# Edinburgh Research Explorer

# Signatures of TOP1 transcription-associated mutagenesis in cancer and germline

# Signatures of TOP1 transcription-associated mutagenesis in cancer and germline

Martin A.M. Reijns[1,14*], David A. Parry[1,14], Thomas C. Williams[1,2,14], Ferran Nadeu[3,4], Rebecca L. Hindshaw[5], Diana O. Rios Szwed[1], Michael D. Nicholson[6], Paula Carroll[1], Shelagh Boyle[7], Romina Royo[8], Alex J. Cornish[9], Hang Xiang[10], Kate Ridout[11], The Genomics England Research Consortium[+], Colorectal Cancer Domain UK 100,000 Genomes Project[+], Anna Schuh[11], Konrad Aden[10], Claire Palles[5], Elias Campo[3,4,12,13], Tatjana Stankovic[5], Martin S. Taylor[2*], Andrew P. Jackson[1*]

[1] Disease Mechanisms, MRC Human Genetics Unit, Institute of Genetics and Cancer, The University of Edinburgh, Edinburgh, UK

[2] Biomedical Genomics, MRC Human Genetics Unit, Institute of Genetics and Cancer, The University of Edinburgh, Edinburgh, UK

[3] Institut d'Investigacions Biomèdiques August Pi i Sunyer (IDIBAPS), Barcelona, Spain

[4] Centro de Investigación Biomédica en Red de Cáncer (CIBERONC), Madrid, Spain

[5] Institute of Cancer and Genomic Sciences, University of Birmingham, Edgbaston, UK

[6] Cancer Research UK Edinburgh Centre, Institute of Genetics and Cancer, The University of Edinburgh, Edinburgh, UK

[7] Genome Regulation, MRC Human Genetics Unit, Institute of Genetics and Cancer, The University of Edinburgh, Edinburgh, UK

[8] Barcelona Supercomputing Center (BSC), Barcelona, Spain

[9] The Institute of Cancer Research, London, UK

[10] Institute of Clinical Molecular Biology, Christian-Albrechts-University and University Hospital Schleswig-Holstein, Kiel, Germany

[11] Department of Oncology, University of Oxford, Oxford, UK

[12] Hospital Clínic of Barcelona, Barcelona, Spain

[13] Departament de Fonaments Clínics, Universitat de Barcelona, Barcelona, Spain

[14] These authors contributed equally: Martin A.M. Reijns, David A. Parry, Thomas C. Williams

[+] A list of authors and their affiliations appears at the end of the paper.

[*] Correspondence to MAMR, MST, APJ

## Abstract

The mutational landscape is shaped by many processes, with genic regions vulnerable to mutation but preferentially protected by transcription-coupled repair[1]. In microbes, transcription has been demonstrated to be mutagenic[2,3]; however, the impact of transcription-associated mutagenesis remains to be established in higher eukaryotes[4]. Here we show that ID4, an indel cancer signature of unknown aetiology[5] characterised by short deletions (2-5 bp), is due to a transcription-associated mutagenesis process. We demonstrate defective ribonucleotide excision repair in mammals to be associated with the ID4 signature, with mutations occurring at a TNT sequence motif, implicating Topoisomerase 1 (TOP1) activity at sites of genome-embedded ribonucleotides as a mechanistic basis. Such TOP1-mediated deletions occur somatically in cancer, and the ID-TOP1 signature is also found in physiological settings, contributing to genic *de novo* indel mutations in the germline. Hence, while topoisomerases protect against genome instability by releasing topological stress[6], their activity may also be an important source of mutations in the human genome.

## Introduction

45    Eukaryotic cells employ many strategies to ensure integrity of their genomes, with high-fidelity DNA

46    replication[7] and DNA repair processes countering exogenous and endogenous DNA lesions[8]. The

47    process of transcription targets DNA repair machinery to expressed genes, preferentially reducing

48    their mutation rate following DNA damage[1]. Despite this targeted repair, in micro-organisms the

49    process of transcription itself is mutagenic; a phenomenon referred to as transcription-associated

50    mutagenesis (TAM)[2,3]. In yeast, Topoisomerase 1 (Top1) activity is a major source of TAM and results

51    in a distinctive transcription-dependent signature of 2-5 bp deletions at tandem repeat sequences[9-

52    11]. Genome-embedded ribonucleotides have been established as a cause of Top1-TAM deletions in

53    yeast[12]. Such ribonucleotides are frequently incorporated by DNA polymerases during replication,

54    and represent the most prevalent aberrant nucleotides in the eukaryotic genome[13,14]. These

55    genome-embedded ribonucleotides are normally removed by ribonucleotide excision repair (RER), a

56    process initiated by the heterotrimeric Ribonuclease H2 enzyme[15]. However, when Top1 cleaves at

57    embedded ribonucleotides instead of RNase H2 this can result in small deletions[16,17].

58    In the last decade, widespread use of genome sequencing has enabled unbiased sampling of human

59    mutations, substantially advancing understanding of mutagenesis in the germline[18] and in

60    neoplasia[19]. Multiple mutational processes act during cancer evolution, and mathematical methods

61    have been developed to define signatures that may correspond to individual mutagenic

62    mechanisms, through decomposition of tumour mutational profiles[19]. This has successfully defined

63    cell-intrinsic, environmental and treatment-related origins for many base substitution signatures in

64    cancer[20-22]. However, the origin of a substantial number of signatures remains unknown, and some

65    may be artefactual. Recently, cancer signature analysis has been extended to indels[5], small (1-49 bp)

66    insertions and deletions. Such indels are an important class of mutations that contribute

67    substantially to disease-causing germline variants (>20%) and human variation[23] .

68    Here we investigate an indel signature of unknown cause, ID4. We show experimentally that ID4

69    deletions are increased in RNase H2 deficient cell lines and cancers and delineate a human TOP1-

70    mediated TAM signature (ID-TOP1) relevant to both somatic and germline mutagenesis.

71

## 72 Results

### 73 ID4, a distinct cancer indel signature

74    The ID4 cancer signature, as categorised by COSMIC[24], comprises 2-5 bp deletions, often with loss of

75    a single repeat unit at short repeat sequences[5]. Most commonly these occur where the deleted

76    sequence is repeated one, two or three times in tandem (Fig. 1a). Hereafter we use the term SSTRs

77    (short-short tandem repeats) to distinguish such short tandem repeats (STRs) with less than 5

78    repeats (i.e. <6 repeat units) from microsatellite STRs with many repeats. In addition to these SSTR

79    deletions, ID4 is characterised by small deletions at sequences with microhomology (MH), in

80    particular 2 bp deletions with single nucleotide microhomology (SNMH). Both features are distinct

81    from cancer deletion signatures resulting from other well-recognised mechanisms like replication

82    slippage and non-homologous/microhomology-mediated end joining (NHEJ/MMEJ) (Extended Data

83    Fig. 1a,b). In support of a distinct aetiology, SSTR and SNMH deletions are not apparent in cancer

84    associated with homologous recombination (HR) or mismatch repair (MMR) deficiency, which are

85    expected to have higher levels of MMEJ and replication slippage mutagenesis, respectively

86    (Extended Data Fig. 1c,d).

87

### 88 ID4 resembles a yeast mutation signature

89    Noting similarities to a Top1-induced transcription-associated mutagenesis (Top1-TAM) in

90    *Saccharomyces cerevisiae*, we re-analysed published genome-wide mutation accumulation

91    experiments performed with *rnh201Δ pol2-M644G* yeast[25]. This strain is particularly susceptible to

92     Top1-TAM as it accumulates genome-embedded ribonucleotides at high levels due to RNase H2/RER

93     deficiency and enhanced ribonucleotide incorporation by a steric-gate mutation in the catalytic site

94     of the replicative polymerase Pol ε[26]. Similarities to the ID4 signature were apparent with a

95     comparable pattern of small deletions at SSTRs, although mutational events at sites of SNMH were

96     not evident in the yeast data (Fig. 1b). As over one million ribonucleotides are incorporated by DNA

97     polymerases per replicating mouse cell[14], we reasoned that genome-embedded ribonucleotides

98     might cause similar mutational events in mammalian cells. To experimentally address whether TAM

99     contributes to indel formation in human RER-deficient cells, we developed a novel reporter to

100    enable sensitive and specific detection of mutational events arising from TOP1 activity in both yeast

101    and mammals.

102

103    **Top1-dependent deletions in yeast**

104    Mutation rates are routinely measured in *S. cerevisiae* using well characterised but species-specific

105    selectable markers (LYS2, URA3, CAN1). Therefore, to establish a system that could be transferred

106    between yeast and mammalian cells, we used an approach inspired by the Traffic Light reporter

107    assay[27], incorporating both positive and negative selection cassettes in a single transcriptional unit

108    (Fig. 1c). The hygromycin resistance gene (HygroR) served both as the mutational target and

109    negative selection marker. Indels causing a 2 bp frameshift within HygroR, including 2 bp deletions,

110    result in translation of an otherwise out-of-frame P2A self-cleaving peptide and the neomycin

111    resistance (NeoR) gene, permitting positive selection of mutated colonies with neomycin (Extended

112    Data Fig. 2a). To enrich the target for 2 bp tandem repeats, *in silico* re-design incorporated

113    synonymous substitutions such that SSTRs accounted for >50% of the HygroR open reading frame.

114    For validation, the reporter was inserted into the *S. cerevisiae* genome and fluctuation assays used

115    to assess mutation rates in strains deficient for RER and/or Top1. A 37-fold increase in mutation rate

116    was seen for the *rnh201Δ* (RNase H2 null) strain compared to wild type (Fig. 1d), with a mutation

117    rate of 6.1x10$^{-9}$ per bp per generation (95% CI, 5.4-6.9x10$^{-9}$), whereas the increased mutation rate

118    was abolished in the *rnh201Δ top1Δ* double mutant strain, in keeping with Top1-dependent

119    mutagenesis at genome-embedded ribonucleotides[12,28]. Notably, there was a 10-fold decrease in the

120    mutation rate for *top1Δ* compared to the wild-type strain, and a 35-fold decrease in 2 bp SSTR

121    deletions (Extended Data Fig. 2b), consistent with previous reports[10,11]. In addition, the observed

122    mutational spectrum was most similar for wild-type and *rnh201Δ* strains, but substantially different

123    compared to *top1Δ* and *rnh201Δ top1Δ* strains (Fig. 1e; Extended Data Fig. 2c-f). Taken together, we

124    conclude that the same Top1-mediated mutations occur, albeit at different frequencies, in wild-type

125    cells when RER is functional and in RNase H2 deficient strains when elevated levels of

126    ribonucleotides are present in the genome.

127

**TOP1-mediated mutations in human cells**

129    Having validated the reporter in yeast, the same 2 bp repeat-enriched HygroR sequence was used to

130    address whether TOP1-mediated mutagenesis at embedded ribonucleotides is conserved in human

131    cells (Fig. 2; Extended Data Fig. 2g). NeoR was replaced by the puromycin resistance (PuroR) gene,

132    with reporter expression driven from the mammalian ubiquitous CAG promoter, permitting rapid

133    antibiotic selection in mammalian cells. This modified reporter was inserted at the *AAVS1* safe

134    harbour locus in HeLa cells (Fig. 2a; Extended Data Fig. 3a-e). CRISPR/Cas9-mediated genome

135    editing, targeting the catalytic site of *RNASEH2A*, was then used to generate two independent

136    knockout (KO) reporter clones, alongside a control clone that had also been taken through editing

137    and clonal selection steps (Fig 2b,c; Extended Data Fig. 3). The control clone retained RNase H2

138    activity, while there was complete loss of cellular RNase H2 activity in KO clones, accompanied by

139    high levels of ribonucleotides in genomic DNA (Fig. 2b,c; Extended Data Fig. 3f,g).

140    In fluctuation assays, RNase H2 null clones demonstrated a significant 3.1-fold increase in mutation

141    rate (Fig. 2d) and 5.2-fold more 2 bp SSTR deletions (Extended Data Fig. 3h) compared with RNase

142    H2 proficient cells (RNASEH2A+), consistent with conservation of TOP1-directed mutagenesis in

143    human cells. As in yeast (Fig. 1e), the overall mutational profile of reporter mutations was similar

144    between RNase H2 proficient and null HeLa cells (cosine similarity 0.89, $P < 10^{-4}$), predominantly

145    comprised of 2 bp SSTR deletions (Fig. 2e).

146    The mutation rate for RNase H2 null HeLa cells ($8.0 \times 10^{-9}$ per bp per generation; 95% CI, $6.7\text{-}9.5 \times 10^{-9}$)

147    was similar to that seen for *rnh201Δ* yeast (Fig. 1d), whereas the rate was substantially higher for

148    RNASEH2A+ control cells compared to wild-type yeast. However, the increased mutation rate in

149    RNase H2 null HeLa cells likely underestimates the true impact of RER deficiency in human cells, as

150    despite the control RNASEH2A+ HeLa reporter cells retaining protein expression (Fig. 2b), the clone

151    had also acquired mutations at the CRISPR editing site that reduced enzymatic activity (Fig 2c),

152    causing a moderate increase in genomic ribonucleotide content (Extended Data Fig. 3f,g).

153    To confirm these findings we used a complementary approach to establish the relevance of such

154    mutational events genome-wide, performing mutation accumulation experiments using hTERT-RPE1

155    ($TP53^{-/-}$) diploid cell lines.  Ancestral populations for RNase H2 wild-type and null cells (RNASEH2A-

156    KO or RNASEH2B-KO; Extended Data Fig. 4a-d) were established after initial single cell sorting, and

157    clones then grown for approximately 100 generations. Single cell sorting was performed every 25

158    generations, creating bottlenecks to capture accumulating mutations (Fig. 3a). Combined variant

159    calling on whole genome sequencing (WGS) from paired ancestral and endpoint cultures identified a

160    total of 1,698 acquired high confidence indel mutations, captured by at least 3 out of 4 variant

161    callers. Consistent with TOP1-mediated mutagenesis, among all indel categories, only 2-5 bp

162    deletions were found to be substantially (7.4-fold) and significantly enriched in RNase H2 null RPE1

163    cells compared to wild-type (Fig. 3b; Extended Data Fig. 4e,f), with an estimated rate of $1.1 \times 10^{-10}$ 2-5

164    bp deletions per generation per bp for KO and $1.4 \times 10^{-11}$ for WT. Of these deletions in RNase H2 null

165    cells, 82% were 2 bp deletions, of which 48% were at SSTRs (Extended Data Fig. 4g). Furthermore,

166    signature decomposition using SigProfilerExtractor[5] reported a 21% ID4 contribution in RNase H2

167    null cells, that increased to 61% when subtraction of background mutation patterns was performed

168    to identify RER-deficiency specific mutation signatures (Fig. 3c,d; Extended Data Fig. 5). The ID4

169    signature was substantially enriched in transcribed genomic regions (Extended Data Fig. 5e). ID5, a

170    clock-like signature[5], was also enriched in KO cells, likely due to slower growth and longer culture

171    time needed to achieve the same number of doublings for RNase H2 null cells[14].

172

173    **MH deletions specific to mammalian TOP1**

174    Small deletions at sequences with microhomology are an additional feature of ID4 (Fig. 1a), not

175    observed in *rnh201Δ pol2-M644G* yeast (Fig. 1b). However, consistent with a ribonucleotide-induced

176    mutational origin in mammalian cells, they are observed frequently in RNase H2-deficient RPE1 cells,

177    in which SNMH sites account for 31% of 2 bp deletions, indicating that in humans they share the

178    same aetiology as those occurring at SSTRs. Taken together, our reporter and mutation

179    accumulation experiments demonstrate that genome-embedded ribonucleotides cause a similar

180    mutational signature in yeast and mammalian cells. Therefore Topoisomerase 1-mediated

181    mutagenesis likely also occurs in humans and is associated with 2-5 bp deletions at SSTR and SNMH

182    sequences.

183

184    **ID4 mutations in a murine cancer model**

185    To determine if TOP1-induced mutations resulting in the ID4 signature can be detected *in vivo*, we

186    next studied an RER-deficient murine cancer model in which Villin-Cre conditional deletion of

187    *Rnaseh2b* and *Tp53* results in intestinal malignancy[29]. Whole genome sequencing of paired tumour-

188    normal tissue samples from 6 mice, identified a total of 989 high-confidence tumour-specific somatic

189    indels. Analysis of the resulting mutational signature established that ID4 substantially contributed in

190    all tumours (Fig. 4a,b and Extended Data Fig. 6a), accounting for 32% of acquired indels. Consistent

191 with a transcription-associated process, the ID4 signature was again most evident in transcribed

192 genomic regions (Fig. 4b). Commonly occurring cancer signatures[5] ID1, ID2 and ID5 were also seen,

193 in line with expectations of multiple mutational processes active in neoplasia.

194 The observed ID4 mutation spectrum corresponded closely to that seen in the RPE1 mutation

195 accumulation experiment: 28% of indels were at 2-5 bp deletions, of which the majority were again

196 2 bp deletions (82%) predominantly at SSTRs (51%) and sites of SNMH (34%) (Extended Data Fig.

197 6b,c). This is consistent with the occurrence of TOP1-induced somatic mutations at genome-

198 embedded ribonucleotides *in vivo*, conserved across different tissue and cellular contexts, and

199 shows that it can be detected in a cancer setting.

200

201 **A sequence motif for ID4 mutations**

202 While COSMIC defines the ID4 signature on the basis of indel size and repeat/microhomology

203 context (Fig. 1a), the number of indels in the murine RER-deficient tumour model permitted us to

204 further investigate the characteristics of mammalian Topoisomerase 1-induced mutations. We

205 focussed our analysis on 2 bp deletions, as such events represented 81% of >1 bp deletions in the

206 context of tandem repeats and 85% of deletions in sequences with microhomology.

207 First, we classified all 2 bp deletions at STR/SNMH sequences into 6 non-redundant dinucleotide

208 classes, grouping together complementary sequences (Fig. 4c). We noted that the deleted

209 sequences substantially deviated from genome-wide frequencies, with a complete absence of CC/GG

210 and CG/GC deletions, as well as an overrepresentation of the CT category (containing CT, TC, GA and

211 AG deletions). All observed deletions therefore included at least one thymidine (T), which

212 functionally could be accounted for by the very strong preference of mammalian Topoisomerase 1

213 to cleave at a phosphodiester bond with a T immediately upstream[30].

214    Next, to investigate the wider sequence context, we aligned sequences containing all 228 two bp

215    deletions (Extended Data Fig. 6e), which indicated that deletions preferentially occur when T

216    nucleotides are spaced at a 2-base interval. Indeed, this TNT motif was present in 100% of SNMH

217    (*n*=77) and STR sites (*n*=124), providing a common unifying sequence context for both deletion types

218    (Fig. 4d), a finding replicated in both our RPE1 (Extended Data Fig. 6e) and yeast datasets (Extended

219    Data Fig. 7). We found TNT to be substantially overrepresented at deletions sites compared to the

220    genome-wide null expectation. Furthermore, while the TNT motif is common at tandem repeat

221    sequences, 2 bp deletions at this motif are still significantly enriched when considering the

222    occurrence of 2 bp STR and SNMH sequences in mouse and human genomes (Fig. 4e; Extended Data

223    Fig. 6f), and STR sequences in the yeast genome (Extended Data Fig. 7).

224    To account for thymidines spaced at a 2-base interval and the occurrence of mammalian SNMH

225    deletions, we developed a revised model based on the established strand realignment model for

226    yeast Top1-mediated mutagenesis[12,16,17]. In this "TNT model", TOP1 cleaves preferentially 3' of an

227    embedded ribouridine, with nucleophilic attack by the 2'-OH of the ribose ring resulting in TOP1

228    release and formation of a non-ligatable nick with a terminal 2',3'-cyclic phosphate (Fig. 4f, i-iii). This

229    then provides a substrate for TOP1 cleavage 2 bp or more upstream[17], preferentially at a

230    thymidine[30]. When this second cleavage event happens at a base identical to that of the first cleaved

231    nucleotide, an event more likely at STR and microhomology sequences, strand realignment can then

232    occur, resulting in a nick permissive to religation and TOP1cc reversal (Fig. 4f, iv-vi). An alternative

233    mechanism of sequential Top1 cleavage, in which double-strand breaks occur due to nicking of

234    opposite strands[31] could not be reconciled with our TNT model, but may account for deletions

235    occurring at non-STR/SNHM sites. Within the TNT motif, deletions were most common at CT and GT

236    dinucleotides in both mammals and yeast (Fig. 4c; Extended Data Fig. 6 and 7b,e), which may be

237    explained, at least in part, by preferential incorporation of ribouridine at CT and GT dinucleotides

238    (Extended Data Fig. 7f and [32]).

239

240      Implicating TOP1-TAM as the cause of the ID4 signature permits us to include additional features in

241      the definition of this COSMIC signature, namely preference for a TNT sequence motif at 2 bp

242      deletion sites and enrichment in transcribed genes. Hereafter, we refer to this extended definition as

243      ID-TOP1. To establish the relevance of the ID-TOP1 signature for human disease and genetic

244      variation, we next examined publicly available datasets.

245

246      **ID-TOP1 in human cancer**

247      *RNASEH2B* is frequently deleted in human cancer, in particular in chronic lymphocytic leukaemia

248      (CLL) given its proximity to a tumour suppressor locus, the *DLEU2-mir-15-16* microRNA cluster[33].

249      Such RNase H2 deficient human cancers should therefore be enriched for the ID4/ID-TOP1 signature.

250      We analysed whole genome sequencing data for 348 CLL patients from two independent

251      cohorts[34,35], stratified on *RNASEH2B* deletion status. Somatic variant calling identified a significant

252      increase in 2-5 bp deletions in RNase H2 null tumours (Fig. 5a), while other indels were equally

253      represented across wild-type, heterozygous and null categories (Extended Data Fig. 8a). Of the 2-5

254      bp deletions in tumours with biallelic *RNASEH2B* loss more than half (57%) were 2 bp deletions,

255      which were predominantly at STR and SNMH sequences and substantially enriched for the TNT motif

256      (Extended Data Fig. 8b,c), consistent with the ID-TOP1 mutational signature. Furthermore,

257      mutational signature decomposition for RNase H2 null CLL cases confirmed the presence of the ID4

258      signature, most apparent in genic regions (Extended Data Fig. 8d). We therefore conclude that the

259      ID-TOP1 signature is present in human cancer and enriched in tumours that are RNase H2 deficient.

260      Topoisomerase 1 also causes mutations in RER proficient cells (Fig. 1d-f and [10,11]), and therefore is

261      likely to cause mutations in other cancers, with deletions expected to occur most frequently in highly

262      transcribed genes[4]. Accordingly, analysis of WGS data across cancer types (ICGC/PCAWG)

263      demonstrated that the 2-5 bp deletion rate correlates with expression levels of ubiquitously

264    expressed genes (Pearson's r = 0.86; *P* = 0.0014), with deletions markedly elevated in the most highly

265    expressed genes (Fig. 5b), in line with previous reports of such deletions in certain cancer genes[36,37].

266    Examination of 2 bp deletions (42% of 2-5 bp deletions) across cancer types also demonstrated them

267    to be predominantly in STR and SNMH contexts (Extended Data Fig. 8f) and enriched for the TNT

268    sequence motif (Fig. 5c). Furthermore, using a dataset of TOP1 cleavage events captured by TOP1-

269    seq[38], we found 2-5 bp deletions increase in frequency with TOP1 enzymatic activity, with such

270    deletions more prevalent in regions of high TOP1 activity (Fig. 5d). Likewise, TOP1-ID deletion rates

271    also corresponded to TOP1 activity and transcription level, in contrast to all other deletions

272    (Extended Data Fig. 8g,h). Taken together, this establishes a significant role for TOP1-mediated

273    mutagenesis in the generation of somatic deletions.

274    To further explore the role of transcription in deletion mutagenesis of cancer genomes, we identified

275    genes that are highly expressed, but only in certain tissues. For prostate adenocarcinoma, highly

276    expressed prostate-restricted genes were significantly enriched for 2-5 bp deletion mutations

277    compared to other genes in this cancer type, as well as the same genes in other cancers (two-tailed

278    Fisher's exact test, OR 3.5, *P* = 2.5x10$^{-8}$ after Bonferroni correction; Extended Data Fig. 8i).

279    Importantly, this analysis considers the same sets of genes between cancer types and therefore rules

280    out sequence composition biases as a confounder for elevated ID-TOP1 mutagenesis in highly

281    expressed genes. Extending this approach in an all-versus-all comparison between 8 cancer types

282    and 17 tissues demonstrated specificity between high expression in a tissue of origin and enrichment

283    for 2-5 bp deletions (Fig. 5e). These results extend the relevance of TOP1-mediated mutagenesis to

284    other cancers, confirms the ID-TOP1 mutational signature to be transcription-associated, and

285    supports the occurrence of TAM in humans.

286

287    **TOP1-mediated deletions in the germline**

288    TOP1 is ubiquitously expressed, so we reasoned that it could cause germline as well as somatic

289    mutations. To investigate this possibility we examined mutations from parent-child trio WGS studies

290    in the Gene4Denovo database[39]. *De novo* mutations identified in such datasets represent germline

291    events, as they occur in germ cells or during early embryonic cell divisions. Strikingly, 2-5 bp

292    deletions were the largest category identified, accounting for 33% of the 40,936 *de novo* indels (Fig.

293    5f), and the majority of these were compatible with the ID-TOP1 signature. Analysis of 2 bp deletions

294    (41% of 2-5 bp deletions) demonstrated that most occur at SSTR or MH sites (Extended Data Fig.

295    9a,b), with enrichment of the TNT sequence motif both genome wide and in the context of

296    STR/SNMH sites (Fig. 5g; Extended Data Fig. 9c). Likewise for 3 and 4 bp deletions, TNNT and TNNNT

297    motifs respectively were significantly overrepresented compared to genome-wide expectation

298    (Extended Data Fig. 9d), in support of sequential TOP1 cleavage and strand realignment as the

299    underlying cause. Consistent with TOP1-TAM aetiology, 2-5 bp deletion and ID-TOP1 deletion

300    frequency correlated with transcript expression in male germ cells (Fig. 5h and Extended Data Fig.

301    9e). We therefore conclude that the ID-TOP1 mutational signature also occurs in the human

302    germline, implicating TOP1-induced strand realignment mutagenesis as an important mutational

303    process in mammalian cells.

304

## Discussion

306    Here we establish a biological basis for the ID4 cancer signature[5], experimentally demonstrating it to

307    occur in RNase H2 deficient cells both *in vitro* and *in vivo*. This implicates TOP1-mediated cleavage at

308    genome-embedded ribonucleotides as its cause. TOP1 is cell-essential in mammals, and it is

309    therefore not possible to similarly confirm a genetic dependency on TOP1 in human cells, as has

310    been done in yeast[12]. However, conservation of this mechanism across eukaryotes is supported by

311    us finding a Topoisomerase 1 dependent TNT deletion motif present in both yeast and humans, and

312    demonstrating that deletion frequency is dependent on human TOP1 activity levels. Previously

313     published work also provides evidence for TOP1-mutagenesis at ribonucleotide sites in humans. The

314     reversible transesterification reaction of Type 1 Topoisomerases is conserved from yeast to

315     humans[6], and human TOP1 has site-specific activity for ribonucleotides[40], causing DNA breaks in

316     mammalian RNase H2 deficient cells[33]. Furthermore, generation of 2 bp deletions through sequential

317     TOP1 cleavage at embedded ribonucleotides has been biochemically reconstituted with both human

318     and yeast enzymes[17,31].

319     We define additional features of this ID-TOP1 mutational signature, with deletions strongly enriched

320     at TNT motifs in both yeast and mammals, a sequence context specific to Topoisomerase 1

321     (Extended Data Fig. 7g,h), and deletions most frequent in highly transcribed regions. Consequently,

322     we show that a transcription-associated mutagenesis process first identified in yeast[10,12,41] is relevant

323     to higher eukaryotes, establishing TOP1-induced mutagenesis as an important process for human

324     variation and disease. Additional signatures associated with topoisomerases or indeed RNase H2

325     may be identified in future, particularly given that ID17 has been recently been linked to TOP2A$^{K743N}$

326     cancers[42].

327     The substantial contribution of ID-TOP1 deletions to germline mutagenesis has particular

328     significance given that such deletions will be disproportionately disruptive, particularly in transcribed

329     regions. Notably, such deletions occur in the context of normal RER function, consistent with the

330     mutagenic potential of Topoisomerase 1 in physiological wild-type settings (Fig. 1d and [10,11]). Given

331     that genome-embedded ribonucleotides are the most common endogenous lesion in replicating

332     mammalian cells[14], they are the most likely sites of TOP1-TAM mutagenesis, where TOP1 could

333     cleave before their removal by RNase H2-dependent RER. Processing of TOP1cc may be an

334     alternative, less frequent source of 2-5 bp deletions[41], but we did not detect ID4 in Topoisomerase 1

335     inhibitor treated cancers (Extended Data Fig. 8j). TOP1 canonical function is to relieve DNA

336     topological stress, arising during both transcription and replication[6] (Extended Data Fig. 10). Hence

337     TOP1-mediated deletions are not restricted to transcribed regions of the genome, with deletions

338    also evident in non-genic regions with high TOP1 activity (Extended Data Fig. 8k). However, overall,

339    enhanced TOP1 activity associated with transcription accounts for more frequent mutagenesis

340    within genes.

341    Given the essential nature of topoisomerase activity across tissues and cell states, TOP1-mediated

342    mutagenesis is likely to occur in many contexts. Frequent TOP1-mediated human germline

343    mutations (Fig. 5i-k) and identification of ID4 at early embryonic stages[43] suggest developmental

344    vulnerability to TOP1-TAM. In addition, 2-5 bp somatic deletions at SSTRs are also observed at high

345    frequency in non-dividing neurons[36], and ID4 has been identified in multiple tumour types[5]. As such,

346    this mutational process is likely to be significant not only in cancers with RER deficiency, but also

347    those with high TOP1 activity and tumours with defects in relevant repair mechanisms, such as

348    enzymes that process TOP1cc[6] or non-ligatable TOP1-induced nicks[44-46]. In addition, alternative RER

349    pathways may exist[47] that could reduce TOP1-mutagenesis. The ID-TOP1 signature may provide a

350    useful biomarker with potential future diagnostic and therapeutic utility[48], for instance as an

351    indicator of TOP1-induced genome instability targetable by PARP or ATR inhibitors[33,49].

352    In conclusion, alongside its essential role in relieving DNA torsional stress, TOP1 also drives

353    mutagenesis in somatic and germline contexts, relevant to neoplasia, inherited disease and human

354    variation.

355

## References

357    1    Hanawalt, P. C. & Spivak, G. Transcription-coupled DNA repair: two decades of progress and
358         surprises. *Nat Rev Mol Cell Biol* **9**, 958-970, doi:10.1038/nrm2549 (2008).
359    2    Datta, A. & Jinks-Robertson, S. Association of increased spontaneous mutation rates with
360         high levels of transcription in yeast. *Science* **268**, 1616-1619, doi:10.1126/science.7777859
361         (1995).
362    3    Herman, R. K. & Dworkin, N. B. Effect of gene induction on the rate of mutagenesis by ICR-
363         191 in Escherichia coli. *J Bacteriol* **106**, 543-550, doi:10.1128/JB.106.2.543-550.1971 (1971).
364    4    Jinks-Robertson, S. & Bhagwat, A. S. Transcription-associated mutagenesis. *Annu Rev Genet*
365         **48**, 341-359, doi:10.1146/annurev-genet-120213-092015 (2014).

| 366 | 5 | Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer. *Nature* **578**, |
| 367 | | 94-101, doi:10.1038/s41586-020-1943-3 (2020). |
| 368 | 6 | Pommier, Y., Sun, Y., Huang, S. N. & Nitiss, J. L. Roles of eukaryotic topoisomerases in |
| 369 | | transcription, replication and genomic stability. *Nat Rev Mol Cell Biol* **17**, 703-721, |
| 370 | | doi:10.1038/nrm.2016.111 (2016). |
| 371 | 7 | Kunkel, T. A. Evolving views of DNA replication (in)fidelity. *Cold Spring Harb Symp Quant Biol* |
| 372 | | **74**, 91-101, doi:10.1101/sqb.2009.74.027 (2009). |
| 373 | 8 | Ciccia, A. & Elledge, S. J. The DNA damage response: making it safe to play with knives. *Mol* |
| 374 | | *Cell* **40**, 179-204, doi:10.1016/j.molcel.2010.09.019 (2010). |
| 375 | 9 | Lippert, M. J., Freedman, J. A., Barber, M. A. & Jinks-Robertson, S. Identification of a |
| 376 | | distinctive mutation spectrum associated with high levels of transcription in yeast. *Mol Cell* |
| 377 | | *Biol* **24**, 4801-4809, doi:10.1128/MCB.24.11.4801-4809.2004 (2004). |
| 378 | 10 | Lippert, M. J. *et al.* Role for topoisomerase 1 in transcription-associated mutagenesis in |
| 379 | | yeast. *Proceedings of the National Academy of Sciences* **108**, 698-703, |
| 380 | | doi:10.1073/pnas.1012363108 (2011). |
| 381 | 11 | Takahashi, T., Burguiere-Slezak, G., Auffret Van Der Kemp, P. & Boiteux, S. Topoisomerase 1 |
| 382 | | provokes the formation of short deletions in repeated sequences upon high transcription in |
| 383 | | Saccharomyces cerevisiae. *Proceedings of the National Academy of Sciences of the United* |
| 384 | | *States of America* **108**, 692-697, doi:10.1073/pnas.1012582108 (2011). |
| 385 | 12 | Kim, N. *et al.* Mutagenic processing of ribonucleotides in DNA by yeast topoisomerase I. |
| 386 | | *Science (New York, N.Y.)* **332**, 1561-1564, doi:10.1126/science.1205016 (2011). |
| 387 | 13 | Nick McElhinny, S. A. *et al.* Abundant ribonucleotide incorporation into DNA by yeast |
| 388 | | replicative polymerases. *Proceedings of the National Academy of Sciences of the United* |
| 389 | | *States of America* **107**, 4949-4954, doi:10.1073/pnas.0914857107 (2010). |
| 390 | 14 | Reijns, M. A. M. *et al.* Enzymatic removal of ribonucleotides from DNA is essential for |
| 391 | | mammalian genome integrity and development. *Cell* **149**, 1008-1022, |
| 392 | | doi:10.1016/j.cell.2012.04.011 (2012). |
| 393 | 15 | Sparks, J. L. *et al.* RNase H2-Initiated Ribonucleotide Excision Repair. *Molecular Cell* **47**, 980- |
| 394 | | 986, doi:10.1016/j.molcel.2012.06.035 (2012). |
| 395 | 16 | Huang, S. Y., Ghosh, S. & Pommier, Y. Topoisomerase I alone is sufficient to produce short |
| 396 | | DNA deletions and can also reverse nicks at ribonucleotide sites. *J Biol Chem* **290**, 14068- |
| 397 | | 14076, doi:10.1074/jbc.M115.653345 (2015). |
| 398 | 17 | Sparks, J. L. & Burgers, P. M. Error-free and mutagenic processing of topoisomerase 1- |
| 399 | | provoked damage at genomic ribonucleotides. *The EMBO journal* **34**, 1259-1269 (2015). |
| 400 | 18 | Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat Genet* **48**, 126- |
| 401 | | 133, doi:10.1038/ng.3469 (2016). |
| 402 | 19 | Alexandrov, L. B. & Stratton, M. R. Mutational signatures: the patterns of somatic mutations |
| 403 | | hidden in cancer genomes. *Curr Opin Genet Dev* **24**, 52-60, doi:10.1016/j.gde.2013.11.014 |
| 404 | | (2014). |
| 405 | 20 | Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, |
| 406 | | 415-421, doi:10.1038/nature12477 (2013). |
| 407 | 21 | Nik-Zainal, S. *et al.* The genome as a record of environmental exposure. *Mutagenesis* **30**, |
| 408 | | 763-770, doi:10.1093/mutage/gev073 (2015). |
| 409 | 22 | Pich, O. *et al.* The mutational footprints of cancer therapies. *Nat Genet* **51**, 1732-1740, |
| 410 | | doi:10.1038/s41588-019-0525-5 (2019). |
| 411 | 23 | Eichler, E. E. Genetic Variation, Comparative Genomics, and the Diagnosis of Disease. *New* |
| 412 | | *England Journal of Medicine* **381**, 64-74, doi:10.1056/NEJMra1809315 (2019). |
| 413 | 24 | Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res* |
| 414 | | **47**, D941-D947, doi:10.1093/nar/gky1015 (2019). |
| 415 | 25 | Conover, H. N. *et al.* Stimulation of Chromosomal Rearrangements by Ribonucleotides. |
| 416 | | *Genetics* **201**, 951-961, doi:10.1534/genetics.115.181149 (2015). |

417 26   Nick McElhinny, S. A. *et al.* Genome instability due to ribonucleotide incorporation into DNA.
418       *Nature chemical biology* **6**, 774-781, doi:10.1038/nchembio.424 (2010).
419 27   Certo, M. T. *et al.* Tracking genome engineering outcome at individual DNA breakpoints. *Nat*
420       *Methods* **8**, 671-676, doi:10.1038/nmeth.1648 (2011).
421 28   Williams, J. S. *et al.* Topoisomerase 1-mediated removal of ribonucleotides from nascent
422       leading-strand DNA. *Molecular cell* **49**, 1010-1015, doi:10.1016/j.molcel.2012.12.021 (2013).
423 29   Aden, K. *et al.* Epithelial RNase H2 Maintains Genome Integrity and Prevents Intestinal
424       Tumorigenesis in Mice. *Gastroenterology* **156**, 145-159.e119,
425       doi:10.1053/j.gastro.2018.09.047 (2019).
426 30   Tanizawa, A., Kohn, K. W. & Pommier, Y. Induction of cleavage in topoisomerase I c-DNA by
427       topoisomerase I enzymes from calf thymus and wheat germ in the presence and absence of
428       camptothecin. *Nucleic Acids Res* **21**, 5157-5166, doi:10.1093/nar/21.22.5157 (1993).
429 31   Huang, S. N., Williams, J. S., Arana, M. E., Kunkel, T. A. & Pommier, Y. Topoisomerase I-
430       mediated cleavage at unrepaired ribonucleotides generates DNA double-strand breaks.
431       *EMBO J* **36**, 361-373, doi:10.15252/embj.201592426 (2017).
432 32   Balachander, S. *et al.* Ribonucleotide incorporation in yeast genomic DNA shows preference
433       for cytosine and guanosine preceded by deoxyadenosine. *Nat Commun* **11**, 2447,
434       doi:10.1038/s41467-020-16152-5 (2020).
435 33   Zimmermann, M. *et al.* CRISPR screens identify genomic ribonucleotides as a source of
436       PARP-trapping lesions. *Nature* **559**, 285-289, doi:10.1038/s41586-018-0291-z (2018).
437 34   Consortium, T. G. E. R. The 100,000 Genomes Project | Genomics England. (2020).
438 35   Puente, X. S. *et al.* Non-coding recurrent mutations in chronic lymphocytic leukaemia.
439       *Nature* **526**, 519-524, doi:10.1038/nature14666 (2015).
440 36   Abascal, F. *et al.* Somatic mutation landscapes at single-molecule resolution. *Nature*,
441       doi:10.1038/s41586-021-03477-4 (2021).
442 37   Rheinbay, E. *et al.* Analyses of non-coding somatic drivers in 2,658 cancer whole genomes.
443       *Nature* **578**, 102-111, doi:10.1038/s41586-020-1965-x (2020).
444 38   Baranello, L. *et al.* RNA Polymerase II Regulates Topoisomerase 1 Activity to Favor Efficient
445       Transcription. *Cell* **165**, 357-371, doi:10.1016/j.cell.2016.02.036 (2016).
446 39   Zhao, G. *et al.* Gene4Denovo: an integrated database and analytic platform for de novo
447       mutations in humans. *Nucleic acids research* **48**, D913--D926 (2020).
448 40   Sekiguchi, J. & Shuman, S. Site-specific ribonuclease activity of eukaryotic DNA
449       topoisomerase I. *Mol Cell* **1**, 89-97, doi:10.1016/s1097-2765(00)80010-6 (1997).
450 41   Cho, J. E., Kim, N., Li, Y. C. & Jinks-Robertson, S. Two distinct mechanisms of Topoisomerase
451       1-dependent mutagenesis in yeast. *DNA Repair (Amst)* **12**, 205-211,
452       doi:10.1016/j.dnarep.2012.12.004 (2013).
453 42   Boot, A. & Rozen, S. G. Recurrent mutations in topoisomerase 2a cause a novel mutator
454       phenotype in human cancers. *bioRxiv* (2020).
455 43   Park, S. *et al.* Clonal dynamics in early human embryogenesis inferred from somatic
456       mutation. *bioRxiv*, 2020.2011.2023.395244, doi:10.1101/2020.11.23.395244 (2020).
457 44   Alvarez-Quilon, A. *et al.* Endogenous DNA 3' Blocks Are Vulnerabilities for BRCA1 and BRCA2
458       Deficiency and Are Reversed by the APE2 Nuclease. *Mol Cell* **78**, 1152-1165 e1158,
459       doi:10.1016/j.molcel.2020.05.021 (2020).
460 45   Li, F. *et al.* Apn2 resolves blocked 3' ends and suppresses Top1-induced mutagenesis at
461       genomic rNMP sites. *Nat Struct Mol Biol* **26**, 155-163, doi:10.1038/s41594-019-0186-1
462       (2019).
463 46   Potenski, C. J., Niu, H., Sung, P. & Klein, H. L. Avoidance of ribonucleotide-induced mutations
464       by RNase H2 and Srs2-Exo1 mechanisms. *Nature* **511**, 251-254 (2014).
465 47   Riva, V. *et al.* Novel alternative ribonucleotide excision repair pathways in human cells by
466       DDX3X and specialized DNA polymerases. *Nucleic Acids Res* **48**, 11551-11565,
467       doi:10.1093/nar/gkaa948 (2020).

468    48    Davies, H. *et al.* HRDetect is a predictor of BRCA1 and BRCA2 deficiency based on mutational
469        signatures. *Nat Med* **23**, 517-525, doi:10.1038/nm.4292 (2017).
470    49    Wang, C. *et al.* Genome-wide CRISPR screens reveal synthetic lethality of RNASEH2
471        deficiency and ATR inhibition. *Oncogene* **38**, 2451-2463, doi:10.1038/s41388-018-0606-4
472        (2019).
473    50    Consortium, T. I. T. P.-C. A. o. W. G. Pan-cancer analysis of whole genomes. *Nature* **578**, 82-
474        93, doi:10.1038/s41586-020-1969-6 (2020).

475

476 **Figure legends**

477 **Fig. 1 | Top1-dependent deletions in *S. cerevisiae* resemble ID4, a cancer mutational signature of**

478 **unknown aetiology. a,** The ID4 signature comprises small deletions (typically 2, 3 or 4 bp in size) of

479 one repeat unit at short short tandem repeat (SSTR) and microhomology (MH) sites. i-vi; repeated

480 sequence in bold colour; deletions in red; SNMH, single nucleotide MH. **b,** Indel mutations similar to

481 those detected in ID4 accumulate genome-wide in yeast with high levels of genome-embedded

482 ribonucleotides. Reanalysis of WGS for *rnh201Δ pol2-M644G* yeast[25]. **c,** Schematic of novel

483 frameshift mutation reporter containing many 2 bp SSTRs. Frameshift mutations in HygroR result in

484 neomycin resistant yeast colonies. $P_{TEF}$, TEF promoter; HygroR/NeoR, hygromycin/neomycin

485 resistance genes; P2A, self-cleaving peptide. **d,e,** Fluctuation assays demonstrate that Top1-

486 mediated 2 bp SSTR mutations occur in wild-type and RNase H2 deficient (*rnh201Δ*) backgrounds.

487 Mutation rates, median ± 95% confidence intervals for *n*=16 independent cultures per strain (**d**). WT

488 and *rnh201Δ* have similar indel mutation spectra, and differ from *top1Δ* strains. Spectra of neomycin

489 resistant colonies. n, number of independent indels detected (**e**). Cosine similarity P-values

490 empirically determined, Extended Data Fig. 2e,f.

491

492 **Fig. 2 | Two bp SSTR deletions are increased in RNase H2 null HeLa cells. a,** Schematic of reporter

493 targeting to *AAVS1* safe harbour locus to generate reporter cells. PuroR, puromycin resistance; $P_{CAG}$,

494 CAG promoter; HA, homology arm (L, left; R, right). Also see Extended Data Fig. 3. **b,c,** Validation of

495 *RNASEH2A* knockout reporter clones. Immunoblot of cell lysates detecting the three RNase H2

496 subunits. GAPDH, loading control. For gel source data, see Supplementary Fig. 1 (**b**). Cellular RNase

497     H2 enzyme activity. Bars, mean; error bars, s.d.; *n*=3 technical replicates. HeLa, no modification;

498     Parental, HeLa with reporter (grey); RNASEH2A+, CRISPR-edited reporter clone retaining RNase H2

499     activity (green); KO1, KO2, CRISPR-mediated *RNASEH2A* knockout clones (red) (**c**). **d,** Fluctuation

500     assays establish a significantly increased mutation rate in RNase H2 null (KO) cells. Median ± 95%

501     confidence intervals. Data points, rates for independent cultures (RNase H2 proficient, RNASEH2A+,

502     *n*=9; knockout, KO1, open circles, *n*=10; KO2, open squares, *n*=6). **e,** 2 bp SSTR and SNMH deletions

503     are frequent in both RNASEH2A+ and KO cells. Indel mutation spectra. n, number of indels identified

504     by sequencing colonies from independent cultures.

505

506     **Fig. 3 | ID4 SSTR and MH mutations are increased genome-wide in RNase H2 deficient RPE1 cells.**

507     **a,** Schematic of mutation accumulation experiment. Long-term culture of hTERT-RPE1 *TP53*$^{-/-}$ RNase

508     H2 wildtype (WT) and null cell lines (AKO, BKO: RNASEH2A, RNASEH2B knockout respectively)

509     bottlenecked every 25 doublings by single cell sorting. **b,** Mutations acquired during long-term

510     culture were significantly enriched for 2-5 bp deletions in RNase H2 null cells, but not other mutation

511     categories (also see Extended Data Fig 4e). Mean ± s.d.; P-value, two-sided Fisher's exact test with

512     Bonferroni correction, WT (*n*=3 independent clones) vs KO (*n*=2 independent clones) for 2-5 bp

513     deletions vs all other indel types. **c,d,** ID4 occurs in RNase H2 null cells (**c**), and is the major signature

514     once background mutations observed in wildtype cells are subtracted (**d**).

515

516     **Fig. 4 | RER-deficient tumours have an ID4 signature, associated with transcription and a TNT**

517     **sequence motif. a,** ID4 contributes substantially to the mutational spectrum of Rnaseh2b-KO murine

518     intestinal tumours (WGS, paired tumour-normal samples from *n*=6 mice). **b,** ID4 contribution is

519     greater in transcribed regions of the genome. Two-sided Fisher's exact test, ID4 vs other indels.

520     *n*=969 indels from 6 biologically independent tumours. **c,** 2 bp STR/SNMH deletions have biased

521     sequence composition. Genome, frequency of dinucleotides in STR/SNMH sequences in the

522  mappable genome. Deletions (bold red), right aligned. **d,e,** A TNT sequence motif is present at all 2

523  bp STR and SNMH deletions. Sequence logo: Two-bit representation of the sequence context of 2 bp

524  deletions at STR and SNMH sequences (**d**). Deletion sites are significantly enriched for the TNT

525  sequence motif compared to genome-wide occurrence, for all genome sequence, as well as STR and

526  SNMH sites. P-values, Two-sided Fisher's exact test, observed vs expected. $n$=228 (all; $P$=1.7x10$^{-28}$),

527  124 (STR; $P$=0.0008), 77 (SNMH; $P$=1.4x10$^{-8}$) deletions in 6 biologically independent tumours (**e**). **f,**

528  Model for TOP1-mediated mutations at TNT sequences containing embedded ribonucleotides, in

529  which strand realignment results in 2 nt deletion (description in main text).

530

531  **Fig. 5 | TOP1-mediated deletions in human cancer and germline. a,** 2-5 bp deletions are

532  significantly increased in CLL with biallelic *RNASEH2B* deletions (null). Box, 25-75%; line, median;

533  whiskers 5-95%; data points, values outside range. WT, $n$=116, 85; het (heterozygous), $n$=72, 59; null,

534  $n$=10, 6 tumours (GEL, ICGC respectively). Multiple-testing corrected q-values, 2-sided Mann-

535  Whitney. **b-d,** ID-TOP1 deletions are frequent somatic mutations in cancer. Indels per expression

536  stratum of ubiquitously expressed genes (defined in Extended Data Fig. 8e). Dotted line, genome-

537  wide rate (**b**). Two bp deletions preferentially occur at TNT motifs. P-values, two-sided Fisher's exact

538  test, observed vs expected. $n$=11,853 (all; $P$<10$^{-200}$), 6,699 (STR; $P$=1.9x10$^{-60}$), 2,872 (SNMH; $P$=1.5x10$^{-}$

539  $^{51}$) deletions (**c**). 2-5 bp deletions increase with TOP1 cleavage activity in ID4-positive PCAWG

540  tumours (**d**). Solid lines, relative deletion rate. Shading, 95% confidence intervals from 100 (**b**) or

541  1,000 bootstrap replicates (**d**), n= 11,853 biologically independent tumours[50] (**b-d**). **e,** 2-5 bp

542  deletions are enriched at tissue-specific highly transcribed genes in associated cancers. Heatmap of

543  significant Odds Ratio scores (2-5 bp deletions in top 10% tissue-restricted genes vs deletions in

544  other genes, relative to expected frequency from all other tissues) for tissue-tumour pairs. Two-

545  sided Fisher's exact test. **f-h,** ID-TOP1 deletions are frequent human *de novo* mutations enriched in

546  highly transcribed germ cell genes. 2-5 bp deletions are the most common indels in the human

547 germline. Gene4Denovo WGS data ([39] ; $n$=40,936 indels) (**f**). TNT sequence motif is significantly

548 enriched in *de novo* 2 bp deletions (**g**). P-values, two-sided Fisher's exact test, observed vs expected;

549 $n$=5,569 two bp deletions ($P<10^{-200}$), at STR ($n$=3,294; $P$=5.2x10$^{-47}$) and SNMH sequences ($n$=1,093;

550 $P$=2.9x10$^{-26}$). 2-5 bp deletion frequency correlates with gene transcription level in germ cells (**h**).

551 Solid lines, Gene4Denovo indel mutations per individual per Mbp. Shading, 95% confidence intervals,

552 100 bootstrap replicates.

553

554 **Methods**

555 **Plasmids**

556 A description of all plasmids used in this work can be found in Supplementary Table 1. The *S.*

557 *cerevisiae* reporter was generated by DNA synthesis (GeneArt Gene Synthesis, Thermo Fisher

558 Scientific; gBlocks Gene Fragments, IDT) and conventional cloning (restriction, ligation and

559 Quikchange site-directed mutagenesis). The final construct (pTCW12) was used for *S. cerevisiae*

560 reporter strain construction and fluctuation assays. A Gateway compatible reporter construct for

561 mammalian cells (pTCW14) was similarly generated using a combination of DNA synthesis and

562 conventional cloning strategies. Gateway cloning was then used to move the reporter cassette into

563 pAAVS-Nst-CAG-Dest (a gift from Knut Woltjen; Addgene plasmid # 80489; [51]) to generate pTCW15

564 for targetting it to the human AAVS1 locus.

565

566 *In silico* **re-design of the Hygromycin resistance gene**

567 To increase the frequency of 2-bp tandem repeats, synonymous substitutions were introduced in the

568 1 kbp *hph* coding sequence, the *Klebsiella pneumoniae* hygromycin resistance gene (HygroR)[52]. Using

569 Python, a 5-codon (15-base) sliding window was moved one codon at a time, to identify all possible

570 synonymous permutations. Permutations were ranked on the basis of tandem dinucleotide repeat

571   sequence length, with the highest ranking sequences used to replace whole codons, prioritising

572   dinucleotide couplet repeats over mononucleotide repeats. Edited codons were then censored from

573   subsequent permutation. Subsequently, to eliminate stop codons that would arise after a 2 bp

574   deletion or equivalent frameshift mutations, further synonymous changes were made, where

575   possible preserving tandem repeat sequences.

576

577   **Yeast strains and growth conditions**

578   All *S. cerevisiae* strains used in this work (Supplementary Table 2) are isogenic with BY4741 [53] and

579   were grown at 30°C. *TOP1* and *RNH201* open reading frames (ORFs) were deleted using 1-step allele

580   replacement using PCR products generated from plasmid templates with selection cassettes

581   (Supplementary Table 2) and primers containing 60 nt homology directly up and downstream of the

582   ORF. Gene deletions were confirmed by PCR. The 2 bp deletion reporter was inserted at the *AGP1*

583   locus using a PCR product amplified from pTCW12 using primers AGP1-MX6-F and AGP1-MX6-R

584   (Supplementary Table 3). Correct reporter insertion was confirmed by PCR and Sanger sequencing.

585   Growth under selection was on YPD (10 g/l yeast extract, 20 g/l bactopeptone, 20 g/l dextrose, 20 g/l

586   agar) supplemented with hygromycin B (300 mg/l), nourseothricin (100 mg/l) and/or G418 (1 g/l), or

587   on Synthetic Defined medium (6.7 g/l yeast nitrogen base without amino acids, complete

588   supplement single dropout mixture (Formedium), 20 g/l dextrose, 20 g/l agar).

589

590   **Fluctuation assays (yeast)**

591   Fluctuation assays were performed as previously described[54]. Yeast was grown overnight in YPD with

592   hygromycin B (300 mg/l), plated on YPD and grown at 30°C to obtain individual colonies derived

593   from a single cell without HygroR mutations. For each strain, 16 independent colonies were then

594   used to inoculate 5 ml YPD, and grown for 3 days at 30°C with shaking at 250 rpm. Cells were

595    pelleted by centrifugation and resuspended in 1 ml of $H_2O$. Undiluted suspensions for each culture

596    were plated (100 μl per plate) on 2 YPD plates supplemented with 1 g/l G418, with the exception of

597    *rnh201Δ* for which a $10^{-2}$ dilution was used. In addition, each suspension was serially diluted to $10^{-6}$

598    of which 100 μl per plate was spread on 2 YPD plates to estimate the total number of viable cells per

599    culture. Plates were incubated at 30°C for 2-3 days, and colonies counted. Mutation rates were

600    determined in Microsoft Excel 2016 for each individual culture, and an overall rate for each strain

601    calculated using the Lea Coulson method of the median[55]. The number of mutants for each culture

602    were ranked, and those ranked 4[rd] and 13[th] used to calculate the rates that define the lower and

603    upper limits of the 95% confidence interval[56]. A single G418-resistant colony for each independent

604    culture was used to determine the spectrum of frame shift mutations. A 1.3 kbp region including

605    HygroR was amplified in two overlapping amplicons (primers S297F and S1113R; S752 and S1658R)

606    using FastStart PCR Master Mix (Roche) and direct colony PCR (5 min 95˚C; 35 cycles 30 s 95˚C, 30 s

607    58˚C, 45 s 72˚C; 45 s 72˚C). Each amplicon was Sanger sequenced using primers described in

608    Supplementary Table 3, and analysed using Sequencher 5.4.6 (Gene Codes Corporation) and/or

609    Mutation Surveyor V3.30 (SoftGenetics). Mutation rates (per bp) were calculated for 1,032 bp of

610    sequence in which productive frameshift mutations can occur.

611

612    **Cell lines**

613    Human cell lines used in this work are summarised in Supplementary Table 4. All cells were grown at

614    37°C and 5% $CO_2$, authenticated using STR DNA profiling in the labs of origin and shown to be

615    mycoplasma negative through routine testing. HeLa cells (a gift from G. Stewart, University of

616    Birmingham, UK; originally purchased from ATCC) were grown in Dulbecco's Modified Eagle Medium

617    (DMEM; Gibco/Thermo Fisher Scientific) supplemented with 10% fetal bovine serum (FBS), 100 U/ml

618    penicillin and 100 μg/ml streptomycin. hTERT-RPE1 cells (a gift from D. Durocher, University of

619    Toronto, Canada; originally purchased from ATCC) were grown in DMEM/F12 medium mixture

620    (Gibco/Thermo Fisher Scientific) supplemented with 10% FBS, 100 U/ml penicillin and 100 µg/ml

621    streptomycin. The 2-bp deletion reporter was integrated at the *AAVS1* safe harbour locus in HeLa

622    cells using a published CRISPR/Cas9 targeting protocol [51]. HeLa cells were transfected with pXAT2

623    and pTCW15 in Opti-MEM reduced-serum medium using Invitrogen Lipofectamine 2000 (Thermo

624    Fisher Scientific). After 48 h cells were re-plated in medium containing 500 µg/ml G418, and after

625    another 48 h and a second round of re-plating in selective medium, single cells were sorted into 96-

626    well plates using a BD FACSJazz instrument (BD Biosciences). Resulting G418-resistant clones were

627    screened by PCR for reporter integration at the correct locus, retention of integration-free *AAVS1*

628    and Sanger sequencing of resulting PCR products. Single-locus integration was confirmed by FISH as

629    previously described[57], using pTCW16 to generate a fluorescently labelled probe. The full reporter

630    sequence of selected clones was checked, with amplification of a 1.9 kbp fragment using Prime Star

631    Max PCR Master Mix (Takara Bio) with primers HygroR_up and PuroR_rev (40 cycles 10 s 98˚C, 15 s

632    70˚C, 2 min 72˚C), followed by Sanger sequencing with additional primers (Supplementary Table 3).

633    To generate RNASEH2A-KO reporter cells, the selected parental HeLa reporter clone was transfected

634    with pMAR526 and pMAR527 (Supplementary Table 1), using Lipofectamine 2000. Forty-eight hours

635    after transfection, single EGFP-expressing cells were sorted into 96-well plates and grown until

636    colonies formed. Initial screening was based on PCR amplification (primers RNASEH2A-ex1F and

637    RNASEH2A-ex1R) of the CRISPR/Cas9-targeted region of *RNASEH2A* with mutations present in

638    selected clones determined by Sanger sequencing. The cellular RNase H2 status was then confirmed

639    by immunoblotting, RNase H2 enzymatic activity assay, and alkaline gel electrophoresis to determine

640    ribonucleotide content of genomic DNA (detailed methods below).

641

642    **Fluctuation assays (human)**

643    Hygromycin resistant HeLa reporter cells (400 µg/ml hygromycin B) were recovered from frozen

644    stocks in the absence of selection. The following day, 10 wells of a 96-well plate were seeded with

645    2,000 cells per well for each line. The experiment was performed with the operator blinded to the

646    identity of the cell lines. Cells were cultured under non-selective conditions and re-plated

647    subsequently in 24-well, 6-well plates and ultimately T75 flasks, in which they were grown to

648    confluence. Cells were then dissociated using Gibco TrypLE (Thermo Fisher Scientific) and cells

649    counted using a Moxi Z automated cell counter. After serial dilution 1,000 cells were plated into two

650    10-cm plates for each culture and grown for 14 days to determine plating efficiency. All other cells

651    were plated into two 10-cm plates, 0.5 μg/ml puromycin added after 4 h, with medium subsequently

652    changed every 2-3 days for 14 days to remove dead cells and maintain a puromycin concentration of

653    0.5 μg/ml.

654    To establish mutation spectra, colonies were removed by scraping and then cultured in a 96-well

655    plate. When confluent, cells were lysed with 75 μl DirectPCR Lysis Reagent (Viagen Biotech) and 0.4

656    mg/ml PCR-grade Proteinase K (Roche), heating overnight at 55 °C followed by 45 min at 85 °C. Only

657    one sample per independent culture was used for PCR amplification and Sanger sequencing to

658    determine the nature of mutations in the HygroR coding sequence. A 1.24 kbp region including

659    HygroR was amplified with Prime Star Max PCR Master Mix (Takara Bio), HygroR_up and H1327R

660    primers (40 cycles 10 s 98˚C, 15 s 70˚C, 2 min 72˚C). Sanger sequencing was then performed with

661    additional primers (Supplementary Table 3) and mutations identified using Mutation Surveyor V3.30

662    (SoftGenetics). All mutants showed double traces of equal height from the point of indel mutations,

663    consistent with the presence of two copies of the reporter in all reporter lines. As FISH indicated

664    presence of the reporter at a single AAVS1 locus, we inferred that two copies of the reporter were

665    inserted in tandem at this locus. As a 2 bp deletion or equivalent frameshift mutation in either

666    HygroR copy would bring the associated PuroR coding sequence into the translated reading frame,

667    we corrected mutation rate calculations (per bp) for the presence of 2 copies.

668    To determine colony numbers, plates were washed with PBS, fixed with 2% formaldehyde in PBS for

669    10 min, rinsed with water, and colonies stained with 0.1% crystal violet solution for 10 min. Plates

670    were then washed with water and left to dry before counting colonies. After counting the

671    experiment was unblinded. Mutation rates were determined for each individual culture in Microsoft

672    Excel 2016, and an overall rate for WT and KO strains calculated using the Lea Coulson method of

673    the median. The number of mutants for each culture were ranked, and appropriate ranks[56] used to

674    calculate the rates that define the lower and upper limits of the 95% confidence interval.

675

676    **Immunoblotting**

677    Whole-cell extracts (WCE) to determine protein levels of RNase H2 subunits by immunoblotting and

678    for RNase H2 activity assays were prepared as previously described[58]. Equal amounts of protein from

679    WCE were separated by SDS-PAGE on 4-12% NuPAGE gels and transferred to PVDF. Membranes

680    were probed in 5% milk (w/v; Marvel Original Dried Skimmed), TBS+0.2% Tween-20 (v/v) with the

681    following antibodies: sheep anti-RNase H2 (raised against human recombinant RNase H2, 1:1,000)[14];

682    mouse anti-RNASEH2A G-10 (Santa Cruz Biotechnologies sc-515475, lot #A1416, 1:1,000); rabbit

683    anti-GAPDH (Abcam ab9485, 1:2,000, lot #GR3380498-1). For detection we Rabbit Anti-Sheep

684    Immunoglobulins/HRP (Dako, P04163, lot #00047199, 1:2,000); Goat Anti-Mouse

685    Immunoglobulins/HRP (Dako, P0447, lot #20039214, 1:10,000); Anti-rabbit IgG, HRP-linked Antibody

686    (Cell Signaling Technologies, 7074S, lot #29, 1:10,000); Amersham ECL Prime Western Blotting

687    Detection Reagent (GE Healthcare Life Sciences) and an ImageQuantLAS4000 device, or IRDye

688    secondary antibodies and an Odyssey CLx Imaging System (LI-COR Biosciences). Uncropped

689    immunoblots are presented in Supplementary Fig. 1.

690

691    **RNase H2 activity assays**

692    To assess cellular RNase H2 activity, a FRET-based fluorescent substrate release assay was

693    performed as previously described[14]. Briefly, RNase H2-specific activity was determined by

694  measuring the cleavage of double-stranded DNA substrate containing a single embedded

695  ribonucleotide. Activity against a DNA-only substrate of the same sequence was used to correct for

696  background activity. Substrates were formed by annealing a 3'-fluorescein-labeled oligonucleotide

697  (GATCTGAGCCTGGGaGCT or GATCTGAGCCTGGGAGCT; uppercase DNA, lowercase RNA) to a

698  complementary 5'-DABCYL-labelled DNA oligonucleotide (Eurogentec). Reactions were performed in

699  100 µl reaction buffer (60 mM KCl, 50 mM Tris–HCl pH 8.0, 10 mM $MgCl_2$, 0.01% BSA, 0.01% Triton

700  X-100) with 250 nM substrate in black 96-well flat-bottomed plates (Costar) at 24°C. WCE was

701  prepared as described above, protein concentrations determined using a Bio-Rad Bradford Protein

702  Assay, and the final protein concentration per reaction was 50 ng/µl. Fluorescence was read (100

703  ms) every 5 min for up to 90 min using a VICTOR2 1420 multilabel counter (Perkin Elmer), with a

704  480-nm excitation filter and a 535-nm emission filter. Initial substrate conversion after background

705  subtraction was used to calculate RNase H2 enzyme activity.

706

707  **Alkaline gel electrophoresis**

708  To determine the presence of excess genome-embedded ribonucleotides in nuclear DNA, alkaline

709  gel electrophoresis of RNase H2 treated genomic DNA was performed as previously described[58].

710  Briefly, total nucleic acids were isolated from pellets from ~1 million cells by incubation in ice-cold

711  buffer (20 mM Tris–HCl pH 7.5, 75 mM NaCl, 50 mM EDTA) with 200 µg/ml proteinase K (Roche) for

712  10 min on ice, followed by addition of N-lauroylsarcosine sodium salt (Sigma) to a final concentration

713  of 1%. Nucleic acids were phenol:chloroform-extracted, isopropanol precipitated and dissolved in

714  nuclease-free water. For alkaline gel electrophoresis, 500 ng of total nucleic acids was incubated

715  with 1 pmol of purified recombinant human RNase H2 (isolated as previously described[59]) and 0.25

716  µg of DNase-free RNase (Roche) for 30 min at 37°C in 100 µl reaction buffer (60 mM KCl, 50 mM

717  Tris–HCl pH 8.0, 10 mM $MgCl_2$, 0.01% Triton X-100). Nucleic acids were ethanol precipitated,

718  dissolved in nuclease-free water and 250 ng separated on 0.7% agarose gels in 50 mM NaOH, 1 mM

719    EDTA. After overnight electrophoresis, the gel was neutralised in 0.7 M Tris–HCl pH 8.0, 1.5 M NaCl

720    and stained with SYBR Gold (Invitrogen). Imaging was performed on a FLA-5100 imaging system

721    (Fujifilm), and densitometry plots were generated using AIDA Image Analyzer v3.44.035 (Raytest).

722

723    **Mutation accumulation experiment**

724    TP53-KO hTERT-RPE1 cells without and with loss-of-function mutations in *RNASEH2A* or *RNASEH2B*,

725    introduced by CRISPR/Cas9 genome editing, a gift from D. Durocher (The Lunenfeld–Tanenbaum

726    Research Institute, Toronto), have been previously described[33]. RNase H2 proficient (WT),

727    RNASEH2A-KO and RNASEH2B-KO cells were single cell sorted into 96-well plates using a BD

728    FACSJazz instrument (BD Biosciences). Multiple individual clones for each were expanded to

729    confluent T75 flasks for cryopreservation and genomic DNA isolation of these ancestral populations.

730    In addition, lines were again single cell sorted into 96-well plates to start the mutation accumulation

731    experiment. Cultures were expanded by subsequent growth in 24-well, 6-well plates and T75 flasks

732    until confluent (approximately 25 population doublings), and this process of single cell sorting and

733    expansion was repeated 4 more times providing bottlenecks to capture mutations that occurred

734    since the previous sort. From the first to the last single cell sort a total of approximately 100

735    population doublings occurred and the final culture was expanded for cryopreservation and genomic

736    DNA isolation of these end-point populations.

737    Genomic DNA was isolated using phenol extraction as previously described[58], for alkaline gel

738    electrophoresis and whole genome sequencing. Library preparations and sequencing were

739    performed by Edinburgh Genomics. Libraries were prepared using Illumina SeqLab specific TruSeq

740    PCRFree High Throughput library preparation kits as per manufacturer's instructions, with DNA

741    samples sheared to a 450 bp mean insert size. Libraries were sequenced using paired-end reads on

742    an Illumina HiSeqX instrument using v2.5 chemistry to achieve minimum mean genome-wide

743    sequencing depth of 30x per sample.

744

**Mouse whole genome sequencing**

746 Villin-Cre$^+$ Trp53$^{fl/fl}$ Rnaseh2b$^{fl/fl}$ mice with epithelial-specific deletion of *Trp53* and *Rnaseh2b* on a

747 C57Bl/6J background have been described previously[29]. Animal experiments were conducted with

748 appropriate permission, in accordance with guidelines for animal care of the Christian-Albrechts-

749 University (Kiel, Germany), in agreement with national and international laws and policies. No

750 randomisation or blinding was performed. Paired tumour-normal DNA was isolated from small

751 intestinal tumours (*Trp53$^{-/-}$ Rnaseh2b$^{-/-}$*) and liver tissue (*Trp53$^{+/+}$ Rnaseh2b$^{+/+}$*) from 52-week old

752 females, using a Qiagen DNeasy Blood & Tissue Kit. Library preparations and sequencing were

753 performed by Edinburgh Genomics using Illumina DNA PCR-Free Library Prep as per manufacturer's

754 instructions. Paired end sequencing was performed by Edinburgh Genomics on a NovaSeq 6000

755 using v1.5 chemistry. Mean genome-wide sequencing depth of at least 30x for liver samples and 60x

756 for tumour samples was obtained.

757

758 *S. cerevisiae* **WGS analysis**

759 Whole genome sequencing SRA files for *rnh201Δ pol2-M644G S. cerevisiae*[25] from the NCBI

760 Sequence read archive (SRA) were converted to FASTQ files using SRA Toolkit v2.5.4-1 (SRA Toolkit

761 Development Team; http://ncbi.github.io/sra-tools/). FASTQ reads were aligned to the

762 GSE56939_L03_ref_v2 reference genome ([60]; Supplementary Table 5) and sorted BAM files created

763 using BWA-MEM 0.7.12[61], and deduplicated with SAMBLASTER v0.1.22[62]. To select high quality indel

764 variants, GATK (v3.6-0) Haplotype Caller (without Base Quality Score Recalibration)[63] variant calling

765 was performed with "Hard Filters" (--filterExpression "QD < 2.0 || FS > 200.0 || ReadPosRankSum < -

766 20.0") . Filtering for strain-specific variants was performed as previously described[60], with minor

767 modifications. Filters: 1) eliminate variants shared with an ancestral clone; 2) required ≥ 20 reads for

768    variant allele in descendent; 3) exclusion of repetitive sequences as defined in [60]; 4)

769    reference/variant depth ratio 0.4-0.6; < 0.4 if homozygous variant allele .

770

771    **RPE1 WGS analysis**

772    FASTQs were converted to unaligned BAM format and Illumina adaptors marked using GATK v4.1.9.0

773    FastqToSam and MarkIlluminaAdapters tools[64]. Reads were aligned to the human genome (hg38,

774    including alt, decoy and HLA sequences) using BWA-MEM v0.7.16 [61] and read metadata merged

775    using GATK's MergeBamAlignment tool. PCR and optical duplicate marking and base quality score

776    recalibration were performed using GATK. Variants from NCBI dbSNP build 151 were used as known

777    sites for base quality score recalibration. Post-processed alignments were genotyped using Mutect2,

778    Strelka2, Platypus and SvABA using somatic calling models for each pair of ancestral and endpoint

779    cultures, as detailed below.

780

781    **Mouse WGS analysis**

782    FASTQ processing and alignment was performed as for RPE1 WGS analysis, using the GRCm38 mouse

783    genome reference and known variant sites from the Mouse Genomes Project[65] (REL-1807-

784    SNPs_Indels) for base quality score recalibration. Somatic variant calling of post-processed

785    alignments was performed using Mutect2, Strelka2, Platypus and SvABA for each tumour-liver pair,

786    as detailed below. Somalier v0.2.12 (https://github.com/brentp/somalier) was used to confirm each

787    paired tumour and liver sample originated from the same animal.

788

789    **Human Ethics approval**

790      Data generated from Genomics England 100,000 genomes and ICGC-CLL studies were analysed. In

791      these respective studies, informed consent for participation was obtained. Ethical approval for

792      Genomics England 100,000 genomes project: East of England and South Cambridge Research Ethics

793      Committee; CLL-ICGC: International Cancer Genome Consortium (ICGC) guidelines from the ICGC

794      Ethics and Policy committee were followed and the study was approved by the Research Ethics

795      Committee of the Hospital Clínic of Barcelona.

796

797      **CLL WGS analysis**

798      *Genomics England*: CLL tumour-normal pairs (*n*=198) were processed as part of the 100,000

799      Genomes Project (pilot and main programme v8). Samples were sequenced using the Illumina HiSeq

800      X System with 150 bp paired-end reads at a minimum of 75x coverage for tumours and 30x coverage

801      for germline samples. Reads were mapped to GRCh38 using ISAAC aligner v03.16.02.19 [66]. SNVs and

802      indels were called using Strelka v2.4.7 using somatic calling mode. Structural and copy number

803      variants were called using Manta v0.28.0 and Canvas v1.3.1 [67], respectively. Samples with a tumour

804      purity estimate from Canvas of less than 50% were excluded from analysis. *RNASEH2B* copy number

805      was determined using a combination of Canvas, Manta, read depth counts with samtools v1.9 and

806      confirmed by manual inspection using IGV (v2.5.0)[68].

807      *ICGC*: WGS from the ICGC-CLL cohort[35] (*n*=150) was re-analysed. Raw reads were mapped to the

808      human reference genome (GRCh37) using BWA-MEM (v0.7.15)[61]. BAM files were generated, sorted,

809      indexed and optical or PCR duplicates flagged using biobambam2

810      (https://gitlab.com/german.tischler/biobambam2, v2.0.65). Copy number alterations were called

811      from WGS data using Battenberg (cgpBattenberg, v3.2.2)[69], ASCAT (ascatNgs, v4.1.0)[70], and Genome-

812      wide Human SNP Array 6.0 (Thermo Fisher Scientific) data[35] re-analysed using Nexus 9.0

813      Biodiscovery software (Biodiscovery). *RNASEH2B* copy number was established by combining the

814      three analyses and manual review with IGV.

815

**Colorectal Cancer WGS analysis**

817   Irinotecan-treated (*n*=39) and irinotecan-untreated (*n*=78) colorectal cancers from the 100,000

818   Genomes Project Colorectal Cancer Domain were 1:2 matched using a multivariate greedy matching

819   algorithm without replacement, implemented in the Matching R-package[71]. Matching was

820   conducted considering sex, age at sampling, whether a primary tumour or metastasis had been

821   sequenced, microsatellite instability status, and whether the individual have previously received

822   radiotherapy, oxaliplatin, capecitabine or fluorouracil treatment.

823

**Somatic Variant Calling**

825   Somatic variant calling was performed in parallel using four distinct methods: Mutect2 (as part of

826   GATK v4.1.9.0)[72,73], Strelka2 (v2.1.9.10)[74], SvABA (v1.1.3)[75] and Platypus (v0.8.1)[76]. High-confidence

827   indel calls were defined as the intersected output of these four tools, where variants passed all

828   filters for ≥3 of 4 callers. The intersection was performed using the bcftools (v1.10.2)[77] isec function

829   after normalising variant calls and left-aligning ambiguous alignment gaps using the bcftools norm

830   function. For Platypus (v0.8.1)[76], joint calling all samples in each cohort was performed before

831   filtering for somatic variants; the other variant callers were run in paired tumour-normal mode. For

832   the RPE1 mutation accumulation experiment the endpoint and ancestral cultures were defined as

833   "tumour" and "normal" samples respectively. Variant filtering strategies were optimised to both

834   available information on segregating genetic variation for humans and mice, and the functionality of

835   each calling method as detailed below.

836   *Mutect2*: unfiltered genotypes for all normal samples were combined to filter germline variants.

837   Somatic calls were obtained using GATK's FilterMutectCalls command. Human polymorphism data

838    and allele frequencies from, gnomAD[78] were provided to Mutect2 for the filtering of germline

839    variants.

840    *SvABA*: Germline indel and structural variants were filtered using --dbsnp-vcf and --germline-sv-

841    database options. Mouse indels were obtained from Mouse Genomes Project version 5 SNP and

842    (ftp://ftp-mouse.sanger.ac.uk/REL-1505-

843    SNPs_Indels/mgp.v5.merged.indels.dbSNP142.normed.vcf.gz); structural variants from SV release

844    version 5 (ftp://ftp-mouse.sanger.ac.uk/REL-1606-SV/mgpv5.SV_insertions.bed.gz and ftp://ftp-

845    mouse.sanger.ac.uk/REL-1606-SV/mgpv5.SV_deletions.bed.gz). Human indels were extracted from

846    NCBI dbSNP build 151 and common structural variants from dbVAR

847    (https://hgdownload.soe.ucsc.edu/gbdb/hg38/bbi/dbVar/).

848    *Strelka2*: candidate small indels for each pair were first generated by Manta (v1.6.0)[79] in somatic

849    calling mode. Strelka2 was then executed in somatic calling mode for each pair with Manta's

850    candidate small indels output provided to the --indelCandidates option.

851    *Platypus*: Germline variants were filtered on the basis of any "normal" sample with ≥ 2 variant allele

852    reads. Somatic variant calls for each sample pair were retained if "tumour"/endpoint sample > 2

853    variant reads; site depth > 9; and "normal" sample read depth ≥ 20, <2 variant reads. Additionally, a

854    >10x ratio of tumour to normal for variant/total depth was required.

855    For Genomics England CLL tumour-normal pairs, pre-existing Strelka2 calls from the 100,000

856    Genomes Project pipeline were used, while variant calling with Mutect2, Platypus and SvABA was

857    performed as above. Colorectal cancer tumour-normal pairs from Genomics England were processed

858    as for Genomics England CLL but without Mutect2 analysis. For ICGC CLL, somatic indels were called

859    using Mutect2 (GATK v4.0.2.0)[72,73], Strelka2 (v2.8.2)[74], SvABA (v1.1.0)[75], and Platypus (v0.8.1)[76].

860    Candidate small indels generated by Manta (v1.2)[79] were used as input for Strelka2. Mutect2,

861    Strelka2 and SvABA were run in paired tumour-normal mode. somaticMutationDetector.py

862    (https://github.com/andyrimmer/Platypus/tree/master/extensions/Cancer) was used to identify

863    somatic indels called by Platypus with a minimum posterior of 1. SNVs called by Platypus were

864    considered somatic if they had at least 2 alternative reads in the tumour, fewer than 2 alternative

865    reads in the normal, a minimum tumour VAF of 10x the control VAF, and a minimum depth of 10.

866

867    **Germline mutation analysis**

868    *De novo* WGS variants were downloaded from the Gene4Denovo database (Supplementary Table 5).

869    Reference assembly conversion errors were removed by discarding variants where the reference

870    allele did not match the genome reference at the given position or where the variant position was

871    greater than the length of the reference chromosome. In addition, individuals with total *de novo*

872    variants below the 10[th] (*n*=33) or above the 90[th] (*n*=140) percentile were excluded. For germline

873    gene expression we used pre-defined expression groups[80] based on Ensembl release 90 annotation

874    (ftp://ftp.ensembl.org/pub/release-90/gtf/homo_sapiens/Homo_sapiens.GRCh38.90.gtf.gz). Initially

875    stratified as nine expression groups from 1 (=unexpressed) to 9 (=high), we collapsed them into a

876    smaller set of unexpressed (1), low (2, 3, 4), mid (5, 6, 7) and high (8, 9). The annotations were

877    converted to GRCh37 coordinates using liftover (kent source version 417). Genomic segments

878    overlapping multiple distinct expression groups, due to overlapping genes, were assigned to the

879    higher of those expression groups. For each expression group we summed the count (*c*) of *de novo*

880    indels contained within the genomic span of those genes. This was converted to rate estimates by

881    dividing by the union genomic span (*g* nucleotides) of genes in that expression group, and adjusting

882    for the number of mutated genomes considered (*n*); *rate = c/(gn)*. To obtain 95% confidence

883    intervals, gene selection was bootstrapped (sampled to an identical number with replacement) 100

884    times and the 0.025 and 0.975 quantiles of the bootstrapped rate calculation taken as the 95%

885    confidence interval.

886

887    **ICGC Pan-cancer expression analysis**

888    The ICGC PCAWG somatic mutations[50]

889    (https://dcc.icgc.org/api/v1/download?fn=/PCAWG/consensus_snv_indel/final_consensus_passonly

890    .snv_mnv_indel.icgc.public.maf.gz) and ICGC PCAWG "baseline" gene expression[50] were obtained

891    (ArrayExpress https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-5200/). Genomic

892    annotation of gene extents on the GRCh37 reference genome match the Ensembl version 75

893    annotation (http://ftp.ensembl.org/pub/release-

894    75/gtf/homo_sapiens/Homo_sapiens.GRCh37.75.gtf.gz) of the ICGC gene expression calls. Mean,

895    median and maximal gene expression (transcripts per million, TPM) were calculated for each gene

896    across the 76 ICGC baseline gene expression tissues/samples. Only genes annotated on the main

897    autosomal chromosomes, 1 to 22 and the X chromosome were considered. Overlapping genes were

898    removed, keeping only the most abundantly (highest median, then mean in the case of ties)

899    expressed genes from overlapping pairs. This filtering was applied hierarchically, starting with the

900    most abundant. Following [81] genes with housekeeping-like expression were defined as those with

901    maximal expression of less than ten times median expression. Housekeeping-like genes were decile

902    binned into expression groups based on median expression. Mutations were stratified by type (1 bp

903    deletion, 2-5 bp deletion) or by the "TN*T" motif defined below and counted by intersection with

904    the annotated genomic extents of genes in each expression group.

905    For the analysis of tissue-biased gene expression, the 76 ICGC baseline samples were grouped by

906    annotated tissue (e.g. breast, prostate, kidney, liver) and matched where possible to the tissue of

907    origin for ICGC cancer types. For each tissue, the median expression (in TPM) of each gene was

908    calculated for (a) within-tissue samples and (b) for all other samples. The 90th quantile of gene

909    expression (q90, top 10%) within a tissue was set as a threshold for "high" level expression. Genes

910    with high expression in a tissue (a) but a median expression of less than q90*0.1 in the other tissues

911    (b) were considered highly expressed but tissue restricted (HETR). For the set of HETR genes from a

912    tissue, we counted the number of 2-5 bp deletions within the annotated genomic extent of the HETR

913    genes in a cancer type of interest. We similarly counted 2-5 bp deletions in all other genes for that

914    cancer type, and counted both the HETR and non-HETR 2-5 bp deletions from all other cancer types

915    within the ICGC cohort. For each cancer:tissue pair this provided 4 sets of counts, analysed as two-

916    tailed Fisher's exact test using the R fisher.test function. A positive odds-ratio indicating enrichment

917    of 2-5 bp deletions in the HETR genes, compared to a background of the remainder of the ICGC

918    cohort in which HETR genes are not highly expressed. For each cancer type considered, this test was

919    repeated for each tissue type ($n$=17). Analyses were carried out for eight of the ICGC cohort cancer

920    types which met the combined criteria of having a well-matched and known tissue of origin amongst

921    the ICGC baseline samples, and requiring the cancer type cohort to have at least $n$=2,500 2-5 bp

922    deletions in aggregate. This represents $n$=17*8=136 statistical tests, adjusted for by Bonferroni

923    correction. Odds ratios (r) for mutation depletion were transformed to their reciprocal (1/r) for

924    display purposes.

925

926    **ICGC Pan-cancer TOP1-seq analysis**

927    Data corresponding to two replicates of TOP1-seq, a modified ChIP-seq technique to

928    immunoprecipitate only catalytically engaged TOP1[38], were downloaded from the NCBI GEO

929    database (accession code GSE57628, samples GSM1385717 and GSM1385718). Autosomal

930    chromosomes 1 to 22 and the X chromosome were divided into 1 kbp bins and for each bin the

931    amount of mappable sequence was determined using Umap's regions mappable using 36mers[82] to

932    approximate read length of the TOP1-seq data. For each 1 kbp window, the TOP1-seq signal within

933    mappable regions was summed for each replicate and mean signal calculated. This mean was

934    divided by the amount of mappable sequence to calculate the TOP1-seq signal per bp and each 1 kb

935    window was then assigned to decile bins using this value.

936    Somatic deletion calls from ID4-positive PCAWG samples (as defined in

937    https://dcc.icgc.org/api/v1/download?fn=/PCAWG/mutational_signatures/Signatures_in_Samples/S

938    P_Signatures_in_Samples/PCAWG_SigProfiler_ID_signatures_in_samples.csv ) were counted within

939    the same 36-mer mappable regions for each 1kbp window and either stratified by type (1 bp

940    deletion, 2-5 bp deletion) or by the "TN*T" motif defined below. Relative rates of deletions in each

941    category were calculated relative to the first TOP1-seq signal decile.

942

943    **Mutational signatures**

944    *De novo* extraction and decomposition of mutational signatures was performed in Python 3.8.5 using

945    SigProfilerExtractor (v1.1.0)[5], along with SigprofilerMatrixGenerator (v1.1.14/1.1.15)[83] and

946    SigprofilerPlotting (v1.1.27). Recommended default settings (including 500 NMF replicates) were

947    applied (https://github.com/AlexandrovLab/SigProfilerExtractor). Subtraction of mutations in RPE1

948    wildtype cells from those detected in RNase H2 null cells was performed as follows. The average

949    number of indels per line for each of the 83 categories was determined for the three wildtype lines.

950    Counts per category for AKO and BKO lines were subtracted using these averages, with negative

951    values set to 0. SigProfilerExtractor was then performed on the resulting WT-subtracted AKO and

952    BKO ID-83 matrices for both *de novo* signature detection and decomposition analysis.

953

954    **Indel sequence context analysis**

955    WGS indels were categorized based on repeat sequence context. Genome-wide occurrence of short

956    repeats and regions of microhomology were identified and filtered to include only the mappable

957    genome, defined by Umap's regions mappable using 100mers [82]. For both WGS-identified indel

958    variants and genome-wide occurrence, scoring of 2-bp deletions compliant with the "TNT" motif at

959    MH/SSTR sites required the deleted bases to match the sequence NT with a T immediately 5' of the

960    deleted dinucleotides. More generally, for varying sized deletions these were considered to fit a

961    "TN*T" motif if the deletion lay within an SSTR or region of microhomology containing the motif $TN_{(d}$

962    $_{-1)}T$ where $d$ = the length of the deletion. Genome-wide occurrences were estimated from 100,000

963    randomly generated deletions of given lengths within the mappable genome. For SSTRs and MH

964    regions, all regions containing the respective motifs $(TN_{(r-1)})_n$ or $TN_{(r-1)}T$ were identified (where $r$ =

965    the length of the repeat unit and $n > 1$), and the fraction of SSTR/MH sequence containing TNT

966    motifs determined against total SSTR/MH sequence in the mappable genome.

967    To derive a null expectation for *de novo* deletions matching the TNT, TNNT and TNNNT motif for 2, 3,

968    and 4 bp deletions respectively, deletions at repeats from the Gene4Denovo database were first

969    classified by deletion length, repeat type (STR or MH) and repeat length. Bootstrap samples of

970    corresponding repeats from the genome were generated with 1,000 replicates. That is, for each

971    deletion category an equal number of repeats of matching repeat type, repeat unit length and total

972    repeat length were randomly drawn from the genome for each bootstrap sample.

973

974    **Sequence logos**

975    Genomic sequences containing 2 bp deletions were reversed and complemented when the deleted

976    dinucleotide contained an adenosine (A), except when the dinucleotide was AT or TA. For SNMH and

977    STR deletions, the position of the deleted dinucleotide cannot be unequivocally assigned, and

978    therefore the deleted sequence was right aligned in the repeat/microhomology region, either to the

979    most 3' T, where present, or otherwise to the limit of the repeat/microhomology region. Sequences

980    were converted to bit score matrices and logos drawn using Logomaker v0.8[84].

981

982    **Embedded ribonucleotide sequence context analysis**

983    EmRiboSeq data from *rnh201Δ* yeast prepared during mid-log phase growth[85] was obtained

984    (Supplementary Table 5) and aligned to the sacCer3 reference genome as previously described to

985    identify the genomic coordinates of genome embedded ribonucleotides[86]. Bedtools (version v2.30.0,

986    [87]) utilities groupby, slop and getfasta were used to extract and count the sequence context of

987    genome embedded ribonucleotides with downstream analysis and plotting implemented in R

988    (version 4.0.5). Genome sequence composition adjusted relative rates were calculated as previously

989    described[32] such that under the null expectation of no sequence bias in ribonucleotide

990    incorporation, all sequence contexts have an expected relative rate of 1/n where n is the number of

991    contexts considered.

992

993    **Statistical methods**

994    Statistical testing was performed using GraphPad Prism v9.1.1, Python v3.8.5 or R v3.3.1. Two-side

995    non-parametric Mann-Whitney tests were performed for quantitative measurements; multiple

996    testing correction, FDR set at 0.05; and for categorical data Fisher's exact tests were performed in

997    Python using stats.fisher_exact from scipy v1.6.3. Calculation of cosine similarities: Mutations for

998    each strain were converted into a vector, with ordered values representing different mutation

999    categories as a proportion of total mutations. These were then compared in a pairwise fashion.

1000   Given two vectors A and B, the cosine similarity (cos(θ)) was calculated as:

$$\cos(\theta) = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n}(A_i)^2}\sqrt{\sum_{i=1}^{n}(B_i)^2}}$$

1001

1002

1003   Hierachical clustering used the hclust function of R (version 4.1.0) with complete linkage clustering

1004   of pairwise cosine distances (1 - cosine similarity) between ID-83 mutation spectra, with 41

1005   categories of productive reporter frameshift mutations. For bootstrap support, *n*=1,000 bootstrap

1006   datasets were generated by sampling with replacement the mutations observed with a strain, for

1007   each strain; then calculating the cosine distance and hierarchical clustering for each bootstrap

1008    dataset. Reported bootstrap scores are the percentage of bootstrap replicates hierarchical clustering

1009    of which supports the clustering to the right of the indicated position.

1010    To test significance of cosine similarities, we used a null model based on the Dirichlet-multinomial

1011    distribution. Briefly, when comparing two mutation count vectors, with total mutations $m_1$ and $m_2$,

1012    over n mutation classes, we constructed a distribution of cosine values by comparing 10,000

1013    simulated pairs of random vectors generated as follows. For each simulated pair, we sampled from a

1014    Dirichlet-multinomial distribution with the concentration parameters as a vector of ones of

1015    dimension n, and number of trials as $m_1$ for the first vector in the pair, and $m_2$ for the second vector.

1016    The null distribution was obtained by computing the cosine similarity of the 10,000 pairs of mutation

1017    count vectors.

1018

## Additional References

1020    51    Oceguera-Yanez, F. *et al.* Engineering the AAVS1 locus for consistent and scalable transgene
1021          expression in human iPSCs and their differentiated derivatives. *Methods* **101**, 43-55,
1022          doi:10.1016/j.ymeth.2015.12.012 (2016).
1023    52    Gritz, L. & Davies, J. Plasmid-encoded hygromycin B resistance: the sequence of hygromycin
1024          B phosphotransferase gene and its expression in Escherichia coli and Saccharomyces
1025          cerevisiae. *Gene* **25**, 179-188, doi:10.1016/0378-1119(83)90223-8 (1983).
1026    53    Brachmann, C. B. *et al.* Designer deletion strains derived from Saccharomyces cerevisiae
1027          S288C: A useful set of strains and plasmids for PCR-mediated gene disruption and other
1028          applications. *Yeast* **14**, 115-132, doi:10.1002/(SICI)1097-0061(19980130)14:2<115::AID-
1029          YEA204>3.0.CO;2-2 (1998).
1030    54    Spell, R. M. & Jinks-Robertson, S.     3-12 (Springer, 2004).
1031    55    Lea, D. E. & Coulson, C. A. The distribution of the numbers of mutants in bacterial
1032          populations. *Journal of Genetics* **49**, 264-285, doi:10.1007/BF02986080 (1949).
1033    56    Altman, D. G. *Practical statistics for medical research*.  (Chapman and Hall, 1991).
1034    57    Vallot, C., Herault, A., Boyle, S., Bickmore, W. A. & Radvanyi, F. PRC2-independent chromatin
1035          compaction and transcriptional repression in cancer. *Oncogene* **34**, 741-751,
1036          doi:10.1038/onc.2013.604 (2015).
1037    58    Benitez-Guijarro, M. *et al.* RNase H2, mutated in Aicardi-Goutières syndrome, promotes
1038          LINE-1 retrotransposition. *The EMBO journal* **37**, doi:10.15252/embj.201798506 (2018).
1039    59    Reijns, M. A. M. *et al.* The Structure of the Human RNase H2 Complex Defines Key
1040          Interaction Interfaces Relevant to Enzyme Function and Human Disease. *Journal of Biological*
1041          *Chemistry* **286**, 10530-10539, doi:10.1074/jbc.M110.177394 (2011).

1042    60    Lujan, S. A. *et al.* Heterogeneous polymerase fidelity and mismatch repair bias genome
1043          variation and composition. *Genome research* **24**, 1751-1764, doi:10.1101/gr.178335.114
1044          (2014).
1045    61    Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
1046          (2013).
1047    62    Faust, G. G. & Hall, I. M. SAMBLASTER: fast duplicate marking and structural variant read
1048          extraction. *Bioinformatics* **30**, 2503-2505, doi:10.1093/bioinformatics/btu314 (2014).
1049    63    Van der Auwera, G. A. *et al.* From fastQ data to high-confidence variant calls: The genome
1050          analysis toolkit best practices pipeline. *Current Protocols in Bioinformatics* **43**,
1051          doi:10.1002/0471250953.bi1110s43 (2013).
1052    64    Van der Auwera, G. A. & O'Connor, B. D. *Genomics in the Cloud: Using Docker, GATK, and*
1053          *WDL in Terra*.  (O'Reilly Media, 2020).
1054    65    Keane, T. M. *et al.* Mouse genomic variation and its effect on phenotypes and gene
1055          regulation. *Nature* **477**, 289-294, doi:10.1038/nature10413 (2011).
1056    66    Raczy, C. *et al.* Isaac: ultra-fast whole-genome secondary analysis on Illumina sequencing
1057          platforms. *Bioinformatics* **29**, 2041-2043, doi:10.1093/bioinformatics/btt314 (2013).
1058    67    Roller, E., Ivakhno, S., Lee, S., Royce, T. & Tanner, S. Canvas: versatile and scalable detection
1059          of copy number variants. *Bioinformatics* **32**, 2375-2377, doi:10.1093/bioinformatics/btw163
1060          (2016).
1061    68    Robinson, J. T. *et al.* Integrative genomics viewer. *Nat Biotechnol* **29**, 24-26,
1062          doi:10.1038/nbt.1754 (2011).
1063    69    Nik-Zainal, S. *et al.* The life history of 21 breast cancers. *Cell* **149**, 994-1007,
1064          doi:10.1016/j.cell.2012.04.023 (2012).
1065    70    Raine, K. M. *et al.* ascatNgs: Identifying Somatically Acquired Copy-Number Alterations from
1066          Whole-Genome Sequencing Data. *Curr Protoc Bioinformatics* **56**, 15 19 11-15 19 17,
1067          doi:10.1002/cpbi.17 (2016).
1068    71    Sekhon, J. S. Multivariate and Propensity Score Matching Software with Automated Balance
1069          Optimization: The Matching package for R. *Journal of Statistical Software* **42**, 1 - 52,
1070          doi:10.18637/jss.v042.i07 (2011).
1071    72    Benjamin, D. *et al.* Calling Somatic SNVs and Indels with Mutect2. *bioRxiv*, 861054,
1072          doi:10.1101/861054 (2019).
1073    73    McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing
1074          next-generation DNA sequencing data. *Genome Res* **20**, 1297-1303,
1075          doi:10.1101/gr.107524.110 (2010).
1076    74    Kim, S. *et al.* Strelka2: fast and accurate calling of germline and somatic variants. *Nat*
1077          *Methods* **15**, 591-594, doi:10.1038/s41592-018-0051-x (2018).
1078    75    Wala, J. A. *et al.* SvABA: genome-wide detection of structural variants and indels by local
1079          assembly. *Genome Res* **28**, 581-591, doi:10.1101/gr.221028.117 (2018).
1080    76    Rimmer, A. *et al.* Integrating mapping-, assembly- and haplotype-based approaches for
1081          calling variants in clinical sequencing applications. *Nat Genet* **46**, 912-918,
1082          doi:10.1038/ng.3036 (2014).
1083    77    Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**,
1084          doi:10.1093/gigascience/giab008 (2021).
1085    78    Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in
1086          141,456 humans. *Nature* **581**, 434-443, doi:10.1038/s41586-020-2308-7 (2020).
1087    79    Chen, X. *et al.* Manta: rapid detection of structural variants and indels for germline and
1088          cancer sequencing applications. *Bioinformatics* **32**, 1220-1222,
1089          doi:10.1093/bioinformatics/btv710 (2016).
1090    80    Xia, B. *et al.* Widespread Transcriptional Scanning in the Testis Modulates Gene Evolution
1091          Rates. *Cell* **180**, 248-262 e221, doi:10.1016/j.cell.2019.12.015 (2020).

1092 81    Consortium, F. *et al.* A promoter-level mammalian expression atlas. *Nature* **507**, 462-470,
1093       doi:10.1038/nature13182 (2014).
1094 82    Karimzadeh, M., Ernst, C., Kundaje, A. & Hoffman, M. M. Umap and Bismap: quantifying
1095       genome and methylome mappability. *Nucleic Acids Res* **46**, e120, doi:10.1093/nar/gky677
1096       (2018).
1097 83    Bergstrom, E. N. *et al.* SigProfilerMatrixGenerator: A tool for visualizing and exploring
1098       patterns of small mutational events. *BMC Genomics* **20**, 685-685, doi:10.1186/s12864-019-
1099       6041-2 (2019).
1100 84    Tareen, A. & Kinney, J. B. Logomaker: beautiful sequence logos in Python. *Bioinformatics* **36**,
1101       2272-2274, doi:10.1093/bioinformatics/btz921 (2020).
1102 85    Reijns, M. A. M. *et al.* Lagging-strand replication shapes the mutational landscape of the
1103       genome. *Nature* **518**, 502-506, doi:10.1038/nature14183 (2015).
1104 86    Ding, J., Taylor, M. S., Jackson, A. P. & Reijns, M. A. M. Genome-wide mapping of embedded
1105       ribonucleotides and other noncanonical nucleotides using emRiboSeq and EndoSeq. *Nature*
1106       *protocols* **10**, 1433-1444, doi:10.1038/nprot.2015.099 (2015).
1107 87    Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic
1108       features. *Bioinformatics* **26**, 841-842, doi:10.1093/bioinformatics/btq033 (2010).
1109 88    Kornberg, A., Bertsch, L. L., Jackson, J. F. & Khorana, H. G. Enzymatic Synthesis of
1110       Deoxyribonucleic Acid, Xvi. Oligonucleotides as Templates and the Mechanism of Their
1111       Replication. *Proc Natl Acad Sci U S A* **51**, 315-323, doi:10.1073/pnas.51.2.315 (1964).
1112 89    Fan, H. & Chu, J. Y. A brief review of short tandem repeat mutation. *Genomics Proteomics*
1113       *Bioinformatics* **5**, 7-14, doi:10.1016/S1672-0229(07)60009-6 (2007).
1114 90    Stok, C., Kok, Y. P., van den Tempel, N. & van Vugt, M. Shaping the BRCAness mutational
1115       landscape by alternative double-strand break repair, replication stress and mitotic
1116       aberrancies. *Nucleic Acids Res* **49**, 4239-4257, doi:10.1093/nar/gkab151 (2021).
1117 91    Drost, J. *et al.* Use of CRISPR-modified human stem cell organoids to study the origin of
1118       mutational signatures in cancer. *Science* **358**, 234-238, doi:10.1126/science.aao3130 (2017).
1119 92    Hiller, B. *et al.* Mammalian RNase H2 removes ribonucleotides from DNA to maintain
1120       genome integrity. *The Journal of Experimental Medicine* **209**, 1419-1426,
1121       doi:10.1084/jem.20120876 (2012).

1122

## Acknowledgments

1147

## Author Contributions

1148

1149    M.A.M.R, T.C.W., M.S.T. and A.P.J. conceived the project and designed the experiments. T.C.W. and

1150    M.A.M.R, with help from P.C., performed fluctuation assays and sequencing experiments. M.A.M.R.,

1151    with help from P.C., performed the RPE1 mutation accumulation experiment. S.B. performed FISH

1152    experiments. M.A.M.R., T.C.W. and D.O.R.S. performed all other molecular biology experiments. H.X.

1153    and K.A. provided mouse tumour and control tissue samples. D.A.P., T.C.W., M.D.N. and M.S.T.

1154    designed and implemented computational analyses. D.A.P., T.C.W. and M.S.T analysed yeast, mouse,

1155    RPE1 and Gene4Denovo WGS data. D.A.P. and M.S.T. performed pan-cancer analyses. T.G.E.R.C., K.R

1156    and A.S. provided CLL WGS data. A.J.C. provided CRC data. D.A.P., F.N., R.L.H., R.R. and C.P. analysed

1157    CLL data. D.A.P. analysed CRC data. M.A.M.R, C.P., T.S., E.C, M.S.T. and A.P.J. supervised the work.

1158    T.C.W., F.N., E.C., T.S., M.S.T. and A.P.J. funded the work. M.A.M.R. and A.P.J. wrote the manuscript.

1159    All authors had the opportunity to edit the manuscript. All authors approved the final manuscript.

1160

## Data Availability

1172    RPE1 mutation accumulation experiment and mouse tumour WGS data are available from European

1173    Nucleotide Archive accession PRJEB48753 (https://www.ebi.ac.uk/ena/browser/view/PRJEB48753).

1174    All other data were previously published and sources are cited in Supplementary Table 5.

1175

## Code Availability

1177    Code documented in Methods is available at https://git.ecdf.ed.ac.uk/ID-TOP1

1178

## The Genomics England Research Consortium

John C. Ambrose[15], Prabhu Arumugam[15], Roel Bevers[15], Marta Bleda[15], Freya Boardman-Pretty[15,16], Christopher R. Boustred[15], Helen Brittain[15], Mark J. Caulfield[15,16], Georgia C. Chan[15], Greg Elgar[15,16], Tom Fowler[15], Adam Giess[15], Angela Hamblin[15], Shirley Henderson[15,16], Tim J.P. Hubbard[15], Rob Jackson[15], Louise J. Jones[15,16], Dalia Kasperaviciute[15,16], Melis Kayikci[15], Athanasios Kousathanas[15], Lea Lahnstein[15], Sarah E.A. Leigh[15], Ivonne U.S. Leong[15], Javier F. Lopez[15], Fiona Maleady-Crowe[15], Meriel McEntagart[15], Federico Minneci[15], Loukas Moutsianas[15,16], Michael Mueller[15,16], Nirupa Murugaesu[15], Anna C. Need[15,16], Peter O'Donovan[15], Chris A. Odhams[15], Christine Patch[15,16], Mariana Buongermino Pereira[15], Daniel Perez-Gil[15], John Pullinger[15], Tahrima Rahim[15], Augusto Rendon[15], Tim Rogers[15], Kevin Savage[15], Kushmita Sawant[15], Richard H. Scott[15], Afshan Siddiq[15], Alexander Sieghart[15], Samuel C. Smith[15], Alona Sosinsky[15,16], Alexander Stuckey[15], Mélanie Tanguy[15], Ana Lisa Taylor Tavares[15], Ellen R. A. Thomas[15,16], Simon R. Thompson[15], Arianna Tucci[15,16], Matthew J. Welland[15], Eleanor Williams[15], Katarzyna Witkowska[15,16], Suzanne M. Wood[15,16].


[15] Genomics England, London, UK

[16] William Harvey Research Institute, Queen Mary University of London, London, UK


## Colorectal Cancer Domain UK 100,000 Genomes Project

Daniel Chubb[9], Alex Cornish[9], Ben Kinnersley[9], Richard Houlston[9], David Wedge[17], Andreas Gruber[17], Anna Frangou[18], William Cross[19], Trevor Graham[20], Andrea Sottoriva[9], Gulio Caravagna[9], Nuria Lopez-Bigas[21], Claudia Arnedo Pac[21], David Church[18], Richard Culliford[9], Steve Thorn[22], Phil Quirke[23], Henry Wood[23], Ian Tomlinson[22], Boris Noyvert[5]

1202 [17] Manchester Interdisciplinary Biocentre, University of Manchester, Manchester, UK

1203 [18] Wellcome Centre for Human Genetics, Oxford, UK

1204 [19] Cancer Institute, University College London, London, UK

1205 [20] Barts Cancer Institute, Barts and The London School of Medicine and Dentistry, Queen Mary

1206 University of London, London, UK

1207 [21] Institute for Research in Biomedicine (IRB Barcelona), The Barcelona Institute of Science and

1208 Technology (BIST), Barcelona, Spain

1209 [22] Edinburgh Cancer Research Centre, IGC, The University of Edinburgh, Edinburgh, UK

1210 [23] Pathology and Data Analytics, Leeds Institute of Medical Research, St James's University Hospital,

1211 University of Leeds, Leeds, UK

1212

## Extended data figure legends

1214 **Extended Data Fig. 1 | ID4 is distinct from small deletion signatures of known aetiology. a,b,** The

1215 mechanistic basis for many COSMIC indel signatures is unknown, with only 9 out of 18 having a

1216 proposed aetiology. ID2 (**a**) is attributed to DNA polymerase slippage[88,89] and ID6 (**b**) to

1217 microhomology mediated end-joining (MMEJ) activity, associated with HR deficiency[5,90]. **c,d,**

1218 Mechanism for these signatures supported by: impaired MMR promoting replication slippage

1219 mutagenesis in MLH1-/- colonic organoids resulting in ID2 (and ID1) signatures (**c**); ID6 contributing

1220 substantially (along with ID8) to the indel signature in ovarian cancer, in which HR deficiency is

1221 common (**d**). Analysis of data from [91] in **c**; data for 73 ovarian adenocarcinomas with ID6

1222 contribution from ICGC[5,50] in **d**.

1223

1224 **Extended Data Fig. 2 | Yeast and human frameshift mutation reporters detect indels at tandem**

1225 **repeats. a,** Yeast reporter. Synonymous substitutions were made in the hygromycin resistance gene

1226    (HygroR), such that it contained many short 2 bp tandem repeats (SSTRs). Expression from the TEF

1227    promoter (P$_{TEF}$) ensures a constitutive high level of transcription. Mutations within HygroR that

1228    result in a frameshift simultaneously put the HygroR coding sequence out of frame and the

1229    downstream neomycin resistance (NeoR) sequence in frame, allowing antibiotic selection of cells

1230    with such mutations. **b,** Top1-dependent 2 bp SSTR deletions occur in both WT and *rnh201Δ* (RNase

1231    H2 null) yeast, with the highest mutation rate for *rnh201Δ* (related to Fig. 1d). **c-e,** WT and *rnh201Δ*

1232    have similar spectra, and differ from *top1Δ* strains. Mutation spectra of neomycin resistant colonies.

1233    n, number of independent colonies sequenced. Other: complex indels, missense mutations or

1234    mutation not characterised (**c**). Tree for pairwise clustering with percent bootstrap support to the

1235    right of the indicated position, based on cosine scores calculated for mutation spectra (Fig. 1e) of the

1236    41 mutation categories that give productive reporter frameshift mutations (**d**). Matrix of pairwise

1237    cosine similarities and P-values between reporter mutation spectra in different yeast strains. Darker

1238    blue indicates greater similarity; darker grey greater significance. Test statistic is the cosine similarity

1239    value for 41 mutation categories and the null hypothesis is that that the cosine value will be

1240    distributed according to the Dirichlet-multinomial model, as described in Methods. The test is one-

1241    sided and no adjustments were made for multiple comparisons (**e**). **f,** Null distribution for cosine

1242    pairwise vector comparisons for 41 and 83 mutation categories. Plots, cosine values for 10,000

1243    randomly generated pairs of vectors of mutation spectra. Each vector contained 100 randomly

1244    assigned mutations (see Methods for further details). Cosine value thresholds indicated for *P* < 0.05

1245    and *P* < 0.01. **g,** The human reporter is expressed from the ubiquitous CAG promoter (P$_{CAG}$) , and

1246    NeoR is replaced with the puromycin resistance gene (PuroR) to allow more rapid antibiotic

1247    selection in mammalian cell culture.

1248

1249    **Extended Data Fig. 3 | Validation and characterisation of RNASEH2A+ and KO HeLa reporter cells.**

1250    **a-c,** Reporter integration at the *AAVS1* locus and retention of a reporter-free locus with a 200 bp

1251    deletion at the target site was confirmed by PCR and Sanger sequencing. Green arrow head, specific

1252    PCR product. Representative of at least 2 independent experiments. **d,e,** FISH shows integration of

1253    the reporter (**d**) at a single *AVS1* locus (**e**). Representative image of approximately one hundred

1254    mitotic chromosome spreads in 3 independent experiments. SA, splice acceptor; T2A, self-cleaving

1255    peptide; pA, polyadenylation site; also see Fig. 2a. **f,g,** Alkaline gel electrophoresis of RNase H2

1256    treated genomic DNA (**f**) shows a small increase in fragmentation for the RNASEH2A+ control clone

1257    and a more substantial increase in two independent RNASEH2A-KO clones (representative of 4

1258    independent experiments), indicating the presence of more genome-embedded ribonucleotides

1259    compared to HeLa and parental reporter cells (**g**). "Control KO" cells were reported previously[33,58].

1260    RFU, relative fluorescence units. **h,** 2 bp SSTR deletions are frequent in both RNASEH2A+ and KO

1261    cells. Mutation spectra, quantitation of indel type. Relative area of pie charts scaled to mutation

1262    rate. n, number of colonies sequenced from independent cultures. Other: complex indels or

1263    missense mutations.

1264

1265    **Extended Data Fig. 4 | RPE1 RNase H2 null cells accumulate embedded ribonucleotides and 2-5 bp**

1266    **deletions across the genome. a,b,** *RNASEH2A* and *RNASEH2B* KO cells (AKO, BKO, respectively) have

1267    substantially reduced cellular levels of RNase H2 subunits (**a**) and are deficient for RNase H2 enzyme

1268    activity (**b**) at the outset (ancestral) and at the end of the mutation accumulation experiment (end

1269    point). Individual data points, *n*=3 technical replicates; mean ± s.d. For gel source data, see

1270    Supplementary Fig. 1. **c,d,** Alkaline gel electrophoresis of RNase H2 treated genomic DNA (**c**) shows a

1271    substantial increase in fragmentation for *RNASEH2A* and *RNASEH2B* KO clones (representative of 3

1272    independent experiments), indicating the presence of more genome-embedded ribonucleotides

1273    compared to two WT control clones (**d**). Densitometry plots of **c**. RFU, relative fluorescence units. As

1274    RNase H2 deficiency activates the p53 pathway[14,92], experiments were performed in a *TP53* knockout

1275    background. **e,** Only 2-5 bp deletions are significantly increased in RNase H2 null cells. Data points

1276     for acquired indel mutations in individual cell lines after 100 population doublings. Individual data

1277     points, indel counts per cell line; mean ± s.d.; P-values for two-sided Fisher's exact test between WT

1278     (*n*=3 independent clones) and KO (*n*=2 independent clones) for one indel type vs all other indel

1279     types, after Bonferroni correction. **f,** Proportions of acquired indels in WT and KO RPE cells. After

1280     correction for indels occurring in WT, 69% of indels in RNase H2 null cells are 2-5 bp deletions. n,

1281     total indel counts. **g,** Quantification of 2 bp deletions by context. n, total number of 2 bp deletions.

1282     For **f** and **g**, chart areas scaled to mutation counts per line.

1283

1284     **Extended Data Fig. 5 | ID4 occurs in RNase H2 null RPE1 cells, particularly in transcribed regions. a-**

1285     **d,** Mutational spectra detected by WGS after 100 population doublings in RPE1 cells demonstrates

1286     that SSTR and SNMH/MH deletions are enriched in RNase H2 cells. Spectra for combined RNase H2

1287     null and wildtype cell lines (**a**), and individual cell lines (**b**). Mutational signature analysis confirms

1288     ID4 contribution in RNase H2 null (**c**), but not WT cells (**d**). **e,** In RNase H2 null cells, ID4 contributes

1289     significantly more to indel mutations in transcribed genomic regions ($P$=1.3x10$^{-29}$). Two-sided

1290     Fisher's exact test, ID4 indels vs other indels.

1291

1292     **Extended Data Fig. 6 | ID4 mutations in RNase H2 null mouse tumours and RPE1 cells occur at a**

1293     **TNT motif, defining ID-TOP1. a,** Mutation spectra for individual Rnaseh2b-KO murine intestinal

1294     tumours (WGS, paired tumour-normal samples from 6 mice). **b,** Indel classes, detected in murine

1295     Rnaseh2b-KO tumours. n, total indel count for 6 tumours. **c,** Most 2 bp deletions in these tumours

1296     occur at SSTRs and sites of single nucleotide microhomology (SNMH). n, number of 2 bp deletions.

1297     **d,e,** A TNT sequence motif is present at all 2 bp STR and SNMH deletions in RNase H2 null mouse

1298     tumours (**d**) and RPE1 cells (**e**). Related to Fig. 4d and Fig. 3, respectively. Sequence logo: 2-bit

1299     representation of the sequence context of 2 bp deletions. Top, all deletions, with those sequences

1300     containing a deleted adenosine (except AT/TA) reverse complemented, and deletions right-aligned.

1301     Middle, re-aligned on right-hand T. Bottom, aligned on T (STR and SNMH context only). n, number of

1302     deletions. **f,** Deletion sites in RNase H2 null RPE1 cells are significantly enriched for the TNT

1303     sequence motif compared to genome-wide occurrence, for all genome sequence, as well as SNMH

1304     sites. P-values, two-sided Fisher's exact, observed vs expected. $n$=98 (all; $P$=8.3x10$^{-14}$), 54 (STR;

1305     $P$=0.057), 30 (SNMH; $P$=0.0008) deletions.

1306

1307     **Extended Data Fig. 7 | ID4 deletions in RNase H2 null *S. cerevisiae* occur at a TNT motif in a Top1-**

1308     **dependent manner. a,** 2 bp deletion sites in *rnh201Δ pol2-M644G* yeast are significantly enriched

1309     for the TNT sequence motif compared to genome-wide occurrence, for all genome sequence, as well

1310     as STR sites. P-values, two-sided Fisher's exact, observed vs expected. $n$=94 (all; $P$=1.0x10$^{-9}$), 91 (STR;

1311     $P$=0.029), 3 (SNMH; $P$=1) deletions. **b,** A TNT sequence motif is present at all 2 bp STR and SNMH

1312     deletions in *rnh201Δ pol2-M644G* yeast. Sequence logo: 2-bit representation of the sequence

1313     context of 2 bp deletions. Top, all deletions, with those sequences containing a deleted adenosine

1314     (except AT/TA) reverse complemented, and deletions aligned on right-hand T. Bottom, aligned on T

1315     (STR and SNMH context only). n, number of deletions. **c,d,** TN*T motifs extend beyond 2 bp

1316     deletions, with enrichment above expectation for 2 bp deletions at TNT, 3 bp deletions at TNNT and

1317     4 bp deletions at TNNNT motifs in *rnh201Δ pol2-M644G* yeast WGS. Null expectations were

1318     generated by randomly simulating deletions of 2, 3 and 4 bp (**c**) or 2 bp STR sequences (**d**) genome-

1319     wide and scoring those simulated events for TN*T compliance. Each simulated dataset matched the

1320     count of observed mutations for the corresponding deletion class and $n$=1,000 replicate simulated

1321     datasets were produced. The frequency distribution of TN*T compliance in simulations is plotted as

1322     histograms, and comparison to the observed frequency of TN*T compliance (dotted red lines) used

1323     to derive a two-tailed empirical P-value. **e,** 2 bp STR deletions have biased sequence composition.

1324     Deletion observed in *rnh201Δ pol2-M644G* yeast WGS. Genome, frequency of dinucleotides in STR

1325     sequences in mappable genome. **f,** Ribouridine (rU) is more common in a CrU/GrU than in an

1326  ArU/TrU dinucleotide context. Genome-embedded ribonucleotide frequency determined by

1327  emRiboSeq[86]. Dotted line indicates relative rate in absence of bias (=0.25). Horizontal lines, mean;

1328  individual data points, values for *n*=4 independent experiments[85]. **g,h,** 2 bp TNT deletions in wildtype

1329  and RNase H2 null cells are dependent on Topoisomerase 1. Mutation rates for 2 bp deletions at

1330  TNT-compliant SSTRs (**g**). Deletions at TNT motifs are significantly increased above expectation in WT

1331  and *rnh201Δ*, but not in *top1Δ* and *rnh201Δ top1Δ* yeast. Horizontal bars, 95% confidence intervals

1332  for odds ratio estimates (diamonds). P-values, two-sided Fisher's exact after Bonferroni correction;

1333  *n*=86, 28, 103, 19 2-bp deletions, with each deletion from an independent culture, for WT, *top1Δ*,

1334  *rnh201Δ*, *rnh201Δ top1Δ*, respectively. Null expectation, random occurrence of mutations in

1335  reporter target sequence (**h**).

1336

1337  **Extended Data Fig. 8 | TOP1-mediated mutagenesis causes increased 2-5 bp deletions in cancer. a,**

1338  Of all indels, only 2-5 bp deletions are significantly increased in CLL with biallelic *RNASEH2B* loss.

1339  Box, 25-75%; line, median; whiskers 5-95% with data points for values outside this range. WT (2

1340  copies), *n*=201; monallelic loss (1 copy), *n*=131; biallelic loss (0 copies), *n*=16 independent tumours.

1341  Indels as percentage of all variants per sample (GEL and ICGC data combined). q-values, 2-sided

1342  Mann-Whitney test with 5% FDR. **b,c,** In RNase H2 null CLL, 2 bp deletions predominantly occur at

1343  STR and SNMH sequences (**b**), and at the TNT sequence motif (**c**), consistent with TOP-mediated

1344  mutagenesis. Mean ± s.e.m., percentage of all variants per sample. GEL and ICGC data combined. *n*=

1345  1,711; 1,244; 443 2-bp indels identified in 201, 131, 16 biologically independent tumours,

1346  respectively. **d,** ID4 contribution in RNase H2 null CLL is greater in transcribed regions. Two-sided

1347  Fisher's exact test, ID4 indels vs other indels ($P$=9.2x10$^{-16}$). **e,** Pan-cancer transcript expression data

1348  divided into ten expression strata for ubiquitously expressed genes (used in panel **h** and Fig. 5b

1349  analysis). Data points, median/maximum expression across cancer types for individual genes. Genes

1350  with similar median and maximum TPMs were considered to be ubiquitously expressed and divided

1351    into expression groups from low (1) to high (10) expression. **f,** Two bp deletions in cancer

1352    preferentially occur at STRs. **g,** ID-TOP1 deletions increase in frequency with TOP1 cleavage activity

1353    (measured by TOP1-Seq; [38]). Dotted line, relative rate in lowest TOP1-seq category set to 1. Solid

1354    lines, relative deletion rate. ID-TOP1, 2-5 bp MH and SSTR deletions containing the TN*T sequence

1355    motif. **h,** ID-TOP1, but not deletions in other sequence contexts, correlate with transcription. **i,** 2-5

1356    bp deletions from prostate adenocarcinoma are most enriched amongst the top 10% of highly

1357    expressed prostate 'tissue-restricted' genes. Odds ratio (OR): number of 2-5bp deletions in top 10%

1358    tissue restricted genes vs deletions in other genes, relative to expected frequency from all other

1359    tissues. **j,** ID4 is not detected in the indel signature of irinotecan-treated colorectal cancers.

1360    Untreated (*n*=78), treated (*n*=39). **k,** 2-5 bp deletion frequency in cancer corresponds to TOP1

1361    cleavage activity, in both genic and non-genic regions. Data analysed from PCAWG[50], all tumours in

1362    **e**, **h**; ID4 positive tumours in **g**, **k**; Genomics England in **j**. In **g**, **h** and **k**, solid line, relative deletion

1363    rate; shading indicates 95% confidence intervals from 1,000 (**g**,**k**) or 100 (**h**) bootstrap replicates.

1364

1365    **Extended Data Fig. 9 | Human germline *de novo* indels are enriched for ID-TOP1 deletions. a,** Most

1366    *de novo* 2 bp deletions occur at SSTR, STR and SNMH sequences. **b,c,** A TNT sequence motif is

1367    present at the majority of 2 bp STR and SNMH deletions (**b**). Sequence logos: 2-bit representation of

1368    the sequence context of 2 bp deletions. Top, all deletions, with those containing A (except AT/TA)

1369    reverse complemented, and deletions right-aligned on T (where present). Bottom, STR/SNMH

1370    deletions only (**c**). **d,** TN*T motifs extend beyond 2 bp deletions, with enrichment above expectation

1371    for 2 bp deletions at TNT, 3 bp deletions at TNNT and 4 bp deletions at TNNNT motifs ($P < 0.001$;

1372    two-tailed empirical P-value determined for each category). Bootstrap sampling (*n*=1,000) of 2, 3

1373    and 4 bp STR/MH sequences genome-wide to derive expected frequencies of those matching TN*T

1374    motifs. Sampling was performed to match the numbers of deletions at repeats observed in the

1375    Gene4Denovo database for each category defined by repeat type, repeat unit length and total

1376    repeat length. Histograms, distribution of the number of repeats matching TN*T motifs over these

1377    samplings. Solid blue lines, kernel density estimates for these distributions. Dotted red lines, number

1378    of deletions observed in Gene4Denovo matching TN*T motifs for each category. **e,** ID-TOP1

1379    correlates with germline expression level. ID-TOP1, defined as 2-5 bp MH and SSTR deletions

1380    containing the TN*T sequence motif. Shading, 95% confidence intervals from 100 bootstrap

1381    replicates.

1382

1383    **Extended Data Fig. 10 | Topoisomerase 1 causes small deletions while protecting against**

1384    **topological stress. a,** The canonical role of Topoisomerase 1 (TOP1) is to relieve torsional stress (sc,

1385    supercoiling) during replication and transcription. **b,** TOP1 acts by forming ssDNA nicks to release

1386    supercoils and then religates the relaxed DNA. However, TOP1 cleavage at genome-embedded

1387    ribonucleotides (frequently incorporated by replicative polymerases such as Pol ε), can lead to short

1388    deletions that will be most frequent at sites of torsional stress in the genome, such as occurs at

1389    highly transcribed genes. Adapted from [6].

**a**

COSMIC ID4 signature

**b**

*S. cerevisiae rnh201Δ pol2-M644G*  **100% ID4** (0.78 cosine similarity; 148 indels)

**c**

**d**

**e**

**a**

Reporter
HygroR | TAA TT | P2A | PuroR
STOP

Targeting construct
HA-L | SA | T2A | NeoR | pA | P_CAG | Reporter | pA | HA-R

CRISPR/Cas9

HeLa
AAVS1 locus

**b**

| | HeLa | Parental | RNASEH2A+ | KO1 | KO2 | |
|---|---|---|---|---|---|---|
| RNASEH2A | | | | | | 37 kDa |
| RNASEH2B | | | | | | 37 kDa |
| RNASEH2C | | | | | | 20 kDa |
| GAPDH | | | | | | 37 kDa |

**c**

RNase H2 activity (% of HeLa)

HeLa  Parental  RNASEH2A+  KO1  KO2

**d**

RNASEH2A+
KO
3.1x

mutation rate (x 10^-10 per bp per generation)
0  50  100  150  200  250

**e**

1 bp del    1 bp ins    >1 bp del (repeats)    >1 bp ins (repeats)    del (MH)

Proportion of indels

RNASEH2A+          n = 31
SSTR          SNMH

RNASEH2A-KO          n = 46

homopolymer length    Number of repeat units    Number of repeat units    MH length

**a**

WT  AKO  BKO  hTERT-RPE1

initial single cell sort

WT1  WT2  AKO  BKO  ancestral (WGS)

single cell sort
1
2
3
4

4 x 25 doublings

"offspring" after ~100 doublings (WGS)

**b**

Deletion events

$P = 4 \times 10^{-17}$

WT  KO    WT  KO
  1 bp       2-5 bp

**c**

ID4  ID1  ID2  ID5

%ID4

WT        0%
KO        21%
subtracted  61%

0  20  40  60  80  100
% contribution

**d**

1 bp deletion | 1 bp insertion | >1 bp deletions at repeats (Deletion length) | >1 bp insertions at repeats (Insertion length) | Deletions with microhomology (Deletion length)

C  T  C  T  2  3  4  5+  2  3  4  5+  2 3 4  5+

Percentage of indels

16%

**61% ID4**, 39% ID5 (0.89 cosine similarity; 212 indels)

12%

8%

4%

0

Background subtracted

1 2 3 4 5 6+  1 2 3 4 5 6+  0 1 2 3 4 5+  0 1 2 3 4 5+  1 2 3 4 5 6+  1 2 3 4 5 6+  1 2 3 4 5 6+  1 2 3 4 5 6+  0 1 2 3 4 5+  0 1 2 3 4 5+  0 1 2 3 4 5+  0 1 2 3 4 5+  1 1 2 1 2 3  1 2 3 4 5+

homopolymer length | Number of repeat units | Number of repeat units | Microhomology length

**a**

1 bp del    1 bp ins    >1 bp del, repeats (del length)    >1 bp ins, repeats (ins length)    del, MH (del length)

C   T    C   T    2   3   4   5+    2   3   4   5+    2   3   4   5+

Percentage of indels

28%
21%
14%
7%
0

**32% ID4** (0.98 cosine similarity; 989 indels)

homopolymer length     Number of repeat units     MH length

**b**

ID4   ID1   ID2   ID5   ID7    indels (n=)   ID4

Genome — 989 — 32%
Transcribed — 418 — 42%
Untranscribed — 571 — 31%

$P = 4.7 \times 10^{-4}$

% contribution
0   20   40   60   80   100

**c**

% of 2 bp repeat deletions

60
50
40
30
20
10
0

example:

CTCT / GAGA    Deletion / Genome   STR
NTCT / NAGA    Deletion / Genome   SNMH

deleted bases: CT   GT   TT   AT   CC   CG

**d**

n = 201 deletions
STR/SNMH

deletion

bits
2.0
1.0
0.0

-5   -4   -3   -2   -1   0   1   2   3   4   5

deletion site

**e**

$P < 10^{-27}$    $P = 0.0008$    $P < 10^{-7}$

% of 2 bp deletions
100
75
50
25
0

TNT   Other

Deletion site / Genome-wide   All
Deletion site / Genome-wide   STR
Deletion site / Genome-wide   SNMH

**f**

(i) Top1 cleaves at rN

(ii) Nucleophilic attack by 2'OH

(iii) 2',3'-cP formation & Top1 release

(iv) 2 nt release & strand realignment

(v) Religation of single strand nick

(vi) Top1 release — 2 bp (NT) deletion

**a**

q = 0.0013
q = 0.0077

q < 0.0002
q = 0.0023

2-5 bp deletion events

RNASEH2B  WT  het  null    WT  het  null
CLL - GEL              CLL - ICGC

**b**

PCAWG, transcription

Deletion rate relative to genome

— 2–5 bp deletions
— 1 bp deletions

Housekeeping gene expression category
(low to high)

**c**

$P < 10^{-200}$  $< 10^{-59}$  $< 10^{-50}$

% of 2 bp deletions

All    STR    SNMH

Deletion site  Genome-wide  Deletion site  Genome-wide  Deletion site  Genome-wide

■ TNT
■ Other

**d**

PCAWG, TOP1 activity

Relative rate

— 2–5 bp deletions
— 1 bp deletions

Top1-seq signal category
(low to high)

**e**

Breast-AdenoCa
Prost-AdenoCa
Kidney-RCC
Liver-HCC
Panc-AdenoCa
Eso-AdenoCa
Stomach-AdenoCa
Ovary-AdenoCa

breast prostate kidney liver pancreas esophagus stomach ovary bladder brain cervix intestine lung skeletal muscle skin thyroid uterus

**Enriched**
■ $P$<0.01, nominal
■ $P$<0.01, Bonferroni corrected

**Depleted**
■ $P$<0.01, nominal
■ $P$<0.01, Bonferroni corrected

• No significant change

3  5  7  9  11
Odds ratio

**f**

Germline

deletion   1 bp
           2-5 bp    33 %
           > 5 bp

insertion  1 bp
           2-5 bp
           > 5 bp

0    5,000    10,000    15,000
De novo indels

**g**

$P < 10^{-200}$  $< 10^{-46}$  $< 10^{-25}$

% of 2 bp deletions

All    STR    SNMH

Deletion site  Genome-wide  Deletion site  Genome-wide  Deletion site  Genome-wide

■ TNT
■ Other

**h**

x 10$^{-3}$ indels per Mbp
per individual

Indel Type
— 1 bp deletion
— 1 bp insertion
— 2-5 bp deletion
— 2-5 bp insertion
— >5 bp deletion
— >5 bp insertion

None  Low  Mid  High
Germline expression

**a**

**Replication slippage**

| 1 bp deletion | | 1 bp insertion | | >1 bp deletions at repeats | | | | >1 bp insertions at repeats | | | | Deletions with microhomology | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| C | T | C | T | 2 | 3 | 4 | 5+ | 2 | 3 | 4 | 5+ | 2 | 3 | 4 | 5+ |

ID2

Percentage of indels: 104%, 78%, 52%, 26%, 0

homopolymer length          Number of repeat units          Microhomology length

**b**

**Microhomology-mediated end joining (MMEJ)**

ID6

Percentage of indels: 32%, 24%, 16%, 8%, 0

**c**

**MMR deficiency**

Percentage of indels: 92%, 69%, 46%, 23%, 0

**91% ID2**, 9% ID1 (0.998 cosine similiarity; 21,457 indels)

**d**

**HR deficiency**

Percentage of indels: 24%, 18%, 12%, 6%, 0

**48% ID6**, 32% ID8 (0.995 cosine similarity; 55,595 indels)

**a**

$P_{TEF}$ >50% SSTRs STOP
HygroR TAA TT P2A NeoR

↓ mutagenesis

2 bp del GO
Hygro TA ATT P2A NeoR

2+(3n) bp del GO
TA ATT P2A NeoR

1 bp or 1+(3n) bp ins GO
gibberish N TA ATT P2A NeoR

missense GO
ATG gibberish TA ATT P2A NeoR

**b**

2 bp deletion rate, at SSTRs
(x $10^{-12}$ per bp, per generation)

top1Δ, WT (35x), rnh201Δ, rnh201Δ top1Δ (52x)

**c**

WT — 67%, n = 112
rnh201Δ — 92%, n = 105
rnh201Δ top1Δ — 15%, n = 101
top1Δ — 20%, n = 97

- 2 bp del (SSTR)
- 2 bp del (SNMH)
- 2 bp del (no repeat)
- >2 bp deletion
- 1 bp insertion
- >1 bp insertion
- Other

**d**

WT — 100%
rnh201Δ
top1Δ — 97.9%
rnh201Δ top1Δ

pairwise cosine distance
(1 - cosine similarity)
0.8 0.6 0.4 0.2 0.0

**e**

cosine similarities

|          | WT   | rnh201Δ | top1Δ |
|----------|------|---------|-------|
| rnh201Δ  | 0.98 |         |       |
| top1Δ    | 0.31 | 0.21    |       |
| rnh201Δ top1Δ | 0.46 | 0.38 | 0.82 |

P-values

|          | WT      | rnh201Δ | top1Δ   |
|----------|---------|---------|---------|
| rnh201Δ  | <0.0001 |         |         |
| top1Δ    | 0.90    | 0.99    |         |
| rnh201Δ top1Δ | 0.38 | 0.67 | <0.0001 |

**f**

100 mutations, 41 categories
0.61 — P < 0.05
0.68 — P < 0.01
Percentage / Cosine value

100 mutations, 83 categories
0.5 — P < 0.05
0.56 — P < 0.01
Cosine value

**g**

$P_{CAG}$ >50% SSTRs STOP
HygroR TAA TT P2A PuroR

↓ mutagenesis

2 bp del GO
Hygro TA ATT P2A PuroR

**a**

left arm | internal | wildtype

**b**

right arm

**c**

**left arm - reporter**

dna803

HA-L | SA T2A | neo

dna804

|← ---- 1.2 kb ---- →|

**AAVS1 - wildtype**

dna803

HA-L | HA-R

dna183

|← ---- 1.4 kb ---- →|

**reporter - right arm**

puro-F | puro-3F
pA

HygroR | PuroR | HA-R

P2A | dna183 | HA-doR

|← -- 1.4 kb --→| internal

|← ---- 1.7 kb ---- →| right arm

**d**

**Reporter construct**

HA-L | SA | T2A | neo | pA | CAG | REPORTER | pA | HA-R

|← ----- 5.3 kb ----- →|

FISH probe

**e**

Reporter probe
DAPI

10 μm

**f**

**g**

— HeLa
— Parental
— RNASEH2A+
— KO1
— KO2
— Control KO
---- Ladder

RFU

kb (gel migration)

**h**

RNASEH2A +   RNASEH2A-KO

47%   79%

n = 34   n = 46

■ 2 bp del (SSTR)
■ 2 bp del (SNMH)
□ 2 bp del (no repeat)
■ >2 bp deletion
■ 1 bp insertion
■ >1 bp insertion
■ Other

**a**

| | Ancestral | | | | End point | | | |
|---|---|---|---|---|---|---|---|---|
| | WT1 | WT2 | AKO | BKO | WT1 | WT2 | AKO | BKO |

RNASEH2A — 37 kDa

RNASEH2B — 37 kDa

RNASEH2C — 20 kDa

GAPDH — 37 kDa

**b**

RNase H2 activity (% of WT1 ancestral)

WT1 WT2 AKO BKO — Ancestral
WT1 WT2 AKO BKO — End point

**c**

| | Ancestral | | | | End point | | | |
|---|---|---|---|---|---|---|---|---|
| | WT1 | WT2 | AKO | BKO | WT1 | WT2 | AKO | BKO |

kb
48.5
20
15
10
8
6
5
4
3
2
1.5
1
0.5

**d**

RFU

48.5   15   8   5   3   1.5   0.5
kbp (gel migration)

— WT1
— WT2
— RNASEH2A-KO
— RNASEH2B-KO
---- Ladder

**e**

Number of indel events

$P = 4 \times 10^{-17}$

$P = 0.17$

WT KO — 1 bp deletion
WT KO — 1 bp insertion
WT KO — 2-5 bp deletion
WT KO — 2-5 bp insertion
WT KO — >5 bp deletion
WT KO — >5 bp insertion

**f**

WT
3%

KO
16%

KO (WT-subtracted)
69%

n = 933 (3 lines)    n = 765 (2 lines)

- 1 bp deletion
- 1 bp insertion
- 2-5 bp deletion
- 2-5 bp insertion
- > 5 bp deletion
- > 5 bp insertion

**g**

WT

KO
48%   31%

n = 13 (3 lines)    n = 98 (2 lines)

2 bp deletion context
- SSTR (<5 repeats)
- STR ≥5 repeats
- microhomology (SNMH)
- no repeat

**a**

Tumour 1 (160 indels)
Tumour 2 (469 indels)
Tumour 3 (19 indels)
Tumour 4 (136 indels)
Tumour 5 (173 indels)
Tumour 6 (32 indels)

1 bp del    1 bp ins    >1 bp del at repeats (del length)    >1 bp ins at repeats (ins length)    del with MH (del length)

Percentage of indels

homopolymer length    Number of repeat units    Number of repeat units    MH length

**b**

n = 989

- 1 bp deletion (49%)
- 1 bp insertion
- 2-5 bp deletion (28%)
- 2-5 bp insertion
- > 5 bp deletion
- > 5 bp insertion

**c**

n = 228

2 bp deletion context
- SSTR (<5 repeats) (51%)
- STR ≥5 repeats
- microhomology (SNMH) (34%)
- no repeat

**d**

Mouse tumours 2 bp deletions

bits

n = 228    All
n = 228    All (T-aligned)
n = 201    STR/SNMH

deletion site

**e**

RPE1 2 bp deletions

n = 98    All
n = 98    All (T-aligned)
n = 84    STR/SNMH

deletion site

**f**

% of 2 bp deletions

TNT
Other

$P = 8.3 \times 10^{-14}$    $P = 0.057$    $P = 0.0008$

All    STR    SNMH

Deletion site    Genome-wide

**a**

P = 10⁻⁹    0.029    1.0

**b** *S. cerevisiae* 2 bp deletions

deletion

All (T-aligned)    n = 94

STR/SNMH    n = 91

deletion site

**c**

sampling, genome-wide occurrence

deletions observed

91 of 94 deletions
P < 0.0001

2 bp deletions

20 of 24 deletions
P < 0.0001

3 bp deletions

8 of 8 deletions
P = 0.0005

4 bp deletions

**d**

sampling, genome-wide occurrence

deletions observed

88 of 88 deletions
P = 0.0088

2 bp STR deletions

**e**

STR
Deletions observed
Genome

CT GT TT AT CC CG

**f**

Relative rate

ArU CrU GrU TrU

rU context

**g**

WT
*top1Δ*   78x
*rnh201Δ*
*rnh201Δ top1Δ*   1190x

TNT mutation rate
(x 10⁻¹² per bp, per generation)

**h**

WT    P = 6.4 x 10⁻⁹
*top1Δ*    P = 1.0
*rnh201Δ*    P = 2.5 x 10⁻³²
*rnh201Δ top1Δ*    P = 0.96

Odds ratio (TNT observed/expected)

**a**

2 bp deletion context
- SSTR (<5 repeats)
- STR ≥5 repeats
- microhomology (SNMH)
- no repeat

26%
45%

n = 5,569

**b**

2 bp deletion events

TNT
Other

STR    SNMH    No repeat

**c**

2 bp deletions
n = 5,569

deletion

bits

deletion site

All

n = 4,387

STR/SNMH

-5 -4 -3 -2 -1 0 1 2 3 4 5

**d**

sampling, genome-wide occurrence    deletions observed

2 bp          3 bp          4 bp

STR

counts

3075 3100 3125 3150 3175 3200 3225 3250
1820 1830 1840 1850 1860
1230 1235 1240 1245 1250 1255

MH

counts

700 750 800 850 900 950
events (TNT)

780 800 820 840 860
events (TNNT)

1650 1700 1750 1800 1850
events (TNNNT)

**e**

Deletions
- ID-TOP1
- Other

Relative rate

1.4
1.2
1.0
0.8

None    Low    Mid    High
Germline expression

a

**Replication**

Pol ε

Sc

Origin

Pol ε

Genome-embedded
ribonucleotide

TOP1

Sc

Sc

RNA

Pol II

TOP1

Sc

Gene

Sc

TOP1

**Transcription**

b

TOP1

Sc

**Resolution of topological stress**

Rotation

**TOP1-mediated mutagenesis**

ID-TOP1
2-5 bp deletions

Genome-embedded
ribonucleotide