

# **Accurate functional classification of thousands of *BRCA1* variants with saturation genome editing**

Gregory M. Findlay<sup>1</sup>, Riza M. Daza<sup>1</sup>, Beth Martin<sup>1</sup>, Melissa D. Zhang<sup>1</sup>, Anh P. Leith<sup>1</sup>, Molly Gasperini<sup>1</sup>, Joseph D. Janizek<sup>1</sup>, Xingfan Huang<sup>1</sup>, Lea M. Starita<sup>1,2\*</sup>, Jay Shendure<sup>1,2,3\*</sup>

<sup>1</sup>Department of Genome Sciences, University of Washington, Seattle, WA 98195

<sup>2</sup>Brotman Baty Institute for Precision Medicine, Seattle, WA 98195

<sup>3</sup>Howard Hughes Medical Institute, University of Washington, Seattle, WA 98195

\*Correspondence to Jay Shendure ([shendure@uw.edu](mailto:shendure@uw.edu)) or Lea Starita ([Istarita@uw.edu](mailto:Istarita@uw.edu))

Running title: Saturation genome editing of *BRCA1*

Keywords: BRCA1, functional assay, genome editing, VUS

1 Variants of uncertain significance (VUS) fundamentally limit the utility of genetic  
2 information in a clinical setting. The challenge of VUS is epitomized by *BRCA1*, a tumor  
3 suppressor gene integral to DNA repair and genomic stability. Germline *BRCA1* loss-of-  
4 function (LOF) variants predispose women to early-onset breast and ovarian cancers.  
5 Although *BRCA1* has been sequenced in millions of women, the risk associated with most  
6 newly observed variants cannot be definitively assigned. Data sharing attenuates this  
7 problem but it is unlikely to solve it, as most newly observed variants are exceedingly rare.  
8 In lieu of genetic evidence, experimental approaches can be used to functionally  
9 characterize VUS. However, to date, functional studies of *BRCA1* VUS have been  
10 conducted in a *post hoc*, piecemeal fashion. Here we employ saturation genome editing to  
11 assay 96.5% of all possible single nucleotide variants (SNVs) in 13 exons that encode  
12 functionally critical domains of *BRCA1*. Our assay measures cellular fitness in a haploid  
13 human cell line whose survival is dependent on intact *BRCA1* function. The resulting  
14 function scores for nearly 4,000 SNVs are bimodally distributed and almost perfectly  
15 concordant with established assessments of pathogenicity. Sequence-function maps  
16 enhanced by parallel measurements of variant effects on mRNA levels reveal mechanisms  
17 by which loss-of-function SNVs arise. Hundreds of missense SNVs critical for protein  
18 function are identified, as well as dozens of exonic and intronic SNVs that compromise  
19 *BRCA1* function by disrupting splicing or transcript stability. We predict that these  
20 function scores will be directly useful for the clinical interpretation of cancer risk based on  
21 *BRCA1* sequencing. Furthermore, we propose that this paradigm can be extended to  
22 overcome the challenge of VUS in other genes in which genetic variation is clinically  
23 actionable.

24  
25 Despite our rapidly advancing knowledge of the genetic underpinnings of human disease, our  
26 ability to predict the phenotypic consequences of an arbitrary genetic variant in a human genome  
27 remains poor. This problem manifests most poignantly in the large numbers of ‘variants of  
28 uncertain significance’ (VUS) identified in ‘clinically actionable’ genes, *i.e.* genes that are  
29 already etiologically linked with a specific disease, and for which a definitive interpretation of  
30 the variant as benign or pathogenic would significantly impact clinical care<sup>1,2</sup>.

31  
32 The gene that perhaps best highlights the challenge of VUS is *BRCA1*. Germline variants that  
33 disrupt *BRCA1* function are associated with a hereditary predisposition to breast and ovarian  
34 cancer<sup>3-6</sup>. Functionally disruptive germline variants in *BRCA1* are clinically actionable, *e.g.* by  
35 more aggressive screening or prophylactic surgery, interventions which lead to improved  
36 outcomes<sup>7,8</sup>. Furthermore, functionally disruptive somatic *BRCA1* mutations influence how  
37 tumors respond to specific therapeutic agents, *e.g.* PARP inhibitors<sup>9-11</sup>. Clinical sequencing of  
38 *BRCA1*, as well as many other genes linked to cancer predisposition such as *BRCA2*, *PALB2*,  
39 *BARD1*, *ATM*, etc., has the potential to implicate specific variants in disease<sup>12</sup>. Documented  
40 pathogenic *BRCA1* variants in the ClinVar database include complete or partial gene deletions,  
41 frameshifting insertions and deletions (indels), nonsense SNVs, missense variants detrimental to  
42 protein stability and function, and both intronic and exonic variants that perturb splicing<sup>13</sup>.  
43 However, as of January 2018, over half of *BRCA1* SNVs in ClinVar are classified as VUS. VUS  
44 are typified by rare missense SNVs, but also include variants potentially affecting mRNA  
45 production, such as SNVs near splice junctions. Further indicative of the challenge of variant  
46 interpretation, ClinVar is replete with *BRCA1* variants that have received conflicting

47 interpretations from different experts. Of 3,936 germline *BRCA1* SNVs currently represented in  
48 ClinVar, only 983 are classified by an expert panel as ‘benign’ or ‘pathogenic’ without  
49 conflicting interpretations.

50  
51 There are two major approaches for resolving VUS. The first approach, data sharing, relies on  
52 the expectation that as *BRCA1* is sequenced in increasing numbers of individuals<sup>14</sup>, the recurrent  
53 observation of a specific variant in multiple individuals who either have or have not developed  
54 breast and/or ovarian cancer will enable the definitive interpretation of that variant. However,  
55 although this may be possible for some variants, given that the vast majority of potential SNVs  
56 in *BRCA1* are exceedingly rare<sup>15,16</sup> and that the phenotype is incompletely penetrant, it may be  
57 decades or centuries before sufficient numbers of humans are included in genotype-phenotype  
58 studies to accurately quantify cancer risk for each individual rare variant.

59  
60 The second approach, functional assessment, has spurred the development of diverse *in vitro*  
61 assays for *BRCA1*<sup>17</sup>. As the homology-directed DNA repair (HDR) function of BRCA1 is key  
62 for tumor suppression, one commonly used assay involves expressing a BRCA1 variant in cells  
63 and assessing the integrity of the cells’ HDR pathway via inducing repair of a double strand  
64 DNA break in a fluorescent reporter construct<sup>18,19</sup>. Other approaches include assays for  
65 embryonic stem cell viability<sup>20</sup>, cell sensitivity to chemotherapeutic drugs<sup>20</sup>, binding to known  
66 partners such as BARD1<sup>18,21</sup>, and minigene-based splicing assays<sup>22,23</sup>. Computational tools can  
67 predict variant effects based on features such as amino acid conservation. However, although  
68 many such metrics correlate with pathogenicity, at present no computational tool is sufficiently  
69 accurate to be used for the clinical interpretation of newly observed *BRCA1* variants in the  
70 absence of genetic or experimental evidence<sup>24,25</sup>.

71  
72 Functional assessment of *BRCA1* variants has historically been limited in several ways. Chiefly,  
73 experimental studies are *post hoc* and have not kept pace with the scaling of *BRCA1* sequencing  
74 and the accumulation of VUS. Additionally, assays that express variants as cDNA-based  
75 transgenes removed from their genomic context<sup>18,21</sup> fail to assess effects on splicing or transcript  
76 stability, as well as potential artifacts of overexpression<sup>26</sup>. Genome editing technologies provide  
77 a means to overcome these challenges. Yet to our knowledge, genome editing has not yet been  
78 applied to functionally characterize VUS in *BRCA1* or other genes similarly linked to cancer  
79 predisposition.

80  
81 Here we set out to apply genome editing to measure the functional consequences of all possible  
82 SNVs in *BRCA1*, regardless of whether they have been previously observed in a human. Given  
83 *BRCA1*’s immense size, this initial study focuses on 13 exons that encode the functionally  
84 critical RING and BRCT domains. In each experiment, a single exon is subjected to ‘saturation  
85 genome editing’<sup>27</sup>, wherein all possible SNVs are simultaneously introduced to a haploid human  
86 cell line in which *BRCA1* is essential. Consequently, *BRCA1* variants that result in nonfunctional  
87 alleles are depleted over time, a selection that is quantified by deep targeted sequencing. We  
88 optimized this method to obtain function scores for 3,893 SNVs, comprising 96.5% of all  
89 possible SNVs in the targeted exons. These function scores are bimodally distributed and nearly  
90 perfectly concordant with expert-based assessments of pathogenicity. We predict that our  
91 functional classifications will be of immediate clinical utility, and argue that the scaling of this

92 approach to additional clinically actionable genes will substantially enhance the utility of genetic  
93 testing.

94

## 95 RESULTS

96

### 97 Saturation genome editing of *BRCA1* exons

98

99 Many genes in the HDR pathway, including those associated with hereditary cancer  
100 predisposition such as *BRCA1*, *BRCA2*, *PALB2* and *BARD1*<sup>12</sup>, were recently identified in a gene  
101 trap screen as being essential in the human haploid cell line HAP1<sup>28</sup> (**Fig. 1a**). To validate this  
102 finding, we designed guide RNAs (gRNAs) to target exons of each of these genes and assessed  
103 HAP1 cell viability after transfecting each gRNA on a plasmid co-expressing Cas9 and a  
104 puromycin resistance cassette<sup>29</sup>. High cell death was evident by light microscopy (**Fig. 1b**), and a  
105 luminescence-based survival assay established that targeting any of these genes substantially  
106 reduces viability of HAP1 cells within one week (**Extended Data Fig. 1**). Deep sequencing of  
107 the edited loci of *BRCA1*-targeted cells confirmed that cell death was consequent to mutations, as  
108 there was widespread selection against frameshifting indels in favor of unedited loci and some  
109 in-frame indels (**Fig. 1c**). Overall, these results confirm the essentiality of HDR pathway  
110 components in HAP1 cells and establish targeted sequencing as a strategy to distinguish  
111 functional vs. non-functional *BRCA1* variants in a population of edited HAP1 cells.

112

113 We next designed and optimized experiments for saturation genome editing (SGE)<sup>27</sup> (**Fig. 1d**) .  
114 We chose to focus on the thirteen exons of *BRCA1* encoding the RING (exons 2-5) and BRCT  
115 domains (exons 15-23) because these domains are essential for the protein's role as a tumor  
116 suppressor<sup>30-32</sup> and harbor missense variants known to be pathogenic or benign, as well as ~400  
117 VUS or variants with conflicting reports of pathogenicity<sup>13,33,34</sup>. To create a library of repair  
118 templates, we used array-synthesized oligo pools containing all possible SNVs spanning each  
119 exon and ~10 bp of adjacent intronic sequence. Oligo pools for each exon were PCR-amplified  
120 and cloned into plasmids with homology arms to mediate genomic integration and make 'SNV  
121 libraries'. Each SNV library molecule also included a fixed synonymous substitution at the target  
122 site to reduce re-cutting by Cas9 after successful HDR<sup>27</sup>. Each SGE experiment targeted a single  
123 exon. In brief, a population of 20 million HAP1 cells was co-transfected on day 0 with the  
124 exon's corresponding SNV library and Cas9/gRNA plasmid. Successfully transfected cells were  
125 selected with puromycin (days 1-4), expanded, and sampled on day 5 and day 11. Variant  
126 frequencies were quantified by targeted amplification and sequencing of the edited exon from  
127 genomic DNA (gDNA) harvested on day 5 and day 11. Negative controls were used to confirm  
128 that PCR amplicons were not derived from the plasmid DNA of the SNV library.

129

130 We initially performed SGE experiments in replicate for each exon in wild-type (WT) HAP1  
131 cells. In each of the 13 exons, we observed depletion of frameshifting indels, confirming  
132 intolerance to loss of *BRCA1* function (**Extended Data Fig. 2**). However, towards achieving  
133 more robust data, we optimized SGE in HAP1 cells in two ways. First, to increase HDR rates in  
134 HAP1 cells, we generated a monoclonal *LIG4* knockout HAP1 line (HAP1-Lig4KO) (**Extended**  
135 **Data Fig. 3a-b**). *LIG4* acts in the non-homologous end joining (NHEJ) pathway, and its  
136 depletion can increase the proportion of cells with HDR-mediated repair of double-stranded  
137 breaks<sup>35,36</sup>. We observed a median 3.6-fold increase in HDR rates on day 5 in HAP1-Lig4KO

138 relative to WT HAP1 (**Fig. 2a**). Second, it is known that HAP1 cells can spontaneously revert to  
139 diploidy<sup>37</sup>. Simply sorting HAP1 cells for 1N ploidy prior to editing improved reproducibility  
140 (**Extended Data Fig. 3c-e**).

141  
142 We next performed optimized SGE experiments for each of the 13 targeted exons in 1N-sorted  
143 HAP1-Lig4KO cells, testing nearly every possible SNV per exon in replicate (**Fig. 2b**).  
144 Functional effects of SNVs on survival were determined by targeted DNA sequencing of each  
145 SNV library as well as the edited exon in gDNA harvested on day 5 and day 11 (**Fig. 2c-e**).  
146 Additionally, targeted RNA sequencing of day 5 samples was used to determine how abundant  
147 exonic SNVs were in *BRCA1* mRNA (**Fig. 2f**). Because these optimizations resulted in greater  
148 reproducibility (**Extended Data Fig. 4**), we moved forward with data from the 1N-sorted HAP1-  
149 Lig4KO cells only.

### 150 151 **Function scores for 3,893 *BRCA1* SNVs**

152  
153 We sought to calculate function scores for each SNV in a way that accurately quantified  
154 selection throughout the experiment while also minimizing experimental biases. First, we  
155 calculated the log<sub>2</sub> ratio of the SNV's frequency on day 11 vs. its frequency in the original  
156 plasmid library. Second, positional biases in editing rates were modeled (using day 5 SNV  
157 frequencies) and subtracted (**Extended Data Fig. 5**). Third, to enable comparisons between  
158 exons, we normalized function scores such that each experiment's median synonymous and  
159 nonsense SNV matched global medians. Finally, a small number of SNVs were filtered out that  
160 could not confidently be scored (e.g. SNVs poorly represented on day 5; **Extended Data Fig. 6**).  
161 Altogether, we obtained function scores for 3,893 SNVs within or immediately intronic to these  
162 exons (**Fig. 2e, Supplementary Table 1, <https://sge.gs.washington.edu/BRCA1>**). This  
163 corresponds to 96.5% of all possible SNVs in these regions.

164  
165 Function scores for SNVs in these 13 *BRCA1* exons were bimodally distributed (**Fig. 2g**). All  
166 nonsense SNVs scored below -1.25 (N = 138, median = -2.12), whereas 98.7% of synonymous  
167 SNVs >3 bp from splice junctions scored above -1.25 (N = 544, median = 0.00). We classified  
168 all SNVs as 'functional', 'non-functional', or 'intermediate' by fitting a two-component  
169 Gaussian mixture model in which the parameters of the 'non-functional' distribution were based  
170 on all nonsense SNVs and the 'functional' distribution based on synonymous SNVs not depleted  
171 in RNA (**Extended Data Fig. 7**). We then used this model to estimate the probability of each  
172 SNV's score being drawn from the non-functional distribution ( $P_{nf}$ ). SNVs with  $P_{nf} < 0.01$  were  
173 categorized as functional (72.5%); SNVs with  $P_{nf} > 0.99$  were categorized as non-functional  
174 (21.1%); and SNVs with  $0.01 < P_{nf} < 0.99$  (6.4%) were categorized as intermediate.

175  
176 Rare missense variants in *BRCA1* are particularly challenging to interpret clinically. Of the  
177 missense SNVs that we scored here, 21.1% (441 of 2,086) scored as non-functional (**Fig. 2h**).  
178 Although most of the remaining missense SNVs were functional (70.6%), there was an  
179 enrichment for missense SNVs with intermediate effects (8.1%, compared to 4.4% of all other  
180 SNVs; Fisher's exact  $P = 2.7 \times 10^{-6}$ ).

181  
182 An advantage of assaying variants by genome editing is that their impact on native regulatory  
183 mechanisms such as RNA splicing can be ascertained<sup>27</sup>. Whereas SNVs disrupting canonical

184 splice sites (the two intronic positions immediately flanking each exon) overwhelmingly scored  
185 as non-functional (89.5%) or intermediate (5.5%) ('CS' in **Fig. 2h**). SNVs positioned 1-3 bp into  
186 the exon or 3-8 bp into the intron had variable effects. We defined SNVs in these regions that did  
187 not alter the amino acid sequence as 'splice region' variants, of which 22.9% were non-  
188 functional ('SR' in **Fig. 2h**), on par with missense SNVs (21.2% non-functional). SNVs  
189 positioned more deeply in introns or in the 5' UTR were similar to non-splice-region  
190 synonymous SNVs, in that they were much less likely to score as non-functional (intronic: 1.8%  
191 non-functional; 5' UTR: 0.0% non-functional; synonymous: 1.3% non-functional).

192

### 193 **Function scores are nearly perfectly concordant with ClinVar**

194

195 We next asked how well our function scores agreed with expert-based clinical variant  
196 interpretations, where available in ClinVar. Of 169 SNVs deemed 'pathogenic' in ClinVar that  
197 overlapped with our classifications, 162 were designated 'non-functional', 2 'functional', and the  
198 remaining 5 'intermediate'. In contrast, of 22 SNVs deemed 'benign' in ClinVar that overlapped  
199 with our classifications, 1 was designated 'non-functional', 1 'intermediate', and 20 'functional'  
200 (**Fig. 3a**). The three SNVs for which our function scores are unambiguously discordant with  
201 ClinVar are discussed further below. A ROC curve showed a sensitivity of 96.7% at 98.2%  
202 specificity when we treat 'likely pathogenic' and 'likely benign' ClinVar annotations as  
203 pathogenic and benign, respectively (**Fig. 3b**). Importantly, our assay accurately predicts ClinVar  
204 interpretations independent of mutational consequence; sensitivity and specificity are high for  
205 both missense and splice site SNVs when these are considered separately from nonsense SNVs  
206 (**Extended Data Fig. 7f**). We find 64 of 256 (25.0%) VUS and 60 of 122 (49.2%) SNVs with  
207 conflicting interpretations to be non-functional in our assay (**Fig. 3c**). Missense VUS from  
208 ClinVar were significantly more likely to score as non-functional compared to missense SNVs  
209 absent from ClinVar (25.9% vs. 17.2%,  $P = 0.002$ ). Apart from largely corroborating established  
210 ClinVar annotations, our scores also provide functional classifications for an additional 3,140  
211 SNVs, the vast majority of which have yet to be publicly reported in clinical sequencing. Of  
212 these SNVs, 498 (15.9%) are classified as non-functional.

213

214 We also investigated the relationship between our function scores and SNV frequencies in large-  
215 scale databases of human genetic variation. Of 302 assayed SNVs that overlap with the Genome  
216 Aggregation Database (gnomAD)<sup>16</sup>, higher allele frequencies were associated with higher  
217 function scores (**Fig. 3d**). For instance, 33 of 166 (19.9%) of singleton gnomAD variants were  
218 non-functional, whereas only 8 of 136 SNVs (5.9%) seen in multiple individuals were non-  
219 functional (Fisher's exact  $P = 3 \times 10^{-4}$ ). A similar trend was observed with the Bravo database  
220 (**Extended Data Fig. 8a**). The FLOSSIES database contains *BRCA1* variants observed in women  
221 over seventy years old who have not developed breast or ovarian cancer<sup>38</sup>. Of 39 intersecting  
222 SNVs, only one scored as non-functional (**Extended Data Fig. 8b**). Collectively, these  
223 observations show that *BRCA1* SNVs with higher allele frequencies are more likely to be  
224 functional, as expected. However, the fact that >70% of ClinVar variants and >95% of non-  
225 ClinVar variants that we assayed here have not been observed even once in sequencing of  
226 >120,000 humans illustrates the challenges facing observational approaches to variant  
227 interpretation.

228

229 Several computational metrics are currently used to assess deleteriousness of variants and  
230 often included in genetic testing reports. Although our function scores correlate with metrics  
231 such as CADD<sup>39</sup>, phyloP<sup>40</sup>, and Align-GVGD<sup>41</sup>, which are largely based on evolutionary  
232 conservation and biochemical properties of missense variants, the modesty of these correlations  
233 underscores the value of functional assays (**Fig. 3e, Extended Data Fig. 9a-g**). ROC curve  
234 analysis restricted to missense variants reveals that SGE-based function scores outperform these  
235 metrics at predicting pathogenicity status in ClinVar (**Extended Data Fig. 9h-l**). This  
236 outperformance is likely underestimated because some of these metrics (*e.g.* Align-GVGD) or  
237 their correlates (*e.g.* evolutionary conservation) informed the ClinVar classifications of  
238 pathogenicity in the first place.

239

### 240 **Mechanisms of *BRCAl* loss-of-function**

241

242 To gain insights into the various mechanisms by which SNVs compromise function, we  
243 performed targeted RNA sequencing of *BRCAl* transcripts from day 5 cells. We normalized  
244 SNV frequencies in cDNA to their frequency in gDNA to produce mRNA expression scores  
245 ('RNA scores') for 96% of the functionally characterized exonic SNVs. Together with function  
246 scores, RNA scores enable fine mapping of molecular consequences of SNVs (**Fig. 4**). For  
247 instance, regions of exons 2 and 15 that respectively code for RING and BRCT domain residues  
248 contain numerous loss-of-function missense variants. This contrasts with coding sequence in the  
249 same exons that fall outside of the boundaries of these protein domains. Overall, 89% of non-  
250 functional missense SNVs did not reduce RNA levels substantially, suggesting that their effects  
251 are likely mediated at the protein level (**Fig. 5a**). Many residues that are sensitive to missense  
252 SNVs *not* impacting RNA levels map to buried hydrophobic residues or to the zinc-coordinating  
253 loops that are required for proper RING domain folding (**Fig. 5b-c**). However, 11% of non-  
254 functional missense SNVs are depleted from RNA by 4-fold or more. Many of these SNVs map  
255 outside of key protein-protein interfaces and rather in unstructured loops, suggesting that they  
256 cause loss-of-function by lowering mRNA expression levels. Consistent with this, the 12  
257 synonymous SNVs classified as non-functional also tended to markedly reduce mRNA levels  
258 (median 5.4-fold reduction).

259

260 How do these exonic SNVs cause reductions in mRNA levels? Although other mechanisms  
261 cannot be ruled out, many of the variants depleted in mRNA are likely impacting RNA splicing.  
262 This is evidenced by an overrepresentation of non-functional SNVs near splice junctions,  
263 including low scores for many SNVs at terminal G nucleotides of exons (**Fig. 4**), non-functional  
264 exonic SNVs with low mRNA levels that create new acceptor or donor sequences (SNVs  
265 annotated with asterisks in **Fig. 5d**), and the presence of short regions (~6-8 bp) in which many  
266 SNVs have moderate-to-strong effects on RNA levels, suggestive of exonic splice enhancers<sup>42</sup>  
267 (**Fig. 5e**). Certain exons appeared particularly prone to harbor non-functional SNVs with low  
268 RNA scores. In exon 16, for instance, 46 of 244 SNVs (excluding nonsense) were non-functional  
269 (**Fig. 5e**). Of these, more than half ( $n = 26$ ) reduced RNA levels by more than 2-fold, and nearly  
270 a third ( $n = 15$ ) by more than 4-fold. In contrast, in exon 19, of 55 of 234 SNVs (excluding  
271 nonsense) that were non-functional, none lowered expression by more than 2-fold (**Fig. 5f**). Exon  
272 19 also completely lacks non-functional SNVs in its flanking intronic regions (apart from the  
273 acceptor and donor sites), suggesting the exon is robustly spliced compared to other exons.

274

## 275 **Discordances with ClinVar Interpretations**

276

277 We leveraged sequence-function maps in reviewing the evidence around the three SNVs for  
278 which our classifications were clearly discordant with ClinVar. Discordant SNVs assayed in our  
279 preliminary experiments in WT HAP1 cells had similar scores, suggesting their classifications  
280 are not secondary to noise in our assay (**Extended Data Fig. 10**). One missense SNV designated  
281 ‘pathogenic’ in ClinVar that we scored as functional, c.5359T>A (C1787S), was identified  
282 through segregation with disease. However, in each case, it was seen in *cis* with a second SNV at  
283 the neighboring amino acid position<sup>43</sup>. Our data as well as data from other functional assays<sup>44</sup>  
284 suggest c.5359T>A on its own is functional. The linked SNV c.5363G>T (G1788D), however,  
285 scored as non-functional, calling into question the ClinVar annotation (**Extended Data Fig.**  
286 **10c**).

287

288 A second disagreement was identified in the exon 2 splice acceptor, c.-19-2A>G. This SNV was  
289 annotated as ‘pathogenic’ in ClinVar based on its occurrence at a splice acceptor site<sup>45</sup>, rather  
290 than from having been associated with disease. Exon 2 contains the *BRCA1* translation initiation  
291 codon, meaning that alternate splice forms may preserve the complete open reading frame. Of  
292 note, CADD scores for SNVs across the exon 2 acceptor site were much lower than for SNVs in  
293 other canonical splice sites (**Extended Data Fig. 10d**), and none of the 6 SNVs that we  
294 introduced here scored as non-functional. Further supporting that this splice site is not essential  
295 for *BRCA1* function, RNA sequencing from breast and ovarian tissue in the GTEx database<sup>46</sup>  
296 shows this exon junction is poorly represented among *BRCA1* transcripts (**Extended Data Fig.**  
297 **10e**). This suggests that this acceptor site is likely dispensable both in our assay and in tissues  
298 relevant to disease, again calling the ClinVar annotation into question.

299

300 Exon 16 harbored the third discordantly classified SNV, the ‘benign’ c.5044G>A (E1682K)  
301 variant, which scored as non-functional in our assay. Of note, c.5044G>A resides in a predicted  
302 exonic splice enhancer (ESE)<sup>42</sup>, and its low function score was substantiated by a reduction in  
303 RNA levels of over 90% (**Fig. 5e**). Neighboring SNVs in the predicted ESE also reduced RNA  
304 expression, corroborating the element’s importance. Although this missense SNV is rare (absent  
305 from gnomAD and Bravo), reports indicate it was designated as benign based on being observed  
306 in *trans* with a variant considered pathogenic<sup>33</sup>, as biallelic *BRCA1* loss-of-function mutations  
307 are thought to be embryonic lethal. The underlying data supporting this finding are not publically  
308 available, and previous assays of this variant did not measure splicing consequences<sup>44</sup>.

309

## 310 **DISCUSSION**

311

312 Here we applied saturation genome editing to the 13 exons that encode functionally critical  
313 domains of the cancer risk gene, *BRCA1*, characterizing the functional consequences of nearly  
314 4,000 SNVs in their native genomic context. Specifically, we used CRISPR/Cas9 to introduce  
315 hundreds of SNVs per experiment, followed by deep sequencing to measure the functional  
316 consequences of each SNV in parallel. Because we measured cell survival, the effects of SNVs  
317 on multiple layers of gene function (*e.g.* RNA splicing, translation, protein function, protein  
318 stability) are effectively integrated. The approach is validated by nearly perfect concordance of  
319 function scores with available evidence for clinical pathogenicity.

320



321 Our experimental approach has several caveats. First, the exact requirements for *BRCA1* function  
322 essential to maintaining *in vitro* viability and growth of HAP1 cells, as opposed to mediating *in*  
323 *vivo* tumor suppression, are not known. For instance, we cannot rule out, differences in splicing  
324 or dosage requirements between our *in vitro* model vs. *in vivo* physiology. Second, we are not  
325 currently able to interrogate every possible SNV. Of note, most of the 3.5% of SNVs for which  
326 we do not provide function scores were excluded by factors related to genome editing, rather  
327 than because of sampling (**Extended Data Fig. 6**). Lastly, as these experiments were designed to  
328 measure loss-of-function in a haploid cell line, we are unable to detect all types of functional  
329 effects (e.g. dominant negative variants).

330  
331 Notwithstanding these limitations, we achieved nearly comprehensive coverage of the targeted  
332 regions and our functional classifications are nearly perfectly concordant with current clinical  
333 interpretations. As such, we anticipate that our results will be clinically useful, both for  
334 adjudicating hundreds of observed variants whose interpretation is currently ambiguous, as well  
335 as for providing immediate functional assessments for variants newly observed. Therefore, the  
336 pressing question becomes how to best to integrate this functional data within existing clinical  
337 variant classification schemes<sup>47</sup>.

338  
339 A benefit of functional data is that measurements are systematically derived, independent of  
340 prior expectation<sup>48</sup>. As such, function scores add an additional layer of evidence to support  
341 interpretations of variants made through segregation with disease. However, for the large number  
342 of VUS for which genetic evidence is insufficient, the predictive power demonstrated here  
343 suggests function scores can be used to classify variants with >95% accuracy. As current  
344 standards for defining ‘likely pathogenic’ and ‘likely benign’ variants accept a comparable level  
345 of uncertainty<sup>49</sup>, we argue that a failure to use appropriately validated functional data to inform  
346 clinical care would be a missed opportunity. There is precedent for incorporating functional data  
347 in interpretation guidelines<sup>24</sup>, but the breadth and predictive power demonstrated by SGE calls  
348 for an increased role. Indeed, given the low likelihood that observational approaches will ever be  
349 sufficient to classify variants not yet seen once in humans, we believe that there is a strong  
350 argument to be made for using highly predictive function scores, where available, to inform  
351 initial interpretations of newly observed variants.

352  
353 The orthologous nature of SGE data also presents an opportunity for integration with other data  
354 sources. For example, a multiplex reporter assay for HDR activity strengthens the functional  
355 evidence presented here for *BRCA1* missense variants (see accompanying manuscript from  
356 Starita *et al.*). Integration and optimal weighting of experimental and computational approaches  
357 may also further improve classification of variants lacking genetic evidence. In cases where  
358 evidence is contradictory, functional data may yield specific hypotheses to test. For example,  
359 c.5044G>A, for which our data contradicts the ClinVar interpretation (**Fig. 5e**), would be  
360 disambiguated by testing *BRCA1* mRNA levels in individuals harboring this SNV. Similar  
361 approaches should be taken to more confidently resolve unlikely functional classifications, such  
362 as synonymous SNVs with low function scores and canonical splice SNVs deemed functional.  
363 Furthermore, the ~6% of SNVs exhibiting intermediate function scores remain beyond definitive  
364 interpretation. The fact that we observe an excess of missense SNVs with intermediate scores  
365 suggests that some of these may be hypomorphic *BRCA1* alleles<sup>50-52</sup>. Further studies will be  
366 necessary to quantify the penetrance of intermediately functional variants.

367

368 Moving forward, our study provides a blueprint for comprehensive functional analysis of all  
369 potential SNVs in clinically actionable genes for which appropriate assays can be developed.  
370 Here, we prioritized *BRCA1* exons encoding the RING and BRCT domains, but SGE of the  
371 entire coding sequence and promoter are also well motivated. Furthermore, the essentiality of  
372 *BRCA2*, *PALB2*, *BARD1*, and *RAD51C* in HAP1 suggests that these genes are assayable by the  
373 same method. For genes in other pathways, assays that are compatible with saturation genome  
374 editing (*e.g.* drug selection, FACS on phenotypic markers, etc.) may need to be developed and  
375 validated. For any gene tested, it is critical that functional measurements be calibrated to clinical  
376 evidence of pathogenicity. Given that SGE tests variants in their endogenous genomic context,  
377 the scaling of SGE to many loci promises to improve our understanding of how diverse  
378 biological functions are encoded by the genome.

379

380 Delivering on the promise of genomic medicine requires that we not only be able to cost-  
381 effectively ascertain genetic variation, but also accurately and definitively interpret it. Presently,  
382 interpretation is the rate limiting step. As a potential path forward, we show that saturation  
383 genome editing is a viable strategy for functionally classifying thousands of variants in a  
384 clinically actionable gene, most of which have yet to be observed in a human. With further  
385 scaling, we anticipate that this paradigm will substantially improve the utility of genetic  
386 information in clinical decision making.

387 **DATA AVAILABILITY**

388

389 Saturation genome editing data is available at: <https://sge.gs.washington.edu/BRCA1>.

390

391 **ACKNOWLEDGEMENTS**

392

393 We thank Malte Spielmann, Daniela Witten, Aaron McKenna, Martin Kircher, Max Dougherty,  
394 John Lazar, Yi Yin, and Brian Shirts for insights on data analysis and/or comments on the  
395 manuscript, Jacob Kitman for sharing reagents and protocols, Rocío Acuña-Hidalgo and  
396 Jennifer Milbank for experimental assistance and the Feng Zhang lab for sharing Cas9/gRNA  
397 plasmids. This work was supported by an NIH Director's Pioneer Award (DP1HG007811 to J.S.)  
398 and a training award from the National Cancer Institute (F30CA213728 to GMF). JS is an  
399 Investigator of the Howard Hughes Medical Institute.

## 400 METHODS

401

### 402 HDR pathway essentiality analysis in HAP1 cells

403 HAP1 cells were derived from KBM7 cells (a near-haploid immortalized chronic  
404 myelogenous leukemia line) by introduction of induced pluripotent stem cell factors<sup>56</sup>. HAP1  
405 gene essentiality scores were obtained<sup>28</sup> and filtered on genes with greater than 20 mapped gene-  
406 trap insertions (N = 14,306). Of 78 HDR genes defined by the GO term ‘double-strand break  
407 repair via homologous recombination’ (GO:0000724), 66 were among the 14,306 genes included  
408 in analysis. To rank genes by essentiality, they were first ordered by q-value (low to high) and  
409 second by the proportion of gene-trap insertions in the sense orientation (low to high). HDR  
410 pathway genes implicated in cancer (labelled in Fig. 1) were defined as those included on the  
411 University of Washington BROCA sequencing panel<sup>57</sup>.

412

### 413 gRNA design and cloning

414 All CRISPR gRNAs used in SGE and essentiality experiments were cloned into pX459<sup>29</sup>.  
415 This plasmid expresses the gRNA from a U6 promoter, as well as a Cas9-2A-puromycin  
416 resistance (puroR) cassette. *S. pyogenes* Cas9 target sites were chosen for SGE experiments on  
417 multiple criteria, assessed in the following order: 1.) To induce cleavage within *BRCAl* coding  
418 sequence, 2.) To target a genomic site permissive to synonymous substitution within the guanine  
419 dinucleotide of the PAM or the protospacer, 3.) To have minimal predicted off-target activity<sup>58</sup>,  
420 4.) To have maximal predicted on-target activity<sup>59</sup>.

421 Complementary oligos ordered from Integrated DNA Technologies (IDT) were annealed,  
422 phosphorylated, diluted and ligated into BbsI-digested and gel-purified pX459, as described<sup>29</sup>.  
423 Ligation reactions were transformed into *E. coli* (Stellar competent cells, Takara), which were  
424 plated on ampicillin. Colonies were cultured and Sanger sequenced to confirm correct gRNA  
425 sequences. Purification of sequence-verified plasmids for transfection was performed with the  
426 ZymoPure Maxiprep kit (ZymoResearch). For targeting *LIG4* in HAP1 cells, pX458<sup>29</sup> was used  
427 instead of pX459, which expresses EGFP in lieu of puroR.

428

### 429 HDR library design and cloning

430 Array-synthesized oligos were designed as follows for each saturation genome editing  
431 region (*i.e.* a *BRCAl* exon). The sequence to be mutated (~100bp) was obtained from the human  
432 genome (hg19) and a synonymous substitution was introduced at the chosen Cas9 target site (*e.g.*  
433 a substitution at the PAM site). This ‘fixed’ substitution in the library was included in design to  
434 serve multiple purposes: 1.) plasmid library molecules harboring the substitution are predicted to  
435 be cleaved less frequently by Cas9:gRNA complexes, 2.) SNVs introduced to cells are predicted  
436 to be depleted via Cas9 re-cutting less frequently as a consequence of the fixed substitution, and  
437 3.) sequencing reads can be filtered on the fixed substitution to distinguish true SNVs introduced  
438 via HDR from sequencing errors. A second synonymous substitution at an alternative CRISPR  
439 target site was introduced to the sequence as well, such that each exon’s SNV library would be  
440 compatible with multiple gRNAs. Next, a sequence was created for every possible single  
441 nucleotide substitution on this template. For all sequences, adapters were added to both ends to  
442 enable PCR amplification from the oligo pool. For each SGE region, the total number of oligos  
443 designed was three times the length of the region, plus the oligo template without any SNV (*e.g.*  
444 for a 100 bp SGE region, 301 total oligos were designed).

445 Pooled oligos were synthesized (Agilent Technologies). Primers designed to amplify the  
446 subset of oligos corresponding to a single exon's region were used to perform PCR with Kapa  
447 HiFi Hot-start Ready Mix ('Kapa HiFi', Kapa Biosystems). PCR products were purified with  
448 Ampure beads (Agencourt) to be used in subsequent library cloning reactions.

449 Homology arms were cloned into pUC19 by PCR-amplifying (Kapa HiFi) regions  
450 surrounding each targeted exon from HAP1 gDNA. Primers for these reactions were designed  
451 such that homology arms would be between 600 and 1,000 bp on both sides of the targeted  
452 region. Adapters homologous to pUC19 were added to primers to facilitate NEBuilder HiFi  
453 Assembly cloning (NEB) into a linearized pUC19 vector. Cloning reactions were transformed  
454 into Stellar competent cells and selected with ampicillin. Plasmid DNA was isolated from  
455 colonies (Qiagen MiniPrep kit) and sequence-verified.

456 To make the HDR library, homology arm plasmids were linearized via PCR using  
457 primers that conferred 15-20 bp of terminal overlap with the adapter sequences flanking each  
458 PCR-amplified oligo pool. This sequence overlap enabled cloning via the NEBuilder HiFi  
459 Assembly Cloning Kit (NEB). Cloning reactions were transformed into Stellar competent cells,  
460 and a small proportion (1%) of the transformation was plated on ampicillin-containing plates to  
461 assess efficiency. All remaining transformed cells were grown directly in 100 ml of media with  
462 ampicillin for 16-18 hours, and plasmid DNA from the culture was isolated (ZymoPure  
463 Maxiprep kit) to produce each final HDR library.

464

#### 465 HAP1 cell culture

466 Quality-controlled WT HAP1 cells were purchased (Haplogen/Horizon Discovery) and  
467 cultured in media comprising Iscove's Modified Dulbecco's Medium (IMDM) with L-glutamine  
468 and 25 mM HEPES (GIBCO) supplemented with 10% fetal bovine serum (Rocky Mountain  
469 Biologicals) and 1% penicillin-streptomycin (GIBCO). Cells were grown on plates at 37C with  
470 5% CO<sub>2</sub>, and passaged prior to becoming confluent. For routine passaging, cells were washed  
471 once with 1x phosphate buffered saline (PBS, Gibco), trypsinized with 0.25% trypsin with  
472 EDTA (Gibco), resuspended in media, centrifuged for 5 min at 300 rcf, and then resuspended  
473 and plated.

474 A monoclonal *LIG4* knock-out HAP1 line (HAP1-Lig4KO) was generated by  
475 transfecting a plasmid expressing a Cas9-2A-GFP cassette and a gRNA targeting the human  
476 *LIG4* coding sequence (gRNA sequence: 5'-GCATAATGTCACACTACAGATC) into WT HAP1  
477 cells. Single GFP-expressing HAP1 cells were sorted into wells of a 96-well plate and cultured.  
478 After two weeks, gDNA was harvested and Sanger sequencing was performed to assess *LIG4*  
479 editing. A clone with a 4bp deletion was identified and expanded further for use in saturation  
480 genome editing experiments.

481 HAP1 cells can spontaneously revert to a diploid state in cell culture. Therefore, to sort a  
482 1N-enriched population of cells prior to transfection, cells were stained for DNA content with  
483 Hoechst 34580 (BD Biosciences) at 5 ug/ml media for 1h at 37C. FACS was performed to  
484 isolate 1-2x10<sup>6</sup> cells from the lowest intensity Hoechst peak, corresponding to 1N ploidy. These  
485 cells were expanded for seven days prior to transfection.

486

#### 487 Transfection of HAP1 cells

488 For all experiments, HAP1 cells were transfected using TurboFectin 8.0 (Origene)  
489 according to manufacturer's protocol. A 2.5x volume of Turbofectin was added to the  
490 transfection mix for each ug of plasmid DNA in Opti-Mem (Life Technologies). For each SGE

491 transfection, 10 million cells were passaged to a 10 cm dish. The next day (day 0), cells were co-  
492 transfected with 12 ug of the Cas9/gRNA plasmid (pX459) and 3 ug of the SGE library  
493 corresponding to a single exon. For negative control transfections, a pX459 vector targeting  
494 *HPRT1* was used instead. On day 1, cells were passaged into media supplemented with  
495 puromycin (1 ug/ml) to select for successfully transfected cells. On day 4, cells were washed  
496 twice and passaged to 6 cm plates in regular media.

497 Cell populations were sampled on day 5 and day 11 for all SGE experiments. On day 5,  
498 half of the cells were pelleted and frozen and the other half passaged. The cells were passaged on  
499 day 8 into 15 cm dishes and then harvested on day 11. Negative control transfections were  
500 harvested on day 5.

501 For the luminescence-based viability assay, HAP1 cells were plated at ~35-40%  
502 confluency in a 6-well dish (approximately 1.2 million cells per well per target) then transfected  
503 with 1.5 ug Cas9/gRNA plasmid targeting coding exons of HDR genes or controls the following  
504 day. 24 hours after transfection the cells were plated in time-point triplicates at 20,000 cells per  
505 well in 96-well clear bottom plates in media with and without puromycin. Cells without  
506 puromycin were assessed 4 hours after plating to establish baseline absorbance for each target.  
507 Cell survival was assessed at day 2, day 5, and day 7 post-transfection using the CellTiterGlow  
508 reagent (Promega, 1:10 dilution of suggested reagent). Luminescence at 135 nm absorbance was  
509 measured using a Synergy plate reader (Biotek Instruments).

510

#### 511 Nucleic acid sampling and sequencing library production

512 For obtaining WT HAP1 genomic DNA for cloning homology arms and for genotyping  
513 the HAP1-Lig4KO cell line, DNA was isolated using the DNeasy kit (Qiagen). For each SGE  
514 experiment, DNA and total RNA were purified using the AllPrep kit (Qiagen). DNA samples  
515 were quantified with the Qubit dsDNA Broad Range kit (Thermo Fisher) and RNA samples by  
516 UV spectrometry (Nanodrop). PCR primers for genomic DNA were designed such that one  
517 primer would anneal outside of the homology arm sequence, thereby selecting for amplicons  
518 derived from gDNA and not plasmid DNA. PCR conditions were optimized using gradient qPCR  
519 on WT HAP1 gDNA.

520 All gDNA harvested from the population of day 5 cells was sampled by performing many  
521 PCR reactions in parallel on a 96-well plate, using 250 ng of gDNA per 50 ul reaction such that  
522 all day 5 gDNA was used in PCR (Kapa HiFi). At least as many PCR reactions were performed  
523 for day 11 samples (which yielded more gDNA) to ensure adequate sampling. PCRs were  
524 performed for the minimal number of cycles needed to complete amplification, with cycling  
525 conditions as specified in the Kapa HiFi protocol. An additional PCR was performed using day 5  
526 gDNA from negative control transfections for each exon.

527 After PCR, multiple wells of amplicons from the same sample were pooled and purified  
528 using Ampure beads. Next, a nested qPCR was performed using the first reaction as template to  
529 produce a smaller amplicon with custom sequencing adapters ('PU1L' and 'PU1R'), which was  
530 likewise purified with Ampure beads. The SGE libraries were also PCR-amplified at this step,  
531 starting from 50 ng of plasmid DNA. Lastly, a final qPCR was performed using purified  
532 products from the second reaction as template to add dual sample indexes and flow cell adapters.

533 RNA was sampled from day 5 HAP1-Lig4KO cells (AllPrep, Qiagen). Reverse  
534 transcription followed by RNase H treatment was performed on all RNA harvested or a  
535 maximum of 5 ug per sample (Superscript IV Kit, Life Technologies). This reaction was primed  
536 with a gene-specific primer complementary to the 3' UTR in exon 23 of *BRCA1*. Primers were

537 designed for each exon to amplify across exon junctions, and reaction conditions were optimized  
538 using gradient PCR. cDNA was distributed into 5 equal PCR reactions, which were run on a  
539 qPCR machine and then pooled in equal ratios. Flow cell adapters and sample indexes were  
540 added in an additional reaction (as for gDNA samples).

541 All sequencing libraries were purified with Ampure beads, quantified with the Qubit  
542 dsDNA High Sensitivity kit (Life Technologies), diluted and denatured for sequencing in  
543 accordance with protocols for the Illumina NextSeq or MiSeq machines.

544

#### 545 Sequencing and data analysis

546 Sequencing was performed on an Illumina NextSeq or MiSeq instrument, allocating  
547 about 3 million reads to each gDNA and cDNA sample, 1 million reads for each HDR library,  
548 and 500,000 reads for each negative control sample. gDNA samples for individual exons were  
549 sequenced on the same run. 300 cycle kits were used, with 150 cycles for read 1 and read 2 each,  
550 and 19 cycles for dual index reads. Custom sequencing primers and indexing primers are  
551 provided in Supplementary Table 2. Illumina PhiX control DNA was added to each sequencing  
552 run (~10% MiSeq, ~30-40% NextSeq) to improve base calling.

553 Illumina's bcl2fastq 2.16 was used to call bases and perform sample demultiplexing and  
554 fastqc 0.11.3 was run on all samples to assess sequencing quality. SeqPrep was used with the  
555 following parameters to perform adapter trimming and to merge perfectly matched overlapping  
556 read pairs: '-A GGTGGAGCGAGATTGATAAAGT -B  
557 CTGAGCTCTCTCACAGCCATTTAG -M 0.1 -m 0.001 -q 20 -o 20'. Merged reads containing  
558 'N' bases were removed. Reads from cDNA samples were removed if they contained indels or  
559 did not perfectly match transcript sequence flanking each targeted exon. Remaining cDNA reads  
560 were processed to match genomic DNA amplicons by removing flanking exonic sequence and  
561 replacing it with the exon's corresponding intronic sequence. All reads were then aligned to  
562 reference gDNA amplicons for each exon using the needleall command in the EMBOSS 6.4.0  
563 package with the following parameters: '-gapopen 10 -gapextend 0.5 -aformat sam'. Reads not  
564 aligning to the reference amplicon (alignment score < 300) were removed from analysis. To  
565 analyze indels, unique cigar counts were quantified from day 5 and day 11 samples using a  
566 custom Python script. Reads were classified as HDR events for rate calculations if the  
567 programmed edit or edits to the PAM or protospacer (HDR marker edits) were observed in the  
568 alignment. Variants without identifiable markers of HDR were not used. Abundances of SNVs  
569 were quantified only from aligned reads that had no other mismatches or indels, with the  
570 exception of the HDR markers. SNV reads with only the cut-site proximal HDR marker were  
571 summed with reads that had both HDR markers to get total abundances for each SNV in each  
572 sample, to which a pseudocount of 1 was added to all variants present in either the library, day 5  
573 or day 11 sample. Frequencies for each SNV were calculated as SNV reads over total reads.  
574 SNV measurements from WT HAP1 cells and HAP1-Lig4KO cells were processed separately at  
575 all steps.

576

#### 577 Modeling positional biases of library integration

578 Positional biases in editing rates were modeled for each SNV by using a LOESS  
579 regression to fit the log2 day 5 over library ratios as a function of chromosomal position. To  
580 avoid modeling biological effects instead of positional effects, the model was fit only on the  
581 subset of SNVs that were not substantially depleted between any two timepoints in the  
582 experiment (*i.e.* SNVs with day 5 over library ratios > 0.5 and day 11 over d5 ratios > 0.8.). The

583 regression was performed for each exon replicate, using the ‘loess’ function in R with span =  
584 0.15. Each model was extended flatly outward to include any positions not fit (a total of 22  
585 nucleotides of sequence on the edges of the edited regions). We subtracted each SNV’s  
586 positional fit (e.g. the model’s output) from the SNV’s log<sub>2</sub> day 11 over library ratio to get  
587 position-adjusted ratios for each SNV.  
588

#### 589 Normalizing scores within and across exons

590 Position-adjusted log<sub>2</sub> day 11 over library ratios were normalized first across exon  
591 replicates, and then across all exons assayed. Scores from within each replicate were linearly  
592 scaled such that the median synonymous and median nonsense SNVs within the replicate were  
593 set to the median synonymous and median nonsense SNV values averaged across replicate  
594 experiments. The ensuing SNV scores for each replicate were then normalized across exons in  
595 the same way by again using median synonymous and median nonsense SNVs.  
596

#### 597 SNV functional class assignment

598 Function scores were averaged across replicates and a mixture model was used to  
599 estimate the probability that each SNV’s score was drawn from the non-functional distribution of  
600 scores. The non-functional distribution was defined as nonsense SNVs across all exons. The  
601 functional distribution was defined as exonic synonymous SNVs not within 3 bp of splice  
602 junctions and with RNA scores within 1 standard deviation of the median synonymous SNV.  
603 This definition does not fully guarantee that these SNVs have no functional consequence. The  
604 means and variances of the ‘non-functional’ and ‘functional’ groups were fixed and a model was  
605 fit using the normalmixEM function of the mixtools package in R, with starting component  
606 proportions set to 0.5. The posterior probabilities generated from the model were used as point  
607 estimates of the probability of drawing each SNVs score from the non-functional distribution  
608 ( $P_{nf}$ ). Functional classifications were made by setting thresholds for  $P_{nf}$  as follows:  $P_{nf} > 0.99 =$   
609 ‘non-functional’,  $0.01 < P_{nf} < 0.99 =$  ‘intermediate’,  $P_{nf} < 0.01 =$  ‘functional’.

610 Independent of mixture modelling, ROC curves were used to assess performance of SGE  
611 data and other metrics’ ability to predict assigned ClinVar classifications. These analyses were  
612 performed with the plotROC package in R, and Youden’s J-statistic was calculated (sensitivity  
613 plus specificity minus 1) to determine optimal values reported in text.  
614

#### 615 Variant filtering

616 A small minority of SNVs that could not be accurately scored were removed from  
617 analysis. If a SNV was not present in the HDR library at a frequency over 1 in  $10^4$ , it was  
618 presumed to have been lost in oligo synthesis or cloning and was removed. Additionally, if a  
619 SNV was not observed with complete HDR markers at a frequency over over 1 in  $10^5$  in day 5  
620 genomic DNA samples from both replicate experiments, it was removed. SNVs introduced near  
621 the CRISPR recognition site have the potential to facilitate Cas9 recutting of the locus (e.g. by  
622 replacing the PAM edit or introducing an alternative PAM site). Because these SNVs are likely  
623 to score lower consequent to Cas9 editing biases and not their effects on gene function, SNVs  
624 were filtered that created increased potential for re-cutting as follows: When an HDR marker  
625 mutation used to disrupt editing occurred at position 2 of the PAM (e.g. ‘NGG’ to ‘NCG’),  
626 SNVs that replaced this marker with an alternate base were removed to prevent biases introduced  
627 by recutting non-canonical *S. pyogenes* Cas9 PAMs (e.g. ‘NAG’, ‘NTG’). Additionally, variants  
628 that created a new PAM 1 bp 3’ of the mutated PAM were excluded due to the potential for



629 recutting (*e.g.* unedited PAM: 5'-NGGA, edited PAM with HDR marker: 5'-NCGA, filtered out  
630 SNV that creates *new PAM* +1bp 3': 5'-NCGG). (Extended Data Fig. 6 describes recutting  
631 observed at alternative PAMs.) To prevent misinterpretation, we also removed SNVs that created  
632 amino acid changes specific to the context of the library's fixed edits (*e.g.* if in the unedited  
633 background, the SNV causes an X to Y change, but with a fixed edit in the same codon, the SNV  
634 causes an X to Z change). We also applied this logic to remove SNVs that introduced splice  
635 donor sites only in the context of the edited PAM, and SNVs that create splice donor sites in the  
636 unedited context but not in the context of the edited PAM.

637 The RNA scores for exon 18 samples were neither well correlated across replicates nor  
638 with SNV abundances in genomic DNA, indicating likely bottlenecking in library preparation.  
639 Therefore, RNA data from exon 18 was excluded. WT HAP1 function scores from exon 22 were  
640 excluded because there was an unusually high correlation between SNV frequencies sampled  
641 from the plasmid library and from day 5 gDNA, suggesting plasmid contamination in gDNA  
642 sequencing. This problem was fixed by designing a new primer to prepare gDNA sequencing  
643 samples from HAP1-Lig4KO cells.

644

#### 645 External data sources

646 Variant annotations were downloaded from CADD<sup>39</sup> version 1.3  
647 (<http://cadd.gs.washington.edu/download>). This included the following scores: mammalian  
648 phyloP, Grantham deviation, SIFT, Polyphen-2, and CADD. Align-GVGD scores were obtained  
649 by running the Align-GVGD program on BRCA1 sequences conserved to sea urchin. ClinVar  
650 data were downloaded on 1/2/2018 for all germline SNVs with at least a 1-star annotation. SNVs  
651 annotated as 'Benign/Likely benign' were grouped with 'Likely benign' SNVs and SNVs  
652 classified 'Pathogenic/Likely pathogenic' were grouped with 'Likely pathogenic' SNVs. SNV  
653 allele frequencies were obtained from <http://gnomad.broadinstitute.org/> on 12/26/2017 for  
654 gnomAD<sup>16</sup>, from <https://bravo.sph.umich.edu/freeze5/hg38/> on 11/19/2017 for Bravo, and from  
655 <https://whi.color.com/> on 10/9/2017 for FLOSSIES data. Transcript data was obtained from  
656 GTEx on 1/3/2018. Throughout this study, *BRCA1* exons, coding nucleotide positions, and  
657 amino acid positions are referenced by the ClinVar transcript annotation for *BRCA1*, transcript  
658 NM\_007294.3 (NCBI).

659

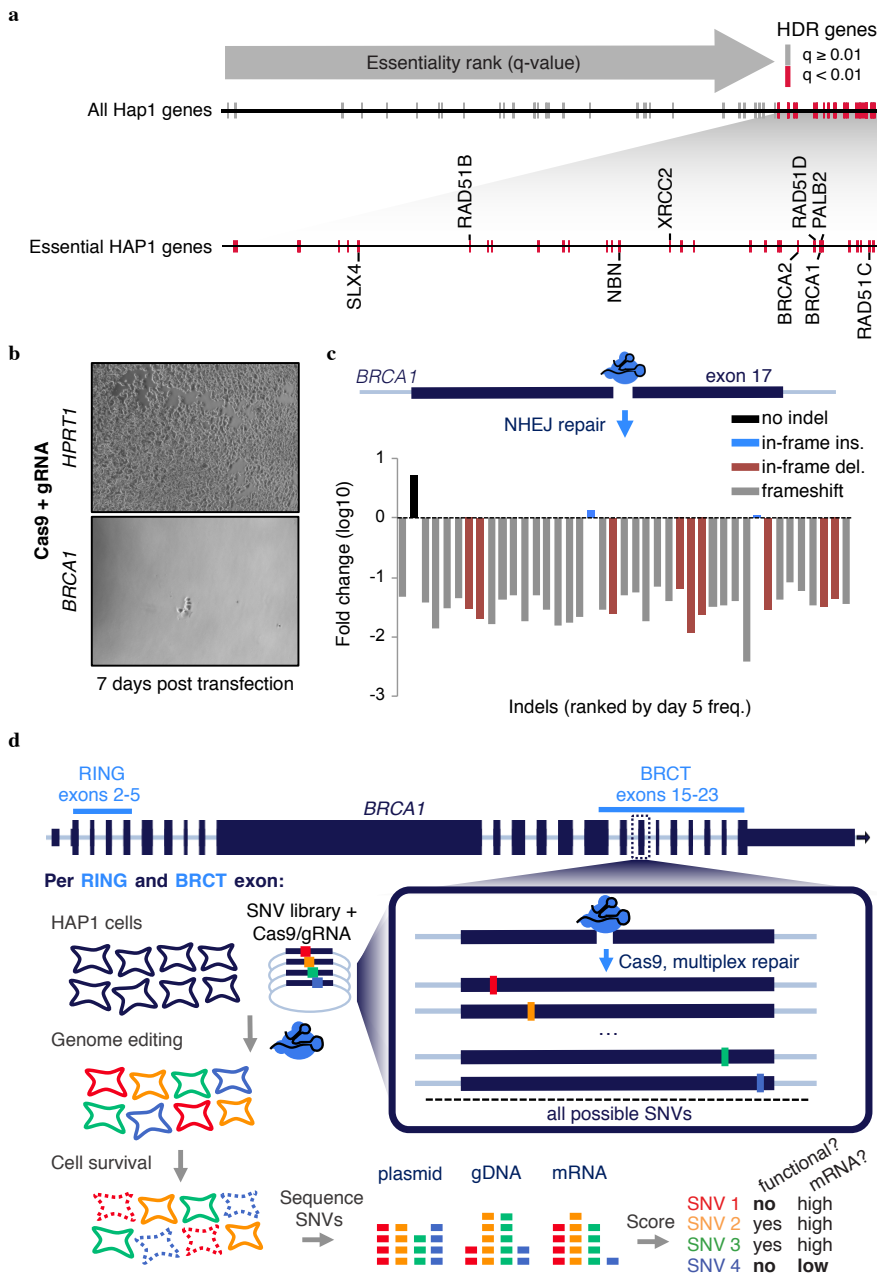
#### 660 Statistical reporting

661 All statistical tests described were performed as two-tailed tests using the R software  
662 package.

663

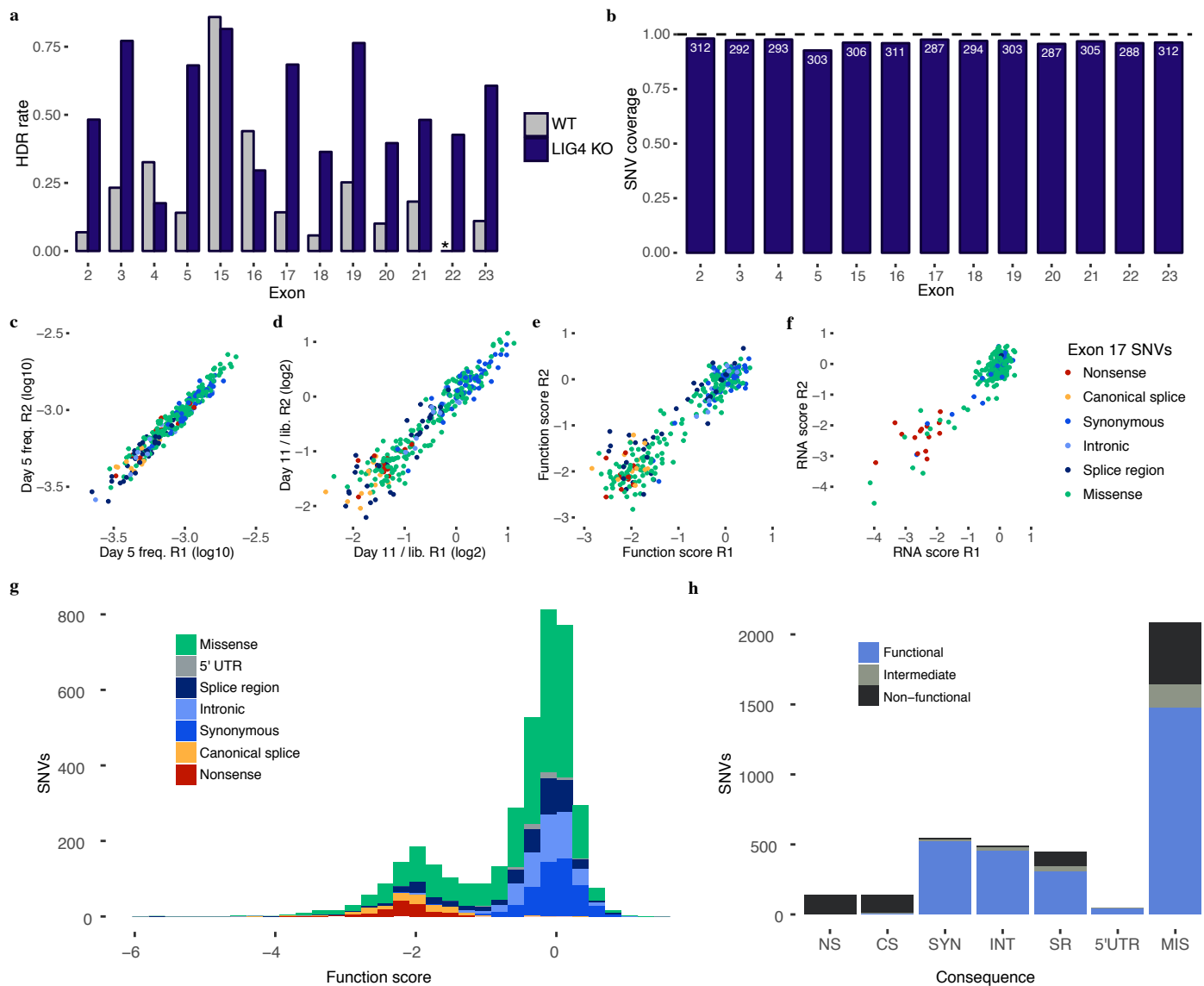
#### 664 Code availability

665 Custom scripts for analyzing sequencing data were written in Python and R. All code will  
666 be made available upon request.



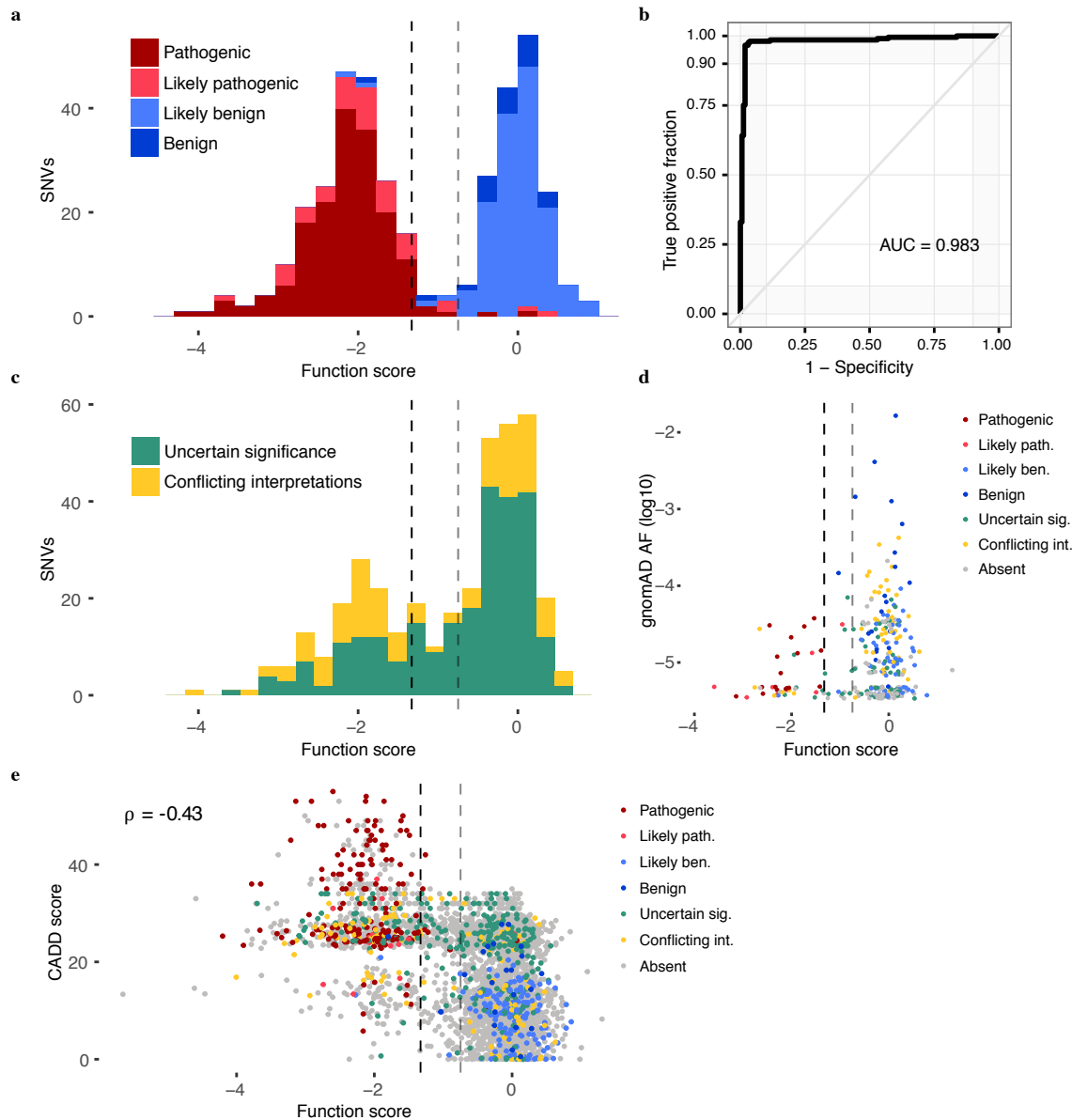
**Figure 1 | *BRCA1* and other HDR pathway genes are essential in HAP1 cells.** **a**, The q-value rankings of HDR pathway genes ( $N = 66$ , defined by Gene Ontology) among 14,306 genes scored in a HAP1 gene trap screen for essentiality<sup>28</sup> are indicated with tick marks. Essential HDR genes are colored red and those implicated in cancer predisposition are labelled in the enlargement below. Of the 66 HDR pathway genes scored, 34 including *BRCA1* were ‘essential’, a 3.4-fold enrichment compared to non-HDR genes (Fisher’s exact  $P = 6.1 \times 10^{-12}$ ). **b**, HAP1 cell populations were transfected with a Cas9/gRNA plasmid either targeting the non-essential gene *HPRT1* (control) or exon 17 of *BRCA1* on day 0. Successfully transfected cells were selected with puromycin (days 1-4) and cultured until day 7, at which point cells were washed prior to imaging. Images are representative of two transfection replicates. **c**, The targeted *BRCA1* exon 17 locus was deeply sequenced from a population of transfected cells sampled on day 5 and day 11. The fold-change from day 5 to day 11 for each editing outcome observed at a frequency over 0.001 in day 5 sequencing reads is plotted. All alleles but indel-free sequences and two in-frame insertions were depleted. **d**, Saturation genome editing experiments were designed to introduce all possible SNVs across thirteen *BRCA1* exons encoding the protein’s RING (exons 2-5) and BRCT domains (exons 15-23). For each exon, a Cas9/gRNA construct was designed to be transfected with a library of plasmids containing all SNVs across ~100 bp of genomic sequence (the ‘SNV library’). SNV libraries were designed to saturate a total of 1,345 bp of genomic sequence, spanning BRCT and RING domain coding regions and adjacent intronic sequences. SNV library plasmids contain homology arms to mediate genomic integration, as well as fixed synonymous variants within the CRISPR target site to prevent Cas9 re-cutting. Upon HAP1 cell transfection of each Cas9/gRNA plasmid / SNV library pair, successfully edited cells harbor a single *BRCA1* SNV from the library. Cells are sampled 5 and 11 days after transfection and targeted gDNA and RNA sequencing is performed to quantify SNV abundances. SNVs compromising *BRCA1* function are selected against, manifesting in reduced gDNA representation, and SNVs impacting mRNA production are depleted in RNA samples relative to gDNA.

**Figure 2**



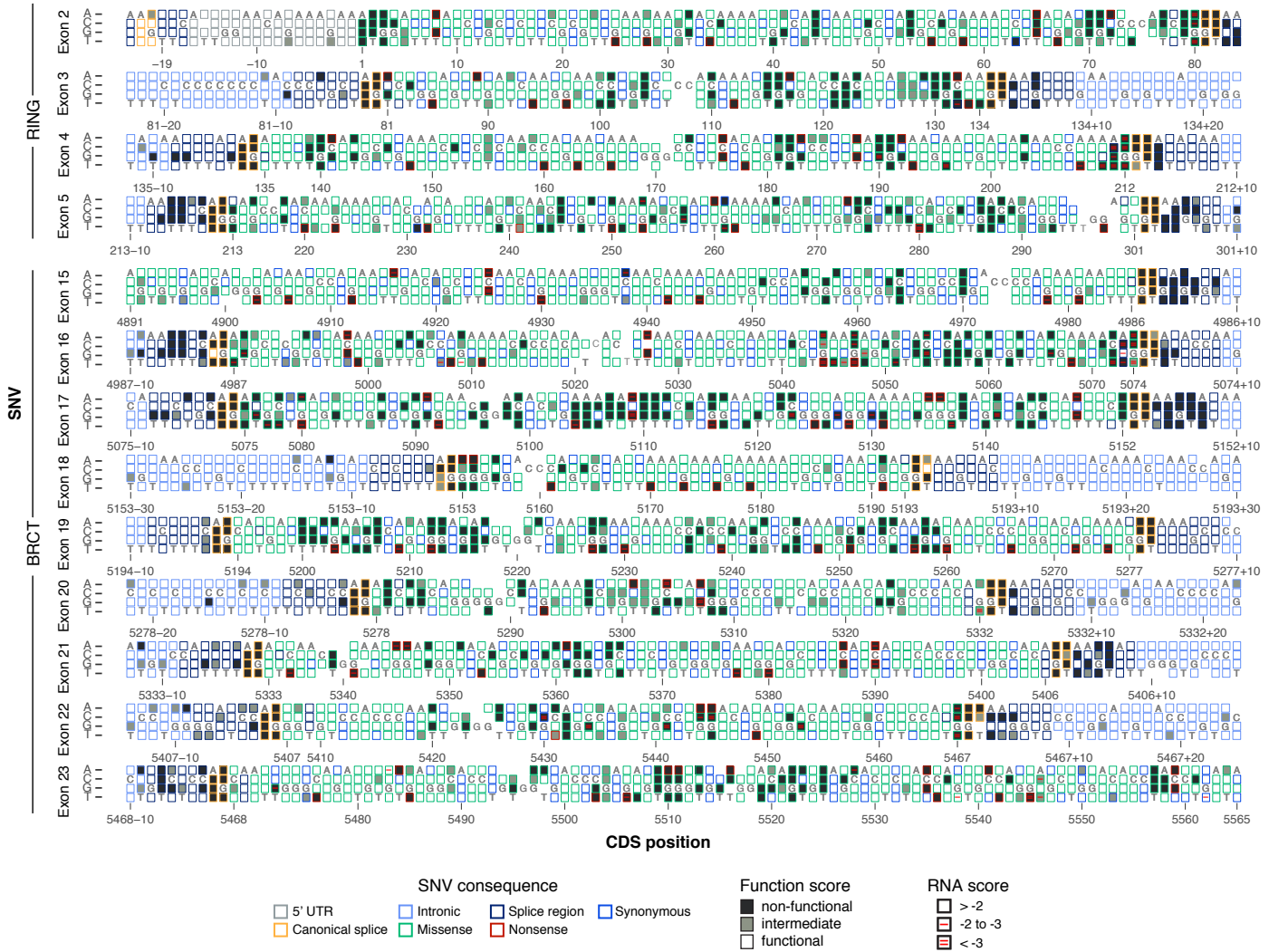
**Figure 2 | Saturation genome editing enables functional classification of 3,893 *BRCA1* SNVs.** **a**, HDR editing rates were calculated for each exon as the fraction of day 5 reads containing the SNV library's fixed synonymous variant (*i.e.* an 'HDR marker' edit). The average of two WT HAP1 replicates and two HAP1-Lig4KO replicates is plotted for comparison. (Asterisk denotes missing exon 22 data.) **b**, The fraction of all possible SNVs scored is shown for each exon. SNVs were excluded mainly due to proximity to the HDR marker and/or poor sampling (Extended Data Fig. 6 and Methods). **c-f**, Reproducibility was assessed across all exon replicates (Extended Data Fig. 5). Measurements for exon 17 SNVs assayed in HAP1-Lig4KO cells are plotted to show correlations of day 5 frequencies (**c**,  $\rho = 0.97$ ), day 11 over library ratios (**d**,  $\rho = 0.95$ ), function scores (**e**,  $\rho = 0.88$ ), and RNA expression scores (**f**,  $\rho = 0.61$ ). **g**, A histogram of 3,893 SNV function scores (averaged across replicates and normalized across exons) shows how each category of mutation compares to the overall distribution. **h**, The number of SNVs within each category of mutation is plotted and colored by functional classification determined by SGE. (NS = nonsense, CS = canonical splice, SYN = synonymous, INT = intronic, SR = splice region, 5'UTR = 5' untranslated region, MIS = missense.)

**Figure 3**



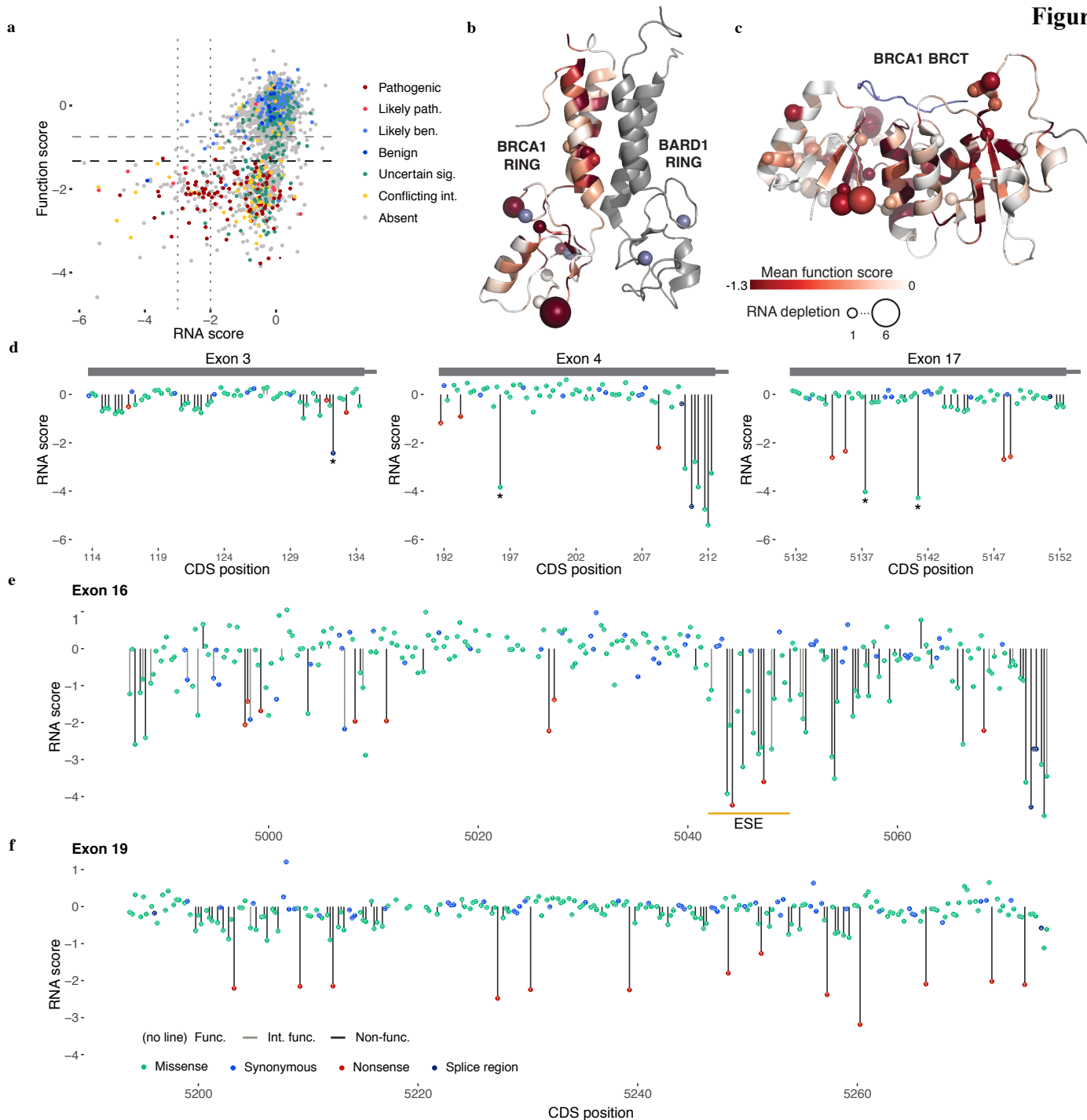
**Figure 3 | SGE function scores are highly accurate at predicting clinical interpretations of *BRCA1* SNVs.** **a**, The distribution of SNV function scores colored by ClinVar interpretation. Scores are shown for the 375 SNVs with at least a ‘1-star’ review status in ClinVar and either a ‘pathogenic’ or ‘benign’ interpretation (including ‘likely’). The dashed lines indicate the functional classification thresholds determined by mixture modeling (gray = intermediate, black = non-functional). **b**, An ROC curve reveals optimal sensitivity and specificity for classifying the same 375 SNVs in **a** at SGE function score cutoffs from -1.03 to -1.22. **c**, The distribution of scores plotted as in **a** for the 378 SNVs annotated as variants of uncertain significance or with conflicting interpretations. 91.3% of such variants are classified as ‘functional’ or ‘non-functional’ by SGE. **d,e**, SNVs are colored by ClinVar annotation. **d**, Among the 302 SNVs assayed also present in gnomAD, higher allele frequencies associated with higher function scores (Wilcoxon Signed Rank Test,  $P = 3.7 \times 10^{-12}$ ). **e**, CADD scores (which predict deleteriousness) inversely correlate with function scores.

Figure 4



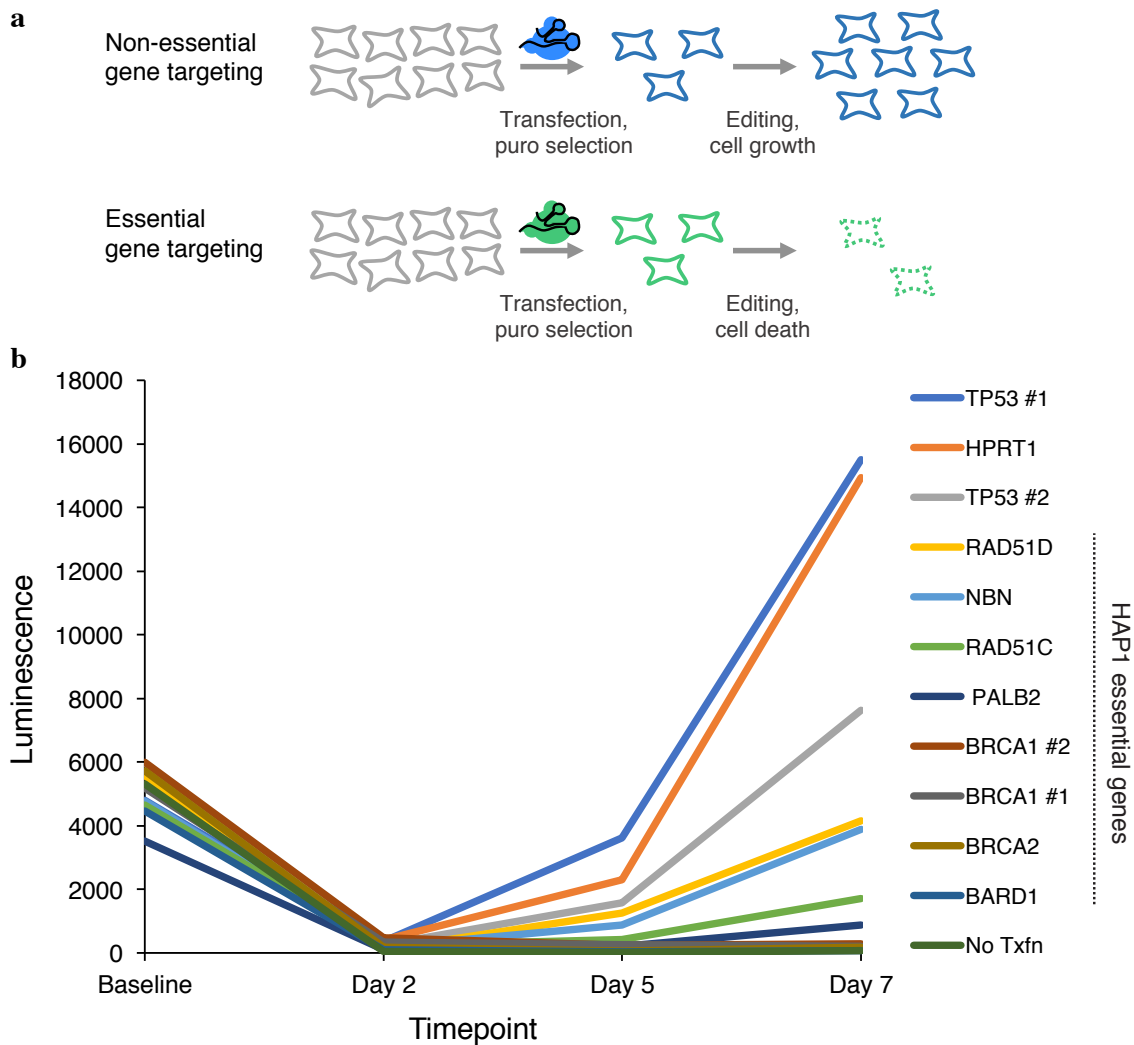
**Figure 4 | Sequence-function maps for 13 *BRCA1* exons.** The 3,893 SNVs scored with SGE are each represented by a box corresponding to coding sequence position (NCBI transcript ID: NM\_007294.3) and nucleotide identity. Boxes are filled corresponding to functional class, and outlined corresponding to the SNV's mutational consequence. Red lines within boxes mark SNVs depleted in RNA; one line indicates an RNA score between -2 and -3 (log<sub>2</sub> scale) and two lines indicate a score below -3. RNA measurements were determined only for exonic SNVs, excluding exon 18. Reference nucleotides are indicated by dark gray letters; blank boxes indicate missing data.

**Figure 5**



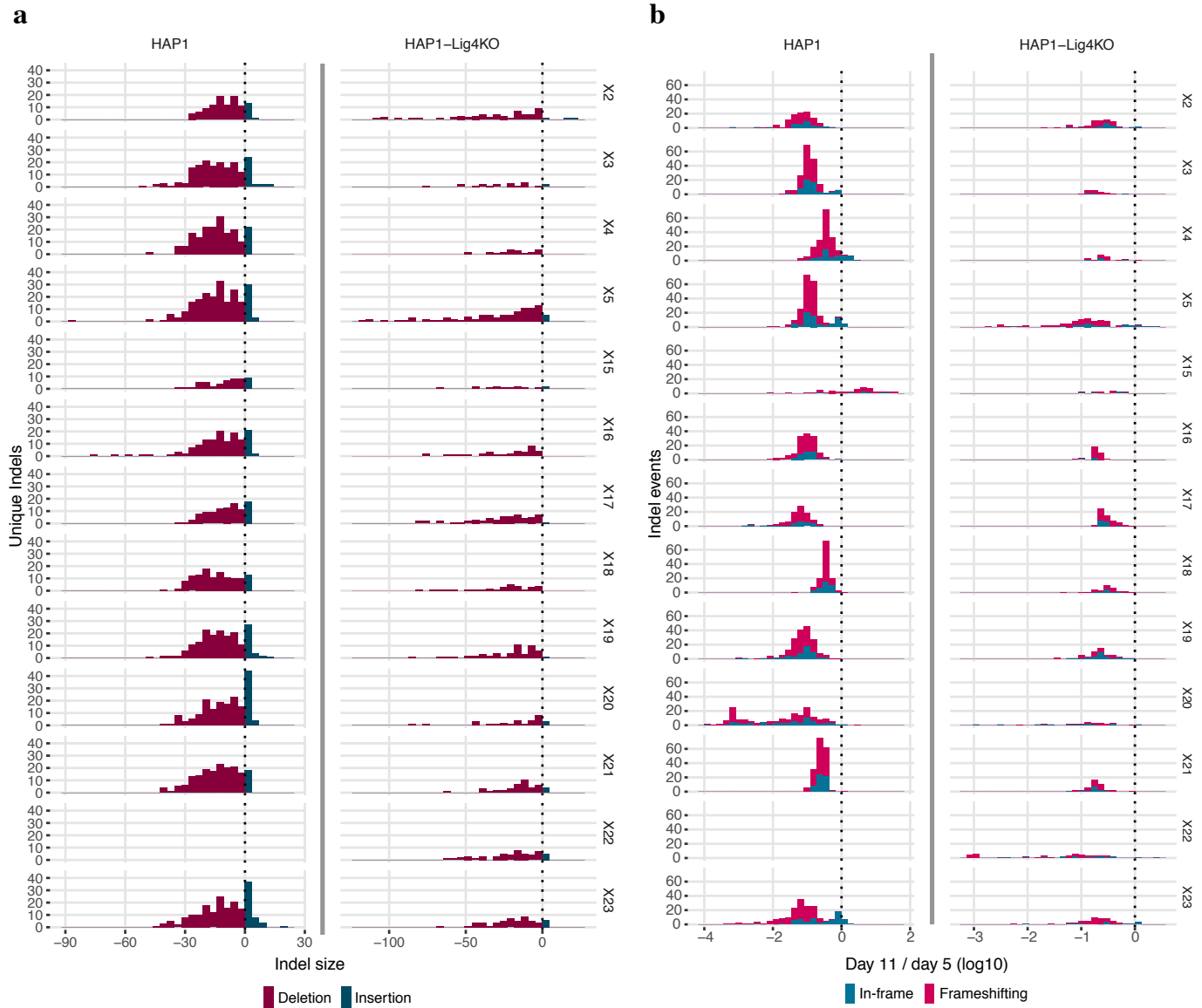
**Figure 5 | Measuring SNV mRNA abundance and function in parallel delineates mechanisms of variant effect.** **a**, Function scores are plotted against RNA scores for all exonic synonymous and missense SNVs scored ( $N = 2,646$ ). Horizontal dashed lines indicate functional thresholds, and vertical dotted lines mark RNA scores of -2 and -3. **b,c**, Function scores for all SNVs were mapped onto the structures of the RING (**b**, pdb 1JM7) and BRCT (**c**, pdb 1T29) domains in shades of red by averaging missense SNV scores at each amino acid position. The number of SNVs that cause >75% reduction in RNA levels at each amino acid position is represented by the size of the sphere at the alpha-carbon at each residue. Grey denotes residues not assayed and the BACH1 peptide bound to the BRCT structure is colored slate blue. **d,e,f**, SNV RNA scores are plotted by transcript position, with lines denoting SNV functional classification. **d**, Examples of non-functional SNVs with low RNA scores that create new 5'-GU splice donor motifs are shown. Complete maps of RNA scores for exons 16 (**e**) and exon 19 (**f**) reveal highly variable sensitivity to RNA depletion. The location of the strongest predicted exonic splice enhancer in exon 16<sup>42</sup> is indicated by the orange line (**e**).

## Extended Data Figure 1



**Extended Data Figure 1 | CRISPR targeting of HDR pathway genes to confirm essentiality in HAP1 cells.** **a**, Schematic; HAP1 cells are transfected with a plasmid expressing a gRNA and a Cas9-2A-puromycin cassette<sup>29</sup>. Due to low transfection rates for HAP1 cells, puromycin selection reduces viable cells in all transfections. Over time, however, CRISPR targeting of non-essential genes leads to increased cell growth compared to CRISPR targeting of essential genes. **b**, Cell viability of HAP1 cells transfected with Cas9/gRNA constructs targeting different HDR genes and controls (*HPRT1*, *TP53*) was measured using the CellTiterGlow assay. Luminescence is proportional to the number of living cells in each well when the assay is performed. Triplicate wells for each gRNA at each time point were processed, quantified on a plate reader and averaged. gRNA sequences are included in Supplementary Table 2.

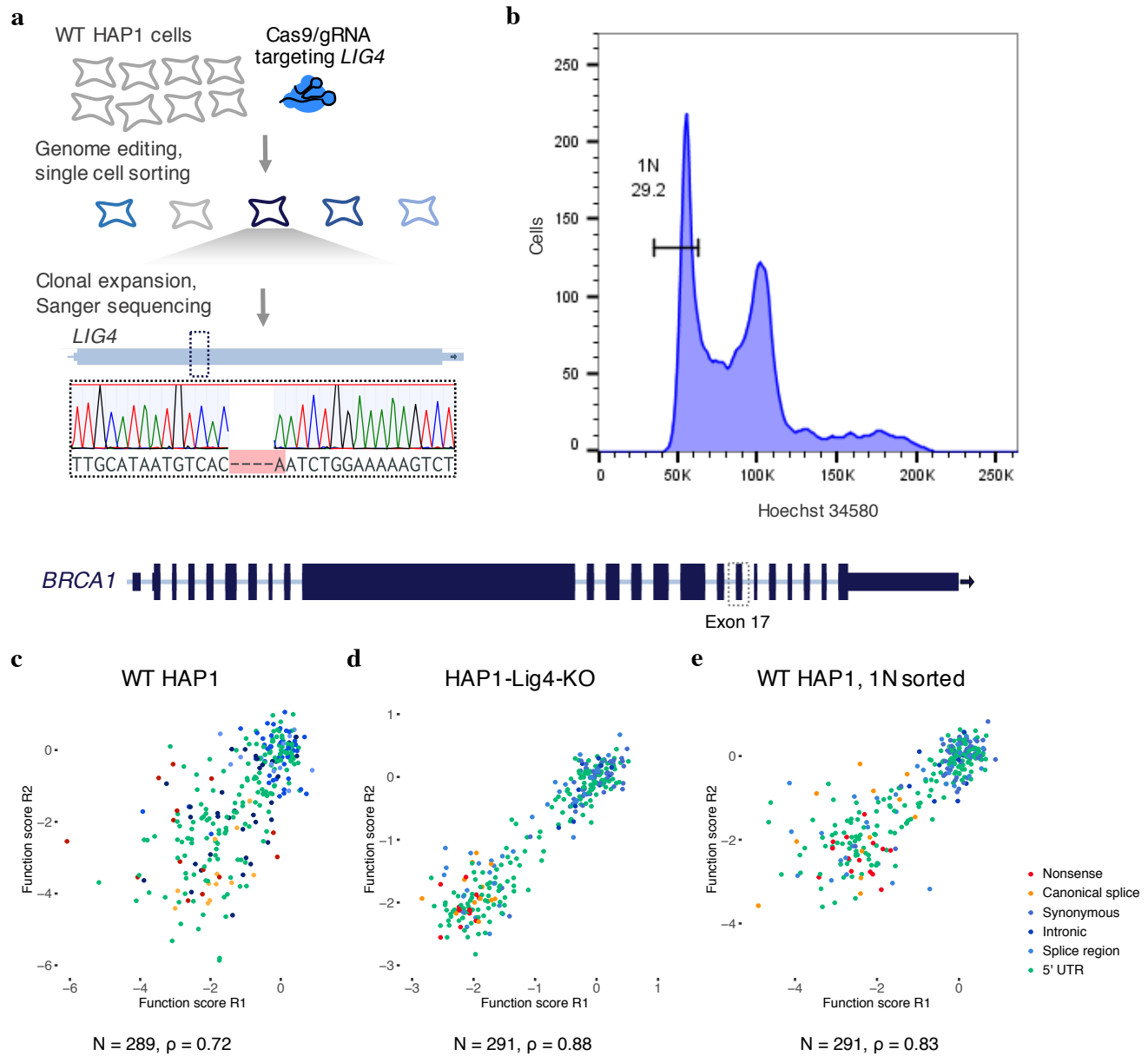
## Extended Data Figure 2



**Extended Data Figure 2 | Analysis of Cas9-induced indels observed in *BRCA1* SGE experiments.** Variants observed in gDNA sequencing were included in this analysis if i) they aligned to the reference with either a single insertion or deletion within 15 bp of the predicted Cas9 cleavage site and ii) were observed at a frequency greater than 1 in 10,000 reads in both replicates. **a**, Histograms show the number of unique indels observed of each size, with negative sizes corresponding to deletions. More unique indels were observed in WT HAP1 cells compared to HAP1-Lig4KO cells for exons compared (WT data for exon 22 was excluded). **b**, Day 11 over day 5 indel frequencies were normalized to the median synonymous SNV in each replicate and then averaged across replicates to measure selection on each indel. The distribution of selective effects is shown for each experiment as a histogram, in which indels are colored by whether their size was divisible by 3 (*i.e.* ‘in-frame’ vs. ‘frameshifting’). Whereas frameshifting variants were consistently depleted, some exons were tolerant to in-frame indels.

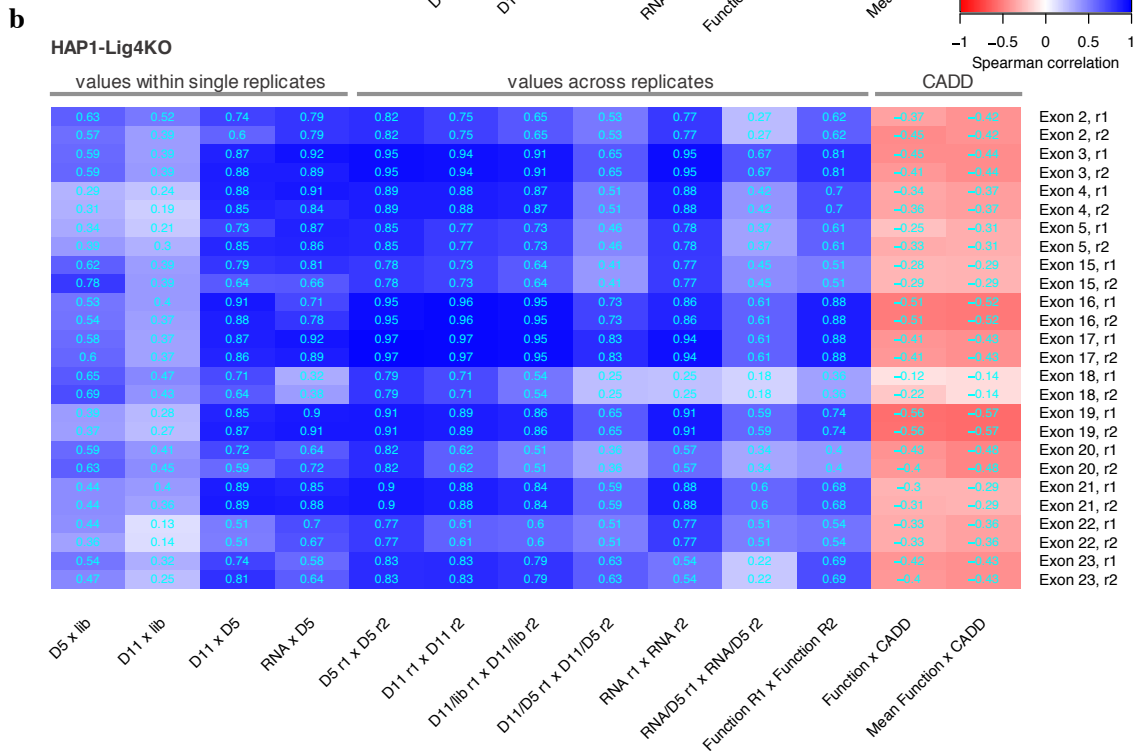
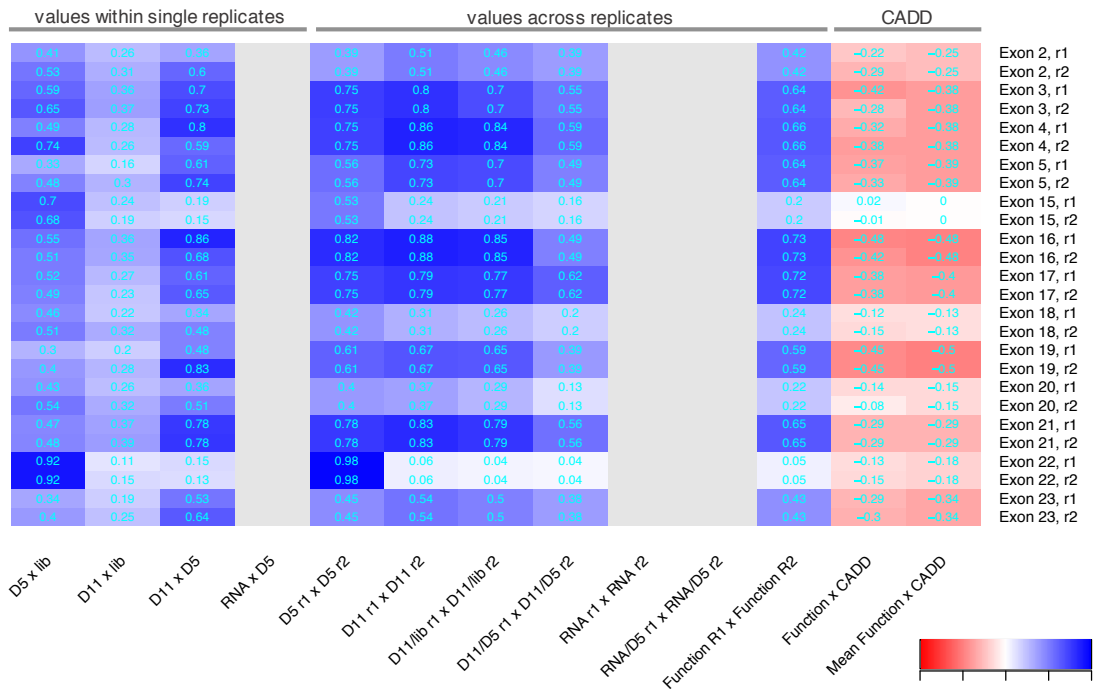


### Extended Data Figure 3



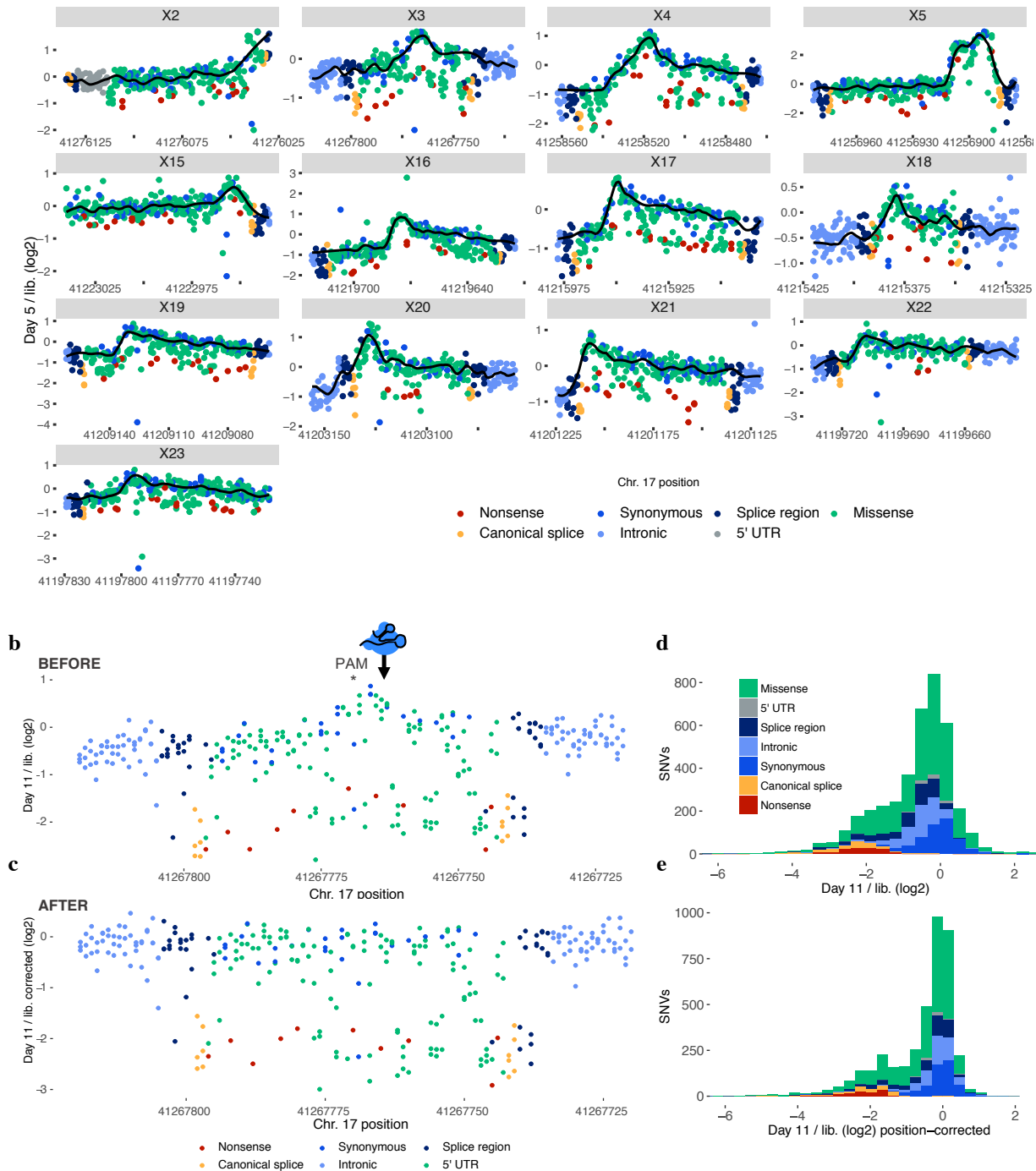
**Extended Data Figure 3 | HAP1 cell line optimizations for saturation genome editing to assay essential genes.** **a**, A gRNA targeting Cas9 to the coding sequence of *LIG4*, a gene integral to the non-homologous end-joining pathway, was cloned into a vector co-expressing Cas9-2A-GFP<sup>29</sup>. WT HAP1 cells were transfected, and single GFP-expressing cells were sorted into wells of a 96-well plate. Eight monoclonal lines were grown out over a period of three weeks and screened using Sanger sequencing for frameshifting indels in *LIG4*. The Sanger trace shows the frameshifting deletion present in the clonal line chosen for subsequent experiments, referred to as ‘HAP1-Lig4KO’. **b**, To purify HAP1 cells for haploid cells, live cells were stained for DNA content with Hoechst 34580 and sorted using a gate to select cells with the lowest DNA content, corresponding to 1N cells in G1. **c-e**, Plots comparing SNV function scores across replicate experiments for exon 17 saturation genome editing experiments performed in unsorted WT HAP1 cells (**c**), HAP1-Lig4KO cells (**d**), and WT HAP1 cells sorted on 1N ploidy (**e**). Both *LIG4* knockout and 1N-sorting improved replicate correlations.

**a** WT HAP1 **Extended Data Figure 4**



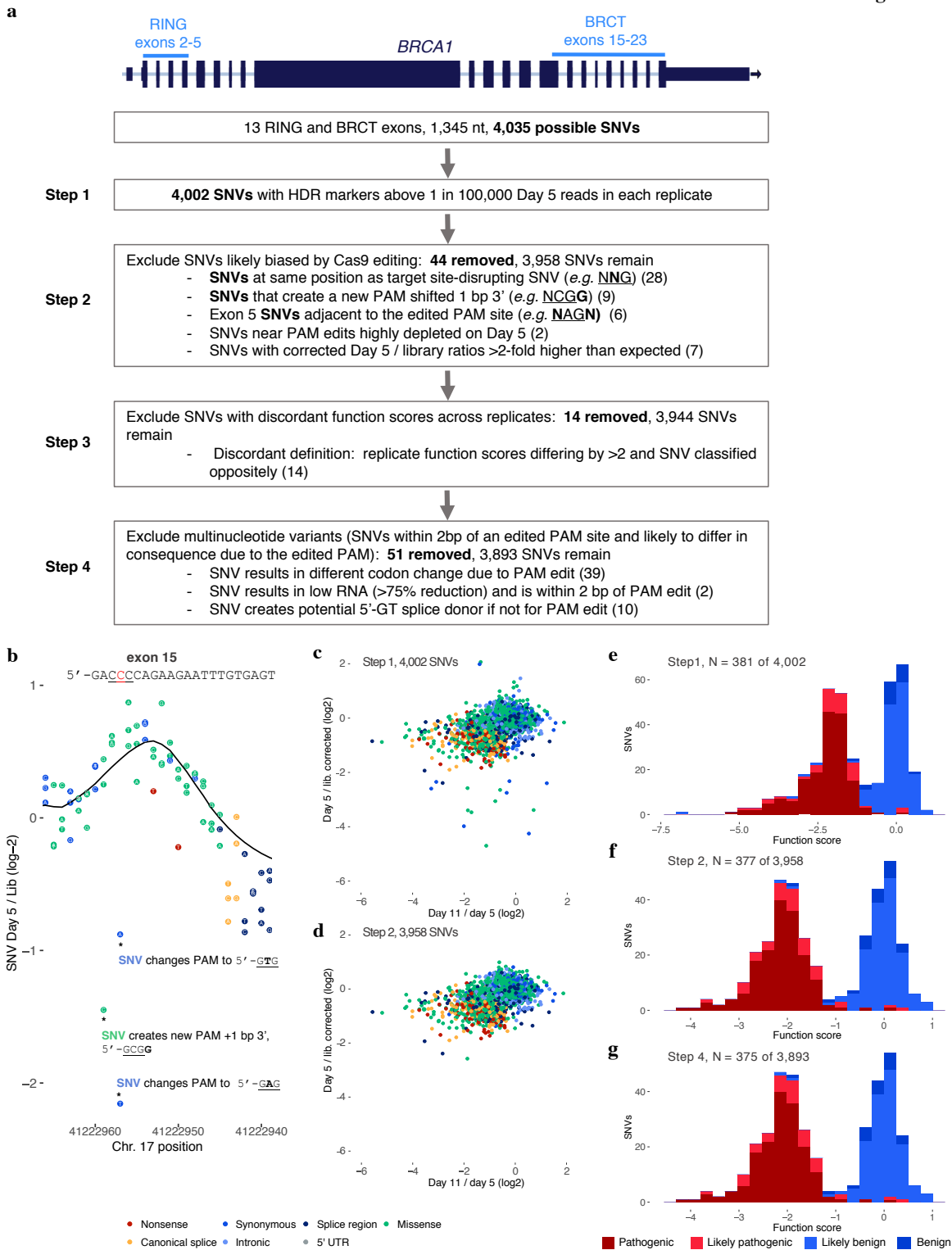
**Extended Data Figure 4 | Correlations for SNV measurements within single experiments, across transfection replicates, and to CADD scores for all SGE experiments.** Heatmaps indicate Spearman correlation coefficients for SNV measurements from experiments in WT HAP1 cells (**a**) and in HAP1-Lig4KO cells (**b**). Gray boxes indicate absent RNA data from WT HAP1 cells. The four leftmost columns show how SNV frequencies correlate between samples from within a single replicate experiment. The unusually high correlations between exon 22 SNV frequencies in the plasmid library and in day 5 gDNA samples from WT HAP1 cells suggests plasmid contamination in gDNA. Indeed, primer homology to a repetitive element in the exon 22 library was identified. Consequently, the WT HAP1 exon 22 data was removed from analysis and a different primer specific to gDNA was used to prepare exon 22 sequencing amplicons from HAP1-Lig4KO cells. The low HAP1-Lig4KO correlations between exon 18 SNV frequencies in day 5 gDNA and RNA and between RNA replicates suggests RNA sample bottlenecking consequential to low RNA yields. Therefore, exon 18 RNA was also excluded from analysis. Consistent with the higher rates of HDR-mediated genome editing (**Fig. 2a**), replicate correlations (middle columns) were generally higher in HAP1-Lig4KO cells than WT HAP1 cells. CADD scores predict the deleteriousness of each SNV, and are therefore negatively correlated with function scores (rightmost columns).

## Extended Data Figure 5

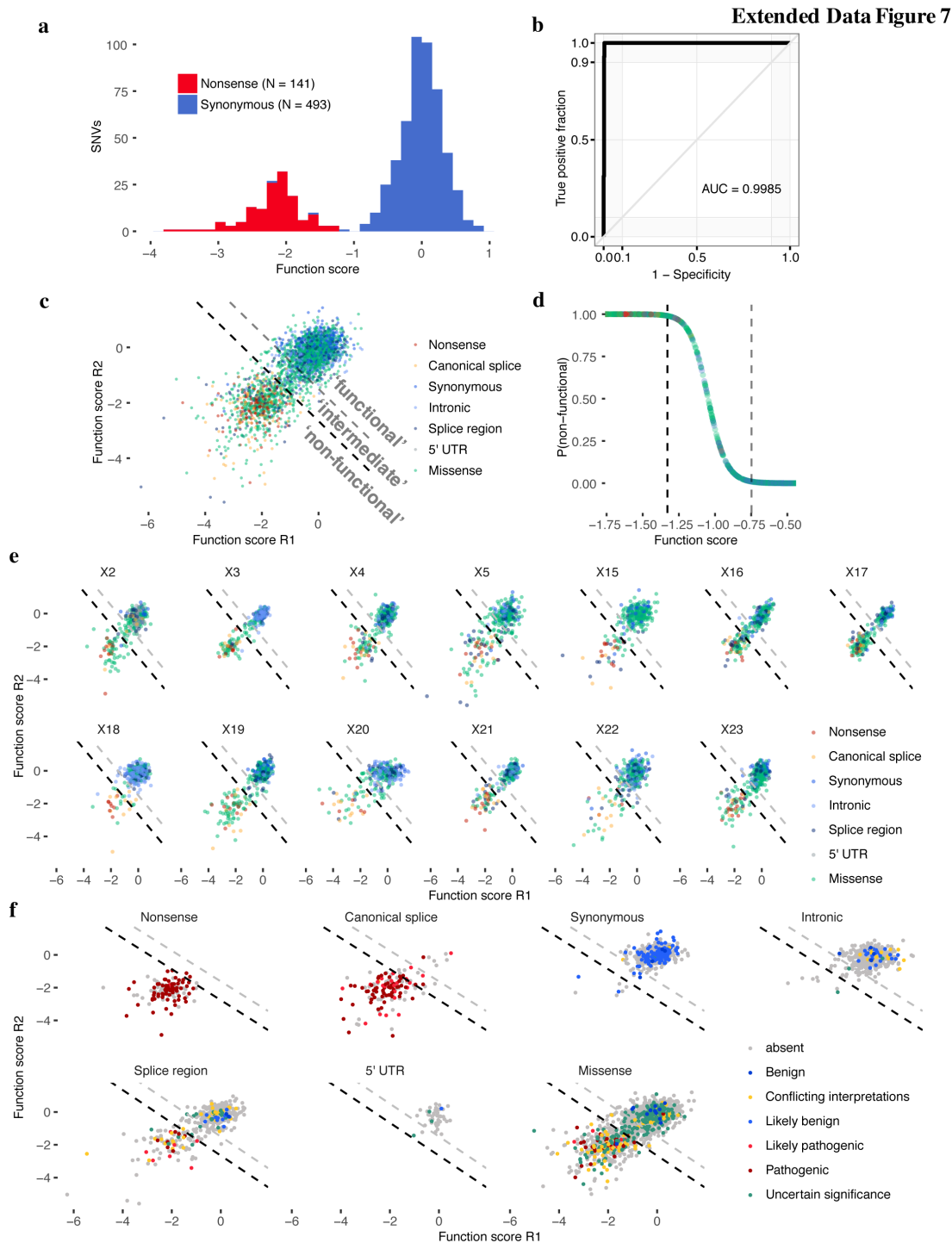


**Extended Data Figure 5 | Models of SNV editing rates across *BRCA1* exons account for positional biases.** **a**, Gene conversion tracts arising during HDR in human cells are short such that library SNVs are introduced to the genome more frequently near the CRISPR target site. We modelled this positional effect in our data using a LOESS regression fit on day 5 over library SNV ratios. Plots shown here are of the average of two replicate experiments per exon, with the black line indicating the LOESS regression. By day 5 sampling, selective effects on gene function are evidenced by nonsense SNVs (red) appearing at lower frequencies compared to neighbouring SNVs. Therefore, to best approximate the SNV editing rate as a function of position alone (*i.e.* the ‘baseline’), the regression excluded SNVs that were selected against between day 11 and day 5 (see Methods). **b,c**, Day 11 over library SNV ratios were adjusted by the positional fit for each experiment in calculating function scores. This adjustment is illustrated here for an exon 3 replicate by plotting the ratio as a function of position before (**b**) and after (**c**) adjustment. The elevated day 11 over library ratios for SNVs near the CRISPR target site are corrected to achieve a more uniform baseline across the mutagenized region. **d,e**, The distributions of SNV day 11 over library ratios before and after accounting for positional effects are shown, colored by mutational consequence (pre-filtering, N = 4,002).

Extended Data Figure 6

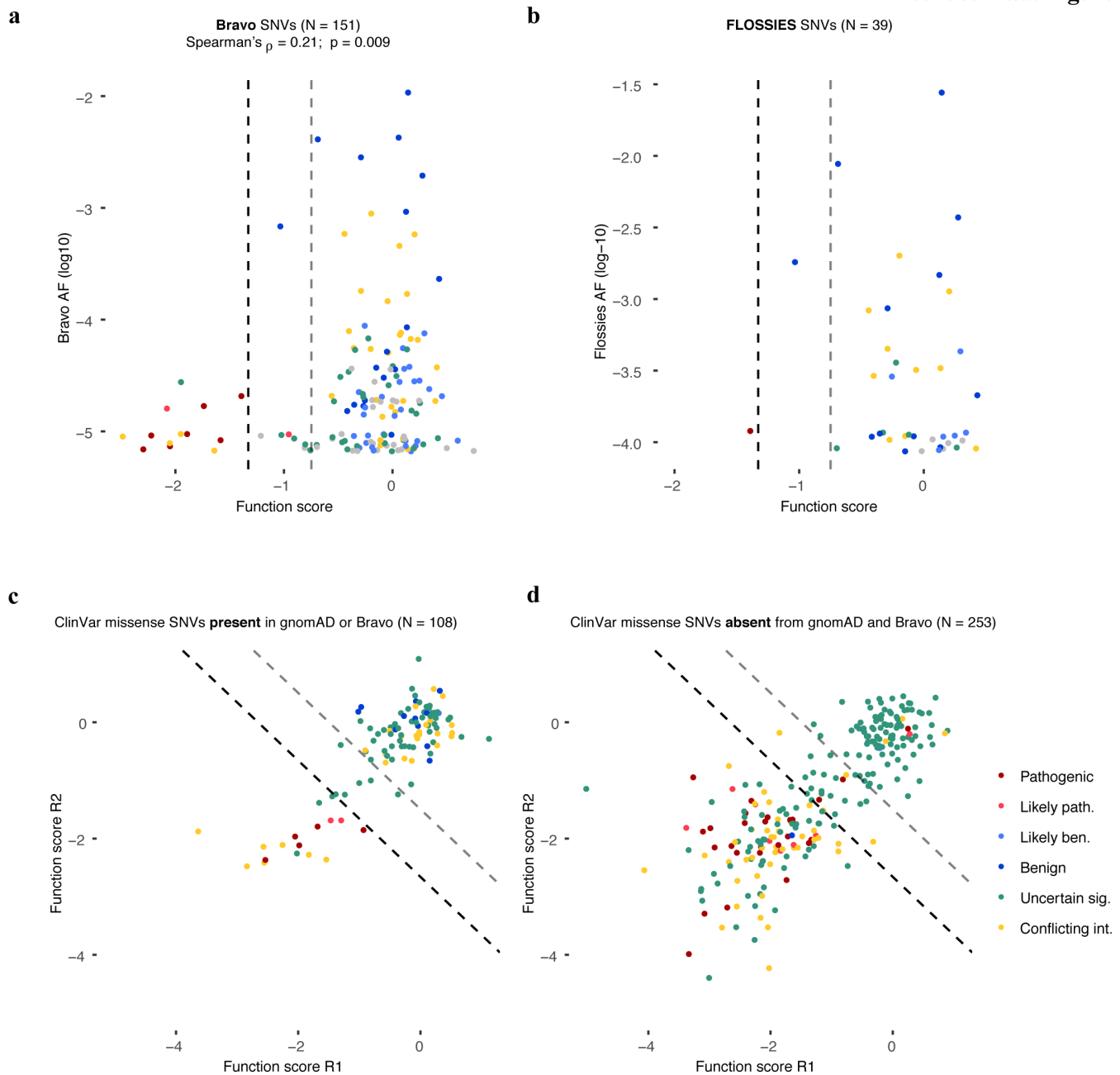


**Extended Data Figure 6 | SNV filtering to prevent erroneous functional classification.** **a**, The flow chart describes filters used to produce the final SNV data set and shows how many SNVs were removed at each step. **b**, Raw day 5 over library SNV ratios are shown for a portion of exon 15 to illustrate how re-editing biases necessitate filtering. The three depleted SNVs marked with asterisks create alternative PAM sequences that likely allow the Cas9:gRNA complex to re-cut the locus and cause their removal. For other SNVs, the fixed PAM edit (a GGG to GCG synonymous change) minimalizes re-editing. The location of the target PAM is underlined and each indicated SNV is bolded in the annotations. The LOESS regression curve is shown in black. **c,d**, Plots show the relationship between day 5 over library and day 11 over day 5 ratios before (**c**) and after (**d**) filtering steps 1 and 2. Filtering removes outliers because editing biases primarily affect the day 5 over library ratio. **e-g**, Histograms show the distributions of function scores for SNVs deemed ‘pathogenic’ or ‘benign’ in ClinVar at different stages of filtering. Scores in **e** are derived prior to normalization across exons.



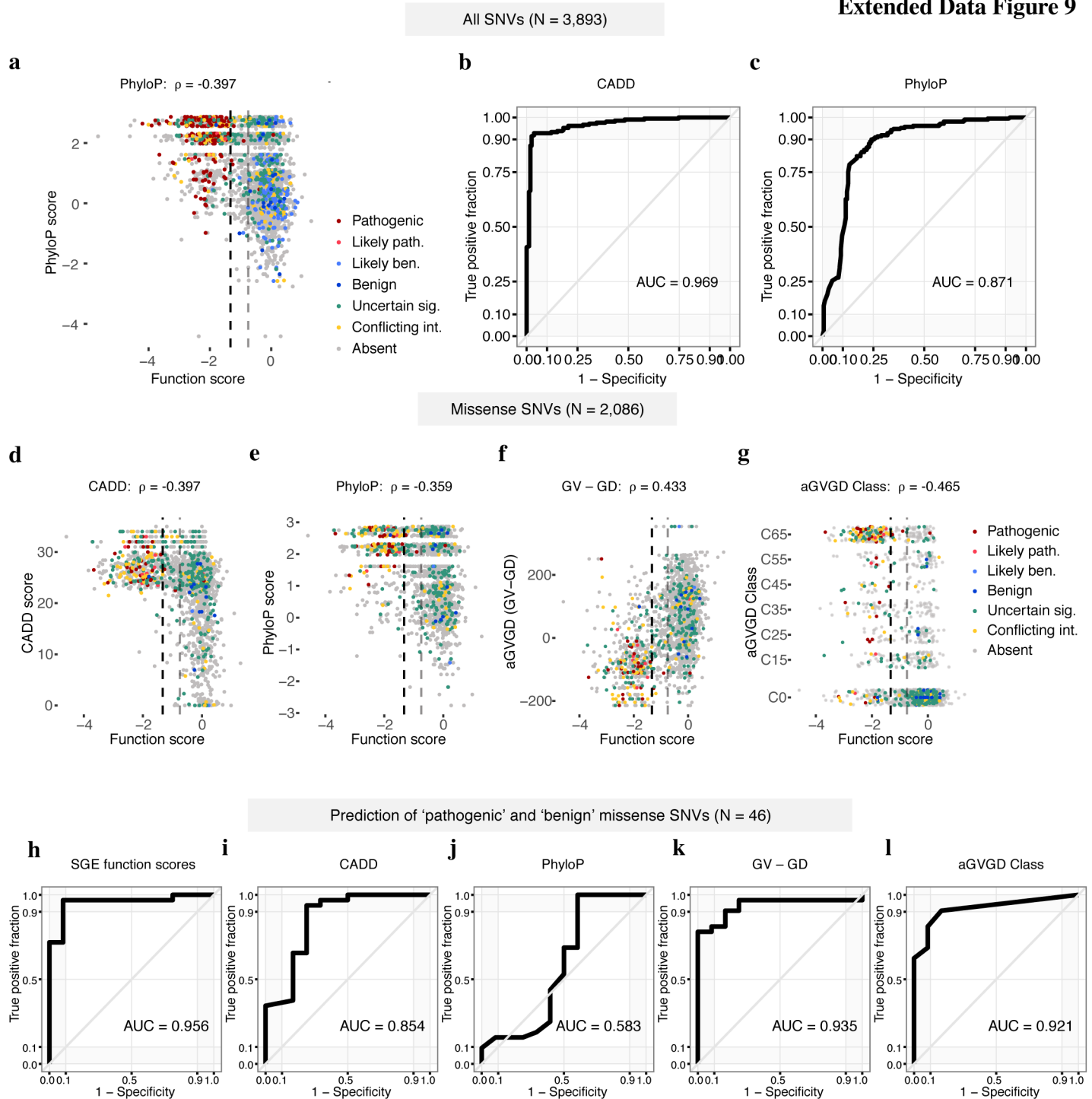
**Extended Data Figure 7 | Mixture modeling of scores to classify SNVs by functional effect.** **a**, Distributions of ‘non-functional’ and ‘functional’ SNVs plotted here were defined respectively as all nonsense SNVs and all synonymous SNVs with RNA scores within 1 SD of the median synonymous SNV. **b**, An ROC curve was generated using SGE function scores to distinguish the 634 ‘functional’ and ‘non-functional’ SNVs defined in **a**. **c**, A two-component Gaussian mixture model was used to produce point estimates of the probability that each SNV was ‘non-functional’,  $P(\text{nf})$ , given its average function score across replicates. These P-values are plotted in **d** against function scores for a subset of the data. Thresholds were set such that  $P(\text{nf}) < 0.01$  corresponds to ‘functional’, and  $P(\text{nf}) > 0.99$  corresponds to ‘non-functional’, and  $0.01 < P(\text{nf}) < 0.99$  corresponds ‘intermediate’ classification. Functional classification thresholds are drawn as dashed lines; black denotes the non-functional threshold and gray the intermediate threshold. **e,f**, SNV function scores across replicates are plotted for each exon with SNVs colored by mutational consequence (**e**), and for each type of mutational consequence with SNVs colored by ClinVar status (**f**). Using the optimal function score cutoff for all SNVs tested (Fig. 3b), sensitivities and specificities for distinguishing ‘Pathogenic’/‘Likely pathogenic’ from ‘Benign’/‘Likely benign’ ClinVar annotations for each type of mutation are as follows: 92.7% and 92.9% for missense SNVs (N = 55), 100% and 100% for splice region SNVs (N = 23), and 95.2% sensitivity for canonical splice site SNVs (N = 83; specificity not calculable).

## Extended Data Figure 8



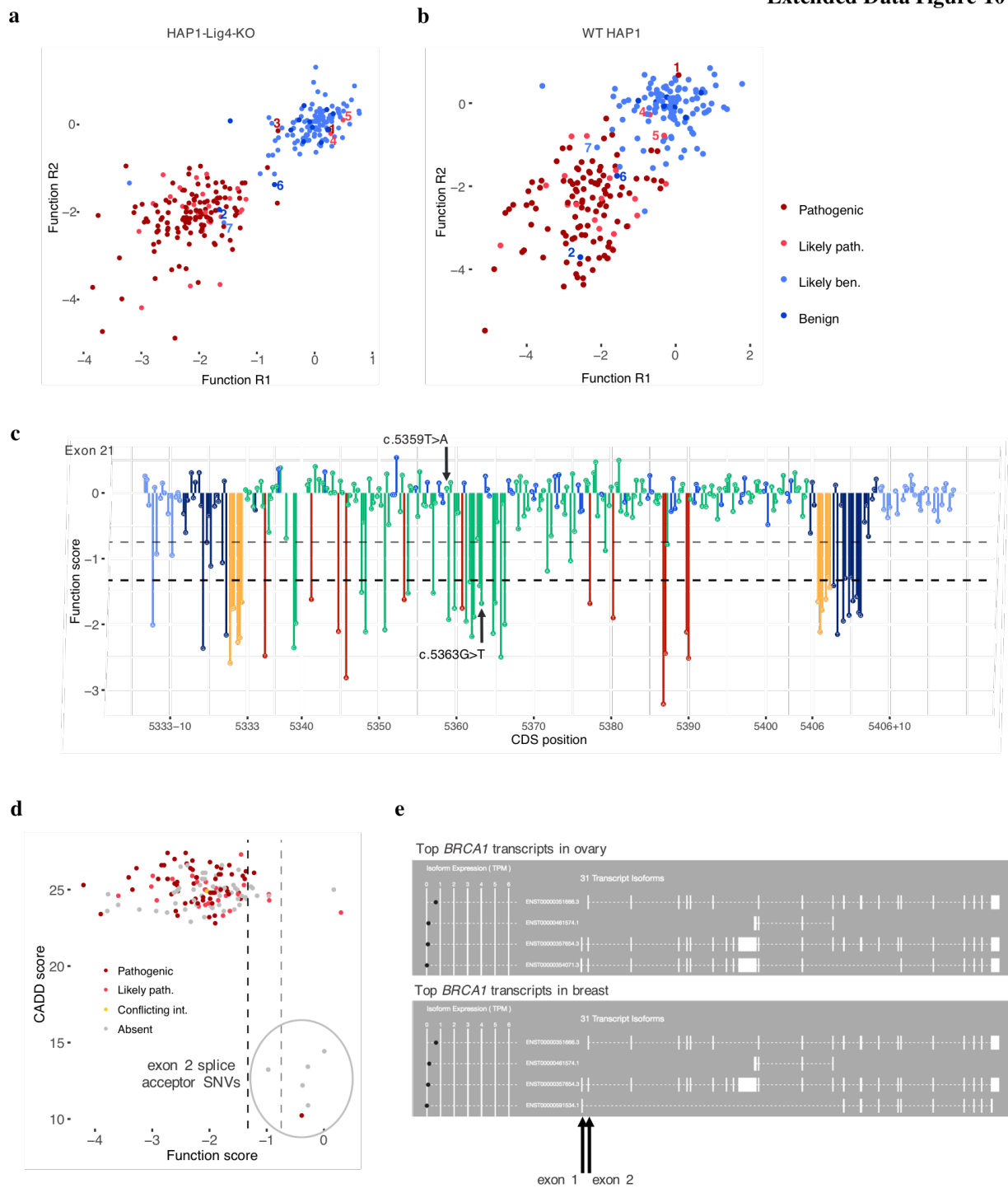
**Extended Data Figure 8 | *BRCA1* SNVs observed more frequently in large-scale population sequencing are more likely to score as functional.** SNV function scores are plotted against Bravo allele frequencies (**a**) and FLOSSIES allele frequencies (**b**). **a**, Bravo is a collection of whole genome sequences ascertained from 62,784 individuals through the NHLBI TOPMed program. Similarly to SNVs present in gnomAD (Fig. 3d), higher allele frequencies of SNVs in Bravo correlate with higher function scores. **b**, FLOSSIES is a database of variants seen in targeted sequencing of breast cancer genes sampled from approximately 10,000 cancer-free women at least 70 years old. Only 1 of 39 SNVs observed in FLOSSIES scored as non-functional. **c,d**, Missense SNVs in ClinVar are separated by whether they have (**c**) or have not (**d**) been seen in either gnomAD or Bravo and function scores across replicates are plotted, with dashed lines demarcating functional classes. A higher proportion of ClinVar missense SNVs absent from gnomAD and Bravo score as non-functional (50.6% vs. 15.7%, Fisher's exact  $P = 1.80 \times 10^{-17}$ ).

Extended Data Figure 9



**Extended Data Figure 9 | SGE function scores correlate with computational metrics and perform favorably at predicting ClinVar annotations.** **a**, SNV function scores are plotted against mammalian phyloP scores, with colors indicative of ClinVar status. **b,c**, ROC curves show the performance of CADD scores and phyloP scores for discriminating ClinVar 'pathogenic' and 'benign' SNVs (including 'likely'), as described in Fig. 3b for SGE data. **d-g** Plots as in **a**, but for missense SNVs only, showing correlations between SGE function scores and CADD<sup>39</sup> scores, phyloP scores<sup>40</sup>, Grantham differences (Grantham amino acid variation minus Grantham amino acid deviation; GV - GD), and align-GVGD classifications<sup>53</sup>. Missense SNV function scores also correlate with SIFT scores<sup>54</sup> ( $\rho = 0.363$ ) and PolyPhen-2 scores<sup>55</sup> ( $\rho = -0.277$ ). ( $P < 1 \times 10^{-37}$  for all correlations.) **h-l**, ROC curves assess the performance of SGE function scores and each indicated metric at distinguishing firmly 'pathogenic' and 'benign' missense SNVs. (*i.e.* not including 'likely').

Extended Data Figure 10



**Extended Data Figure 10 | Evidence supporting SNV scores in discordance with ClinVar classifications.** Function scores of SNVs classified as ‘benign’ or ‘pathogenic’ (including likely’s) are shown across replicates for experiments using HAP1-Lig4KO cells (**a**) and for preliminary experiments using WT HAP1 cells (**b**). Plots exclude exons with low overall reproducibility in WT HAP1 cells (replicate correlations < 0.4: exons 15, 18, 20 and 22). The three SNVs firmly discordant with ClinVar are labelled 1-3 in **a**, corresponding to c.5359T>A (dark red 1), c.5044G>A (dark blue 2), and c.-19-2A>G (dark red 3), respectively. The same filtering criteria were applied to both sets of experiments, which led to the removal of SNV 3 from the WT HAP1 data due to disagreement of scores between replicates. Discordant ‘likely pathogenic’ SNVs (4,5), an intermediate scoring ‘benign’ SNV (6) and a discordant ‘likely benign’ SNV (7) are also labelled for comparison. **c**, The sequence-function map of exon 21 is shown with the function scores for the two ‘pathogenic’ SNVs observed in linkage indicated. Dashed lines demarcate functional classifications. **d**, Function scores are plotted against CADD scores for all canonical splice SNVs assayed, colored by ClinVar status. The six possible exon 2 splice acceptor SNVs (circled) have the lowest CADD scores among all canonical splice SNVs assayed, and none score as ‘non-functional’. **e**, GTEx browser shots show that many of the most common *BRCA1* transcripts mapped from ovarian and breast tissues lack the exon 1 / exon 2 junction.



## REFERENCES

1. Cooper, G. M. Parlez-vous VUS? *Genome Res.* **25**, 1423–1426 (2015).
2. Rehm, H. L. *et al.* ClinGen — The Clinical Genome Resource. *N. Engl. J. Med.* **372**, 2235–2242 (2015).
3. Hall, J. M. *et al.* Linkage of early-onset familial breast cancer to chromosome 17q21. *Science* **250**, 1684–1689 (1990).
4. Miki, Y. *et al.* A strong candidate for the breast and ovarian cancer susceptibility gene BRCA1. *Science* **266**, 66–71 (1994).
5. Friedman, L. S. *et al.* Confirmation of BRCA1 by analysis of germline mutations linked to breast and ovarian cancer in ten families. *Nat. Genet.* **8**, 399–404 (1994).
6. Kuchenbaecker, K. B. *et al.* Risks of Breast, Ovarian, and Contralateral Breast Cancer for BRCA1 and BRCA2 Mutation Carriers. *JAMA* **317**, 2402–2416 (2017).
7. Olopade, O. I. & Artioli, G. Efficacy of risk-reducing salpingo-oophorectomy in women with BRCA-1 and BRCA-2 mutations. *Breast J.* **10 Suppl 1**, S5–9 (2004).
8. Rebbeck, T. R. *et al.* Bilateral prophylactic mastectomy reduces breast cancer risk in BRCA1 and BRCA2 mutation carriers: the PROSE Study Group. *J. Clin. Oncol.* **22**, 1055–1062 (2004).
9. Chan, S. L. & Mok, T. PARP inhibition in BRCA-mutated breast and ovarian cancers. *Lancet* **376**, 211–213 (2010).
10. Hollis, R. L., Churchman, M. & Gourley, C. Distinct implications of different BRCA mutations: efficacy of cytotoxic chemotherapy, PARP inhibition and clinical outcome in ovarian cancer. *Onco. Targets. Ther.* **10**, 2539–2551 (2017).
11. Farmer, H. *et al.* Targeting the DNA repair defect in BRCA mutant cells as a therapeutic strategy. *Nature* **434**, 917–921 (2005).

12. Easton, D. F. *et al.* Gene-Panel Sequencing and the Prediction of Breast-Cancer Risk. *N. Engl. J. Med.* **372**, 2243–2257 (2015).
13. Landrum, M. J. *et al.* ClinVar: public archive of interpretations of clinically relevant variants. *Nucleic Acids Res.* **44**, D862–8 (2016).
14. Cook-Deegan, R., Conley, J. M., Evans, J. P. & Vorhaus, D. The next controversy in genetic testing: clinical data as trade secrets? *Eur. J. Hum. Genet.* **21**, 585–588 (2013).
15. Yang, S., Cline, M., Zhang, C., Paten, B. & Lincoln, S. E. DATA SHARING AND REPRODUCIBLE CLINICAL GENETIC TESTING: SUCCESSES AND CHALLENGES. *Pac. Symp. Biocomput.* **22**, 166–176 (2017).
16. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
17. Millot, G. A. *et al.* A guide for functional analysis of BRCA1 variants of uncertain significance. *Hum. Mutat.* **33**, 1526–1537 (2012).
18. Ransburgh, D. J. R., Chiba, N., Ishioka, C., Toland, A. E. & Parvin, J. D. Identification of breast tumor mutations in BRCA1 that abolish its function in homologous DNA recombination. *Cancer Res.* **70**, 988–995 (2010).
19. Pierce, A. J., Hu, P., Han, M., Ellis, N. & Jasin, M. Ku DNA end-binding protein modulates homologous repair of double-strand breaks in mammalian cells. *Genes Dev.* **15**, 3237–3242 (2001).
20. Bouwman, P. *et al.* A high-throughput functional complementation assay for classification of BRCA1 missense variants. *Cancer Discov.* **3**, 1142–1155 (2013).
21. Starita, L. M. *et al.* Massively Parallel Functional Analysis of BRCA1 RING Domain Variants. *Genetics* **200**, 413–422 (2015).

22. Steffensen, A. Y. *et al.* Functional characterization of BRCA1 gene variants by mini-gene splicing assay. *Eur. J. Hum. Genet.* **22**, 1362–1368 (2014).
23. de la Hoya, M. *et al.* Combined genetic and splicing analysis of BRCA1 c.[594-2A>C; 641A>G] highlights the relevance of naturally occurring in-frame transcripts for developing disease gene variant classification algorithms. *Hum. Mol. Genet.* **25**, 2256–2268 (2016).
24. Richards, S. *et al.* Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* **17**, 405–424 (2015).
25. Ghosh, R., Oak, N. & Plon, S. E. Evaluation of in silico algorithms for use with ACMG/AMP clinical variant interpretation guidelines. *Genome Biol.* **18**, 225 (2017).
26. Gibson, T. J., Seiler, M. & Veitia, R. A. The transience of transient overexpression. *Nat. Methods* **10**, 715 (2013).
27. Findlay, G. M., Boyle, E. A., Hause, R. J., Klein, J. C. & Shendure, J. Saturation editing of genomic regions by multiplex homology-directed repair. *Nature* **513**, 120–123 (2014).
28. Blomen, V. A. *et al.* Gene essentiality and synthetic lethality in haploid human cells. *Science* **350**, 1092–1096 (2015).
29. Ran, F. A. *et al.* Genome engineering using the CRISPR-Cas9 system. *Nat. Protoc.* **8**, 2281–2308 (2013).
30. Moynahan, M. E., Chiu, J. W., Koller, B. H. & Jasin, M. Brca1 controls homology-directed DNA repair. *Mol. Cell* **4**, 511–518 (1999).
31. Drost, R. *et al.* BRCA1 RING function is essential for tumor suppression but dispensable for therapy resistance. *Cancer Cell* **20**, 797–809 (2011).
32. Shakya, R. *et al.* BRCA1 Tumor Suppression Depends on BRCT Phosphoprotein Binding,

- But Not Its E3 Ligase Activity. *Science* **334**, 525–528 (2011).
33. Easton, D. F. *et al.* A Systematic Genetic Assessment of 1,433 Sequence Variants of Unknown Clinical Significance in the BRCA1 and BRCA2 Breast Cancer–Predisposition Genes. *Am. J. Hum. Genet.* **81**, 873–883 (2007).
  34. Vega, A. *et al.* The R71G BRCA1 is a founder Spanish mutation and leads to aberrant splicing of the transcript. *Hum. Mutat.* **17**, 520–521 (2001).
  35. Beumer, K. J. *et al.* Efficient gene targeting in Drosophila by direct embryo injection with zinc-finger nucleases. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 19821–19826 (2008).
  36. Ma, Y. *et al.* Increasing the efficiency of CRISPR/Cas9-mediated precise genome editing in rats by inhibiting NHEJ and using Cas9 protein. *RNA Biol.* **13**, 605–612 (2016).
  37. Essletzbichler, P. *et al.* Megabase-scale deletion using CRISPR/Cas9 to generate a fully haploid human cell line. *Genome Res.* **24**, 2059–2065 (2014).
  38. whi.color.com. *FLOSSIES* Available at: <https://whi.color.com/gene/ENSG00000012048>. (Accessed: 9th October 2017)
  39. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
  40. Pollard, K. S., Hubisz, M. J., Rosenbloom, K. R. & Siepel, A. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res.* **20**, 110–121 (2010).
  41. Tavtigian, S. V., Byrnes, G. B., Goldgar, D. E. & Thomas, A. Classification of rare missense substitutions, using risk surfaces, with genetic- and molecular-epidemiology applications. *Hum. Mutat.* **29**, 1342–1354 (2008).
  42. Desmet, F.-O. *et al.* Human Splicing Finder: an online bioinformatics tool to predict splicing signals. *Nucleic Acids Res.* **37**, e67 (2009).

43. Goldgar, D. E. *et al.* Integrated Evaluation of DNA Sequence Variants of Unknown Clinical Significance: Application to BRCA1 and BRCA2. *Am. J. Hum. Genet.* **75**, 535–544 (2004).
44. Woods, N. T. *et al.* Functional assays provide a robust tool for the clinical annotation of genetic variants of uncertain significance. *Npj Genomic Medicine* **1**, 16001 (2016).
45. Spurdle, A. B. *et al.* ENIGMA - evidence-based network for the interpretation of germline mutant alleles: an international initiative to evaluate risk and clinical significance associated with sequence variation in BRCA1 and BRCA2 genes. *Hum. Mutat.* **33**, 2–7 (2012).
46. GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–585 (2013).
47. Starita, L. M. *et al.* Variant Interpretation: Functional Assays to the Rescue. *Am. J. Hum. Genet.* **101**, 315–325 (2017).
48. Gasperini, M., Starita, L. & Shendure, J. The power of multiplexed functional analysis of genetic variants. *Nat. Protoc.* **11**, 1782–1787 (2016).
49. Plon, S. E. *et al.* Sequence variant classification and reporting: recommendations for improving the interpretation of cancer susceptibility genetic test results. *Hum. Mutat.* **29**, 1282–1291 (2008).
50. Lovelock, P. K. *et al.* Identification of BRCA1 missense substitutions that confer partial functional activity: potential moderate risk variants? *Breast Cancer Res.* **9**, R82 (2007).
51. Spurdle, A. B. *et al.* BRCA1 R1699Q variant displaying ambiguous functional abrogation confers intermediate breast and ovarian cancer risk. *J. Med. Genet.* **49**, 525–532 (2012).
52. Domchek, S. M. *et al.* Biallelic Deleterious BRCA1 Mutations in a Woman with Early-Onset Ovarian Cancer. *Cancer Discov.* **3**, 399–405 (2013).
53. Tavtigian, S. V. *et al.* Comprehensive statistical study of 452 BRCA1 missense substitutions

- with classification of eight recurrent substitutions as neutral. *J. Med. Genet.* **43**, 295–305 (2006).
54. Kumar, P., Henikoff, S. & Ng, P. C. Predicting the effects of coding non-synonymous variants on protein function using the SIFT algorithm. *Nat. Protoc.* **4**, 1073–1081 (2009).
  55. Adzhubei, I. & Jordan, D. M. Predicting functional effect of human missense mutations using PolyPhen-2. *Current protocols in* (2013).
  56. Carette, J. E. *et al.* Ebola virus entry requires the cholesterol transporter Niemann-Pick C1. *Nature* **477**, 340–343 (2011).
  57. Walsh, T. *et al.* Detection of inherited mutations for breast and ovarian cancer using genomic capture and massively parallel sequencing. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 12629–12633 (2010).
  58. Hsu, P. D. *et al.* DNA targeting specificity of RNA-guided Cas9 nucleases. *Nat. Biotechnol.* **31**, 827–832 (2013).
  59. Doench, J. G. *et al.* Optimized sgRNA design to maximize activity and minimize off-target effects of CRISPR-Cas9. *Nat. Biotechnol.* **34**, 184 (2016).