

1
2 **Targeted genome fragmentation with CRISPR/Cas9 improves hybridization capture, reduces PCR**
3 **bias, and enables efficient high-accuracy sequencing of small targets**
4

5 Daniela Nachmanson¹, Shenyi Lian^{1†}, Elizabeth K. Schmidt¹, Michael J. Hipp¹, Kathryn T. Baker¹,
6 Yuezheng Zhang¹, Maria Tretiakova¹, Kaitlyn Loubet-Seneor¹, Brendan F. Kohn¹, Jesse J. Salk^{2‡}, Scott
7 R. Kennedy^{1*}, Rosa Ana Risques^{1*}
8

9 ¹Department of Pathology, University of Washington, Seattle, WA 98195, USA. ²Department of
10 Medicine, Division of Hematology and Oncology, University of Washington, Seattle, WA 98195, USA.
11

12 †Current address: Key laboratory of Carcinogenesis and Translational Research (Ministry of
13 Education/Beijing), Department of Pathology, Peking University Cancer Hospital & Institute, Beijing, PR,
14 China.

15 ‡ Current address: TwinStrand Biosciences, Seattle, WA 98121, USA.
16

17 *These authors contributed equally
18

19 **Corresponding authors:**

20
21 Scott Kennedy, PhD
22 Dept. of Pathology
23 University of Washington
24 Box 357470
25 1959 NE Pacific Ave.
26 Seattle, WA 98195-7705
27 (206) 543-5452
28 scottrk@uw.edu
29
30

31 Rosa Ana Risques, PhD
32 Dept. of Pathology
33 University of Washington
34 Box 357470
35 1959 NE Pacific Ave.
36 Seattle, WA 98195-7705
37 (206) 616 4976
38 rrisques@uw.edu
39
40
41
42
43
44
45

46 **ABSTRACT**

47 Current next-generation sequencing techniques suffer from inefficient target enrichment and frequent
48 errors. To address these issues, we have developed a targeted genome fragmentation approach based
49 on CRISPR/Cas9 digestion. By designing all fragments to similar lengths, regions of interest can be size-
50 selected prior to library preparation, increasing hybridization capture efficiency. Additionally,
51 homogenous length fragments reduce PCR bias and maximize read usability. We combine this novel
52 target enrichment approach with ultra-accurate Duplex Sequencing. The result, termed CRISPR-DS, is
53 a robust targeted sequencing technique that overcomes the inherent challenges of small target
54 enrichment and enables the detection of ultra-low frequency mutations with small DNA inputs.

55

56 **Key words:** Target enrichment, Duplex Sequencing, Next-Generation Sequencing, NGS, CRISPR/Cas9

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72 **BACKGROUND**

73 In the past decade, NGS has revolutionized the fields of biology and medicine. However, standard
74 NGS suffers from two major problems that negatively impact multiple applications: the limited efficiency
75 of the current target selection methods and the high error rate of the sequencing process. Targeted
76 genome enrichment is essential to many applications that do not require whole genome sequencing and
77 it is performed either by PCR or by hybridization capture. PCR is simple and efficient but does not scale
78 well and suffers from biases that result in uneven coverage and false mutation calls [1, 2]. Hybridization
79 capture improves coverage uniformity and mutation call accuracy but has low recovery, especially when
80 the target region is small, which leads to the requirement of larger amounts of DNA [2]. An additional
81 complication is that DNA is typically fragmented by sonication which introduces DNA damage resulting
82 in sequencing errors [3]. Moreover, the heterogeneous fragment sizes generated by sonication are
83 subject to PCR bias and contribute to uneven coverage. An alternative option to sonication is enzymatic
84 fragmentation. This method resolves some issues but introduces different artifacts that also result in
85 sequencing errors [4]. Thus, at the library preparation step, both methods of target selection suffer
86 important limitations that lead to non-optimal sequencing outcomes, including uneven coverage,
87 introduction of false mutations, and low recovery.

88 The second major problem of NGS is the high error rate inherent to the sequencing process.
89 Illumina offers the most accurate sequencing platform with an estimated error rate of 10^{-3} [5]. This error
90 rate, however, translates into millions of false calls in each sequencing run and precludes the detection
91 of low frequency mutations (Additional file 1: Figure S1), which is critical for applications such as
92 forensics, metagenomics, and oncology [6]. While the accuracy of NGS can be improved by repair of the
93 DNA prior to sequencing [7, 8] and by computational error correction [9, 10], these strategies do not
94 remove all potential artifacts. An alternative approach employs unique molecular identifiers, also called
95 molecular barcodes or molecular tags, to identify the reads derived from an original DNA molecule and
96 use their redundant information to create a consensus sequence [11]. The unique random shear points
97 generated at sonication can be used as “endogeneous barcodes”, but only when sequencing depth is

98 low (~10x) to avoid overlapping of sharing points between independent DNA molecules [12]. To enable
99 higher sequencing depth, exogenous barcodes are necessary. Exogenous barcodes are random DNA
100 sequences attached to the original DNA molecules before or during PCR. Single-stranded molecular
101 barcodes produce a consensus with the reads derived from one DNA strand [11], whereas double-
102 stranded molecular barcodes introduce an additional level of correction by allowing the comparison of
103 independent consensus sequences derived from the two complementary strands of the original DNA
104 molecule [13]. This additional level of correction is essential for removing polymerase errors occurring in
105 the first round of PCR and subsequently propagated to all reads derived from a given DNA strand [7].
106 Polymerase errors caused by DNA damage are one of the most pervasive problems of NGS [8] but can
107 be successfully addressed with double-strand barcodes given the extremely low probability that the same
108 error occurs in the same position on both strands of DNA. Duplex Sequencing (DS), the method that
109 pioneered double-strand molecular barcodes [13, 14], has an estimated error rate $<10^{-7}$, four orders of
110 magnitude less than single-strand molecular barcode methods. This level of accuracy allows for very
111 sensitive ultra-deep sequencing (Additional file 1: Figure S1) and has been employed in a variety of
112 applications including the detection of very low frequency somatic mutations in cancer and aging [15-18].

113 DS successfully addresses the problem of sequencing errors, but it suffers from the limitations of
114 hybridization capture, which is required to perform target selection while preserving the strand recognition
115 of molecular barcodes. As described above, hybridization capture is highly inefficient when selecting
116 small target sizes [19]. It is estimated that for targets $<50\text{Kb}$ only 5-10% of reads are on-target after
117 hybridization capture [20]. In DS, as well as in other panel-based sequencing approaches, the region of
118 interest is usually small as a cost-effective trade-off for higher sequencing depth. In this situation, a
119 successful approach for target enrichment is to perform two consecutive rounds of capture [20]. However,
120 this approach results in a time consuming, costly, and inefficient protocol that requires large amounts of
121 DNA [14]. For example, in DS at least $1\mu\text{g}$ of DNA was historically needed to produce depths $>3,000\text{x}$
122 [17], which is prohibitive in many applications that rely on small samples.

123 Here we present CRISPR-DS, a new method that addresses the two main problems of NGS:
124 limited efficiency of target selection and high error rate. Target selection is facilitated by an enrichment
125 of the regions of interest using the CRISPR/Cas9 system. *In vitro* digestion with CRISPR/Cas9 has been
126 proven to be a useful tool for multiplexed excision of large megabase fragments and repetitive sequence
127 regions for PCR-free NGS [21, 22]. We reasoned that targeted *in vitro* CRISPR/Cas9 digestion could be
128 used to excise similar length fragments covering the area of interest, which could then be enriched by
129 size selection prior to library preparation. We designed this method to enable target enrichment while
130 simultaneously eliminating sonication-related errors and biases arising from random genome
131 fragmentation. In addition, by pairing this method with double-strand molecular barcoding, we aimed to
132 produce a method that preserves the sequencing accuracy of DS, while increasing the recovery rate,
133 enabling low DNA input and a simplified protocol for translational applications.

134

135 RESULTS

136 Design of CRISPR-DS based on CRISPR/Cas9 target fragmentation and double strand molecular 137 barcodes

138 CRISPR-DS is based on *in vitro* CRISPR/Cas9 excision of target sequences to generate DNA
139 molecules of uniform length which are then enriched by size selection. The versatility, specificity, and
140 multiplexing capabilities of the CRISPR/Cas9 system enable its application for the excision of any target
141 region of interest by simply designing guide RNAs (gRNA) to the desired cutting points. As a proof of
142 principle, we developed the method for sequencing the exons of *TP53*. Further, in order to achieve high
143 recovery as well as high sequencing accuracy, we combined it with DS. The main steps of the protocol
144 are illustrated in Figure 1. First, target regions are excised from genomic DNA by multiplexed *in vitro*
145 CRISPR/Cas9 digestion (Fig. 1a), followed by enrichment of the excised fragments by size-selection
146 using SPRI beads (Fig. 1b). The selected fragments are then coupled with the double-strand molecular
147 barcodes used in DS (Fig. 1c) [14]. These fragments are then amplified and captured with biotinylated
148 hybridization probes as previously described for DS [14], with the exception that only one round of

149 hybridization capture is required due to the prior enrichment of target fragments (see below). Finally, the
150 library is sequenced and the resulting reads are analyzed to perform error correction based on the
151 consensus sequences of both strands of each DNA molecule (Fig. 1d) [14]. Due to the requirement of
152 only one round of hybridization capture, the workflow of CRISPR-DS is almost one day shorter than
153 standard-DS (Fig. 2, Additional file 1: Figure S2), enabling a more cost-efficient and applicable method.

154

155 **CRISPR/Cas9 cut fragments can be designed to be of homogenous length, reducing PCR bias** 156 **and producing uniform coverage**

157 Typically, genome fragmentation is performed with sonication, which generates randomly sized
158 fragments that have different amplification efficiencies [23]. Short fragments are preferentially amplified,
159 resulting in uneven coverage of the regions of interest and decreased recovery. In DS, amplification bias
160 introduces an additional problem because short fragments produce an excess of PCR copies that do not
161 further aid error reduction. To produce a consensus, only three PCR copies of the same molecule are
162 required. Additional copies waste resources because they produce sequencing reads but do not generate
163 additional data. By using CRISPR/Cas9, gRNA can be designed such that restriction with Cas9 produces
164 fragments of predefined, homogeneous size. We reasoned that these fragments would eliminate PCR
165 bias, leading to homogeneous sequencing coverage and minimizing wasted reads that are PCR copies
166 of the same original molecule.

167 To test this approach, we designed gRNAs to specifically excise the coding regions and their
168 flanking intronic sequence of *TP53* (Fig. 1a). Fragment length was designed to be ~500bp in order to
169 maximize read space of an Illumina MiSeq v3 600 cycle kit while allowing for sequencing of the molecular
170 barcode (10 bp) and 3'-end clipping of 30bp to remove low-quality bases produced in the later sequencing
171 cycles. gRNAs were selected based on the highest specificity score that produced appropriate fragment
172 length (Additional file 2: Table S1, Additional file 3: Data S1) [24]. The fragment comprising exon 7 was
173 designed shorter than the rest (336 bp) to avoid a homopolymeric run of T's in the flanking intronic region
174 which induced poor base quality in reads that span this region (Additional file 1: Figure S3).

175 We performed a side-by-side comparison of library performance (Fig. 3a-c) and sequencing
176 coverage (Fig. 3d) of a sample DNA processed with CRISPR-DS vs standard-DS (see Material and
177 Methods). Standard-DS for *TP53* had been previously performed using sonication and published
178 protocols [14, 17]. Visualization of the resulting sequencing library by gel electrophoresis showed that
179 CRISPR restriction produced distinct bands/peaks (Fig. 3a-b) corresponding to the predesigned size of
180 target fragments as opposed to the diffuse “smear” characteristic of libraries prepared by sonication. The
181 discrete peaks allow confirmation of correct library preparation and target enrichment, preventing the
182 sequencing of suboptimal libraries. Sequencing and mapping of the libraries demonstrated that targeted
183 Cas9 restriction results in well-defined DNA fragments corresponding to the expected size (Fig. 3d).
184 Importantly, these fragments exhibited extremely uniform sequencing depth. In contrast, sonicated DNA
185 fragments resulted in significant variability in depth across target regions. Because DS reads correspond
186 to individual DNA molecules, the uniform depth achieved by CRISPR-DS indicates a homogenous
187 representation of the original genomic DNA in the final sequencing output, confirming the proper excision
188 of all fragments.

189 The ability to uniformly control the DNA insert size should not only provide homogenous depth,
190 but should also produce a more uniform number of copies of each molecule, minimizing the waste of
191 unnecessary reads to produce a consensus sequence. To test this possibility, we counted the number of
192 PCR copies for each molecular barcode and plotted it as a function of the DNA fragment size (Fig. 3c).
193 Sonicated DNA exhibited a strongly negative association between DNA fragment size and the number of
194 PCR copies, as expected due to the fact that small DNA fragments are preferentially amplified (Fig. 3c,
195 *blue*). In contrast, targeted fragmentation produced a consistent number of PCR copies for all fragments,
196 including the smaller exon 7 fragment (Fig. 3c, *red*).

197

198 **CRISPR/Cas9 cut fragments can be designed to be of optimal length to maximize read usage**

199 An additional disadvantage of the variable fragment size produced by sonication is inefficient read
200 usage: fragments that are too short generate overlapping reads that waste sequencing space, whereas

201 fragments that are too long get sequenced on the ends, leaving captured but un-sequenced DNA in the
202 middle (Fig. 4a). The programmable nature of Cas9 can be leveraged to reduce the amount of data “lost”
203 by generating optimal length fragments tailored to the preferred number of sequencing cycles. To
204 illustrate the improvement in read usage, we quantified the amount of deviation from the optimal fragment
205 size (defined as the total number of sequencing cycles minus the total length of the molecular barcodes
206 and 3'-end clipping) of seven samples independently processed with sonication and targeted
207 fragmentation. Sonication produced significant variability in the amount of deviation from the optimal
208 fragment size with a large fraction of fragments being twice the optimal size for one of the samples (Fig.
209 4b,c; Additional file 1: Figure S4). Indeed, only $9.1\pm 4.2\%$ of reads had inserts that were within 10%
210 deviation from the optimal fragment length. Even samples with more stringent size selection had only
211 ~61% of reads within the 10%-deviation window (Fig. 4c; Additional file 1: Figure S4). In contrast, the
212 same samples fragmented with Cas9 had $71.0\pm 3.2\%$ of reads within the same window range, with the
213 vast majority of the reads outside the window being due the purposefully shorter Exon 7 fragment (Fig.
214 4b,c; Additional file 1: Figure S3, S4). Exclusion of exon 7 from this analysis improved the percent of
215 reads within the 10%-deviation window to $94.3\pm 2.1\%$. These data indicate that targeted fragmentation
216 can tightly control the fragment size to optimize read usage, thereby increasing the efficiency of
217 sequencing.

218

219 **CRISPR/Cas9 fragmentation enables target enrichment by size selection, eliminates one round of** 220 **hybridization capture, and increases sequencing yield**

221 While performing two rounds of capture substantially increases the number of on-target reads for
222 standard-DS and other small target applications, the process is time consuming, expensive, and requires
223 additional PCR steps that introduce further bias [20]. We hypothesized that target enrichment via size
224 selection of CRISPR/Cas9 digested fragments would sufficiently enrich for on-target DNA fragments and
225 eliminate the need for a second capture. To test this hypothesis, we performed CRISPR/Cas9 digestion
226 of targeted *TP53* exons (Fig. 1a) on a range of DNA input amounts (10-250ng) followed by SPRI size

227 selection to remove undigested high molecular weight DNA fragments (> 1kb in size). The selected DNA
228 fragments were ligated to DS adapters, PCR amplified, and sequenced (see Material and Methods). No
229 hybridization capture or any other type of target enrichment was performed. Mapping of raw reads
230 revealed between 0.2% to 5% reads on-target, corresponding to ~2,000x to 50,000x fold enrichment
231 given the fact that our target region only amounted for 0.000101% of the human genome (Table 1). This
232 level of enrichment matches or exceeds what is typically achieved with solution based hybridization for
233 small target size [19, 20]. Notably, lower DNA inputs showed the highest enrichment, potentially reflecting
234 more efficient digestion or improved removal of off-target, high molecular weight DNA fragments when
235 they are in lower abundance. These results suggested that a simple size selection step can be used in
236 lieu of a targeted hybridization enrichment step.

237 To test this possibility, we performed a side by side comparison of standard-DS [14] with one and
238 two rounds of hybridization capture vs. CRISPR-DS with only one round of hybridization capture. Three
239 input amounts of the same control DNA extracted from normal human bladder tissue were sequenced in
240 parallel for each of the methods. CRISPR-DS with one round of capture achieved >90% raw reads on-
241 target (e.g. covering *TP53*) (Fig. 5a), a significant improvement over standard-DS which only achieved
242 ~5% raw reads on-target with a single capture, consistent with prior work [20]. In an independent
243 experiment, we tested the reproducibility of this result with three different DNA samples that were
244 sequenced with CRISPR-DS using one and two rounds of capture (Additional file 1: Figure S5).
245 Confirming the prior result, the three samples produced >90% raw reads on target using only one round
246 of capture. The second round of capture only minimally increased raw reads on-target and is
247 unnecessary.

248 The side-by-side comparison of CRISPR-DS vs standard-DS also demonstrated a substantial
249 increase in recovery using CRISPR-DS. Sequencing recovery, also referred to as yield, is typically
250 measured as the fraction or percentage of sequenced genomes compared to input genomes. Consistent
251 with prior studies[13, 17], standard-DS produced a recovery rate of ~1% across the different inputs, while
252 CRISPR-DS recovery rate ranged between 6 and 12% (Fig. 5b). Notably, 25ng of DNA prepared with

253 CRISPR-DS produced a post-processing depth comparable to 250ng with standard-DS. This indicates that
254 size selection for CRISPR/Cas9 excised fragments not only removes a step from the library preparation
255 but, most importantly, increases the recovery of input DNA enabling deep sequencing with greatly reduced
256 DNA requirements.

257

258 **Validation of CRISPR-DS recovery in an independent set of samples, including low quality DNA**

259 We further confirmed the performance of CRISPR-DS in an independent set of 13 DNA samples
260 extracted from bladder tissue (Additional file 2: Table S3). We used 250ng and obtained a median DCS
261 depth of 6,143x, corresponding to a median recovery rate of 7.4% in agreement with the prior experiment.
262 Reproducible performance was demonstrated with technical replicates for two samples (B2 and B4,
263 Additional file 2: Table S3). All samples had >98% reads on-target after consensus making, but the
264 percentage of on-target raw reads ranged from 43% to 98%. We noticed that the low target enrichment
265 corresponded to samples with DNA Integrity Number (DIN) <7. DIN is a measure of genomic DNA quality
266 ranging from 1 (very degraded) to 10 (not degraded) [25]. We reasoned that degraded DNA compromises
267 enrichment by size selection, and the poor yield could be mitigated by removing low molecular weight DNA
268 prior to CRISPR/Cas9 digestion. To test this hypothesis, we used the pulse-field feature of the BluePippin
269 system to select high molecular weight DNA (> 8kb) from two samples with degraded DNA (DINs 6 and 4).
270 This pre-enrichment resulted in successful removal of low molecular weight products and increased on-
271 target raw reads by 2-fold and DCS depth by 5-fold (Additional file 1: Figure S6). These results indicate that
272 enrichment of high molecular weight DNA could be used as a solution for successful CRISPR-DS
273 performance in partially degraded DNA.

274

275 **Validation of CRISPR-DS for the detection of low-frequency mutations**

276 To validate the ability of CRISPR-DS to detect low-frequency mutations, we analyzed four peritoneal
277 fluid samples collected during debulking surgery from women with ovarian cancer and previously analyzed
278 for *TP53* mutations using the standard-DS protocol [17]. The tumor mutation was previously identified in

279 the four samples: in one sample at a high frequency (68.5%) and at a very low frequency (around or below
280 1%) in the remaining 3 samples. CRISPR-DS detected the tumor mutation in all samples at frequencies
281 comparable to what was reported in the original study (Table 2) [17]. In addition to the tumor mutation,
282 standard-DS also revealed the presence of additional exonic *TP53* mutations in these samples which were
283 at an extremely low frequency (<0.1%) in all cases. These mutations are considered “biological
284 background” mutations to distinguish them from the tumor-derived mutations [17]. Standard-DS revealed
285 between 1 to 5 biological background mutations in each of the samples, representing an overall mutation
286 frequency of about $\sim 1 \times 10^{-6}$. Similarly, CRISPR-DS identified biological background mutations in the 4
287 samples at a comparable overall mutation frequency (Additional file 1: Figure S7). These results indicate
288 that CRISPR-DS preserves the sequencing accuracy and sensitivity for mutation detection previously
289 described for DS [13, 17].

290 Table 2 also illustrates a critical advantage of CRISPR-DS compared to standard-DS in terms of
291 translational applicability: the reduced requirement of input DNA as a result of a more efficient library
292 preparation method that enables higher recovery. Standard-DS of these peritoneal fluid samples required
293 between 3-10 μg of DNA to compensate for the $\sim 1\%$ recovery rate of standard-DS and to achieve the high
294 depth necessary to detect low frequency tumor mutations. With CRISPR-DS we only used 100ng of DNA
295 (30-100 fold less than what was used for standard-DS), and we obtained comparable DCS depth to
296 standard-DS (Table 2). Recovery rates ranged between 6 and 12%, as in prior experiments (Fig. 5 and
297 Additional file 2: Table S3). These results represent an efficiency increase of 15x-200x compared to
298 standard-DS with the same DNA. Notably, CRISPR-DS not only preserved sensitivity for mutation
299 detection, increased sequencing recovery, and reduced DNA input, but also shortened the protocol by
300 nearly one day (Additional file 1: Figure S2), making it a more cost effective option for accurate deep
301 sequencing of samples with limited DNA amounts.

302

303 **DISCUSSION**

304 While CRISPR-based target enrichment can be applied to any sequencing method that requires
305 hybridization capture of small targets, here we have leveraged its qualities for the optimization of DS,
306 producing a new method called CRISPR-DS. CRISPR-DS merges the increases in efficiency provided
307 by CRISPR-based targeted genome fragmentation with the high accuracy of sequencing provided by
308 double strand molecular barcodes, thus enabling ultra-accurate sequencing of small target regions using
309 minimal DNA inputs. In addition to CRISPR-DS, the CRISPR-based target enrichment approach can be
310 used in combination with other methods for targeted sequencing to improve recovery of small targets and
311 to reduce PCR bias and uneven coverage arising from random fragment sizes.

312 Targeted sequencing remains a cost effective alternative to whole genome-sequencing,
313 especially when high depth is desired [2]. In multiple applications, such as oncology, the goal is to
314 sequence a small panel of relevant genes with high accuracy in order to find low frequency mutations.
315 While the selected target panel can be amplified by PCR, this method creates uneven coverage and false
316 mutations, thus hybridization capture is typically preferred [2]. Hybridization capture improves coverage
317 uniformity and removes certain artifactual mutations but does not resolve these issues completely. A
318 major disadvantage in hybridization-based sequencing methods is the reliance on sonication for genome
319 fragmentation which generates DNA fragments of random size. We have demonstrated that this size
320 heterogeneity generates two problems that can be solved by replacing sonication with CRISPR-based
321 genome fragmentation. The first problem is PCR bias, which results in the preferential amplification of
322 short DNA fragments. PCR bias leads to wasted reads that contain an excess of PCR copies of the same
323 molecule. While these reads can be removed bioinformatically [26], the amplification advantage of certain
324 molecules can lead to uneven coverage and reduced recovery [27]. In methods that employ molecular
325 barcodes, such as DS, three PCR copies are typically sufficient to generate a consensus sequence [14].
326 Thus, additional sequencing of PCR copies does not produce additional data and only wastes resources.
327 We have demonstrated that with CRISPR-based fragmentation all fragments amplify similarly. This
328 homogeneous amplification translates into uniform coverage across all targeted regions, a critical feature
329 when the goal is to detect low frequency mutations in selected panel of genes.

330 The second problem associated with the heterogeneous fragment sizes relates to reduced data
331 yield at the read level. Because sonication allows minimal control over fragment size, a large proportion
332 of fragments are typically too short or too long compared to the optimal length size determined by the
333 number of sequencing cycles. When reads are too short, paired-end reads overlap and the middle region
334 is double-sequenced. Conversely, when reads are too long, the middle part of the DNA fragment, which
335 may contain a variant or region of interest, remains un-sequenced. This inefficient read usage is solved
336 with CRISPR-based target selection because the fragments are tailored to the desired read length.

337 CRISPR-based target fragmentation also offers two additional advantages. First, homogeneously
338 sized DNA fragments can be visualized to confirm library target enrichment prior to sequencing. In
339 sonication-based hybridization capture, the gel electrophoresis for a target-enriched library looks identical
340 to a library with no target enrichment. This issue can result in the costly waste of a sequencing run where
341 the majority of reads are in off-target regions. We show that the defined fragment lengths created by
342 CRISPR-based digestion produce distinct peaks which are easily visualized and confirm that the
343 sequencing library is target-enriched. A second advantage of Cas9 digestion over sonication is the
344 elimination of sonication-induced sequencing errors [3] and the preservation of double stranded DNA at
345 the ends of fragments. Sonication produces ssDNA at the end of molecules which is susceptible to
346 damage and converted into “pseudo-dsDNA” by end repair. This process has the potential to introduce
347 false variant calls, but it is prevented by CRISPR-DS because Cas9 produces blunt ends which do not
348 require end repair.

349 In the context of small target sequencing by hybridization-capture, the major advantage
350 introduced by CRISPR-based target enrichment is increased recovery, that is, percentage of input
351 genomes that produce sequencing data. Hybridization capture is notably inefficient, especially for small
352 target regions [19, 20]. As demonstrated with our experiments and in agreement with prior studies, the
353 average recovery rate of DS is ~1% which translates to at least 1 μ g of DNA being needed to produce
354 an average depth of ~3,000x. This recovery is improved 10-fold by the addition of CRISPR-based target
355 enrichment and the elimination of one round of capture. We have demonstrated that by simply excising

356 the genomic regions of interest and performing size selection, we can achieve a level of enrichment
357 comparable to a single round of capture. By performing this step prior to library preparation, only one
358 round of hybridization capture is needed, greatly minimizing DNA loss and increasing recovery.
359 Therefore, using CRISPR-based target enrichment prior to DS achieves the same depth with 10 times
360 less DNA.

361 To take advantage of the accuracy of DS while enabling low DNA inputs, several groups have
362 developed DS-based approaches that combine endogenous and exogenous barcodes. Yet each comes
363 with its own set of compromises. BotSeqS, iDES, and SIP-HAVA-Seq all use DS-based error correction
364 and require little DNA input [12, 28, 29], but the reliance on endogenous barcodes means that depth is
365 limited in order to keep shearing points unique. BiSeqS uses chemical conversions to distinguish one
366 strand from another in combination with molecular barcodes [30] which allows for an increased recovery
367 and high sequencing depth. However, as a consequence of the chemical conversions, it is unable to
368 detect all mutation types. In contrast, CRISPR-DS preserves the sequencing accuracy of DS because it
369 relies exclusively on exogenous double strand molecular barcodes, and the error correction method and
370 analytical algorithms remain identical to standard-DS. We have demonstrated that CRISPR-DS identified
371 very low frequency mutations previously detected by DS, confirming its sensitivity. Remarkably this
372 validation experiment was performed with 10 to 100 times less DNA than the original standard-DS
373 experiment, illustrating a significant improvement in recovery that will enable the use of the extreme
374 sensitivity of DS for mutation detection in samples with low input DNA.

375 Though CRISPR-DS addresses several needs in targeted NGS, it could still benefit from
376 optimizations. First, improvements could be made to increase the recovery of degraded samples.
377 Currently, in order to perform efficient target enrichment with CRISPR/Cas9 digestion and size selection,
378 degraded samples must be pre-processed to remove low molecular weight fragments. We performed
379 this pre-processing using electrophoretic size selection with the BluePippin system. However, to minimize
380 loss of DNA, high molecular weight DNA could be selected with alternative methods such as micro-
381 column filters. Second, we noticed that the best recovery was achieved with smaller inputs of DNA. Since

382 our goal was to achieve higher depth with smaller amounts of input DNA, this was not problematic.
383 However, further efforts should be directed to improve recovery from larger DNA inputs as well. Lastly,
384 although CRISPR-DS provides an effective solution for small-target region deep sequencing, the method
385 becomes costly for deep sequencing of large genomic regions, an inherent problem of deep sequencing.
386 Nevertheless, fragmentation by CRISPR/Cas9 followed by size selection for fragments as a generic
387 target enrichment technique can easily be scaled to many genomic regions as each region only requires
388 the addition of the appropriate gRNAs for target excision. Thus, CRISPR-DS is ideal for small to moderate
389 size panels (1-100Kb) that require ultra-sensitive mutation detection with minimal DNA inputs.

390

391 **CONCLUSION**

392 We have demonstrated that CRISPR/Cas9 fragmentation followed by size selection enables efficient
393 target enrichment, increasing the recovery of hybridization capture and eliminating the need for a second
394 round of capture for small target regions. In addition, it eliminates PCR bias, maximizes the use of
395 sequencing resources, and produces homogeneous coverage. This fragmentation method can be
396 applied to multiple sequencing modalities that suffer from these problems. Here we have applied it to DS
397 in order to produce CRISPR-DS, an efficient, highly accurate sequencing method with significantly
398 reduced input DNA requirements. CRISPR-DS has broad application for the sensitive identification of
399 mutations in situations in which samples are DNA-limited, such as forensics and early cancer detection.

400

401 **METHODS**

402

403 **Samples.** The samples analyzed included de-identified human genomic DNA from peripheral blood,
404 bladder with and without cancer, and peritoneal fluid DNA from a prior study [17]. Only peritoneal fluid
405 samples had patient information available, which was necessary to confirm the tumor mutation. The
406 remainder of the study samples were used solely to illustrate technical aspects of the technology, no
407 patient information was available, and interpretation of the mutational status of *TP53* is not reported.

408 Frozen bladder samples were obtained from unfixed or frozen autopsy tissue. DNA was extracted with
409 the QIAamp DNA Mini kit (Qiagen, Inc., Valencia, CA, USA) and it had never been denatured, which is
410 essential to preserve the double-stranded nature of each DNA molecule prior to ligation of DS adapters.
411 DNA was quantified with a Qubit HS dsDNA kit (ThermoFisher Scientific). DNA quality was assessed
412 with Genomic TapeStation (Agilent, Santa Clara, CA) and DNA integrity numbers (DIN) were recorded.
413 Peripheral blood DNA and peritoneal fluid DNA had $DIN > 7$ reflecting good quality DNA with no
414 degradation. Bladder samples, however, were purposely selected to include different levels of DNA
415 degradation. Samples B1 to B13 had DINs between 6.8 and 8.9 and were successfully analyzed by
416 CRISPR-DS (Additional file 2: Table S3). Samples B14 and B16 had DINs of 6 and 4, respectively, and
417 were used to demonstrate pre-enrichment of high molecular weight DNA with the BluePippin system (see
418 below and Additional file 1: Figure S6).

419

420 **CRISPR guide design.** CRISPR/Cas9 uses a gRNA to identify the site of cleavage. gRNAs are
421 composed of a complex of CRISPR RNA (crRNA), which contains the ~20bp unique sequence
422 responsible for target recognition, and a trans-activating crRNA (tracrRNA), which has a universal
423 sequence [31]. To select the best gRNAs to excise *TP53* exons we used the CRISPR MIT design website
424 (<http://CRISPR.mit.edu>). The selection criteria were: (1) production of fragments of ~500bp covering
425 exons 2-11 of *TP53* and (2) highest MIT website score (Additional file 2: Table S1 and Additional file 3:
426 Data S1). For exon 7, a smaller size fragment was required in order to avoid a proximal poly-T repeat
427 (Additional file 1, Figure S3) . We designed a total of 12 gRNA, which excised *TP53* into 7 different
428 fragments (Figure 1a). All gRNA had scores > 60 . 10 gRNAs were successful with the first chosen
429 sequence and 2 had to be redesigned due to poor cutting. Initially, the quality of the cut was assessed
430 by reviewing the alignment of the final DCS reads with Integrative Genomics Viewer [32]. Successful
431 guides produced a typical coverage pattern with sharp edges in region boundaries and proper DCS depth
432 (Figure 3d). Unsuccessful guides led to a drop in DCS depth and the presence of long reads that spanned
433 beyond the expected cutting point. In order to simplify and speed up the assessment of guides, we

434 designed a synthetic GeneBlock DNA fragment (IDT, Coralville, IA) that included all gRNA sequences
435 interspaced with random DNA sequences (Additional file 3: Data S2). 3ng of GeneBlock DNA were
436 digested with each of the gRNAs using the CRISPR/Cas9 *in vitro* digestion protocol described below.
437 After digestion, the reactions were analyzed by TapeStation 4200 (Agilent Technologies, Santa Clara,
438 CA, USA) (Additional file 1: Figure S9). The presence of predefined fragment lengths confirms: (1) Proper
439 gRNA assembly (2) The ability of the gRNA to cleave the designed site.

440

441 **CRISPR/Cas9 *in vitro* digestion of genomic DNA.** The *in vitro* digestion of genomic DNA with *S.*
442 *pyogenes* Cas9 Nuclease requires the formation of a ribonucleoprotein complex, which both recognizes
443 and cleaves a pre-determined site. This complex is formed with gRNAs (crRNA + tracrRNA) and Cas9.
444 For multiplex cutting, the gRNAs can be complexed by pooling all the crRNAs, then complexing with
445 tracrRNA, or by complexing each crRNA and tracrRNA separately, then pooling. The second option is
446 preferred because it eliminates competition between crRNAs. gRNAs are at risk of quick degradation and
447 repeated cycles of freeze-thawing should be avoided. crRNAs and tracrRNAs (IDT, Coralville, IA) were
448 complexed into gRNAs and then 30nM of gRNAs were incubated with Cas9 nuclease (NEB, Ipswich,
449 MA) at ~30nM, 1x NEB Cas9 reaction buffer, and water in a volume of 23-27 μ L at 25°C for 10 min. Then,
450 10-250ng of DNA was added for a final volume of 30 μ L. The reaction was incubated overnight at 37°C
451 and then heat shocked at 70°C for 10 min to inactivate the enzyme.

452

453 **Size Selection.** Size selection for the predetermined fragment length is critical for target enrichment prior
454 to library preparation. AMPure XP Beads (Beckman Coulter, Brea, CA, USA) were used to remove off-
455 target, un-digested high molecular weight DNA. After heat inactivation, the reaction was combined with
456 a 0.5x ratio of beads, briefly mixed, and then incubated for 3 min to allow the high molecular weight DNA
457 to bind. The beads were then separated from the solution with a magnet and the solution containing the
458 targeted DNA fragment length was transferred into a new tube. This was followed by a standard AMPure

459 1.8x ratio bead purification eluted into 50 μ L of TE Low to exchange the buffer and remove small DNA
460 contaminants.

461

462 **A-tailing, and ligation.** The fragmented DNA was A-tailed and ligated using the NEBNext Ultra II DNA
463 Library Prep Kit (NEB, Ipswich, MA) according to manufacturer's protocol. The NEB end-repair and A-
464 tailing (ERAT) reaction was incubated at 20°C for 30 min and 65°C for 30 min. Note that end-repair is not
465 needed for CRISPR-DS because Cas9 produces blunt ends, but the ERAT reaction was used for
466 convenient A-tailing. The NEB ligation mastermix and 2.5 μ l of DS adapters at 15 μ M were added and
467 incubated at 20°C for 15 min according to the manufacturer's instructions. Instead of relying on in-house
468 manufactured adapters using previously published protocols [13, 14], which tend to exhibit substantial
469 batch-to-batch variability, we used a commercial adapter prototype of the structure shown in Fig. 1c that
470 were synthesized externally through arrangement with TwinStrand Biosciences. The two differences from
471 the previous adapters are: (1) 10bp random double stranded molecular tag instead of 12bp and (2)
472 substitution of the previous 3' 5bp conserved sequence by a simple 3'-dT overhang to ligate onto the 5'-
473 dA-tailed DNA molecules. Upon ligation, the DNA was cleaned by a 0.8X ratio AMPure Bead purification
474 and eluted into 23 μ L of nuclease free water.

475

476 **PCR.** The ligated DNA was amplified using KAPA Real-Time Amplification kit with fluorescent standards
477 (KAPA Biosystems, Woburn, MA, USA). 50 μ l reactions were prepared including KAPA HiFi HotStart
478 Real-time PCR Master Mix, 23 μ l of previously ligated and purified DNA and DS primers MWS13, 5'-
479 AATGATACGGCGACCACCGAG-3', and MWS20, 5'- GTGACTGGAGTTCAGACGTGTGC-3' [13, 14] at
480 a final concentration of 2 μ M. The reactions were denatured at 98°C for 45 sec and amplified with 6-8
481 cycles of 98°C for 15 sec, 65°C for 30 sec, and 72°C for 30 sec, followed by final extension at 72°C for 1
482 min. Samples were amplified until they reached Fluorescent Standard 3, which typically takes 6-8 cycles
483 depending on the amount of DNA input. Reaching Fluorescent Standard 3 produces a sufficient and

484 standardized number of DNA copies into capture across samples and prevents over-amplification. A 0.8X
485 ratio AMPure Bead wash was performed to purify the amplified fragment and eluted into 40µL of nuclease
486 free water.

487

488 **Capture and post-capture PCR.** *TP53* xGen Lockdown Probes (IDT, Coralville, IA) were used to
489 perform hybridization capture for *TP53* exons as previously reported with minor modifications. From the
490 pre-designed IDT *TP53* Lockdown probes, we selected 21 probes that cover the entire *TP53* coding
491 region (exon 1 and part of exon 11 are not coding) (Additional file 2: Table S2). Each CRISPR/Cas9
492 excised fragment was covered by at least 2 probes and a maximum of 5 probes (Additional file 3: Data
493 S1). To produce the capture probe pool, each of the probes for a given fragment was pooled in equimolar
494 amounts, producing 7 different pools, one for each fragment. The pools were mixed again in equimolar
495 amounts, except for the pools for exon 7 and exons 8-9, which were represented at 40% and 90%
496 respectively. The decrease of capture probes for those exons was implemented after observing
497 consistent overrepresentation of these exons at sequencing. The final capture pool was diluted to 0.75
498 pmol/µl. Of note, it is essential to dilute the capture pool in low TE (0.1 mM EDTA) and to aliquot it in
499 small volumes suitable for 2-3 uses. Excessive rounds of freeze-thaw severely impact the efficiency of
500 the protocol. Hybridization capture was performed according to the IDT protocol, except for 3
501 modifications. First, we used blockers MWS60, 5'-
502 AATGATACGGCGACCACCGAGATCTACACTCTTCCCTACACGACGCTCTTCCGATCTIIIIIIIIITGA
503 CT-3' and MSW61, 5'-GTCAIIIIIIIIIIAGATCGGAAGAGCACACGTCTGAACTCCAGTCAC-3', which are
504 specific to DS adapters. Second, we used 75µl of Dynabeads M-270 Streptavidin beads instead of 100µl.
505 Third, the post-capture PCR was performed with the KAPA Hi-Fi HotStart PCR kit (KAPA Biosystems,
506 Woburn, MA, USA) using MWS13 and indexed primer MWS21 at a final concentration of 0.8 µM. The
507 reaction was denatured at 98°C for 45 sec and then amplified for 20 cycles at 98°C for 30 sec, 60°C for

508 45 sec, and 72°C for 45 sec, followed by extension at 72°C for 60 sec. The PCR product was purified
509 with a 0.8X AMPure Bead wash.

510

511 **Sequencing.** Samples were quantified using the Qubit dsDNA HS Assay Kit, diluted, and pooled for
512 sequencing. The sample pool was visualized on the Agilent 4200 TapeStation to confirm library quality.
513 The TapeStation electropherogram should show sharp, distinct peaks corresponding to the fragment
514 length of the designed CRISPR/Cas9 cut fragments (Fig. 3b-c). This step can also be performed for each
515 sample individually, prior to pooling, to verify the performance of each individual sample. The final pool
516 was quantified using the KAPA Library Quantification kit (KAPA Biosystems, Woburn, MA, USA). The
517 library was sequenced on the MiSeq Illumina platform using a v3 600 cycle kit (Illumina, San Diego, CA,
518 USA) as specified by the manufacturer. For each sample, we allocated ~7-10% of a lane corresponding
519 to ~2 million reads. Each sequencing run was spiked with approximately 1% PhiX control DNA.

520

521 **Standard-DS experiments.** Three amounts of DNA (25ng, 100ng, and 250ng) from normal human
522 bladder sample B9 were sequenced with standard-DS with one round and two rounds of capture to
523 provide direct comparison with CRISPR-DS. Standard-DS was performed as previously described [14],
524 with the exception that the KAPA Hyperprep kit (KAPA Biosystems, Woburn, MA, USA) was used for
525 end-repair and ligation and the KAPA Hi-Fi HotStart PCR kit (KAPA Biosystems, Woburn, MA, USA) was
526 used for PCR amplification. Hybridization capture was performed with xGen Lockdown probes that
527 covered *TP53* exons 2-11, the same that were used for CRISPR-DS. Samples were sequenced on ~10%
528 of a HiSeq 2500 Illumina platform to accommodate shorter fragment lengths. Data analysis was perform
529 with the standard-DS analysis pipeline (<https://github.com/risqueslab/DuplexSequencingScripts>).

530

531 **CRISPR-DS target enrichment experiments.** Two different experiments were performed to
532 characterize CRISPR-DS target enrichment. The first experiment consisted of comparing one vs. two
533 rounds of capture. Three DNA samples were processed for CRISPR-DS and split in half after one

534 hybridization capture. The first half was indexed and sequenced and the second half was subject to an
535 additional round of capture, as required in the original DS protocol. Then the percentage of raw reads
536 on-target (covering *TP53* exons) was compared for one vs. two captures. The second experiment
537 assessed the percentage of raw reads on-target without performing hybridization capture to determine
538 the enrichment produced exclusively by size selecting CRISPR excised fragments. Fold-enrichment was
539 calculated as the fraction of on-target raw reads divided over the expected fraction of on-target reads
540 given the size of the target region (bases in the target region/total genome bases). Different DNA amounts
541 (from 10ng to 250ng) of three different samples were processed with the protocol described above until
542 first PCR, that is, prior to hybridization capture. Then the PCR product was indexed and sequenced. The
543 percentage of raw reads on-target was calculated and the fold enrichment was estimated considering the
544 size of the targeted region, which is 3,280bp.

545

546 **Pre-enrichment for high molecular weight DNA.** Selection of high molecular weight DNA improves the
547 performance of degraded DNA in CRISPR-DS. We performed this selection using a BluePippin system
548 (Sage Science, Beverly, MA). Two bladder DNAs with DINs of 6 and 4 were run using a 0.75% gel
549 cassette and high-pass setting to obtain >8kb fragments. Size selection was confirmed by TapeStation
550 (Additional file 1: Figure S6a). Then 250ng of DNA before BluePippin and 250ng of DNA after BluePippin
551 were processed in parallel with CRISPR-DS. The percentage of raw reads on-target as well as average
552 DCS depth was quantified and compared (Additional file 1: Figure S6b.). Alternative methods for size
553 selection such as AMPure beads might be suitable to perform this enrichment.

554

555 **Data processing.** A custom bioinformatics pipeline was created to automate analysis from raw FASTQ
556 files to text files (Additional file 1: Figure S8). This pipeline includes two major modifications compared to
557 the previously described method for DS analysis [13, 14]: (1) the retention of paired read information and
558 (2) consensus-making performed prior to alignment. Paired-end reads are essential to the analysis of
559 CRISPR-DS data, but are also an important improvement for the analysis of DS in general, as they allow

560 critical quality control of fragment size and removal of potential technical artifacts related to short
561 fragments. In this pipeline, consensus is executed by a custom python and bash scripts. After consensus
562 calling, the resulting processed FASTQ files are aligned to the reference genome of interest, in this case
563 human reference genome v38, using bwa-mem v.0.7.4[33] with default parameters. Mapped reads are
564 re-aligned with GATK Indel-Realigner and low quality bases are clipped from the ends with GATK Clip-
565 Reads (<https://software.broadinstitute.org/gatk/>). Because of the expected decrease in read quality in the
566 latest cycles of sequencing, we performed a conservative clipping of 30 bases from the 3' end and
567 another 7 bases from 5' end were clipped to avoid the occasional extra overhang left by incorrectly
568 synthesized adapters. In addition, overlapping areas of read-pairs, which in our *TP53* design spanned
569 ~80bp, are trimmed back using fgbio ClipOverlappingReads (<https://github.com/fulcrumgenomics/fgbio>).
570 Software for CRISPR-DS is available at <https://github.com/risqueslab/CRISPR-DS>.

571

572 **Data analysis.** Recovery rate (also called fractional genome-equivalent recovery) was calculated as
573 average DCS depth (sequenced genomes) divided by number of input genomes (1ng of human genomic
574 DNA corresponds to ~330 haploid genomes). The number of on-target raw reads was calculated by
575 counting the number of reads within 100bp window on either side of the CRISPR/Cas9 cut sites. Optimal
576 fragment size (Fig. 4b-c and Additional file 1: Figure S4) was calculated as the sequencing read length
577 minus the barcode sequence and minus clipped off bases for poor quality at the ends of reads. For
578 peritoneal fluid samples sequenced with both CRISPR-DS and standard-DS, *TP53* biological background
579 mutation frequency was calculated as the number of *TP53* mutations in *TP53* exons 4 to 10 (excluding
580 the tumor mutation) divided by the total number of nucleotides sequenced in those exons. The 95%
581 confidence intervals were calculated in R using the Clopper-Pearson 'exact' method for binomial
582 distributions.

583

584

585

586 **ABBREVIATIONS**

587 **DS:** Duplex Sequencing **DCS:** Double-stranded consensus sequence **SSCS:** Single-stranded consensus
588 sequence **gRNA:** Guide RNA **crRNA:** CRISPR RNA **tracrRNA:** Trans-activating crRNA **NGS:** Next-
589 generation Sequencing **ng:** Nanogram **bp:** Basepair **ssDNA:** Single-stranded DNA **dsDNA:** Double-
590 stranded DNA **DIN:** DNA integrity number

591

592 **DECLARATIONS**

593

594 **Ethics approval and consent to participate**

595 Samples in these studies were obtained from: (1) the University of Washington Gynecologic Oncology
596 Tissue Bank, which collected specimens and clinical information after informed consent under protocol
597 number 27077 approved by the University of Washington Human Subjects Division institutional review
598 board; (2) the University of Washington Genitourinary Cancer Specimen Biorepository and from not
599 previously fixed or frozen autopsy tissue with waiver of consent under protocol number 52389 approved
600 by the Fred Hutchinson Cancer Research Center Human Subjects Division institutional review board.

601 **Consent for publication**

602 All the samples in the study were de-identified. Consent for publication is included under the informed
603 consent for research described above.

604 **Availability of data and material**

605 Sequencing data that supports the findings of this study have been deposited in the Sequence Read
606 Archive (BioProject ID: [PRJNA412416](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA412416)). Software for CRISPR-DS data analysis is available at
607 <https://github.com/risqueslab/CRISPR-DS>.

608 **Competing interests**

609 SRK is a consultant and equity holder for TwinStrand Biosciences Inc. JJS is a founder and equity holder
610 in TwinStrand Biosciences Inc. RAR is the principal investigator on a NIH SBIR R44CA221426
611 subcontract research agreement with TwinStrand Biosciences Inc.

612 **Funding**

613 Research reported in this publication was supported by grants from the NIH under award numbers
614 R01CA160674 and R01CA181308 to RAR; Mary Kay Foundation grant 045-15 to RAR. Cooperative
615 Agreement Number W911NF-15-2-0127 from the Department of Defense Army Research Office/Defense
616 Forensic Science Center(DFSC), as well as grant W81XWH-16-1-0579 from the Department of Defense
617 Congressionally Directed Medical Research Program to SRK. The views and conclusions contained in
618 this document are those of the authors and should not be interpreted as representing the official policies,
619 either expressed or implied, of the Army Research Office, DFSC, or the U.S. Government. The U.S.
620 Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding
621 any copyright notation hereon.

622 **Authors' contributions**

623 S.R.K. conceived the idea; D.N. S.R.K., and R.A.R. developed the method; D.N. and R.A.R. designed
624 the experiments; D.N., S.L., E.K.S., M.J.H., K.B., K.L.S., B.F.K., R.A.R, and S.R.K. carried out
625 experiments and/or performed data analysis; M.T. provided samples and scientific input; Y.Z., and J.S.
626 contributed to assay development and provided invaluable critical discussion; D.N. S.R.K. and R.A.R.
627 wrote the paper.

628 **Acknowledgements**

629 We thank Shilpa Kumar for assistance with computational analysis, Emily Kohlbrenner for technical
630 support and helpful discussions, Penny Faires for critical reading and copy editing of the manuscript, and
631 the Genitourinary Cancer Specimen Biorepository for providing access to bladder cases (Director Dr
632 Colm Morrissey, PhD), We thank the University of Washington Gynecologic Oncology Tissue Bank for
633 providing peritoneal fluid DNA and the Brigham and Women's Hospital/Harvard Cohorts Biorepository for
634 sending archived samples from the Nurses' Health Study for pilot testing.

635

636

637

638 **TABLES**

639

Table 1: Target enrichment due to size selection

Sample	DNA Input (ng)	Reads On Target (%)	Fold Enrichment
B9	25	0.76%	7,527
	200	0.25%	2,452
	250	0.21%	2,037
PF1	10	2.85%	28,139
	25	1.99%	19,583
	100	0.68%	6,667
	250	0.70%	6,878
PF5	10	5.05%	49,794
	25	0.96%	9,456
	100	0.34%	3,321
	250	0.22%	2,217

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

Table 2. Comparison of Standard-DS vs CRISPR-DS for four different samples with *TP53* mutations.

Method	Sample	Input DNA (ng)	Raw Reads On Target	Median Final Depth*	Recovery (%)	Tumor Mutation	Mutant Allele Fraction
Standard-DS	PF1	9,196	92.4%	2742	0.09%	chr17:g.7578275G>A	68.5%
	PF2	3,000	92.8%	5381	0.54%	chr17:g.7577548C>T	1.2%
	PF3	10,186	95.9%	1866	0.06%	chr17:g.7578403C>T	1.6%
	PF4	7,436	95.4%	2029	0.08%	chr17:g.7578526C>T	0.6%
CRISPR-DS	PF1	100	76.6%	2039	6.18%	chr17:g.7578275G>A	68.4%
	PF2	100	94.3%	2831	8.58%	chr17:g.7577548C>T	1.0%
	PF3	100	87.6%	3801	11.52%	chr17:g.7578403C>T	0.4%
	PF4	100	96.5%	2194	6.65%	chr17:g.7578526C>T	0.1%

*After final Duplex Sequencing data processing is performed

657 **FIGURES**

658

659 Figure 1. Schematic representation of key aspects of CRISPR-DS. (a) CRISPR/Cas9 digestion of *TP53*.
660 Seven fragments containing all *TP53* coding exons were excised via targeted cutting using gRNAs. Dark
661 grey represents reference strand and light grey represents the anti-reference strand. (b) Size selection
662 using 0.5x SPRI beads. Uncut, genomic DNA binds to the beads and allows the recovery of the
663 homogenously sized excised fragments in solution. (c) Double-stranded DNA molecule fragmented and
664 ligated with DS-adapters. Adapters are double-stranded and contain 10-bp of random, complementary
665 nucleotides and a 3'-dT overhang. (d) Error correction by DS. Reads derived from the same strand of
666 DNA are compared to form a Single-Strand Consensus Sequence (SSCS). Then both strands of the
667 same original DNA molecule are compared with one another to create a Double-Strand Consensus
668 Sequence (DCS). Only mutations found in both SSCS reads are counted as true mutations in DCS reads.

669

670 Figure 2. Comparison of library preparation protocols for standard-DS vs. CRISPR-DS. The primary
671 differences between the CRISPR-DS and standard-DS library preparation are the fragmentation
672 techniques and the number of hybridization capture steps. Instead of fragmentation by sonication as
673 performed in standard-DS, CRISPR-DS relies on an *in vitro* excision of target regions by CRISPR/Cas9
674 followed by size selection for the excised fragments. The size selection eliminates the need for a second
675 round of hybridization capture which is required for sufficient target enrichment in the standard-DS
676 protocol. CRISPR-DS reduces the workflow by nearly a day. Colored boxes represent 1h of time.

677

678 Figure 3. Visualization of sequencing libraries and data prepared with CRISPR-DS and standard-DS. (a)
679 TapeStation gels show distinct bands for CRISPR-DS as opposed to a smear for standard-DS. The size
680 of bands corresponds to the CRISPR/Cas9 cut fragments with adapters. (b) CRISPR-DS
681 electropherograms allow visualization and quantification of peaks for quality control of the library prior to
682 sequencing. Standard-DS electropherograms show a diffuse peak that harbors no information about the

683 specificity of the library. (c) Dots represent original barcoded DNA molecules. Each DNA molecule has
684 multiple copies generated at PCR (x-axis). In CRISPR-DS, all DNA molecules (red dots) have preset
685 sizes (y-axis) and generate similar number of PCR copies. In standard-DS, sonication shears DNA into
686 variable fragment lengths (blue dots). Smaller fragments amplify better and generate an excess of copies
687 that waste sequencing resources. (d) Integrative Genomics Viewer of *TP53* coverage with DCS reads
688 generated by CRISPR-DS and standard-DS. CRISPR-DS shows distinct boundaries that correspond to
689 the CRISPR/Cas9 cutting points and an even distribution of depth across positions, both within a fragment
690 and between fragments. Standard-DS shows the typical 'peak' pattern generated by random shearing of
691 fragments and hybridization capture, which leads to variable coverage.

692

693 Figure 4. CRISPR/Cas9 fragmentation produces optimal fragment lengths. (a) Sonication produces
694 fragments that are either too short or too long, corresponding to redundant or lost information,
695 respectively. CRISPR-DS produces optimally sized fragments which are perfectly covered by the
696 sequencing reads. (b-c) Comparison of histograms of the insert sizes of two samples prepared with
697 standard-DS (*blue*, left panels), which uses sonication for fragmentation, and CRISPR-DS (*red*, right
698 panels), which uses CRISPR/Cas9 digestion for fragmentation, The x-axis represents the percent
699 difference from the optimally sized fragment, e.g. fragment size that matches the sequencing read length
700 after adjustments for molecular barcodes and clipping. Yellow shading highlights range of fragment sizes
701 which are within 10% difference from optimal size.

702

703 Figure 5. Technical comparison of 250ng, 100ng and 25ng of DNA sequenced with both standard-DS
704 and CRISPR-DS. Measurements were obtained by sequencing samples prepared with standard-DS
705 (*blue*) using one and two rounds of hybridization capture and CRISPR-DS (*red*) with only one round of
706 hybridization capture. (a) The percentage of raw sequencing reads on-target (covering *TP53*) was
707 comparable between Standard-DS with two rounds of capture and CRISPR-DS with one round of
708 capture, demonstrating the target enrichment efficiency of the novel method. (b) Percentage recovery

709 was calculated as the percentage of genomes in input DNA that produced DCS reads. CRISPR-DS
710 increases recovery thanks to the initial CRISPR-based target enrichment, which eliminates one round of
711 hybridization capture. (c) After creating DCS reads, the median DCS depth across all targeted regions
712 was calculated for each input amount. The increased recovery enabled by CRISPR-DS translates into 5-
713 10 times more sequencing depth for the same input DNA.

714

715

716 **ADDITIONAL FILE 1:**

717 **SUPPLEMENTARY FIGURES**

718 Figure S1. Comparison of mutation limit detection by sequencing accuracy

719 Figure S2. Timeline of library preparation for CRISPR-DS and standard-DS

720 Figure S3. Homopolymer region produces suboptimal sequencing near *TP53* exon 7

721 Figure S4. Fraction of reads within 10% of optimal insert size: CRISPR-DS vs standard-DS

722 Figure S5. Target enrichment for CRISPR-DS with one vs. two captures

723 Figure S6. Pre-enrichment for high molecular weight DNA with BluePippin

724 Figure S7. Comparison of *TP53* biological background mutation frequency measured by Standard-DS
725 and CRISPR-DS

726 Figure S8. Overview of CRISPR-DS data processing

727 Figure S9. Control CRISPR/Cas9 digestion of *TP53* gRNAs

728

729 **ADDITIONAL FILE 2:**

730 **SUPPLEMENTARY TABLES**

731 Table S1. crRNA sequences for *TP53* CRISPR/Cas9 digestion

732 Table S2. *TP53* hybridization capture probes

733 Table S3. CRISPR-DS sequencing results for 15 samples processed with 250ng input DNA

734

735 **ADDITIONAL FILE 3:**

736 **SUPPLEMENTARY DATA**

737 Data S1. *TP53* sequence with crRNA and capture probes

738 Data S2. GeneBlock sequence

739

740 **REFERENCES**

- 741 1. Kebschull JM, Zador AM: Sources of PCR-induced distortions in high-throughput sequencing data
742 sets. *Nucleic Acids Res* 2015, 43:e143.
- 743 2. Samorodnitsky E, Jewell BM, Hagopian R, Miya J, Wing MR, Lyon E, Damodaran S, Bhatt D,
744 Reeser JW, Datta J, Roychowdhury S: Evaluation of Hybridization Capture Versus Amplicon-
745 Based Methods for Whole-Exome Sequencing. *Hum Mutat* 2015, 36:903-914.
- 746 3. Park G, Park JK, Shin SH, Jeon HJ, Kim NKD, Kim YJ, Shin HT, Lee E, Lee KH, Son DS, et al:
747 Characterization of background noise in capture-based targeted sequencing data. *Genome Biol*
748 2017, 18:136.
- 749 4. Knierim E, Lucke B, Schwarz JM, Schuelke M, Seelow D: Systematic comparison of three
750 methods for fragmentation of long-range PCR products for next generation sequencing. *PLoS*
751 *One* 2011, 6:e28240.
- 752 5. Goodwin S, McPherson JD, McCombie WR: Coming of age: ten years of next-generation
753 sequencing technologies. *Nat Rev Genet* 2016, 17:333-351.
- 754 6. Fox EJ, Reid-Bayliss KS, Emond MJ, Loeb LA: Accuracy of Next Generation Sequencing
755 Platforms. *Next Gener Seq Appl* 2014, 1.
- 756 7. Arbeithuber B, Makova KD, Tiemann-Boege I: Artfactual mutations resulting from DNA lesions
757 limit detection levels in ultrasensitive sequencing applications. *DNA Res* 2016, 23:547-559.
- 758 8. Chen L, Liu P, Evans TC, Jr., Ettwiller LM: DNA damage is a pervasive cause of sequencing
759 errors, directly confounding variant identification. *Science* 2017, 355:752-756.

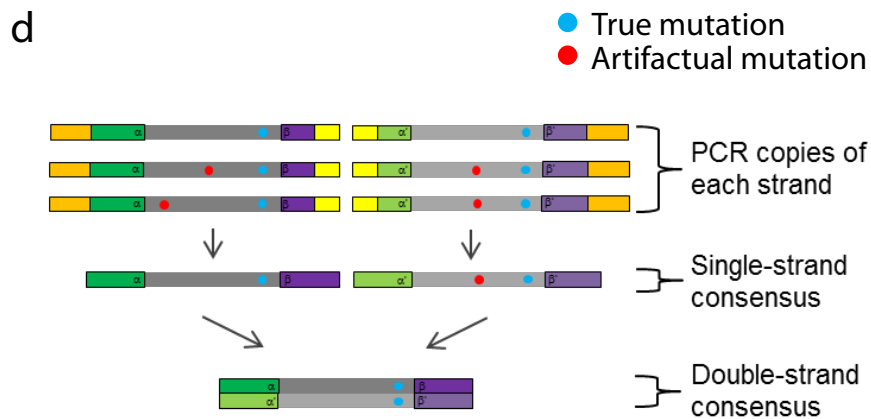
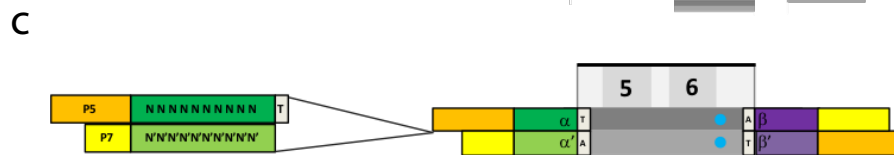
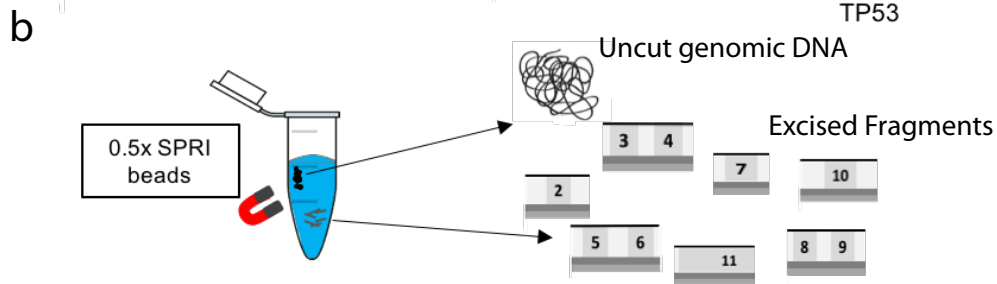
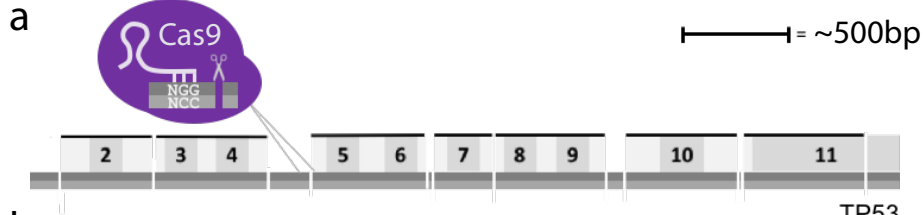
- 760 9. Akogwu I, Wang N, Zhang C, Gong P: A comparative study of k-spectrum-based error correction
761 methods for next-generation sequencing data analysis. *Hum Genomics* 2016, 10 Suppl 2:20.
- 762 10. Laehnemann D, Borkhardt A, McHardy AC: Denoising DNA deep sequencing data-high-
763 throughput sequencing errors and their correction. *Brief Bioinform* 2016, 17:154-179.
- 764 11. Kinde I, Wu J, Papadopoulos N, Kinzler KW, Vogelstein B: Detection and quantification of rare
765 mutations with massively parallel sequencing. *Proc Natl Acad Sci U S A* 2011, 108:9530-9535.
- 766 12. Hoang ML, Kinde I, Tomasetti C, McMahon KW, Rosenquist TA, Grollman AP, Kinzler KW,
767 Vogelstein B, Papadopoulos N: Genome-wide quantification of rare somatic mutations in normal
768 human tissues using massively parallel sequencing. *Proc Natl Acad Sci U S A* 2016, 113:9846-
769 9851.
- 770 13. Schmitt MW, Kennedy SR, Salk JJ, Fox EJ, Hiatt JB, Loeb LA: Detection of ultra-rare mutations
771 by next-generation sequencing. *Proc Natl Acad Sci U S A* 2012, 109:14508-14513.
- 772 14. Kennedy SR, Schmitt MW, Fox EJ, Kohn BF, Salk JJ, Ahn EH, Prindle MJ, Kuong KJ, Shen JC,
773 Risques RA, Loeb LA: Detecting ultralow-frequency mutations by Duplex Sequencing. *Nat Protoc*
774 2014, 9:2586-2606.
- 775 15. Ahn EH, Lee SH, Kim JY, Chang CC, Loeb LA: Decreased Mitochondrial Mutagenesis during
776 Transformation of Human Breast Stem Cells into Tumorigenic Cells. *Cancer Res* 2016, 76:4569-
777 4578.
- 778 16. Kennedy SR, Salk JJ, Schmitt MW, Loeb LA: Ultra-sensitive sequencing reveals an age-related
779 increase in somatic mitochondrial mutations that are inconsistent with oxidative damage. *PLoS*
780 *Genet* 2013, 9:e1003794.
- 781 17. Krimmel JD, Schmitt MW, Harrell MI, Agnew KJ, Kennedy SR, Emond MJ, Loeb LA, Swisher EM,
782 Risques RA: Ultra-deep sequencing detects ovarian cancer cells in peritoneal fluid and reveals
783 somatic TP53 mutations in noncancerous tissues. *Proc Natl Acad Sci U S A* 2016, 113:6005-
784 6010.

- 785 18. Reid-Bayliss KS, Arron ST, Loeb LA, Bezrookove V, Cleaver JE: Why Cockayne syndrome
786 patients do not get cancer despite their DNA repair deficiency. *Proc Natl Acad Sci U S A* 2016,
787 113:10151-10156.
- 788 19. Winters M, Monroe C, Barta JL, Kemp BM: Are we fishing or catching? Evaluating the efficiency
789 of bait capture of CODIS fragments. *Forensic Science International-Genetics* 2017, 29:61-70.
- 790 20. Schmitt MW, Fox EJ, Prindle MJ, Reid-Bayliss KS, True LD, Radich JP, Loeb LA: Sequencing
791 small genomic targets with high efficiency and extreme accuracy. *Nat Methods* 2015, 12:423-425.
- 792 21. Bennett-Baker PE, Mueller JL: CRISPR-mediated isolation of specific megabase segments of
793 genomic DNA. *Nucleic Acids Res* 2017, 45:e165.
- 794 22. Shin G, Grimes SM, Lee H, Lau BT, Xia LC, Ji HP: CRISPR-Cas9-targeted fragmentation and
795 selective sequencing enable massively parallel microsatellite analysis. *Nat Commun* 2017,
796 8:14291.
- 797 23. Dabney J, Meyer M: Length and GC-biases during sequencing library amplification: a comparison
798 of various polymerase-buffer systems with ancient and modern DNA sequencing libraries.
799 *Biotechniques* 2012, 52:87-94.
- 800 24. Hsu PD, Scott DA, Weinstein JA, Ran FA, Konermann S, Agarwala V, Li Y, Fine EJ, Wu X, Shalem
801 O, et al: DNA targeting specificity of RNA-guided Cas9 nucleases. *Nat Biotechnol* 2013, 31:827-
802 832.
- 803 25. Jung H, Ji S, Song S, Park Y, Yang J, Schmidt E: The DNA Integrity Number (DIN) provided by
804 the genomic DNA ScreenTape assay allows for streamlining of NGS on FFPE tissue samples.
805 *Application Note Nucleic Acid Analysis* 2014
- 806 26. Li H: A statistical framework for SNP calling, mutation discovery, association mapping and
807 population genetical parameter estimation from sequencing data. *Bioinformatics* 2011, 27:2987-
808 2993.
- 809 27. Kozarewa I, Armisen J, Gardner AF, Slatko BE, Hendrickson CL: Overview of Target Enrichment
810 Strategies. *Curr Protoc Mol Biol* 2015, 112:7 21 21-23.

- 811 28. Newman AM, Lovejoy AF, Klass DM, Kurtz DM, Chabon JJ, Scherer F, Stehr H, Liu CL, Bratman
812 SV, Say C, et al: Integrated digital error suppression for improved detection of circulating tumor
813 DNA. *Nat Biotechnol* 2016, 34:547-555.
- 814 29. Taylor PH, Cinquin A, Cinquin O: Quantification of in vivo progenitor mutation accrual with ultra-
815 low error rate and minimal input DNA using SIP-HAVA-seq. *Genome Res* 2016, 26:1600-1611.
- 816 30. Mattox AK, Wang Y, Springer S, Cohen JD, Yegnasubramanian S, Nelson WG, Kinzler KW,
817 Vogelstein B, Papadopoulos N: Bisulfite-converted duplexes for the strand-specific detection and
818 quantification of rare mutations. *Proc Natl Acad Sci U S A* 2017, 114:4733-4738.
- 819 31. Ran FA, Hsu PD, Wright J, Agarwala V, Scott DA, Zhang F: Genome engineering using the
820 CRISPR-Cas9 system. *Nat Protoc* 2013, 8:2281-2308.
- 821 32. Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP:
822 Integrative genomics viewer. *Nat Biotechnol* 2011, 29:24-26.
- 823 33. Li H, Durbin R: Fast and accurate long-read alignment with Burrows-Wheeler transform.
824 *Bioinformatics* 2010, 26:589-595.

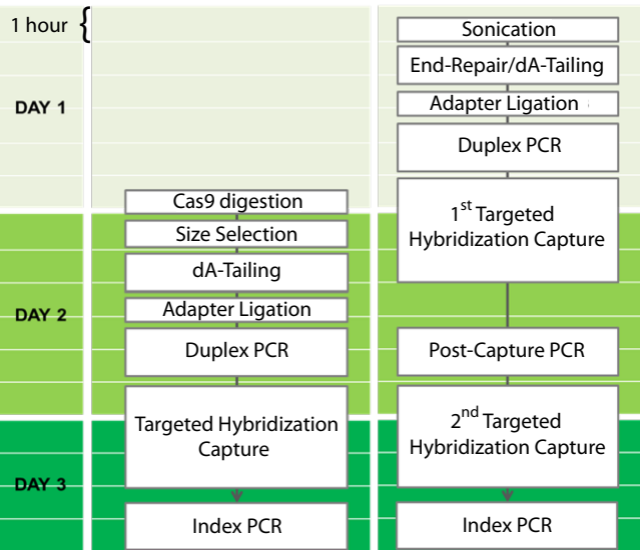
825

826

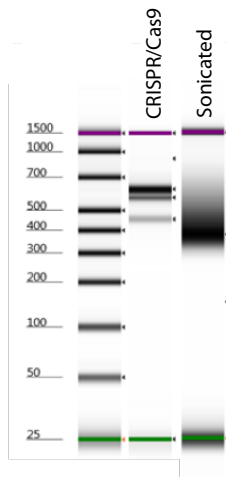


CRISPR-DS

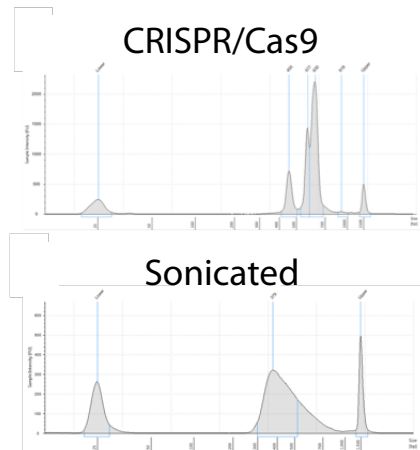
Standard-DS



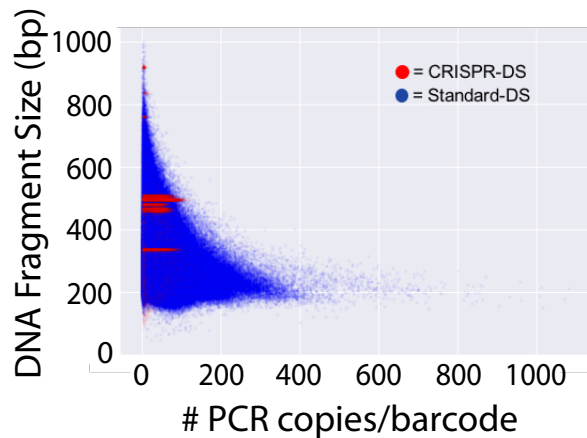
a



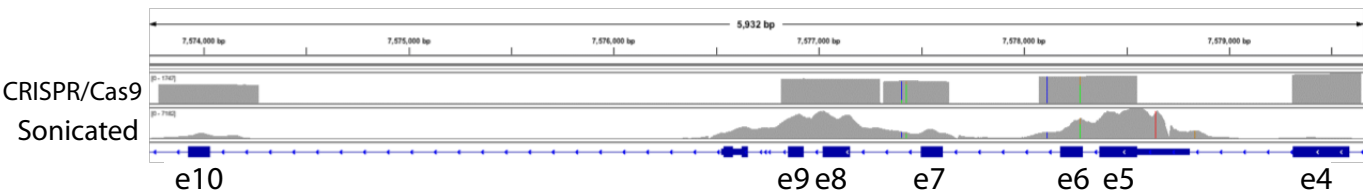
b

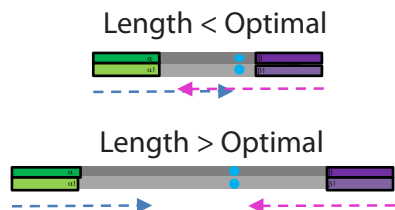


c

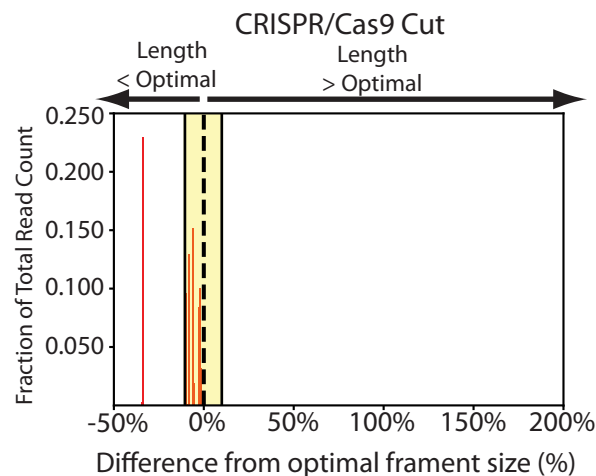
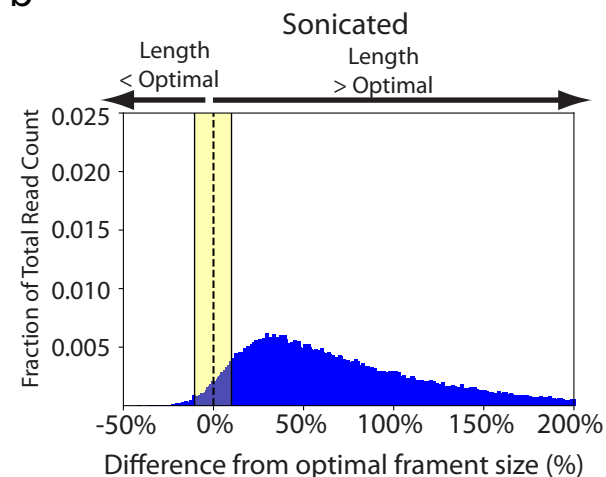


d



a

Length = Optimal

**b****c**