



35 **ABSTRACT**

36 Xeroderma pigmentosum (XP) is a genetic disorder caused by mutations in genes of the  
37 Nucleotide Excision Repair (NER) pathway (groups A-G) or in Translesion Synthesis (TLS) DNA  
38 polymerase  $\eta$  (V). XP is associated with an increased skin cancer risk, reaching, for some groups,  
39 several thousand-fold compared to the general population. Here, we analyzed 38 skin cancer  
40 genomes from five XP groups. We found that the activity of NER determines heterogeneity of the  
41 mutation rates across skin cancer genomes and that transcription-coupled NER extends beyond  
42 the gene boundaries reducing the intergenic mutation rate. Mutational profile in XP-V tumors and  
43 experiments with *POLH*-KO cell line revealed the role of polymerase  $\eta$  in the error-free bypass of  
44 (i) rare TpG and TpA DNA lesions, (ii) 3' nucleotides in pyrimidine dimers, and (iii) TpT  
45 photodimers. Our study unravels the genetic basis of skin cancer risk in XP and provides insights  
46 into the mechanisms reducing UV-induced mutagenesis in the general population.

## 47 INTRODUCTION

48 Xeroderma Pigmentosum (XP) is a group of eight rare hereditary recessive disorders  
49 caused by mutations in seven nucleotide excision repair (NER) pathway genes (groups A-G) or  
50 in the *POLH* gene coding the translesion synthesis (TLS) DNA polymerase  $\eta$  (XP-V)<sup>1</sup>. XP is  
51 characterized by up to a 10000-fold increased risk of non-melanoma skin cancers and 2000-fold  
52 increased risk of melanoma<sup>2</sup>. Moreover, epidemiological studies revealed a 34-fold increased risk  
53 of internal tumors in XP patients which was associated with characteristic mutation signature and  
54 accelerated accumulation of mutations<sup>3,4</sup>.

55 Nucleotide excision repair (NER) is the main pathway that removes bulky DNA lesions in  
56 the genome in an error-free manner<sup>5</sup>. NER can be initiated by two sub-pathways: global genome  
57 repair (GG-NER) and transcription-coupled repair (TC-NER), while the downstream mechanism  
58 of lesion removal is shared between the two and involves recruitment of the TFIIH complex and  
59 XPA, which unwind the DNA helix at the lesion site, and XPG and XPF-ERCC1, which excise the  
60 fragment containing the damaged nucleotides<sup>6</sup>. GG-NER operates genome-wide; it recognizes  
61 UV-induced bulky lesions with the XPE/DBB2 protein or helix distortions caused by those lesions  
62 with XPC protein. TC-NER is initiated by lesion-stalled RNA polymerase II and operates mainly  
63 on the transcribed strand of active genes.

64 The photoproducts which have not been removed by NER may block the progression of  
65 replicative polymerases during DNA replication. The TLS polymerase  $\eta$  is a DNA polymerase that  
66 bypasses the UV-induced photoproducts, thus preventing replication fork stalling<sup>7,8</sup>.

67 Skin cancer predisposition among XP groups is highly heterogenous, and there is an  
68 inverse relationship between the level of sunburn sensitivity and skin cancer incidence between  
69 the groups<sup>9</sup>. The most skin cancer-prone groups are XP-C and XP-E, with impaired GG-NER,  
70 and XP-V - with the deficiency of polymerase  $\eta$ <sup>9</sup>. Other XP groups with the deficiency in both GG-  
71 NER and TC-NER are also associated with considerable skin cancer risk and demonstrate a high  
72 level of sunburn sensitivity and neurological symptoms<sup>10</sup>. A current model posits that UV exposure  
73 in the context of TC-NER deficiency may cause an impairment of transcription and result in  
74 decreased cell fitness, whereas defective GG-NER results in tumor-proneness<sup>11</sup>.

75 The process of UV-induced mutagenesis depends on several major factors, including  
76 DNA lesion generation, removal by NER, and bypass by TLS polymerases. Skin cancers from  
77 XP groups differ from each other and sporadic skin cancers by the ability to repair or bypass DNA  
78 lesions, but not by the sources of DNA damage. Thus, analysis of XP tumors with defects in GG-  
79 NER alone, both TC- and GG-NER, or TLS, enables the disentanglement of the contribution of  
80 those components to mutagenesis in a natural physiological system. Moreover, an extreme skin  
81 cancer susceptibility in XP patients points to vulnerabilities in the mechanisms of protection from  
82 excessive mutation accumulation in normal skin cells.

83 In this study, we assembled a unique collection of 38 skin cancers from 5 xeroderma  
84 pigmentosum groups (XP-A, C, D, E, V). We used Whole Genome Sequencing (WGS) to study  
85 the role of defects in the major components of NER and translesion synthesis on tumor mutation

86 burden, mutation profiles, genomic landscape, and protein-damaging effects of mutagenesis in  
87 human skin cancers.

88

89

90

## 91 **RESULTS**

### 92 **Samples and clinical characteristics**

93 We collected and sequenced genomes of 33 skin cancers from 21 patients representing  
94 5 out of 8 xeroderma pigmentosum (XP) groups (3 XP-A, 4 XP-C, 2 XP-D, 10 XP-E, and 14 XP-  
95 V tumors; **Supplementary Table 1**). Causative homozygous (n=12) or compound heterozygous  
96 (n=8) germline variants were identified in 20 patients, 13 of which had known causative germline  
97 mutations (**Supplementary Table 1**), while the 7 others – had novel germline mutations  
98 compatible with the diagnosis. The mean tumor purity and sequence coverage were 41% and  
99 40X (30X for normal tissue), respectively. Additionally, we sequenced genomes of 6 sporadic  
100 cutaneous squamous cell carcinoma samples (SCC). This newly generated data was combined  
101 with WGS data from four previously published XP-C tumors<sup>12,13</sup>, one XP-D<sup>14</sup>, as well as 25  
102 sporadic cutaneous Squamous Cell Carcinomas (SCC)<sup>12,15,16</sup>, 8 Basal Cell Carcinomas<sup>15</sup> (BCC)  
103 and 113 Melanomas<sup>17</sup> (MEL) from individuals not affected with XP. The resulting cohort of XP  
104 tumors included 17 BCCs, 15 SCCs, five melanomas, and one angiosarcoma. The mean age at  
105 biopsy in XP-cohort was 33 years old (ranging from 25 years old in the XP-C group to 48 years  
106 old in the XP-V group) while in sporadic skin cancer group it was 65 years old (**Table 1**,  
107 **Supplementary Table 1**).

108

### 109 **XP groups demonstrate different mutation burden and mutation profiles**

110 We assessed the Tumor Mutation Burden (TMB) and mutation profiles of skin cancer  
111 genomes from 5 sequenced XP groups and compared them with the three types of sporadic skin  
112 cancer including BCC, SCC and MEL (**Fig. 1a**). The mean TMB of single base substitutions (SBS)  
113 was significantly higher in 3 XP groups: XP-E (350 mut/Mb,  $P = 0.0241$ ), XP-V (248 mut/Mb,  
114  $P=0.0014$ ) and XP-C skin cancers (162 mut/Mb,  $P=0.0220$ ), than the dataset weighted average  
115 (130 mut/Mb, global  $P < 2.2e-16$ ; Kruskal–Wallis H test; **Fig. 1a**). We also observed a striking  
116 difference in the TMB and the proportion of CC>TT tandem base substitutions (DBS)  
117 characteristic of UV-induced mutagenesis between the different XP groups and sporadic cancers  
118 (**Fig. 1a**). The highest proportion of CC>TT DBS from UV-induced SBS in pyrimidine dimers (C>T  
119 in YpC or CpY contexts; Y denotes a pyrimidine) was observed in XP-C and XP-D tumors (0.2  
120 and 0.17, respectively), which was 6 times higher than in sporadic skin cancers (0.03,  $P = 4.7e-$   
121 08, Mann–Whitney U test, two-sided).

122

123 The mutation profiles of skin cancers in all XP groups were dominated by C>T  
124 substitutions at pyrimidine dimers, as also found in sporadic skin cancers. However, some XP  
125 groups demonstrated marked differences for C>T mutations in specific contexts, such as  
enrichment at TCA in XP-E, TCW in XP-C, or NCY in XP-D (where W denotes A or T; N: A, C, G,

126 or T; Y: C or T). Moreover, in XP-V skin cancers, we report abundant mutations, namely C:G>A:T,  
127 T:A>A:T, and T:A>C:G, which were not previously seen to a significant degree in skin cancer  
128 (**Fig. 1b, Supplementary Fig. 1**). XP tumors formed clusters by XP group, which were non-  
129 overlapping with the cluster of sporadic skin cancers based on SBS mutation profiles and  
130 multidimensional scaling analysis (MDS; **Fig. 1c,d,e**). XP-V, XP-C, and XP-A clusters were  
131 located distantly, while the XP-E / XP-D cluster was closer to the cluster of sporadic skin cancers.

132 Among 78 COSMIC mutation signatures<sup>18</sup> (v3.2) extracted from the pan-cancer dataset,  
133 four mutation signatures (SBS7a/b/c/d) are associated with UV irradiation. We investigated  
134 whether these signatures could explain the observed mutation profiles in XP skin cancer with an  
135 accuracy comparable to sporadic skin cancers. For that, we compared observed and  
136 reconstructed mutation profiles for each sample in our cohort. The mean Cosine dissimilarity  
137 distance was small for sporadic skin cancers (0.004) but increased drastically for all the XP groups  
138 (0.16) and particularly for XP-C (0.237), XP-A (0.1957), and XP-V (0.222, **Fig. 1f**).

139

## 140 **Nucleotide excision repair efficiency determines mutation load distribution along the** 141 **genome**

142 Strong heterogeneity in the mutation rate across the genome is an important fundamental  
143 feature of mutagenesis, which has several clinical implications, for example, the discovery of  
144 cancer driver genes. We investigated the distribution of typical UV mutations (YC>YT or CY>TY)  
145 in XP and sporadic skin cancers in relation to replication timing (RT), active and inactive  
146 topologically-associated domains (TAD), and markers of chromatin states. These analyses  
147 revealed a major role for NER in shaping the heterogeneity of local rates of UV-induced mutations  
148 across the genome. A maximal 5.2-fold difference was observed between the earliest and the  
149 latest replicating bins in sporadic skin cancers (average for BCC, cSCC, and MEL) with a  
150 monotonal decrease of mutation load from late to early replicating genomic regions (**Fig. 2a**). This  
151 effect was much weaker in GG-NER deficient XP-C genomes (2.4-fold) and almost disappeared  
152 in GG-NER and TC-NER deficient XP-A (1.5-fold) and XP-D (0.99-fold) genomes (**Fig. 2a**).  
153 Interestingly, the distribution of UV-induced SBS by RT in XP-E and XP-V genomes was not very  
154 different from sporadic skin cancer genomes, 4.6-fold and 5.4-fold, respectively.

155 It has been recently shown that TAD boundaries between active and inactive chromatin  
156 domains strongly delineate the transition between regions with low and high mutation load in  
157 different human cancers<sup>19</sup>. Indeed, in our cohort, we found a 2.2-fold difference in mutation load  
158 between active and inactive TADs in sporadic cancers, but it was noticeably decreased in XP-C  
159 (1.4-fold) cancers and was virtually absent in XP-A (1.05-fold) and XP-D (1.09-fold; **Fig. 2b**).  
160 Similarly, the mutation load in XP-A and XP-D tumors was independent of chromatin states, the  
161 XP-C group demonstrated a mild dependence, while the XP-E and the XP-V groups were not  
162 different from sporadic cancers (**Fig. 2c**).

163 CPD and 6-4PP DNA lesions occur on pyrimidine bases, which enabled us to identify the  
164 strand on which the lesion underlying a UV-induced mutation occurred. In order to separately  
165 investigate the genomic targets of GG-NER and TC-NER, we split the genome into intergenic,

166 transcribed, and untranscribed strands of genic regions. A strong decrease mutation rate in the  
167 early RT regions in groups proficient in GG-NER (sporadic cancers and XP-V), and surprisingly  
168 in GG-NER deficient XP-E, was observed in intergenic regions and untranscribed strands of  
169 genes. Whereas XP-A, XP-D, and XP-C groups, which lack GG-NER, had flat slopes compatible  
170 with the lack of repair in the open chromatin of early RT regions (**Fig. 2d, Supplementary Fig.**  
171 **2**).

172

### 173 **Transcriptional bias is different between the XP groups**

174 TC-NER removes UV-induced bulky DNA lesions on the transcribed strand of expressed  
175 genes more efficiently than GG-NER on the untranscribed strand resulting in a decrease of  
176 mutations on the transcribed versus untranscribed strand, a phenomenon called transcriptional  
177 bias (TRB)<sup>20</sup>. In skin tumors with proficient NER, the TRB ranged between 1.3 and 1.6-fold for  
178 sporadic cancers and was 1.7 in XP-V. In the GG-NER-deficient TC-NER proficient groups, TRB  
179 was particularly high, ranging between 1.77-fold (XP-E) and 2.42-fold (XP-C), which is compatible  
180 with defects in the repair of the untranscribed strand. In contrast, in XP-A and XP-D groups with  
181 defects of both TC-NER and GG-NER TRB was minimal or absent: 1.17-fold and 0.97-fold,  
182 respectively (**Fig. 3a**).

183

### 184 **TC-NER removes DNA lesions downstream of genes and influences intergenic mutation** 185 **load**

186 Since early RT regions are particularly gene-rich<sup>21</sup>, we hypothesized that in GG-NER  
187 deficient XP groups, decreased mutation load in early RT intergenic regions might be associated  
188 with the TC-NER activity beyond gene boundaries. Indeed, in GG-NER deficient XP-C tumors,  
189 we revealed a significant TRB up to 50kb downstream of the furthest annotated transcriptional  
190 end sites (TES) of genes with decreased mutation frequency on the transcribed strand of nearby  
191 genes (**Fig. 3b, Supplementary Fig. 3**). The same effect was observed in XP-E and even in NER  
192 proficient skin cancers although with a lower magnitude. As expected, we did not observe TRB  
193 downstream of genes in XP-A and XP-D samples being deficient for both TC-NER and GG-NER  
194 (**Fig. 3b, Supplementary Fig. 3**). To validate TC-NER activity downstream of gene TES, we used  
195 previously published XR-seq data from XPC-deficient cell lines. It is expected that in XPC-  
196 deficient cells, XR-seq data, representing the sequencing of lesion-containing DNA fragments  
197 excised by NER<sup>22</sup>, is produced exclusively by TC-NER. An XR-seq signal was observed up to  
198 50kb downstream of TES on a transcribed strand of a nearby gene, mirroring mutation asymmetry  
199 in the same regions in XP-C tumors (**Fig. 3c**) and was well correlated with the transcriptional  
200 intensity of nascent RNA, which was retrieved from an independent study<sup>23</sup> (**Fig. 3d**). This  
201 suggests that, in some cases, the RNA polymerase might continue transcription after TES and  
202 recruit TC-NER at lesion sites. We identified XR-seq signal in 21% of the cumulative length of  
203 intergenic regions and 14% - of untranscribed strands of genes in XPC-deficient cell line<sup>22</sup>,  
204 suggesting ubiquitous extended TC-NER activity. Analysis of transcriptional bias and relative  
205 mutation rate in intergenic regions of XP-C tumors (**Fig. 3e, f**) revealed strong dependence on

206 the intensity of XR-seq outside the annotated genic regions. This extended TC-NER activity  
207 outside of the transcribed strand of genes is especially strong in early replicating regions with a  
208 high density of active genes (**Fig. 3c**). It may explain the decrease of the mutation density in  
209 intergenic regions and on the untranscribed strands of genes in early replicating genomic regions  
210 of GG-NER deficient XP-C samples (**Fig. 2d, Supplementary Fig. 2**).

211

### 212 **XP-E demonstrates reduced GG-NER activity**

213 The sensors of UV-induced DNA lesions in GG-NER, XPC, and DDB2 (XPE) are thought  
214 to work in tandem when DDB2 binds directly to a lesion and facilitates recruitment of XPC, which  
215 in turn initializes the repair process with the TFIIH complex<sup>24</sup>. We next decided to compare the  
216 features of UV-induced mutagenesis in XP-E resulting from the loss of DDB2 with XP-C and  
217 sporadic tumors.

218 MDS plot based on SBS revealed three different well delineated clusters of XP-C, XP-E,  
219 and sporadic tumors (**Fig. 4a**). At the same time, the proportion of CC>TT DBS was much  
220 increased in XP-C (0.21) versus sporadic cSCC (0.064), but significantly decreased in XP-E  
221 cSCC (0.034,  $P = 0.0003$ ; Mann–Whitney U test, two-sided), confirming qualitative differences of  
222 mutagenesis in XP-E. Unlike XP-C, the distribution of the mutational load in intergenic and  
223 untranscribed strand gene regions by RT in XP-E was very close to that of sporadic cSCC,  
224 suggesting that repair in early RT regions was functional in XP-E (**Fig. 2d, Supplementary Fig.**  
225 **2**). Similarly, the MDS based on the local mutation load in 2684 1Mb-long intervals along the  
226 genome, revealed no difference between XP-E and sporadic samples, while XP-C formed a  
227 separate cluster irrespective the tumor type (**Fig. 4b**).

228 The XP-E group demonstrated a strong TRB (1.77-fold), which was intermediate between  
229 sporadic cSCC (1.33) and the XP-C group (2.47) (**Fig. 3a, Fig. 4c,d**). Given that TC-NER is  
230 functional in XP-E, XP-C, and sporadic samples, and assuming that GG-NER is fully abrogated  
231 in XP-C, we can estimate the relative efficiency of GG-NER in XP-E tumors. Providing all else is  
232 equal, GG-NER is 64% less efficient in XP-E than in sporadic cancers.

233 To provide a more detailed view of the mutation difference between XP-E, XP-C, and  
234 sporadic tumors, we compared the association of mutation load in each group with the core  
235 epigenetic marks from primary keratinocyte cell line<sup>25</sup> using only cSCC samples (**Fig. 4e**). Unlike  
236 XP-C, XP-E tumors did not show strong and significant differences from sporadic cSCC in the  
237 dependence of mutagenesis on the majority of epigenetic covariates except for the histone  
238 modification marks H3K36me3, H3K27ac and H3K9me3 on the transcribed strand of gene  
239 regions (**Fig. 4e**).

240 Taking these observations together, we can speculate that in XP-E tumors, there is a  
241 residual activity of GG-NER associated with the ability of XPC to find a fraction of DNA lesions  
242 and initiate NER. This correlates with the clinical observation that XP-E patients develop less and  
243 later skin tumors than XP-C patients.

244

245

## 246 **Polymerase $\eta$ deficiency causes a unique mutation profile in skin cancers**

247 The analysis of XP-V skin cancers revealed that an average of 27% (15-42%) of SBS  
248 were represented by C:G>A:T mutations with a highly specific 3-nt context (NCA) and a strong  
249 and homogeneous TRB (**Fig. 5a, Fig. 1b, Supplementary Fig. 1**). Similar mutation contexts and  
250 a TRB was observed for a part of T:A>A:T mutations, which represented 8.7% of SBS. In sporadic  
251 skin cancers, C:G>A:T and T:A>A:T mutations represented only 2.5% and 4.6%, respectively,  
252 and had different broad 3-nt contexts without a strong TRB (**Fig. 1b, Supplementary Fig. 1**).  
253 Enrichment of these types of mutations in XP-V suggests that they might originate from lesions  
254 that are bypassed by polymerase  $\eta$  in an error-free manner in sporadic skin cancer, but XP-V  
255 cells have to use an alternative polymerase(s) to bypass these lesions.

256 The direction of TRB for these types of mutations indicates a decrease in mutations from  
257 lesions involving purines on the transcribed strand (**Fig. 5a**). Furthermore, comparison of  
258 C:G>A:T mutation frequencies on the transcribed and untranscribed strands with the proximal 5'  
259 intergenic regions confirmed that TRB is indeed associated with a decrease of C:G>A:T mutations  
260 on the transcribed strand (**Fig. 5b**). This suggests that mutations occur due to lesions involving  
261 purines, which are NER substrates and are effectively repaired by TC-NER on the transcribed  
262 strand (**Fig. 5b**). Interestingly, C:G>A:T mutations had stronger TRB than YC>YI or CY>IY UV-  
263 induced mutations in all bins of genes grouped by the expression level (**Fig. 5c**). This observation  
264 might indicate that those lesions produce a smaller helix distortion and are less visible to GG-  
265 NER than UV-induced pyrimidine lesions.

266 C:G>A:T and T:A>A:T mutations occurred in a very specific dinucleotide context, where  
267 a purine is always preceded by a thymine base (TA/G > TT), suggesting that causative DNA  
268 lesions might be thymine-purine dimers (**Fig. 5d**). The number of mutations in a TG context was  
269 strongly correlated with the number of mutations in a TA context ( $R=0.98$ ; **Supplementary Fig.**  
270 **4**) in our XP-V skin cancer cohort suggesting coordinated mutation processes.

271 To assess the possibility that these lesions were generated directly or indirectly due to  
272 UV-irradiation, we measured a Pearson correlation of TG > TT or TA > TT mutations with typical  
273 UV-induced (YC>YI or CY>IY) mutations and observed strong correlations in both cases,  
274  $R=0.78$  ( $P = 0.001$ ) and  $R=0.99$  ( $P = 1e-10$ ), respectively (**Supplementary Fig. 4**).

275 To further understand the nature of TG > TT and TA > TT mutations we established a  
276 *POLH* knock out of the RPE-1 *TP53*-KO cell line and sequenced whole genomes of the *POLH* wt  
277 and *POLH*-KO clones both without treatment and with treatment with KbrO<sub>3</sub> (to induce reactive  
278 oxygen species), UV-A and UV-C (**Supplementary Fig. 5**). There were no major differences in  
279 the number of mutations and mutational profiles between *POLH*-wt and *POLH*-KO for untreated  
280 cells and KbrO<sub>3</sub>-treated (**Fig. 5e,f,g**). UV-A and UV-C exposures greatly increased number of  
281 SBS in the *POLH*-KO cells (6.4 and 11.7-folds respectively) and dramatically changed the  
282 mutational profiles in comparison with *POLH*-wt clones (**Fig. 5e,f,g**). UV-A-treated *POLH*-KO  
283 clone had 15% of TG > TT mutations and 12% of the TA > TT mutations with specific XP-V context  
284 and strong transcriptional bias while in the UV-C-treated clone these percentages were 10% and  
285 4% respectively (**Fig. 5f,g**). UV-treated *POLH*-KO cells demonstrated a distinct pattern of TG>VT



286 DBS substitutions (V – A, C or G). Interestingly, a similar DBS pattern was also visible in XP-V  
287 tumors (**Supplementary Fig. 6**).

288 Another feature of the XP-V skin cancer profile was the presence of 15% (range 11% -  
289 23%) of mutations originating from TT pyrimidine dimers. Such mutations are very rare in sporadic  
290 cancer (4.8%) because TT pyrimidine dimers are bypassed by polymerase  $\eta$  in a relatively error-  
291 free manner. Two predominant types of mutations at TT were  $\text{TT}\underline{\text{T}}>\text{TA}$  and  $\text{TT}\underline{\text{T}}>\text{TC}$ , and they, as  
292 expected for mutations from pyrimidine lesions, demonstrated strong TRB and were correlated  
293 with the typical UV-induced  $\text{YC}\underline{\text{T}}>\text{YT}$  or  $\text{CY}\underline{\text{T}}>\text{TY}$  mutations (**Fig. 5a, Supplementary Fig. 4**).  
294 Interestingly, UV-A treated *POLH*-KO cells harbored 52% of T>A and T>C mutations, while in  
295 case of treatment with UV-C, it was only 22% (**Fig. 5g**).

296

### 297 **In the absence of polymerase $\eta$ , error-prone bypass of 3' nucleotides in pyrimidine dimers** 298 **shapes the mutation profile of XP-V tumors**

299 The 3-nt context of C>T substitutions in XP-V skin cancers differed from sporadic skin  
300 cancers and other XP groups (**Fig. 1b,d**). Previously it was shown that in the absence of  
301 polymerase  $\eta$ , the bypass of CPD photoproducts can be performed in two steps by two TLS  
302 polymerases, one of which inserts a first nucleotide opposite to a 3' nucleotide of the lesion  
303 ("inserter"), and then is replaced by another TLS polymerase, which performs the extension  
304 opposite to the 5' nucleotide of the lesion ("extender")<sup>26</sup>. We hypothesized that loss of polymerase  
305  $\eta$  in skin cancer might change the probabilities of mutations at 3' versus 5' nucleotides in  
306 pyrimidine dimers and thereafter contribute to the observed differences of the mutation profiles  
307 for C>T SBS in XP-V versus sporadic skin cancer.

308 To test this hypothesis, we first estimated the relative number of mutations arising at 3'  
309 and 5' cytosines in the tetranucleotide ACCA, where we could unambiguously allocate a  
310 pyrimidine dimer (**Fig. 6a**). In sporadic skin cancers, the probabilities of mutations at 3' and 5'  
311 cytosines were similar, with only a slight increase of mutagenesis from the 3'C (55%), while in  
312 XP-V skin cancers 97% of the mutations were from the 3'C (**Fig. 6b**). This bias towards 3'  
313 pyrimidine mutations was also much stronger in XP-V versus other groups of skin cancer for the  
314 CT, TC, and TT pyrimidine dimers. For example,  $\text{AT}\underline{\text{C}}\text{A} > \text{AT}\underline{\text{T}}\text{A}$  mutations were 9.17-fold more  
315 frequent than  $\text{A}\underline{\text{C}}\text{TA} > \text{A}\underline{\text{T}}\text{TA}$  mutations in XP-V than in the other groups (normalized to the  
316 corresponding 4-nt frequencies in the human genome). A similar effect was observed for T>A and  
317 T>C mutations in ATTA context (**Fig. 6c**).

318 These results demonstrate that mutations at pyrimidine dimers in XP-V occur  
319 predominantly at the 3' nucleotide, which might be associated with the error-prone activity of the  
320 inserter polymerase which replaces polymerase  $\eta$ , and modulate the mutational profile of C>T  
321 substitutions. *POLH*-KO cells treated with UV-C conversely demonstrated a very strong bias in  
322 CC pyrimidine dimers towards mutations at 3'C (99%) (**Fig. 6d**).

323

324

325

## 326 **Mutation properties of XP groups modulate protein-damaging effects of mutagenesis.**

327 High mutation rates in cells increase cancer risk and intensify tumor evolution, while the  
328 topography of mutagenesis and mutation signatures can impact the probability of damaging or  
329 driver mutations.<sup>27</sup> In our dataset of skin cancers, the number of oncogenic mutations in the  
330 cancer genome was strongly correlated with the total mutation burden (**Fig. 7a**).

331 Active DNA repair in open chromatin regions decreases the accumulation of mutations  
332 in the early replicating gene-rich regions of cancer genomes (**Fig. 2a**). We estimated a fraction of  
333 mutations per genome falling in the exonic regions across the studied skin cancer groups and  
334 found in XP-A and XP-D tumors a significant enrichment of exonic mutations in comparison with  
335 the other groups (**Fig. 7b**). The effect was caused by the redistribution of mutations from late to  
336 early RT regions of a genome (**Fig. 2a**).

337 C>T transitions, which are the most prevalent UV mutations, have relatively low protein-  
338 damaging effect in the human genome and their damaging/silent mutation ratio is 1.8, while other  
339 types of mutations, such as C:G>A:T transversions or CC>TT DBS are more damaging with a  
340 damaging/silent mutation ratio of 3.4 and 29.5, respectively (**Fig. 7c**). To better understand how  
341 the NER deficiency modulates the protein-damaging effect of UV irradiation we grouped protein-  
342 damaging mutations into 5 categories: C>T mutations on the transcribed and untranscribed  
343 strand, CC>TT double base substitutions on the transcribed and untranscribed strands, and other  
344 SBSs (**Fig. 7d**). The largest fraction of protein-damaging mutations was accounted for by C>T  
345 substitutions in all cancer groups except XP-V where other mutation classes play a more  
346 important role.

347 Contribution of damaging C>T mutations from transcribed and untranscribed strands of  
348 genes (measured as untranscribed/transcribed ratio) differed between groups. It was balanced  
349 between strands in sporadic skin cancers (1.02-fold); at the same time the majority of damaging  
350 mutations in GG-NER deficient XP-E and XP-C groups were attributed to the untranscribed strand  
351 (1.36 and 1.82-fold, respectively), while in GG- and TC- NER deficient XP-D and XP-A groups -  
352 to the transcribed strand (0.77-fold and 0.65-fold, respectively, **Fig. 7d**). These results can be  
353 explained by the fact that UV-induced C>T SBS, which originate from the lesions on the  
354 transcribed strand, are 1.88-fold more protein-damaging as compared to the untranscribed strand  
355 of genes; thereafter, active lesion removal by TC-NER from the transcribed strand of genes  
356 results not only in reduction of a total number of mutations from UV lesions, but is particularly  
357 important for the reduction of the burden of protein-damaging mutations.

358

359

360

361

362

363

364

365

## 366 DISCUSSION

367 Our results indicate that XP skin cancers deficient in GG-NER (XP-C, XP-E) or  
368 polymerase  $\eta$  (XP-V) harbor 3.6-fold more mutations than sporadic skin cancers, on average. The  
369 mutation profiles in all XP groups were dominated by C>T mutations in pyrimidine dimers;  
370 however, they differed from sporadic skin cancers and each other. These differences can be  
371 partially explained by the increased contribution to mutagenesis of early replicating GC-rich  
372 regions in XP groups with significantly impaired NER (XP-C, XP-A, XP-D). Mutational differences  
373 in XP-E might be further explained by the important role of XPE (DDB2) protein in removal of  
374 CPD lesions rather than of 6-4PP lesions, which can be recognized by XPC directly<sup>28</sup>. Thus,  
375 observed distinct mutational profiles in XP-E and XP-C might be associated with the different  
376 relative contributions of CPD and 6-4PP lesions to mutagenesis. Moreover, CC>TT double base  
377 substitutions, a characteristic feature of skin cancers, which is particularly enriched in XP-C (but  
378 depleted in XP-E), could be associated with 6-4PP photolesions. In addition, CC>TT DBS were  
379 not strongly enriched or depleted in XP-V, which might suggest that the occurrence of these  
380 mutations does not depend exclusively on polymerase  $\eta$ .

381 The current knowledge about TLS across UV lesions posits that CPD are bypassed by  
382 polymerase  $\eta$  alone, while bulkier 6-4PP require two TLS polymerases, one of which performs  
383 insertion and the other an extension<sup>29</sup>. In XP-V cancer genomes, we demonstrated a striking  
384 increase in the mutagenicity of the 3' nucleotide of the pyrimidine dimer, a phenomenon observed  
385 before using lacZ mutational reporter gene in polymerase  $\eta$  - deficient mice.<sup>30</sup> This might be  
386 explained by the model where a CPD in the absence of polymerase  $\eta$  is bypassed in two steps  
387 instead of one, by an error-prone inserter polymerase followed by an error-free extender  
388 polymerase. We propose to name the effect of differential mutagenicity of 3' and 5' nucleotides in  
389 the intrastrand crosslinked DNA dimers as "dimer translesion bias". In this study, we presented  
390 an illustrative example of XP-V, where a dimer translesion bias significantly alters the relatively  
391 conserved UV-induced mutation profiles for C>T mutations and drives the mutator phenotype of  
392 XP-V tumors and *POLH*-KO cell line.

393 It is well known that the local mutation rate in cancers and in the germline strongly  
394 correlates with epigenetic features of the genomic regions<sup>31,32</sup>. The most striking associations  
395 were observed for replication timing, chromatin accessibility, active and non-active topologically-  
396 associated domains, and some chromatin marks such as H3K9me3, H3K27me3, and  
397 H3K9me2<sup>32</sup>. In NER-deficient XP-A and XP-D tumors, we observed weak heterogeneity of the  
398 mutation frequency across the genome depending on the chromatin status. This finding  
399 demonstrates that the decreased mutagenesis in sporadic skin cancers in large open chromatin  
400 regions is driven by their accessibility to NER.

401 We observed that intergenic genomic regions and untranscribed strands of genes in GG-  
402 NER-deficient XP-C samples had a decreased mutation load in early-replicating genomic regions,  
403 which are enriched in genes with high levels of transcription. We have shown that it can partially  
404 be explained by the activity of TC-NER up to 50 Kb beyond the annotated TES and propose to  
405 name this phenomenon "extended TC-NER". In the early-replicating genomic regions, genes are

406 densely located and transcribed colinearly or in opposite orientations. Extended TC-NER can  
407 contribute substantially to the DNA repair independently of GG-NER, thus lowering mutation load  
408 in intergenic regions and on the untranscribed strands of closely located or overlapping genes.

409 Properties of C:G>A:T and T:A>A:T mutations, characteristic of XP-V skin tumors, such  
410 as the specific dinucleotide context (TA/G), strong transcription bias, and correlation with the  
411 number of C>T UV-induced mutations, enabled us to speculate that these mutations are  
412 associated with lesions that directly or indirectly induced by UV. C:G>A:T mutations are unlikely  
413 to be associated with 7,8-dihydro-8-oxoguanine (8oxoG) DNA adducts as mutation signatures of  
414 8oxoG (COSMIC signatures SBS18 and SBS36) have different 3-nt context (lack of 5' thymine  
415 specificity) and do not demonstrate transcriptional bias  
416 (<https://cancer.sanger.ac.uk/signatures/sbs/>). Furthermore, rare photolesions in a TA context  
417 have been previously reported, and their chemical structure was described as Thymidylyl-(3'-5')-  
418 Deoxyadenosine "TA photoproducts"<sup>33-35</sup>. More recently, a study dedicated to the discovery of  
419 atypical photolesions reported rare mutagenic TA photoproducts and, to a lesser extent, TG  
420 photoproducts as the most associated with UV-irradiation after CPD and 6-4PP photolesions<sup>36</sup>.  
421 The high mutagenesis in TA/G contexts in XP-V tumors uncovers a critical and non-redundant  
422 function of polymerase  $\eta$  in the error-free bypass of highly mutagenic but poorly studied DNA  
423 lesions, probably induced by UV. The WGS analysis of the RPE-1 *POLH*-KO clones confirmed  
424 that TG>TT and TA>TT mutations are greatly increased after UV-A and UV-C exposure, and  
425 have shown that they are most prevalent after UV-A. However, after KbrO<sub>3</sub> treatment which  
426 induces reactive oxygen species, in *POLH*-KO cells SBS mutational pattern did not change as  
427 compared to the wild type. Presence of TG>VA DBS in XP-V tumors and in *POLH*-KO cells after  
428 UV-exposures further suggests that their origin might depend on thymine-guanine dimers.

429 Another important peculiarity of the XP-V mutation profile is a high fraction of mutations  
430 (on average 15%) originating from a TT dinucleotide (TT>TA/C). It is known that the majority of  
431 CPD lesions occur in a TT context<sup>37</sup> and polymerase  $\eta$  is the main polymerase to bypass them in  
432 an error-free manner. In the absence of polymerase  $\eta$ , other TLS polymerases perform bypass  
433 of TT lesions introducing more errors. Even though TT CPD represents near 50% of UV-induced  
434 lesions, the proportion of TT>TA/C mutations is only 15% in XP-V. This means that even in the  
435 absence of polymerase  $\eta$ , TT CPD is not a highly mutagenic lesion, probably because other  
436 replacing TLS polymerases insert predominantly adenines opposite to the lesion following the "A  
437 rule"<sup>38,39</sup>.

438 In XP skin cancers, UV irradiation results in different mutation profiles and topography of  
439 mutagenesis, which are associated with a variation in the probability of protein-damaging and  
440 oncogenic mutations. We revealed three main factors contributing to heterogeneity in the  
441 proportion of protein-damaging mutations in XP skin cancers, which were associated with the  
442 differences in mutation profiles (fractions of transversions and DBS), the activity of TC-NER, and  
443 mutation distribution between open and closed chromatin.

444 The observed differences in mutation burden and mutation profiles might also reflect the  
445 differences in clinical manifestations between XP groups. XP-A and XP-D patients show severe

446 sunburn reactions and are diagnosed early. Thereafter, those patients are rarely exposed to UV  
447 and have rather few tumors. This might partly explain rather low TMB in XP-A and XP-D skin  
448 tumors. In tumor-prone XP-C, XP-E, and XP-V groups, the amount and mode of exposure to UV  
449 again may be different. XP-C patients do not experience sunburns but develop other skin  
450 symptoms resulting in an early diagnosis and sun protection. XP-E and XP-V patients, on the  
451 contrary, do not have any symptoms till their 20s or 30s, by which time they may have had a lot  
452 of sun exposure, and in subsequent decades, they develop many skin tumors. This might be in  
453 line with the observation that mutational profiles in XP-E and XP-V in some features resemble  
454 sporadic cancers, while XP-C is the most different.

455 Overall, our analysis of rare skin cancers with deficient NER or translesion DNA synthesis  
456 has revealed how the absence of different NER components modulates mutation burden, profiles,  
457 and topography of mutagenesis after UV irradiation. We have attempted to provide mechanistic  
458 explanations for the mutation consequences of DDB2 (XPE) loss in the XP-E group and the  
459 polymerase  $\eta$  deficiency in the XP-V group for UV mutagenesis in skin cancer. Further mutation  
460 studies on experimental cell lines from XP patients can extend our knowledge of the role of major  
461 and rare photoproducts in skin cancer pathogenesis and biological mechanisms supporting  
462 genome stability.

463

#### 464 **ACKNOWLEDGMENTS**

465 S.I.N. was supported by grant Foundation ARC 2017, Foundation Gustave Roussy, and the The  
466 French National Cancer Institute - RPT21145LLA. P.L.K and S.I.N. were supported by grant from  
467 the Foundation ARC-ARCPGA12019120001055\_1578 (P.L.K. and S.N.). This work was also  
468 supported by Prism – National Precision Medicine Center in Oncology funded by the France 2030  
469 program and the French National Research Agency (ANR) under grant number ANR-18-IBHU-  
470 0002. The authors are very thankful to Xiaole Xu (BGI) for the management of sequencing.

471

#### 472 **AUTHOR CONTRIBUTION**

473 S.I.N. and A.A.Y. designed the study. A.A.Y. performed the data analysis and prepared  
474 figures. A.A.Y. and S.I.N. drafted the manuscript. A.S., A.L., and C.F.M.M commented on the  
475 manuscript. F.R. handled biopsies, performed QC of the samples and DNA extraction. F.R.  
476 performed cell line experiments. P.L. participated in the DNA extraction and sample  
477 handling. K.G. participated in the data analysis. I.P. and L.P. performed data preprocessing.  
478 T.B.P., C.F.M.M., H.F., A.L., C.N., P.L.K, and A.S. collected the samples.

479

#### 480 **COMPETING INTERESTS**

481 The authors declare no competing interests.

482

#### 483 **DATA AVAILABILITY**

484 Experimental data generated in this study have been deposited in the European Genome-  
485 phenome Archive (EGA) under accession XXX.

486

487

## 488 REFERENCES

- 489 1. Lehmann, A. R., McGibbon, D. & Stefanini, M. Xeroderma pigmentosum. *Orphanet J Rare*  
490 *Dis* **6**, 1–6 (2011).
- 491 2. Bradford, P. T. *et al.* Cancer and neurologic degeneration in xeroderma pigmentosum:  
492 Long term follow-up characterises the role of DNA repair. *J Med Genet* **48**, 168–176  
493 (2011).
- 494 3. Yurchenko, A. A. *et al.* XPC deficiency increases risk of hematologic malignancies through  
495 mutator phenotype and characteristic mutational signature. *Nat Commun* (2020)  
496 doi:10.1038/s41467-020-19633-9.
- 497 4. Nikolaev, S., Yurchenko, A. A. & Sarasin, A. Increased risk of internal tumors in DNA  
498 repair-deficient xeroderma pigmentosum patients: analysis of four international cohorts.  
499 *Orphanet J Rare Dis* **17**, (2022).
- 500 5. Spivak, G. Nucleotide excision repair in humans. *DNA Repair (Amst)* **36**, 13–18 (2015).
- 501 6. Marteijn, J. A., Lans, H., Vermeulen, W. & Hoeijmakers, J. H. J. Understanding nucleotide  
502 excision repair and its roles in cancer and ageing. *Nat Rev Mol Cell Biol* **15**, 465–481  
503 (2014).
- 504 7. Yang, W. & Gao, Y. Translesion and Repair DNA Polymerases: Diverse Structure and  
505 Mechanism. *Annual Review of Biochemistry* vol. 87 Preprint at  
506 <https://doi.org/10.1146/annurev-biochem-062917-012405> (2018).
- 507 8. Barnes, R. P., Tsao, W., Moldovan, G. & Eckert, K. A. DNA Polymerase Eta Prevents Tumor  
508 Cell Cycle Arrest And Cell Death During Recovery from Replication Stress. (2018)  
509 doi:10.1158/0008-5472.CAN-17-3931.
- 510 9. Sethi, M. *et al.* Patients with xeroderma pigmentosum complementation groups C, e and  
511 v do not have abnormal sunburn reactions. *British Journal of Dermatology* **169**, 1279–  
512 1287 (2013).
- 513 10. Fassihi, H. *et al.* Deep phenotyping of 89 xeroderma pigmentosum patients reveals  
514 unexpected heterogeneity dependent on the precise molecular defect. (2016)  
515 doi:10.1073/pnas.1519444113.
- 516 11. Reid-Bayliss, K. S., Arron, S. T., Loeb, L. A., Bezrookov, V. & Cleaver, J. E. Why Cockayne  
517 syndrome patients do not get cancer despite their DNA repair deficiency. *Proc Natl Acad*  
518 *Sci U S A* **113**, (2016).
- 519 12. Zheng, C. L. *et al.* Transcription Restores DNA Repair to Heterochromatin, Determining  
520 Regional Mutation Rates in Cancer Genomes. *Cell Rep* **9**, 1228–1234 (2014).
- 521 13. Momen, S. *et al.* Dramatic response of metastatic cutaneous angiosarcoma to an  
522 immune checkpoint inhibitor in a patient with xeroderma pigmentosum: Whole-genome  
523 sequencing AIDS treatment decision in end-stage disease. *Cold Spring Harb Mol Case*  
524 *Stud* **5**, (2019).
- 525 14. Cho, R. J. *et al.* APOBEC mutation drives early-onset squamous cell carcinomas in  
526 recessive dystrophic epidermolysis bullosa. *Sci Transl Med* **10**, eaas9668 (2018).
- 527 15. Priestley, P. *et al.* Pan-cancer whole-genome analyses of metastatic solid tumours.  
528 *Nature* **575**, (2019).
- 529 16. Mueller, S. A. *et al.* Mutational Patterns in Metastatic Cutaneous Squamous Cell  
530 Carcinoma. *Journal of Investigative Dermatology* (2019) doi:10.1016/j.jid.2019.01.008.
- 531 17. Hayward, N. K. *et al.* Whole-genome landscapes of major melanoma subtypes. *Nature*  
532 **545**, 175–180 (2017).

- 533 18. Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer. *Nature*  
534 (2020) doi:10.1038/s41586-020-1943-3.
- 535 19. Akdemir, K. C. *et al.* Somatic mutation distributions in cancer genomes vary with three-  
536 dimensional chromatin structure. *Nat Genet* **52**, (2020).
- 537 20. Haradhvala, N. J. *et al.* Mutational Strand Asymmetries in Cancer Genomes Reveal  
538 Mechanisms of DNA Damage and Repair. *Cell* **164**, 538–549 (2016).
- 539 21. Woodfine, K. *et al.* Replication timing of the human genome. *Hum Mol Genet* **13**, (2004).
- 540 22. Hu, J., Adar, S., Selby, C. P., Lieb, J. D. & Sancar, A. Genome-wide analysis of human  
541 global and transcription-coupled excision repair of UV damage at single-nucleotide  
542 resolution. *Genes Dev* **29**, 948–960 (2015).
- 543 23. Barbieri, E. *et al.* Rapid and Scalable Profiling of Nascent RNA with fastGRO. *Cell Rep* **33**,  
544 (2020).
- 545 24. Scrima, A. *et al.* Structural Basis of UV DNA-Damage Recognition by the DDB1-DDB2  
546 Complex. *Cell* **135**, (2008).
- 547 25. Kundaje, A. *et al.* Roadmap Epigenomics Consortium: Integrative analysis of 111  
548 reference human epigenomes. *Nature* (2015) doi:10.1038/nature14248.
- 549 26. Livneh, Z., Ziv, O. & Shachar, S. Multiple two-polymerase mechanisms in mammalian  
550 translesion DNA synthesis. *Cell Cycle* vol. 9 Preprint at  
551 <https://doi.org/10.4161/cc.9.4.10727> (2010).
- 552 27. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues. *Cell*  
553 **171**, 1029-1041.e21 (2017).
- 554 28. Oh, K. S. *et al.* Nucleotide excision repair proteins rapidly accumulate but fail to persist in  
555 human XP-E (DDB2 mutant) cells. *Photochem Photobiol* **87**, (2011).
- 556 29. Quinet, A. *et al.* Translesion synthesis mechanisms depend on the nature of DNA damage  
557 in UV-irradiated human cells. *Nucleic Acids Res* **44**, (2016).
- 558 30. Ikehata, H., Chang, Y., Yokoi, M., Yamamoto, M. & Hanaoka, F. Remarkable induction of  
559 UV-signature mutations at the 3'-cytosine of dipyrimidine sites except at 5'-TCG-3' in the  
560 UVB-exposed skin epidermis of xeroderma pigmentosum variant model mice. *DNA*  
561 *Repair (Amst)* **22**, (2014).
- 562 31. Supek, F. & Lehner, B. Scales and mechanisms of somatic mutation rate variation across  
563 the human genome. *DNA Repair (Amst)* 102647 (2019)  
564 doi:10.1016/j.dnarep.2019.102647.
- 565 32. Schuster-Böckler, B. & Lehner, B. Chromatin organization is a major influence on regional  
566 mutation rates in human cancer cells. *Nature* (2012) doi:10.1038/nature11273.
- 567 33. Bose, S. N., Davies, R. J. H., Sethi, S. K. & McCloskey, J. A. Formation of an adenine-  
568 thymine photoadduct in the deoxydinucleoside monophosphate d(TpA) and in DNA.  
569 *Science (1979)* **220**, (1983).
- 570 34. Zhao, X. & Taylor, J. S. Mutation spectra of TA, the major photoproduct of thymidylyl-(3'-  
571 5')-deoxyadenosine, in *Escherichia coli* under SOS conditions. *Nucleic Acids Res* **24**,  
572 (1996).
- 573 35. Asgatay, S. *et al.* UV-induced TA photoproducts: Formation and hydrolysis in double-  
574 stranded DNA. *J Am Chem Soc* **132**, (2010).
- 575 36. Laughery, M. F. *et al.* Atypical UV Photoproducts Induce Non-canonical Mutation Classes  
576 Associated with Driver Mutations in Melanoma. *Cell Rep* **33**, (2020).
- 577 37. Cadet, J., Grand, A. & Douki, T. Solar uv radiation-induced dna bipyrimidine  
578 photoproducts: Formation and mechanistic insights. *Top Curr Chem* **356**, (2015).
- 579 38. Taylor, J. S. New structural and mechanistic insight into the A-rule and the instructional  
580 and non-instructional behavior of DNA photoproducts and other lesions. *Mutation*

- 581            *Research - Fundamental and Molecular Mechanisms of Mutagenesis* vol. 510 Preprint at  
582            [https://doi.org/10.1016/S0027-5107\(02\)00252-X](https://doi.org/10.1016/S0027-5107(02)00252-X) (2002).
- 583    39.    Strauss, B. S. The 'A rule' of mutagen specificity: A consequence of DNA polymerase  
584            bypass of non-instructional lesions? *BioEssays* vol. 13 Preprint at  
585            <https://doi.org/10.1002/bies.950130206> (1991).
- 586    40.    Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
587            *ArXiv* **1303**, (2013).
- 588    41.    Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler  
589            transform. *Bioinformatics* **25**, 1754–1760 (2009).
- 590    42.    van der Auwera, G. A. *et al.* GATK Best Practices. *Current protocols in bioinformatics /*  
591            *editorial board, Andreas D. Baxevanis ... [et al.]* (2002) doi:10.1002/0471250953.
- 592    43.    Depristo, M. A. *et al.* A framework for variation discovery and genotyping using next-  
593            generation DNA sequencing data. *Nat Genet* (2011) doi:10.1038/ng.806.
- 594    44.    Ramos, A. H. *et al.* Oncotator: Cancer variant annotation tool. *Hum Mutat* (2015)  
595            doi:10.1002/humu.22771.
- 596    45.    Shen, R. & Seshan, V. E. FACETS: Allele-specific copy number and clonal heterogeneity  
597            analysis tool for high-throughput DNA sequencing. *Nucleic Acids Res* **44**, (2016).
- 598    46.    Andrews, S. FASTQC A Quality Control tool for High Throughput Sequence Data.  
599            *Babraham Institute* (2015).
- 600    47.    Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,  
601            2078–2079 (2009).
- 602    48.    Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, (2021).
- 603    49.    Ewels, P., Magnusson, M., Lundin, S. & Källér, M. MultiQC: Summarize analysis results for  
604            multiple tools and samples in a single report. *Bioinformatics* (2016)  
605            doi:10.1093/bioinformatics/btw354.
- 606    50.    Köster, J. & Rahmann, S. Snakemake-a scalable bioinformatics workflow engine.  
607            *Bioinformatics* (2012) doi:10.1093/bioinformatics/bts480.
- 608    51.    David Meyer *et al.* e1071: Misc Functions of the Department of Statistics, Probability  
609            Theory Group. *CRAN Repository* Preprint at (2021).
- 610    52.    Khanna, A. *et al.* Bam-readcount - rapid generation of basepair-resolution sequence  
611            metrics. *J Open Source Softw* **7**, (2022).
- 612    53.    Garrison, E., Kronenberg, Z. N., Dawson, E. T., Pedersen, B. S. & Prins, P. Vcflib and tools  
613            for processing the VCF variant call format. *bioRxiv* (2021).
- 614    54.    Bergstrom, E. N. *et al.* SigProfilerMatrixGenerator: a tool for visualizing and exploring  
615            patterns of small mutational events. *BMC Genomics* **20**, 1–12 (2019).
- 616    55.    Blokzijl, F., Janssen, R., van Boxtel, R. & Cuppen, E. MutationalPatterns: Comprehensive  
617            genome-wide analysis of mutational processes. *Genome Med* (2018)  
618            doi:10.1186/s13073-018-0539-0.
- 619    56.    Manders, F. *et al.* MutationalPatterns: the one stop shop for the analysis of mutational  
620            processes. *BMC Genomics* **23**, (2022).
- 621    57.    Tate, J. G. *et al.* COSMIC: The Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids*  
622            *Res* **47**, (2019).
- 623    58.    Hansen, R. S. *et al.* Sequencing newly replicated DNA reveals widespread plasticity in  
624            human replication timing. *Proc Natl Acad Sci U S A* (2010) doi:10.1073/pnas.0912402107.
- 625    59.    Neph, S. *et al.* BEDOPS: High-performance genomic feature operations. *Bioinformatics*  
626            (2012) doi:10.1093/bioinformatics/bts277.
- 627    60.    Frankish, A. *et al.* GENCODE 2021. *Nucleic Acids Res* **49**, (2021).
- 628    61.    Quinlan, A. R. BEDTools: The Swiss-Army tool for genome feature analysis. *Curr Protoc*  
629            *Bioinformatics* (2014) doi:10.1002/0471250953.bi1112s47.



630 62. Chakravarty, D. *et al.* OncoKB: A Precision Oncology Knowledge Base. *JCO Precis Oncol*  
631 (2017) doi:10.1200/po.17.00011.  
632

633

## 634 **MATERIALS AND METHODS**

635

### 636 **Studied samples**

637 The samples were collected from patients with a confirmed XP diagnosis. Informed signed  
638 consents were obtained from patients and/or their parents per the Declaration of Helsinki and the  
639 French law. This study was approved by the French Agency of Biomedicine (Paris, France), the  
640 Ethics Committee from the CPP of the University Hospital of Bordeaux (Bordeaux, France), the  
641 Institutional Review Board of Gustave Roussy (CSET: 2018-2820; Gustave Roussy, Villejuif,  
642 France), the Research Ethics Committee of Guy's and St Thomas' Foundation Trust, London  
643 (reference 12/LO/0325), and the CONEP (Brazil), Number CAAE 48347515.3.0000.5467.

644 The tumor samples were collected from patients during surgery. The tumors were stored  
645 in liquid nitrogen or allprotect tissue reagent, and 16 in FFPE. Normal control samples were  
646 represented by blood (4 patients), saliva (7 patients), fresh skin (2 patients), or FFPE (6 patients).  
647 DNA from non-FFPE tissues was extracted using AllPrep DNA/RNA/miRNA Universal Kit (Cat.  
648 No. / ID: 80224, Qiagen) according to the manufacturer's instructions. DNA from FFPE blocks  
649 was extracted after examination and dissection by a pathologist. Tumor DNA was extracted from  
650 parts of FFPE containing a high fraction of tumor cells using Maxwell® RSC DNA FFPE Kit  
651 (Catalog number: AS1450, Promega) according to the manufacturer's instructions. Non-tumoral  
652 DNA was extracted from FFPE blocks that did not contain tumor cells if available, or from parts  
653 of tumor cell-containing FFPE blocks free from tumor cells. DNA quantity and quality were  
654 assessed using the NanoDrop-ND-1000 (Nanodrop Technologies).

### 655 **Genome sequencing and variant calling**

656 The genomes were sequenced using BGISEQ-500 in BGI (Shenzhen) according to the  
657 manufacturer's protocols to the mean coverage after deduplication equal to 40X for tumor and  
658 30X for normal DNA (100 bp paired-end reads). Reads were mapped using BWA-MEM<sup>40,41</sup>  
659 (v0.7.12) software to the GRCh37 human reference genome, and then we used the standard  
660 GATK best practice pipeline<sup>42</sup> to process the samples and call somatic and germline genetic  
661 variants. PCR duplicates were removed, and the base quality score recalibrated using GATK<sup>43</sup>  
662 (v4.0.10.1), MarkDuplicates, and BaseRecalibrator tools. Somatic variants were called and  
663 filtered using GATK tools Mutect2, FilterMutectCalls, and FilterByOrientationBias and annotated  
664 with oncotator<sup>44</sup> (v1.9.9.0). SCNAs calling was done with FACETS<sup>45</sup> (v 0.5.14). Quality controls  
665 of FASTQ files and mapping were done with FASTQC<sup>46</sup> (v0.11.7), samtools<sup>47,48</sup> (v1.9), GATK  
666 HSMetrics, and MultiQC<sup>49</sup> (v1.5). All processing steps were combined in a pipeline built with  
667 snakemake<sup>50</sup> (v5.4.0).

### 668 **Filtration of somatic variants in tumor samples**

669 Only PASS-filtered somatic variants supported by at least one read from each strand and  
670 at least three reads in total with variant allele frequency higher than 0.05 and POPAF filter > 5  
671 (negative log 10 population allele frequencies of alt alleles; probability of the mutation to be a  
672 germline polymorphism) were used for the analysis. Additionally, all used VCF files were filtered  
673 based on the alignability map of the human genome from the UCSC browser  
674 (<https://genome.ucsc.edu/cgi-bin/hgFileUi?db=hg19&g=wgEncodeMapability>) with the length of  
675 K-mer equal to 75 bp (wgEncodeCrgMapabilityAlign75mer, mutations overlapped regions with  
676 score <1 were filtered out) and UCSC Browser blacklisted regions (Duke and DAC).

677 To filter out the FFPE artefacts, we employed Support Vector Machine-based (SVM)  
678 methodology with the e1071 R library.<sup>51</sup> For each sample separately, each variant in the  
679 prefiltered VCF file (the same filters as for the fresh non-FFPE samples) was annotated with  
680 additional quality information specific for the alternative allele from the BAM file using bam-  
681 readcount utility<sup>52</sup>. This additional BAM-derived information in the form of a table was merged with  
682 the quality annotations from the VCF file (VCF was parsed into a table with vcf2tsv from vcflib  
683 library<sup>53</sup>) which included CONTQ (Phred-scaled qualities that alt allele are not due to  
684 contamination), SEQQ (Phred-scaled quality that alt alleles are not sequencing errors),  
685 STRANDQ (Phred-scaled quality of strand bias artifact), TLOD (Log 10 likelihood ratio score of  
686 variant existing versus not existing). The typically UV-induced double base substitutions  
687 (CC:GG>TT:AA) were considered true positive variants, while abundant FFPE artefacts  
688 TG:CA>CA:TG were considered false positive variants during the training of SVM. To tune the  
689 SVM parameters we subset 25% of the TG:CA>CA:TG and CC:GG>TT:AA variants and run  
690 tune() command (cost=c(0.001,0.01,0.1, 1,5,10,100)). Then the best tuning parameters for the  
691 model were chosen (tune.out\$best.model) and applied to the training dataset of 50% of the  
692 TG:CA>CA:TG and CC:GG>TT:AA variants using svm() command with 10 k-fold cross  
693 validations (cross=10) and probabilistic assignment of the classification (type="C-classification",  
694 probability = TRUE, scale=T) to build the SVM classification model. Finally, the SVM classification  
695 model was applied to the whole dataset of variants to classify them as true positive in a  
696 probabilistic manner (command predict(), probability = TRUE). We extracted for the downstream  
697 analysis only the variants with a probability of being true positive > 0.95.

### 698 **Mutation spectrum, MDS, and comparison with known signatures**

699 To convert the VCF files into a catalog of mutation matrices, we used  
700 SigProfilerMatrixGenerator v.1.0 software<sup>54</sup>. Before the profiling, VCF files were split into separate  
701 files with single base substitutions and other variants to avoid splitting double base substitutions  
702 into single base substitutions by the software. To construct the multidimensional scaling plots  
703 (MDS), we computed pairwise Cosine similarity distance between all the samples using  
704 MutationalPatterns R package<sup>55,56</sup> (cos\_sim\_matrix()) and then processed the matrix of distances  
705 between the samples in the prcomp() function in R.

706 To understand whether known UV signatures can explain the mutational profiles of XP  
707 and sporadic datasets, we extracted four SBS mutational signatures previously associated with  
708 UV irradiation (SBS 7a,b,c,d) from the COSMIC database<sup>57</sup> (V3.2,

709 <https://cancer.sanger.ac.uk/signatures/sbs/>) and then reconstructed observed mutational profiles  
710 of the studied samples using these four UV-associated mutational signatures (`fit_to_signatures()`,  
711 `MutationalPatterns` R package). The Cosine dissimilarity of the observed and reconstructed  
712 mutational profiles was calculated for each sample as 1-Cosine distance. The procedure was  
713 performed separately for all 96 trinucleotide mutational contexts and only 12 mutational contexts  
714 of the UV-induced spectra (`NCY>NTY` or `YCN>YIN`).

### 715 **Replication timing, TADs, epigenetic marks, and mutational load along the genome**

716 We used Repli-Seq data from 11 cell lines<sup>58</sup> (BG02, BJ, GM0699, HeLa, HEPG2, HUVEC,  
717 IMR90, K562, MCF7, NHEK, SK-N-SH) to identify conservative replication timing regions. For  
718 each 1-kb region, we calculated weighted mean replication timing and then its standard deviation  
719 between all the cell lines and removed all the regions with a standard deviation higher than 15.  
720 For the rest of consistent regions across different cell lines, we calculated the mean values and  
721 used them during analysis. The genome was divided into five or eight bins according to the  
722 replication timing values, and mutation density was calculated for each bin, adjusting for  
723 trinucleotide contexts. Additionally, we computed the dependence of mutation density on  
724 replication timing separately for intergenic and genic regions (splitting mutations on the  
725 transcribed and untranscribed strands).

726 The genomic location of the 1MB borders between topologically-associated domains  
727 (TADs) was downloaded from the recent publication exploring mutation rate dependency on TAD  
728 structures<sup>19</sup>. The border regions were spitted into 1-kb intervals and separated into four bins (two  
729 for active and two for inactive TADs). Then the fraction of mutations per each sample fallen into  
730 each bin was calculated, adjusting for the trinucleotide composition. A similar procedure was  
731 performed for the consensus chromatin states of the genome from the same publication.

732 To calculate the slopes of the mutation load over replication timing (or other epigenetic  
733 marks) bins per sample, the logarithm of the normalized fraction of mutations in each bin was  
734 fitted into a linear model (`lm()`) with the number of each bin (1 to 8).

735 To investigate the relationships between mutation density and intensity of various  
736 epigenetic marks (DNase, H3K36me3, H3K27ac, H3K4me1, H3K27me3, H3K9me3, methylation  
737 level from whole genome bisulfite sequencing), we downloaded bigwig files of the Roadmap  
738 Epigenomics Project<sup>25</sup> and converted them to wig and then bed files (tissue E058, keratinocyte).  
739 The mean intensity of each mark was calculated for 1-kb non-overlapping windows across  
740 autosomes with BEDOPS v2.4.37 (`bedmap`) software.<sup>59</sup> The mark intensities were normalized to  
741 the 1–100 range, and we used only genomic windows with high alignability (equal to 1) along at  
742 least 90% of a window. For each window, we split mark intensities into 5 bins (`cut2()` function in  
743 R) and calculated the trinucleotide-adjusted fraction of mutations per sample per bin for each  
744 mark separately for intergenic regions, transcribed and untranscribed strands of genes.

745 To assess the mutation load distribution along the genome between groups of samples  
746 and irrespective of the epigenetic features, we split the genome into 1MB-long nonoverlapped  
747 intervals and excluded all the intervals with a mappability score less than 1 over 80% of the  
748 interval. For the resulting dataset of 2684 intervals, we calculated the mutation density of C>T

749 substitutions in each interval per sample (with at least 50000 mutations) and then normalized the  
750 mutation density. Finally, the principal component analysis was performed on the resulting matrix.

### 751 **Transcriptional bias and XR-seq**

752 Transcriptional strand bias (TRB) was quantified for each sample based on the stranded  
753 mutation matrixes generated by SigProfilerMatrixGenerator.<sup>54</sup> We computed inequality between  
754 mutations from pyrimidines (C > A/T/G; T > A/C/G) to mutations from purines (G > A/C/T;  
755 A > C/G/T) for genes located on the sense and antisense strands of DNA relative to the reference  
756 human genome.

757 To compute TRB between genes expressed with different levels, we used RPKM values  
758 of RNA-seq from Epigenetic Roadmap Project<sup>25</sup> represented by keratinocytes (E058) and only  
759 samples represented by BCC and cSCC. For each gene, mutations were separated as located  
760 on transcribed or untranscribed strands, and genes were divided into six bins by the level of  
761 expression.

762 Following the hypothesis that cytosine-containing DNA lesions caused the majority of  
763 mutations, we were also able to compute strand-specific mutation densities around transcription  
764 end sites (TESs), and transcription start sites (TSSs). Transcribed and untranscribed strands of  
765 genes and adjacent to TES/TSS intergenic regions were treated separately. TESs/TSSs of all  
766 annotated genes (GENECODE<sup>60</sup> v38) were retrieved using BEDTools v2.30.0<sup>61</sup>, and then regions  
767 located  $\pm 50$  kb of TESs/TSSs were split into 1-kb intervals. The 1-kb intervals that overlapped  
768 with other intergenic or genic intervals (represented mainly by overlapped or closely located  
769 genes) were removed for this analysis, and the rest were aggregated into 10 bins. We then  
770 separately calculated the trinucleotide context-adjusted fraction of mutations per bin per sample  
771 for transcribed and untranscribed strands.

772 XR-seq profiles for XP-C cell lines (XP4PA-SV-EB, GM15983) and nascent RNA-seq data  
773 from the HeLa cell line were downloaded from the previous works of Hu et al. 2015<sup>22</sup> and Barbieri  
774 et al. 2020<sup>23</sup>, respectively. The mean intensity of tracks was calculated for binned 1-kb intervals  
775 along the genome and  $\pm 50$  kb around the TESs.

### 776 **Dimer translesion bias**

777 To calculate the relative amount of mutations arising from 5' and 3' sides of pyrimidine  
778 dimers, we extracted mutations from C>T located in the ACCT context, mutations T>C/A located  
779 in the ATTA context, and calculated the ratio of such mutations originating from 3' C/T to 5' C/T  
780 separately for each mutation type with the corresponding 4-nucleotide context. Additionally, we  
781 calculated the ratio between the number of ATCA > ATTA and ACTA > ATTA mutations per  
782 sample, adjusting for the different fractions of ATCA and ACTA four-nucleotides.

### 783 **Protein-damaging effects of mutagenesis**

784 To assess the protein-damaging effect of different substitutions, we annotated the VCF  
785 files using oncotator<sup>44</sup> software and classified exonic mutations into protein-damaging (missense,  
786 nonsense, splice-site) and silent. For the C>T and CC>TT mutations, we separated them by  
787 strands and calculated the protein-damaging effect separately for the transcribed and

788 untranscribed strands. The number of putative oncogenic drivers per sample was calculated using  
789 the OncoKb<sup>62</sup> database (oncogenic and likely oncogenic events).

#### 790 **Cell culture**

791 RPE-1 *TP53*-KO cell line is obtained as a gift from Dr. Olivier Gavet lab. RPE-1 *POLH*-wt  
792 and RPE1 *POLH*-KO cell lines were cultured in DMEM/F-12 (gibco; life technologies, Ref:  
793 11320033) at 37 °C in a humidified atmosphere containing 5% CO<sub>2</sub>., supplemented with 10%  
794 (v/v) fetal bovine serum (FBS; NB-26-00009).

#### 795 **Generating *POLH*-KO cell line**

796 *POLH*-KO cell lines are obtained from Synthego company. sgRNA was used to generate  
797 the *POLH*-KO cell line. Homozygote knock out was verified by sanger sequencing showing 4  
798 nucleotide insertion.

#### 799 **UV exposure**

800 Cells were irradiated with 10 J/m<sup>2</sup> UV-C (200-280 nm) or UV-A (320-400 nm) for 4  
801 sequential exposures both for *POLH*-KO and *POLH*-wt cell lines. Irradiation were performed every  
802 4 days.

#### 803 **KbrO<sub>3</sub> treatment:**

804 IC50 values for KbrO<sub>3</sub> was identified as following protocol. 5,000 Cells per each well were  
805 plated and grown for 24 hours in 96-well plates. Cells were treated in serial diluted concentrations  
806 of KbrO<sub>3</sub> (500mM-10uM). Treatment was last for 96 hours. After 4 days, density of cells in each  
807 well was quantified using Methylene Blue staining. In the first step, cells (wells) were washed with  
808 PBS 1X. Then 100µl absolute methanol is added to each well and plate was incubated for 1 hour  
809 at room temperature. Then the wells were let to be dried and 100 µl methylene blue solution  
810 (concentration 1gr/L) was added to each well and followed by 1-hour incubation at room  
811 temperature. Following the staining step, wells were rinsed with water for 2 times and then let the  
812 wells to dry. Washing step followed by solubilization of stain by adding 200 µl HCL (0.1N) in each  
813 well and incubation at 60°C for 30 minutes. In the last step O.D of each well was measured at  
814 630 nm using BMG FLUOstar OPTIMA plate reader.

815 *POLH*-wt and *POLH*-KO cells were treated for 8 weeks using 300uM KbrO<sub>3</sub>. Treatment  
816 was refreshed every 48 hours.

#### 817 **Single-cell cloning and DNA extraction**

818 Single cell sorting was performed by Flow Cytometry Cell Sorting (FACS) in P96-well plate  
819 upon the completion of treatment period and reaching the sufficient number of cells. 5-6 hours  
820 after sorting the wells were monitored to confirm the presence of a single cell in the well. Three  
821 clones per condition were randomly selected to pass to P24-well plate to propagate the cells and  
822 extract DNA 18-21 days after cell sorting. Genomic DNA was extracted using the Qiamp DNA  
823 mini kit (QIAGEN) according to the manufacturer's instructions.

#### 824 **Sequencing and bioinformatic analysis of the cell line experiments**

825 The RPE-1 clonal cell populations were sequenced in BGI, Shenzhen (15x coverage,  
826 BGISEQ-500 instrument) and bioinformatically processed in the similar way with tumor samples.  
827 The nontreated samples were used as "normal" and treated as "tumor" during GATK mutect2

828 calling of somatic mutations (and vice versa). Then we removed all the nonunique mutations  
829 between the clones (module *bcftools iseq*) as well as supported by less than 3 reads in total and  
830 at least one read from each strand. Finally, only strictly clonal mutations with VAF > 0.3 were  
831 used for the analysis.

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

## 870 TABLES

871

872 **Table 1. The studied dataset of XP and sporadic skin cancers with WGS data**

873

Group	Tumors (n)	Patients (n)	Mean age at biopsy (years)	cSCC (n)	BCC (n)	Melanoma (n)
XP-E	10	3	28	7	2	1
XP-C	8*	8	25	5	1	1
XP-A	3	1	32	1	2	-
XP-D	3	3	30	2	1	-
XP-V	14	9	48	-	11	3
Sporadic SCC	31	31	73	31	-	-
Sporadic BCC	8	8	66	-	8	-
Sporadic MEL	113	113	57	-	-	113

874 \*- one XP-C patient with angiosarcoma

875

## 876 **Supplementary Table 1. Xeroderma Pigmentosum tumors used in the analysis.**

877

878

879

## 880 **FIGURE LEGENDS**

881

### 882 **Figure 1. Mutation landscape of the studied cancers.**

883 **a** Tumor mutation burden of single base substitutions (SBS; left panel), double base substitutions (CC>TT; middle panel) per group and a proportion of CC>TT DBS relative to C>T SBS in pyrimidine context (right panel). *P*-values from nonparametric ANOVA are indicated. **b** Trinucleotide-context mutation profiles of SBS (left panel) and tetranucleotide-context mutation profiles of CC>TT DBS (right panel) per group. **c** Multidimensional scaling (MDS) plot based on the Cosine similarity distance between the SBS trinucleotide-context mutation profiles of the samples. **d** MDS plot based on the Cosine similarity distance between the trinucleotide-context mutation profiles of the samples using only C>T mutations with an adjacent pyrimidine (YC>YT or CY>TY), the typical UV mutation context. **e** MDS plot based on the Cosine similarity distance between the tetranucleotide-context mutation profiles of the samples using only CC>TT double base substitutions. **f** Mean Cosine dissimilarity (1-Cosine distance) between original and reconstructed trinucleotide-context mutation profiles using only SBS7a/b/c/d COSMIC mutation signatures for all SBS (upper panel) and C>T mutations with adjacent pyrimidine only (lower panel).

896

### 897 **Figure 2. Genomic topography of mutagenesis in the skin cancers.**

899 **a** Fraction of C>T mutations from pyrimidine dimers in genomic regions grouped in 8 equal size  
900 bins by replication timing (RT) for XP groups and sporadic skin cancers. The box contains the  
901 slope values from linear regressions across 8 RT bins. **b** Fraction of C>T mutations from  
902 pyrimidine dimers per group in 1Mb regions centered at the boundary between active and inactive  
903 topologically-associated domains (split into two bins each). **c** Fraction of C>T mutations from  
904 pyrimidine dimers per group across different chromatin states (R - repressed, A and A2 - active,  
905 H - heterochromatin, I - inactive). **d** Fractions of C>T mutations from pyrimidine dimers in  
906 intergenic regions (left panel), on the untranscribed (middle panel) and transcribed (right panel)  
907 DNA strands of gene regions grouped in 5 equal size bins by replication timing (RT) for XP groups  
908 and sporadic skin cancers. The boxes contain the slope values from linear regressions across 5  
909 RT bins. I: intergenic regions; NTR: untranscribed strand of genes; TR: transcribed strand of  
910 genes.

911

### 912 **Figure 3. TC-NER activity behind transcription end sites (TES) of genes.**

913 **a** The transcriptional bias (TRB) per group (ratio between untranscribed and transcribed strand)  
914 for C>T mutations with adjacent pyrimidines for XP groups and sporadic skin cancers. *P*-values  
915 from nonparametric ANOVA are indicated **b** Fractions of C>T mutations with adjacent pyrimidines  
916 separated by strands in the TES-centered 100kb region (binned by 10kb intervals). **c** DNA  
917 context-normalized XR-seq density from XP-C cell line on untranscribed and transcribed gene  
918 strands in the TES-centered 100kb region (binned by 10kb intervals; left panel). DNA context-  
919 normalized XR-seq density from XP-C cell line by replication timing for the transcribed and  
920 untranscribed DNA strands of genes and intergenic regions. I: intergenic regions; NTR:  
921 untranscribed strand of genes; TR: transcribed strand of genes (right panel). **d** Correlation  
922 between XR-seq intensity from XP-C cell line and nascent RNA-seq for genic regions (left panel)  
923 and intergenic regions 50kb downstream of TES (right panel). **e** Transcriptional bias of C:G>T:A  
924 mutations on intergenic regions of XP-C tumors depending on the XR-seq intensity of XP-C cell  
925 line. **f** Relative mutation rate of C:G>T:A mutations in intergenic regions of XP-C tumors  
926 depending on the XR-seq intensity in XP-C cell line.

927

### 928 **Figure 4. Comparison of genomic mutagenesis between sporadic cancers, XP-E and XP-C** 929 **groups.**

930 **a** Multidimensional scaling (MDS) plot based on the Cosine similarity distance between the SBS  
931 trinucleotide-context mutation profiles of the samples (Dimensions 1 and 2 - left panel,  
932 Dimensions 1 and 3 - right panel). **b** PCA plot based on the density of mutations in 2684 1Mb-  
933 long windows along the genome (only for samples with more than 50k mutations belong to  
934 sporadic, XP-C and XP-E groups). **c** The transcriptional bias (TRB; ratio between untranscribed  
935 and transcribed strand mutation number) for C>T mutations from pyrimidine dimers in genes  
936 grouped in 6 bins by gene expression level. Only cSCC tumors were used for XP-C and XP-E  
937 groups **d** Fractions of C>T mutations from pyrimidine dimers separated by strands in the TSS-  
938 centered 100kb region (binned by 10kb intervals). **e** The slope values from linear regressions



939 across C>T mutations from pyrimidine dimers over binned epigenetic features for the whole  
940 genome (left panel), intergenic regions (left middle panel), untranscribed (right middle panel), and  
941 transcribed (right panel) strands of genes separately (only cSCC from sporadic, XP-E and XP-C  
942 groups were used in the analysis). *P*-values based on the Student's t-test pairwise comparisons  
943 between sporadic cSCC and cSCC from XP-C or XP-E groups are indicated.

944

945 **Figure 5. Mutation profiles of XP-V skin cancers and *POLH*-KO clones.**

946 **a** Trinucleotide-context mutation profile of genomic SBS (upper panel) and genic SBS (lower  
947 panel) separated by transcribed (TR) and untranscribed (NT) strands in XP-V tumors. **b** Fractions  
948 of C>A mutations separated by gene strands in the TSS-centered 100kb region of XP-V tumors  
949 (binned by 10kb intervals). Blue – transcribed strand for mutations from purines and untranscribed  
950 strand for mutations from pyrimidines; red - untranscribed strand for mutations from purines and  
951 transcribed strand for mutations from pyrimidines **c** The transcriptional bias (ratio between  
952 transcribed and untranscribed strand) for C>A and C>T mutations per bin of gene expression  
953 level (only XP-V samples represented by SCC and BCC). **d** Trinucleotide-context mutation  
954 profiles of SBS separated by strands in XP-V tumors for C>A and T>A mutations. **e** Mutations per  
955 megabase in the *POLH* wt and *POLH*-KO clones in nontreated cells (NT), treated with KbrO<sub>3</sub>,  
956 UV-A and UV-C. **f** Mutational specificity of the  $\underline{\text{TG}}>\underline{\text{TI}}$  mutations in XP-V tumors and *POLH*-KO  
957 UV-A- and UV-C-treated cell lines. X-axis: log<sub>2</sub>-transformed transcriptional bias of the  $\underline{\text{TG}}>\underline{\text{TI}}$   
958 mutations per genome. Y-axis: Fraction of the mutations in the  $\underline{\text{TG}}>\underline{\text{TI}}$  context from the total  
959 number of C:G>A:T substitutions per genome. *POLH*-KO and *POLH*-wt clones are specifically  
960 indicated with their corresponding treatment (KbrO<sub>3</sub>, UV-A and UV-C) as well as COSMIC SBS18  
961 and SBS36 mutational signatures associated with oxidative DNA damage (black dots).  
962 **g** Mutation profiles of the *POLH*-wt and *POLH*-KO clones for nontreated cells (NT), treated with  
963 KbrO<sub>3</sub>, UV-A and UV-C.

964

965 **Figure 6. Dimer translesion bias in XP-V skin cancers.**

966 **a** Schematic representation of the putative CC photodimer in ACCA context and resulting  
967 mutations analyzed in the panel b. **b** Fraction of C>T mutations from 5' and 3' cytosines of the  
968 dimer in the <sup>5</sup>ACCA<sup>3</sup> context per group of tumors. **c** "Dimer translesion bias" for different  
969 sequence contexts per group of tumors. Comparison of C>T mutation frequency in CT and TC  
970 pyrimidine dimers was performed after normalization to the number of such contexts in the  
971 genome (upper right panel). **d** Fraction of C>T mutations from 5' and 3' cytosines of the dimer in  
972 the <sup>5</sup>ACCA<sup>3</sup> context in the RPE-1 *POLH*-wt and *POLH*-KO clones.

973

974 **Figure 7. Protein-damaging effect of mutation contexts.**

975 **a** Correlations between tumor mutation burden and number of oncogenic and likely oncogenic  
976 mutations in the studied skin cancer samples according to the OncoKB database. **b** Mean fraction  
977 of exonic mutations from all the mutations per sample. **c** Protein-damaging/silent mutation ratio  
978 per substitution type in our skin cancer cohort. Damaging mutations - all non-silent exonic

979 (missense, truncating) and splice site mutations. **d** Mean fraction of protein-damaging mutations  
980 originating from the main mutation classes split by gene strand per group.

## 981 **SUPPLEMENTARY FIGURE LEGENDS**

982

983 **Supplementary Figure 1.** Trinucleotide-context mutation profiles of SBS for each tumor from XP  
984 patients.

985

986 **Supplementary Figure 2.** Fractions of C>T mutations from pyrimidine dimers in intergenic  
987 regions (INT, grey color), on the untranscribed (NTR, red color) and transcribed (TR, blue color)  
988 DNA strands of gene regions grouped in 5 equal size bins by replication timing (RT) for XP groups  
989 and sporadic skin cancers.

990

991 **Supplementary Figure 3.** Transcriptional bias in the TES-centered 100kb region (binned by 10kb  
992 intervals).

993

994 **Supplementary Figure 4.** Correlations between different types of substitutions in specific  
995 contexts in XP-V tumors.

996

997 **Supplementary Figure 5.** Scheme of the mutation accumulation experiment with *POLH* KO and  
998 *POLH* wt cell lines.

999

1000 **Supplementary Figure 6.** Double base substitution (DBS) profiles of XP and sporadic skin  
1001 tumors from fresh-frozen samples (upper panel) and RPE-1 mutation accumulation experiment  
1002 (lower panel). Only fraction from 0 to 0.3 is shown.

1003

1004

1005

1006

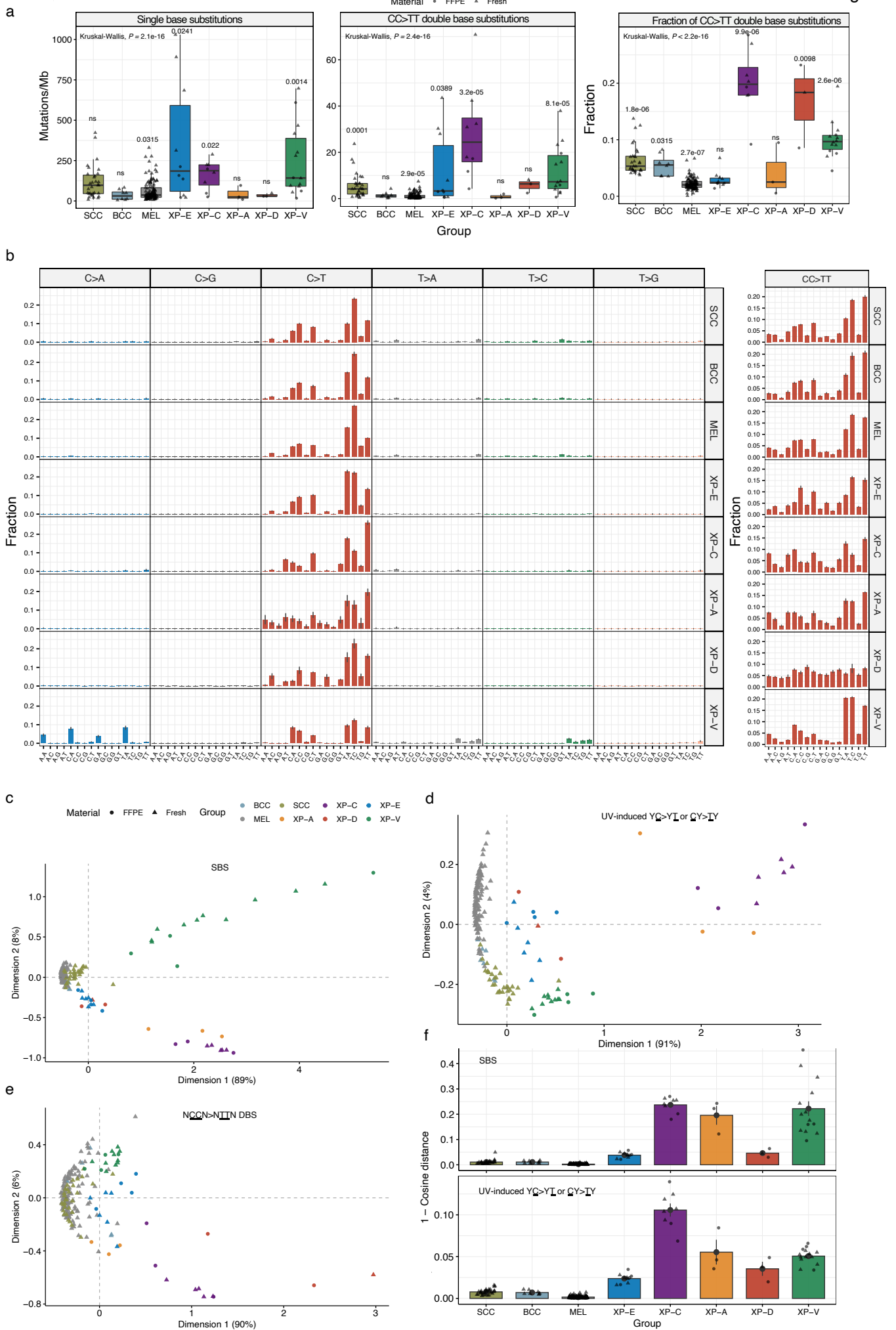


Figure 2

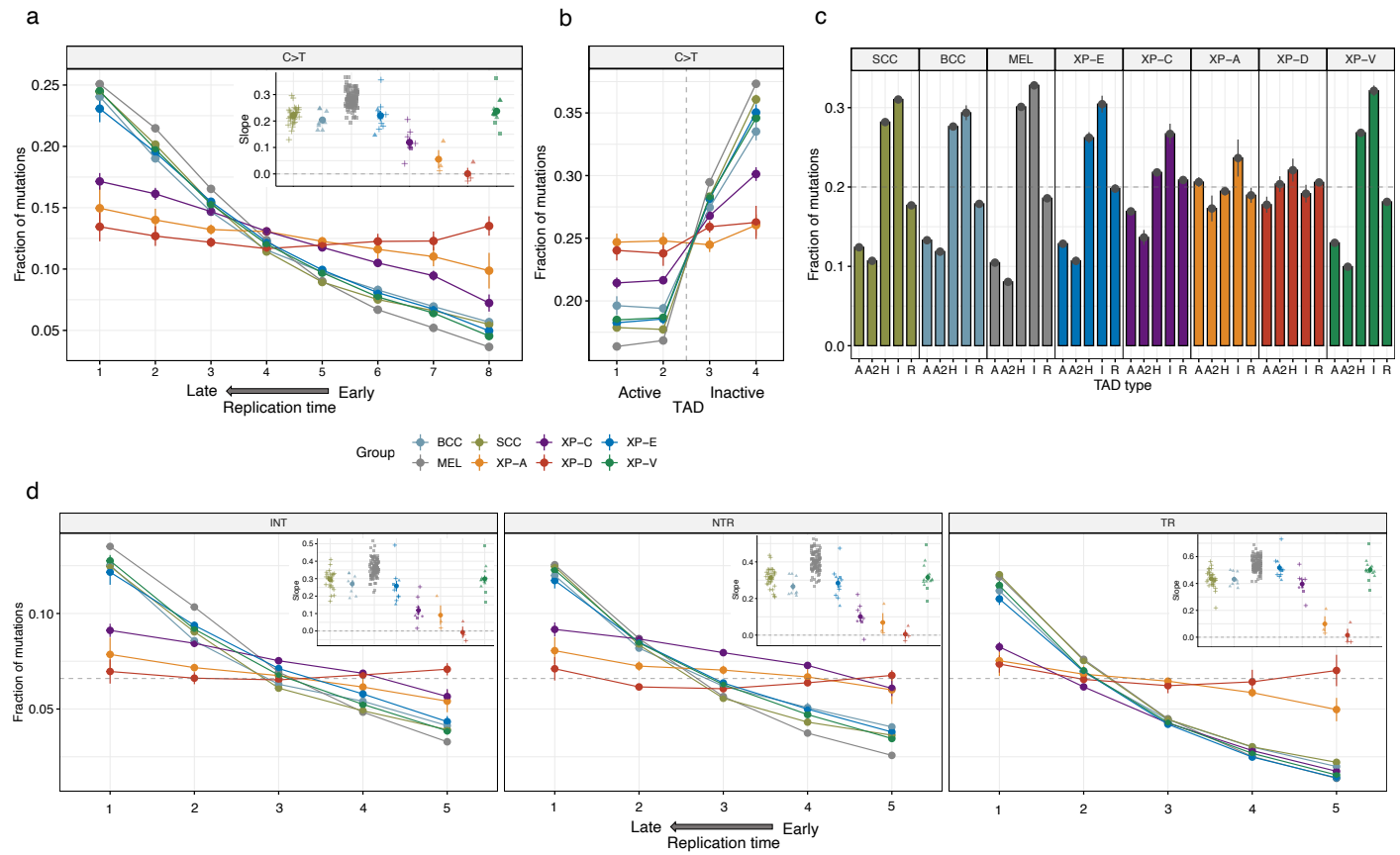


Figure 3

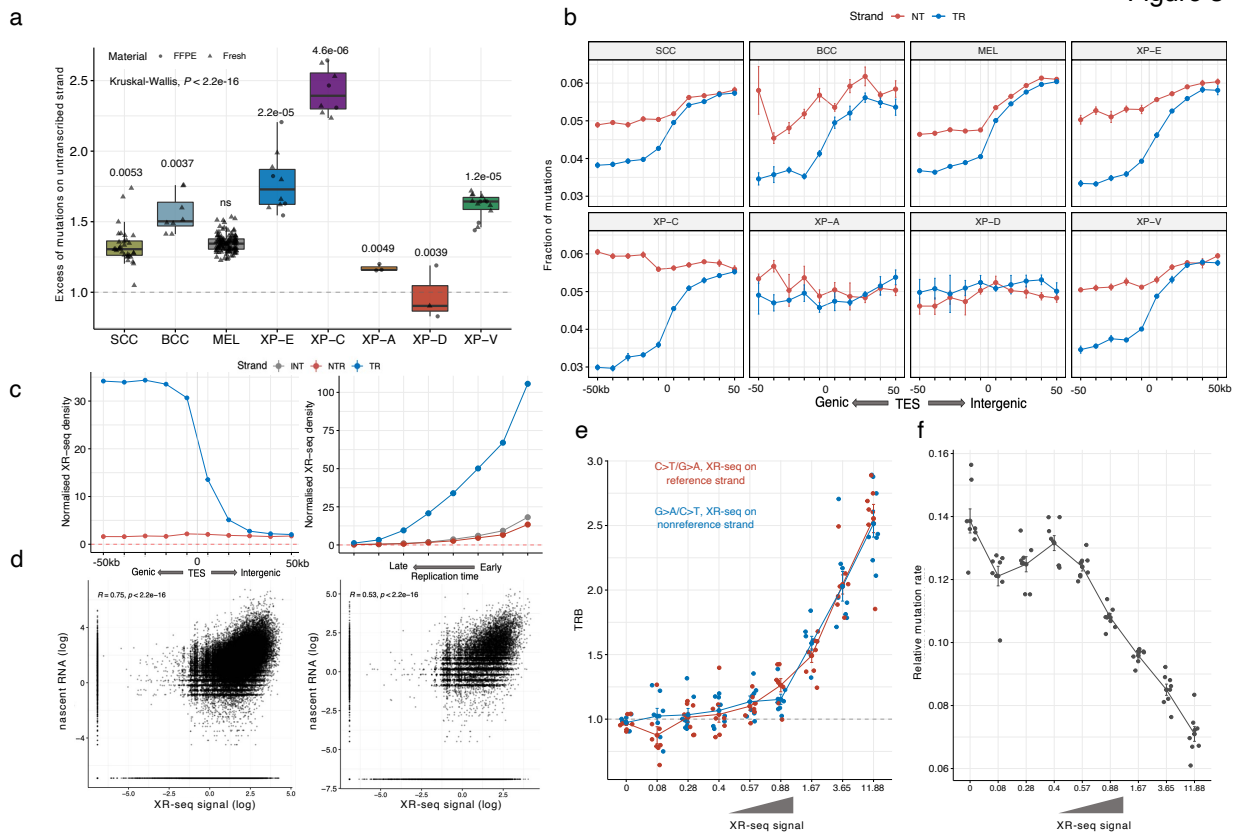


Figure 4

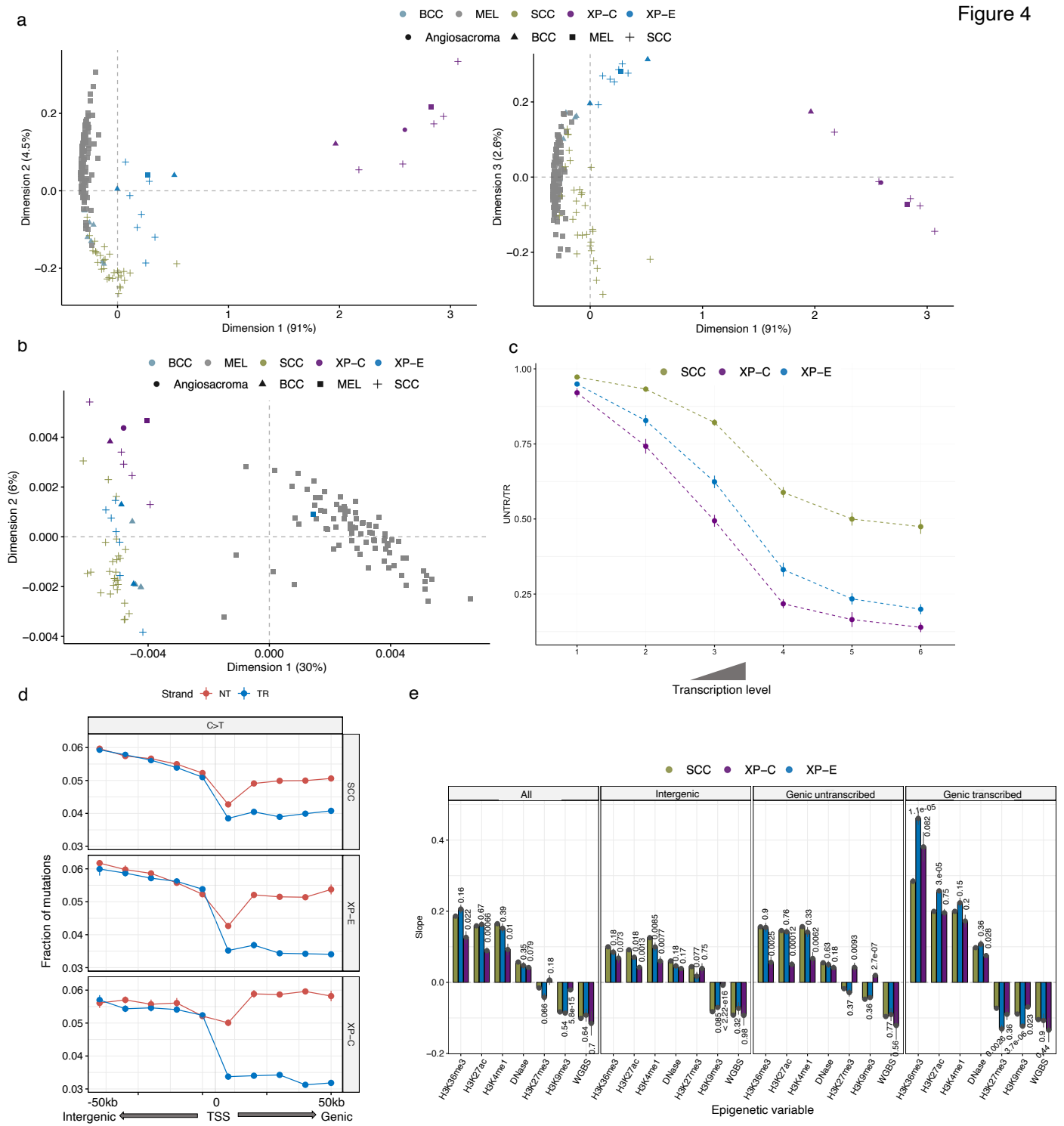


Figure 5

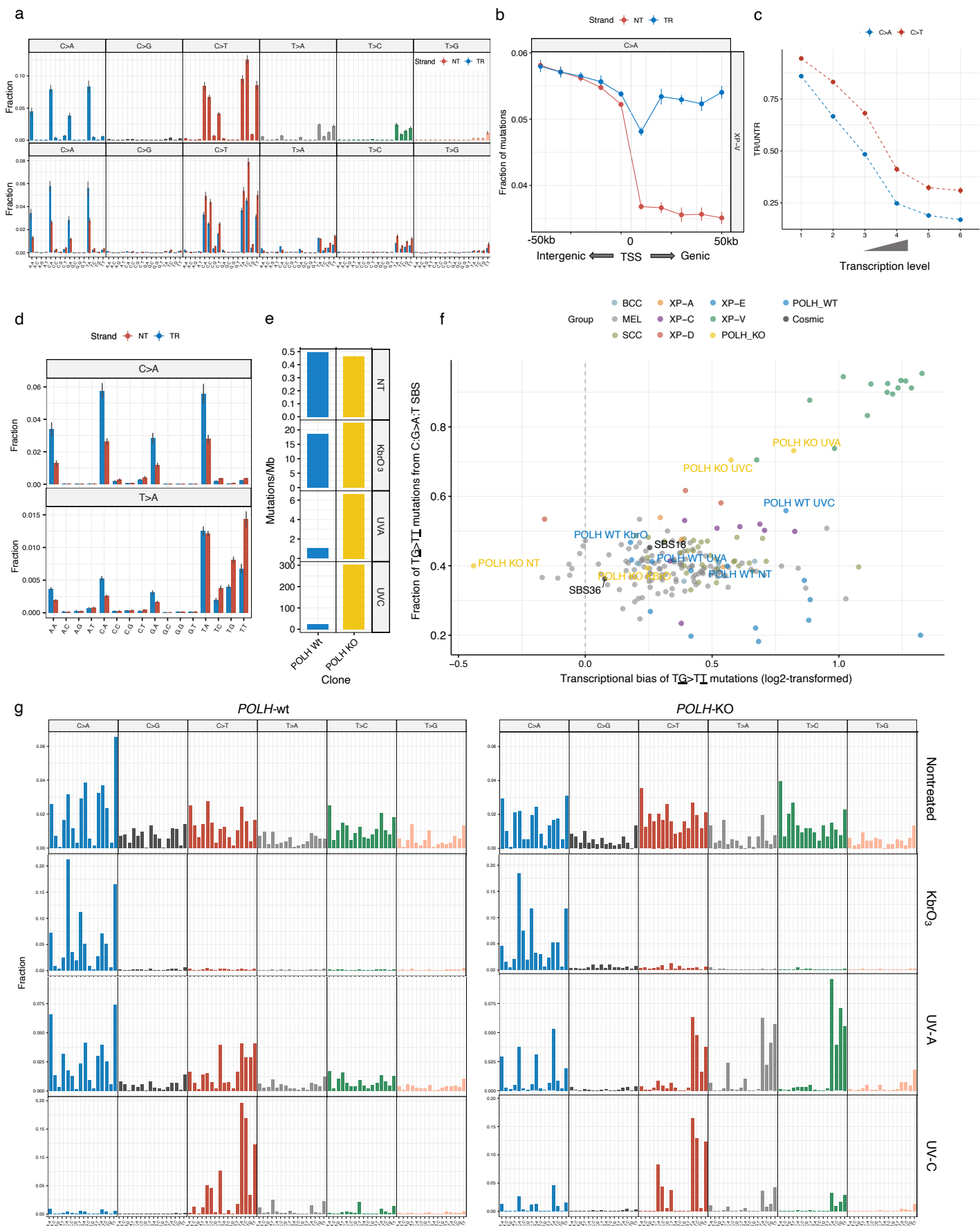


Figure 6

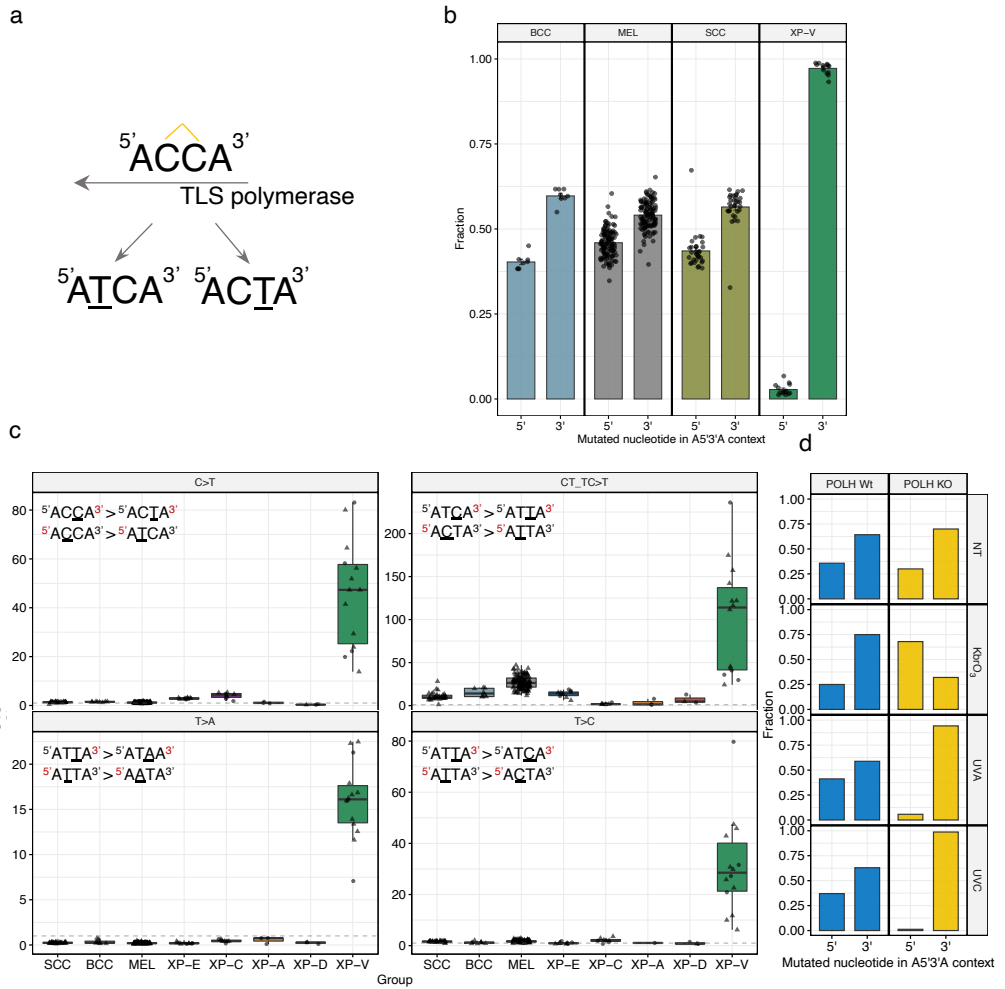


Figure 7

