCGTCACCATCCTCATCTACCTACTCATGGGTGAG
1681 ATGCACATGGTAGACACCGTGACATGGCACAAGATGCGAGGGGCACAGATGATTG
1682 TGGCCATGAAGGCAGTGTCTCTGGGCTTCGACCTGGACCGGGGCGAGGTGGGTA
1683 CGGTGCCCTCGCCAGTGGAGTTCATGGGCTACCTCTACTTCGTGGGCACCATCGT
1684 CTTGGGCCCTGGATATCCTTCCACAGCTACCTACAAGCTGTCCAAGGCCGCCCA
1685 CTGAGCTGCCGGTGGCTGCAGAAGGTGGCCCGGAGCCTGGCACTGGCCCTGCTG
1686 TGCCTTGTGCTGTCCACTTGCGTGGGCCCTACCTCTTCCCGTACTTCATCCCCCT
1687 CAACGGTGACCGCCTCCTTCGCAAGGGCACCATGGTAAGGTGGCTGCGAGCCTA
1688 CGAGAGTGCTGTCTCCTTCCACTTCAGCAACTATTTTGTGGGCTTTCTTTCCGAGG
1689 CCACGGCCACGTTGGCGGGGGCTGGCTTTACCGAGGAGAAGGATCACCTGGAAT
1690 GGGACCTGACGGTGTCCAAGCCACTGAATGTGGAGCTGCCTCGGTCAATGGTGG
1691 AAGTTGTCACAAGCTGGAACCTGCCCATGTCTTATTGGCTAAATAACTATGTTTTCA
1692 AGAATGCTCTCCGCCTGGGGACCTTCTCGGCTGTGCTGGTCACCTATGCAGCCAG
1693 CGCCCTCCTACATGGCTTCAGTTTCCACCTGGCTGCGGTCCTGCTGTCCCTGGCT
1694 TTTATCACTTACGTGGAGCATGTCCTCCGGAAGCGCCTGGCTCGGATCCTCAGTG
1695 CCTGTGTCTTGTCAAAGCGGTGCCCGCCAGACTGTTCCGACCAGCATCGCTTGGG
1696 CCTGGGGGTGCGAGCCTTAAACTTGCTCTTTGGAGCTCTGGCCATCTTCCACCTG
1697 GCCTACCTGGGCTCCCTGTTTGATGTCGATGTGGATGACACCACAGAGGAGCAGG
1698 GCTACGGCATGGCATACTGTCCACAAGTGGTCAGAGCTCAGCTGGGCCAGTCA
1699 CTGGGTCACTTTTGGATGCTGGATCTTCTACCGTCTCATAGGC

1700

1701 **PORCNDup20**

1702 ATGGCCACCTTTAGCCGCCAGGAATTTTTCCAGCAGCTACTGCAAGGCTGTCTCCT
1703 GCCTACTGCCCAGCAGGGCCTTGACCAGATCTGGCTGCTCCTTGCCATCTGCCTC
1704 GCCTGCCGCTCCTCTGGAGGCTCGGGTTGCCATCCTACCTGAAGCATGCAAGCA
1705 CCGTGGCAGGCGGGTTCTTCAGCCTCTACCACTTCTTCCAGCTGCACATGGTTTG
1706 GGTCGTGCTGCTCAGCCTCCTGTGCTACCTCGTGCTGTTCCCTCTGCCGACATTCC
1707 CCCATCGAGGCGTCTTCCTATCCGTCACCATCCTCATCTACCTACTCATGGGTGAG
1708 ATGCACATGGTAGACACCGTGACATGGCACAAGATGCGAGGGGCACAGATGATTG
1709 TGGCCATGAAGGCAGTGTCTCTGGGCTTCGACCTGGACCGGGGCGAGGTGGGTA
1710 CGGTGCCCTCGCCAGTGGAGTTCATGGGCTACCTCTACTTCGTGGGCACCATCGT
1711 CTTGGGCCCTGGATATCCTTCCACAGCTACCTACAAGCTGTCCAAGGCCGCCCA
1712 CTGAGCTGCCGGTGGCTGCAGAAGGTGGCCCGGAGCCTGGCACTGGCCCTGCTG
1713 TGCCTTGTGCTGTCCACTTGCGTGGGCCCTACCTCTTCCCGTACTTCATCCCCCT

1714 CAACGGTGACCGCCTCCTTCGCAAGGGCACCATGGTAAGGTGGCTGCGAGCCTA
1715 CGAGAGTGCTGTCTCCTTCCACTTCAGCAACTATTTTGTGGGCTTTCTTTCCGAGG
1716 CCACGGCCACGTTGGCGGGGGCTGGCTTTACCGAGGAGAAGGATCACCTGGAAT
1717 GGGACCTGACGGTGTCCAAGCCACTGAATGTGGAGCTGCCTCGGTCAATGGTGG
1718 AAGTTGTCACAAGCTGGAACCTGCCCATGTCTTATTGGCTAAATAACTATGTTTTCA
1719 AGAATGCTCTCCGCCTGGGGACCTTCTCGGCTGTGCTGGTCACCTATGCAGCCAG
1720 CGCCCTCCTACATGGCTTCAGTTTCCACCTGGCTGCGGTCCTGCTGTCCCTGGCT
1721 TTTATCCCTGGCTTTTATCACTTACGTGGAGCATGTCCTCCGGAAGCGCCTGGCTC
1722 GGATCCTCAGTGCCTGTGTCTTGTCAAAGCGGTGCCCGCCAGACTGTTTCGCACCA
1723 GCATCGCTTGGGCCTGGGGGTGCGAGCCTTAAACTTGCTCTTTGGAGCTCTGGCC
1724 ATCTTCCACCTGGCCTACCTGGGCTCCCTGTTTGTATGTCGATGTGGATGACACCAC
1725 AGAGGAGCAGGGCTACGGCATGGCATACTGTCCACAAGTGGTCAGAGCTCAG
1726 CTGGGCCAGTCACTGGGTCACTTTTGGATGCTGGATCTTCTACCGTCTCATAGGC
1727

1728 **References**

- 1729 1. Ceccaldi, R., Rondinelli, B. & D'Andrea, A. D. Repair Pathway Choices and Consequences at
1730 the Double-Strand Break. *Spec. Issue Qual. Control* **26**, 52–64 (2016).
- 1731 2. DiCarlo, J. E., Chavez, A., Dietz, S. L., Esvelt, K. M. & Church, G. M. Safeguarding
1732 CRISPR-Cas9 gene drives in yeast. *Nat. Biotechnol.* **33**, 1250 (2015).
- 1733 3. McVey, M. & Lee, S. E. MMEJ repair of double-strand breaks (director's cut): deleted
1734 sequences and alternative endings. *Trends Genet.* **24**, 529–538 (2008).
- 1735 4. Yu, A. M. & McVey, M. Synthesis-dependent microhomology-mediated end joining accounts
1736 for multiple types of repair junctions. *Nucleic Acids Res.* **38**, 5706–5717 (2010).
- 1737 5. Davis, A. J. & Chen, D. J. DNA double strand break repair via non-homologous end-joining.
1738 *Transl. Cancer Res.* **2**, 130–143 (2013).
- 1739 6. Heidenreich, E., Novotny, R., Kneidinger, B., Holzmann, V. & Wintersberger, U. Non-
1740 homologous end joining as an important mutagenic process in cell cycle-arrested cells. *EMBO*
1741 *J.* **22**, 2274 (2003).
- 1742 7. Pfeiffer, P., Goedecke, W. & Obe, G. Mechanisms of DNA double-strand break repair and
1743 their potential to induce chromosomal aberrations. *Mutagenesis* **15**, 289–302 (2000).
- 1744 8. Brown, J. S. *et al.* Neddylation Promotes Ubiquitylation and Release of Ku from DNA-
1745 Damage Sites. *Cell Rep.* **11**, 704–714 (2015).
- 1746 9. Landrum, M. J. *et al.* ClinVar: Public archive of interpretations of clinically relevant variants.
1747 *Nucleic Acids Res.* **44**, D862–D868 (2016).
- 1748 10. Stenson, P. D. *et al.* Human Gene Mutation Database: towards a comprehensive central
1749 mutation database. *J. Med. Genet.* **45**, 124 (2008).

1750
1751
1752