

# 1 **Redistribution of mutation rates across chromosomal domains in human** 2 **cancer genomes**

3

4 Marina Salvadores<sup>1</sup>, Fran Supek<sup>1,2</sup> \*

5

6 <sup>1</sup> Genome Data Science, Institute for Research in Biomedicine (IRB Barcelona), Barcelona  
7 Institute of Science and Technology, 08028 Barcelona, Spain.

8 <sup>2</sup> Catalan Institution for Research and Advanced Studies (ICREA), 08010 Barcelona, Spain.

9 \* correspondence to: [fran.supek@irbbarcelona.org](mailto:fran.supek@irbbarcelona.org)

10

## 11 **Abstract**

12

13 Somatic mutations in human have a heterogeneous genomic distribution, with increased  
14 numbers of mutations in late-replication time (RT), heterochromatic domains of chromosomes.  
15 While this regional mutation rate density (RMD) landscape is known to vary between tissues  
16 and due to deficiencies in DNA repair, we asked whether it varies between individual tumors  
17 and what would be the mechanisms underlying such variation. Here, we identified 13 RMD  
18 signatures that describe mutation redistribution across megabase-scale domains in ~4200  
19 tumors. Of those, 10 RMD signatures corresponded to groupings or subdivisions of cancerous  
20 tissues and cell types. We further identified 3 global RMD signatures of somatic mutation  
21 landscapes that transcended cancer types. One is a known general loss of RMD variation,  
22 previously associated with DNA mismatch repair failures, and was here additionally linked with  
23 homologous recombination (HR) repair deficiencies. Next, we identified a global RMD signature  
24 affecting facultative heterochromatin domains. This RMD signature strongly reflects regional  
25 variation in DNA replication time and in heterochromatin across state tumor samples, and is  
26 associated with altered cell cycle control. Finally we identified a global RMD signature  
27 associated with *TP53* loss-of-function, mainly affecting the very late RT regions. The local  
28 mutation rates in 26%-75% of cancer genes are notably changed in the tumors affected by  
29 these three global RMD signatures of mutation redistribution. Our study highlights how the  
30 plasticity of chromatin states and the RT program in cancers bears upon the regional somatic  
31 mutation rate landscape, and the downstream consequences on mutation supply to disease  
32 genes.

33

34

## 35 **Introduction**

36

37 During cancer evolution, somatic cells accumulate a number of mutations, most of them non-  
38 selected “passengers”. These somatic mutations are caused by different mutagenic processes,  
39 many of which generate higher mutation rates in late DNA replication time (RT), inactive,  
40 heterochromatic DNA. This is likely due to higher activity and/or accuracy of DNA repair in early-  
41 replicating, active chromosomal domains<sup>1,2</sup>.

42

43 These chromosomal segments are defined roughly at the megabase scale, and tend to  
44 correspond to topologically associating domains (TADs) and RT domains<sup>3-5</sup>. Regional mutation  
45 density (RMD) of mutations in megabase-sized domains in the human genome correlates with  
46 domain RT, local gene expression levels, chromatin accessibility (as DNase hypersensitive  
47 sites (DHS)), density of inactive histone marks such as H3K9me3 and inversely with density of  
48 active marks such as H3K4me3<sup>1,6-8</sup>. The RMD signatures have been shown to be tissue-  
49 specific, and can be used to predict cancer type, and potentially subtype at high accuracy<sup>9,10</sup>.  
50 The tissue-specificity of RMD is paralleled in the tissue-specificity of active or inactive domains.  
51 For instance, the domain that switches from late-RT to early-RT, or where genes increase in  
52 expression levels, or that gets more accessible chromatin in a particular tissue, also exhibits a  
53 reduced rate of somatic mutations in that tissue<sup>1,6</sup>; this property may help identify the cell-of-  
54 origin of some cancers<sup>11</sup>.

55  
56 Apart from variation in active chromatin and gene expressions between tissues, recent work  
57 suggests existence of gene expression programs that are variably active between tumors  
58 originating from the same tissue (and also between individual cells), but are recurrently seen  
59 across many different tissues<sup>12,13</sup>. Such programs may conceivably drive, or be driven by  
60 chromatin remodeling that activates or silences chromosomal domains. Indeed, chromatin  
61 remodeling was widely reported to occur during tumor evolution, and this can manifest as  
62 changes in RT between normal and cancerous cells, loss of DNA methylation in some  
63 chromosomal domains with cell cycling, as well as a generalized loss of heterochromatin upon  
64 transformation<sup>14-18</sup>. These changes in RT, DNA methylation and heterochromatin occurring in  
65 cancer cells may plausibly affect chromosomal stability, given the links of various DNA damage  
66 and repair processes and chromatin organization<sup>1,2,16,19-21</sup>.

67  
68 Here, we hypothesized that chromatin remodeling that occurs variably between tumors may  
69 generate inter-individual variation in regional mutation rates, beyond the tissue identity or cell-of-  
70 origin identity effects on mutagenesis.

71  
72 We study the RMD profiles at the megabase scale of somatic mutations from tumor whole-  
73 genome sequences, modeling this mutational portrait as a mixture of several underlying regional  
74 distributions, which may correspond to different mechanisms that produce or prevent mutations  
75 preferentially in some genomic domains. To disentangle these distributions, we apply an  
76 unsupervised factorization approach and extract RMD signatures from ~4200 whole genome  
77 sequenced human tumors. Some of these RMD signatures represent the expected differences  
78 between tissues/cell types, or they may represent consequences of common DNA repair  
79 failures. However others are novel and are associated with RT variation and with chromatin  
80 remodeling upon cell cycle disturbances. We characterize the differences between individuals in  
81 the usage of these different RMD distributions of mutations, suggesting that the chromatin  
82 remodeling RMD signatures are ubiquitous amongst human cancers. They reflect wide-spread  
83 mutation redistribution across domains and affect mutation supply to regions harboring cancer  
84 genes.

85  
86

## 87 Results

88

### 89 Inter-individual variability in megabase-scale regional mutation density in human tumors

90

91 We hypothesize that, in addition to the variability between cancer types, the RMD patterns  
92 encompass variability between individuals that is observed across many tissues. To test this, we  
93 performed a global unsupervised analysis of diversity in one-megabase (1 Mb) RMD patterns  
94 across 4221 whole-genome sequenced tumors that had a mutation burden >3 single-nucleotide  
95 variants (SNV) per Mb. To prevent confounding by the variable SNV mutational signatures  
96 across tumors<sup>22</sup> we controlled for trinucleotide composition across the 1 Mb windows  
97 (Methods). We additionally normalized the RMDs at chromosome arm-level to control for  
98 possible confounding of large-scale copy-number alterations (CNA) on mutation rates. Finally  
99 we removed known mutation hotspots (e.g. CTCF binding sites, see Methods), and also exons  
100 of all protein-coding genes to reduce effects of selection.

101

102 To quantify the systematic variability contained within tumor RMD landscapes, we applied a  
103 Principal Component (PC) analysis on the RMD profiles across all tumor samples (n=4221).  
104 Expectedly, most of the 22 relevant PCs (those with a % of variance explained higher than a  
105 random baseline (Fig 1a)), separated different tissues (Fig 1b, Fig S1a). However, we found  
106 some PCs that captured variability between individuals but not between the known tissue-of-  
107 origin of tumors (Fig S1a). Serving as a positive control, the PC1 separated the canonical RMD  
108 landscape with increased mutation rates in late-replicating DNA versus the known “flat”  
109 landscape of tumors with failed DNA mismatch repair (MMR)<sup>1</sup> (Fig 1c). Next, we observed that  
110 PC7 separates lymphoid tumors with higher somatic hypermutation (SHM) activity (Fig 1c),  
111 using the exposure of mutational signature SBS9 as a proxy for prior activity of SHM in that  
112 lymphoma sample<sup>22</sup>. Reassuringly, we observed that the PC7 1 Mb window weights are  
113 strongest in known SHM regions containing antibody genes (Fig 1d). In summary, our RMD  
114 features were able to capture two known examples of regional redistribution of mutations: one  
115 affects specific sites (SHM regions in B-lymphocytes) and the other causes a global ‘flattening’  
116 of mutation rate landscape along the genome in MMR-deficient samples, supporting the utility of  
117 our RMD profiling method.

118

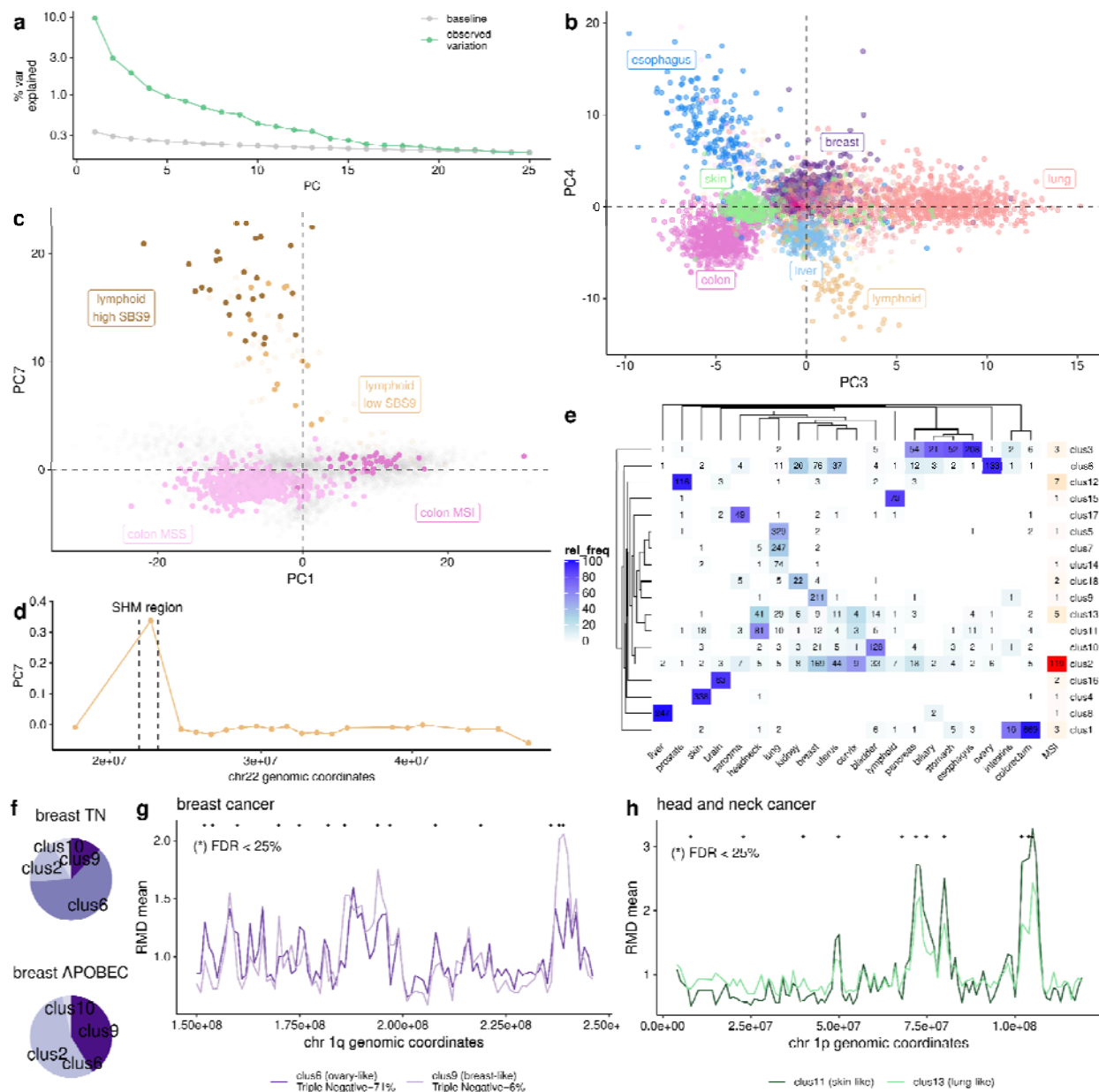
119 Next, we asked if clustering the tumor samples by their RMD feature vectors would reveal, in  
120 addition to an expected grouping by tissue, also other sources of inter-individual variability in  
121 RMD (Fig 1e). For clustering, we selected first 22 RMD PCs based on the amount of variance  
122 explained (Fig 1a), based on the PC window weights’ autocorrelation with neighboring windows  
123 (Fig S1b) indicating a nonrandom organization of the RMD pattern along the chromosomes; and  
124 based on additional criteria (Fig S1c-d). This revealed there exist 3 different types of clusters.  
125 On the one extreme, for example the RMD\_cluster2 contained samples from almost all cancer  
126 types, and was very enriched with MSI (MMR-deficient) samples. On the other extreme, there  
127 were tissue-specific clusters that contained only a single cancer type (e.g. the liver  
128 RMD\_cluster8) (Fig 1e). Interestingly, there was also the third, intermediate case with clusters  
129 that spanned several, apparently similar cancer types (e.g. RMD\_cluster3 with various digestive  
130 tract cancers, or the squamous-like RMD\_cluster11, with head-and-neck cancers, the non-

131 melanoma skin cancers and some esophagus and lung cancers) (Fig 1e). Therefore, there is  
132 information in the RMD feature vector that can transcend the tissue-of-origin, in this case uniting  
133 similar tissues or cell types.

134  
135 In addition to RMDs bridging cancer types, conversely RMD profiles can be used to subdivide  
136 some cancer types such as breast cancer. Breast cancers in RMD\_cluster2 have high APOBEC  
137 mutagenesis (Fig 1f), thus sharing cluster with MSI samples; APOBEC mutagenesis has been  
138 reported to change the regional mutational landscape by preferring early replicating regions<sup>23,24</sup>.  
139 similarly as in MSI tumors<sup>1</sup>. Furthermore, most breast cancer samples in RMD\_cluster6  
140 (ovarian-like), which have visually distinct RMD profiles from the typical breast-like  
141 RMD\_cluster9 (Fig 1f-g), are from the triple negative breast subtype, which was reported to be  
142 more similar to ovarian cancer by gene expression<sup>25</sup>. Another example of how RMD profiles  
143 can be used for subtyping is the head-and-neck cancer, which is split into RMD\_cluster11  
144 (squamous-like, includes non-melanoma skin cancers) and RMD\_cluster13 (also contains some  
145 lung cancers) (Fig 1h).

146  
147 Overall, even though the RMD profiles are tissue specific, there is systematic RMD variability in  
148 certain tumor genomes observed apparently independently of the tissue. This motivated us to  
149 devise a method that is able to extract this inter-individual variability from genomic RMD profiles,  
150 while robustly accounting for the strong tissue-specific signal in RMD.

151



**Figure 1. Chromosomal domain RMD variability across tissues and individuals.** **a)** Variance explained for the first 25 PCs of a PCA on the RMD matrix (4221 samples x 2542 one-megabase windows), and a baseline (by the broken stick rule). **b)** PC3 and 4 separate various cancer types. **c)** As controls, the PC1 separates MSI versus MSS tumors, and PC7 separates lymphoid samples according to their level of the SHM mutational signature (SBS9). **d)** PC7 window weights for chromosome 22 agree with the known SHM region. **e)** Number of tumor samples from each cancer type that are assigned to each RMD-based cluster (Methods). **f)** Cluster assignment for breast cancer samples of triple negative (TN) subtype and samples with high APOBEC (>25% of mutations are in APOBEC contexts) **g)** Mean RMD profiles for breast cancer samples in cluster 6 (n = 76) and cluster 9 (n = 211), shown for chr 1q. **h)** Mean RMD profiles for head and neck squamous samples in cluster 11 (n = 81) and cluster 13 (n = 41), for chr 1p.

152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165

## 166 **A methodology to detect inter-individual variation in regional mutation density**

167  
168 To separate the inter-individual RMD variability from the tissue-specific variability we applied a  
169 methodology analogous to that recently used for extracting trinucleotide SNV mutational  
170 signatures<sup>22,26,27</sup> however here applied to megabase-sized domains. In brief, non-negative  
171 matrix factorization (NMF) is repeatedly applied to bootstrapped mutational data, to find  
172 solutions (sets of factors) that are consistent across bootstrap runs. These solutions contain  
173 multiple RMD signatures (factors), each with RMD window weights (all 1 Mb windows with  
174 varying contributions) and RMD sample ‘exposures’ or activities (the weight of each tumor for  
175 that signature).

176  
177 To test whether our NMF method is sufficiently powered to capture RMD inter-individual  
178 variability, we simulated cancer genomes containing known, ground-truth patterns of RMD that  
179 affected a variable number of windows, being present in variable number of tumor samples, and  
180 at variable intensity (fold-increase over canonical mutation rates) (Fig S2a, see detailed  
181 description in Methods). We ran our NMF methodology for these different scenarios  
182 independently. We selected the number of factors and clusters based on a clustering quality  
183 measure, the silhouette index (SI), over multiple runs of NMF (Fig S2b) and matching the known  
184 ground-truth signatures (Methods, Fig S3). We show an example of an extracted RMD signature  
185 compared to its matching ground-truth signature in Fig 2a.

186  
187 By comparing the different scenarios (Fig S4), encouragingly, we observed that even with a  
188 small fraction of samples affected (5%), the ground-truth RMD signatures can be identified  
189 reliably, as long as the contribution of the RMD signature to the total mutation burden is high  
190 ( $\geq 20\%$ ). In addition, we observed that the NMF setup is very robust to the number of windows  
191 affected and is usually able to recover RMD signatures that affect as little as 10% of all  
192 windows. Out of other characteristics that may affect power to recover RMD signatures, we  
193 identified the signature strength/exposure (fold-enrichment) as showing the highest effect, thus  
194 the signatures with subtle effects on RMD might not be recovered (Fig S4). In summary, our  
195 simulations support that our NMF-based methodology can recover the genome-wide RMD  
196 signatures in a wide variety of tested scenarios.

## 197 198 199 **Three prevalent patterns of megabase-scale mutation rate variation observed across** 200 **most somatic tissues**

201  
202 We applied the NMF methodology to the somatic RMD profiles of 4221 tumor WGS, here  
203 requiring a minimum of 3 mutations/Mb per sample thus restricting to tumors with less noisy  
204 RMD profile (as a limitation, we note that this may exclude samples from some cancer types  
205 preferentially). In total, we extracted a total of 13 RMD signatures based on the silhouette index  
206 that scores the reproducibility of solutions upon 100 bootstraps (Fig 2bc, Fig S5).

207  
208 In accordance with the above RMD clustering analysis (Fig. 1e), we observed that the RMD  
209 signatures from NMF span a continuum from very tissue specific (high Gini index, Fig 2c), to



210 global signatures (low Gini index). We named ten signatures according to the tissue or tissues  
211 they affect (e.g. RMD\_upper-GI, RMD\_liver), while the three global signatures that affect many  
212 cancer types were named RMDglobal1, RMDglobal2 and RMDflat (Fig 2c, Fig S5) (the latter is  
213 named by the visually recognizable pattern, and also has in part known mechanisms; see  
214 below).

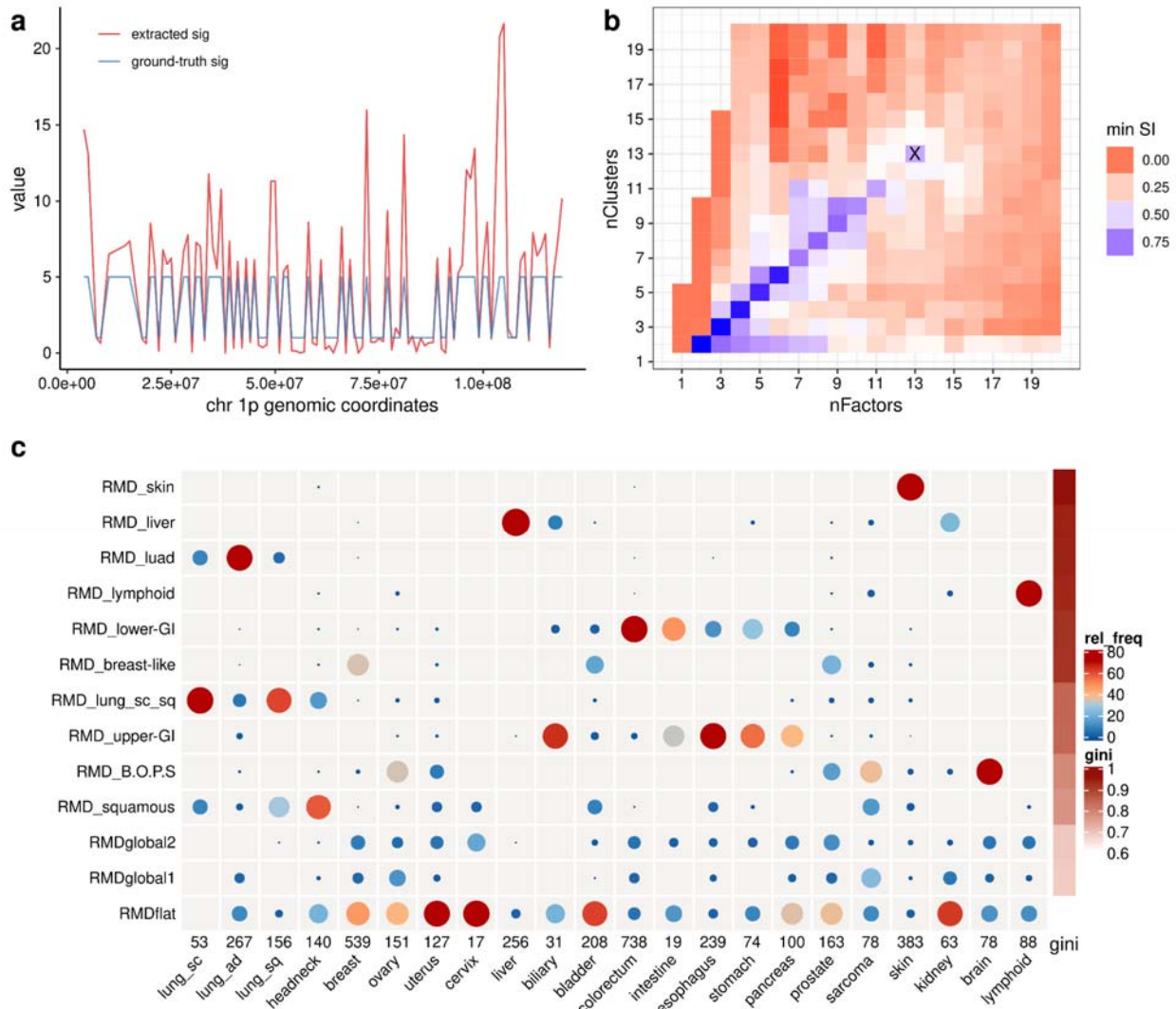
215  
216 In one extreme, there are tissue-specific signatures (e.g. RMD\_skin, RMD\_liver) which capture  
217 the genomic regions with an increase of mutations only in that particular cancer type (e.g. skin  
218 in RMD\_skin, or liver and some biliary and some kidney cancers in RMD\_liver) (Fig 2c, Fig S5).  
219 Windows in these RMD signatures could be used to improve cancer-type classification of  
220 tumors based on regional mutation density from WGS data <sup>9</sup>.

221  
222 Between the two extremes, there are signatures present in several cancer types which are  
223 apparently similar (Fig 2c, Fig S5). For instance, RMD\_upper-GI signature is present in most  
224 esophagus, stomach, pancreas and biliary tumor samples, and some intestine tumors. The  
225 RMD\_lower-GI, in turn, contains mainly the colorectal and most of the intestinal tumors, broadly  
226 consistent with the subdivision by developmental origin into the foregut (RMD\_upper-GI) and  
227 the midgut/hindgut (RMD\_lower-GI; Fig 2c). The RMD\_squamous signature spans some  
228 squamous lung cancers, head-and-neck cancers, some bladder cancers (consistent with reports  
229 based on gene expression data <sup>28</sup>, also expectedly some cervical and esophageal tumors, and  
230 surprisingly some sarcomas and uterus cancers. Interestingly, one signature, provisionally  
231 named “RMD\_B.O.P.S.”, spans brain (B), ovarian (O), prostate (P), sarcomas (S), and uterus  
232 cancers and so probably reflects a convergent phenotype rather than a common cell-of-origin.  
233 These examples suggest that there are commonalities in mutation rates, probably reflecting  
234 chromatin organization in the cell-of-origin of tumor types. These commonalities usually reflect  
235 anatomical subdivisions or cell type similarity, and shape the RMD profiles of those samples.  
236 Our RMD signatures support the proposed uses of RMD profiles for elucidating cell-of-origin  
237 and cancer development trajectories (e.g. metaplasia and/or invasion) <sup>11</sup> by matching to  
238 chromatin profiles.

239  
240 In the other extreme, we identified 3 global RMD signatures, which capture the inter-individual  
241 RMD variability within most cancer types (Fig 2c, Fig S5). While the profile of RMDflat captures  
242 the known “flat” RMD landscape (i.e. a low variation in mutation rates between segments) profile  
243 associated with MMR and NER failures <sup>1,2</sup>, RMDglobal1 and RMDglobal2 profiles have an  
244 apparently complex pattern with their peaks appearing distributed throughout the chromosomes.  
245 We can rule out that RMDglobal1 and 2 are due to random noise, because (a) the silhouette  
246 index of RMDglobal1 and 2 (measuring robustness of their profile to noise that is introduced in  
247 repeated NMF runs) is comparable to the other RMD signatures, and (b) the autocorrelation of  
248 their profiles (measuring similarity in weights of consecutive 1 Mb windows) is comparable to the  
249 other, tissue-associated RMD signatures (Fig S6a-b).

250  
251 In addition to the pan-cancer analysis, we ran NMF for each cancer type independently, for the  
252 12 cancer types with more than 100 genomes meeting criteria (Fig S7). All three global  
253 signatures can be found also in the per-cancer-type NMF runs (Fig S8). We found signatures in

254 breast, lung and esophagus with a cosine similarity > 0.84 with RMDglobal1, and in colon,  
 255 uterus and breast with a cosine similarity > 0.89 with RMDglobal2, supporting that the global  
 256 RMD signatures capture inter-individual RMD variation recurrently observed in various somatic  
 257 tissues.  
 258



259  
 260 **Figure 2. Identifying RMD signatures by an application of a NMF-based methodology to WGS of**  
 261 **human tumors. a)** Example signature from a simulation study, comparing window weights for an  
 262 extracted NMF signature and its matching simulated ground-truth signature along chr 1p. See  
 263 Supplementary Figs 2-4 for additional simulation data. **b)** NMF run on data from 4221 human tumors.  
 264 Minimum silhouette index (SI) across clusters (RMD signatures) for different numbers of NMF factors and  
 265 clusters. Selected case (nFactor=13, nCluster=13) is marked with a cross. **c)** Overview of the 13 RMD  
 266 signatures extracted (rows) and their distribution across different cancer types (columns). The circle size  
 267 and, equivalently, color corresponds to the fraction of samples from a specific cancer type exhibiting a  
 268 specific signature (signature exposure  $\geq 0.177$ ). Total number of samples per cancer type written  
 269 beneath table. The Gini index quantifies the distribution of the signature across different cancer types;  
 270 higher index means more specificity to few cancer types.



271  
272 **Homologous recombination-deficient tumors show lower regional mutation rate**  
273 **variability**

274  
275 The RMDflat global signature we extracted from the NMF analysis captures the ‘flat’ distribution  
276 of mutations that was reported for MSI tumors, which are deficient in MMR <sup>1</sup>. The window  
277 weights of RMDflat correlated with the average RT (Fig 3a), opposite of the canonical RMD  
278 landscape (which has few mutations in early-replicating DNA), thus the additive combination of  
279 the two results in a flat, low-variation landscape.

280  
281 As expected, MSI samples showed high exposures to this RMD signature (Fig 3b), as well as  
282 bladder samples with mutations in the *ERCC2* gene, participating in the NER pathway (Fig 3b),  
283 consistent with previous reports <sup>1,29</sup>. In addition, we observed that tumor samples with high  
284 APOBEC signature mutagenesis also showed high exposures to the RMDflat signature (Fig 3b),  
285 solidifying prior reports of APOBEC mechanisms being enriched in early-replicating DNA,  
286 possibly via their association with DNA repair activity providing ssDNA substrate for APOBECs  
287 <sup>23,24,30</sup>.

288  
289 Based on these known associations involving MMR, NER, and APOBEC activity, we  
290 hypothesized that some of the remaining unexplained cases of RMDflat-high tumors (total 52%  
291 were explained) may be associated with deficiencies in another DNA repair pathway. In  
292 particular, we considered homologous recombination (HR) repair deficient samples, as  
293 ascertained by the CHORD method based on SNV and CNA (but not RMD) mutational  
294 signatures <sup>31</sup>. The HR deficient samples also presented higher RMDflat exposures, both for  
295 BRCA1 and BRCA2 subtypes (Fig 3b). When HR is deficient, there is an increase in the  
296 spectrum of the trinucleotide mutational signature SBS3 <sup>22,31</sup> may result from activity of error-  
297 prone DNA polymerases <sup>32</sup>. We observed that in HR-deficient tumor samples, the SBS3-like  
298 mutational spectrum [mutation types with high weights in SBS3, such as C>G mutations]  
299 accumulate more in early replicating DNA (i.e. opposite to canonical RMD pattern) (Fig S9),  
300 thus contributing to the “flatness” of the RMD landscape.

301  
302 Thus various DNA repair related mechanisms converge onto the RMDflat phenotype, with  
303 considerable variation in prevalence depending on the tissue: in colorectal tumors the main  
304 mechanism is the MMR deficiency, while in ovary and pancreas it is the HR deficiency, and  
305 APOBEC mutagenesis is the main mechanism in bladder and lung (Fig 3c).

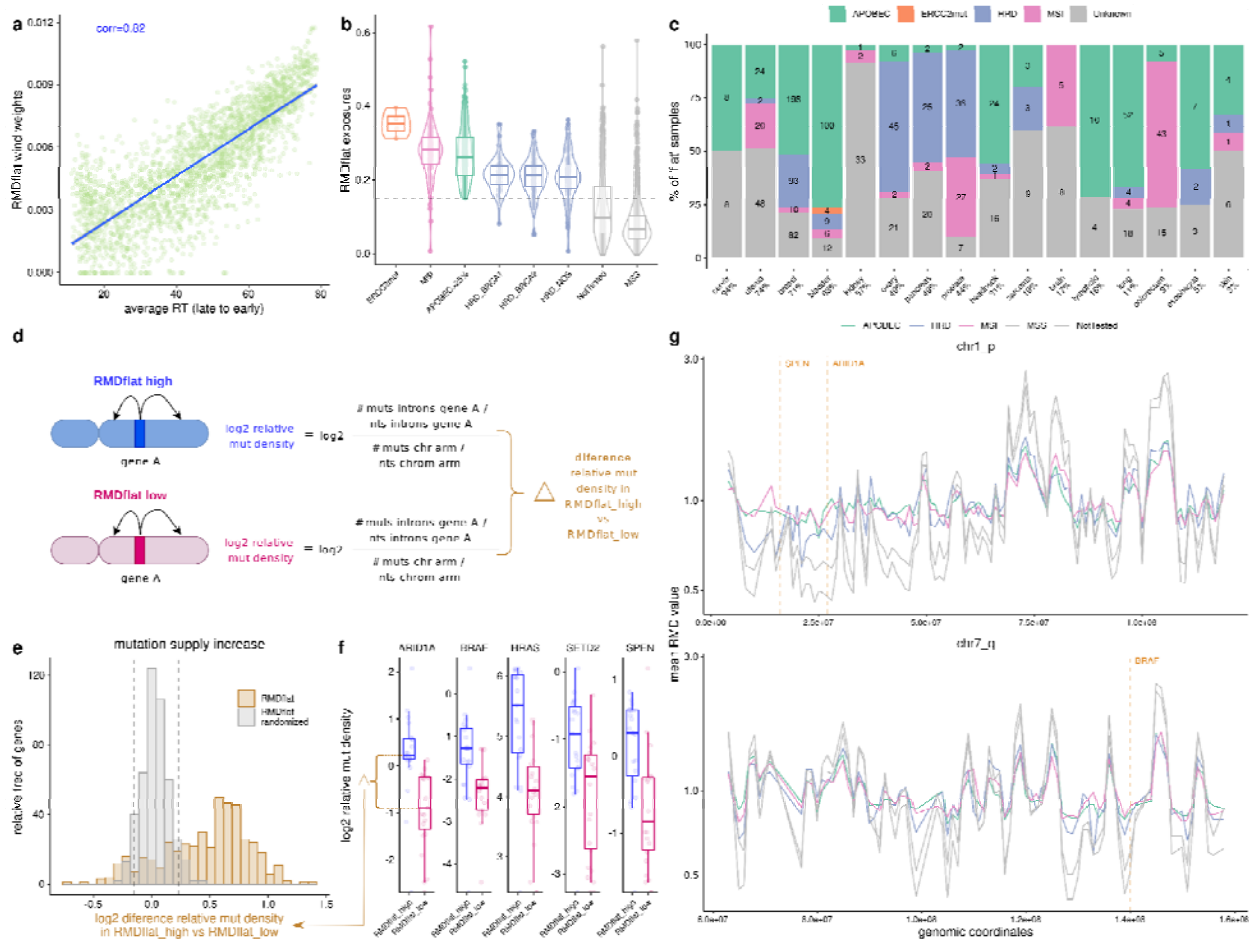
306  
307 For the final 28% of RMDflat tumor samples that are unexplained, we suggest this is unlikely to  
308 be due to false negatives in the MMR or HR deficiency tests, since these tumors had, on  
309 average, an indel spectrum (in microsatellite loci and elsewhere) not obviously different from the  
310 general indel spectrum of same cancer types (Fig S10). Thus we suggest there are other  
311 mechanism(s) involved, for instance in kidney cancer there were many unexplained RMDflat  
312 signature samples, however this cancer type very rarely has known MMR, HR deficiencies or  
313 APOBEC mutagenesis; a possible explanation is a particular mutational process in kidney <sup>33</sup>  
314 that may evades DNA repair mechanisms operative in early-RT domains.

315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344

## Mutation supply towards major cancer genes is altered by global RMDflat signatures

Tumors with RMDflat undergo an increase in mutation rates in early replicating, euchromatic regions<sup>23,34</sup>. These regions also have a higher gene density, so we quantified how RMDflat affects the mutation supply to cancer driver genes. In particular, we tested whether there is a difference in mutation density in cancer genes (considering intronic mutations, to avoid effects of selection, and further normalizing to the mutation burden of that chromosome arm to avoid effects of gross CNA; [Methods](#)), between tumor samples with a high RMDflat exposure (top tertile) versus low RMDflat exposure tumors (bottom tertile). 75% of the 460 tested cancer genes<sup>35</sup> undergo an increase in mutation supply from RMDflat-low to RMDflat-high tumors, when compared to the 95th percentile of a randomized distribution ([Fig 3d](#)). Conversely, few cancer genes decreased in mutation supply in RMDflat-high tumors (9% are below the 5th percentile of the random distribution). We considered the mutation supply density for 5 examples of common driver genes, for which mutation supply is increased 1.8-2.5 fold between RMDflat-high and RMDflat-low tumors ([Fig 3e](#)). Considering for instance the *ARID1A* tumor suppressor gene, located in a lowly-mutated region in chromosome 1p, its mutation supply increased 1.8-fold, 2.1-fold and 2.4-fold in MSI, HRD and APOBEC tumors (all RMDflat-high), respectively, compared to the *ARID1A* baseline mutation supply in tumors without DNA repair deficiencies ([Fig 3f](#)). Similarly, the *BRAF* oncogene (where causal mutations are known to be highly enriched in MSI compared to MSS colorectal tumors<sup>36</sup>) has considerably increased mutation supply in the RMDflat-high tumors ([Fig 3f](#)).

In summary, we detected the three known mechanisms that cause RMDflat (APOBEC mutagenesis, MMR and *ERCC2* deficiency) and we found an additional cause (HR deficiency) of this phenotype ([Fig 3d](#)). The consequence for tumors with RMDflat is an increase in the mutation supply for three-quarters of all cancer genes.



**Figure 3. Characterization of the RMDflat RMD signature, which represents a loss of mutation rate heterogeneity.** **a)** Correlation between RMDflat signature NMF window weights and the DNA replication timing (RT) (Repli-Seq average across 10 cell lines). **b)** RMDflat signature exposures (i.e. activities) for groups of tumor samples with various DNA repair failures (MSI, microsatellite instable, indicating MMR failure; HRD, deficient homologous recombination, by the BRCA1 type or BRCA2 type<sup>31</sup> or not otherwise specified), or high levels of APOBEC mutation signatures. **c)** Percentage of tumor samples with ‘flat’ mutation rate landscapes (RMDflat exposure>0.177, a threshold that recovers 95% of MSI samples) belonging to each of the DNA repair categories, stratified by cancer type. The percentage of ‘flat’ samples in each cancer type is indicated in x-axis labels. **d)** Schematic of the mutation supply analysis in panels e-g. **e)** Distribution of the log2 difference in the relative mutation density (intronic) for 460 cancer genes, comparing between RMDflat high tumors and RMDflat low tumors, using the actual values (“RMDflat” histogram) and randomized values (“RMDflat randomized” histogram). **f)** Log2 relative mutation density (normalized to flanking DNA in same chromosome arm, see panel d) for RMDflat-high versus RMDflat-low for 5 example genes (common drivers across >=4 cancer types and with highest effect sizes in this test). Each dot is a cancer type. **g)** Mean RMD profile across the DNA repair groups, shown examples for chr 1p and chr 7q. Vertical lines mark the position for three example genes from panel f.

345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364

## 365 **RMDglobal1 signature increases mutation rate in regions with variable replication timing**

366  
367 We were interested in the mechanism behind the RMDglobal1 signature. To elucidate this, we  
368 first tried to predict RMDglobal1 signature (the megabase window weights) from epigenomic  
369 features previously reported to associate with megabase mutation rates (reviewed in <sup>34</sup>):  
370 replication timing (RT), density of accessible chromatin (DNase hypersensitive sites, DHS) and  
371 ChipSeq data from a variety of histone marks (Fig 4a). We first tried to predict the chromosome-  
372 wide profile of the RMDglobal1 signature using the average of each feature across many  
373 epigenomic datasets, which failed to predict (Fig 4a). Predicting RMDglobal1 from each  
374 RT/DHS/ChipSeq dataset individually fared slightly better, with moderate associations ( $R^2 \sim$   
375 0.2) for certain datasets with regional density of facultative heterochromatin (H3K27me3) and  
376 constitutive heterochromatin (H3K9me3) marks (Fig 4a), suggesting a role of heterochromatin  
377 organization in determining RMDglobal1.

378  
379 Remarkably, we observed that RMDglobal1 signature can be highly accurately predicted ( $R^2$  up  
380 to 0.7) from certain features most prominently RT, DHS, and the two heterochromatin marks  
381 above, however only when predicting using multiple samples jointly (but not when predicting  
382 from the averaged feature across the samples (Fig 4a)). This suggests that RMDglobal1  
383 signature is explained by the variation between the samples for one feature, e.g. differences  
384 between the individual RT profiles. We observed the same trend using regional density of  
385 chromHMM segmentation states (Fig S11).

386  
387 The features that best predicted RMDglobal1 were the three RT datasets (Fig 4a): (i) a  
388 collection of RT profiles from experiments [RepliChip or RepliSeq] in multiple cell types (expRT,  
389  $n = 158$  samples), (ii) predicted RT in a collection of noncancerous tissues, cultured primary  
390 cells and cell lines including cancer and stem cell lines (predRT,  $n = 597$  samples), and (iii)  
391 predicted RT in human tumors (predRT-TCGA,  $n = 410$  samples, majority measured in technical  
392 duplicate). For the latter two RT datasets, we predicted RT from DHS <sup>37</sup> or ATAC-seq data <sup>38</sup>,  
393 respectively, using the Replicon software, which predicts RT profiles from local distributions in  
394 chromatin accessibility at very high accuracy <sup>39</sup> (see Methods).

395  
396 Next, we aimed to characterize the variability in RT across individuals that predicts RMDglobal1.  
397 By calculating the difference in window-wise RT for each pair of RT samples, and correlating  
398 this difference with RMDglobal1 window weights (Fig 4b, Fig S12), we observed that often only  
399 two RT profiles can be enough to predict RMDglobal1 using either expRT (max  $R=0.47$ ),  
400 predRT (max  $R=0.49$ ) and predRT-TCGA (max  $R=0.62$ ). In predRT, the best correlations are  
401 obtained when the difference in RT is when contrasting a pair that consist of one RT from  
402 (noncancerous) intact tissue *versus* one RT from primary cultured cells (Fig 4b), This suggests  
403 that selection for proliferation-capable stem-like cells when introducing cells into culture may  
404 alter RT, and that this altered RT is reflected in mutation rates in RMDglobal1 (see below for  
405 further discussion). As an illustrative example in a classification analysis using two selected RT  
406 profiles, one from a primary cell culture (ENCFF145RIZ) and one from an intact tissue  
407 (ENCFF315RKI), we observed that while the RT profile of each sample alone does not  
408 accurately identify RMDglobal1-high windows (Fig 4c), the difference in RT of windows between

409 these two RT samples can classify the windows with high RMDglobal1 weights from those  
410 windows with low RMDglobal1 weights (AUC = 0.82) (Fig 4d).

### 411 412 **Cell cycling gene expression-associated changes in RT are relevant for RMDglobal1**

413  
414 To further characterize the source of variability within RT profiles that explains RMDglobal1  
415 signature, we applied a PCA with the predRT-TCGA dataset of RT in TCGA tumors (Fig 4e),  
416 and correlated each RT-PC with the profile of RMDglobal1 signature across megabase  
417 windows. We observed that the strongest PCs, RT-PC3 and RT-PC4 are either tissue-  
418 associated, separating breast from kidney and brain tumors, or represent the average RT profile  
419 (RT-PC1, RT-PC2) (Fig S13). However, the RT-PC5 does not have a strong tissue bias but  
420 correlates strongly with RMDglobal1 ( $R=-0.49$ ) (Fig 4e, Fig S13c). Indeed, when we checked  
421 the RT profiles for the top RT-PC5 positive and RT-PC5 negative tumor samples, we observed  
422 RT differences in the RMDglobal1-relevant windows (Fig S14). The next best correlation was  
423 with RT-PC6 ( $R=0.35$ ).

424  
425 These RT-PCs summarize the global variation in the RT program between tumors of the same  
426 cancer type, and also predict RMDglobal1 global variation in mutation distribution. To interpret  
427 the RT-PCs, we asked how gene expression changed between the TCGA tumor samples with  
428 high values of a RT-PC versus tumors with low values. We considered the RHP gene sets,  
429 representing gene expression programs that are variable in a coordinated manner between  
430 individual cells, and that were recurrently observed across different cancer cell lines<sup>12</sup>.

431  
432 In particular, RT-PC5 correlates strongly with gene expression of cell cycle genes from the RHP  
433 gene sets<sup>12</sup> (there are two RHP sets of cell cycle genes, the G2/M and G1/S, and both correlate  
434 at  $p=9e-40$  and  $1e-14$ , respectively) (Fig 4f, Fig S13d), while the other RHP gene sets correlate  
435 less (next strongest is  $p\text{-value}=5e-05$ ) Consistently, also the RT-PC6, which has a more subtle  
436 correlation to RMDglobal1, also correlates with the cell-cycle RHP gene expression programs  
437 (Fig 4f, Fig S13d). This suggests that the RMDglobal1 regional mutability pattern reflects the RT  
438 program alterations associated with variable speed of cell cycling across different tumors.

439  
440 To further understand the biology of the systematic variation in RT captured by the RT-PCs  
441 relevant to mutation rates, we projected the predRT and expRT data into the existing predRT-  
442 TCGA PC coordinate system. RT-PC5 separated tissues *versus* cultured primary cells in  
443 predRT samples (Fig 4g). One possible interpretation is that RT-PC5 captures the effect of  
444 tissue culture conditions on RT profiles, however this is unlikely because there is a considerable  
445 spread within the cultured cells group, which span across the tissue-side of RT-PC5 on the one  
446 extreme of RT-PC5 and cell line side on the other extreme of RT-PC5 (Fig 4g). The other  
447 interpretation is that RT-PC5 captures the RT program of proliferation-capable, stem-like cells,  
448 which are normally a minority in an intact tissue, but are selected during establishment of cell  
449 culture; this is consistent with the above-mentioned cell cycling RHP gene expression program  
450 association with RT-PC5 and RT-PC6 and so we favor this interpretation. Next, the RT-PC5  
451 also separated healthy *versus* cancerous cells in the expRT samples (considered for blood  
452 cells, where both healthy and tumor was available (Fig 4h)). This suggests that this property



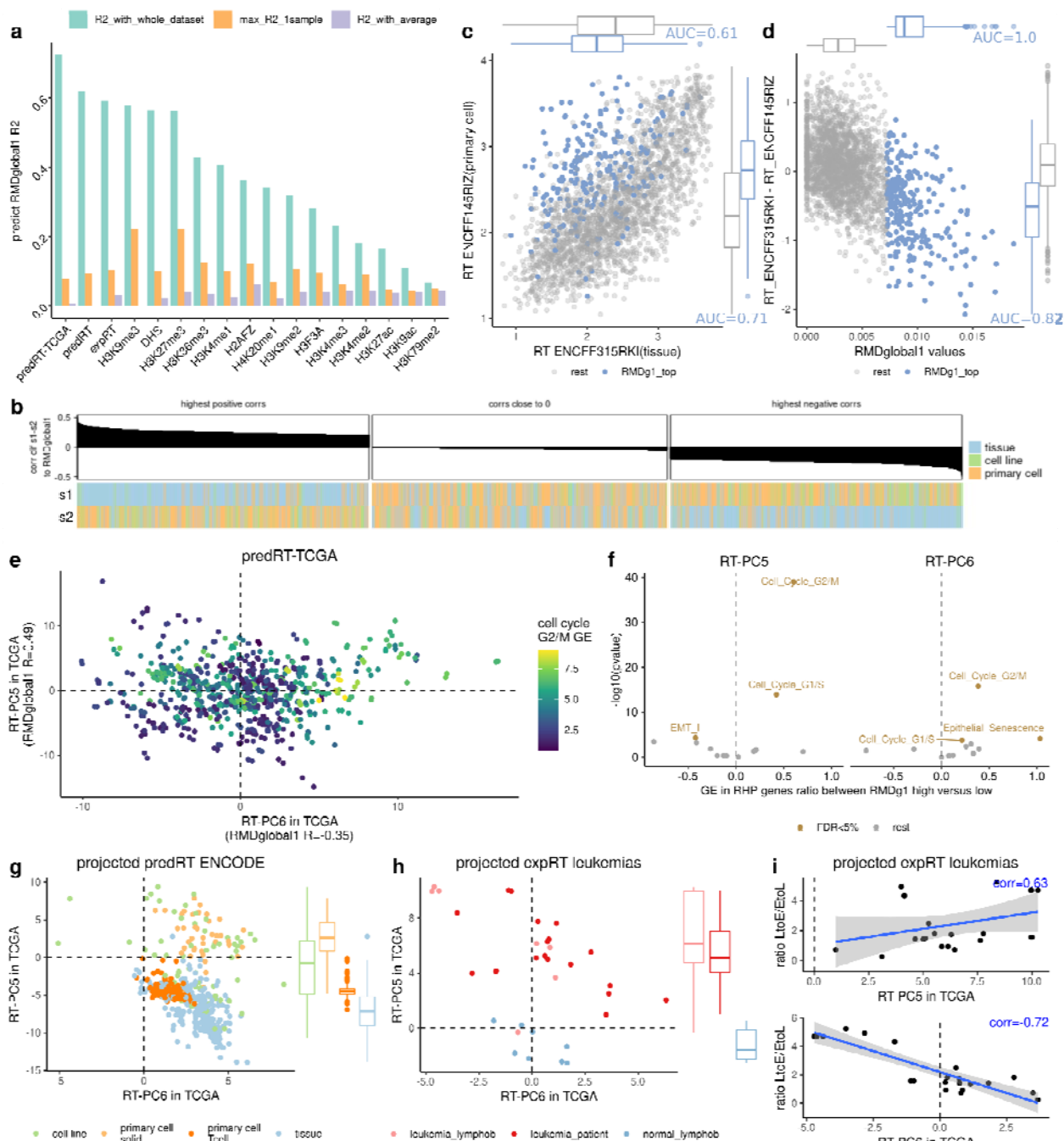
453 captured in RT-PC5 is more prominent in cancerous cells than in normal cells, again consistent  
454 with the property being related with cell cycling, which is often unchecked and accelerated in  
455 cancer cells.

456  
457 In summary, RT-PC5 separates intact tissue samples or tumors with lower cell cycle gene  
458 expression on one side, and cultured primary cells or tumors with higher cell cycle gene  
459 expression on the other side. This suggests that the windows with high RMDglobal1 weights are  
460 those that undergo changes in RT in more proliferative/stem cell-like samples compared to less  
461 proliferative/differentiated cell-like samples.

462  
463 Within a subset of the expRT data, changes in RT were studied previously<sup>40</sup>, reporting late-to-  
464 early (LtoE) and early-to-late (EtoL) RT changes between noncancerous samples  
465 (lymphoblastoid cell lines) and cancers (leukemias and cell lines). Interestingly, their pre-  
466 calculated ratio of LtoE/EtoL strongly correlate with our RT-PC6 ( $R=-0.72$ ) and RT-PC5  
467 ( $R=0.63$ ) (Fig 4i), adding evidence that RMDglobal1 is linked to the genome-wide changes in  
468 RT that occur during cancerous transformation.

469  
470 As a validation, we saw the same trends when we performed the PCA in predRT initially (i.e.  
471 using a mix of tissues and cell types, rather than only tumors in predRT-TCGA), and then  
472 projected the expRT into it (Fig S15). Of note, the expRT-PC that reflects developmental  
473 changes as reported earlier<sup>41</sup> does not correlate with RMDglobal1 (Fig S16), meaning that  
474 RMDglobal1 mutagenesis pattern does not relate to embryonal-characteristic patterns of RT.

475  
476  
477



478  
479

480 **Figure 4. RMDglobal1 signature is linked to regional variability in replication timing.** **a)** Adjusted  $R^2$   
481 of a regression predicting RMDglobal1 window weights from various epigenomic features (x axis) using  
482 either the whole dataset jointly, or selecting the maximum  $R^2$  of each sample in the dataset individually,  
483 or using the average values of the feature across the samples in the dataset. **b)** Correlation between  
484 RMDglobal1 signature, and the difference between each pair of RT profiles (all combinations tested).  
485 Panel shows 1st decile (highest positive R), 5th decile (R close to 0) and 10th decile (highest negative R)  
486 deciles ordered by correlation. **c)** RT profiles for two selected samples, where dots are megabase  
487 windows, colored by their weight in RMDglobal1 signature (top decile in blue). RT of each sample  
488 individually is modestly predictive of RMDglobal1 (AUCs for discriminating top-decile windows are listed

489 next to boxplots of RTs). **d)** Difference between the two RT profiles in panel c (on y axis) is predictive of  
490 the RMDglobal1 signature (see AUC for discriminating top-decile windows). **e)** A PCA on a matrix of  
491 predicted RT from TCGA tumor samples. RT-PC5 and RT-PC6 are shown because of their correlation  
492 with RMDglobal1 ( $R=0.49$  and  $-0.35$ , respectively). Points are colored by the mean gene expression in  
493 the cell cycle G2/M RHP module, in each TCGA tumor sample. **f)** Association between RT-PC5-high (top  
494 tertile) versus RT-PC5-low (bottom tertile) with the expression of genes in various RHP programs<sup>12</sup>, and  
495 same for RT-PC6. **g)** Predicted RT from ENCODE data with tissues, primary cells and cell lines (predRT-  
496 ENCODE) was projected into PCs of the tumor predRT-TCGA data. **h)** Projection of experimentally  
497 determined RT data for leukemias and normal blood cells into the same PCs of predRT-TCGA data. **i)**  
498 Correlation between the projection of expRT leukemia samples in RT-PC5 and RT-PC6, and the ratio of  
499 late-to-early and early-to-late regional RT changes reported previously<sup>40</sup>.

500

## 501 **RMDglobal1 signature associates with RB1 loss and affects regions that undergo** 502 **chromatin remodeling**

503

504 To identify events that may drive the changes in RT we found linked with cell cycling, we  
505 performed an analysis to detect somatic copy number alteration (CNA) events and deleterious  
506 point mutations are associated with RMDglobal1 exposure, while adjusting for cancer type and  
507 for confounding between linked CNAs (qq-plots in [Fig S17](#); [Methods](#) for details). Here, we  
508 considered 1543 chromatin modifier genes, cell cycle genes, DNA replication and repair genes  
509 and cancer genes, compared against a background of 1000 control genes ([Methods](#)).

510

511 For CNA, we found a strong positive association of RMDglobal1 with *RB1* deletion  
512 (FDR=0.05%, and better p-value than all control genes) ([Fig 5a-b](#), [Fig S18a](#)). Because CNA  
513 often affects large segments, we also checked associations with *RB1* neighbors ([Fig 5c](#)), noting  
514 that *RB1* is at the CNA peak (by mean estimated copy-number across tumor samples), meaning  
515 it is the likely causal gene. RMDglobal1 association with *RB1* is gene dosage dependent ([Fig](#)  
516 [S18b](#)). Consistently, we see that the effect of *RB1* mutations shows a trend in the same  
517 direction as *RB1* deletions, even though it is nonsignificant (*RB1* mutations are rarer) ([Fig](#)  
518 [S18c](#)). As independent supporting evidence, we identified deletions in *CDK6*, a negative  
519 regulator upstream of *RB1*, negatively associated with RMDglobal1, also having a stronger p-  
520 value than any of the control genes considered ([Fig 5a](#)).

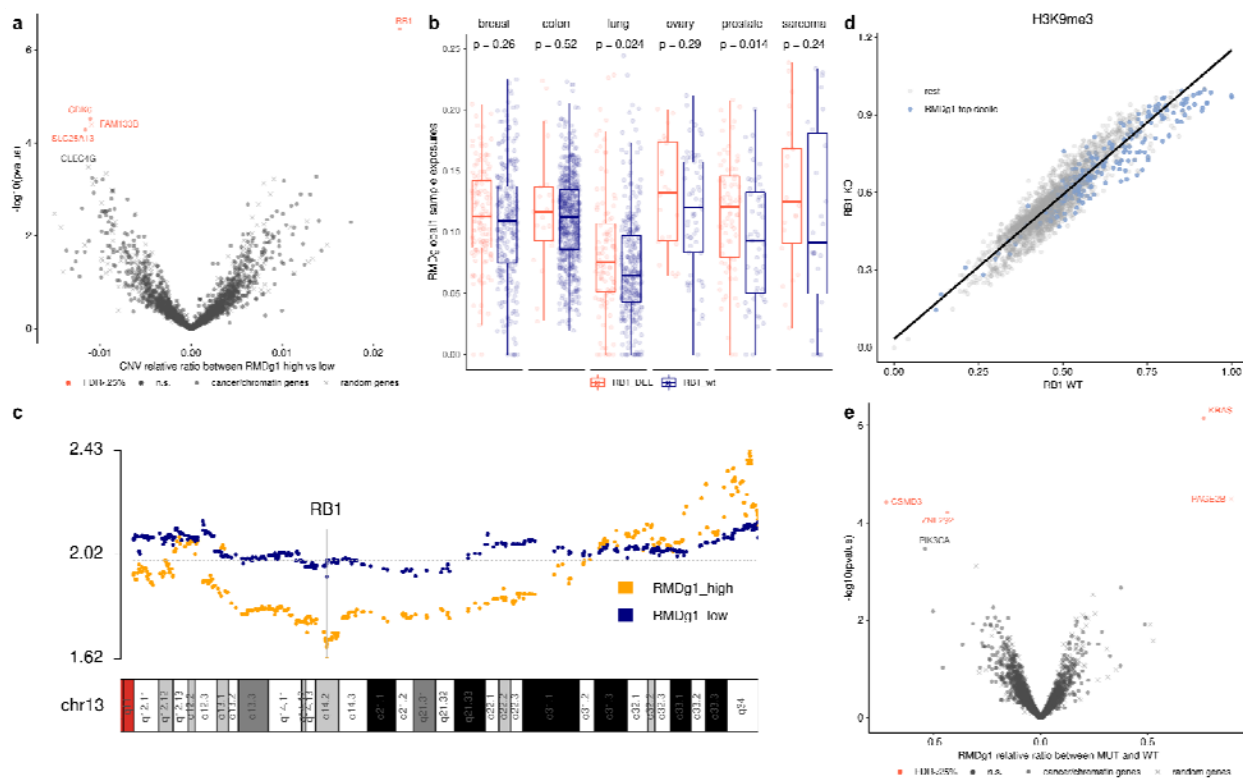
521

522 In addition to its effects on cell cycle regulation, *RB1* has additional important roles in chromatin  
523 organization<sup>19,42-44</sup>. In specific, *RB1* deletions were reported to change heterochromatin marks  
524 H3K9me3 and H3K27me3 in regions enriched at subtelomeres, and this associates with  
525 regional propensity to DNA damage<sup>19</sup>. These same two marks we found to be more highly  
526 correlated to RMDglobal1 than other tested marks ([Fig 4a](#)), and interestingly we find that  
527 RMDglobal1 window weights are also strongly enriched approximately 5 Mb nearby telomeres  
528 ([Fig 6h](#)). Notably, the changes in regional H3K9me3 profile when *RB1* is wild-type *versus* in  
529 isogenic *RB1* k.o. cells<sup>19</sup> predicted RMDglobal1 signature (adjusted  $R^2=0.29$ ), and so did  
530 changes in regional H3K27me3 albeit subtly (adj $R^2=0.18$ ) ([Fig 5d](#), [Fig S19ab](#)). The genome  
531 regions with 10% highest weights in RMDglobal1 are the ones where the level of H3K9me3  
532 heterochromatin mark is more likely to be asymmetrically altered upon *RB1* disruption<sup>19</sup> (off-  
533 diagonal dots in [Fig 5d](#)). This indicates that loci where heterochromatin is remodeled upon *RB1*

534 loss-of-function<sup>19</sup> significantly overlap with loci where RMDglobal1 mutation rates change in  
 535 many tumors, further implicating RB1 in shaping the mutation rate landscape.

536  
 537 As for CNA, we also tested associations between the presence of deleterious point mutations in  
 538 cancer and chromatin and DNA repair genes and the exposure to the RMDglobal1 mutagenic  
 539 pattern. Here, we found the *KRAS* mutation to strongly positively associate with RMDglobal1, at  
 540 FDR=0.1% (Fig 5e, Fig S19c), and this is observed consistently across individual cancer types  
 541 (Fig S19c) and significantly in colon, uterus and bladder (see Fig S19d-e legend for comment on  
 542 lung adenocarcinoma). Of note, the *KRAS* gene was reported to act downstream of *RB1* loss-  
 543 of-function with *RB1* in developmental and in tumor mouse phenotypes<sup>45,46</sup>. Consistently, *KRAS*  
 544 mutation and *RB1* loss (either deletion or mutation) are mutually exclusive in our tumor dataset  
 545 (chi-square  $p < 2.2e-16$ ), supporting that the driver alterations in *RB1* and *KRAS* may converge  
 546 onto the same mutation rate redistribution phenotype, RMDglobal1.

547



548  
 549 **Figure 5. Genetic alterations associated with the activity RMDglobal1 signature.** a) Associations  
 550 between CNA deletions and tumors with higher RMDglobal1 exposures in a pan-cancer analysis,  
 551 adjusting for cancer type and for global CNA patterns (Methods). N=1543 cancer genes and chromatin-  
 552 related genes are shown (dots), as well as a 1000 set of randomly chosen genes (crosses). b)  
 553 Differences in RMDglobal1 exposures between RB1 deletion (-1 or -2 deletion) and wt for several cancer  
 554 types (those with the highest number of samples with *RB1* deletion); remainder in Fig S18. c) Mean local  
 555 CN profile in groups of tumors, grouped by RMDglobal1 high and low, of the segment of chromosome 13  
 556 containing the gene *RB1*. Each dot represents one gene. d) Correlation between the H3K9me3  
 557 heterochromatin profiles for samples with *RB1* knock-out (“KO”) versus wild-type (“WT”). Each dot  
 558 represents a window, colored by RMDglobal1 window weight top decile versus the rest of the windows. e)

559 Associations between deleterious SNV and indel mutations in the same sets of genes as in panel **a**, and  
560 the RMDglobal1-high versus RMDglobal1-low activity of tumor samples, in a pan-cancer analysis.

561  
562 Motivated by these associations between RB1-caused regional heterochromatin changes<sup>19</sup> and  
563 the RMDglobal1 regional mutation rates, we further investigated the variation in the H3K27me3  
564 and H3K9me3 marks across datasets in ENCODE (Fig 6a). To characterize the regional  
565 heterochromatin variability, we performed a PCA on the profiles of the two marks and the  
566 predRT together. The resulting heterochromatin-PC4 (het-PC4) correlated best with RMDglobal1  
567 window weights ( $R=0.53$ ). As above, the difference in the three features (H3K9me, H3K27me3,  
568 RT) separated the proliferative, putatively stem-like samples (het-PC4 positive) *versus* the rest  
569 of the samples (het-PC4 negative) (Fig 6b). The stem-like samples (het-PC4 positive) are later  
570 replicating (relative increase in RT [het-PC4 high vs low] = 61%) and have higher H3K27me3  
571 and H3K9me3 in RMDglobal1 top windows (relative increase of 55% and 78% respectively) (Fig  
572 6c). In summary, our analyses suggest that the chromosomal domains with highest RMDglobal1  
573 weights become later-replicating and more heterochromatic in more stem-like cells (cell  
574 lines/primary cells), associated with an increase of relative mutation rates in these domains.

575  
576  
577 **Gene regulation and chromatin compartments associated with the RMDglobal1 mutation**  
578 **rate phenotype**

579  
580 The regional variability in RT and heterochromatin marks, as reflected in variable somatic RMD  
581 in tumors, suggests the existence of concomitant changes of regional gene expression,  
582 because early RT was reported to be broadly associated with higher gene expression<sup>1</sup>.  
583 Therefore we asked if there are coordinated changes in gene expression levels in certain  
584 windows between the RMDg1-high and RMDg1-low cancers. Indeed, we found several windows  
585 with gene expression upregulation and downregulation between RMDg1-high and low cancers  
586 ( $FDR < 25\%$ ). The windows with coordinated gene expression downregulation are enriched in  
587 higher values of RMDglobal1 window weights, compared to the windows with non-coordinated  
588 gene expression changes (Wilcoxon test, greater; downregulation p-value = 0.03; there is a  
589 nonsignificant trend for coordinated upregulation) (Fig S20a). These regional changes in gene  
590 expression are consistent with chromatin remodeling affecting various chromosomal domains,  
591 which is also mirrored in regional mutation rates (Fig S20b).

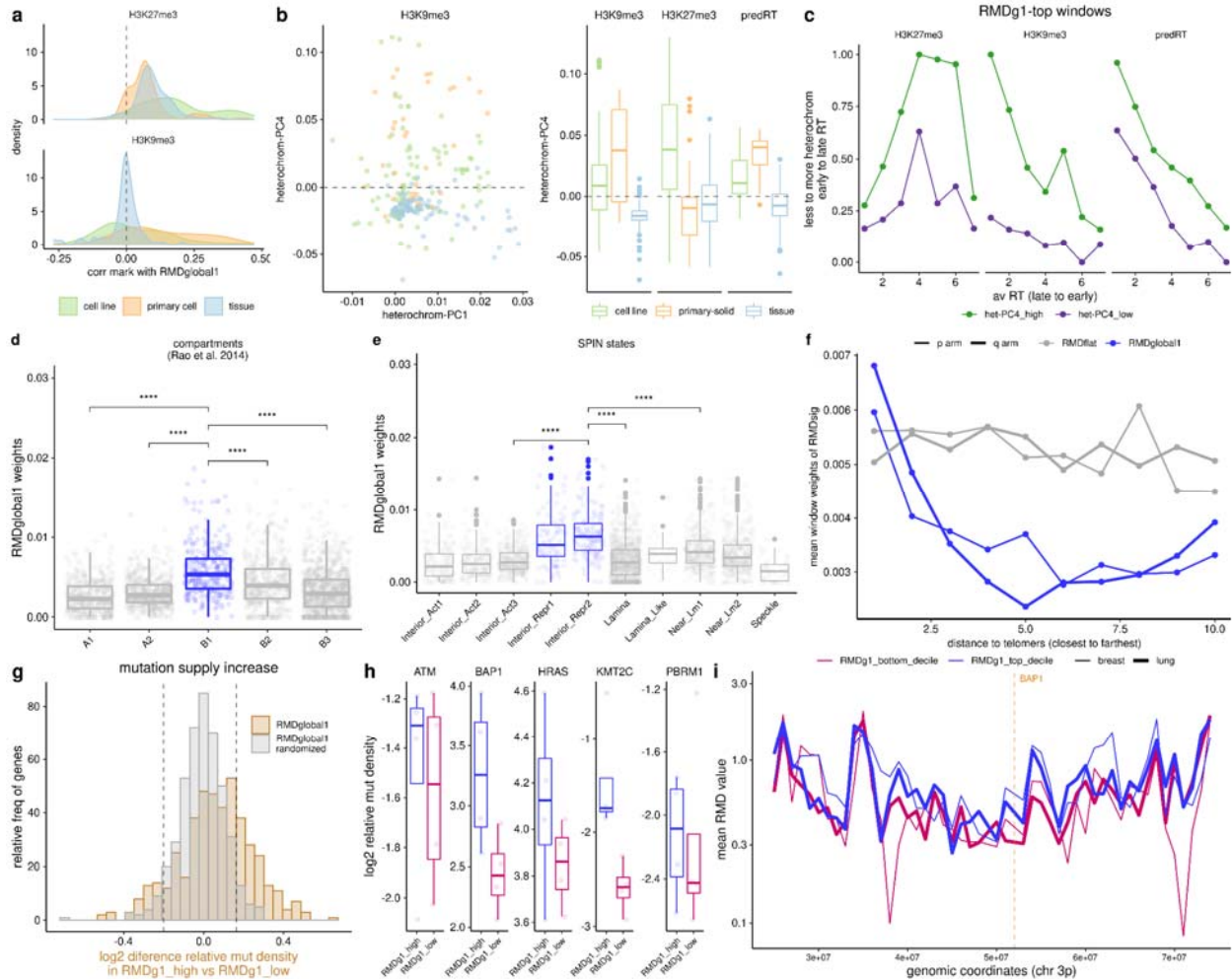
592  
593 To additionally characterize the regions affected by the chromatin remodelling, we analyzed  
594 data from diverse types of genomic assays from various studies (Table S3) that reported some  
595 correlations with RT. We compared the regional density of these various features with our  
596 RMDglobal1 window weights (Table S3). Correlations were noted with Hi-C subcompartments  
597 (Fig 6d), inferred from long-range chromatin interactions at fine resolution (25 kb)<sup>47</sup>. In  
598 particular, the B1 subcompartment was associated with RMDglobal1; this subcompartment  
599 replicates during middle S phase, and correlates positively with the Polycomb H3K27me3 mark  
600 (Fig S21) and negatively with H3K36me3 suggesting that it represents facultative  
601 heterochromatin<sup>47</sup>. Next, we observed a correlation with two SPIN states (Fig 6e), derived by  
602 integrating nuclear compartment mapping assays and chromatin interaction data<sup>48</sup>.



603 RMDglobal1 signature regions are enriched in the two “Interior repressed” SPIN states, marking  
604 regions that are inactive, however unlike other inactive heterochromatic regions they are located  
605 centrally in the nucleus, rather than peripherally (next to the lamina). Additionally, RMDglobal1  
606 important windows are enriched in subtelomeric regions (Fig 6f). In sum, the windows with  
607 higher weights in RMDglobal1 signature are enriched in subtelomeric regions, the B1 facultative  
608 heterochromatin subcompartment, and nuclear interior located, repressed chromatin states.

609  
610 Since RMDglobal1 captures a redistribution of mutation rates genome-wide, we predicted that  
611 this will affect the supply of mutations to some cancer genes. To quantify this, we performed a  
612 similar analysis as for the MSI-associated RMDflat above; for RMDglobal1 shown in (Fig 6g-i).  
613 When compared to a randomized baseline (95th percentile of the random distribution used as  
614 cutoff), 28% of the 460 cancer genes suffer a significant increase of mutation supply when  
615 comparing RMDglobal1-low (bottom tertile) and -high (top tertile) tumor samples. Regarding the  
616 effect size of increase, these genes increase mutation rates on average by 1.21-fold between  
617 the RMDglobal1-low *versus* high tertile tumors. The mutation rate density is shown for 5  
618 example genes with high fold-difference in Fig 6h, where for instance the median mutation rate  
619 for the *ATM* tumor suppressor increases by 1.18-fold, and for the *KMT2C* tumor suppressor by  
620 1.79-fold, in the top tertile by RMDglobal1 signature of mutation redistribution.

621



622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642

**Figure 6. RMDglobal1 mutation rate redistribution is linked with chromatin remodeling.** **a)** Heterochromatin marks correlation with RMDglobal1. Distribution of the correlations between H3K9me3 and H3K27me3 profiles in various ENCODE datasets of healthy tissues/primary cells/cell lines, and the RMDglobal1 signature extracted from tumor mutations. This shows a wide spread of correlations, with some examples of cell lines or primary cells with high correlations of heterochromatin profiles with RMDglobal1. **b)** A PCA was performed on the predicted RT and the heterochromatin marks from ENCODE data. Left panel shows the heterochromatin PC4 (het-PC4) selected for its high correlation with RMDglobal1 ( $R=0.52$ ) and het-PC1 shown for highest amount of variance explained for H3K9me3. Right panel shows the het-PC4 distribution across different cell types for the 3 features. **c)** Mean predicted RT, H3K27me3 and H3K9me3 across PC4-high versus PC4-low groups in ENCODE data, split by RT bins. **d)** RMDglobal1 signature across different Hi-C nuclear subcompartments from reference<sup>47</sup>. **e)** RMDglobal1 signature across different SPIN nuclear compartmentalization states from<sup>48</sup>. **f)** RMDglobal1 and RMDflat signature window weights compared to distance to telomeres. **g)** Distribution for the difference in mutation density (see Fig. 2d), shown for 460 cancer genes, comparing between RMDglobal1-high and low tumors, using the actual values of RMDglobal1 and as a baseline randomized of RMDglobal1. Vertical lines show 5th and 95th percentile of the randomized distribution. **h)** Mutation density for RMDglobal1-high versus low tumor samples (here, top tertile versus bottom tertile) for 5 example genes (drivers in  $\geq 4$  cancer types and with the highest effect size); dots are cancer types. **i)** Mean RMD profile on chromosome 3p across the RMDglobal1-high versus low tumor groups (here, top and bottom decile by

643 RMDglobal1), for two example cancer types. Vertical lines mark the position for the *BAP1* tumor  
644 suppressor gene (example gene in panel h).

645

646

### 647 **A TP53-associated RMDglobal2 signature reduces relative mutation rates in late** 648 **replicating regions**

649

650 RMDglobal2 signature mutations follow a distribution similar to the canonical RMD landscape,  
651 increasing mutation density in late replication, except for a set of very late RT windows , which  
652 acquire fewer mutations than expected from RT (Fig 7ab). Mutation density increases near  
653 linearly with RT bins in tumors with high RMDglobal2, while in tumors with a low RMDglobal2  
654 exposure the RT relationship to mutation rates is better described by a quadratic fit (Fig 7c, Fig  
655 S22). Therefore, qualitatively the canonical RT-associated RMD landscape is preserved  
656 regardless of RMDglobal2 being low or high. However RMDglobal2 changes the shape of  
657 association to RT, by exaggerating (or suppressing) the more prominent peaks in regional  
658 mutation rates, but not affecting the minor peaks.

659

660 We aimed to identify driver event behind this redistribution of mutations by testing for  
661 associations of RMDglobal2 high (top tertile) versus low (bottom tertile) samples, and genetic  
662 events (CNAs, deleterious mutations) in cancer driver, DNA repair and chromatin modifier  
663 genes. Strikingly, we found *TP53* mutation to be uniquely strongly associated with RMDglobal2  
664 signature (effect size = 1.27, FDR = 9e-10) (Fig 7d). As supporting evidence, we found that  
665 *TP53* deletions also positively associated (Fig 7e) and, independently, the known amplifications  
666 that phenocopy *TP53* loss (*MDM2*, *MDM4* and *PPM1D* oncogenes) are also positively  
667 associated with RMDglobal2 RMD signature exposures (Fig 7e, Fig S23). This rules out that the  
668 *TP53* driver mutation is merely the consequence of RMDglobal2 redistribution, and provides  
669 evidence for a causal effect of *TP53* inactivation.

670

671 Since *TP53* mutations were reported to be associated with increased burdens of CNA events<sup>49</sup>,  
672 we tested whether RMDglobal2 RMD signature could be due to confounding from a multiplicity  
673 of focal CNA events (we note our method for RMD analysis does control for confounding by  
674 arm-level CNAs, Methods), which can modify apparent local mutation rates. However, there is  
675 only a weak correlation between the CNA burden and RMDglobal2 signature levels upon  
676 stratifying for *TP53* status ( $R \leq 0.11$ ), suggesting that RMDglobal2 likely does not simply reflect  
677 changes in local DNA copy number (Fig S24).

678

679 RMDglobal2 signature describes variation in certain genome regions, which may affect mutation  
680 supply to genes therein. We tested whether there is a difference in mutation rate in the cancer  
681 genes for RMDglobal2-high (top tertile) versus low (bottom tertile) tumor samples (Fig 7f). When  
682 compared to randomized data (5th percentile), 26% of cancer genes exhibited decreased  
683 mutation supply; only 6% genes exhibit an increased mutation supply with high RMDglobal2  
684 (Fig 7f-g). As an example, we show the mutation density of *ARID1A* and *GATA3*, which  
685 decreased in mutation supply (as above, measured using intronic rates; the decrease implies  
686 they are below the 5th percentile of the randomized distribution) with high RMDglobal2 (Fig 7h).

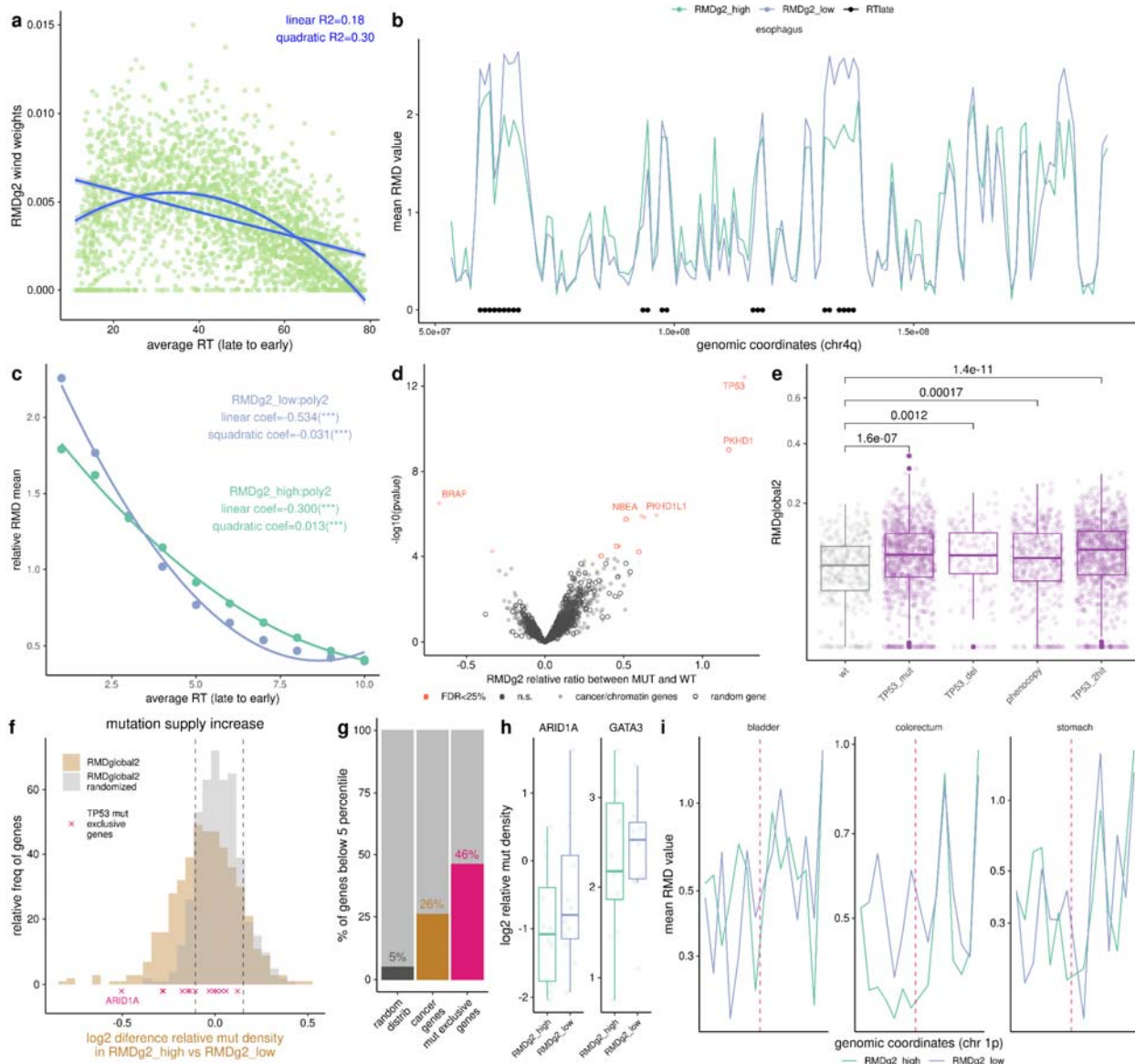
687 We hypothesized that apparent genetic interactions like this mutual exclusivity example might  
688 arise due to redistribution of mutations altering local mutation supply to genes. We thus  
689 considered 13 genes mutually exclusive with *TP53* mutations<sup>50</sup> (note that *TP53* loss is strongly  
690 associated with RMDglobal2) and found that nearly half (6/13) of these genes were below the  
691 5th percentile of the random distribution (Fig 7f-g). Upon inspection of the raw RMD profiles for  
692 RMDglobal2 high and low tumors for several cancer types we noted a difference in the region  
693 where *ARID1A* resides (Fig 7i). Overall, this illustrates how a global redistribution of mutation  
694 rates, here mediated by *TP53* loss, can create apparent genetic interactions that may not  
695 indicate selection on functional effects of the genetic interaction. Thus, regional mutation rates,  
696 which vary extensively between tumors, should be explicitly controlled for in statistical studies of  
697 epistasis in cancer genomes.

698

699

700

701



702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717

**Figure 7. A TP53-associated mechanism underlies the RMDglobal2 mutation rate redistribution pattern.** **a)** A quadratic association of RMDglobal2 signature with the average replication timing. **b)** Mean RMD profiles in chromosome 4q for the RMDglobal2-high versus low tumor samples in esophagus cancer. Latest RT windows (avRT<20) marked with black dots. **c)** Relative RMD mean profile across 10 RT bins for tumors that are RMDglobal2-high (RMDglobal2 exposures > 0.17) versus RMDglobal2-low (RMDglobal2 exposures < 0.01), showing a linearization of the link between RT and mutation rates in RMDglobal2-high. **d)** Associations between deleterious mutations in known cancer genes and chromatin-related genes (dots) and a control set of randomly chosen genes (hollow circles), and RMDglobal2 exposures in samples (p-values from Z-test on regression coefficient). **e)** RMDglobal2 signature exposures of tumor samples stratified by: wild-type for *TP53* (wt), *TP53* with 1 mutation (*TP53\_mut*), *TP53* with 1 deletion (*TP53\_del*), *TP53* loss phenocopy via a amplification in *MDM2*, *MDM4* or *PPM1D* (*TP53\_pheno*), or *TP53* with any two hits of the previously mentioned alteration (*TP53\_2hit*). **f)** Distribution of the log<sub>2</sub> difference in the relative mutation density (intronic) for 460 cancer genes, comparing between RMDglobal2 high tumors and RMDglobal2 low tumors, using the actual values



718 (“RMDglobal2” histogram) and randomized values (“RMDglobal2 randomized” histogram). Position of the  
719 genes mutually exclusive with TP53 marked with crosses. **g)** Percentage of genes above 95 percentile of  
720 a random distribution for the random distribution, cancer genes and TP53 mutually exclusive genes. **h)**  
721 Log2 relative mutation density (normalized to flanking DNA in same chromosome arm, see [Fig 3d](#)) for  
722 RMDglobal2-high versus RMDglobal2-low for 2 example genes (TP53 mutually exclusive genes above  
723 the 95 percentile). Each dot is a cancer type. **i)** Mean RMD profile across the RMDglobal2-high versus low  
724 groups in a region of chr 1p. Vertical lines mark the position for the *ARID1A* gene.

725

726

## 727 Discussion

728 Even though there clearly exists a common, canonical regional mutation rate landscape shared  
729 across human cells, there are RMD patterns superimposed that are tissue-specific. This is  
730 consistent with the fact that tissues have different RT programs and chromatin landscapes.  
731 Here we systematically characterized patterns of regional redistribution of somatic mutations  
732 independent of tissue, identifying two new global RMD patterns and possible underlying  
733 mechanisms.

734

735 We demonstrate how an NMF-based approach can deconvolute the RMD mutational patterns  
736 that compose the final regional mutagenesis ‘portrait’ of each individual tumor. The method is  
737 related for those applied to trinucleotide mutational signatures, however it also rigorously  
738 adjusts for the confounding by these trinucleotide signatures. Of 13 RMD signatures, expectedly  
739 the majority were tissue-associated.

740

741 Some of the tissue-related RMD signatures may bridge various cancer types, usually reflecting  
742 known biology, however the surprising RMD signature (“B.O.P.S.”) has high activity in many  
743 brain samples but also in ovary, uterus, prostate and sarcoma cancers. Since these 4 cancer  
744 types do not have their own specific tissue-RMD signature this can indicate that RMD “B.O.P.S.”  
745 is a residual RMD signature that collects diverse RMD patterns that current NMF methodology  
746 does not resolve well, possibly due to lack of power. However, this RMD B.O.P.S. pattern was  
747 similarly robust (by autocorrelation across windows) as the other tissue-specific patterns ([Fig](#)  
748 [S6](#)) and so this RMD may reflect some commonalities in chromatin organization and gene  
749 regulation connecting those cancer types. For instance, an analysis of transcriptome-based cell  
750 states across tumor types, based on single-cell gene expression data <sup>13</sup> suggested a module of  
751 cilium/cytoskeleton-related genes common to some ovarian cancers, glioblastoma, uterine  
752 cancer and lung adenocarcinoma, thus the tissue spectrum corresponds to B.O.P.S. (we do  
753 note the lung tissue is separate in our RMD analysis). We acknowledge that, as has recently  
754 occurred with the trinucleotide mutational signatures <sup>51</sup>, some of the initially proposed signatures  
755 such as RMD\_B.O.P.S. may, with arrival of more data, be able to be ‘split’ into component RMD  
756 signatures that more precisely match tissue identity. Overall, the RMD features provide an  
757 important tool for understanding the relatedness of cancer types and the chromatin organization  
758 in the cell-of-origin of cancers <sup>11,52</sup>.

759

760 Here, we identified 3 robust ‘global’ (i.e. largely independent of tissue) RMD signatures of  
761 redistribution of mutation rates occurring in human cancer. As expected, we recovered the

762 redistribution of mutations towards a flatter mutational landscape (“RMDflat”) that was  
763 previously described for DNA MMR failures in tumors, cell lines and xenografts<sup>1,53,54</sup>, validating  
764 our methodology. A simulation study further supports the broad adequacy of the NMF method,  
765 with the number of mutations generated by a RMD signature being a limiting element in  
766 identifying the signature. With a higher number of tumor samples, additional RMD patterns may  
767 become sufficiently represented to be recognized by NMF, and our RMD catalog will likely be  
768 extended with rarer RMD patterns, as was the case for trinucleotide mutational signatures<sup>27,51</sup>.  
769 A limitation of the current implementation of our method is that low mutation burden tumors are  
770 not analyzed; increasing WGS sequencing depth and so power to detect subclonal variants may  
771 alleviate this constraint.

772  
773 Of the widespread global RMD signatures, RMDglobal1 causes a genome-wide redistribution of  
774 mutations predominantly in inactive regions enriched in Polycomb facultative heterochromatin  
775 mark (H3K27me3), and centrally located in the nucleus (i.e. not lamina associated  
776 heterochromatin). These regions showed variable RT programs and variable heterochromatin  
777 state, comparing a more proliferative/stem-like group of samples versus a less stem-like group  
778 of samples. Consistently, the corresponding RMDglobal1 mutational pattern was associated  
779 with genetic alterations implicated in cell cycle disturbances. Of note, *RB1* is involved in cell  
780 cycle but also has important roles in chromatin organization<sup>19</sup> and so may affect this RMD  
781 pattern in multiple ways. Together, our analyses converge onto a model where due to more  
782 rapid cycling in tumor cells (e.g. caused by oncogenic drivers) and/or loss of cell cycle control  
783 (e.g. caused by *RB1* loss-of-function), chromatin is remodeled in facultative heterochromatin  
784 locations and RT program changed, and as a consequence the mutation rates in those regions  
785 are altered.

786  
787 One question that arises is why those specific regions undergo the chromatin remodeling and  
788 RT change. We found these RMDglobal1 regions to be enriched in the B1 Hi-C  
789 subcompartment, which is the most dynamic (less conserved) subcompartment across cell lines  
790<sup>55</sup>. Additionally, chromatin remodelling and increase of risk to DNA damage was reported to  
791 affect subtelomeric regions upon *RB1* disruption<sup>19</sup>. Collectively, this suggests that chromatin  
792 state in those facultative heterochromatin regions are likely more malleable and prone to  
793 change upon different processes, either developmental or cancerous in nature.

794  
795 The second global change in the mutational patterns we identified (RMDglobal2) occurs  
796 independently of the above and can be described as a sharp relative reduction of mutations in  
797 latest RT regions, associated with loss of TP53 activity via mutation, CNA or phenocopying  
798 events. Since *TP53* mutations are very common in tumors, the impact of this redistribution of  
799 mutations to many other cancer genes may be widespread. This study provides examples of  
800 how an alteration in one gene -- here, deletions in *RB1* or mutations in *TP53* -- can affect future  
801 evolutionary scenarios: by ‘redirecting’ the regional mutation supply away from one set of genes  
802 and towards another.

803  
804 In conclusion, our large-scale analysis recovered the known differences in mutation density  
805 between tissues and identified three robust global RMD signatures of mutation rate variability

806 across chromosomal domains. The global redistribution of mutations can have an important  
807 impact in mutation supply on cancer genes at their affected regions, increasing their likelihood  
808 to acquire a deleterious mutation.

809  
810

## 811 **Methods**

### 812 **WGS mutation data collection and processing**

813 We collected whole genome sequencing (WGS) somatic mutations from 6 different cohorts  
814 ([Table S1](#)). First, we downloaded 1950 WGS somatic single-nucleotide variants (SNVs) from  
815 the Pan-cancer Analysis of Whole Genomes (PCAWG) study at the International Cancer  
816 Genome Consortium <sup>56</sup> Data portal (<https://dcc.icgc.org/pkawg>). Second, we obtained 4823  
817 WGS somatic SNVs from the Hartwig Medical Foundation (HMF) project <sup>57</sup>  
818 (<https://www.hartwigmedicalfoundation.nl/en/>). Third, we downloaded 570 WGS somatic SNVs  
819 from the Personal Oncogenomics (POG) project <sup>58</sup> from BC Cancer  
820 (<https://www.bcgsc.ca/downloads/POG570/>). Fourth, we obtained 724 WGS somatic SNVs from  
821 The Cancer Genome Atlas (TCGA) study as in <sup>9</sup>; we used QSS\_NT $\geq$ 12 mutation calling  
822 threshold in this study.

823 Finally, we downloaded bam files for 781 WGS samples from the Clinical Proteomic Tumor  
824 Analysis Consortium (CPTAC) project <sup>59,60</sup> and bam files for 758 tumor samples from the MMRF  
825 COMMPASS project <sup>61</sup> from the GDC data portal (<https://portal.gdc.cancer.gov/>). Somatic  
826 variants were called using Illumina's Strelka2 caller <sup>62</sup>, using the variant calling threshold  
827 SomaticEVS  $\geq$ 6. Additionally, for these samples we performed a liftOver from GRCh38 to the  
828 hg19 reference genome.

829 We collected the samples' metadata (MSI status, purity, ploidy, smoking history, gender) from  
830 data portals and/or from the supplementary data of the corresponding publications. Additionally,  
831 we harmonized the cancer type labels across cohorts. Here, since lung tumors in HMF data are  
832 not divided into lung squamous cell carcinoma (LUSC) and lung adenocarcinoma (LUAD) types,  
833 we used a CNA-based classifier to tentatively annotate them in the HMF data. We downloaded  
834 copy number alteration data from HMF and TCGA for lung tumor samples and adjusted for  
835 batch effects between cohorts using ComBat as described in our previous work <sup>63</sup>. We trained a  
836 Ridge regression model with TCGA data to discriminate between LUSC and LUAD and applied  
837 the model to predict LUSC/LUAD in the HMF lung samples. We did not assign a label to  
838 samples with an ambiguous prediction score between 0.4 and 0.6.

839 Similarly, since POG breast cancer (BRCA) samples are not divided into subtypes (luminal A,  
840 luminal B, HER2+ and triple-negative) we used a gene expression classifier to annotate them.  
841 We downloaded gene expression data for TCGA and POG breast tumors and adjusted the data  
842 for batch effect using ComBat as previously described <sup>63</sup>. We trained a Ridge regression model  
843 with TCGA data to discriminate between the breast cancer subtypes (one-versus-rest) and  
844 applied the model to the POG breast samples to assign them to a subtype. We did not assign  
845 23 samples that are predicted as two subtypes and 8 that are not predicted as any subtype.

## 846 **Defining windows and filtered regions**

847 We divided the hg19 assembly of the human genome into 1 Mb-sized windows. These divisions  
848 are performed on each chromosome arm separately. To minimize errors due to misalignment of  
849 short reads, we masked out all regions in the genome defined in the 'CRG Alignability 75' track  
850 <sup>64</sup> with alignability <1.0. In addition, we removed the regions that are unstable when converting  
851 between GRCh37 and GRCh38 <sup>65</sup> and the ENCODE blacklist of problematic regions of the  
852 genome <sup>66</sup>.

853 Additionally, to minimize the effect of known sources of mutation rates variability at the sub-  
854 gene scale we removed CTCF binding site regions (downloaded from the Table Browser), ETS  
855 binding regions (downloaded from <http://funseq2.gersteinlab.org/data/2.1.0>) and APOBEC  
856 mutagenized hairpins downloaded from <sup>67</sup>. Finally, we removed all coding exon regions (+-2nts,  
857 downloaded from the Table Browser) to minimize the effect of selection on mutation rates.

## 858 **Matching trinucleotide composition across megabase windows**

859 To minimize the variability in mutational spectra confounding the analyses, we accounted for  
860 the trinucleotide composition of each window. For this, we removed trinucleotide positions from  
861 the genome in an iterative manner to reduce the difference in trinucleotide composition across  
862 windows. We selected 800,000 iterations that reach a tolerance >0.0005 (difference in relative  
863 frequency of trinucleotides between the windows). After the matching, we removed all windows  
864 that end up with less than 500,000 usable bps. The final number of analyzed windows is 2,540.

## 865 **Calculating the Regional Mutation Density (RMD) of each window**

866 For our WGS tumor sample set (n=9,606 WGS) we counted the number of mutations in the  
867 above-defined windows. We required a minimum number of mutations per sample of 5,876,  
868 which corresponds to 3 muts/Mb (total genome = 1,958,707,652 bp). In total, 4221 tumor  
869 samples remain, which we use for the downstream analyses.

870 To calculate the RMD, we normalized the counts of each window by: (i) the nt-at-risk available  
871 for analysis in each window and (ii) the sum of mutation densities in each chromosome arm. To  
872 control for whole arm copy number alterations.

873 To calculate the RMD applied to NMF analysis, we first subsample mutations from the few  
874 hyper-mutator tumors, to prevent undue influence on overall analysis. We allow a maximum of  
875 20 muts/Mb that is 39,174 muts. If the tumor mutation burden is higher we subsample the  
876 mutations to reduce it to that maximum value. Then, as above, we normalized the RMD by: (i)  
877 the nt-at-risk in each window [  $RMD = counts * average\_nt\_risk / nt\_at\_risk$  ] and (ii) the sum of  
878 mutation density in each chromosome arm [  $RMD * row\_mean\_WG / rowMeans$  by chr arm ].  
879 We multiply by the average nucleotides at risk and the mean whole genome to maintain the  
880 values range of each sample for the bootstrapping.

## 881 **Applying NMF to extract RMD signatures**

882 We applied bootstrap resampling (R function `UPmultinomial` from package `sampling`) to the  
883 RMD scores that we calculated for NMF as above. The result for each tumor sample is a vector  
884 of counts with a tumor mutation burden close to the original one but normalized by the  
885 nucleotides at risk by window and for the possible chromosome arm copy number alterations  
886 (CNA). Then, we applied NMF (R function: `nmf`) to the bootstrapped RMD matrices, testing  
887 different values of the rank parameter (1 to 20), herein referred to as `nFact`.

888  
889 We repeated the bootstrapping and NMF 100 times for each `nFact`. We pooled all the results by  
890 `nFact` and performed a k-medoids clustering (R function `pam`), with different number-of-clusters  
891 `k` values (1 to 20). We calculated the silhouette index value, a clustering quality score (which  
892 here measures, effectively, how reproducible are the NMF solutions across runs), for each  
893 clustering to select the best `nFact` and `k` values.

894  
895 Additionally, we also applied the same NMF methodology to each cancer type separately (`n` =  
896 12 cancer types that had >100 samples available).

897

### 898 **Simulated data with ground-truth RMD signatures**

899 For each cancer type, we calculated a vector of RMD values (i.e. regional mutation density  
900 mean of all samples from that cancer type) based on observed data, and super-imposed  
901 simulated ground-truth signatures onto these cancer type-derived canonical RMD patterns. We  
902 generated 9 simulated ground-truth RMD signatures with different characteristics, varying the  
903 number of windows affected by the signature (10, 20 or 50% of 2540 windows total) and the  
904 fold-enrichment of mutations in those windows (x2, x3 or x5) over the RMD window value in the  
905 canonical RMD pattern for that tissue.

906

907 In particular, we tested 9 different scenarios, varying the signature contribution to the total  
908 mutation burden (10, 20 or 40%) and the number of tumor samples affected by the signature (5,  
909 10 or 20%). We randomly assigned the ground-truth signatures to be super-imposed onto each  
910 tumor sample (e.g. sample A will be affected by RMD signature 1 and 3 while sample B will be  
911 affected by signature 4). In total, we have simulated genomes for 9 different scenarios (different  
912 RMD signature contributions and number of tumor samples affected), each of them containing  
913 the 9 simulated ground-truth RMD signatures.

914 We applied the NMF methodology for the 9 different scenarios independently and obtained NMF  
915 signatures. For each case, we selected an NMF `nFact` and k-medoids clustering `k`, based on the  
916 minimum cluster silhouette index (SI) quality score. To assess the method, we compared the  
917 extracted NMF signatures with the ground-truth simulated signatures. In particular, we  
918 considered that an extracted NMF signature matches the ground-truth simulated signatures  
919 when the cosine similarity is  $\geq 0.75$  only for that ground-truth simulated signature and  $< 0.75$  for  
920 the rest.

921

### 922 **Analysis of differential mutation supply towards cancer genes.**

923 For 460 cancer genes from the MutPanning list <sup>35</sup> (<http://www.cancer-genes.org/>), we tested if  
924 they are enriched in intronic mutations in tumor samples with high `RMDflat`, `RMDglobal1` or



925 RMDglobal2. An enrichment will mean that there is a higher supply of mutations in the intron  
926 regions of those genes when the RMDsignature is high. For this, we considered the counts of  
927 mutations in the intronic regions of the gene, normalized to the number of mutations in the  
928 whole chromosome arm, comparing pooled tumor samples with RMD signatures high or low, by  
929 tissue. Note that the possibly different number of eligible nucleotides-at-risk in the central  
930 window, nor the length of the flanking chromosome arm are relevant in this analysis, because  
931 they cancel out when comparing one group of tumor samples (split by the RMD signature) to  
932 another group of tumor samples. We binarized the tumor samples by RMDflat, RMDglobal1 and  
933 RMDglobal2 by dividing each of them into tertiles, and keeping 1st tertile *versus* 3rd tertile for  
934 further analysis. We applied a Poisson regression with the following formula:

935  $Count\_gene\_intron \sim offset(count\_chr\_arm) + RMDflat + RMDglobal1 + RMDglobal2 + tissue$

936 where “count” refers to mutation counts. By including the tissue as a variable in the regression,  
937 we controlled for possible confounding by cancer type. The log fold-difference in mutation  
938 supply between RMD signature high versus low tumor samples is estimated by the regression  
939 coefficients for RMDflat, RMDglobal1 and RMDglobal2 variables. As a control, we repeated the  
940 exact same analysis but randomizing the tertile assignment for the three RMD signatures prior  
941 to the regression.

942

#### 943 **Association analysis of gene mutations with RMD global signatures.**

944 We created a subset of 1543 relevant genes: cancer genes from the MutPanning list<sup>35</sup> and  
945 Cancer Gene Census list<sup>68</sup>, and furthermore we included genes associated with chromatin and  
946 DNA damage<sup>69</sup>. As control, we used a subset of 1000 random genes selected as in<sup>69</sup>.

947

948 We applied the analysis for two different features: copy number alterations (CNA) and  
949 deleterious point mutations. For CNA, we use the CN values by gene, using a score of -2, -1, 0,  
950 1 or 2 for each gene. We considered a gene to be amplified if CNA value was +1 or +2 and  
951 deleted if the CNA value was -1 or -2. For deleterious mutations, we selected mutations  
952 predicted as moderate or high impact in the Hartwig (HMF) variant calls,  
953 (<https://github.com/hartwigmedical/hmftools>). We binarized the feature into 1 if the sample has  
954 the feature (CNA, or deleterious mutations present) or 0 if it has not. We considered CNA  
955 deletions and amplifications as two independent features. We binarized RMDflat, RMDglobal1  
956 and RMDglobal2 by dividing each of them in tertiles and comparing tumor samples in 1st tertile  
957 *versus* 3rd tertile, by tissue.

958

959 We fit a linear model to test whether the binary genetic feature (amplification CNA, deletion CNA  
960 or deleterious mutation in a particular gene) can be explained by the RMD signatures activity  
961 being high *versus* low (i.e. upper tertile *versus* lower tertile). We controlled for tissue by  
962 including it as covariate. The regression formula was:

963  $genetic\_feature \sim RMDflat + RMDglobal1 + RMDglobal2 + tissue$

964 We used the regression coefficients, and p-values (according to the R function “summary”) from  
965 the variables RMDflat, RMDglobal1 and RMDglobal2 to identify genetic events associated with

966 high levels of each RMD global signatures, suggesting possible RMD signature generating  
967 events. In the case of CNAs, to adjust for the linkage between CNA resulting in confounding, we  
968 added to the regression the PCs from a PCA on the CNA landscape across all genes. We  
969 calculated the lambda (inflation factor) for the p-value distribution of associations, while  
970 including PCs from 1 to 100 to decide the best number of PCs to include so as to minimize  
971 lambda. We included the first 55 PCs for the deletion CNA and the first 63 PCs for the  
972 amplification CNA association study.

### 973 **Epigenomic and related data sources**

974 ENCODE data. We downloaded from ENCODE (<https://www.encodeproject.org/>) all data  
975 available for *Homo sapiens* in the genome assembly hg19 for DHS, H3F3A, H3K27me3,  
976 H3K4me1, H3K4me3, H3K9ac, H3K9me3, HiC, DNA methylation (WGBS), H2AFZ, H3K27ac,  
977 H3K36me3, H3K4me2, H3K79me2, H3K9me2 and H4K20me1 marks. Data is described in  
978 [Table S2](#). For each of these features, we downloaded the narrow peaks, calculated their  
979 weighted density for each 1Mb window as the width of the peak multiplied by the peak value.

980  
981 ChromHMM chromatin states. We downloaded the 25 ChromHMM states segmented files  
982 (“imputed12marks\_segments”) for the 129 cell types available from Roadmap epigenomics  
983 <sup>70</sup>(<http://compbio.mit.edu/ChromHMM/>). We calculated the density of each state for each 1Mb  
984 window as the fraction of the window covered by the chromatin state.

985  
986 Other epigenomic data. We downloaded RT variability genomic data describing RT  
987 heterogeneity <sup>71</sup>, Constitutive and Developmental RT domains <sup>72</sup>, RT changes upon  
988 overexpression of the oncogene *KDM4A* <sup>73</sup>, RT signatures of replication stress <sup>74</sup>, RT signatures  
989 of tissues <sup>41</sup>, RT states <sup>75</sup>, changes in RT upon *RIF1* knock-out <sup>76</sup> and RT changes due to RT  
990 QTLs <sup>77</sup>. In addition, we downloaded data for variability in DNA methylation <sup>15,78</sup>, HMD and PMD  
991 regions <sup>16</sup>, CpG density, gene density, lamina associated domains (LADs), asynchronous  
992 replication domains <sup>79</sup>, early replicating fragile sites <sup>80</sup>, SPIN states <sup>48</sup>, A/B subcompartments <sup>47</sup>,  
993 DHS signatures <sup>81</sup> and H3K27me3 and H3K9me profiles for *RB1* wild-type and knock-out <sup>19</sup>.  
994 Data described in [Table S3](#). We calculated the density for each feature for each 1 Mb window,  
995 and correlated this with the RMDglobal1 signature windows weights.

996

### 997 **Replication timing data sources and generation**

998 We downloaded experimental RT data, from RepliChip or RepliSeq assays, from the Replication  
999 Domain database (<https://www2.replicationdomain.com/index.php>) <sup>72</sup> in multiple human cell  
1000 types (n = 158 samples). In addition, we predicted RT using the Replicon software <sup>39</sup> from two  
1001 type datasets: (i) in noncancerous tissues, cultured primary cells and cell lines including cancer  
1002 and stem cells (n = 597 samples) using the DHS chromatin accessibility data downloaded from  
1003 ENCODE; and (ii) in human tumors (n = 410 samples, most of them with technical replicates)  
1004 using ATAC-seq data of TCGA tumors downloaded from <sup>38</sup>. We used Replicon tool with the  
1005 default settings.

1006

### 1007 **Analysis of coordinated gene expression changes**

1008 For the genomes from the HMF data set, we downloaded gene expression data (as adjusted  
1009 TPM values) from Hartwig<sup>57</sup>, available for a subset of samples for which we derived the RMD  
1010 signatures. In total, we had gene expression data for 1534 samples and 18889 protein coding  
1011 genes. We tested whether the gene expression values of the genes within one window show an  
1012 increase or decrease compared to their flanking windows in RMDglobal1 high (exposure  $\geq$   
1013 0.13) versus RMDglobal1 low (exposure  $<$  0.06) tumor samples using a linear regression model:  
1014  $gene\_expression (adjTPM) \sim is\_RMDglobal1 + is\_window + is\_RMDglobal1:is\_window + tissue$   
1015

1016 In this analysis, we removed from the datasets samples with high RMDflat or with high  
1017 RMDglobal2 value (exposure  $>$  0.15). We used samples from breast, colorectum, lung, ovary  
1018 and skin because they had  $\geq 5$  samples in both categories (RMDglobal1 high and low). To  
1019 analyze the coordinated changes in gene expression we checked the coefficient and p-values of  
1020 the interaction term  $is\_RMDglobal1:is\_window$ .

1021 For the genomes from the TCGA data set, we downloaded gene expression data (as TPM  
1022 values) from the Genomic Data Commons data portal (<https://dcc.icgc.org/pcawg>) for the same  
1023 TCGA samples for which we predicted RT. In total, we have gene expression data for 399  
1024 overlapping samples and 20092 genes. We compared the gene expression between RT-PC5  
1025 (and RT-PC6) high and low for a group of pathways which has been reported to be related with  
1026 recurrent heterogeneity across cell types<sup>12</sup> using a regression model. We binarized RT-PC5  
1027 (and RT-PC6) by dividing each into tertiles and keeping the samples in the 1st tertile to be  
1028 compared versus the samples in the 3rd tertile. We applied a regression for all the genes in  
1029 each RHP gene set separately. We controlled for tissue by including it as covariate. The  
1030 regression formula is:

1031  $gene\_expression (TPM) \sim is\_RT-PC5 + tissue$

1032 We considered the regression coefficient and its p-value of the variable  $is\_RT-PC5$ . We applied  
1033 the same analysis for RT-PC6.

### 1034 **Clustering of RMD profiles**

1035 For RMD profiles we applied a PCA to the centered data, where rows were tumor samples and  
1036 the columns were megabase windows. Next, we applied a clustering on the PC1 to PC21 using  
1037 the R function tclust for robust clustering. We tested different numbers of clusters and alpha  
1038 value (number of outliers removed). In addition, we tested the clustering using all PCs (PC1 to  
1039 PC21) and without PC1 (PC2 to PC21), selecting the clustering for  $k=18$  and  $\alpha = 0.02$   
1040 without PC1 based on the log likelihood measurement.

1041

1042

1043

### 1044 **Acknowledgements**

1045 M.S. was funded by a FPU fellowship of the Spanish government, Ministry of Universities. Work  
1046 in the lab of F.S. is supported by an ERC StG “HYPER-INSIGHT” (757700), Horizon2020  
1047 project “DECIDER” (965193), Spanish government project “REPAIRSCAPE”, CaixaResearch  
1048 project “POTENT-IMMUNO” (HR22-00402), an ICREA professorship to F.S., the SGR funding

1049 of the Catalan government, and the Severo Ochoa centers of excellence award of the Spanish  
1050 government to the hosting institution.

1051 This publication and the underlying research are partly facilitated by Hartwig Medical  
1052 Foundation and the Center for Personalized Cancer Treatment (CPCT) which have generated,  
1053 analysed and made available data for this research. In addition, data used in this publication  
1054 were generated by the Clinical Proteomic Tumor Analysis Consortium (NCI/NIH). We  
1055 acknowledge that the results published here are in part based upon data generated by the  
1056 TCGA Research Network: <http://cancergenome.nih.gov/>.

1057  
1058

## 1059 References:

- 1060 1. Supek, F. & Lehner, B. Differential DNA mismatch repair underlies mutation rate variation across the  
1061 human genome. *Nature* **521**, 81–84 (2015).
- 1062 2. Zheng, C. L. *et al.* Transcription Restores DNA Repair to Heterochromatin, Determining Regional  
1063 Mutation Rates in Cancer Genomes. *Cell Rep.* **9**, 1228–1234 (2014).
- 1064 3. Pope, B. D. *et al.* Topologically associating domains are stable units of replication-timing regulation.  
1065 *Nature* **515**, 402–405 (2014).
- 1066 4. Akdemir, K. C. *et al.* Somatic mutation distributions in cancer genomes vary with three-dimensional  
1067 chromatin structure. *Nat. Genet.* **52**, 1178–1188 (2020).
- 1068 5. Hansen, R. S. *et al.* Sequencing newly replicated DNA reveals widespread plasticity in human  
1069 replication timing. *Proc. Natl. Acad. Sci.* **107**, 139–144 (2010).
- 1070 6. Polak, P. *et al.* Cell-of-origin chromatin organization shapes the mutational landscape of cancer.  
1071 *Nature* **518**, 360–364 (2015).
- 1072 7. Makova, K. D. & Hardison, R. C. The effects of chromatin organization on variation in mutation rates  
1073 in the genome. *Nat. Rev. Genet.* **16**, 213–223 (2015).
- 1074 8. Schuster-Böckler, B. & Lehner, B. Chromatin organization is a major influence on regional mutation  
1075 rates in human cancer cells. *Nature* **488**, 504–507 (2012).
- 1076 9. Salvadores, M., Mas-Ponte, D. & Supek, F. Passenger mutations accurately classify human tumors.  
1077 *PLOS Comput. Biol.* **15**, e1006953 (2019).
- 1078 10. Jiao, W. *et al.* A deep learning system accurately classifies primary and metastatic cancers using  
1079 passenger mutation patterns. *Nat. Commun.* **11**, 1–12 (2020).
- 1080 11. Kübler, K. *et al.* Tumor mutational landscape is a record of the pre-malignant state. 517565 Preprint  
1081 at <https://doi.org/10.1101/517565> (2019).
- 1082 12. Kinker, G. S. *et al.* Pan-cancer single cell RNA-seq uncovers recurring programs of cellular  
1083 heterogeneity. *Nat. Genet.* **52**, 1208–1218 (2020).
- 1084 13. Barkley, D. *et al.* Cancer cell states recur across tumor types and form specific interactions with the  
1085 tumor microenvironment. *Nat. Genet.* **54**, 1192–1201 (2022).
- 1086 14. Du, Q. *et al.* Replication timing and epigenome remodelling are associated with the nature of  
1087 chromosomal rearrangements in cancer. *Nat. Commun.* **10**, 416 (2019).
- 1088 15. Du, Q. *et al.* DNA methylation is required to maintain both DNA replication timing precision and 3D  
1089 genome organization integrity. *Cell Rep.* **36**, 109722 (2021).
- 1090 16. Zhou, W. *et al.* DNA methylation loss in late-replicating domains is linked to mitotic cell division. *Nat.*  
1091 *Genet.* **50**, 591–602 (2018).
- 1092 17. Brinkman, A. B. *et al.* Partially methylated domains are hypervariable in breast cancer and fuel  
1093 widespread CpG island hypermethylation. *Nat. Commun.* **10**, 1749 (2019).
- 1094 18. Gurrion, C., Uriostegui, M. & Zurita, M. Heterochromatin Reduction Correlates with the Increase of  
1095 the KDM4B and KDM6A Demethylases and the Expression of Pericentromeric DNA during the  
1096 Acquisition of a Transformed Phenotype. *J. Cancer* **8**, 2866–2875 (2017).
- 1097 19. Wong, K. M., King, D. A., Schwartz, E. K., Herrera, R. E. & Morrison, A. J. Retinoblastoma protein  
1098 regulates carcinogen susceptibility at heterochromatic cancer driver loci. *Life Sci. Alliance* **5**,  
1099 e202101134 (2022).
- 1100 20. Huang, Y., Gu, L. & Li, G.-M. H3K36me3-mediated mismatch repair preferentially protects actively  
1101 transcribed genes from mutation. *J. Biol. Chem.* **293**, 7811–7823 (2018).



- 1102 21. Poetsch, A. R., Boulton, S. J. & Luscombe, N. M. Genomic landscape of oxidative DNA damage and  
1103 repair reveals regioselective protection from mutagenesis. *Genome Biol.* **19**, 215 (2018).
- 1104 22. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415–421  
1105 (2013).
- 1106 23. Mas-Ponte, D. & Supek, F. DNA mismatch repair promotes APOBEC3-mediated diffuse  
1107 hypermutation in human cancers. *Nat. Genet.* 1–11 (2020) doi:10.1038/s41588-020-0674-6.
- 1108 24. Seplyarskiy, V. B. *et al.* APOBEC-induced mutations in human cancers are strongly enriched on the  
1109 lagging DNA strand during replication. *Genome Res.* **26**, 174–182 (2016).
- 1110 25. Jönsson, J.-M. *et al.* Molecular Subtyping of Serous Ovarian Tumors Reveals Multiple Connections to  
1111 Intrinsic Breast Cancer Subtypes. *PLOS ONE* **9**, e107643 (2014).
- 1112 26. Supek, F. & Lehner, B. Clustered Mutation Signatures Reveal that Error-Prone DNA Repair Targets  
1113 Mutations to Active Genes. *Cell* **170**, 534-547.e23 (2017).
- 1114 27. Degasperi, A. *et al.* Substitution mutational signatures in whole-genome–sequenced cancers in the  
1115 UK population. *Science* **376**, abl9283 (2022).
- 1116 28. Hoadley, K. A. *et al.* Multiplatform analysis of 12 cancer types reveals molecular classification within  
1117 and across tissues of origin. *Cell* **158**, 929–944 (2014).
- 1118 29. Kim, J. *et al.* Somatic ERCC2 mutations are associated with a distinct genomic signature in urothelial  
1119 tumors. *Nat. Genet.* **48**, 600–606 (2016).
- 1120 30. Chen, J., Miller, B. F. & Furano, A. V. Repair of naturally occurring mismatches can induce mutations  
1121 in flanking DNA. *eLife* **3**, e02001 (2014).
- 1122 31. Nguyen, L., W. M. Martens, J., Van Hoeck, A. & Cuppen, E. Pan-cancer landscape of homologous  
1123 recombination deficiency. *Nat. Commun.* **11**, 5584 (2020).
- 1124 32. Chen, D. *et al.* BRCA1 deficiency specific base substitution mutagenesis is dependent on translesion  
1125 synthesis and regulated by 53BP1. *Nat. Commun.* **13**, 226 (2022).
- 1126 33. Franco, I. *et al.* Whole genome DNA sequencing provides an atlas of somatic mutagenesis in healthy  
1127 human cells and identifies a tumor-prone cell type. *Genome Biol.* **20**, 285 (2019).
- 1128 34. Supek, F. & Lehner, B. Scales and mechanisms of somatic mutation rate variation across the human  
1129 genome. *DNA Repair* **81**, 102647 (2019).
- 1130 35. Dietlein, F. *et al.* Identification of cancer driver genes based on nucleotide context. *Nat. Genet.* **52**,  
1131 208–218 (2020).
- 1132 36. Yaeger, R. *et al.* Clinical Sequencing Defines the Genomic Landscape of Metastatic Colorectal  
1133 Cancer. *Cancer Cell* **33**, 125-136.e3 (2018).
- 1134 37. Moore, J. E. *et al.* Expanded encyclopaedias of DNA elements in the human and mouse genomes.  
1135 *Nature* **583**, 699–710 (2020).
- 1136 38. Corces, M. R. *et al.* The chromatin accessibility landscape of primary human cancers. *Science* **362**,  
1137 eaav1898 (2018).
- 1138 39. Gindin, Y., Meltzer, P. S. & Bilke, S. Replicon: a software to accurately predict DNA replication timing  
1139 in metazoan cells. *Front. Genet.* **5**, (2014).
- 1140 40. Ryba, T. *et al.* Abnormal developmental control of replication-timing domains in pediatric acute  
1141 lymphoblastic leukemia. *Genome Res.* **22**, 1833–1844 (2012).
- 1142 41. Rivera-Mulia, J. C. *et al.* Dynamic changes in replication timing and gene expression during lineage  
1143 specification of human pluripotent stem cells. *Genome Res.* **25**, 1091–1103 (2015).
- 1144 42. Gonzalo, S. *et al.* Role of the RB1 family in stabilizing histone methylation at constitutive  
1145 heterochromatin. *Nat. Cell Biol.* **7**, 420–428 (2005).
- 1146 43. Krishnan, B. *et al.* Active RB causes visible changes in nuclear organization. *J. Cell Biol.* **221**,  
1147 e202102144 (2022).
- 1148 44. Dick, F. A., Goodrich, D. W., Sage, J. & Dyson, N. J. Non-canonical functions of the RB protein in  
1149 cancer. *Nat. Rev. Cancer* **18**, 442–451 (2018).
- 1150 45. Takahashi, C., Contreras, B., Bronson, R. T., Loda, M. & Ewen, M. E. Genetic Interaction between  
1151 Rb and K-ras in the Control of Differentiation and Tumor Suppression. *Mol. Cell. Biol.* **24**, 10406–  
1152 10415 (2004).
- 1153 46. Lee, K. Y., Ladha, M. H., McMahon, C. & Ewen, M. E. The Retinoblastoma Protein Is Linked to the  
1154 Activation of Ras. *Mol. Cell. Biol.* **19**, 7724–7732 (1999).
- 1155 47. Rao, S. S. P. *et al.* A 3D Map of the Human Genome at Kilobase Resolution Reveals Principles of  
1156 Chromatin Looping. *Cell* **159**, 1665–1680 (2014).
- 1157 48. SPIN reveals genome-wide landscape of nuclear compartmentalization | Genome Biology | Full Text.



- 1158 <https://genomebiology.biomedcentral.com/articles/10.1186/s13059-020-02253-3>.
- 1159 49. Cramer, D., Serrano, L. & Schaefer, M. H. A network of epigenetic modifiers and DNA repair genes  
1160 controls tissue-specific copy number alteration preference. *eLife* **5**, e16519 (2016).
- 1161 50. Donehower, L. A. *et al.* Integrated Analysis of TP53 Gene and Pathway Alterations in The Cancer  
1162 Genome Atlas. *Cell Rep.* **28**, 1370-1384.e5 (2019).
- 1163 51. Campbell, P. J. *et al.* Pan-cancer analysis of whole genomes. *Nature* **578**, 82–93 (2020).
- 1164 52. Nguyen, L., Van Hoeck, A. & Cuppen, E. Machine learning-based tissue of origin classification for  
1165 cancer of unknown primary diagnostics using genome-wide mutation features. *Nat. Commun.* **13**,  
1166 4013 (2022).
- 1167 53. Zou, X. *et al.* Validating the concept of mutational signatures with isogenic cell models. *Nat.*  
1168 *Commun.* **9**, 1–16 (2018).
- 1169 54. Drost, J. *et al.* Use of CRISPR-modified human stem cell organoids to study the origin of mutational  
1170 signatures in cancer. *Science* **358**, 234–238 (2017).
- 1171 55. Xiong, K. & Ma, J. Revealing Hi-C subcompartments by imputing inter-chromosomal chromatin  
1172 interactions. *Nat. Commun.* **10**, 5069 (2019).
- 1173 56. Hudson (Chairperson), T. J. *et al.* International network of cancer genome projects. *Nature* **464**, 993–  
1174 998 (2010).
- 1175 57. Priestley, P. *et al.* Pan-cancer whole-genome analyses of metastatic solid tumours. *Nature* **575**, 210–  
1176 216 (2019).
- 1177 58. Pleasance, E. *et al.* Pan-cancer analysis of advanced patient tumors reveals interactions between  
1178 therapy and genomic landscapes. *Nat. Cancer* **1**, 452–468 (2020).
- 1179 59. Ellis, M. J. *et al.* Connecting genomic alterations to cancer biology with proteomics: the NCI Clinical  
1180 Proteomic Tumor Analysis Consortium. *Cancer Discov.* **3**, 1108–1112 (2013).
- 1181 60. Edwards, N. J. *et al.* The CPTAC Data Portal: A Resource for Cancer Proteomics Research. *J.*  
1182 *Proteome Res.* **14**, 2707–2713 (2015).
- 1183 61. Walker, B. A. *et al.* A high-risk, Double-Hit, group of newly diagnosed myeloma identified by genomic  
1184 analysis. *Leukemia* **33**, 159–170 (2019).
- 1185 62. Kim, S. *et al.* Strelka2: fast and accurate calling of germline and somatic variants. *Nat. Methods* **15**,  
1186 591–594 (2018).
- 1187 63. Salvadores, M., Fuster-Tormo, F. & Supek, F. Matching cell lines with cancer type and subtype of  
1188 origin via mutational, epigenomic, and transcriptomic patterns. *Sci. Adv.* **6**, eaba1862 (2020).
- 1189 64. Derrien, T. *et al.* Fast Computation and Applications of Genome Mappability. *PLOS ONE* **7**, e30377  
1190 (2012).
- 1191 65. Ormond, C., Ryan, N. M., Corvin, A. & Heron, E. A. Converting single nucleotide variants between  
1192 genome builds: from cautionary tale to solution. *Brief. Bioinform.* **22**, bbab069 (2021).
- 1193 66. Amemiya, H. M., Kundaje, A. & Boyle, A. P. The ENCODE Blacklist: Identification of Problematic  
1194 Regions of the Genome. *Sci. Rep.* **9**, 9354 (2019).
- 1195 67. Buisson, R. *et al.* Passenger hotspot mutations in cancer driven by APOBEC3A and mesoscale  
1196 genomic features. *Science* **364**, (2019).
- 1197 68. Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res.* **47**,  
1198 D941–D947 (2019).
- 1199 69. Vali-Pour, M., Lehner, B. & Supek, F. The impact of rare germline variants on human somatic  
1200 mutation processes. *Nat. Commun.* **13**, 3724 (2022).
- 1201 70. Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330  
1202 (2015).
- 1203 71. Zhao, P. A., Sasaki, T. & Gilbert, D. M. High-resolution Repli-Seq defines the temporal choreography  
1204 of initiation, elongation and termination of replication in mammalian cells. *Genome Biol.* **21**, 76  
1205 (2020).
- 1206 72. Sima, J. *et al.* Identifying cis Elements for Spatiotemporal Control of Mammalian DNA Replication.  
1207 *Cell* **176**, 816-830.e18 (2019).
- 1208 73. Van Rechem, C. *et al.* Collective regulation of chromatin modifications predicts replication timing  
1209 during cell cycle. *Cell Rep.* **37**, 109799 (2021).
- 1210 74. Sarni, D. *et al.* Replication Timing and Transcription Identifies a Novel Fragility Signature Under  
1211 Replication Stress. 716951 Preprint at <https://doi.org/10.1101/716951> (2019).
- 1212 75. Poulet, A. *et al.* RT States: systematic annotation of the human genome using cell type-specific  
1213 replication timing programs. *Bioinformatics* **35**, 2167–2176 (2019).

- 1214 76. Klein, K. N. *et al.* Replication timing maintains the global epigenetic state in human cells. *Science*  
1215 **372**, 371–378 (2021).  
1216 77. Ding, Q. *et al.* The genetic architecture of DNA replication timing in human pluripotent stem cells. *Nat.*  
1217 *Commun.* **12**, 6746 (2021).  
1218 78. Gunasekara, C. J. *et al.* A genomic atlas of systemic interindividual epigenetic variation in humans.  
1219 *Genome Biol.* **20**, 105–105 (2019).  
1220 79. Mukhopadhyay, R. *et al.* Allele-Specific Genome-wide Profiling in Human Primary Erythroblasts  
1221 Reveal Replication Program Organization. *PLoS Genet.* **10**, e1004319 (2014).  
1222 80. Barlow, J. H. *et al.* Identification of Early Replicating Fragile Sites that Contribute to Genome  
1223 Instability. *Cell* **152**, 620–632 (2013).  
1224 81. Meuleman, W. *et al.* Index and biological spectrum of human DNase I hypersensitive sites. *Nature*  
1225 **584**, 244–251 (2020).