

1 **Repeat Expansions Confers WRN Dependence in Microsatellite Unstable**
2 **Cancers**

3

4 **Niek van Wietmarschen^{1,*}, Sriram Sridharan^{1,*}, William J. Nathan^{1,2*}, Anthony**
5 **Tubbs^{1,*}, Edmond M. Chan^{3,4*}, Elsa Callen¹, Wei Wu¹, Frida Belinky¹, Veenu**
6 **Tripathi¹, Nancy Wong¹, Kyla Foster⁴, Javad Noorbakhsh⁴, Kiran Garimella⁴, Abimael**
7 **Cruz-Migoni², Joshua A. Sommers⁵, Yongqing Huang⁴, Ashir A. Borah⁴, Jonathan T.**
8 **Smith⁴, Jeremie Kalfon⁴, Nikolas G. Kesten³, Kasper Fugger⁶, Robert L. Walker⁷, Egor**
9 **Dolzhenko⁸, Michael A. Eberle⁸, Bruce E. Hayward⁹, Karen Usdin⁹, Catherine H.**
10 **Freudenreich¹⁰, Robert M. Brosh, Jr.⁵, Stephen C. West⁶, Peter J. McHugh², Paul S.**
11 **Meltzer⁷, Adam J. Bass^{3,4} & André Nussenzweig^{1,†}**

12

13 ¹ **Laboratory of Genome Integrity, National Cancer Institute, NIH, Bethesda, MD, USA**

14 ² **Department of Oncology, MRC Weatherall Institute of Molecular Medicine, University**
15 **of Oxford, John Radcliffe Hospital, Oxford, United Kingdom**

16 ³ **Department of Medical Oncology, Dana-Farber Cancer Institute, Harvard Medical**
17 **School, Boston MA, USA**

18 ⁴ **Broad Institute of Harvard and MIT, Cambridge, MA, USA**

19 ⁵ **Laboratory of Molecular Gerontology, National Institute on Aging, NIH, Baltimore,**
20 **Maryland, United States of America**

21 ⁶ **DNA Recombination and Repair Laboratory, The Francis Crick Institute, London,**
22 **U.K.**

23 ⁷ **Genetics Branch, National Cancer Institute, NIH, Bethesda, MD, USA.**

24 ⁸ **Illumina Incorporated, San Diego CA, USA**

25 ⁹ **Laboratory of Cell and Molecular Biology, National Institute of Diabetes, Digestive and**
26 **Kidney Diseases, NIH, Bethesda MD, USA**

27 ¹⁰ **Department of Biology, Tufts University, Medford, MA, USA.**

28

29

30 ***Authors contributed equally to this work**

31

32 **†Correspondence: andre_nussenzweig@nih.gov**

33

34

35

36 **The RecQ DNA helicase WRN is a synthetic lethal target for cancers with microsatellite**
37 **instability (MSI), a form of genetic hypermutability arising from impaired mismatch**
38 **repair¹⁻⁴. WRN depletion induces widespread DNA double strand breaks (DSBs) in MSI**
39 **cells, leading to cell cycle arrest and/or apoptosis. However, the mechanism by which**
40 **WRN protects MSI cancers from DSBs remains unclear. Here, we demonstrate that TA-**
41 **dinucleotide repeats are highly unstable in MSI cells and exhibit surprisingly large-scale**
42 **expansions, distinct from previously described insertion/deletion mutations of a few**
43 **nucleotides⁵. We show that expanded TA repeats form non-B DNA secondary structures**
44 **that stall replication forks, activate the ATR checkpoint kinase, and necessitate**
45 **unwinding by the WRN helicase. In the absence of WRN, the expanded TA-dinucleotide**
46 **repeats are susceptible to MUS81 nuclease cleavage, leading to massive chromosome**
47 **shattering. Thus, our study uncovers a distinct biomarker within MSI tumors that**
48 **underlies the synthetic lethal dependence on WRN, thereby supporting the development**
49 **of WRN-based therapeutics.**

50

51

52

53 MSI is characterized by hypermutability of short repetitive DNA sequences scattered
54 throughout the genome. MSI arises from deficiency in DNA mismatch repair (MMR) and
55 contributes to the formation of many types of cancers including colorectal cancer (15%),
56 endometrial (20-30%), gastric cancers (15%), and ovarian cancers (12%)⁶. Recent studies
57 demonstrate that multiple cancer types with MSI are reliant on the WRN helicase activity for
58 survival¹⁻⁴. WRN is a member of the RecQ family of DNA helicases including WRN, BLM,
59 and RECQL4, which when mutated manifest as the distinct chromosome instability disorders
60 Werner syndrome, Bloom syndrome, and Rothmund-Thomson syndrome, respectively⁷. RecQ
61 helicases do not display sequence specificity but resolve non-canonical secondary DNA
62 structures such as bubbles, Holliday junctions, and G-quadruplexes that may be encountered
63 during replication and recombination. However, the mechanism by which WRN helicase is
64 required to protect chromosomal integrity of MSI, but not microsatellite stable (MSS),
65 cancers is not understood.

66 We previously demonstrated that MMR restoration only partially rescued MSI cells
67 from WRN depletion². Recent studies have reported discrepant results on the effect of acute
68 *MLH1* silencing to sensitize MSS cells to WRN depletion^{1,3}. We therefore evaluated *WRN*
69 dependency in human primary stomach epithelial cells following knockout of the MMR genes
70 *MLH1* or *MSH2*. After 4 months of culture, *MLH1* or *MSH2* knockout cells failed to develop
71 a dependency upon *WRN* for survival (Extended Data Fig. 1a, b). These data suggest that
72 instead of WRN loss being simply synthetic lethal with impaired MMR, a “genomic scar”
73 may gradually accumulate in MSI cancers that requires WRN’s role as a structure-specific
74 helicase.

75

76 **WRN deficiency generates recurrent DSBs and extensive end-resection in MSI cells**

77 Loss of viability in MSI cells upon *WRN* silencing is associated with a decrease in
78 proliferation and accumulation of DNA double strand breaks (DSBs)¹⁻⁴. Consistently, we
79 found that *WRN* depletion using a doxycycline-inducible *WRN* shRNA expressed in the MSI
80 KM12 cell line² resulted in decreased DNA synthesis and high levels of the DSB marker,
81 KAP1 phosphorylation at Ser 824 (pKAP1)⁸, predominantly in the G2/M phase of the cell
82 cycle (Extended Data Fig. 1c-e). In contrast, expression of a seed control *WRN.C911* shRNA²
83 did not substantially induce DSBs in KM12 (Extended Data Fig. 1f). Analysis of mitotic
84 spreads (n=100) revealed that all chromosomes were shattered in approximately 35% of
85 *WRN*-depleted KM12 cells (Fig. 1a and Fig. 2a). In contrast, chromosome shattering was not
86 evident in microsatellite stable (MSS) colon cancer line SW837 upon *WRN* depletion (Fig. 1a
87 and Extended Data Fig. 1e, n=100 metaphases).

88 To determine whether *WRN* is necessary to unwind specific regions of the genome to
89 prevent DNA breakage, we queried recurrent sites of DSBs by END-seq⁹. MSI cell lines
90 KM12 and HCT116 show little endogenous DNA breakage; however, upon *WRN* depletion
91 by either shRNA or siRNA, recurrent DSBs were detected at specific locations throughout the
92 genome (Fig. 1b, Extended Data Fig. 2a). Moreover, END-seq peak intensities were highly
93 reproducible among different experiments (Extended Data Fig. 2b). Using the criteria of 20-
94 fold enrichment over non-treated cells (NT; *WRN* proficient), END-seq peaks overlapped
95 significantly between HCT116 and KM12 (Extended Data Fig. 2a,c). In contrast, *WRN*
96 depletion in SW837 did not substantially induce DSBs (Extended Data Fig. 2d). Thus, loss of
97 *WRN* induces DSBs at recurrent genomic loci that are reproducible across distinct MSI
98 cancer cell lines.

99

100 **DNA breaks accumulate around TA dinucleotide repeats**

101 END-seq peaks displayed the characteristic pattern of positive and negative strand reads
102 representing the right and left end of DSBs, respectively (Fig. 1c). END-seq reads spread
103 outwards from DSB sites in a pattern consistent with DNA end resection^{9,10}. To confirm this,
104 we mapped sites of single-strand DNA (ssDNA) bound by Replication Protein A (RPA)^{11,12}.
105 Strikingly, 79% of ssDNA peaks overlapped with END-seq peaks, with resection lengths
106 averaging 500 bp and extending up to 5 kb (Fig. 1b-c, Extended Data Fig. 2e-f). Moreover,
107 polarity of RPA binding was indicative of the accumulation of 3'-overhangs (Extended Data
108 Fig. 2f). Thus, loss of WRN in MSI cells leads to DSBs with extensive 5'-3' end-processing.

109 The left- and right-end of the DSBs were separated from each other by a variable
110 distance (Fig. 1c, Extended Data Fig. 2g). Interestingly, this “gap” harbored a major reduction
111 in sequencing reads, suggesting that DSBs occur at the borders of these regions. We therefore
112 searched for specific DNA motifs in the gap region in the hg19 reference genome, which
113 revealed a dominant TA dinucleotide repeat motif (Fig. 1d) with a median repeat length of 51
114 bp (Extended Data Fig. 2h). Nearly all breaks associated with WRN-deficiency in both KM12
115 and HCT116 cells occurred at TA dinucleotide repeats (Fig. 1e). However, only about 8% of
116 all TA repeats in the reference genome (5,400 out of 66,644) were associated with DSB
117 (Extended Data Fig. 2h). We conclude that DSBs flank a fraction of TA repeats. Hereafter, we
118 refer to these sites as “broken (TA)_n repeats,” and to the DSBs themselves as “TA breaks”.

119 Cruciform structures form at (TA)_n repeats in plasmids in *E. coli* and yeast when their
120 length exceeds roughly 20-22 repeat units¹³⁻¹⁵. Long (TA)_n tracts also cause replication fork
121 stalling and chromosome fragility at late replicating common fragile sites (CFS)¹⁶⁻¹⁸.
122 Exogenous replication stress further enhances replication fork collapse at simple repeats^{12,19},
123 including CFSs. Accordingly, we found that WRN depletion induces DNA breakage precisely
124 at (TA)_n repeats within multiple CFSs (Extended Data Fig. 3a) and at palindromic TA-rich
125 repeats (Extended Data Fig. 3b), which have been proposed to form cruciform

126 structures²⁰. These data suggest that (TA)_n repeats at these sites might fold into secondary
127 structures that are targeted by WRN.

128

129 **MUS81-EME1 shatters chromosomes in WRN-depleted MSI cells**

130 MUS81-EME1 is a structure-specific endonuclease that processes late recombination
131 intermediates at CFSs²¹. MUS81-EME1 forms a complex with the scaffolding protein SLX4
132 that hyperactivates it at the G2/M boundary²². Recent studies indicate that the yeast MUS81-
133 EME1 homolog (Mus81-Mms4) causes DSBs at (TA)_n repeats when the tract exceeds a
134 threshold length for forming cruciform structures¹⁷. To determine whether the DSBs that
135 accumulate in WRN-depleted MSI cells are MUS81- and SLX4-dependent, we depleted these
136 factors prior to WRN depletion (Extended Data Fig. 4a). MUS81 or SLX4 depletion
137 dramatically reduced chromosome shattering (Fig. 2a, Extended Data Fig. 4b). Consistent
138 with this result, depleting MUS81 and SLX4 strongly reduced pKAP1 signaling (Extended
139 Data Fig. 4c) and DSB formation at (TA)_n repeats (Fig. 2b, Extended Data Fig. 4d). Thus,
140 MUS81 and SLX4 nucleases induce toxic chromosome breakage when WRN is depleted from
141 MSI cells.

142

143 **MSI correlates with susceptibility to recombinant MUS81-EME1 cleavage**

144 To test whether MUS81-EME1 acts directly on secondary structures in MSI cells, we treated
145 agarose-embedded KM12 DNA with recombinant MUS81-EME1 *in situ* prior to performing
146 END-seq. MUS81-EME1 promotes the resolution of cruciform structures by a nick and
147 counter-nick mechanism (Extended Data Fig. 4e), and we found that MUS81-EME1
148 generated recurrent and reproducible DSBs (Fig. 2c-d , Extended Data Fig. 4f). These
149 overlapped strikingly with DSBs generated from WRN depletion in these cells (Extended
150 Data Fig. 4g-h). These data show that secondary structures accumulate and can be cleaved by

151 MUS81-EME1 even at baseline conditions in MSI cells. In contrast to DSBs that were highly
152 resected after WRN depletion, *in situ* cleavage by MUS81-EME1 led to accumulation of
153 reads precisely adjacent to the border of (TA)_n repeats (Fig. 2d, Extended Data 4i). Thus, the
154 broad distribution of END-seq reads in the WRN-depleted samples are indeed caused by 5'
155 end-processing *in vivo* post MUS81 cleavage.

156 We hypothesized that WRN is recruited to unwind DNA secondary structures before
157 they can be cleaved by physiologically active MUS81. To test this, we incubated agarose-
158 embedded KM12 DNA with recombinant human WRN prior to MUS81-EME1 treatment *in*
159 *situ*, followed by END-seq detection of DSBs. While incubation with WRN alone did not
160 result in DSB formation, WRN pre-treatment substantially decreased MUS81-EME1-
161 cleavage (Fig. 2e, Extended Data Fig. 4j). These results suggest that WRN melts secondary
162 DNA structures at (TA)_n repeats in MSI cells.

163 If structure-forming (TA)_n repeats are responsible for massive breakage in MSI lines,
164 there should be fewer such structures in MSS cells. To test this, we performed *in situ* MUS81-
165 EME1 cleavage assay in two MSS cell lines (SW837 and RPE-1) and compared these to MSI
166 cell lines KM12 and HCT116. While the MSI cells displayed an overlapping set of strong
167 MUS81-EME1 cleavage sites (Extended Data Fig. 5a, b), the substrates for MUS81-EME1
168 cleavage were significantly reduced in the MSS cells. Thus, structure-forming (TA)_n repeats
169 accumulate in much higher abundance in MSI compared to MSS cells.

170

171 **Structure forming (TA)_n repeats are susceptible to replication fork stalling and collapse**

172 DNA polymerase stalling generates RPA-bound ssDNA, which activates checkpoint kinase
173 ataxia telangiectasia and Rad3 related (ATR) to protect the replication fork. RPA ChIP-seq in
174 WRN-proficient KM12 cells revealed an enrichment of RPA in the vicinity of broken (TA)_n
175 sites, suggesting that stalled forks spontaneously accumulate at these sites (Fig. 3a). WRN is

176 recruited to stalled replication forks in a manner that requires ATR phosphorylation²³.
177 Consistent with this, HCT116 cells expressing WRN mutants WRN^{3A} or WRN^{6A} containing
178 alanine substitutions at identified ATR phosphorylation sites²³ showed an increase in KAP1
179 phosphorylation when endogenous WRN was ablated by an siRNA targeting the 5'
180 untranslated region (Extended Data Fig. 5c). We therefore hypothesized that stalled forks at
181 structure-forming (TA)_n repeats would be susceptible to fork collapse upon ATR kinase
182 inhibition (ATRi). To test this, we treated MSI and MSS cells with ATRi as well as low dose
183 aphidicolin (APH) to partially inhibit DNA polymerase elongation. As determined by END-
184 seq, replication forks collapsed into DSBs preferentially at (TA)_n repeats, including those
185 within CFSs (Fig. 3b, c, Extended Data Fig. 5d). The frequency fork collapse at (TA)_n repeats
186 in MSI cells was at least 30-fold higher than in MSS cells (Fig. 3b, c), and these sites largely
187 overlapped with DSBs generated in the absence of WRN (Extended Data Fig. 5e). Thus,
188 secondary structure-forming (TA)_n repeats are associated with replication fork stalling which
189 is overcome by WRN, possibly through activation by ATR.

190

191 **(TA)_n repeats inhibit DNA synthesis *in vitro***

192 In MSI cancers, insertion/deletions of a few nucleotides are commonly found at
193 mononucleotide repeats⁵. Among dinucleotides, (TA)_n repeats are reported to be the least
194 frequently altered (Extended Data Fig. 6a)²⁴. To determine why (TA)_n repeats are susceptible
195 to breakage, we performed a PCR-based size analysis of several (TA)_n repeats in a panel of
196 MSI and MSS cell lines (Fig. 4a, Extended Data Fig. 6b, Supplementary Table 1). We
197 examined eight (TA)_n repeats that were broken and three that were not broken upon WRN
198 depletion in MSI cells. All 11 sites that we queried harbored uninterrupted stretches of more
199 than 20 (TA)_n units (in the hg19 reference genome), and were amplified with primers
200 flanking the (TA)_n repeat.

201 While MSS cell lines showed the predicted PCR products at broken and non-broken
202 sites (Fig. 4a, Extended Data Fig. 6b), many recurrently broken sites failed to be amplified in
203 KM12 cells (Fig. 4a). These same sites frequently could not be amplified with genomic DNA
204 from other MSI cell lines, highlighting that identical sites are affected across distinct cancers
205 (Extended Data Fig. 6b). A notable exception was HCT116, which showed the expected
206 banding pattern. However, as shown below, this likely reflects that only one of the two alleles
207 is of normal size and can therefore be amplified.

208 The absence of PCR products at multiple broken (TA)_n repeats in MSI cell genomic
209 DNA was unlikely to be caused by large deletions that included primer binding sites, as these
210 sites were highly covered in END-seq and whole genome sequencing data (below). We
211 therefore speculated that these sites contain large, structure-forming expansions that might
212 inhibit polymerase extension during PCR amplification. To test this hypothesis, we performed
213 PCR-free whole genome sequencing of KM12 and HCT116 cells, enabling us to reach an
214 average sequencing depth greater than 100x. By inspection, it was immediately clear that
215 broken (TA)_n repeats displayed markedly lower read depth compared to flanking regions and
216 non-broken (TA)_n repeats (Extended Data Fig. 6c). Indeed, broken (TA)_n repeats (mapped in
217 WRN-depleted MSI cells) displayed the most significant drop in coverage relative to other
218 classes of mono- and di-nucleotide repeats (Fig. 4b, Extended Data Fig. 6d). Finally, as
219 assessed by CRISPR-Cas9 and RNA interference fitness screens², the degree of dependence
220 on WRN for human cancer cell survival was inversely correlated with sequencing coverage at
221 broken (TA)_n sites (Fig. 4c). Thus, structure forming (TA)_n repeats are highly recalcitrant to
222 sequencing. While not previously documented as mutated in MSI cancers, (TA)_n repeats are
223 nevertheless predictive of WRN dependency.

224

225 **MSI cells accumulate large-scale (TA)_n repeat expansions**

226 Since the identification of large-scale changes in short tandem repeats is challenging with
227 short read sequencing, we utilized the ExpansionHunter algorithm²⁵, which has been used to
228 detect various large pathogenic repeat expansions. This analysis revealed that broken (TA)_n
229 repeats exhibit large-scale expansions compared to non-broken (TA)_n repeats (Extended Data
230 Fig. 7a, b). Consistent results were obtained with the expanded short tandem repeat algorithm
231 (exSTRa)²⁶ (Extended Data Fig. 7c).

232 Since large expansions should increase the size of (TA)_n repeats relative to the
233 reference genome, we performed Southern blotting with probes spanning two broken (B1 and
234 B2) and two non-broken sites (NB1 and NB3) previously analyzed by PCR (Fig. 4d,
235 Extended data Fig. 7d-e). For this analysis, we compared two MSI (KM12 and HCT116) and
236 two MSS (RPE-1 and eHAP) cell lines. At the non-broken (TA)_n repeats, we detected bands
237 corresponding to the expected sizes in all cell lines (Extended Data Fig. 7d). Additionally,
238 neither broken (TA)_n region was altered in control RPE-1 or eHAP cell lines (Extended Data
239 Fig. 7d). For broken site B1, we detected one allele of the expected size, but also one
240 expanded allele in both KM12 and HCT116 (Fig. 4d, Extended Data Fig. 6b). The shorter,
241 non-expanded allele likely corresponds to the PCR product detected at B1 in both cell lines.
242 At broken site B2, we detected one expanded allele in HCT116, while the other allele was the
243 expected size (Fig. 4d), corresponding to the product detected by PCR (Extended Data Fig.
244 6b). In KM12 cells, broken site B2 exhibited distinct expansions on the two alleles (Fig. 4d,
245 Extended Data Fig. 7e), and no PCR product was detectable at this site (Fig. 4a). Taken
246 together, our data suggest that (TA)_n repeats undergo large-scale expansion, and that
247 expanded alleles inhibit DNA synthesis during PCR and short read whole genome
248 sequencing.

249 Long-read sequencing technology has successfully characterized expansions of simple
250 tandem repeats responsible for various diseases²⁷. To determine whether a similar approach

251 could uncover variations in secondary structure-forming (TA)_n repeats in MSI cells, we
252 utilized the Pacific Biosciences (PacBio) continuous long-read (CLR) sequencing platform²⁷.
253 With this platform, MSI (HCT116 and KM12) and MSS (SW620 and SW837) colon cancer
254 cell lines displayed similar coverage at broken (TA)_n and other simple repeats (Extended Data
255 Fig. 7f). We detected a large range of expansions at different broken (TA)_n sites in MSI cells,
256 with median repeat length expanding from 54bp in the hg19 reference genome to 91bp and
257 125bp in HCT116 and KM12, respectively (Fig. 4e). In contrast, the median lengths of non-
258 broken (TA)_n repeats in HCT116 (75bp) and KM12 (74bp) were much more similar to the
259 reference genome (72bp). MSS cell lines did not show substantial expansions at either broken
260 or non-broken (TA)_n sites (Fig. 4e). Motif analysis within the broken sites revealed that the
261 expanded alleles consisted of almost pure (TA)_n (Extended Data Fig. 7g). Thus, MSI
262 predisposes (TA)_n repeat tracts to undergo large-scale expansions.

263

264 **Length of uninterrupted TA repeats and replication timing contribute to breakage**

265 Of all annotated (TA)_n repeats in the reference genome, only 8% display DSBs upon WRN
266 depletion in MSI cells (Extended Data Fig. 2h). Thus, it remains unclear what characterizes
267 the underlying expansion and breakage specifically at these sites and why susceptible sites are
268 recurrent across multiple distinct MSI cell lines. Changes in the length of microsatellite DNA
269 are thought to arise from replication slippage caused by the transient dissociation of the
270 polymerase from replicating DNA strands followed by self-annealing and misaligned
271 reassociation²⁸. Since the probability of self-annealing increases with the purity and length of
272 the repeat sequence and decreases with repeat interruptions²⁸, we speculated that
273 predisposition to large-scale expansions of (TA)_n repeats, and therefore breakage, would be
274 influenced by the exact sequence composition of the repeat.

275 Indeed, we found broken sites have a higher TA content, fewer interruptions in the
276 (TA)_n repeat, and longer uninterrupted (TA)_n sequences relative to non-broken sites as
277 assessed both in the reference genome (Extended Data Fig. 8a-c) and in long-read sequencing
278 reads from MSI cells (Extended Data Fig. 8d-f). Moreover, longer pure (TA)_n repeats were
279 much more likely to expand (Extended Data Fig. 8g). Quantitative modeling showed that
280 three features of (TA)_n repeats were predictive of breakage in WRN-deficient cells: the
281 probability that they form secondary structures as measured by MUS81-EME1 cleavage *in*
282 *situ*, the sizes of expansions determined by long-read sequencing, and the likelihood that they
283 occur in late replicating regions (Extended Data Fig. 8h, j). While these features cannot fully
284 predict the propensity of (TA)_n repeats to break upon WRN deficiency, the analysis
285 demonstrates that longer, uninterrupted (TA)_n repeats are more likely to undergo expansion in
286 MSI cell lines, where they form secondary structures and disrupt DNA replication. In the
287 absence of WRN, (TA)_n repeats that are replicated in late S phase are less likely to be
288 resolved prior to mitosis, when they are recognized and cleaved by MUS81 to generate DSBs.

289

290 **Relationship between expandable repeats and genome stability**

291 Secondary structure forming (TA)_n tracts within CFSs cause replication fork stalling and
292 chromosome breakage. CFSs are also associated with deletions in multiple tumor types²⁹,
293 likely due to their susceptibility to replication fork collapse. Since broken (TA)_n sites in MSI
294 cells are also susceptible to replication fork stalling and collapse (Fig. 3), we hypothesized
295 they could be hotspots for deletions in MSI cancers. Although genomic instability in MSI
296 cancers is most frequently associated with small insertions and deletions, large Kb-Mb scale
297 deletions of unknown etiology have been detected³⁰. We therefore analyzed data from pan-
298 cancer whole genome sequencing cohorts in uterine corpus endometrial carcinoma (UCEC),
299 colon adenocarcinoma (COAD) and gastric adenocarcinoma (STAD)³¹. Breakpoints

300 associated with deletions (ranging in size from 459bp to 176Mb) in MSI (24) and MSS (93)
301 cancers were identified⁵. In several cases, one or both deletion breakpoints mapped precisely
302 to broken (TA)_n repeats (Extended Data Fig. 9a, b). We then calculated the enrichment of
303 different annotated repeats at the tumor breakpoints compared to their enrichment at random
304 breakpoints with comparable chromosome- and size- distributions, and found that MSI tumor
305 breakpoints showed the greatest enrichment at broken (TA)_n repeats (Extended Data Fig. 9c).
306 In contrast, these sites were not enriched for deletion breakpoints in MSS tumors (Extended
307 Data Fig. 9c). Thus, in MSI cancers, large scale deletions are frequently associated with
308 broken (TA)_n sites, suggesting that they are inherently fragile.

309 Based on our results, we propose the following model (Extended Data Fig. 8k): MMR
310 deficient cells undergo microsatellite instability, which gradually manifests as large-scale
311 expansions at (TA)_n repeats. The most susceptible sites are those that contain pure,
312 uninterrupted TA-dinucleotide repeats. Over the course of months to years, (TA)_n repeats
313 reach a threshold length above which they extrude into cruciform-like structures, perhaps as a
314 result of negative supercoiling during nucleosome removal ahead of the replication fork.
315 These structures would stall replication forks and trigger ATR-dependent WRN
316 phosphorylation, which promotes unwinding of the secondary structure to complete
317 replication. In the absence of WRN, the structure-specific MUS81-EME1 endonuclease, and
318 its scaffold SLX4, cleave these structures in an attempt to salvage the replication fork.
319 However, thousands of concerted MUS81 cleavage events lead to extensive DNA end-
320 resection, RPA exhaustion (Fig. 1b-c)⁸, chromosomal fragmentation, and cell death.

321 Thirty years ago, Vogelstein and colleagues discovered the *DCC* (Deleted in Colon
322 Cancer) gene, whose expression is greatly reduced in colon cancer³². By Southern blot
323 analysis, they observed that several MSI tumor cell lines harbored “insertions” up to 300 bp at
324 a locus containing an uninterrupted (TA)₂₂ repeat just downstream of *DCC* exon 7. However,

325 numerous attempts to clone or amplify alleles with insertions failed, leading them to conclude
326 that the inserted sequence might form an unusual DNA structure. Strikingly, we detected a
327 MUS81-EME1 sensitive DNA structure precisely at the same (TA)_n repeat (Extended Data
328 Fig. 10a), and long-read sequencing confirmed a (TA)_n expansion at that locus (Extended
329 Data Fig. 10b, right highlighted sequence). Thus, our data demonstrate that rather than a non-
330 templated insertion, DCC contains a structure-forming, MSI-expanded (TA)_n repeat.

331 Anti-PD-1 antibodies have been approved for use in patients with cancers with
332 mismatch repair deficiency or MSI, independent of the cancer lineage. The therapeutic
333 response is correlated with the number indels in coding regions, which can generate
334 immunogenic neoantigens^{33,34}. (TA)_n expansions are mostly localized outside of coding
335 regions, and therefore might not generate immunogenic neoantigens. However, DSBs at
336 (TA)_n motifs are capable of triggering innate cytosolic DNA- and RNA-dependent sensing
337 and signaling pathways^{35,36}, indicating that WRN inhibition has the potential for further
338 stimulating immune responses. Although indel mutations of a few nucleotides and large-scale
339 (TA)_n expansions are both features of MSI, further studies will be necessary to determine
340 whether these defects arise through different mechanisms, in distinct MSI tumors, or within
341 different subclones. In summary, our findings provide a mechanistic explanation for WRN
342 dependence with MSI and identify a novel biomarker to guide selection of patients where
343 WRN inhibition may be effective in combination with immune checkpoint blockade or as an
344 independent line of therapy.

345

346 **References**

- 347
- 348 1 Behan, F. M. *et al.* Prioritization of cancer therapeutic targets using CRISPR-Cas9
349 screens. *Nature* **568**, 511-516, doi:10.1038/s41586-019-1103-9 (2019).
- 350 2 Chan, E. M. *et al.* WRN helicase is a synthetic lethal target in microsatellite unstable
351 cancers. *Nature* **568**, 551-556, doi:10.1038/s41586-019-1102-x (2019).
- 352 3 Kategaya, L., Perumal, S. K., Hager, J. H. & Belmont, L. D. Werner Syndrome
353 Helicase Is Required for the Survival of Cancer Cells with Microsatellite Instability.
354 *iScience* **13**, 488-497, doi:10.1016/j.isci.2019.02.006 (2019).
- 355 4 Lieb, S. *et al.* Werner syndrome helicase is a selective vulnerability of microsatellite
356 instability-high tumor cells. *Elife* **8**, doi:10.7554/eLife.43333 (2019).
- 357 5 Fujimoto, A. *et al.* Comprehensive analysis of indels in whole-genome microsatellite
358 regions and microsatellite instability across 21 cancer types. *Genome Res*,
359 doi:10.1101/gr.255026.119 (2020).
- 360 6 Dudley, J. C., Lin, M. T., Le, D. T. & Eshleman, J. R. Microsatellite Instability as a
361 Biomarker for PD-1 Blockade. *Clin Cancer Res* **22**, 813-820, doi:10.1158/1078-
362 0432.CCR-15-1678 (2016).
- 363 7 Chu, W. K. & Hickson, I. D. RecQ helicases: multifunctional genome caretakers. *Nat*
364 *Rev Cancer* **9**, 644-654, doi:10.1038/nrc2682 (2009).
- 365 8 Toledo, L. I. *et al.* ATR prohibits replication catastrophe by preventing global
366 exhaustion of RPA. *Cell* **155**, 1088-1103, doi:10.1016/j.cell.2013.10.043 (2013).
- 367 9 Canela, A. *et al.* DNA Breaks and End Resection Measured Genome-wide by End
368 Sequencing. *Mol Cell* **63**, 898-911, doi:10.1016/j.molcel.2016.06.034 (2016).
- 369 10 Paiano, J. *et al.* ATM and PRDM9 regulate SPO11-bound recombination
370 intermediates during meiosis. *Nat Commun* **11**, 857, doi:10.1038/s41467-020-14654-
371 w (2020).
- 372 11 Khil, P. P., Smagulova, F., Brick, K. M., Camerini-Otero, R. D. & Petukhova, G. V.
373 Sensitive mapping of recombination hotspots using sequencing-based detection of
374 ssDNA. *Genome Res* **22**, 957-965, doi:10.1101/gr.130583.111 (2012).
- 375 12 Tubbs, A. *et al.* Dual Roles of Poly(dA:dT) Tracts in Replication Initiation and Fork
376 Collapse. *Cell* **174**, 1127-1142 e1119, doi:10.1016/j.cell.2018.07.011 (2018).
- 377 13 Bowater, R., Aboul-ela, F. & Lilley, D. M. Large-scale stable opening of supercoiled
378 DNA in response to temperature and supercoiling in (A + T)-rich regions that promote
379 low-salt cruciform extrusion. *Biochemistry* **30**, 11495-11506,
380 doi:10.1021/bi00113a003 (1991).
- 381 14 Dayn, A. *et al.* Formation of (dA-dT)_n cruciforms in Escherichia coli cells under
382 different environmental conditions. *J Bacteriol* **173**, 2658-2664,
383 doi:10.1128/jb.173.8.2658-2664.1991 (1991).
- 384 15 McClellan, J. A., Boublikova, P., Palecek, E. & Lilley, D. M. Superhelical torsion in
385 cellular DNA responds directly to environmental and genetic factors. *Proc Natl Acad*
386 *Sci U S A* **87**, 8373-8377, doi:10.1073/pnas.87.21.8373 (1990).
- 387 16 Zlotorynski, E. *et al.* Molecular basis for expression of common and rare fragile sites.
388 *Mol Cell Biol* **23**, 7143-7151, doi:10.1128/mcb.23.20.7143-7151.2003 (2003).
- 389 17 Kaushal, S. *et al.* Sequence and Nuclease Requirements for Breakage and Healing of a
390 Structure-Forming (AT)_n Sequence within Fragile Site FRA16D. *Cell Rep* **27**, 1151-
391 1164 e1155, doi:10.1016/j.celrep.2019.03.103 (2019).
- 392 18 Wang, H. *et al.* CtIP maintains stability at common fragile sites and inverted repeats
393 by end resection-independent endonuclease activity. *Mol Cell* **54**, 1012-1021,
394 doi:10.1016/j.molcel.2014.04.012 (2014).

- 395 19 Shastri, N. *et al.* Genome-wide Identification of Structure-Forming Repeats as
396 Principal Sites of Fork Collapse upon ATR Inhibition. *Mol Cell* **72**, 222-238 e211,
397 doi:10.1016/j.molcel.2018.08.047 (2018).
- 398 20 Inagaki, H. *et al.* Chromosomal instability mediated by non-B DNA: cruciform
399 conformation and not DNA sequence is responsible for recurrent translocation in
400 humans. *Genome Res* **19**, 191-198, doi:10.1101/gr.079244.108 (2009).
- 401 21 Minocherhomji, S. & Hickson, I. D. Structure-specific endonucleases: guardians of
402 fragile site stability. *Trends Cell Biol* **24**, 321-327, doi:10.1016/j.tcb.2013.11.007
403 (2014).
- 404 22 Wyatt, H. D., Laister, R. C., Martin, S. R., Arrowsmith, C. H. & West, S. C. The SMX
405 DNA Repair Tri-nuclease. *Mol Cell* **65**, 848-860 e811,
406 doi:10.1016/j.molcel.2017.01.031 (2017).
- 407 23 Ammazalorso, F., Pirzio, L. M., Bignami, M., Franchitto, A. & Pichierri, P. ATR and
408 ATM differently regulate WRN to prevent DSBs at stalled replication forks and
409 promote replication fork recovery. *EMBO J* **29**, 3156-3169,
410 doi:10.1038/emboj.2010.205 (2010).
- 411 24 Cortes-Ciriano, I., Lee, S., Park, W. Y., Kim, T. M. & Park, P. J. A molecular portrait
412 of microsatellite instability across multiple cancers. *Nat Commun* **8**, 15180,
413 doi:10.1038/ncomms15180 (2017).
- 414 25 Dolzhenko, E. *et al.* Detection of long repeat expansions from PCR-free whole-
415 genome sequence data. *Genome Res* **27**, 1895-1903, doi:10.1101/gr.225672.117
416 (2017).
- 417 26 Tankard, R. M. *et al.* Detecting Expansions of Tandem Repeats in Cohorts Sequenced
418 with Short-Read Sequencing Data. *Am J Hum Genet* **103**, 858-873,
419 doi:10.1016/j.ajhg.2018.10.015 (2018).
- 420 27 Mitsushashi, S. & Matsumoto, N. Long-read sequencing for rare human genetic
421 diseases. *J Hum Genet* **65**, 11-19, doi:10.1038/s10038-019-0671-8 (2020).
- 422 28 Khristich, A. N. & Mirkin, S. M. On the wrong DNA track: Molecular mechanisms of
423 repeat-mediated genome instability. *J Biol Chem* **295**, 4134-4170,
424 doi:10.1074/jbc.REV119.007678 (2020).
- 425 29 Glover, T. W., Wilson, T. E. & Arlt, M. F. Fragile sites in cancer: more than meets the
426 eye. *Nat Rev Cancer* **17**, 489-501, doi:10.1038/nrc.2017.52 (2017).
- 427 30 Cancer Genome Atlas, N. Comprehensive molecular characterization of human colon
428 and rectal cancer. *Nature* **487**, 330-337, doi:10.1038/nature11252 (2012).
- 429 31 Consortium, I. T. P.-C. A. o. W. G. Pan-cancer analysis of whole genomes. *Nature*
430 **578**, 82-93, doi:10.1038/s41586-020-1969-6 (2020).
- 431 32 Fearon, E. R. *et al.* Identification of a chromosome 18q gene that is altered in
432 colorectal cancers. *Science* **247**, 49-56, doi:10.1126/science.2294591 (1990).
- 433 33 Ding, L. & Chen, F. Predicting Tumor Response to PD-1 Blockade. *N Engl J Med*
434 **381**, 477-479, doi:10.1056/NEJMcibr1906340 (2019).
- 435 34 Mandal, R. *et al.* Genetic diversity of tumors with mismatch repair deficiency
436 influences anti-PD-1 immunotherapy response. *Science* **364**, 485-491,
437 doi:10.1126/science.aau0447 (2019).
- 438 35 Feng, X. *et al.* ATR inhibition potentiates ionizing radiation-induced interferon
439 response via cytosolic nucleic acid-sensing pathways. *EMBO J*, e104036,
440 doi:10.15252/emboj.2019104036 (2020).
- 441 36 Harding, S. M. *et al.* Mitotic progression following DNA damage enables pattern
442 recognition within micronuclei. *Nature* **548**, 466-470, doi:10.1038/nature23470
443 (2017).
- 444
- 445

446 **Figure Legends**

447

448 **Fig. 1: WRN depletion in MSI cells induces recurrent DSBs at (TA)_n dinucleotide**

449 **repeats.**

450 (a) Representative metaphase spreads from KM12 NT, KM12-shWRN and SW837-shWRN

451 cells. Cells were treated with DMSO (NT) or doxycycline (shWRN) for 48 hrs.

452 (b) Genome browser screenshots displaying END-seq and RPA-ssDNA CHIP-seq profiles as

453 normalized read density (reads per million, RPM) for KM12-shWRN cells. Positive and

454 negative strand END-seq reads are displayed in black and grey, and positive and negative

455 strand RPA-ssDNA CHIP-seq reads in blue and red, respectively.

456 (c) Genome browser screenshot zoomed in on the highlighted region in panel b (orange box).

457 The light blue highlight indicates the gap region between the left and right ends of the

458 DSB.

459 (d) Motif analysis for sequence enrichment in the gap between positive and negative END-

460 seq peaks in KM12-shWRN cells.

461 (e) Fraction of END-seq peaks occurring at sites of (TA)_n repeats for KM12-shWRN and

462 HCT116-shWRN cells treated with either doxycycline (shWRN) or WRN siRNAs

463 (siWRN) for 72 hrs.

464

465 **Fig. 2: TA breaks are dependent on structure-specific endonucleases MUS81-EME1 and**

466 **SLX4.**

467 (a) Quantification of metaphases displaying chromosome shattering (defined as the absence

468 of intact chromosomes) in KM12-shWRN cells treated with doxycycline plus non-

469 targeting (siCTRL), MUS81 (siMUS81), or SLX4 (siSLX4) siRNAs for 72hrs. Student's

470 t-test values are shown.

- 471 (b) Quantification of END-seq signal intensity at broken TA repeats in KM12-shWRN cells
472 with same treatment as in panel (a). *** indicates a p-value $< 10^{-16}$.
- 473 (c) Genome browser screenshot for KM12-shWRN cells treated with doxycycline (shWRN,
474 top), and DMSO-treated cells processed with purified recombinant MUS81-EME1
475 enzyme *in situ* (bottom).
- 476 (d) Zoom-in of TA break at highlighted region in Fig. 2c (orange box).
- 477 (e) Quantification of END-seq peak intensity for DMSO-treated KM12-shWRN cells
478 processed *in situ* with either purified recombinant MUS81-EME1, WRN, or WRN
479 followed by MUS81-EME1. For the latter, Proteinase K digestion was performed between
480 the two enzymatic treatments. *** indicates a p-value $< 10^{-16}$.

481

482 **Fig. 3: Replication stalling and collapse at (TA)_n repeats in MSI cell lines.**

483

- 484 (a) Composite plot of RPA-ssDNA ChIP-seq signal (blue: RPA-ChIP, red: input DNA)
485 around broken (TA)_n repeats in KM12 cells.
- 486 (b) Genome browser screenshot displaying END-seq profiles for KM12-shWRN cells,
487 HCT116-shWRN, RPE-1-shWRN, and eHAP-shWRN cells treated with doxycycline
488 (shWRN) for 72hrs or APH+ATRi for 8 hrs. MSI cells are marked in red, MSS cell lines
489 in blue.
- 490 (c) Barplots indicating percentage of total DSBs located at (TA)_n repeats after APH+ATRi
491 treatment, as shown in panel (b). MSI cell are marked in red, MSS cell lines in blue.

492

493 **Fig. 4: (TA)_n repeats undergo large-scale expansion in MSI cell lines.**

- 494 (a) Agarose gels showing PCR fragments (or lack thereof) at (TA)_n repeats in KM12
495 (MSI) and SW837 (MSS) cells. B1-B8 were chosen based on the presence of END-

496 seq peaks upon WRN depletion in KM12 cells. NB1-NB3 were chosen for similar
497 (TA)_n repeat lengths as broken sites without breakage upon WRN depletion in KM12
498 cells. Ladder fragment sizes (in bp) are displayed. For gel source data, see
499 Supplementary Figure 1.

500 (b) Boxplots displaying coverage at different classes of repeats in WGS data from KM12
501 cells. Dotted red lines indicate the average coverage over the genome.

502 (c) Cell lines plotted by their average WRN dependency score and sequencing coverage
503 of broken (TA)_n loci. fpbm: fragments per base per million.

504 (d) Southern blots for two different genomic regions containing broken (TA)_n repeats
505 corresponding to the same sites in panel (a). Red markers and dotted lines represent
506 expected fragment sizes. For gel source data, see Supplementary Figure 1.

507 (e) Boxplot of long-read sequencing data demonstrating total length of broken and non-
508 broken (TA)_n in indicated cell lines compared to hg19 reference genome. MSI (red),
509 MSS (blue). *** p-value <10⁻⁵. Wilcox testing for the alternative hypothesis that the
510 broken (TA)_n has a greater group mean than the non-broken (TA)_n was used. p values
511 were corrected for multiple testing using the Benjamini-Hochberg method.

512

513

514 **Extended Data Figure Legends**

515

516 **Extended Data Fig. 1: WRN depletion induces DNA damage in different MSI cell lines.**

517 (a) Western Blot analysis of MLH1, MSH2, and GAPDH protein levels in human stomach

518 epithelial cells (HSEC) following CRISPR/Cas9 knockout. *sgLuc*: control sgRNA

519 targeting *luciferase*. *sgMLH1* and *sgMSH2*: sgRNAs targeting *MLH1* and *MSH2*,

520 respectively. For gel source data, see Supplementary Figure 1.

521 (b) Relative viability 7 days after sgRNA transduction in HSEC. Negative controls targeting

522 chromosome 2 intergenic sites: *sgCh2.2* and *sgCh2.4*. Pan-essential control: *sgPolR2D*.

523 Experimental sgRNA targeting *WRN*: *sgWRN2*, *sgWRN3*. Two-tailed Student's t-test

524 values are shown.

525 (c) Example of flow cytometry gating strategy used in Extended Data Fig. 1d, 4c.

526 (d) Flow cytometry profiles for exponentially growing KM12-shWRN cells treated with

527 DMSO (NT) or doxycycline (shWRN) for 72 hrs. EdU was added during the last 30

528 minutes before harvesting cells. Percentage of cells in the gates are indicated.

529 (e) Western blot analysis of WRN protein levels in KM12-shWRN and SW837-shWRN

530 treated with DMSO or doxycycline for 72 hrs. For gel source data, see Supplementary

531 Figure 1.

532 (f) Western blot analysis of WRN and pKAP1 protein levels in KM12-shWRN and KM12-

533 shWRN.C911 (non-targeting shRNA) treated with DMSO or doxycycline for 72 hrs. For

534 gel source data, see Supplementary Figure 1.

535

536 **Extended Data Fig. 2: WRN depletion induces recurrent and overlapping DSBs in MSI**

537 **cells.**

- 538 (a) Genome browser screenshot displaying END-seq profiles as normalized read density
539 (reads per million, RPM) for HCT116-shWRN cells and KM12-shWRN treated with
540 DMSO (NT) or doxycycline (shWRN), or transfected with non-targeting siRNAs
541 (siCTRL) or WRN siRNAs (siWRN) for 72 hrs.
- 542 (b) Scatterplots of END-seq peak intensity between replicate experiments of KM12-
543 shWRN and HCT116-shWRN cells treated with doxycycline for 72 hrs. Pearson
544 correlation coefficients are indicated.
- 545 (c) Venn diagrams showing overlap between peaks detected in HCT116-shWRN and
546 KM12-shWRN cells treated with either doxycycline (shWRN) or WRN siRNAs
547 (siWRN) for 72 hrs (p-value < 10^{-16} for both comparisons).
- 548 (d) Quantification of END-seq peak intensity for KM12-shWRN and SW837-shWRN
549 cells treated with doxycycline for 72 hrs. *** indicates a p-value < 10^{-16} .
- 550 (e) Venn diagram showing overlap between peaks identified from END-seq and RPA-
551 ssDNA ChIP-seq for KM12-shWRN cells treated with doxycycline for 72 hrs (p-value
552 < 10^{-16}).
- 553 (f) Composite plot of END-seq (black: positive strand reads, grey: negative strand reads)
554 and RPA-ssDNA ChIP-seq (blue: positive strand reads, red: negative strand reads)
555 signal around DSB sites in KM12-shWRN cells treated with doxycycline for 72 hrs.
- 556 (g) Heatmap displaying intensity of END-seq signal in KM12-shWRN cells treated with
557 doxycycline for 72 hrs, relative to the center of the gap between positive and negative
558 strand peaks. Sites are ordered by the size of the gap, from smallest to largest.
- 559 (h) Calculated size distribution from the reference genome of (TA)_n repeats either located
560 in gaps between positive and negative END-seq peaks (black, broken sites) or located
561 elsewhere in the genome (grey, non-broken sites), determined from KM12-shWRN
562 cells treated with doxycycline for 72 hrs.

563

564 **Extended Data Fig. 3: WRN depletion induces DNA breakage in common fragile sites**
565 **and palindromic TA-rich repeats in MSI cells.**

- 566 (a) Genome browser screenshot displaying END-seq profiles of common fragile sites
567 FRA16D, FRA3B, FRA10B, and FRA7I as normalized read density (reads per million,
568 RPM) for KM12-shWRN cells treated with DMSO (NT) or doxycycline (shWRN) for 72
569 hrs. Indicated is the number of uninterrupted (TA)_n repeat units in the hg19 reference
570 genome at DSB sites.
- 571 (b) Genome browser screenshot displaying END-seq profiles of PATRRs on chromosomes
572 11 and 22 as normalized read density (reads per million, RPM) for KM12-shWRN cells
573 treated with DMSO (NT) or doxycycline (shWRN) for 72 hrs.

574

575 **Extended Data Fig. 4: (TA)_n repeat-forming repeats in MSI cell lines are substrates for**
576 **MUS81-EME1.**

577

- 578 (a) qRT-PCR quantification of MUS81 and SLX4 mRNA levels in KM12-shWRN cells
579 transfected with non-targeting siRNAs (siCTRL), MUS81 siRNAs (siMUS81), or
580 SLX4 siRNAs (siSLX4).
- 581 (b) Representative images of metaphase spreads from KM12-shWRN cells treated with
582 doxycycline (shWRN) and either non-targeting siRNAs (siCTRL), MUS81 siRNAs
583 (siMUS81), or SLX4 siRNAs (siSLX4) for 48 hrs.
- 584 (c) Flow cytometric profiles for KAP1 phosphorylation in exponentially growing KM12-
585 shWRN cells treated with doxycycline (shWRN), plus non-targeting siRNAs
586 (siCTRL), MUS81 siRNAs (siMUS81), or SLX4 siRNAs (siSLX4) for 72 hrs.

- 587 (d) Genome browser screenshot displaying END-seq profiles as normalized read density
588 (reads per million, RPM) for KM12-shWRN cells treated with DMSO (NT), plus non-
589 targeting siRNAs (siCTRL), MUS81 siRNAs (siMUS81), or SLX4 siRNAs (siSLX4)
590 for 72 hrs.
- 591 (e) Schematic representation of DNA cruciform cleavage by MUS81-EME1 structure-
592 specific endonuclease.
- 593 (f) Venn diagram displaying overlap of END-seq peaks between two replicate
594 experiments of DMSO-treated KM12-shWRN cells processed with purified
595 recombinant MUS81-EME1 enzyme *in situ* (MUS81-EME1) (p-value < 10^{-16}).
- 596 (g) Venn diagram showing overlap in TA breaks between KM12-shWRN treated with
597 doxycycline (shWRN) for 72 hrs, and DMSO-treated cells processed with MUS81-
598 EME1 enzyme *in situ* (MUS81-EME1) (p-value < 10^{-16}).
- 599 (h) Venn diagram displaying overlap between END-seq peaks from DMSO-treated
600 KM12-shWRN and HCT116-shWRN cells processed *in situ* with MUS81-EME1 (p-
601 value < 10^{-16}).
- 602 (i) Genome-wide aggregate analysis of END-seq signal around TA breaks from KM12-
603 shWRN cells treated with doxycycline for 72 hrs (shWRN - black: positive strand
604 reads, grey: negative strand reads), or DMSO-treated KM12-shWRN cells processed
605 with purified recombinant MUS81-EME1 enzyme *in situ* (blue: positive strand reads,
606 red: negative strand reads).
- 607 (j) Genome browser screenshot displaying END-seq profiles for DMSO-treated KM12-
608 shWRN cells (WRN proficient) processed *in situ* with either purified recombinant
609 WRN, MUS81-EME1, or WRN followed by MUS81-EME1. For the latter, Proteinase
610 K digestion was performed between the two enzymatic treatments.

611

612 **Extended Data Fig. 5: Structure-forming repeats in MSI cells activate ATR.**

- 613 (a) Genome browser screenshot displaying END-seq profiles for DMSO-treated KM12-
614 shWRN, HCT116-shWRN, SW837-shWRN and RPE-1-shWRN cells processed *in situ*
615 with purified recombinant MUS81-EME1. Cells are indicated as MSI (red label) or MSS
616 (blue label).
- 617 (b) Quantification of END-seq peak intensity for libraries displayed in panel a.
- 618 (c) Western blot analysis of WRN and pKAP1 levels in HCT116 cells expressing WT WRN,
619 or ATR phosphorylation mutants WRN^{3A} or WRN^{6A}. Endogenous WRN was depleted
620 using an siRNA targeting WRN 5'UTR. For gel source data, see Supplementary Figure 1.
- 621 (d) Genome browser screenshot displaying END-seq profiles within FRA3B on chromosome
622 3 as normalized read density (reads per million, RPM) for KM12-shWRN cells, HCT116-
623 shWRN, RPE-1-shWRN, and eHAP-shWRN cells treated with doxycycline (shWRN) for
624 72hrs or APH+ATRi for 8 hrs.
- 625 (e) Venn diagrams displaying overlap of DSBs detected after WRN depletion or APH+ATRi
626 treatment in KM12 and HCT116 cells.

627

628 **Extended Data Fig. 6: (TA)_n repeat sequences are underrepresented in whole-genome**
629 **sequencing data from MSI cells.**

- 630 (a) Barplots indicating the percentage of recurrent mutations in different classes of repeats
631 (left panel - mono, di, tri and tetra) and a barplot (right panel) showing the number of
632 various dinucleotide repeats in the 1000 altered loci. The plots were based on
633 sequencing analysis from²⁴, which considered microsatellites smaller than 40bp.
- 634 (b) Agarose gels showing PCR fragments (or lack thereof) of sites of different (TA)_n
635 repeats in one MSS and four MSI cell lines. Broken sites B1-B8 were chosen based on
636 the presence of END-seq peaks upon WRN depletion in KM12 cells. Sites NB1-NB3

637 were chosen with similar (TA)_n repeat lengths as broken sites, but were not broken
638 upon WRN depletion in KM12 cells. Fragment sizes (in bp) are displayed. For gel
639 source data, see Supplementary Figure 1.

640 (c) Genome browser screenshots of short read PCR-free whole genome sequencing reads,
641 indicating coverage, in KM12 and HCT116 cell lines. Shown are two regions
642 containing (TA)_n repeats, one that displays END-seq peaks upon WRN depletion in
643 KM12 (left panels), and one that does not (right panels). Regions correspond to
644 equivalent PCR sites in Fig. 4a and Extended Data Fig. 5b.

645 (d) Boxplots displaying coverage at different classes of mono- and di-nucleotide repeats
646 in PCR-free whole-genome sequencing libraries made from HCT116 cells. (TA)_n
647 repeats are split into those that overlap END-seq peaks upon shWRN induction, and
648 those that do not harbor DSBs. Dotted red lines indicate the average coverage over the
649 genome.

650

651

652 **Extended Data Fig. 7: (TA)_n repeats undergo large-scale expansions in MSI cells.**

653 (a) Cumulative fraction of expanded (TA)_n repeats in KM12 and HCT116, based on
654 ExpansionHunter analysis of PCR-free whole genome sequencing data. TA repeats
655 were split into broken (red) and non-broken (black) based on presence or absence of
656 END-seq peaks after WRN depletion in KM12 cells.

657 (b) Graphical representation of a (TA)_n repeat expansion in HCT116. This site has 33
658 (TA)_n repeat units in the reference genome; ExpansionHunter identified an expansion
659 to 86-87 repeat units based on PCR-free whole genome sequencing of HCT116.

660 (c) Empirical cumulative distribution function based on the length by which each read
661 overlaps the (TA)_n repeat shown in panel b as identified by exSTRA.

- 662 (d) Southern blots for two different genomic regions containing non-broken (TA)_n repeats
663 corresponding to the same sites in Fig. 4a and Extended Data Fig. 6b. Red markers
664 and dotted lines represent expected fragment sizes. For gel source data, see
665 Supplementary Figure 1.
- 666 (e) Southern blots for broken (TA)_n repeat B2 (top) and non-broken (TA)_n repeat NB3
667 (bottom) in MSS (blue) and MSI (red) cell lines, confirming expansion of broken
668 (TA)_n repeats in MSI cell lines. Red markers and dotted lines represent expected
669 fragment sizes based on the reference genome. For gel source data, see Supplementary
670 Figure 1.
- 671 (f) Boxplots displaying coverage at different classes of repeats in long-read sequencing
672 libraries made from MSI (red) and MSS (blue) cells.
- 673 (g) Motif analysis for sequence enrichment at broken (TA)_n in the KM12 cell line from
674 long-read sequencing data.

675

676 **Extended Data Fig. 8: Large-scale expansions occur at long, uninterrupted (TA)_n repeat**
677 **sequences.**

678

- 679 (a) Boxplot showing, in the hg19 reference genome, the proportion of (TA)_n repeat units
680 found within the full annotated sequence at broken or non-broken (TA)_n repeats in
681 KM12 cells. *** indicates a p-value < 10⁻¹⁶.
- 682 (b) Boxplot showing, in the hg19 reference genome, the proportion of the longest run of
683 uninterrupted (TA)_n within the full annotated sequence at broken or non-broken (TA)_n
684 repeats in KM12 cells. *** indicates a p-value < 10⁻¹⁶.

- 685 (c) Boxplot showing, in the hg19 reference genome, the length (bp) of the longest
686 uninterrupted (TA)_n dinucleotide repeats within the full annotated sequence at broken
687 or non-broken (TA)_n repeats in KM12 cells. *** indicates a p-value < 10⁻¹⁶.
- 688 (d) Boxplot showing, in long read sequencing data, the proportion of (TA)_n repeat units
689 found within the full sequence at broken or non-broken (TA)_n repeats in KM12 cells.
690 *** indicates a p-value < 10⁻¹⁶.
- 691 (e) Boxplot showing, in long read sequencing data, the proportion of the longest run of
692 uninterrupted (TA)_n within the full sequence at broken or non-broken (TA)_n repeats in
693 KM12 cells. *** indicates a p-value < 10⁻¹⁶.
- 694 (f) Boxplot showing, in long read sequencing data, the length (bp) of the longest
695 uninterrupted (TA)_n dinucleotide repeat within the full sequence at broken or non-
696 broken (TA)_n repeats in KM12 cells. *** indicates a p-value < 10⁻¹⁶.
- 697 (g) Boxplot showing, in the hg19 reference genome, lengths (bp) of the longest
698 uninterrupted (TA)_n repeats for expanded or non-expanded (TA)_n repeats as
699 determined through long read sequencing data from KM12 cells.
- 700 (h) Multiple linear regression model predicting END-seq peak intensity of KM12 shWRN
701 cells treated with doxycycline (shWRN) for 72 hrs derived from END-seq intensity of
702 MUS81-EME1 cleavage *in situ*, replication timing, and expanded length of broken
703 (TA)_n. The Pearson correlation coefficient is indicated (see panel j).
- 704 (i) END-seq intensity of broken (TA)_n repeats in KM12 shWRN cells treated with
705 doxycycline (shWRN) for 72 hrs grouped by replication timing values from late
706 replicating to early replicating.
- 707 (j) Multiple linear regression was performed to predict END-seq peak intensity of KM12
708 shWRN cells treated with doxycycline (shWRN) for 72 hrs based on following
709 parameters: END-seq intensity of MUS81-EME1 cleavage *in situ*, replication timing,

710 and expanded length of broken (TA)_n. END-seq intensity of shWRN END-seq and
711 MUS81-EME1 cleavage were calculated using RPKM in ± 1 kb window around
712 broken (TA)_n. Mean value was used for replication timing quantification. Expanded
713 lengths were identified from long read sequencing data. Estimates of the standardized
714 regression coefficients (β) are shown, along with t statistics and P-values based on the
715 standardized coefficients.

716 (k) Model for MSI cell dependence on WRN. Large scale expansions of (TA)_n repeats are
717 associated with microsatellite instability in MMR deficient cells. When (TA)_n reach
718 above a critical length, they extrude into cruciform-like structures, which stall
719 replication forks and activate ATR kinase, which in turn phosphorylates WRN and
720 other substrates to complete DNA replication. In the absence of WRN, MUS81-
721 EME1/SLX4 cleaves secondary structures at (TA)_n repeats, thereby shattering the
722 chromosomes.

723

724 **Extended Data Fig. 9: Deletion breakpoints in MSI cancers are enriched at (TA)_n**
725 **repeats.**

- 726 a) Genome browser screenshot of a broken (TA)_n (defined from KM12), MSI deletion
727 (derived from a patient sample), and END-seq profile (in WRN-depleted KM12 cells).
728 The sequences around the breakpoints are shown in the inset.
- 729 b) Junctions associated with 6 different MSI deletions from patients. Seq1 represents the
730 sequence from -50bp to left breakpoint and Seq2 represents the sequence from right
731 breakpoint to +50bp.
- 732 c) Enrichment of simple repeats, broken and non-broken (TA)_n, and LINE, SINE and
733 LTR at patient deletion breakpoints relative to their overlap with random deletion
734 breakpoints of the same size (enrichment value=1). B(TA)_n-B(TA)_n represents cases

735 in which both breakpoints overlap with broken (TA)_n repeats; B(TA)_n- represents
736 cases in which only one breakpoint overlaps with a broken (TA)_n repeat. B(TA)_n,
737 broken TA repeat; NB(TA)_n, nonbroken TA repeat; LTR, long terminal repeat; SINE,
738 short interspersed nuclear element; LINE, long interspersed nuclear element.

739

740 **Extended Data Fig. 10: DNA breaks within *DCC* gene body**

741 (a) Genome browser screenshots within *DCC* gene displaying END-seq profiles as
742 normalized read density (reads per million, RPM) for KM12-shWRN cells treated
743 DMSO (NT), doxycycline (shWRN) for 72 hrs, or MUS81-EME1 *in situ*.

744 (b) Zoom-in view of region including exons 6 and 7 of *DCC* gene, containing two (TA)_n
745 repeats displaying END-seq peaks. The highlighted sequences below were extracted
746 from long read sequencing reads in KM12 cells. The (TA)_n repeat in intron 7 is where
747 Vogelstein and colleagues previously detected an insertion.

748

749

750 **Methods**

751 **Cell lines and cell culture**

752 Cell lines containing doxycycline-inducible shWRN cassette (KM12, HCT116, SW837, RPE-
753 1) were generated as previously described². Cell lines were grown in medium supplemented
754 with 10% fetal bovine serum (FBS), penicillin (100 µg/mL) / streptomycin (100 µg/mL), and
755 L- glutamine (Gibco, 292 µg/mL) unless stated otherwise. KM12, SW837, and SW48 were
756 grown in RPMI1640 (Gibco), HCT116 in McCoy's 5A (Gibco), RPE-1 in DMEM/F12
757 (Gibco), SW620 in Leibovitz's L-15 (Gibco), OVK18 and LS180 in MEMα (Gibco)
758 supplemented with 15% FBS. Human primary stomach epithelial cells (HSEC) were obtained
759 from Cell Biologics (H-6039) and cultured with complete human epithelial cell medium
760 (CellBiologics H6621). Independent HCT116 clones were used for whole genome sequencing
761 and continuous long-read sequencing. Cell lines were tested for mycoplasma contamination.

762 **Generation of MMR-deficient HSEC cell lines.**

763 Transduction of HSEC with Cas9 was performed with pLX311-Cas9 (Addgene #118018)
764 followed by blasticidin (4µg/mL) selection. Subsequently, we transduced pXPR_BRD003
765 carrying sgRNAs targeting *MLH1* (target sequence TTTGGCCAGCATAAGCCATG), *MSH2*
766 (CCGGTCGAAAAGGCGCACTG), or *luciferase* (ACAAC TTTACCGACCGCGCC) to
767 generate stable cell lines following puromycin selection (2ug/mL).

768 **Cell viability assay**

769 Cas9-expressing HSEC with sgRNAs targeting *MLH1*, *MSH2*, or *luciferase* were transduced
770 with pXPR_BRD003 harboring the following sgRNAs. sgRNAs targeting *WRN* include
771 sg*WRN2* (ATCCTGTGGAACATACCATG) and sg*WRN3*
772 (GTAGCAGTAAGTGCAACGAT). Two negative controls targeting intergenic sites on
773 chromosome 2 were used: Chr.2-2 sgRNA (GGTGTGCGTATGAAGCAGTG) and Chr.2-4

774 sgRNA (GCAGTGCTAACCTTGCATTG). sgRNA targeting *PolR2D*
775 (AGAGACTGCTGAGGAGTCCA) was used as a pan-essential control. All sgRNAs were
776 inserted in the pXPR_BRD003 lentiviral vector and inserts were verified by Sanger
777 sequencing. Viability was determined using CellTiter-Glo (Promega G7572) 7 days following
778 lentiviral transduction.

779 **Protein depletion and exogenous replication stress**

780 Expression of shWRN was induced by adding 1 µg/mL doxycycline to cell culture medium for
781 indicated time. The following siRNAs were used for experiments: siGenome Human siRNA
782 Smartpool targeting WRN, MUS81, and SLX4, as well as non-targeting control pool (Horizon
783 Discovery). WRN 5'UTR (AAACCCGAGAAGAUAUCCAGUCCAACA) was described in ³
784 and ordered from ThermoFisher. Cells were transfected using RNAiMax Transfection
785 Reagent (ThermoFisher) according to manufacturer's instructions. To induce exogenous
786 replication fork collapse, cells were treated with aphidicolin (0.2 µg/mL, Sigma Aldrich) for
787 24 hrs, and ATR inhibitor AZ20 (10 µM, SelleckChem) was added for the final 8 hrs.

788 **RNA isolation and qRT-PCR**

789 Total RNA was extracted from cells using Trizol Reagent (Invitrogen), and cDNA was made
790 using SuperScript II Reverse Transcriptase (ThermoFisher), according to manufacturer's
791 instructions. RT-PCR was performed using iTaq Universal SyBR Green (BioRad), samples
792 were run and analyzed on a BioRad CFX96 Real-Time PCR detection system. Primer
793 sequences were as follows.

794 MUS81 Fw: 5' GCTGCTCCGAGAGCTACAG 3'

795 MUS81 Rv: 5' CAGGGTTTGCTGGGTCTCTA 3'

796 SLX4 Fw: 5' AGTGTGCTGTGAAGATGGAG 3'

797 SLX4 Rv: 5' CCGTTTCAGACCTCTACTGTG 3'

798 βACTIN Fw: 5' CGTCACCAACTGGGACGACA 3'

799 β ACTIN Rv: 5' CTTCTCGCGGTTGGCCTTGG 3'

800

801 **Metaphase analysis**

802 Cells were arrested at mitosis with 0.04 μ g/ml colcemid (Roche) for 16h and metaphase
803 chromosome spreads were prepared as previously described³⁷.

804 **Western blotting**

805 Cells were lysed in a buffer containing 50 mM Tris-HCl (pH 7.5), 200 mM NaCl, 5% Tween-
806 20, 2% Igepal CA-630, 2 mM PMSF, 50 mM β -glycerophosphate (Merck) and protease
807 inhibitor cocktail tablet (cOmplete Mini, Roche Diagnostics). Equal amounts of lysates were
808 loaded into precast mini-gels (Invitrogen) and resolved by SDS-PAGE. Transfer of proteins
809 onto nitrocellulose membranes and incubation with primary/secondary antibodies were
810 performed according to standard procedures. Visualization of protein bands was achieved by
811 fluorescence imaging on the Odyssey Clx system (LI-COR Biosciences). Antibodies and
812 dilutions used were as follows: Anti-WRN (1:5000, Novus), anti-pKAP1 (Bethyl
813 Laboratories, 1:100), anti-tubulin (1:5,000, sigma), IRDye 680RD Goat anti-Mouse IgG
814 (1:2000, LI-COR Biosciences), IRDye 800CW Goat anti-Rabbit IgG (1:2000, LI-COR
815 Biosciences).

816 **Flow cytometry**

817 For cell cycle analysis, exponentially growing cells were incubated with 10mM (5-ethynyl-
818 2'-deoxyuridine) for 30 min at 37°C and stained using the Click-IT EdU Alexa Fluor 488
819 Flow Cytometry Assay Kit (ThermoFisher) according to the manufacturer's instructions.
820 DNA content was measured by DAPI (4',6-diamidino-2-phenylindole, 0.5 μ g/mL). For
821 pKAP1 staining, fixed and permeabilized cells were washed in PBS+2%FSB, incubated for 1
822 hour with pKAP1 antibody (Bethyl Laboratories, 1:100), washed, stained with anti-rabbit
823 Alexa Fluor 647 antibody (ThermoFisher, 1:2000).

824

825 **Protein purification**

826 Expression and purification of MUS81-_{FLAG}EME1²² and WRN³⁸ were performed as
827 previously described. Human MUS81 (residues 246–551) and double 8xHIS-EME1 (residues
828 246-570) were expressed in a bicistronic expression vector and purified as previously
829 described³⁹, minus the cation exchange chromatography step.

830 **Enzymatic reactions**

831 For enzymatic reactions, END-seq plugs were incubated in the presence of 50 nM of the
832 indicated enzyme for 1.5 hours at 37°C. MUS81/EME1 enzyme reactions were done in buffer
833 consisting of 25 mM Tris-HCl pH 8.0, 30 mM NaCl, 3 mM MgCl₂, 100 ng/μL BSA, 5%
834 glycerol, and 1 mM DTT. WRN enzyme reactions were done in buffer consisting of 30 mM
835 HEPES pH 7.4, 40 mM KCl, 8 mM MgCl₂, 100 ng/μL BSA, 5% glycerol and 2 mM ATP.

836 For *in situ* experiments in which plugs were treated with both WRN and MUS81/EME1, after
837 WRN incubation, agarose plugs were treated with proteinase K in lysis buffer at 50°C for 1
838 hour and then washed in 10 mM Tris-HCl pH 8.0 solution, before incubation with MUS81-
839 EME1.

840 **END-seq**

841 For END-seq, 7-8 million cells were harvested, embedded in 1% agarose plugs, and
842 processed as previously reported^{9,12}.

843 **RPA ChIP-seq**

844 20 million cells were harvested and ChIP-seq was performed as previously reported¹².

845 **PCR-free whole genome sequencing**

846 Genomic DNA was isolated from cultured cells by phenol-chloroform extraction. To provide
847 sufficient mass and maximize library diversity, two identical library preparations were
848 performed in parallel and combined for each sample. DNA (1000 ng) was fragmented on a

849 Covaris S2 sonicator, using 50 µl microTUBE AFA Snap-Caps (Intensity=5, Duty Cycle=5%,
850 Cycles per Burst=200, Time=35 sec). Sequencing libraries were prepared using PCR-free
851 KAPA HyperPrep Kit (Roche) according to the manufacturer's instructions. During ligation,
852 standard adaptors were substituted for xGen Dual Index UMI Adapters (IDT). Following
853 ligation, libraries were size selected using SPRIselect Beads (Beckman Coulter) twice, at
854 concentrations of 0.7x and 0.5x by volume.

855 **Sequencing of END-seq and ChIP seq libraries**

856 END-seq and ChIP-seq libraries were sequenced on the Nextseq 500 or Nextseq 550 platform
857 (Illumina), using 75bp single-end read kits. WGS libraries were sequenced on the Novaseq
858 6000 platform (Illumina), using 250bp paired-end kits.

859 **Continuous long-read sequencing (CLR)**

860 High molecular weight genomic DNA was purified using the MagAttract HMW DNA kit
861 (Qiagen). For CLR library preparation, ≥ 5 ug of high molecular weight genomic DNA (more
862 than 50% of fragments ≥ 40 kb) was sheared to ~ 40 kb using the Megaruptor 3 (Diagenode
863 B06010003), followed by DNA repair and ligation of PacBio adapters using the PacBio
864 SMRTbell Express Template Prep Kit 2.0 (100-938-900). Libraries were then size-selected
865 for >30 kb using a BluePippin instrument with 0.75% agarose cassettes (Sage Science).
866 Following quantification with the Qubit dsDNA High Sensitivity assay (Thermo Q32854),
867 libraries were diluted to 50 pM per SMRT cell, hybridized with PacBio V2 sequencing
868 primer, and bound with SMRT seq polymerase using Sequel II Binding Kit 2.0 (PacBio 101-
869 842-900). CLR sequencing was performed on the Sequel II instrument using 8M SMRT Cells
870 (101-389-001) and Sequel II Sequencing 2.0 Kit (101-820-200), with a 15 hour movie time
871 per SMRT cell. Initial quality filtering, basecalling, and adapter marking was done
872 automatically on board the Sequel II.

873 **Polymerase chain reaction**

874 PCR reactions consisted of 1x KAPA LongRange Buffer, 0.5 U KAPA LongRange HotStart
875 Polymerase (Roche), 1.75 mM MgCl₂, 300 μM dNTPs, 0.4 μM of each primer, and 75 ng of
876 genomic DNA used as template. Samples were initially denatured at 95°C for 3 minutes,
877 underwent 32 cycles of: 1) denaturation at 95°C for 30 seconds and 2) annealing/extension at
878 60°C for 3 minutes, followed by a final extension at 60°C for 10 minutes. PCR products were
879 separated on a 2% agarose gel and visualized by 0.5 μg/mL ethidium bromide staining.
880 PCR primer sequences are as follows, FW/RV indicate forward and reverse primers,
881 respectively:
882 B1 Fw: 5'-GCAACCAGCTGTTTTTGTGA-3'
883 B1 Rv: 5'-GCAATAGTATGCAGCTTGCCC-3'
884 B2 Fw: 5'-TTTGCATCCTGCTTTTCTCATCT-3'
885 B2 Rv: 5'-GAAGAGGTGCCTGGTAGCTG-3'
886 B3 Fw: 5'-TTTGGCTTAGGGGAAGTGTGG-3'
887 B3 Rv: 5'-GTTTTGAGCATGCTGACCTGA-3'
888 B4 Fw: 5'-GCAAGAACCAAATGCTGCAC-3'
889 B4 Rv: 5'-ACTCCTGTTGCTCAGGCAAT-3'
890 B5 Fw: 5'-TGTCCGTGTCTCGAGGAGT-3'
891 B5 Rv: 5'-GTGTCTCCATCCATTGTTCTGC-3'
892 B6 Fw: 5'-TGCTTTCAACCTGCCCAAAC-3'
893 B6 Rv: 5'-GCACTTGAGCCTTGCTGGTA-3'
894 B7 Fw: 5'-TGTGGTTGTCTTCTCCACCC-3'
895 B7 Rv: 5'-AGCTGGGTGTTAAGGGATGAA-3'
896 B8 Fw: 5'-ATGGGATGGCCACACTGAAG-3';
897 B8 Rv: 5'-AACTGCCTTTCACCTGCCT-3';
898 NB1 Fw: 5'-ATAGTCTGTCTCCCGCAGTCT-3';

899 NB1 Rv: 5'-GAGACCGCCGATTAGCATTC-3'

900 NB2 Fw: 5'-TACCTGACAGAACCACTGGC-3'

901 NB2 Rv: 5'-GACAAGGATTCCCCTCCTGC-3'

902 NB3 Fw: 5'-GTGGTGTGGTAAAGGGACCA-3'

903 NB3 Rv: 5'-CCTCTCCCTGTTAAGTCATTACC-3'

904 **Genomic location and reference genome size for TA loci amplified by PCR**

905

906 Genomic location and reference genome size of sites selected for PCR amplification are

907 shown below . Each site amplified contains an uninterrupted stretch of (TA)_n repeats and was

908 amplified with primers localized in flanking areas outside of the (TA)_n repeat. Broken sites

909 B1-B8 were chosen based on the presence of END-seq peaks upon WRN depletion in KM12

910 cells. Sites NB1-NB3 were chosen with similar TA repeat lengths as broken sites but were not

911 broken upon WRN depletion in KM12 cells.

912

913 **Southern blotting**

914 Native (1 % agarose) Southern blot analyses of broken and non broken (TA)_n dinucleotide

915 lengths were carried out using standard Southern blot techniques. Genomic DNA (10 ug)

916 from indicated cell lines were digested with the indicated restriction enzymes: B1: EcoRV;

917 B2: HindIII; NB1: ApaLI; NB3: HindIII. Probes were generated by PCR amplification using

918 genomic DNA from the RPE-1 cell line using the following oligonucleotides:

919 B1 Fw: 5'-GTCAAGGAAAGCCAAGAATTGGAA-3'

920 B1 Rv: 5'-AGAGTTGGAATTTGAACCAAAGC-3'

921 B2 Fw: 5'-GAAACTGCCAAACACAGGGC-3'

922 B2 Rv: 5'-TCTTCAGCACAGGAGCAAGG-3'

923 NB1 Fw: 5'-GTAAGGAGTGTGTGTGGGGG-3'

924 NB1 Rv: 5'-ACTCATGGAGAATAGACTACGACT-3'

925 NB3 Fw: 5'-AGCCATGGTGTGTTTCTGGT-3'

926 NB3 Rv: 5'-GCCGTCTTCGAACCTGTAGA-3'

927 **Quantification and Statistical Analysis**

928 1) Genome Alignment: For Single-end sequencing, NextSeq 500/550 was utilized for
929 single end sequencing using 75bp reads. The sequenced tags were aligned using
930 Bowtie2⁴⁰ with the parameters -N 0 -k 1 -q --local --fast-local. For paired end
931 sequencing: NovoSeq 6000 was utilized for paired-end sequencing using 2x250bp
932 reads. The sequenced tags were aligned using Bowtie2 with the parameters -N 0 –
933 local.

934 2) Peak calling: Peaks were called for single end sequencing experiments using MACS
935 1.4.3⁴¹ using the parameters -p 1e-5 --nolambda --nomodel --keep-dup = all (keep all
936 redundant reads). For each case the experimental sample was doxycycline (shWRN)
937 treated sample compared with DMSO (NT) sample as control. There were some
938 additional filtering criteria as explained below:

939 a) shWRN- ENDseq, RPA ChIPseq

940 The output of the peak calling was filtered by a 20-fold enrichment over background
941 and a minimum size of 1kb. The resulting regions were merged using bedtools⁴²
942 merge -d 1000 function to define the list of regions.

943 b) siWRN- ENDseq, MUS81-EME1- ENDseq

944 The output of the peak calling was filtered by a 20-fold enrichment over background.
945 The resulting regions were merged using bedtools merge -d 1000 function to define
946 the list of regions. For each region defined above the summit in positive and negative
947 strand reads was evaluated using bedtools coverage -d function. The bedtools intersect
948 function was used to filter the regions for which the positive summit was downstream

949 of the negative summit and the distance between the two summits was set to be
950 between 30 to 2,500bp. This defined the final list of peaks for each experiment.

951 3) Motif Finding: For each peak the taller of the positive or negative summit was
952 identified. From this taller summit a 1kb window was generated downstream
953 (upstream) if the summit was on negative (positive) strand. MEME suite⁴³ was used
954 to identify the common sequence motif in these genomic windows.

955 4) For each annotated (TA)_n repeat sequence extracted from the hg19 genome sequence,
956 three parameters were calculated: (a) TA proportion - the cumulative TA length
957 divided by the length of the full annotated (TA)_n repeat; (b) uninterrupted (TA)_n
958 proportion - the length of the longest continuous TA sequence divided by the length of
959 the full annotated (TA)_n repeat; (c) uninterrupted (TA)_n length - the length of the
960 longest continuous TA sequence within the annotated (TA)_n repeat. These three
961 parameters were plotted separately to compare between (TA)_n repeats at broken and
962 non-broken sites of shWRN-KM12 cells. This data was plotted in Fig. 6a-c,
963 respectively. The same was done to calculate the proportion and length of (TA)_n
964 repeats at annotated (TA)_n sequences in the KM12 cell line using long read
965 sequencing data (Extended Data Fig. 7a-c).

966 5) Visualization of Sequencing Data:

967 a) The alignment output sam files were converted and sorted into bam files using
968 samtools⁴⁴ view function. For WGS paired-end data the bam files are directly
969 visualized using IGV⁴⁵ from the Broad Institute. The following steps were
970 carried out for single end sequencing experiments.

971 b) The bamtoBed function in bedtools package was used to convert the bam files
972 to bed files

973 c) The `genomecov` function in `bedtools` package was used to convert the bedfiles
974 into `bedGraph`. This file was converted to `bigwig` file using the UCSC tools
975 package and visualized on the UCSC Genome Browser⁴⁶.

976 6) Expansion detection by ExpansionHunter and exSTRa: Expansions are identified for
977 broken and non-broken (TA)_n from the HCT116 and KM12 whole genome
978 sequencing data by `ExpansionHunter` v3.2.2^{25,47}. The supporting reads for expansions
979 are visualized by `GraphAlignmentViewer`
980 (<https://github.com/Illumina/GraphAlignmentViewer>). Empirical cumulative
981 distribution function for the TA repeat located at hg19 coordinates chr8:106950919-
982 106950985 was generated using `exSTRa` v0.89.0²⁶ and `Bio-STR-exSTRa` v1.1.0 using
983 the default parameter settings.

984 7) Long read sequencing data processing: We aligned CLR reads to the hg19 reference
985 sequence using `minimap2`⁴⁸ (with parameters `-ayYL --MD --eqx -x map-pb`) and
986 called indel/structural variants with the `pbsv` tool using the `--tandem-repeats` argument
987 for increased sensitivity to repeat expansion/contraction variants in accordance with
988 the software's documentation (Pacific Biosciences. `pbsv`
989 <https://github.com/PacificBiosciences/pbsv>). We extracted indel variants passing built-
990 in filters and overlapping with repetitive loci of interest (6,915 “broken” loci, 59,729
991 “non-broken” loci) using the `bedtools closest` utility⁴². We computed coverage over
992 each interval using the `mosdepth`⁴⁹ package and determined the repetitive motifs
993 present using the Tandem Repeat Finder (`trf`) tool⁵⁰, lowering the score threshold in
994 order to see all possible motifs (i.e. with parameters `2 5 7 80 10 0 2000 -l 6 -ngs -`
995 `h`). We quantified the copy number of each motif in the reference haplotype and the
996 alternate haplotype (computed by permuting the reference haplotype with the variants

997 detected in each locus). We report the copy number change for the highest copy
998 number motif per locus.

999 8) Cancer cell line encyclopedia (CCLE) whole genome sequencing analysis: For 321
1000 cancer cell lines from the CCLE2 project⁵¹ that had WGS data we determined the
1001 number of reads at each AT-rich interval using ‘samtools bedcov’⁴⁴. For each interval,
1002 we then defined: fragments per base (fpb) = number of reads covering the interval /
1003 (length of the interval in bases * 2). We then normalized this value across all intervals
1004 within the genome as: fragments per base per million (fpbm) = (fpb / total fpb) * 1
1005 million. This latter quantity (fpbm) was used to determine the read coverage at
1006 broken/non-broken loci. We determined overall read counts at loci for with high AT
1007 repeats

1008 9) MSI deletion analysis: MSI deletions in uterine corpus endometrial carcinoma
1009 (UCEC), colon adenocarcinoma (COAD), and gastric adenocarcinoma (STAD)⁵ were
1010 identified from ICGC (<https://dcc.icgc.org/releases/PCAWG/msi> and
1011 https://dcc.icgc.org/releases/PCAWG/consensus_sv). In total, 396 deletions from 24
1012 MSI samples and 4501 deletions from 93 MSS samples were analyzed. 1000 random
1013 sets were generated using the bedtools⁴² shuffle command in order to estimate
1014 enrichments for different repeat annotations relative to random.

1015 10) Statistical Analysis:

1016 a) Venn Diagrams

1017 For Venn diagrams 10,000 random sites were generated for one of the peak
1018 lists (while keeping others the same) using bedtools shuffle command. The
1019 significance was generated using the pnorm function in R.

1020 b) Boxplots

1021 The statistical significance of the data represented in boxplots are calculated
1022 using the `wilcox.test` function in R.
1023

- 1024 37 Callen, E. *et al.* 53BP1 Enforces Distinct Pre- and Post-resection Blocks on
1025 Homologous Recombination. *Mol Cell* **77**, 26-38 e27,
1026 doi:10.1016/j.molcel.2019.09.024 (2020).
- 1027 38 Palermo, V. *et al.* CDK1 phosphorylates WRN at collapsed replication forks. *Nat*
1028 *Commun* **7**, 12880, doi:10.1038/ncomms12880 (2016).
- 1029 39 Chang, J. H., Kim, J. J., Choi, J. M., Lee, J. H. & Cho, Y. Crystal structure of the
1030 Mus81-Eme1 complex. *Genes Dev* **22**, 1093-1106, doi:10.1101/gad.1618708 (2008).
- 1031 40 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat*
1032 *Methods* **9**, 357-359, doi:10.1038/nmeth.1923 (2012).
- 1033 41 Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137,
1034 doi:10.1186/gb-2008-9-9-r137 (2008).
- 1035 42 Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing
1036 genomic features. *Bioinformatics* **26**, 841-842, doi:10.1093/bioinformatics/btq033
1037 (2010).
- 1038 43 Bailey, T. L. *et al.* MEME SUITE: tools for motif discovery and searching. *Nucleic*
1039 *Acids Res* **37**, W202-208, doi:10.1093/nar/gkp335 (2009).
- 1040 44 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
1041 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).
- 1042 45 Robinson, J. T. *et al.* Integrative genomics viewer. *Nat Biotechnol* **29**, 24-26,
1043 doi:10.1038/nbt.1754 (2011).
- 1044 46 Kent, W. J. *et al.* The human genome browser at UCSC. *Genome Res* **12**, 996-1006,
1045 doi:10.1101/gr.229102 (2002).
- 1046 47 Dolzhenko, E. *et al.* ExpansionHunter: a sequence-graph-based tool to analyze
1047 variation in short tandem repeat regions. *Bioinformatics* **35**, 4754-4756,
1048 doi:10.1093/bioinformatics/btz431 (2019).
- 1049 48 Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**,
1050 3094-3100, doi:10.1093/bioinformatics/bty191 (2018).
- 1051 49 Pedersen, B. S. & Quinlan, A. R. Mosdepth: quick coverage calculation for genomes
1052 and exomes. *Bioinformatics* **34**, 867-868, doi:10.1093/bioinformatics/btx699 (2018).
- 1053 50 Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic*
1054 *Acids Res* **27**, 573-580, doi:10.1093/nar/27.2.573 (1999).
- 1055 51 Ghandi, M. *et al.* Next-generation characterization of the Cancer Cell Line
1056 Encyclopedia. *Nature* **569**, 503-508, doi:10.1038/s41586-019-1186-3 (2019).

1060

1061

1062 **Acknowledgements**

1063 We thank Roackie Awasthi for assistance with Southern Blotting; David Goldstein, Bao Tran
1064 and the CCR Genomics core for sequencing support; Michael Lawrence for computational
1065 assistance. Work in S.C.W.'s laboratory is supported by the Francis Crick Institute (FC10212)
1066 and the European Research Council (ERC-ADG-666400). The Francis Crick Institute receives

1067 core funding from Cancer Research UK, the Medical Research Council, and the Wellcome
1068 Trust. K.F. is the recipient of fellowships from the Benzon Foundation and the Lundbeck
1069 Foundation. The P.J.M laboratory is funded by the MRC MR/R009368/1; A.C-M. is the
1070 recipient of a fellowship from AstraZeneca; E.C.M is supported by the Damon Runyon
1071 Cancer Research Foundation, and E.C.M and A.J.B are supported by a pilot grant from the
1072 Dana-Farber Department of Medical Oncology. The A.N. laboratory is supported by
1073 the Intramural Research Program of the NIH, an Ellison Medical Foundation Senior Scholar
1074 in Aging Award ([AG-SS- 2633-11](#)), the Department of Defense Idea Expansion ([W81XWH-](#)
1075 [15-2-006](#)) and Breakthrough([W81XWH-16-1-599](#)) Awards, the Alex's Lemonade Stand
1076 Foundation Award, and an NIH Intramural FLEX Award.

1077

1078

1079 **Author Contribution**

1080 N.V.W. set up the project, performed END-seq and flow cytometry experiments upon
1081 WRN/MUS81/SLX4 depletion, and performed preliminary analysis of END-seq data; W.J.N.
1082 performed MUS81-EME1 *in situ* END-seq and PCR; A.T. performed END-seq, Southern
1083 Blotting and designed ATR mutant WRN cDNA; E.M.C. generated the inducible shWRN in
1084 KM12 and HCT116, performed and analyzed the HSEC western blot and viability
1085 experiments, long-read sequencing, and analyzed the CCLE/WRN dependency data; E.C.
1086 performed ATRi END-seq experiments, Western blotting, and metaphase analysis. V.T.
1087 performed RPA ChIP-seq; K.F. performed the HSEC and long-read sequencing experiments;
1088 N.W. performed Western blotting and helped generate WRN^{3A} and WRN^{6A} cells; J.N. and
1089 J.K. analyzed the CCLE/WRN dependency data; S.S. analyzed END-seq, RPA ChIP-seq
1090 experiments; W.W. analyzed WGS, PacBio coverage across repeats, deletion breakpoints in
1091 MSI cancers, and performed quantitative modeling; F.B. analyzed nucleotide composition of

1092 broken vs. non broken repeats and replication timing; E.D. performed ExpansionHunter and
1093 exSTRa bioinformatic analysis; M.E. supervised computational work; K.G., Y.H., A.A.B.,
1094 J.T.S., and N.K. analyzed the data and designed bioinformatic pipelines; R.W. prepared WGS
1095 libraries; A.C.M. and K.F. provided recombinant MUS81-EME1; J.S. provided recombinant
1096 WRN; B.E.H. provided advice about PCR across repeats; K.U. provided advice about repeat
1097 expansion biology; C.H.F. provided advice about secondary structure biology; R.M.B.
1098 provided advice about WRN helicase; S.W. provided advice about structure specific nucleases
1099 and recombination intermediates, P.J.M. helped design *in situ* experiments with recombinant
1100 proteins, P.S.M., provided advice on WGS experiments and analyses; A.J.B. and A.N.
1101 supervised the project; N.V.W., W.J.N, A.T., E.M.C., A.J.B. and A.N. wrote the manuscript
1102 with comments from the authors.

1103

1104 **Author Information**

1105 These authors contributed equally: Niek van Wietmarschen, Sriram Sridharan, William J.

1106 Nathan, Anthony Tubbs, Edmond M. Chan

1107 Correspondence and requests from materials: Andre Nussenzweig

1108

1109 **Competing Interest**

1110 Adam J. Bass has research support from Bayer, Merck and Novartis.

1111

1112 **Data Availability**

1113 END-seq, ChIP-seq, WGS and Pacbio CLR data have been deposited in the Gene Expression

1114 Omnibus (GEO) database under the accession number GSE149709. Source data are provided

1115 with this paper.







