

1 **A multiplexed homology-directed DNA repair assay reveals the impact of ~1,700**
2 **BRCA1 variants on protein function.**

3
4 Lea M. Starita^{1,2,7}, Muhtadi M. Islam^{3,7}, Tapahsama Banerjee³, Aleksandra I. Adamovich³,
5 Justin Gullingsrud¹, Stanley Fields^{1,4,5}, Jay Shendure^{*1,2,5}, and Jeffrey D. Parvin^{*3,6}.

6
7 ¹ Department of Genome Sciences, University of Washington, Seattle, Washington, USA.

8 ² Brotman Baty Institute for Precision Medicine, Seattle, Washington, USA.

9 ³ Department of Biomedical Informatics, The Ohio State University, Columbus, Ohio,
10 USA.

11 ⁴ Department of Medicine, University of Washington, Seattle, Washington, USA.

12 ⁵ Howard Hughes Medical Institute, Seattle, Washington, USA.

13 ⁶ The Ohio State University Comprehensive Cancer Center, The Ohio State University,
14 Columbus, Ohio, USA.

15 ⁷ These authors contributed equally to this work.

16
17 * Correspondence should be addressed to J.D.P (Jeffrey.Parvin@osumc.edu) or J.S.
18 (shendure@u.washington.edu).
19

20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Abstract

Loss-of-function mutations in *BRCA1* confer a predisposition to breast and ovarian cancer. Genetic testing for mutations in the *BRCA1* gene frequently reveals a missense variant for which the impact on the molecular function of the BRCA1 protein is unknown. Functional BRCA1 is required for homology directed repair (HDR) of double-strand DNA breaks, a key activity for maintaining genome integrity and tumor suppression. Here we describe a multiplex HDR reporter assay to simultaneously measure the effect of hundreds of variants of BRCA1 on its role in DNA repair. Using this assay, we measured the effects of ~1,700 amino acid substitutions in the first 302 residues of BRCA1. Benchmarking these results against variants with known effects, we demonstrate accurate discrimination of loss-of-function versus benign variants. We anticipate that this assay can be used to functionally characterize BRCA1 missense variants at scale, even before the variants are observed in results from genetic testing.

Introduction

Genetic testing for hereditary breast and ovarian cancer genes often reveals a missense variant in BRCA1 whose impact on the molecular function of the encoded protein, and therefore its contribution to cancer risk, is unknown. These variants are reported as “variants of uncertain significance” (VUS). VUS cause distress for both physicians and patients and can lead to unnecessary surgeries^{1,2}. As genetic testing becomes more common in the clinic, reports of missense VUS are rapidly accumulating in public databases³. Additional VUS arise from the increasingly widespread sequencing of tumor genomes and exomes to guide precision therapy. Therapeutic agents, such as PARP inhibitors, that specifically target BRCA-deficient tumors, are an effective therapy only for the subset of tumors with specific types of *BRCA1* mutation [for example refs. 4,5]. Therefore, knowledge of germline and somatic *BRCA1* mutation status is important to guide both cancer prevention and treatment strategies.

The most commonly reported class of VUS in *BRCA1* are single nucleotide variants (SNVs) that are predicted to result in missense amino acid substitutions. Currently, 1,794 missense VUS in BRCA1 are in the clinical genetics database Clinvar³. An additional 218 missense variants have conflicting interpretation reports, suggesting that clinical testing labs apply discordant classification criteria. Across BRCA1’s 1863 amino acids, there are 12,458 SNVs that are predicted to result in a missense substitution that might or might not affect protein function; these become VUS when identified as a germline or somatic variant in a patient. Most strategies for variant interpretation fail when the variants are sufficiently rare⁶⁻⁹. Computational variant-effect prediction algorithms scale without limit, but these are not accurate enough for routine clinical use¹⁰. Functional assays, on the other hand, are considered by the American College of Medical Genetics (ACMG) guidelines as strong evidence for or against the pathogenicity of missense variants¹⁰. However, performing a *post hoc* functional assay for each BRCA1 SNV as it is discovered is infeasible at their current rate of accumulation.

BRCA1 is required for maintenance of genome integrity via the homology-directed DNA repair (HDR) pathway. The effect of variants in BRCA1 on its HDR function can be

66 determined in tissue culture using a GFP-based reporter assay for intact DNA repair
67 function¹¹. For BRCA1 variants tested thus far, this assay has stratified the protein
68 function of known benign and pathogenic variants with high sensitivity and specificity^{12–}
69 ¹⁵. Here, we describe a multiplex version of this HDR assay that we developed toward the
70 goal of testing all possible protein variants in the N-terminus of BRCA1 (residues 2-302),
71 which includes the RING domain (residues 7-98). Proper folding of the RING domain is
72 required for the stability and function of the full-length protein^{13,16,17}. In addition, missense
73 mutations that cause increased cancer risk frequently map to either the RING or BRCT
74 domains.

75
76 Using the HDR reporter cell line with integrated *BRCA1* variant libraries, we test
77 approximately 600 BRCA1 missense variants per experiment. Damaging mutations are
78 identified by their relative depletion from the subset of cells that are GFP-positive. We
79 show that, as expected, the first 98 amino acids of BRCA1 which encode the RING
80 domain are more sensitive to substitutions than the subsequent 204 residues. The results
81 of the multiplexed assay correlate well with the results from singleton HDR reporter and
82 other functional assays. Furthermore, known pathogenic variants can be separated from
83 known benign variants by their recurrent depletion from the functional population in across
84 replicate experiments. Our results for 66 VUS or variants with conflicting reports of
85 pathogenicity from ClinVar show that seven are nonfunctional for HDR. We suggest
86 improvements to our current protocol that could increase its throughput and accuracy.
87 Finally, we anticipate that this assay can be used to functionally characterize BRCA1
88 missense variants at scale in order to provide the additional information necessary to
89 more definitively interpret VUS in the clinic.

90

91 **Results**

92

93 **A multiplexed assay to measure the effect of protein variants on homology-directed** 94 **DNA repair**

95 The multiplexed HDR reporter assay is based on an assay developed by the Jasin
96 laboratory¹¹. In this reporter assay, a site-specific double-stranded DNA break induced
97 by the I-SceI endonuclease results in conversion of a GFP-negative cell to GFP-positive
98 if the HDR pathway is intact (**Fig. 1a**). Loss of BRCA1 activity (e.g. by depletion with
99 siRNA) results in cells that cannot repair the GFP through HDR from the donor template.
100 We previously used this assay^{13–15,18} to test the capacity of 158 individual BRCA1
101 missense variants to rescue the loss of endogenous BRCA1. BRCA1 HDR function with
102 these individual assays demonstrates 100% specificity and 91.6% sensitivity for
103 predicting the cancer risk associated with variants with established interpretations (n =
104 43, **Supplementary Fig. 1**). The one known pathogenic variant misidentified as benign,
105 R71G, affects splicing, rather than protein function¹⁹, and therefore could not have been
106 classified correctly with an assay that expresses variants within a cDNA copy of *BRCA1*.
107 Given this functional assay's established specificity and sensitivity for predicting clinical
108 pathogenicity, we sought to convert it to a multiplexed format in order to generate high-
109 confidence predictions for hundreds of BRCA1 variants in a single experiment.

110

111 As a pilot version of a multiplexed HDR reporter assay (**Fig. 1b**), we prepared a pool of
112 plasmids containing the wild-type (WT) *BRCA1* and 15 *BRCA1* variants that had been
113 tested individually. We integrated these plasmids *en masse* into the HeLa-derived HDR
114 reporter cell line, which we had engineered to contain a single recombinase-based
115 landing pad (HeLa-DR-FRT, Methods). Though it is low efficiency, we used a
116 recombinase-based system to ensure that only a single variant would be integrated into
117 a common site in each cell, ensuring consistent expression levels and avoiding
118 complications associated with lentiviruses, *e.g.* variable integration sites and template
119 switching²⁰.

120

121 After selecting for cells containing an integrated *BRCA1* variant, we performed the three
122 steps of the HDR reporter assay. First, endogenous *BRCA1* was depleted with an siRNA
123 targeting the 3' untranslated region of the endogenous *BRCA1* mRNA. Second, I-SceI
124 was expressed to generate a double-strand break in the broken GFP reporter. Third, after
125 allowing sufficient time for double-stranded DNA break repair, cells were sorted into GFP-
126 positive and GFP-negative populations. The *BRCA1* variants in each sorted population
127 were PCR amplified and sequenced, and DNA reads for each variant were counted. We
128 calculated a score for each variant by taking the ratio of its frequency in the GFP-positive
129 cells to its frequency in the GFP-negative cells, and normalized to the equivalently
130 calculated ratio for the WT transgene. Scores from this pilot assay were largely consistent
131 with published results for these variants (**Fig. 1c, d**); the sole exception was the *BRCA1*
132 H41R variant, which had an intermediate level of activity in the multiplex assay. Known
133 pathogenic mutants were defective in this pilot experiment and the single known benign
134 missense variant had similar repair activity to the wild-type *BRCA1*. We therefore
135 proceeded to scale up the assay to analyze hundreds of variants per experiment, focusing
136 on *BRCA1* residues 2-302.

137

138 We created three individual pools of barcoded site-saturation mutagenesis libraries^{14,21}.
139 Pool1 contains variants in residues 2-96, pool2 in residues 97-192, and pool3 in residues
140 193-302 (**Supplementary Table 1**). We integrated the plasmids in each pool into the
141 HeLa-DR-FRT cell line. We then performed four replicates of the multiplexed HDR
142 reporter assay on each pool, using either an siRNA against endogenous *BRCA1* mRNA
143 or a control siRNA (**Supplementary Table 2**). The barcodes were amplified and
144 sequenced from the genome of the cells in the GFP-positive and GFP-negative
145 populations. The score for each variant was calculated as described above for the small-
146 scale assay (**Supplementary Table 4**). Variant scores between replicates were well
147 correlated for pool1, but less so for pools 2 and 3, for which most variants scored close
148 to $\ln(0)$, indicating no depletion (**Supplementary Fig. 2**)

149

150 **Functional classification of variants**

151 A characteristic of all sequencing count data, including from multiplexed assays for variant
152 effect, is that the variance of scores increases at low read counts due to error from
153 Poisson-distributed shot noise and stochastic dropout of variants²². This effect is
154 exacerbated in cases like the HDR reporter assay that, despite being the gold standard
155 assay for *BRCA1* missense variants, have a bottleneck. Specifically, only ~10% of the
156 cells at most become GFP-positive; this bottleneck limits the dynamic range of functional

157 scores (**Supplementary Table 2**). Because the score depends on read count, we cannot
158 directly compare the scores of any two variants. Therefore, we chose to construct a binary
159 classifier (“depleted” vs. “not depleted”) by modeling the relationship between the read
160 count and score (Methods). With this model, we tested the significance of each variant’s
161 depletion in the *BRCA1* siRNA experiments compared to variants in the control siRNA
162 experiments at the same read count in the same experiment (**Fig. 2a**, **Supplementary**
163 **Fig. 3**, Methods, and **Supplementary Table 4**). To remove as many variants as possible
164 that could be substantially affected by stochastic dropout, we applied a stringent read
165 count filter to the GFP negative population and removed variants below this threshold
166 from further analysis (**Supplementary Fig. 3**, Methods).

167
168 For each experimental replicate, the number of variants above the read-count threshold
169 in the GFP negative population varied, as did the number of variants classified as
170 depleted (*i.e.* nonfunctional variants, **Table 1**). The Venn diagrams shown in **Fig. 2b**
171 indicate the number of variants found to be depleted among the 4 replicate experiments
172 for each pool of variants. In each set of replicate experiments, some variants were scored
173 as depleted only once. Most of these singletons were depleted due to stochastic dropout,
174 but some were nonsense variants that passed the read-count filter only in one replicate
175 (**Supplementary Table 4**). The variants found to be depleted in 3 or 4 replicates were
176 enriched in nonsense codons and residues with known pathogenic mutations (Fisher’s
177 exact test, 2X enrichment, $p = 2.8 \times 10^{-6}$ and 20X enrichment, $p = 2.0 \times 10^{-29}$, respectively).
178 Pool1 had the highest percentage of nonfunctional variants, which was expected because
179 the structured RING domain is found almost entirely within the amino acids mutated in
180 pool1. Pool2 ranked second, with substitutions at the only two remaining positions of the
181 RING domain (97 and 98) repeatedly depleted from the functional GFP-positive
182 population. In contrast, pool3 had relatively few depleted variants, which were mostly not
183 shared between replicate experiments.

184
185 **Variant effects measured in this multiplexed HDR assay are strongly concordant**
186 **with singleton HDR assays and ClinVar interpretations**

187 In total, we measured the functional effect of 1,696 *BRCA1* missense or nonsense
188 variants in three or more replicates. The great majority of *BRCA1* missense variants were
189 functional for DNA repair, with only 61 variants (3.6%) depleted from the population in 3
190 or 4 replicates and therefore likely to be nonfunctional for DNA repair (**Fig. 3a**). Among
191 the variants successfully tested, results from singleton HDR reporter assays were
192 available for 15 variants^{12,14,15}. All 11 that had been scored as functional in HDR were
193 depleted in 0 replicates in the multiplexed assay, and all 4 that were nonfunctional in
194 singleton assays were all depleted in 3 or 4 replicates of the multiplexed assay (**Fig. 3b**).

195
196 The number of times that variants were depleted from the functional population in the
197 multiplexed HDR reporter assay was strongly correlated with the results of other
198 multiplexed functional assays. Specifically, we previously used *in vitro* ubiquitin ligase and
199 BARD1-binding yeast two-hybrid scores to predict HDR function for variants within the
200 first 102 amino acids of *BRCA1*¹⁴. The 25 variants that overlap in the two assays are
201 highly concordant (**Fig. 3c**). Variants that were never depleted in the multiplexed HDR
202 reporter assay correspond to those in the previous study with higher HDR predictions,

203 whereas variants that were repeatedly depleted correspond to those in the previous study
204 with lower HDR predictions. However, the phage and yeast-based functional assays
205 misidentified the known pathogenic mutant L22S as functional, whereas it was correctly
206 found to be defective in both single¹⁴ and multiplexed HDR experiments in human cells.

207
208 The number of replicates in which variants were depleted from the multiplexed HDR
209 reporter assay was also strongly consistent with functional scores from a multiplex assay
210 for BRCA1 function that relies on survival of a haploid cell line after saturation genome
211 editing (SGE) of exons in *BRCA1* (**Fig. 3d**; Findlay et al, enclosed). To an extent, the
212 exceptions are predictable. The SGE-based functional assay identifies variants that
213 reduce splicing because the variants are edited into their native context in the genome,
214 whereas they are false negatives in the HDR reporter assay in which the BRCA1 ORF
215 is not spliced (**Fig. 3d**, red points represent variants with reduced RNA levels).

216
217 Like the multiplexed HDR reporter assay, the SGE assay is highly accurate for identifying
218 variants that are damaging for protein function. However, there are variants that are found
219 to be damaging to protein function by SGE and not the multiplexed HDR reporter assay
220 (V14G, C44G, E85G, A92G, and A92T). In previous work, we have individually compared
221 the effects of amino acid substitutions on BRCA1 function in HDR versus repair by single-
222 strand annealing (SSA)¹⁵. Though HDR and SSA are both related mechanisms for
223 repairing DNA double-strand breaks dependent on BRCA1, there were different
224 tolerances for seven of the 35 variants tested when comparing the two assays in that
225 study. While all seven of these were functional in HDR, three were non-functional in SSA
226 and four had significantly depleted partial function in SSA¹⁵. When comparing the two
227 multiplexed assays for HDR and SGE, finding five differences from among 1,696
228 measurements of BRCA1 function in HDR is a low number and may reflect appropriate
229 biology of the two assays. This is especially true given that the biochemical mechanism
230 affecting cell growth in the SGE assay is unknown. Of the five differences when
231 comparing the high-throughput HDR and SGE assays, C44G may be a false negative of
232 the multiplexed HDR reporter assay given that all amino acid substitutions in the cysteine
233 and histidine residues that bind zinc ions in the BRCA1 RING domain that have been
234 tested to date in singleton functional assays were damaging¹²⁻¹⁴, although this specific
235 C44G variant has not before been tested for HDR function. Alternatively, C44G and the
236 other discordant variants may be slightly destabilized and rescued by the expression level
237 in the HDR reporter assay, or these residues may be necessary for another function of
238 BRCA1 critical for cell survival in the SGE assay.

239
240 Clinical interpretations for some of the BRCA1 variants that we functionally score here
241 are reported in the ClinVar database (accessed June, 2017)²⁰. Of the ~1,700 variants
242 scored here, seven were classified as benign or likely benign in ClinVar. All seven were
243 either never depleted, or depleted in only one replicate, in the multiplexed HDR assay. In
244 contrast, we tested five BRCA1 variants that are established pathogenic mutations in
245 ClinVar. Of these, four were depleted in 3 or 4 replicates in the multiplexed HDR assay.
246 The fifth pathogenic variant, R71G, had WT-like DNA repair activity in our data, a result
247 consistent with previous HDR reporter assays¹⁵. R71G causes a defect in RNA splicing¹⁸,
248 which as discussed above we do not expect to be detectable in our assay (**Fig. 3e**).

249

250 In summary, we observe strong concordance between variants depleted in 0-1 replicates
251 in our assay and benign status in ClinVar or WT-like function in other assays; as well as
252 between variants depleted in 3-4 replicates in our assay and pathogenic status in ClinVar
253 or loss-of-function in other assays. We therefore concluded that the number of replicates
254 in which a variant is found depleted is a reasonable proxy for its functionality, and term
255 these as “depletion scores” (range 0-4, Supplementary Table 5).

256

257 **Depletion scores identify damaging BRCA1 variants**

258 Depletion scores for all variants passing the read-count threshold in at least three
259 replicates are shown in the form of a sequence–function map in **Fig. 4a**. 100% of the
260 damaging amino acid substitutions (dark red) were observed in the first 98 amino acids
261 which comprises the RING domain, whereas residues 99-302 strongly tended to tolerate
262 amino acid substitutions. Nonsense mutants, indicated with the asterisk on the bottom
263 row were for the most part non-functional. Three of these nonsense mutants were scored
264 as functional (codons 289, 290, and 291), and we consider these to be false negatives.
265 The distribution of variants indicates that the degenerate positions in the mutagenic
266 oligonucleotides used to make the site-saturation mutagenesis libraries contained a much
267 higher fraction of guanines than the 1:1:1:1 ratio of the four nucleotides specified for the
268 synthesis (see Methods). Thus, codons containing guanine in the first or second position
269 (particularly glycine, arginine, alanine and valine) are the most highly represented. Since
270 the highest representation of substitutions was to glycine, we mapped the number of
271 times each substitution to glycine was found depleted in replicate experiments to the
272 solution structure of the BRCA1 and BARD1 RING domain dimer²³ (**Fig. 4b**). Amino acids
273 positions that were the most intolerant to substitution were either buried in the interior of
274 the 4-helix bundle, which acts as the BARD1 interface, or in the loops that coordinate zinc
275 ions.

276

277 We examined the depletion scores of 77 BRCA1 variants that are ambiguously classified
278 in ClinVar (VUS or conflicting reports of pathogenicity) to predict their functional impact.
279 Of 66 variants tested in the multiplexed HDR reporter assay that were classified as VUS
280 in ClinVar, three (V11G, C47G and I68R) had a depletion score of 3-4 and are thus likely
281 nonfunctional for HDR. Of 12 variants that have conflicting interpretations of pathogenicity
282 in ClinVar, four (M18T, T37R, C39F and H41L) were nonfunctional in the DNA repair
283 assay (**Fig. 3e**). These seven variants, all ambiguous in ClinVar and nonfunctional in our
284 assay, are either at positions that coordinate zinc ions (C39F, C47G, H41L), at positions
285 within the zinc-finger loops (T37R and I68R) or to have side chains in the interior of the
286 4-helix bundle (V11G and M18T). Most known pathogenic missense mutations map to
287 these same structural features. Of note, the four variants currently listed in ClinVar as
288 having conflicting interpretations of pathogenicity and nonfunctional here were interpreted
289 in the past as likely pathogenic or pathogenic³.

290

291 The remaining 70 variants that are ambiguously classified in ClinVar (VUS or conflicting
292 reports of pathogenicity) have a low depletion score in our assay. Of these, 60 lie outside
293 of the RING domain where we did not identify any damaging amino acid substitutions.
294 The 10 variants found within the structured RING include V11A, a conservative amino

295 acid substitution inside the 4-helix bundle, and I90T, whose sidechain points out of the
296 helix. T77M and T97A are substitutions to the amino acids that abut the helices and may
297 not affect BARD1 binding; however, other changes at T97 are not tolerated (**Fig. 4a**). I42L
298 and M48V are in the RING loops but are conservative changes and therefore may be
299 tolerated; E29G, E33A, G57R and P58A are also in the loops, but their side chains point
300 away from the interior of the structure which may be why they are tolerated. In studies of
301 ubiquitin ligase activity, substitutions at E29 are damaging to ubiquitin ligase function¹⁴
302 but not to HDR. This finding adds to the growing evidence that ligase activity is not
303 required for BRCA1's HDR function^{24,25}.

304

305 **Discussion**

306 We developed a multiplexed reporter assay to measure the effect of hundreds of amino
307 acid substitutions in BRCA1 on HDR, the molecular function most closely associated with
308 its role in tumor suppression. We reproducibly measured 1,696 of the possible 6,020
309 amino acid or nonsense substitutions (301 x 20). Although these 1,696 constitute only
310 28% of all possible substitutions, our assay is the most high-throughput to date that
311 specifically analyzes the DNA repair function of BRCA1. Of the 12 variants known to be
312 benign or pathogenic that are assayed here, we classified them with 80% sensitivity and
313 100% specificity³. The assay was only 80% sensitive because it misclassified R71G, a
314 mutation that is pathogenic consequent to its impact on splicing, a feature not tested in
315 the DNA repair assay.

316

317 In this multiplex HDR reporter assay we used a cell line (HeLa) that was not derived from
318 breast or ovarian tissue. The requirement for a high transfection efficiency restricted us
319 to only a few choices of human cell lines. The function of BRCA1 in HDR is considered
320 ubiquitous for all human cell types. Although the use of HeLa cells for HDR reporter
321 assays produced results that accurately predict hereditary breast and ovarian cancer risk,
322 this assay does not address the breast and ovarian specificity of loss of BRCA1 activity.

323

324 The non-uniform coverage of amino acid substitutions in our variant libraries was a
325 technical challenge. Two features of the protocol employed here to make variant
326 libraries²¹ contributed to this non-uniformity. First, inverse PCR reactions were performed
327 individually for each amino acid position using individually synthesized oligonucleotides.
328 Therefore, failed reactions and uneven mixing of the PCR products caused loss of all or
329 most substitutions at some positions (**Fig. 4a**). Second, a strong guanine bias at each
330 degenerate nucleotide position during oligonucleotide synthesis led to a bias in the
331 encoded amino acids (**Fig. 4a**). We anticipate that if the distribution of variants in the
332 library were more uniform, we would be able to query more variants per experiment, with
333 a theoretical maximum of ~4,000 variants per experiment under ideal conditions, at our
334 current number of sorted cells. Alternative approaches using array-derived
335 oligonucleotides may create more uniformly distributed variant libraries^{26,27}. It may also
336 be useful to limit the number of amino acid changes to those accessible by single
337 nucleotide changes, as these are the missense variants relevant to human disease.
338 However, this limitation would also reduce the information content of the multiplexed
339 experiments.

340

341 In addition to improving library uniformity, restricting the number of barcoded variants in
342 each multiplexed HDR reporter experiment should result in more variants that pass the
343 read-count filter. For example, a higher percentage of variants passed the read count
344 threshold in pool2 (50%, on average) than pool1 (29%) and pool3 (31%). Pool3 had twice
345 as many barcoded variants as pool2 (50,000 vs. 25,000) and covered more sequence
346 space (106 vs. 96 amino acids positions), which resulted in nearly a 20% reduction in the
347 number of reproducibly queried variants in pool3 compared to pool2. On the other hand,
348 the variant libraries for pool1 and pool2 contained similar numbers of barcoded variants.
349 However, more variants were damaging to BRCA1 HDR function in pool1, and thus only
350 ~7% of the cells in this pool converted to being GFP-positive following the double-strand
351 break, further limiting the dynamic range of pool 1 experiments.

352
353 In summary, we describe the development of the first multiplexed assay for measuring
354 the effects of amino acid substitutions on a protein's function in double-strand DNA break
355 repair in human cells. We analyzed nearly 1,700 amino acid substitutions of BRCA1
356 residues 2-302 and found that this approach yielded results comparable to low-throughput
357 HDR analysis of single variants and other functional assays. More importantly, our results
358 are concordant with known cancer-predisposing mutants of BRCA1, including perfect
359 positive predictive value for identifying known mutations damaging to protein function.
360 This assay can be repurposed to measure the effect of variants in other proteins in the
361 HDR pathway that are also hereditary breast and ovarian cancer tumor suppressors, such
362 as BRCA2²⁸⁻³⁰ and BARD1^{31,32}. As genetic testing for cancer risk becomes more common
363 and additional genes are added to testing panels, the number of rare missense variants
364 that inevitably become VUS will continue to increase. We anticipate that multiplexed
365 functional assays can be used to functionally characterize such variants at scale, even
366 before the variants are observed in results from genetic testing.

367

368 **Table 1**

Pool	rep	Total Variants	Variants above read threshold	Variants $q < 0.05$	% depleted
AA 1-96	1	1197	734	135	18%
	2	1167	793	212	27%
	3	1281	254	47	19%
	4	1292	228	61	27%
AA 97-192	1	1648	784	13	2%
	2	1643	663	63	10%
	3	1649	1091	34	3%
	4	1659	1291	31	2%
AA 193-302	1	1055	747	11	1%
	2	1056	718	3	0%
	3	1112	499	6	1%
	4	1117	670	4	1%

369

370

371 **Figure Legends**

372 **Figure 1 | Overview of the multiplexed-HDR reporter assay**

373 **a**, Schematic of the integrated HDR reporter. **b**, Schematic of workflow for multiplexed
 374 HDR-reporter assay. **c**, Results from the pilot 16-plex HDR assay testing WT and 15
 375 variants in a multiplexed format. WT-normalized, GFP-positive:GFP-negative ratios are
 376 plotted on the y-axis with variant identifications on the x-axis. **d**, The correlation between
 377 scores from the 16-plex experiment and scores from individual HDR-reporter assays.
 378 Spearman rho and p-value are reported. Bar and points are colored according to ClinVar
 379 variant interpretation.

380

381 **Figure 2 | Identifying BRCA1 variants depleted from the GFP-positive population**

382 **a**, The log of the WT-normalized, GFP-positive:GFP-negative ratios are on the y-axis and
 383 \log_{10} read counts are on the x-axis for a single replicate of the HDR-reporter assay for
 384 codons 2-96. Variants from the control (pink) or BRCA1 (black) siRNA conditions are
 385 indicated, variants significantly depleted from the GFP-positive population in the BRCA1
 386 siRNA condition, $q < 0.05$, colored blue. The dashed line represents the read-count
 387 threshold. **b**, Venn diagrams of the number of variants found depleted in multiple or single
 388 replicates.

389

390 **Figure 3 | Comparison of depletion scores to scores from other functional assays
 391 and ClinVar classifications**

392 **a**, Histogram of depletion score for all variant above the read count threshold in at least
 393 three replicates. **b**, Histograms of variant depletion scores that were functional or
 394 nonfunctional as measured by individual HDR assays. **c**, Box and strip plots comparing
 395 variant depletion scores (x-axis) to HDR predictions (y-axis) from ref.16; BRCA1 L22S is
 396 indicated. **d**, Box and strip plots comparing variant depletion scores (x-axis) to SGE

397 functional scores (y-axis; Findlay et al, unpublished), points marking variants with >80%
398 RNA depletion in the SGE assay are colored red; BRCA1 R71G is indicated. **e**,
399 Histograms of variant depletion scores as for each ClinVar classification.

400

401 **Figure 4 | The effect of amino acid substitutions the DNA repair function of BRCA1**
402 **2-302.**

403 **a**, A sequence-function map of the effect of amino acid substitutions in BRCA1 2-302.
404 The depletion score for each variant is the count of experimental replicates in which it
405 was found depleted from the functional population. Each position in BRCA1(2-302) is
406 arranged along the x-axis, the position of the RING domain is diagrammed above. The
407 amino acid substitutions, grouped by side-chain properties, are on the y-axis, nonsense
408 codons are *. The depletion scores range from never depleted, or likely functional
409 (white), to likely nonfunctional in dark red. Black ovals demarcate the wild-type residue
410 and gray, missing data. **b**, The depletion score for amino acid substitutions to glycine
411 are mapped to the solution structure of the BRCA1 1-102, BARD1 26-125 dimer (pdb
412 1JM7). Color scale as in (a) and spheres are shown for side chains at amino acid
413 positions with depletion score of 3 or 4.

414

415

416

417

418

419

420 **Acknowledgements**

421 We thank Jason Underwood and Katy Munson of the University of Washington PacBio
422 Sequencing Services for assistance with long-read sequencing, Ronald Hause and Alan
423 Rubin for helpful suggestions regarding statistical methods, Martin Kircher for help with
424 analysis of long read sequences, Ethan Ahler for assistance with figures and the OSU
425 Comprehensive Cancer Center Analytical Cytometry Shared Resource for sorting of the
426 GFP-positive cells.

427

428 **Funding Statement**

429 This work was supported by National Institutes of Health grants to S.F. (Biomedical
430 Technology Research Resource project #P41GM103533), J.S. (Director's Pioneer Award
431 #DP1HG007811-05), and a Bassar BRCA Innovation Award to J.D.P. S.F. and J.S. are
432 Investigators of the Howard Hughes Medical Institute. M.M.I. was supported by a
433 Pelotonia Cancer Training fellowship.

434

435 **Contributions**

436 L.M.S. and J.G. created, barcoded and assembled the plasmid libraries for the expression
437 of mutant BRCA1. L.M.S. amplified barcodes, sequenced and analyzed results of sorting
438 experiments. M.M.I. established the HeLa-derived cell lines with integrated libraries and
439 with T.B. and A.I.A., performed multiplexed HDR assay, sorted cells and extracted
440 genomic DNA. S.F., J.S., and J.D.P. provided guidance. L.M.S., S.F., J.S., and J.D.P.
441 wrote the manuscript.

442

443 **Methods**

444 All enzymes unless specifically mentioned were purchased from New England Biolabs.
445 Primer sequences can be found in Supplementary Table 3.

446
447 **Creating the HeLa DR-FRT cell line:** A description of the HeLa-DR cell line can be found
448 here¹³. To create a FLP-in version of HeLa DR, we stably integrated into the cells a flipase
449 recognition target (FRT) sequence using the pFRT/lacZeo plasmid (Thermo Fisher).
450 Zeocin resistant clones that had a single integration site detected by Southern blot were
451 tested for high activity integration sites using the mammalian β -galactosidase activity
452 assay (Gal-Screen, ThermoFisher). Clonal expansion of the selected colony established
453 the HeLa DR-FRT cell line.

454
455 **Site-saturation mutagenesis libraries and barcoding:** The HA-tagged BRCA1 N-
456 terminal HindIII-EcoRI fragment containing amino acids 1-302 was cloned into the pUC18
457 plasmid. Three site-saturation mutagenesis libraries of BRCA1 were constructed using a
458 previously reported inverse PCR-based method¹⁹. Pool1 has amino acid substitutions in
459 amino acids 2-96, pool2 97-192 and pool3 193-302. For each codon, 30 base, mutagenic
460 primers were ordered with machine-mixed NNK bases at the 5' end of the sense
461 oligonucleotide (N = ACTG, K = GT). The mutagenized HA-BRCA1 2-302 fragments were
462 ligated into the EcoRI and HindIII sites of the BRCA1 cDNA in pcDNA5/FRT-TO vector
463 that had been modified to remove a second EcoRI site in the gene for hygromycin
464 resistance and a second multiple cloning site added at the MluI/NruI sites for the barcode.
465 A 16 base, degenerate barcode encoded on oligos (pc5_barcode_longer_ W and _C)
466 that had been annealed, extended, and digested with NotI/SbfI was then ligated into the
467 second multiple cloning site. There were ~25,000 barcoded clones for each of the pool1
468 and pool2 libraries and ~50,000 for pool3 as assessed by colony forming units after
469 barcode ligation and transformation. Metrics for the variant library cloning steps can be
470 found in Supplementary Table 1.

471
472 **Assigning barcodes to variants using PacBio long reads:** To prepare the circular
473 SMRT-bell templates for the pool 1 and 2 variable regions and barcode, the intervening
474 sequence between the barcode and the BRCA1 N-terminal variable region was removed
475 by NotI/HindIII restriction digest, followed by end-repair and blunt-end ligation³³. The
476 ligations were transformed into *E. coli* to remove concatamers. The plasmids were then
477 cut with SbfI and EcoRI to release the barcode and BRCA1 N-terminal variable region.
478 For pool3 the entire SbfI and EcoRI fragment including the extra 2 Kb of sequences was
479 released. Custom SMRT bell adapters pb_SbfI and pb_EcoRI were sticky-end ligated to
480 the purified fragment. To make a working stock of 20 μ M SMRT bell adapters in 10 mM
481 Tris, 0.1 mM EDTA, 100 mM NaCl, they were heated to 85°C and snap cooled on ice.
482 The ligation reaction contained 0.3 pmol purified fragment, 0.4 μ M of each adaptor, 0.25
483 μ L of EcoRI-HF, 0.25 μ L of SbfI-HF, 1X ligase buffer, and 0.5 μ L of T4 ligase in a 25 μ L
484 reaction. The ligation was performed at room temperature for 30 minutes, then heat
485 inactivated at 65°C for 20 minutes. To cut destroy SMRT-bells with the remaining plasmid
486 backbone, 0.25 μ L each XhoI and NdeI were added and incubated 15 min 37°C. And
487 finally, to digest noncircular DNA, 0.5 μ L each of ExoIII (Enzymatics) and ExoVII were
488 added and incubated at 37°C for 15 minutes. The final SMRT bell fragments were purified

489 via AmpurePB (Pacific Biosciences) at 1.8X concentration, washed in 70% ethanol,
490 eluted in 15 μ L 10mM Tris pH 8 and quantified by BioAnalyzer (Agilent).

491
492 Each BRCA1 library was sequenced on four SMRT cells on a Pacific Biosciences RS II
493 sequencer. Barcodes and variable regions were identified as in ref. 34 as follows: Base
494 call files were converted from the bax format to the bam format using bax2bam (version
495 0.0.2) and then bam files for each library from separate lanes were concatenated.
496 Consensus sequences for each sequenced molecule in every library were determined
497 using the Circular Consensus Sequencing algorithm (version 2.0.0) with default
498 conditions (ccs and bax2bam can found in the PacBio Github repository,
499 <https://github.com/PacificBiosciences/unanimity/blob/master/doc/PBCCS.md>). Each
500 resulting consensus sequence was then aligned to a BRCA1 reference sequence using
501 Burrows-Wheeler Aligner³⁵ (<http://bio-bwa.sourceforge.net/>). Barcodes and insert
502 sequences were extracted from each alignment using custom scripts that parsed the
503 CIGAR and MD strings. For barcodes sequenced more than once, if barcode-variant
504 sequences differed, the barcode was assigned to the variant that represented more than
505 50% of the sequences. Barcodes lacking a majority variant sequence were assigned the
506 variant sequence with the highest average quality score as determined by the ccs
507 algorithm. The barcode-variant extraction and barcode unification scripts can be found at
508 <https://github.com/shendurelab/AssemblyByPacBio/>. Pool1 had 19,809 barcodes
509 assigned to variants with 0 or 1 amino substitution encoding 1602 unique protein variants.
510 Pool2 had 17,635 barcodes assigned to variants with 0 or 1 amino substitution encoding
511 1695 unique protein variants. Pool3 had 11857 barcodes assigned to variants with 0 or
512 1 amino substitution encoding 1987 unique protein variants. Additional metrics regarding
513 the sequence processing for the barcode-variant assignments can be found in
514 Supplementary Table 1. For all three BRCA1 libraries, a barcode-variant map file was
515 created that contains each barcode and its nucleotide sequence.

516
Integration of libraries into cells: For each of the three BRCA1 plasmid libraries (pools
517 1-3), 70-80 million HeLa DR-FRT cells in ten 10 cm tissue culture plates were transfected
518 with 200 μ g of pOG44 to express a modified flipase enzyme and 100 μ g pcDNA5/FRT
519 BRCA1 variant library. Plasmids are diluted in 10 mL Opti-MEM and incubated for 5
520 minutes. 300 μ l of Lipofectamine 2000 (ThermoFisher) is diluted in 10 mL Opti-MEM for
521 5 minutes. The lipofectamine and plasmid dilutions are then combined and incubated for
522 20 minutes. The mixture was then applied directly to cells. After 24 hours, cells were
523 trypsinized and transferred to ten 15 cm tissue culture dishes. Since the pOG44 flipase
524 has reduced activity at 37°C, four to eight hours post transfer, the cells were moved from
525 a 37°C humidified incubator to a 30°C humidified incubator for 24 hours. The cells were
526 then returned to 37°C for an additional 24 hours. Approximately 72 hours after the initial
527 transfection, the cells were trypsinized and transferred to twenty 15 cm plates containing
528 selection media (50% fresh DMEM supplemented with 10% fetal bovine serum, 50% filter
529 sterilized conditioned media and hygromycin B at 550 μ g/mL). Hygromycin resistant cells
530 are selected at 37°C for 24 hours. The cells are washed with sterile phosphate buffered
531 saline (PBS) and the selection media is replaced after the first 24 hours and again every
532 48 hours until cell colonies are visible without a microscope (about 14 days). Colonies
533 were then counted, trypsinized, resuspended in 20 mL culture media, and mixed
534

535 thoroughly, colony counts can be found in Supplementary Table 1. 15 ml of the
536 resuspended cells were frozen in 1 mL aliquots. 5 mL of resuspended colony mixture
537 was plated onto three 15 cm plates (3 mL, 1.5 mL, 0.5 mL) and incubated for 24 hrs. The
538 plate that was closest to a confluent monolayer of cells was passaged. The cells were
539 passaged for an additional two weeks before performing HDR reporter experiments to
540 assure loss of the unintegrated BRCA1 expression plasmid.

541
542 **HDR reporter assays, sorting, gDNA prep:** HDR reporter assays and FACS sorts were
543 performed for each of the three BRCA1 variant libraries in quadruplicate. A confluent 10
544 cm plate of HeLa BRCA1 variant cell line was trypsinized and resuspended in 10 mL of
545 culture media. 65 μ l of the suspension was plated in each of 48 wells across two 24-well
546 tissue culture plates. 24 hours later, each well was transfected with 30 pmol of siRNA and
547 1.5 μ l Oligofectamine (ThermoFisher). Oligofectamine was diluted with 6 μ l Opti-MEM
548 and siRNA was diluted with 25 μ l Opti-MEM for 5 minutes. The dilutions were then
549 combined and incubated for an additional 30 minutes. The transfection mixture was then
550 applied directly to the cells. 24 hours later, the cells were trypsinized and transferred to
551 four 6-well tissue culture plates, then incubated for another 24 hours. Each well was then
552 transfected with 50 pmol of siRNA, 3 μ g pCBASceI (for I-SceI expression), and 3 μ l
553 Lipofectamine 2000. Plasmid and siRNA were diluted in 125 μ l Opti-MEM and
554 Lipofectamine was diluted in 125 μ l Opti-MEM, then incubated for 5 minutes. The dilutions
555 were combined and incubated for an additional 20 minutes. The transfection mixture was
556 then applied directly to cells. Four to six hours later, the culture media was replaced with
557 fresh media. BRCA1 3'UTR siRNA, BRCA1 coding sequence siRNA and control siRNA
558 were used in both rounds of transfection. After 24 hours, one well of cells treated with
559 each condition was analyzed for GFP expression using flow cytometry to confirm
560 transfection efficiency. If cells treated with control siRNA and 3'UTR siRNA fell within 7-
561 9% and 4-7% GFP+ cells respectively, and cells treated with BRCA1 coding sequence
562 siRNA were 1-2% GFP+, the experiment would proceed. 72 hours post-transfection, the
563 cells were pooled according to treatment and sorted using fluorescent activated cell
564 sorting (FACS) using an Aria IIu instrument. Cells were resuspended and pooled in filter-
565 sterilized sorting buffer containing 1X Phosphate Buffered Saline ($\text{Ca}^{2+}/\text{Mg}^{2+}$ free), 5 mM
566 EDTA, 25 mM HEPES pH 7.0, and 1% heat-inactivated fetal bovine serum dialyzed
567 against $\text{Ca}^{2+}/\text{Mg}^{2+}$ PBS. A minimum of 500,000 GFP+ cells and a maximum of 2 million
568 GFP- cells were collected per pool (Supplementary Table 2). Genomic DNA (gDNA) was
569 extracted from the GFP+ and GFP- cells with a DNeasy Blood & Tissue Kit according to
570 manufacturer instructions (Qiagen). DNAs were eluted in 200 μ l Buffer EB.

571
572 **Barcode amplification and sequencing:** To amplify the barcode from gDNA from the
573 GFP-positive and negative populations were spread over 8 - 16 reactions containing 250
574 μ g of gDNA each (220K - 440K genome equivalents by weight given HeLa triploid
575 genome ~9 pg, see Supplementary Table 2). Reactions also contained Kapa2G Robust
576 Polymerase (Kapa Biosystems), a primer that annealed to the SV40 promoter 5' of the
577 integrated plasmid and a primer adjacent to the barcode (SV40_F and
578 newpc5bc_nexteraR). PCR was performed using the following conditions [95°, 5 min;
579 {95°, 40s, 65°, 30s, 72°, 3 min} x ~28 cycles; 72°, 10 min]. The reactions produce a ~3,700
580 base amplicon specifically from integrated plasmids. The reactions for each sample were

581 combined and the amplicons were purified by 0.5X Ampure and eluted with 10 mM Tris
582 at 10% of the original reaction volume. 10% of the eluted PCR volume was re-amplified
583 with primers containing sample indexes and Illumina cluster generating sequences
584 (pc5bc_p5_F, nextIndex). Reactions also contained Kapa2G Robust Polymerase and
585 0.5X SYBR green II (ThermoFisher). PCR reactions are monitored on a Mini-opticon
586 qPCR machine (Bio-Rad) and removed during exponential amplification using the
587 following conditions [95°, 3 min; {95°, 20s, 65°, 30s, 72°, 20s} x 5-10 cycles]. The PCR
588 reaction produced a 350-base amplicon that was purified using a double-Ampure
589 protocol. First the large DNA fragments were removed by precipitation by addition of 0.6X
590 volume Ampure beads to the reactions. Then the 350-base amplicons were purified from
591 the supernatant using 0.9X volume Ampure beads. The samples were multiplexed and
592 the barcodes and sample indexes were sequenced (single read,
593 pcDNA5_barcodeSeq_F) on a Nextseq 500 High Output 75 base kit, reads per
594 experiment can be found in Supplementary Table 2. For the pilot 16-plex experiment
595 without barcodes, the region of BRCA1 containing the variants was amplified and directly
596 sequenced.

597
598 **Variant scoring, classifications and depletion score:** FASTQ files containing either
599 barcodes for each sample and the barcode-map for each library were used as input for
600 the software package Enrich2²². Enrich2 was used to count the barcodes, associate each
601 barcode with a nucleotide variant, and then translate and count both the unique-
602 nucleotide and unique-amino acid variants. Barcodes assigned to variants containing
603 insertions or deletions were removed from analysis. The counts for each protein variant
604 were converted to frequencies by dividing by the total number of variant counts for each
605 sample. The ratio of the frequency of each variant in the GFP positive population over its
606 frequency in the GFP negative population was calculated. That ratio for each variant was
607 then normalized to the GFP+/GFP- ratio for the wild-type sequence at $\ln(\text{GFP+}/\text{GFP-}) =$
608 0. Variants with multiple amino acid substitutions were removed from further analysis.

609
610 To construct the binary, functional / nonfunctional classifier for variants, we first
611 determined that the standard deviation of the score decays according to read count in a
612 way that can be modeled by a \log_{10} - \log_{10} curve. We then modeled the decay of the
613 standard deviation of scores from the control siRNA experiment. Next, we applied that
614 model to the standard deviation of scores from the BRCA1 siRNA experiment and
615 calculated a p-value and an FDR-adjusted q-value for each variant to determine if it was
616 similar to the control siRNA experiment ($q > 0.05$) or significantly different from control
617 experiment ($q < 0.05$). Finally, we determined where the false positives in the control
618 experiments occurred along the continuum of read counts to assign a read-count
619 threshold for the BRCA1 siRNA condition. A read-count threshold is usually part of the
620 heuristic applied to remove noise due to stochastic dropout from deep mutational
621 scanning data¹⁶. We removed variants with below that read count threshold from further
622 analyses. These analyses are performed using R studio, an R markdown file containing
623 all data manipulations is supplied as Supplementary File 1.

624
625 The depletion score represents the number of replicates in which a variant was found to
626 be depleted from the BRCA1 siRNA GFP-positive population. Only variants that had been

627 present above the read-count threshold in at least three replicates have depletion scores.
628 Depletion scores for those 1699 variants can be found in Supplementary Table 5.

629

630 **Website accessions:**

631 ClinVar, accessed June 2017, minimum 1-star interpretation.

632

633

634 **References**

- 635 1. Murray, M. L., Cerrato, F., Bennett, R. L. & Jarvik, G. P. Follow-up of carriers of
636 BRCA1 and BRCA2 variants of unknown significance: Variant reclassification and
637 surgical decisions. *Genet. Med.* **13**, 998–1005 (2011).
- 638 2. Welsh, J. L. *et al.* Clinical Decision-Making in Patients with Variant of Uncertain
639 Significance in BRCA1 or BRCA2 Genes. *Ann. Surg. Oncol.* **24**, 3067–3072 (2017).
- 640 3. Landrum, M. J. *et al.* ClinVar: public archive of relationships among sequence
641 variation and human phenotype. *Nucleic Acids Res.* **42**, D980–D985 (2014).
- 642 4. Drost, R. *et al.* BRCA1185delAG tumors may acquire therapy resistance through
643 expression of RING-less BRCA1. *J. Clin. Invest.* **126**, 2903–2918 (2016).
- 644 5. Hollis, R. L., Churchman, M. & Gourley, C. Distinct implications of different BRCA
645 mutations: Efficacy of cytotoxic chemotherapy, PARP inhibition and clinical
646 outcome in ovarian cancer. *Onco. Targets. Ther.* **10**, 2539–2551 (2017).
- 647 6. Chenevix-Trench, G. *et al.* Genetic and histopathologic evaluation of BRCA1 and
648 BRCA2 DNA sequence variants of unknown clinical significance. *Cancer Res.* **66**,
649 2019–27 (2006).
- 650 7. Sweet, K., Senter, L., Pilarski, R., Wei, L. & Toland, A. E. Characterization of
651 BRCA1 ring finger variants of uncertain significance. *Breast Cancer Res. Treat.*
652 **119**, 737–43 (2010).
- 653 8. Osorio, A. *et al.* Loss of heterozygosity analysis at the BRCA loci in tumor samples
654 from patients with familial breast cancer. *Int. J. cancer* **99**, 305–9 (2002).
- 655 9. Easton, D. F. *et al.* A Systematic Genetic Assessment of 1,433 Sequence Variants
656 of Unknown Clinical Significance in the BRCA1 and BRCA2 Breast Cancer–
657 Predisposition Genes. *Am. J. Hum. Genet.* **81**, 873–883 (2007).
- 658 10. Richards, S. *et al.* Standards and guidelines for the interpretation of sequence
659 variants: a joint consensus recommendation of the American College of Medical
660 Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.*
661 **17**, 405–423 (2015).
- 662 11. Pierce, A. J., Johnson, R. D., Thompson, L. H. & Jasin, M. XRCC3 promotes
663 homology-directed repair of DNA damage in mammalian cells. *Genes Dev.* **13**,
664 2633–8 (1999).
- 665 12. Lu, C. *et al.* Patterns and functional implications of rare germline variants across 12
666 cancer types. *Nat. Commun.* **6**, (2015).
- 667 13. Ransburgh, D. J. R., Chiba, N., Ishioka, C., Toland, A. E. & Parvin, J. D.
668 Identification of breast tumor mutations in BRCA1 that abolish its function in
669 homologous DNA recombination. *Cancer Res.* **70**, 988–995 (2010).
- 670 14. Starita, L. M. *et al.* Massively parallel functional analysis of BRCA1 RING domain
671 variants. *Genetics* **200**, (2015).
- 672 15. Towler, W. I. *et al.* Analysis of BRCA1 variants in double-strand break repair by

- 673 homologous recombination and single-strand annealing. *Hum. Mutat.* **34**, 439–45
674 (2013).
- 675 16. Drost, R. *et al.* BRCA1 RING function is essential for tumor suppression but
676 dispensable for therapy resistance. *Cancer Cell* **20**, 797–809 (2011).
- 677 17. Wu, W. *et al.* HERC2 is an E3 ligase that targets BRCA1 for degradation. *Cancer*
678 *Res.* **70**, 6384–6392 (2010).
- 679 18. Lu, C. *et al.* Patterns and functional implications of rare germline variants across 12
680 cancer types. *Nat. Commun.* **6**, 10086 (2015).
- 681 19. Vega, A. *et al.* The R71G BRCA1 is a founder Spanish mutation and leads to
682 aberrant splicing of the transcript. *Hum. Mutat.* **17**, 520–521 (2001).
- 683 20. Sack, L. M., Davoli, T., Xu, Q., Li, M. Z. & Elledge, S. J. Sources of Error in
684 Mammalian Genetic Screens. *G3 (Bethesda)*. **6**, 2781–90 (2016).
- 685 21. Jain, P. C. & Varadarajan, R. A rapid, efficient, and economical inverse polymerase
686 chain reaction-based method for generating a site saturation mutant library. *Anal.*
687 *Biochem.* **449**, 90–8 (2014).
- 688 22. Rubin, A. F. *et al.* A statistical framework for analyzing deep mutational scanning
689 data. doi:10.1186/s13059-017-1272-5
- 690 23. Brzovic, P. S., Rajagopal, P., Hoyt, D. W., King, M. C. & Klevit, R. E. Structure of a
691 BRCA1-BARD1 heterodimeric RING-RING complex. *Nat. Struct. Biol.* **8**, 833–837
692 (2001).
- 693 24. Shakya, R. *et al.* BRCA1 Tumor Suppression Depends on BRCT Phosphoprotein
694 Binding, But Not Its E3 Ligase Activity. *Science (80-.)*. **334**, 525–528 (2011).
- 695 25. Reid, L. J. *et al.* E3 ligase activity of BRCA1 is not essential for mammalian cell
696 viability or homology-directed repair of double-strand DNA breaks. *Proc. Natl. Acad.*
697 *Sci.* **105**, 20876–20881 (2008).
- 698 26. Majithia, A. R. *et al.* Prospective functional classification of all possible missense
699 variants in PPARG. *Nat. Genet.* **48**, 1570–1575 (2016).
- 700 27. Kitzman, J. O., Starita, L. M., Lo, R. S., Fields, S. & Shendure, J. Massively parallel
701 single-amino-acid mutagenesis. *Nat. Methods* **12**, (2015).
- 702 28. Wooster, R. *et al.* Identification of the breast cancer susceptibility gene BRCA2.
703 *Nature* **378**, 789–792 (1995).
- 704 29. Moynahan, M. E., Pierce, A. J. & Jasin, M. BRCA2 is required for homology-
705 directed repair of chromosomal breaks. *Mol. Cell* **7**, 263–272 (2001).
- 706 30. Guidugli, L. *et al.* A classification model for BRCA2 DNA binding domain missense
707 variants based on homology-directed repair activity. *Cancer Res.* **73**, 265–75
708 (2013).
- 709 31. Couch, F. J. *et al.* Associations Between Cancer Predisposition Testing Panel
710 Genes and Breast Cancer. *JAMA Oncol.* **3**, 1190 (2017).
- 711 32. Lee, C. *et al.* Functional Analysis of BARD1 Missense Variants in Homology-
712 Directed Repair of DNA Double Strand Breaks. *Hum. Mutat.* **36**, 1205–1214 (2015).
- 713 33. Travers, K. J., Chin, C.-S., Rank, D. R., Eid, J. S. & Turner, S. W. A flexible and
714 efficient template format for circular consensus sequencing and SNP detection.
715 *Nucleic Acids Res.* **38**, e159–e159 (2010).
- 716 34. Matreyek, K. A. *et al.* Multiplex Assessment of Protein Variant Abundance by
717 Massively Parallel Sequencing. *bioRxiv* 211011 (2018). doi:10.1101/211011
- 718 35. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler

719 transform. *Bioinformatics* **26**, 589–595 (2010).
720

Fig. 1

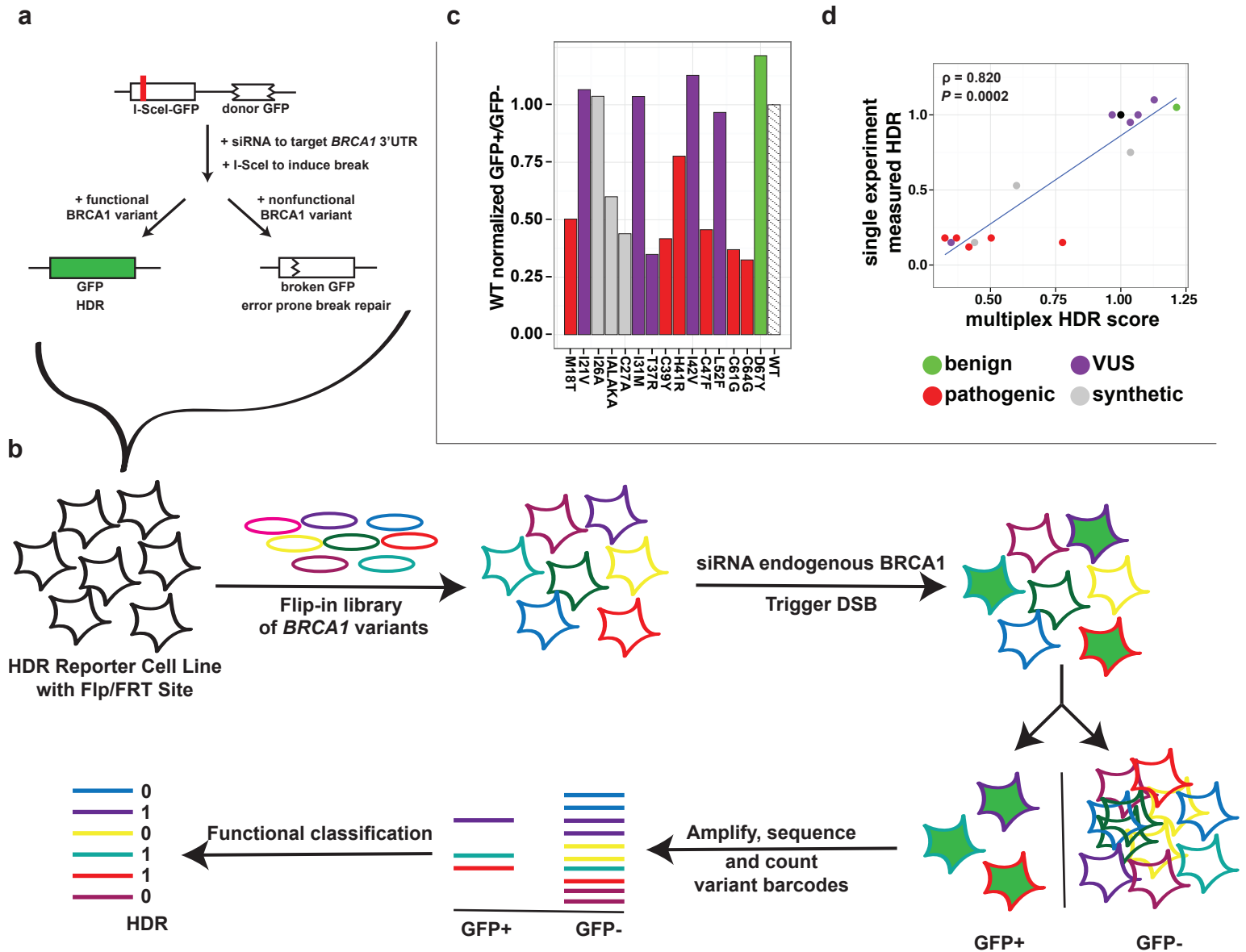


Figure 1 | Overview of the multiplexed-HDR reporter assay.

a, Schematic of the integrated HDR reporter. **b**, Schematic of workflow for multiplexed HDR-reporter assay. **c**, Results from the pilot 16-plex HDR assay testing WT and 15 variants in a multiplexed format. WT-normalized, GFP-positive:GFP-negative ratios are plotted on the y-axis with variant identifications on the x-axis. **d**, The correlation between scores from the 16-plex experiment and scores from individual HDR-reporter assays. Spearman rho and p-value are reported. Bar and points are colored according to ClinVar variant interpretation.

Fig. 2

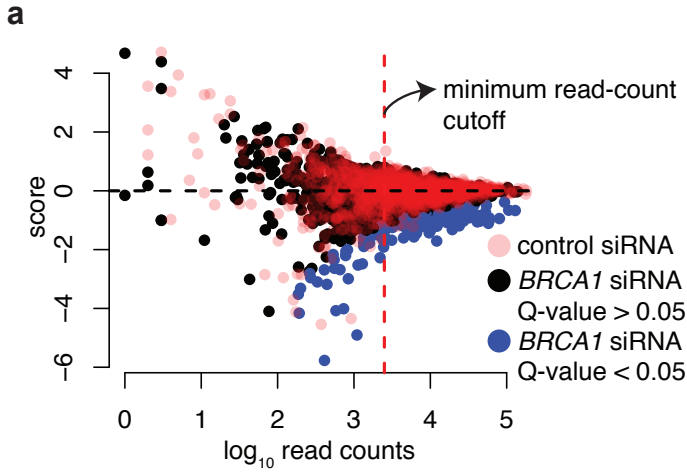


Figure 2 | BRCA1 variants depleted from the GFP-positive population.

a, The log of the WT-normalized, GFP-positive:GFP-negative ratios are on the y-axis and log₁₀ read counts are on the x-axis for a single replicate of the HDR-reporter assay for codons 2-96. Variants from the control (pink) or BRCA1 (black) siRNA conditions are indicated, variants significantly depleted from the GFP-positive population in the BRCA1 siRNA condition, $q < 0.05$, colored blue. The dashed line represents the read-count threshold. **b**, Venn diagrams of the number of variants found depleted in multiple or single replicates.

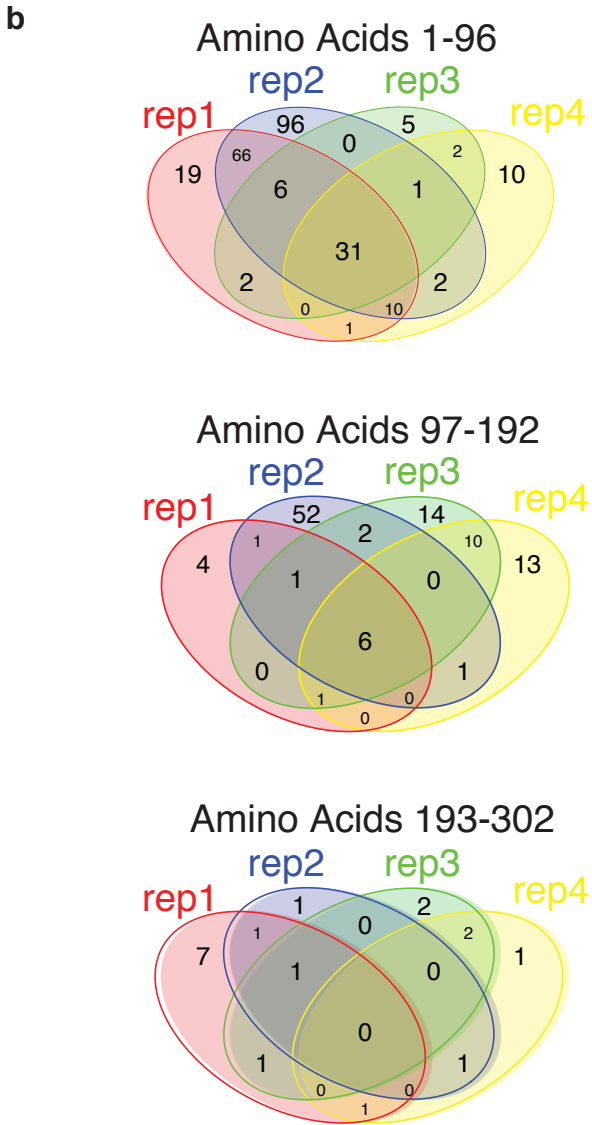


Fig. 3

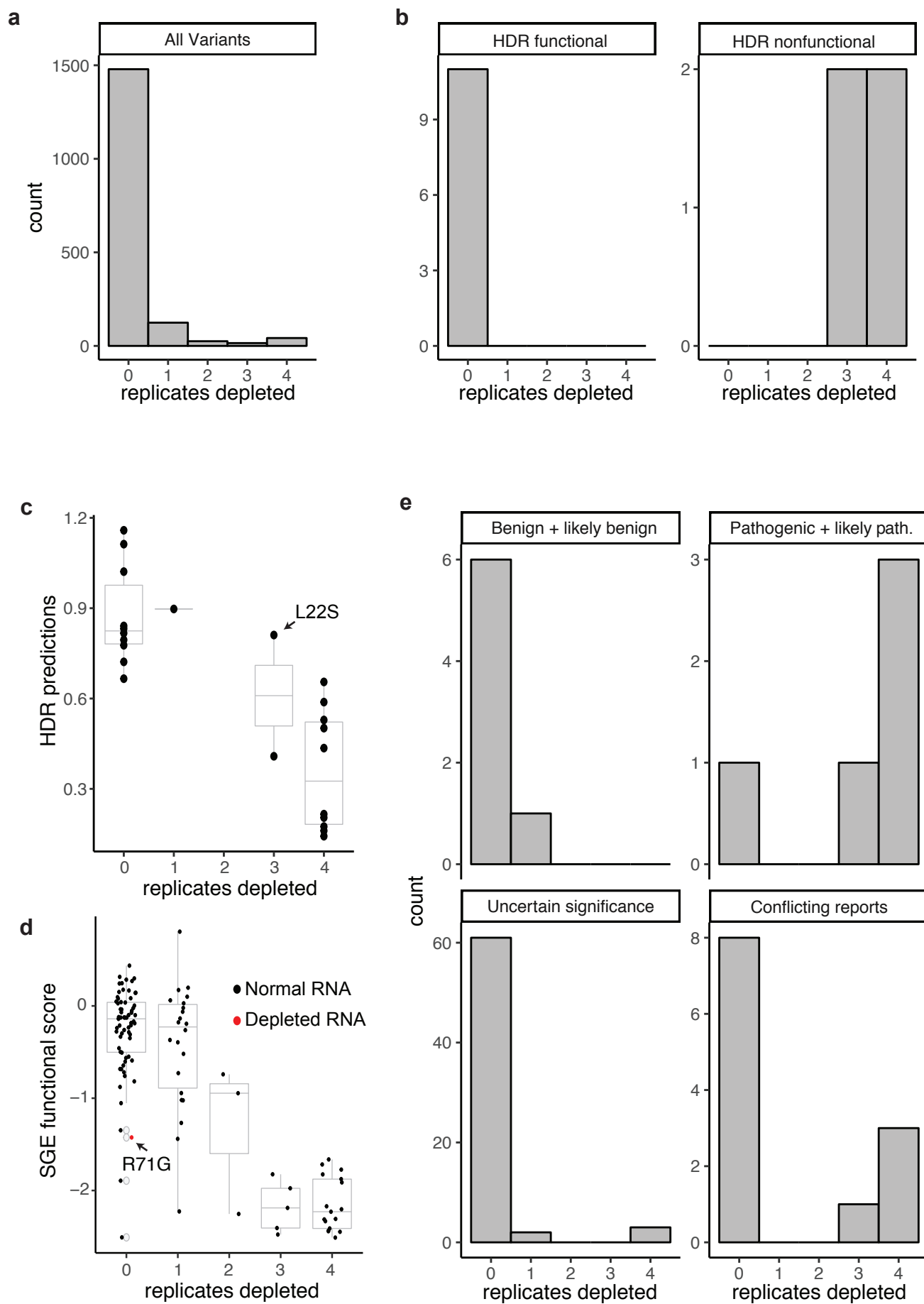


Figure 3 | Comparison of depletion scores to scores from other functional assays and ClinVar classifications.

a, Histogram of depletion score for all variant above the read count threshold in at least three replicates. **b**, Histograms of variant depletion scores that were functional or nonfunctional as measured by individual HDR assays. **c**, Box and strip plots comparing variant depletion scores (x-axis) to HDR predictions (y-axis) from ref.16; BRCA1 L22S is indicated. **d**, Box and strip plots comparing variant depletion scores (x-axis) to SGE functional scores (y-axis; Findlay et al, unpublished), points marking variants with >80% RNA depletion in the SGE assay are colored red; BRCA1 R71G is indicated. **e**, Histograms of variant depletion scores as for each ClinVar classification.

Fig. 4

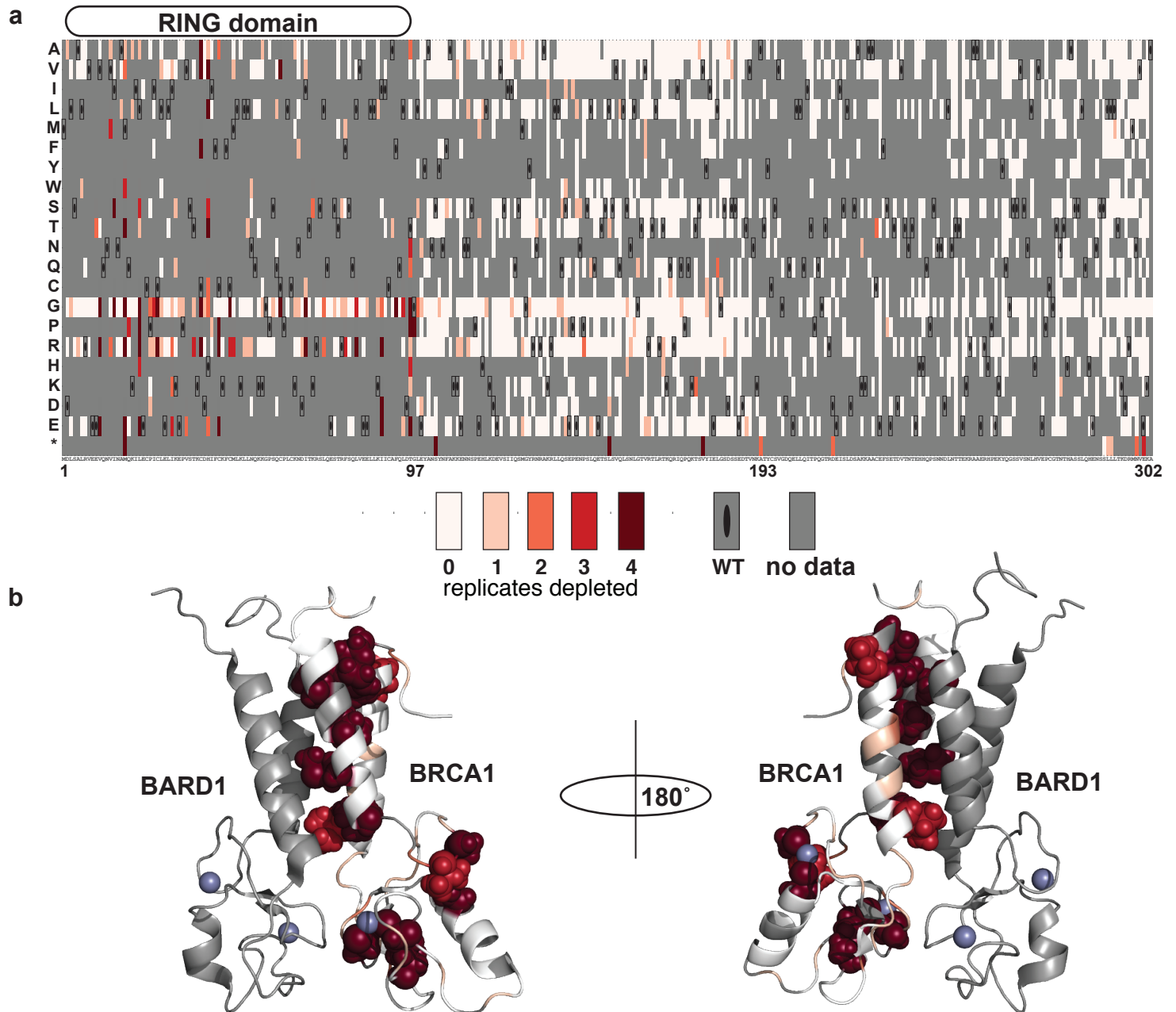


Figure 4 | The effect of amino acid substitutions the DNA repair function of BRCA1 2-302.

a, A sequence-function map of the effect of amino acid substitutions in BRCA1 2-302. The depletion score for each variant is the count of experimental replicates in which it was found depleted from the functional population. Each position in BRCA1(2-302) is arranged along the x-axis, the position of the RING domain is diagrammed above. The amino acid substitutions, grouped by side-chain properties, are on the y-axis. The depletion scores range from never depleted, or likely functional (white), to likely nonfunctional in dark red. Black ovals demarcate the wild-type residue and gray, missing data. **b**, The depletion score for amino acid substitutions to glycine are mapped to the solution structure of the BRCA1 1-102, BARD1 26-125 dimer (pdb 1JM7). Color scale as in (a) and spheres are shown for side chains at amino acid positions with depletion score of 3 or 4.

Supplementary Materials Table of Contents:

In this file:

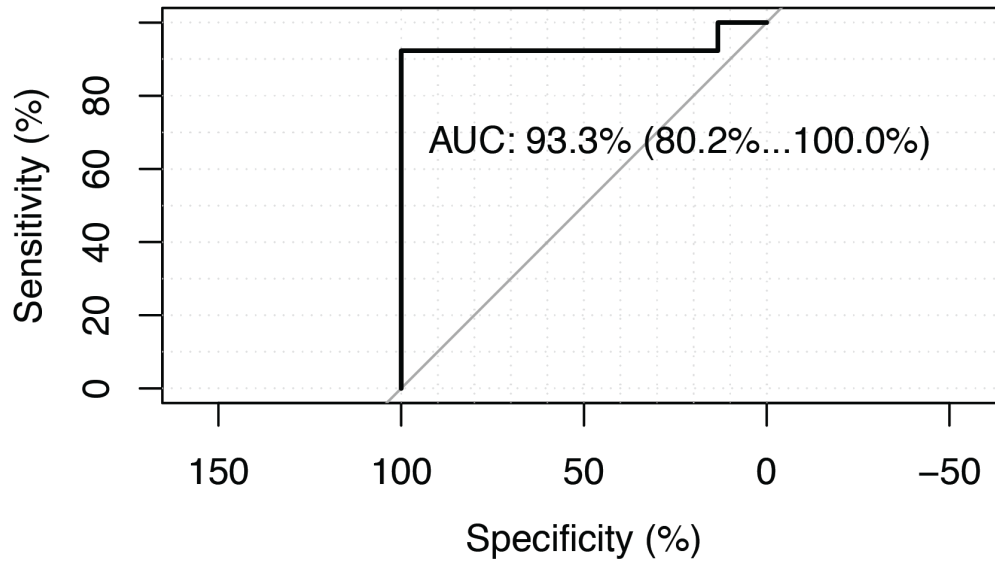
Supplementary Figures 1-3

Supplementary Tables 1-3

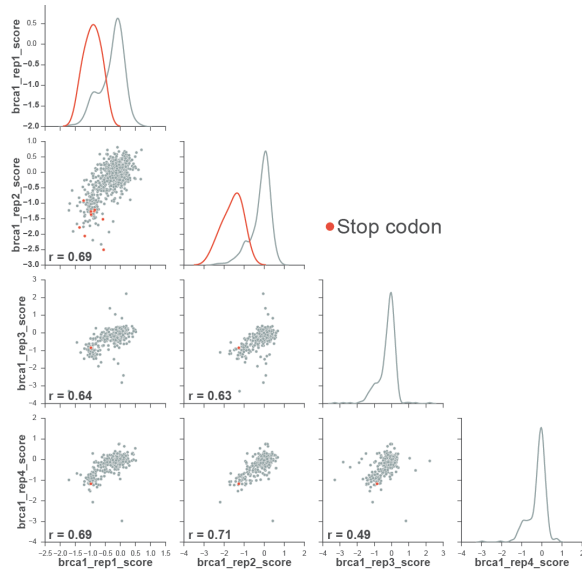
Separate files:

Supplementary Tables 4-5 are tab separated data tables. Descriptions of data and columns names are provided in this document. Supplementary file 1 is uploaded as both rmd (R markdown) and pdf.

Supplementary Figure 1 | Receiver-Operator Characteristic curve for singleton HDR reporter assay results. N = 43, 30 known benign and 13 known pathogenic BRCA1 variants stratified by results from single HDR reporter assays from refs 1–4. Plotted with the pROC r package⁵. 90% confidence interval reported.

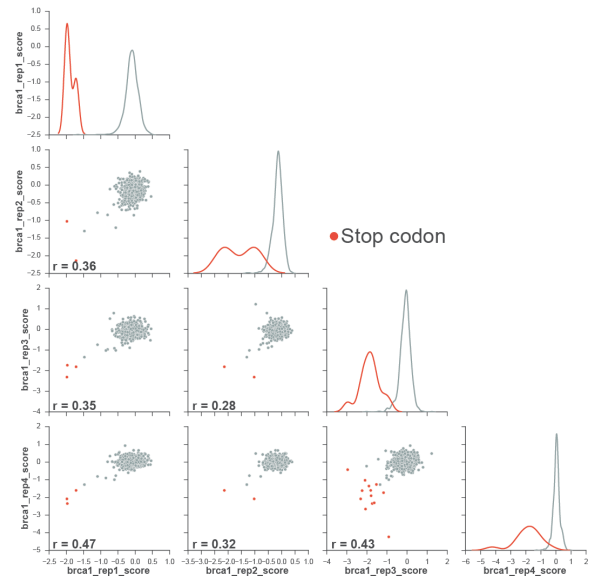


Supplementary Figure 2

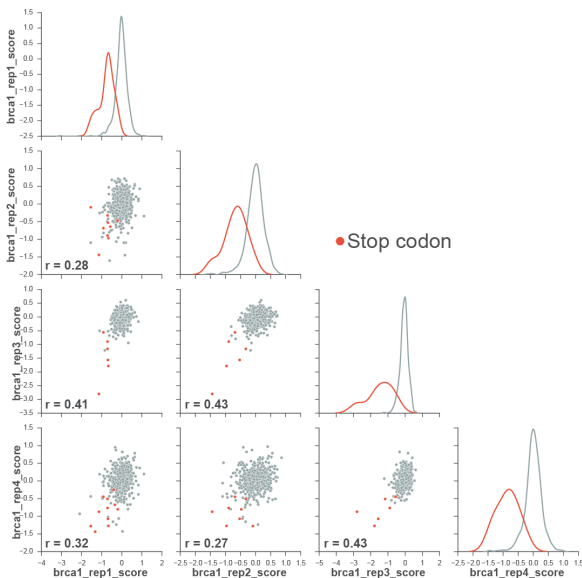


Pool1 replicates correlation of WT-normalized-GFP-positive/GFP-negative ratios

Pool2 replicates correlation of WT-normalized-GFP-positive/GFP-negative ratios

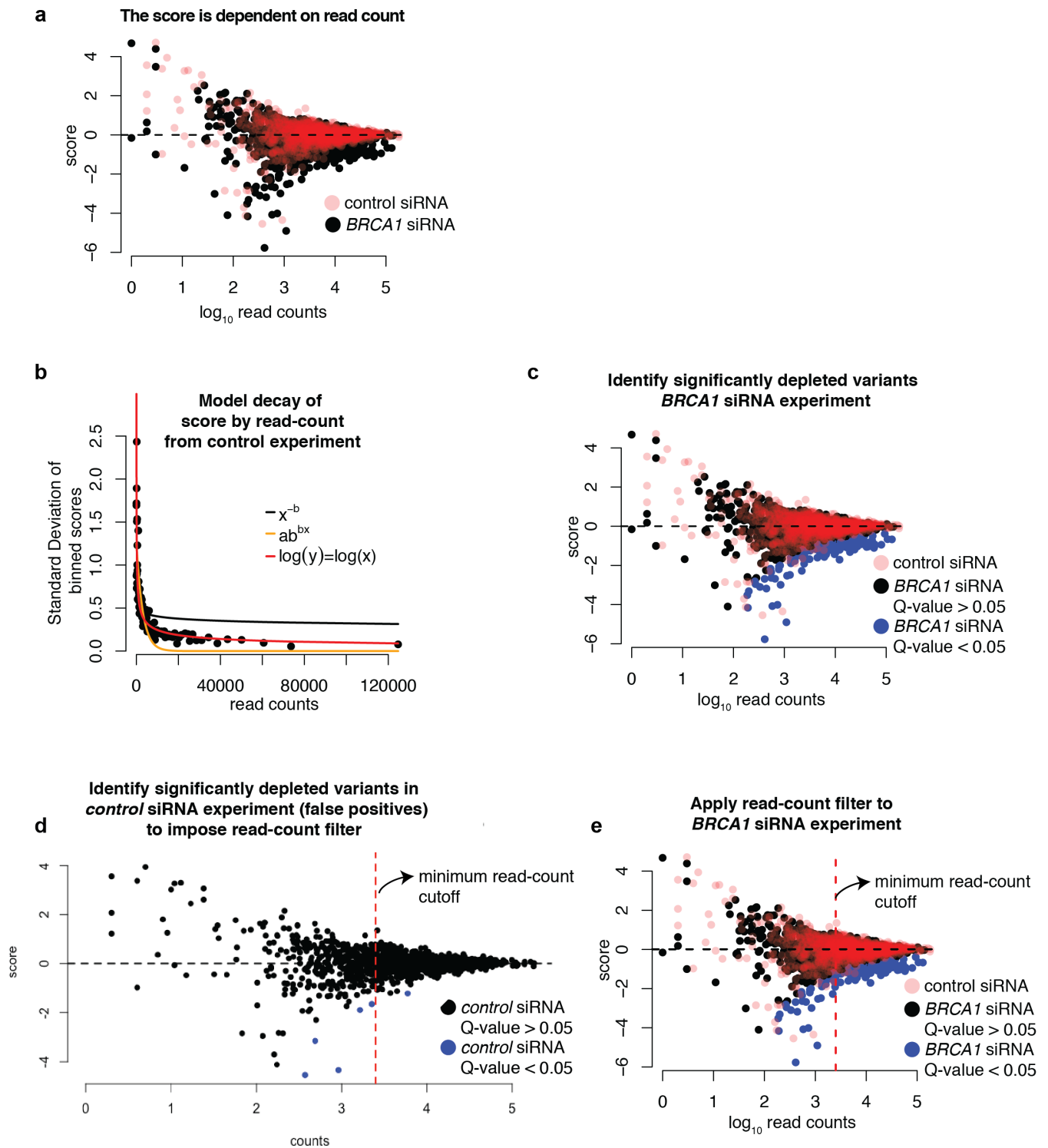


Pool3 replicates correlation of WT-normalized-GFP-positive/GFP-negative ratios



Supplementary Figure 2 | Scatter plots of multiplex HDR reporter assay scores from replicate experiments. Only variants above the read count threshold are included. Missense variants are grey and nonsense variants are orange. Pearson's r reported.

Supplementary Figure 3



Supplementary Figure 3 | Illustration of the depletion classifier and read count thresholds using pool 1, replicate 1.

a, The log of the WT-normalized, GFP-positive:GFP-negative ratios are on the y-axis and \log_{10} read counts are on the x-axis for a single replicate of the HDR-reporter assay for codons 2-96. Variants from the control (red) or BRCA1 (black) siRNA conditions are indicated. **b**, $\text{Log}(x) = \log(y)$ best fits the decay of the standard deviation of the score in 100 variant bins, given the read counts in the control siRNA condition. Lines representing each model are indicated. **c**, The same plot as in (a) with variants significantly depleted from the GFP-positive population in the *BRCA1* siRNA condition, $q < 0.05$, colored blue. **d**, The variants in the control siRNA population (here, colored black) with the false positive ($q < 0.05$) colored blue. The dashed line represents the read-count threshold chosen to minimize the number of false positives. **e**, The same plot as in (a) with variants significantly depleted from the GFP-positive population in the *BRCA1* siRNA condition, $q < 0.05$, colored blue and the dashed line (from d) represents the read-count threshold.

All R code and plots for the remaining pools and replicates can be found in the R Markdown supplied as supplementaryFile1.Rmd. Figures can be remade using the markdown file and data provided in SupplementaryTable4.tsv

Supplementary Table 1 | Metrics for BRCA1 barcoded variant library construction, cell integration, barcode assignments and library quality.

aa = amino acid, CFU = colony forming units

library	pool1	pool2	pool3
BRCA1 variable region aa positions	2-96	97-192	192-302
inverse PCR clones	3.2 M	1.6 M	360 K
pcDNA5 cloning CFU	180 K	80 K	140 K
barcode cloning CFU	25 K	25 K	50 K
HeLa HDR FRT integrations	54 K	150 K	35 K
PacBio SMRT cells	4	4	4
PacBio reads with more than 3 passes	141720	162351	78875
passed alignment filters (min mapQ 30, min quality score per base = 20, no softclipping, correct size barcode)	61906	73254	
passed alignment filters (min mapQ 30, min quality score per base = 0, no softclipping, correct size barcode)			54001
unique barcodes	27433	24497	26200
barcode had one consensus read	10963	6632	11945
barcode had two consensus reads	7484	5710	6259
barcode had three consensus reads	4283	4320	2669
barcode had > 3 consensus reads	4703	7835	5327
all reads associated with a barcode were identical	15966	17266	3193
read assigned by major allele	363	488	1836
read assigned by highest quality score tie breaker	141	111	5471
barcodes associated with WT or single amino acid substitution	19809	17635	11857
barcodes associated with indel	6689	6046	13426
barcodes associated with >1 aa substitution	935	816	917
Unique variants with 0 or 1 aa substitution	1602	1695	1987

Supplementary Table 2 | Metrics for FACS sorts, PCR amplification and DNA sequencing for replicate multiplexed HDR reporter assays.

GEq = genome equivalents (9 pg per triploid HeLa genome)

Pool	1	1	1	1	2	2	2	2	3	3	3	3
Sort	1	2	3	4	1	2	3	4	1	2	3	4
control siRNA % GFP+	0.11	0.113	0.156	0.172	0.158	0.157	0.129	0.146	0.104	0.103	0.155	0.163
BRCA1 siRNA % GFP+	0.07	0.073	0.074	0.074	0.099	0.102	0.087	0.1	0.075	0.075	0.084	0.085
control siRNA GFP-sorted cells	2.5M	3M	2M	2M	2M	2M	2.3M	2.3M	2M	2M	2M	2M
control siRNA GFP+ sorted cells	0.77M	1.14M	0.5M	0.8M	0.5M	0.6M	0.627M	0.75M	0.574M	0.729M	0.614M	0.798M
BRCA1 siRNA GFP-sorted cells	2M	3M	2M	2M	2M	2M	2M	2.6M	2M	2M	2M	2M
BRCA1 siRNA GFP+ sorted cells	0.7M	0.7M	0.5M	0.55M	0.892M	0.8M	0.8M	0.8M	0.663M	0.554M	0.655M	1M
control siRNA GFP- GEq	444444	444444	444444	444444	444444	444444	444444	444444	444444	444444	444444	444444
control siRNA GFP+ GEq	416666	444444	217800	416666	233200	206800	180400	440000	398200	444444	444444	444444
BRCA1 siRNA GFP- GEq	444444	444444	444444	444444	444444	444444	444444	444444	444444	444444	444444	444444
BRCA1 siRNA GFP+ GEq	277777	333333	277988	250000	413600	413600	415800	382800	444444	444444	444444	444444
control siRNA GFP-reads	28102878	16433473	13253616	11499041	6291880	1163211	2899216	2669299	2923995	4100215	4487348	3222144
control siRNA GFP+ reads	25652248	21453958	11548560	9448454	5749200	3057588	1043964	2713935	2777241	3554904	4672780	2978423
BRCA1 siRNA GFP-reads	21022149	36084108	14097458	13104520	4899883	1999841	4951048	3010396	4166809	3407709	2968153	3183004
BRCA1 siRNA GFP+ reads	17916667	24168968	7140895	16885180	3822271	13841768	2335473	3723491	3716272	4378930	6062421	3445223

Supplementary Table 3 | Primer sequences used for this study.

name	primer
pc5_barcode_longer_W	ggccggCCTGCAGNNNNNNNNNNNNNNNNNNCTGATGCGATGACGATGACGATGCATGGG
pc5_barcode_longer_C	catgctgGCGGCCGCCCATCGTCATCGTACTGCATCCCATGCATCGTCATCGTCATCG
pb_EcoRI	/5PHOS/AAT TAT CTC TCT CTT TTC CTC CTC CTC CGT TGT TGT TGT TGA GAG AGA T
pb_SbfI	/5PHOS/AT CTC TCT CTT TTC CTC CTC CTC CGT TGT TGT TGT TGA GAG AGA TTG CA
SV40_F	GAA GTA GTG AGG AGG CTT TTT TGG AGG CTA CC
new_pc5bc_nexteraR	GGCTCGGAGATGTGTATAAGAGACAGCTATGAACTAATGACCCCGTAATTGATTACTA
pc5bc_p5_F	AATGATACGGCGACCACCGAGATCTACACGAGCAAAATTTAAGCTACAACAAGG
nextIndex	CAAGCAGAAGACGGCATAACGAGATNNNNNNNNNGTCTCGTGGGCTCGGAGATGTGTATAAGAGACAG
pcDNA5_barcodeSeq_F	GCT TAG GGT TAG GCG TTT TGC GCT GCT CCT GCA GG

Supplementary Table 4 | Counts, scores and q - values for all variants in all replicates.

Counts: Column headers have the name of the siRNA treatment (brca1 or control), replicate (sort 1-4), and population (noGFP or GFP). Columns appended with _counts are raw sequencing counts for each variant.

Score: Column headers have the name of the siRNA treatment (brca1 or control), replicate (sort 1-4) and are appended with _score. The score calculated by taking the ratio of GFP/noGFP of the frequency of each variant in each population, normalized to that of WT, scores are reported as the natural log.

Pool: The pool represents which amino acids are varied, pool_1 2-96, pool_2 97-192, pool_3 = 192-302.

Variant: Protein variants are listed using HGVS nomenclature.

ProtPos: Amino acid position

mut: mutant amino acid

WT: Wild-type amino acid

pval: Column headers have the name of the siRNA treatment (brca1 or control), replicate (sort 1-4). Columns appended with _pval are the raw p-value describing the significance of the score difference between a variant in the BRCA1 population and variants in the control siRNA population at a given read count.

qval: Column headers have the name of the siRNA treatment (brca1 or control), replicate (sort 1-4). Columns appended with _qval are false discovery rate adjusted p-values⁶.

Sigcol: Column headers have the name of the siRNA treatment (brca1 or control), replicate (sort 1-4). Columns appended with _sigcol determine the color of the points for the variants. $q < 0.05$ = blue, the remainder are black.

Supplementary Table 5 | Depletion scores for 1,700 BRCA1 variants.

Depletion scores for the 1,699 variants that passed the read count threshold in 3 or 4 replicates.

Variant: Protein variants are listed using HGVS nomenclature.

ProtPos: Amino acid position

mut: mutant amino acid

WT: Wild-type amino acid

variantID: Wild-type amino acid concatenated with protPos and mut

mHDR_repPass: Number of replicates in which the variant passed the read count threshold, variants with only 3 or 4 are included.

mHDR_depletionScore: Number of replicates in which the variant was significantly depleted from the BRCA1 siRNA GFP-positive population ($q < 0.05$).

clinvar_clinsig: The ClinVar variant interpretation as of June, 2017.

Parvin_HDR_average: singleton HDR scores from refs. 1–4.

HDR_function_cat: singleton HDR functional category. HDR > 0.5 = high, HDR < 0.5 = low.

Starita_Y2H_score: Multiplexed yeast-two-hybrid scores from ref. 3.

Starita_E3_score: Multiplexed ubiquitin ligase scores from ref. 3.

Starita_HDR_predict: HDR predictions from ref. 3.

Bouwman_class: functional predictions from ref. 7.

SGE_functionScore_aveVar: Saturation Genome Editing (SGE) function score averaged for SNVs that make the same protein variant (Findlay et al. enclosed).

SGE_lowRNAcount: The number of times an SNV at a the codon position causes a >75% reduction in RNA as measured by SGE (Findlay et al. enclosed).