

1           **SEQUENCING CHEMICALLY INDUCED MUTATIONS IN THE MUTAMOUSE LACZ**  
2           **REPORTER GENE IDENTIFIES HUMAN CANCER MUTATIONAL SIGNATURES**

3  
4   Marc A. Beal<sup>1,4\*</sup>, Matt J. Meier<sup>2\*</sup>, Danielle LeBlanc<sup>1</sup>, Clotilde Maurice<sup>1</sup>, Jason O'Brien<sup>3</sup>, Carole L.  
5   Yauk<sup>1</sup>, Francesco Marchetti<sup>1,5</sup>

6  
7   \*These authors contributed equally to this work.

8   <sup>1</sup>Environmental Health Science and Research Bureau, Healthy Environments and Consumer  
9   Safety Branch, Health Canada, Ottawa, Ontario, K1A 0K9, Canada.

10   <sup>2</sup>Science and Technology Branch, Environment and Climate Change Canada, Ottawa, Ontario,  
11   K1A 0H3, Canada.

12   <sup>3</sup>National Wildlife Research Centre, Environment and Climate Change Canada, Ottawa, Ontario,  
13   K1A 0H3, Canada.

14   <sup>4</sup>Present address: Existing Substances Risk Assessment Bureau, Health Canada, Ottawa,  
15   Ontario, Canada

16   <sup>5</sup>Corresponding author: [francesco.marchetti@canada.ca](mailto:francesco.marchetti@canada.ca)

17  
18   Other email addresses: [marc.beal@canada.ca](mailto:marc.beal@canada.ca), [matthew.meier@canada.ca](mailto:matthew.meier@canada.ca),  
19   [danielle.leblanc2@canada.ca](mailto:danielle.leblanc2@canada.ca), [clotilde.maurice@canada.ca](mailto:clotilde.maurice@canada.ca), [jason.obrien@canada.ca](mailto:jason.obrien@canada.ca),  
20   [carole.yauk@canada.ca](mailto:carole.yauk@canada.ca)

21  
22   Running title: Induced lacZ mutations and human cancer signatures

23 **ABSTRACT**

24 Transgenic rodent (TGR) models use bacterial reporter genes to quantify *in vivo*  
25 mutagenesis. Pairing TGR assays with next-generation sequencing (NGS) enables  
26 comprehensive mutation spectrum analysis to inform mutational mechanisms. We used this  
27 approach to identify 2,751 independent *lacZ* mutations in the bone marrow of MutaMouse  
28 animals exposed to four chemical mutagens: benzo[a]pyrene, *N*-ethyl-*N*-nitrosourea,  
29 procarbazine, and triethylenemelamine. We also collected published data for 706 *lacZ*  
30 mutations from eight additional environmental mutagens. We demonstrate that *lacZ* gene  
31 sequencing generates chemical-specific mutation signatures observed in human cancers with  
32 established environmental causes. For example, the mutation signature of benzo[a]pyrene, a  
33 potent carcinogen in tobacco smoke, matched the signature associated with tobacco-induced  
34 lung cancers. Our results show that the analysis of chemically induced mutations in the *lacZ*  
35 gene shortly after exposure provides an effective approach to characterize human-relevant  
36 mechanisms of carcinogenesis and identify novel environmental causes of mutation signatures  
37 observed in human cancers.

38

39

40 Key words: Mutagenesis, COSMIC, cancer, Next Generation Sequencing, Benzo(a)pyrene, N-  
41 ethyl-N-nitrosourea, Procarbazine, Triethylenemelamine

## 42 INTRODUCTION

43 Transgenic rodent (TGR) mutation reporter models have enabled unprecedented  
44 insights into spontaneous and chemically induced mutagenesis<sup>1</sup>. Studies of over 200 chemicals,  
45 including more than 90 carcinogens, have demonstrated that TGR models offer high sensitivity  
46 and specificity for identifying mutagenic carcinogens<sup>1,2</sup>. One of the most commonly used TGR  
47 models is the MutaMouse whose genome was recently sequenced<sup>3</sup>. The MutaMouse harbors  
48 ~29 copies of the bacterial *lacZ* transgene on each copy of chromosome 3<sup>4</sup>. This is a neutral,  
49 transcriptionally-inert reporter gene carried on a shuttle vector that can be recovered from any  
50 cell type and transfected into a bacterial host to detect somatic or germline mutations that  
51 occurred *in vivo*<sup>5,6</sup>. A major advantage of TGR models is the possibility to sequence mutants in  
52 order to characterize mutation spectra. This information is necessary to understand mutational  
53 mechanisms associated with mutagen exposure and response in different tissues, life stages  
54 genetic backgrounds or other contexts. Advances in next-generation sequencing (NGS)  
55 technologies have enabled rapid and accurate characterization of TGR mutants<sup>7,8</sup>, and  
56 integrated TGR-NGS approaches have been used to sequence thousands of mutations<sup>8,9</sup> at a  
57 fraction of the cost of whole genome sequencing. Thus, TGR-NGS approaches currently  
58 provide a unique methodology for simultaneously assessing the magnitude of the mutagenic  
59 response and mutation spectrum to inform underlying mechanisms.

60 Somatic mutation analysis by NGS has greatly advanced our understanding of the  
61 mutational processes operating in human cancers. Algorithms have been developed to mine the  
62 extensive database of single nucleotide variations (SNVs) in cancer genomes to identify  
63 mutational signatures contributing to individual cancers<sup>10,11,12</sup>. These signatures represent a  
64 computationally derived prediction of the relative frequencies of mutation types induced by  
65 processes that contribute to all observed mutations within The Cancer Genome Atlas datasets  
66 (TCGA; <https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga>).

67 As opposed to standard mutation spectrum characterization that simply describes the frequency  
68 of individual nucleotide changes, mutational signatures incorporate flanking nucleotide context.  
69 Originally, 30 mutational signatures from 40 different cancer types were identified and reported  
70 in the Catalogue of Somatic Mutations in Cancer (COSMIC) database<sup>13,14</sup>. This database was  
71 recently expanded to include 71 cancer types and 77 signatures, including 49 single base  
72 substitution (SBS) signatures, 11 doublet base substitution (DBS) signatures, and 17 small  
73 insertion and deletion (ID) signatures<sup>15</sup>. Each signature encompasses 96 possible mutation  
74 types (i.e., 6 possible base pair alterations × 4 different 5' bases × 4 different 3' bases). Many of  
75 these signatures have been attributed to endogenous processes, but chemical mutagens also  
76 play a major contributing role in certain signatures<sup>16</sup>. For example, SBS 4 signature is observed  
77 in lung cancer and has been attributed to tobacco smoke<sup>16,17</sup>. This signature has been  
78 recapitulated by exposing murine embryo fibroblasts to benzo[a]pyrene (BaP)<sup>18,19</sup>, a major  
79 mutagenic component of tobacco smoke. However, several of the mutational signatures  
80 currently have no known endogenous or exogenous causative agents<sup>17</sup>; thus, identification of  
81 exogenous environmental exposures that contribute to these mutational signatures may aid in  
82 elucidating carcinogenic mechanisms.

83 The pattern of mutations observed in a fully developed cancer is a composite of the  
84 signature of the molecular initiating events in the early stages of tumour formation and  
85 signatures arising as a result of genomic instability in the evolving tumour<sup>20</sup>. For example, a  
86 tumour that originates in the lung of a smoker will have a mutational fingerprint that is caused  
87 primarily by DNA damage induced by the many mutagenic compounds found in tobacco  
88 smoke<sup>21</sup>. In addition, the person's age at the time of tumour formation will also determine the  
89 contribution of "clock-like" signatures, caused by lifetime DNA replication, to the fingerprint of  
90 the tumour<sup>22</sup>. There is now compelling evidence that analysis of the spectrum of mutations in a  
91 cancer can provide clues to past environmental exposures that contributed to the development  
92 of the cancer<sup>23,24</sup>. Implicit in this is that the exposure signature should be present in the normal

93 tissue before the carcinogenic process becomes apparent. Indeed, previous studies have  
94 demonstrated that mutational signatures observed in aflatoxin-induced cancers are observed in  
95 normal tissues long before tumour formation<sup>25, 26</sup>. Recent work *in vivo*<sup>27</sup> and *in vitro*<sup>28</sup> has shown  
96 that chemical-specific signatures detected shortly after exposure match signatures seen in  
97 human cancers. Thus, characterization of short-term mutational signatures in non-tumour  
98 tissues is a valuable approach to elucidate human-relevant mechanisms of carcinogenesis.

99 In this study, we used TGR-NGS to characterize mutations induced by four established  
100 mutagens to determine if these mutation profiles inform carcinogenic mechanisms within  
101 COSMIC signatures. For this purpose, we chose four chemicals with varying mutagenic  
102 potencies, mode of action, and carcinogenic classification (as determined by the International  
103 Agency for Research on Cancer): one known class 1 carcinogen, BaP; two probable class 2  
104 carcinogens including *N*-ethyl-*N*-nitrosourea (ENU) and procarbazine (PRC); and one class 3  
105 chemical with inadequate information to be classified, triethylenemelamine (TEM). MutaMouse  
106 males were exposed by gavage to the chemicals or solvent for 28 days and DNA was collected  
107 from bone marrow for analysis. To further compare *lacZ* mutation spectra and COSMIC  
108 signatures, published Sanger sequencing data from 17 studies involving eight mutagens were  
109 also examined (Supplementary Table S1). These studies include data from mice exposed to  
110 electromagnetic radiation<sup>29,30,31,32,33</sup>, alkylating agents and adduct-forming agents<sup>34,35,36,37,38,39,40</sup>,  
111 and a nitrogenous base analog<sup>41</sup>. Data from control animals in these studies and  
112 others<sup>42,43,44,45,46</sup> were also included to generate a background mutation signature. Using *lacZ*-  
113 derived mutation data, we validated COSMIC signatures with proposed aetiologies through the  
114 identification of the expected signatures in the relevant exposure groups. We argue that  
115 analysis of COSMIC signatures observed in exposed animals can be used to generate or test  
116 hypotheses of mutagenic mechanisms associated with human mutational signatures of  
117 unknown etiology.

118

## 119 **RESULTS**

120 We used mutation spectra generated in-house for four chemicals (BaP, ENU, PRC, and  
121 TEM) and vehicle matched controls, and published data from eight agents, including BaP and  
122 ENU (Supplementary Table S1) and their matched controls, to query the COSMIC database  
123 and elucidate the role of environmental mutagens in cancer development. The overall  
124 experimental design is summarized in Figure 1.

125 Mutation spectra were generated from plaques collected during experiments aimed at  
126 evaluating the induction of mutations in the bone marrow of MutaMouse males exposed to  
127 either BaP, ENU, PRC, or TEM using the *lacZ* assay<sup>5</sup>. Mutant frequencies were previously  
128 reported for BaP<sup>8</sup>, PRC<sup>47</sup> and TEM<sup>48</sup>, while mutant frequencies are reported here for the first  
129 time for ENU. All of the exposures caused increases in mutant frequencies relative to vehicle-  
130 matched controls (Supplementary Table 2), and the results were highly significant ( $P < 0.0001$ )  
131 for BaP (122.9-fold), PRC (9.7-fold), and ENU (7.2-fold). TEM exposure also increased mutant  
132 frequency relative to controls (1.6-fold;  $P = 0.048$ ), but it was less potent than the other agents.  
133 The potency ranking of exposures (BaP > PRC/ENU > TEM) was consistent with expectations.

134

### 135 **Mutation Characterization and Spectral Analysis**

136 Sequencing of 5,419 mutant plaques from bone marrow DNA enabled the  
137 characterization of 2,751 independent mutations (Supplementary Table S3). Sequenced  
138 plaques from BaP, ENU, and controls were generated by both NGS and Sanger sequencing.  
139 Specifically, there were 1105, 406, and 438 mutations identified by NGS for BaP, ENU, and  
140 Controls, respectively. The corresponding numbers were 60, 207, and 508 for Sanger  
141 sequencing. The mutation spectra generated by the two sequencing approaches were  
142 consistent for each of the three groups (data not shown). Thus, within each group, the two sets  
143 of mutations were combined. Overall, there were 1,046, 2,914, 129, 902, and 428 mutants

144 sequenced in the Controls, BaP, PRC, ENU, and TEM groups, respectively. These sequenced  
145 mutants represented 512, 1,547, 120, 419, and 153 independent mutations in the five groups,  
146 respectively.

147 In the *lacZ* gene, there are 3,096 positions  $\times$  3 possible substitutions at each position for  
148 a total of 9,288 possible unique SNV events; however, not all of these can be detected using a  
149 functional assay, since many result in silent mutations. Sequencing mutants from the different  
150 groups identified 891 unique SNVs, 338 of which overlapped between two or more groups  
151 (Supplementary Figure S1). Specific to each group, there were 55, 377, 14, 85, and 22 unique  
152 SNVs for Controls, BaP, PRC, ENU, and TEM, respectively (Supplementary Table S3). The  
153 mutations detected in this study are limited almost exclusively to point mutations and small  
154 indels (1-21 bp), as large deletions are infrequently recovered during packaging of the DNA for  
155 the *lacZ* assay<sup>8</sup>.

156 The mutation spectra of the four chemicals were significantly different from the control  
157 mutation spectrum (Figure 2;  $P \leq 0.0008$ ). The COSMIC convention is to represent mutations  
158 based on pyrimidine changes; thus, we present our mutation spectrum using the same  
159 convention. The main spontaneous mutation is represented by C>T transitions, which are  
160 thought to arise through spontaneous mechanisms such as deamination of methylated  
161 cytosines<sup>49</sup>. Although there may be proportional declines in specific mutations relative to  
162 controls (Figure 2), all of the chemicals tested in this study, with the exception of TEM,  
163 increased the mutation frequency of substitutions (e.g., C>T; Supplementary Figure S2).

164 The mutation spectra of BaP and ENU are consistent with previous observations. BaP  
165 exposure caused cytosine transversions and indels (Figure 2), mainly C>A SNVs, consistent  
166 with the formation of bulky DNA adducts mostly at the N2 of guanine<sup>8</sup>. ENU induced T>A  
167 mutations consistent with alkylation of thymine, specifically O<sup>2</sup>- and O<sup>4</sup>-ethyl thymine<sup>50, 51</sup>. We  
168 found that PRC induced T>A mutations and, to a lesser extent T>C mutations, which is  
169 consistent with the pattern of mutations that was observed in an endogenous gene<sup>52</sup>. The

170 mutation spectrum of TEM was significantly different from controls, but there were no significant  
171 changes in specific SNV types. Instead, this effect is mainly driven by the higher proportion of  
172 TEM-induced single nucleotide insertions compared to control animals. TEM also induced the  
173 highest proportion of >1bp indels among all chemicals tested (Figure 2).

174

## 175 **Identification of COSMIC Signatures Using *lacZ* Mutations**

176 We explored the use of the *lacZ* sequence to obtain mutational signatures associated  
177 with human cancers. Although the COSMIC database (version 3) includes also DBS and ID  
178 signatures, we focused on SBS signatures because the *lacZ* assay detects almost exclusively  
179 these types of events. We first divided each trinucleotide frequency in the *lacZ* transgene  
180 (Figure 3) by the respective human genome frequencies (hg38) to create a *lacZ*-normalized set  
181 of the 49 COSMIC SBS signatures (Supplementary Figure S3). We then used the *lacZ*  
182 sequencing data from NGS and Sanger experiments in COSMIC format (Supplementary Figure  
183 S4) to identify which of the normalized signatures were most closely associated with the  
184 mutation spectrum of each agent. This initial analysis showed that mutational signatures in  
185 human cancers that have been associated with specific mutagenic exposures were enriched in  
186 the *lacZ* mutation profiles for the appropriate agent tested in this study (Figure 4). For example,  
187 the UVB<sup>29, 31</sup> and sunlight<sup>30</sup> mutation profiles had very strong correlations (Pearson's coefficient  
188 = 0.93-0.98) with the SBS 7a signature, which is observed in human skin cancers. Similarly, the  
189 BaP mutation profile showed a strong correlation with several signatures including SBS 4  
190 (Pearson's coefficient = 0.76), which is observed in tobacco smoke-induced cancers. In total,  
191 there were six SBS signatures that had a Pearson's coefficient greater than 0.8 with the  
192 mutation spectra generated from sequenced *lacZ* mutations (Figure 4).

193 Next, we used each of the agent-generated mutation data to simultaneously query the  
194 entire COSMIC SBS database to establish which of the signatures contributed to the observed



195 spectra (Supplementary Figure S5). This process was conducted as described in steps 7-9 of  
196 Figure 1. Prior to this analysis, we used control mutation data to generate an *in vivo* background  
197 signature (Figure 5) to account for the fact that some mutations present in the exposure groups  
198 are also spontaneous in origin rather than specific to the mutagen tested. This is especially true  
199 for weak mutagens. As shown in Figure 5, the *in vivo* background signature is enriched primarily  
200 in C>T mutations, and to a lesser extent C>A mutations, and this was consistent among all  
201 tissues that contributed to the control signature (Supplementary Figure S6). Inclusion of the  
202 control signature improved the association between reconstructed signatures and mutation data  
203 by as much as 26% (Figure 6) and eliminated five of the weakly associated signatures  
204 (Supplementary Figure S7). Finally, application of stringent filtering criteria (see Methods)  
205 revealed the association of nine COSMIC SBS signatures with mutation data from the various  
206 exposure groups (Figure 6).

207         The signatures produced by the three electromagnetic radiations (i.e., UVB, sunlight,  
208 and X-rays) appear to be broadly similar when visually assessing individual SBS signature  
209 heatmaps (Figure 4). However, we found that different mutational processes contribute to each  
210 signature. Specifically, SBS 2 and the control signature each explained 33% of the UVB  
211 mutation profile; SBS 2 and SBS 7a each explained 27% of the sunlight data (Figure 6); and the  
212 mutation spectrum associated with X-rays, which induces large deletions rather than point  
213 mutations (56 indels ranging from 1-437 bp vs 35 SNVs<sup>33</sup>), was most associated with the SBS  
214 10b signature (49%).

215         For the bulky adduct group, the mutation spectrum of BaP revealed mutational  
216 processes characteristic of SBS 4 (36%) and SBS 39 (27%) signatures. SBS 4 is most notably  
217 associated with tobacco-smoke induced cancer<sup>16</sup>, while SBS 39 is one of the new signatures  
218 that currently does not have a proposed etiology. No SBS signature was associated with the  
219 mutation profile of NDBzA and the control signature explained 53% of the mutation profile of this  
220 agent.

221 Analysis of the alkylating agent exposure group revealed that SBS 11 and SBS 30  
222 signatures were associated with N-nitrosodimethylamine (NDMA) mutation data<sup>34</sup> and explained  
223 37% and 50% of the mutations, respectively. SBS 11 has previously been linked to exposures  
224 to the methylating agents temozolomide and N-methyl-N'-nitro-N-nitrosoguanidine<sup>17,19</sup>. SBS 30  
225 is hypothesized to be associated with defects in base excision repair. No SBS signatures were  
226 associated with the mutation profiles of ENU or PRC while the control signature explained 36%  
227 and 42% of the mutation profiles of these two chemicals, respectively.

228 There were limited data available for nitrogenous base analogs. Data were only  
229 obtained from mice exposed to 5-(2-chloroethyl)-2-deoxyuridine (CEDU)<sup>41</sup>, a uridine analog.  
230 This included only 14 characterized mutants from bone marrow, 13 of which were T>C  
231 mutations. Consequently, 80% of the data were explained by the SBS 26 signature, which  
232 exhibits a bias for these types of substitutions.

233 TEM had a SNV mutation spectrum that was similar to controls (Figure 2). Nevertheless,  
234 we found that SBS 40 contributed to 32% of the TEM data, which was higher than the 20% that  
235 can be attributed to the control signature. There is currently no known etiology for the SBS 40  
236 signature.

237 Overall, the reconstructed signatures had very strong Pearson's coefficients (0.84 - 0.98)  
238 for six of the agents and strong coefficients (0.63 - 0.74) for four agents with the respective *lacZ*-  
239 generated mutation profiles (Figure 6 and Supplementary Figure S4).

240

## 241 DISCUSSION

242 We show that *in vivo* NGS-TGR data can be used to extract mutagenic mechanisms that  
243 may contribute to human cancers through application of COSMIC signature analysis. We also  
244 show that such analyses are improved through the inclusion of a background mutational  
245 signature (i.e., control signature) that reflects spontaneous mutations resulting from endogenous  
246 processes. Analysis of induced mutations in mouse tissues following exposures to 10 mutagenic  
247 agents (two sequenced by NGS, six sequenced by the Sanger method, and two by both)  
248 revealed high concordance between the expected mutagenic mode of action and the relevant  
249 COSMIC signature. The data suggest that our approach can be used to: (i) test if TGR mutation  
250 spectra support hypotheses that COSMIC signatures are attributed to particular mutagenic  
251 exposures, and (ii) generate hypotheses about the mutagenic mechanisms underlying human  
252 cancers through identifying enriched COSMIC signatures in TGR mutation spectra.

253 A large portion of mutations collected from weak mutagens are spontaneous rather than  
254 chemically induced. Thus, we developed a background signature derived from our empirical  
255 control data that can be integrated with COSMIC signatures to reduce the noise attributable to  
256 spontaneous mutation patterns. This *in vivo* control signature is a unique feature of our study,  
257 as there is currently no 'background' COSMIC signature and no *in vivo* control signature is  
258 reported in a recent study that generated chemical-specific signatures *in vivo* using a different  
259 approach<sup>27</sup>. Our results show that C>T transitions<sup>27</sup> are the most common spontaneous mutations  
260 *in vivo* (Figure 5) and this was consistent among all tissues analyzed (Supplementary Figure 6).  
261 C>T transitions at CpG sites are known hotspots of mutation due to spontaneous deamination  
262 of cytosine<sup>49</sup>. Previous work using bisulfite sequencing has shown that CpG sites in *lacZ* are  
263 heavily methylated, and CpG flanked by a 5' pyrimidine were most likely to have C>T base  
264 substitutions<sup>46</sup>. This is supported by our control data: the most prevalent spontaneous mutations  
265 were C>T at CCG, and, the third most prevalent were C>T mutations at TCG (Figure 5). Thus,  
266 our background control signature is consistent with expectations.

267 An *in vitro* background signature was recently reported<sup>28</sup>; however, the correlation  
268 between the two control signatures is modest (Pearson's coefficient = 0.45) because, at  
269 variance with our results, the *in vitro* control signature is enriched for C>A mutations.  
270 Spontaneous deamination of cytosine is also the most likely reason for C>A transversions and  
271 seem to be the most common spontaneous mutation *in vitro*<sup>53</sup>. This suggests that the same type  
272 of event, i.e., cytosine deamination, can result in different outcomes, i.e., C>T versus C>A  
273 mutations, depending on the physiological context.

274 Application of the control signature (Figure 5) and stringent statistical analysis identified  
275 nine SBS signatures that were associated with the *lacZ* SNVs induced by the investigated  
276 exposures. Two major outcomes from this analysis are: 1) mutation profiles for some of the  
277 tested agents were highly enriched for COSMIC signatures from cancers where the agents are  
278 known etiological factors (e.g., UV for skin cancer and BaP for tobacco-related cancers); and, 2)  
279 a few *lacZ* mutation profiles were associated with a variety of signatures of unknown aetiologies.  
280 This raises the question of whether the mutagenic mechanisms of these prototype agents are  
281 determinants of the signatures.

282 We identified the SBS 2, SBS 7a, and SBS 10b signatures as important contributors to  
283 the mutagenic mechanisms of all three electromagnetic radiation agents investigated (i.e., X-  
284 ray, UVB, and sunlight). SBS 2 has been observed in ~14% of cancer samples and is present in  
285 22 cancer types but is most often found in cervical and bladder cancers<sup>14,17</sup>. In this study, the  
286 signature was most strongly associated with UV skin exposure, representing 33-27% of  
287 mutations in exposed animals. Mechanistically, cytosine deamination is accelerated by UV  
288 exposure<sup>54</sup>; thus, it is possible that we observed SBS 2 in this study because of UV-dependent  
289 cytosine deamination. However, SBS 2 is not observed in skin cancers<sup>17</sup>. This suggests that  
290 mutations arising from UV-dependent cytosine deamination are not the primary drivers of the  
291 surveyed human skin cancers in the COSMIC database, and that other lesions (e.g., various  
292 types of photodimers) are the main contributors to the mutation catalogue of UV-induced skin

293 cancers. Another possible explanation is that with a small sample size of mutations, the high  
294 degree of similarity in the SBS 2 and SBS 7a signatures confounds this analysis. By this logic,  
295 some portion of the mutational signature identified as SBS 2 in our study may be the result of  
296 the mutational processes associated with SBS 7a, which is found in multiple cancer types but is  
297 most pronounced in skin cancers<sup>14,17</sup>. Indeed, the SBS 7a signature contributes to 27% of the  
298 mutations observed after sunlight exposure.

299       Activation of error-prone polymerases has been attributed to SBS 10b<sup>14</sup>, a signature that  
300 is mostly found in colorectal and uterine cancers. In the present study, this signature was only  
301 associated with X-ray mutations (49%). X-ray mutations show a high proportion of C>T  
302 substitutions at the TCG motif (Supplementary Figure S4), which is characteristic of the *lacZ*  
303 normalized SBS 10b signature (Supplementary Figure S3). It is possible that there is an ionizing  
304 radiation component to this signature. However, given previous work in this area, it is more  
305 likely that the association between SBS 10b and X-ray SNVs is a result of error-prone  
306 replication occurring in response to DNA damage.

307       The analysis of mutational signatures for the electromagnetic radiation agents provide  
308 support for the ability of the expanded repertoire of COSMIC signatures to exploit subtle  
309 differences in the mutation profiles to extract different mutational mechanisms. Using the  
310 previous version of the COSMIC database, all three radiation types had a comparable  
311 contribution from signature 7 (21-33%; data not shown). However, there are now four SBS  
312 signatures (7a-7d) derived from the original signature 7 in the latest COSMIC database<sup>15</sup>, and of  
313 these, only the SBS 7a signature contributes significantly to the mutation profile of sunlight.

314       Tobacco smoking is strongly associated with SBS 4, and this signature is commonly  
315 found in the lung tumors of smokers. BaP is a major mutagenic component in tobacco smoke<sup>21</sup>  
316 and as expected, SBS 4 contributed the highest percentage (36%) to the mutation profile of  
317 BaP. Interestingly, SBS4 was the only signature that contributed to the mutation profile of BaP  
318 and accounted for 60% of the observed profile when using the previous version of the COSMIC

319 database (data not shown). However, using version 3 of the COSMIC database<sup>15</sup>, the  
320 contribution of SBS 4 declined while we identified a second signature that contributed to the  
321 BaP mutation profile. Specifically, we detected a significant contribution (27%) of the SBS 39  
322 signature, which is one of the new signatures and currently has no known etiology. These  
323 results suggest an overlap in the aetiology of these two signatures and that SBS 39 may be  
324 associated with exposure to chemicals that induce bulky adducts.

325 The BaP mutation profile that we derived using our approach is consistent with previous  
326 work *in vivo*<sup>27</sup> and *in vitro*<sup>28</sup> that demonstrated the presence of SBS 4 after exposure to BaP.  
327 Indeed, the BaP mutation profile is consistent among the three studies (Pearson's correlation of  
328 0.80 and 0.71 with the *in vivo* and *in vitro* profile, respectively). Remarkably, signatures SBS 24,  
329 which has been associated with aflatoxin adducts, and SBS 29, which has been associated with  
330 tobacco chewing, are strikingly similar to SBS 4 (Supplementary Figure S3). However, only SBS  
331 4 strongly correlates with the BaP mutation data. This demonstrates the robustness of the  
332 mutational signatures and the ability of TGR-NGS to correctly discriminate between similar  
333 signatures that have different aetiologies. It also emphasizes the importance of the flanking  
334 nucleotides to increasing the specificity of the signatures; this work demonstrates that 96-bp  
335 signatures provide superior mechanistic information to standard mutation spectrum analysis.

336 NDMA was the only alkylating agent among those investigated that was associated with  
337 an established COSMIC signature. About 50% of the NDMA mutation profile was explained by  
338 the SBS 30 signature that has been associated with a deficiency in base excision repair. NDMA  
339 is known to induce mostly O<sup>6</sup>- and N<sup>7</sup>-methyl guanine adducts<sup>34</sup>, thus, a role of base excision  
340 repair in the response to this chemical is expected. NDMA exposure was also enriched for SBS  
341 11 (37%), inducing primarily C>T mutations at CpC motifs (Supplementary Figure S4). SBS 11  
342 has been detected in melanomas and glioblastomas, and the mutation pattern of this signature  
343 has been attributed to alkylating agent exposures, such as temozolomide and N-methyl-N'-nitro-  
344 N-nitrosoguanidine<sup>17,19</sup>. These alkylating agents induce C>T mutations, mostly at CpC motifs,

345 and mutations at this motif are the four most common in the SBS 11 signature. The TGR  
346 mutation data from our study are consistent with this expected mutation spectrum.

347         The SBS 11 signature was not enriched within the mutation spectrum of the two other  
348 alkylating agents (i.e., ENU or PRC) in our mutation database. This is expected because these  
349 compounds induce a very different mutation spectrum, causing primarily T>A mutations. These  
350 differences demonstrate that SBS 11 is specific to a particular mechanism of alkylation (i.e.,  
351 target sites for the alkylation events) and that there is currently no COSMIC signature for  
352 alkylating agents that target thymine. Further TGR-NGS analyses of alkylating agents may  
353 refine our understanding regarding which specific alkylating agents or defective  
354 alkyltransferases underlie the mechanisms associated with SBS 11.

355         The mutation profile obtained with ENU, demonstrating a slight preponderance of T>A  
356 mutations over T>C mutations, is consistent (Pearson's coefficient = 0.71) with that obtained in  
357 the bone marrow of *gpt* delta mice<sup>27</sup>, although the correlation is reduced when expanding the six  
358 possible base pair alterations to the 96 possible mutation types (Pearson's coefficient = 0.49).  
359 This is mostly due to a deficiency of T>C mutations at CTN motifs with respect to *gpt* delta mice.  
360 Nevertheless, the similarity with the ENU mutation profile from *gpt* delta mice is greater than  
361 that obtained *in vitro* with an induced pluripotent stem cell (iPSC) line (Pearson's coefficient =  
362 0.30) where the ENU signature is dominated by T>C mutations<sup>28</sup>. These authors speculate that  
363 the preponderance of T>C mutations after *in vitro* exposure to ENU is driven by the intrinsic  
364 characteristics of DNA repair processes in iPSCs.

365         The mutation profile of CEDU, a nitrogenous base analog, closely matched the SBS 26  
366 signature, which contributed to 80% of mutations in exposed animals. This is the highest  
367 contribution of a COSMIC signature to any of the mutation profiles generated in this study. SBS  
368 26 is one of the seven SBS signatures associated with defective mismatch repair, which is one  
369 of the major repair pathways that deals with base analogs<sup>55</sup>. Due to the limited number of  
370 mutations recovered in the CEDU study, the association between SBS 26 and CEDU should be

371 further tested. Also, considering that CEDU is similar in structure to existing halogenated uracil  
372 analogs that serve as therapeutics (e.g., fluorouracil), attention should be given to these  
373 compounds as possible contributors to the SBS 26 signature and associated cancers.

374         Among the agents tested in this study, TEM is the only one that is more effective at  
375 inducing chromosomal structural aberrations than mutations. TEM is a trifunctional alkylating  
376 agent that induced a strong micronucleus response while eliciting a weak mutagenic response  
377 in the hematopoietic system<sup>48</sup>. Our analysis identified SBS 40 signature as a strong contributor  
378 (32%) to the mutation profile of TEM. SBS 40 is one of those signatures that is not dominated  
379 by any specific type of base pair alteration and does not have a proposed aetiology. Further  
380 studies are needed to confirm whether SBS 40 signature is an indicator of a clastogenic mode  
381 of action.

382         In summary, we demonstrated that *lacZ* transgene sequence data can be used, in  
383 conjunction with established mutation signatures derived from COSMIC cancer data sets, to test  
384 the hypothesis that a given class of mutagenic agents is linked with specific human cancers.  
385 Moreover, COSMIC signature mining based on TGR mutation data sets can be used to  
386 generate new hypotheses regarding the mutagenic mechanisms associated with human  
387 cancers. This study presents a potential avenue through which mutation signature analysis can  
388 be applied to *in vivo* experimental models, and the analyses employed to improve  
389 understanding of mode of action. The analyses can also generate hypotheses regarding the  
390 mutational mechanisms of uncharacterized chemicals. The *in vivo* TGR-NGS approach has  
391 comparable sensitivity to whole genome approaches used for investigating the mutational  
392 landscape of environmental agents<sup>18, 19, 26, 28, 56</sup>. However, by avoiding the orders-of-magnitude  
393 higher cost of whole-genome sequencing, the *in vivo* TGR-NGS approach offers much higher-  
394 throughput for the testing of chemical mutagens. Overall, these results highlight that some  
395 mutational signatures may have large environmental components and contribute to the growing



396 body of evidence that analyses of mutations spectra shortly after exposure has bearing on the  
397 carcinogenic mechanism and the mutational profile observed in fully developed cancers.  
398

## 399 **MATERIALS AND METHODS**

### 400 **Animal Treatment**

401 Male MutaMouse animals (6-15 weeks old; 6-8 per group) were exposed daily to either  
402 100 mg/kg BaP, 5 mg/kg ENU, 25 mg/kg PRC or 2 mg/kg TEM by oral gavage for 28 days as  
403 per the Organisation for Economic Co-operation and Development (OECD) test guideline 488<sup>57</sup>.  
404 All doses were selected based on pilot studies conducted to identify the maximum tolerated  
405 dose as per TG 488 guidance. The BaP<sup>8</sup>, PRC<sup>47</sup> and TEM<sup>48</sup> data are the same presented their  
406 respective reference. Matched controls received the solvent (olive oil or water) by oral gavage  
407 during the same period. Three days after the last daily exposure, mice were anaesthetized with  
408 isoflurane and euthanized via cervical dislocation. Bone marrow cells were isolated by flushing  
409 femurs with 1X phosphate-buffered saline. After brief centrifugation, the supernatant was  
410 discarded, and the pellet was flash-frozen in liquid nitrogen prior to storage at -80 °C. All animal  
411 procedures were carried out under conditions approved by the Health Canada Ottawa Animal  
412 Care Committee.

413

### 414 ***lacZ* Mutant Quantification, Collection, and Sequencing**

415 The experimental protocol for enumerating *lacZ* mutants followed OECD guideline 488<sup>57</sup>.  
416 Briefly, bone marrow was thawed and digested overnight with gentle shaking at 37 °C in 5 mL of  
417 lysis buffer (10 mM Tris-HCl, pH 7.6, 10 mM ethylenediaminetetraacetic acid (EDTA), 100 mM  
418 NaCl, 1 % sodium dodecyl sulfate (w/v), 1 mg/mL Proteinase K). High molecular weight  
419 genomic DNA was isolated using phenol/chloroform extraction as described previously<sup>42,58</sup>. The  
420 isolated DNA was dissolved in 100 µL of TE buffer (10 mM Tris pH 7.6, 1 mM EDTA) and stored  
421 at 4 °C for several days before use. The phenyl-β-D-galactopyranoside (P-gal) positive selection  
422 assay<sup>59</sup> was used to identify *lacZ* mutants present in the DNA. Briefly, the *λgt10lacZ* construct  
423 present in the genomic DNA was isolated and packaged into phage particles using the

424 Transpack™ lambda packaging system (Agilent, Mississauga, Ontario, Canada). The phages  
425 were then mixed with *E. coli* (*lacZ*<sup>-</sup>, *galE*<sup>-</sup>, *recA*<sup>-</sup>, *pAA119*<sup>-</sup> with *galT* and *galk*)<sup>58</sup> in order to  
426 transfect the cells with the *lacZ* construct. *E. coli* were then plated on a selective media  
427 containing 0.3% P-gal (w/v) and incubated overnight at 37 °C. Only *E. coli* receiving a mutant  
428 copy of *lacZ* where the gene function is disrupted can form plaques on the P-gal medium,  
429 because P-gal is toxic to *galE*<sup>-</sup> strains with a functional *lacZ* gene product<sup>1</sup>. Packaged phage  
430 particles were concurrently plated on plates without P-gal (titre plates) to quantify the total  
431 plaque-forming units to be used as the denominator in the mutant frequency calculation.

432         After enumeration, plaques from each individual sample were collected and pooled  
433 together in microtubes containing autoclaved milliQ water (0.3 plaques/μL; mutants from 1  
434 sample per tube). Mutant amplification and sequencing were done as described previously<sup>8</sup>.  
435 Briefly, the mutant pools were boiled for 5 minutes and transferred to a PCR mastermix  
436 containing a final concentration of 1X Q5 reaction buffer, 200 μM dNTPs, 0.5 μM Forward  
437 primer (GGCTTTACACTTTATGCTTC), 0.5 μM Reverse Primer  
438 (ACATAATGGATTTCTTACG), and 1U Q5 enzyme (New England BioLabs Ltd., Whitby,  
439 Ontario, Canada); the final volume of each PCR was 50 μL. To control for errors introduced  
440 during PCR, each mutant pool was amplified twice as two separate technical replicates. The  
441 following thermocycle program was used for amplification: 95 °C for 3 min; 30 cycles of 95 °C  
442 for 45 s, 50 °C for 1 min, 72 °C for 4 min; final extension at 72 °C for 7 min. PCR products were  
443 purified using the QIAquick PCR purification kit (Qiagen, Montreal, Quebec, Canada).

444         NGS libraries were built using the NEBNext® Fast DNA Library Prep Set for Ion  
445 Torrent™. Each technical replicate had a unique barcoded adaptor ligated to the *lacZ* DNA  
446 fragments allowing for many samples to be sequenced simultaneously (up to 96 libraries per  
447 NGS run). Sequencing was performed using the Ion Chef™ workflow and Ion Proton™ system  
448 with P1 chips. NGS reads were aligned to the *lacZ* gene using bowtie 2 (version 2.1.0) and read  
449 depths for every possible mutation were quantified using samtools (version 0.1.19). Mutations

450 were called if, after background correction (determined by sequencing non-mutants), both  
451 technical replicates had mutation read depths above threshold values (equal to at least  
452  $1/\text{number of plaques in pool}$ )<sup>8</sup>. To further filter the data in this study, if the mutation read depths  
453 between two technical replicates varied by  $\geq 50\%$  then that mutation was removed from analysis.  
454 Clonally expanded mutants were only counted as one mutation.

455

## 456 **Published Sanger Sequencing Data**

457 Published data came from studies where *lacZ* transgene mutants were sequenced and  
458 the position and type of each mutation was reported (summarized in Supplementary Table S1).  
459 Mutants were characterized from MutaMouse or *LacZ* Plasmid mice<sup>60</sup>. Some studies reported  
460 the position of the mutation in the plasmid construct, while others reported the position in the  
461 coding sequence. For consistency, the positional information was adjusted to reflect the position  
462 of the mutation in the coding sequence of the *lacZ* gene. Furthermore, the reference sequence  
463 of *lacZ* used for NGS has four variations<sup>38</sup> relative to the *E. coli lacZ* coding sequence  
464 (Genbank: V00296.1)<sup>61</sup>, including a 15 bp insertion into codon 8. Thus, mutation positions were  
465 also adjusted to reflect this where applicable (e.g., if *LacZ* Plasmid mice were used instead of  
466 MutaMouse). No mutations were detected at or next to the variant positions in the *LacZ* Plasmid  
467 motif. In contrast to NGS work, different tissues were used for these analyses (i.e., bone  
468 marrow, brain, colon, germ cells, kidney, liver, skin, spleen, and stomach). Tissue sources are  
469 noted in the results with the accompanying data.

470

## 471 **Signature Analyses**

472 The workflow used to do signature analyses are available as an RShiny web-application  
473 ([https://github.com/MarcBeal/HC-MSD/tree/master/lacZ\\_Mutations\\_COSMIC\\_Signatures](https://github.com/MarcBeal/HC-MSD/tree/master/lacZ_Mutations_COSMIC_Signatures)).  
474 Mutations for control and exposed samples (see metadata in Supplemental Material) were

475 imported into the R console<sup>62</sup> as VRanges using the package “VariantAnnotation”<sup>63</sup> with the *lacZ*  
476 coding sequence as the reference FASTA file. To determine which of the COSMIC mutation  
477 signatures best explained the observed *lacZ* mutant spectrum, the COSMIC mutation signature  
478 weights, which are derived from human mutation data, were first normalized to *lacZ* trinucleotide  
479 frequencies. This was done using the ratio of trinucleotide frequencies in *lacZ* to the  
480 trinucleotide frequencies in the human genome (Figure 3; the normalized signatures are shown  
481 in Supplementary Figure S3 and the raw numbers in Supplementary Material). Analysis was  
482 done this way (as opposed to converting *lacZ* mutation data themselves to human trinucleotide  
483 frequencies) because the COSMIC signature are based on a much larger database, and  
484 therefore, represent a more robust signal with less variance. Following normalization, each of  
485 the 96 trinucleotide substitutions within each signature were represented as the relative  
486 frequency (i.e., all values in a signature sum to 1) by dividing each normalized value by the sum  
487 of all values for that signature. The trinucleotide mutation context (i.e., the nucleotide  
488 immediately upstream and downstream of the mutation) was obtained with the  
489 “mutationContext” function and converted to a motif matrix using the “motifMatrix” function (both  
490 in the “SomaticSignatures” package<sup>64</sup>). The motif matrix was then transposed to obtain the  
491 required format, and finally decomposed into the constituent *lacZ*-normalized signatures using  
492 the “whichSignatures” function from “deconstructSigs”<sup>65</sup>. The contribution of each identified  
493 signature to the mutation data was reported as a fraction. If the sum of each signature did not  
494 account for 100% of the mutation data, then the remainder was reported as the “residual”.

495       In order to account for spontaneous mutations often present alongside induced  
496 mutations, which is especially true for weak mutagens, we generated a signature for the  
497 spontaneous mutation background using the mutations observed in control animals. This  
498 included all control mutations characterized by NGS and Sanger sequencing. However,  
499 spontaneous SNVs characterized by Sanger sequencing were heavily biased towards positions  
500 1072, 1090, 1187, 1627, and 2374. Therefore, Sanger sequencing data at these 5 positions

501 were not used for deriving the control mutation signature. Signatures were plotted using  
502 ggplot2<sup>66</sup>.

503 “Signature reconstruction” was then used to determine how well the combination of  
504 normalized signatures, identified using the “whichSignatures” function, explain the mutation data  
505 from the respective exposure groups. For example, if signatures 3 and 4 contributed 40% and  
506 60% to the mutation profile of a compound, respectively, then the motif matrices for signatures 3  
507 and 4 were multiplied by 0.4 and 0.6, respectively, and summed together. The reconstructed  
508 signature was then compared against the motif matrices of the compound using Pearson  
509 correlation.

510 Lastly, the contribution of individual signatures was further validated using Pearson  
511 correlation. Specifically, each signature was compared against the respective 96-base context  
512 mutation spectra from which the signature was identified. In the final results, COSMIC  
513 signatures were only reported if the contribution was greater than the largest residual, and the  
514 Pearson coefficient with the reconstructed signature was greater than 0.5.

515

## 516 **Statistics**

517 Statistical analyses were done using the R programming language<sup>62</sup>. Mutant frequencies  
518 were compared between exposure groups and controls using generalized estimating equations  
519 assuming a Poisson distribution for the error, as done previously<sup>8</sup>, using the geepack library<sup>67</sup>  
520 with outliers (1 in control, 1 in TEM) removed. Bonferonni correction for multiple comparisons  
521 was used to adjust the threshold of significance. Mutation spectra of the chemical exposure  
522 groups were compared against controls using mutation proportions. The standard error for the  
523 mutation spectra was determined using error propagation. Significant differences in mutation  
524 spectra between chemically induced mutants and spontaneous control mutants were  
525 determined using Fisher’s exact tests with Bonferonni correction for multiple comparisons (i.e.,  
526 across different chemical groups). To compare whole mutation spectra between control and

527 exposed groups, Fisher's exact tests were performed with Monte Carlo simulation with 10,000  
528 replicates. Fisher's exact tests were also performed on 2 × 2 sub-tables for each mutation type.  
529

## 530 **FUNDING**

531 Funding for this research was provided for by Health Canada's Chemicals Management  
532 Plan and Genomics Research and Development Initiative.  
533

## 534 **ACKNOWLEDGEMENTS**

535 We would like to thank Angela Dykes, Lynda Soper, and John Gingerich for their  
536 technical contributions to this research. We are grateful for advice provided by Dr. Ludmil  
537 Alexandrov and Andrew Williams.

538

## 539 **COMPETING FINANCIAL INTERESTS**

540 The authors declare that there are no competing financial and non-financial interests.

541

## 542 **AUTHOR CONTRIBUTIONS**

543 MAB, CM, MJM and JOB conducted the MutaMouse animal studies and collected  
544 samples. MAB and MJM sequenced plaques. MAB, MJM and DL conducted the COSMIC  
545 analyses. CY and FM secured funding for the study and were responsible for study conception  
546 and design. All authors contributed to data analysis, interpretation, paper writing and approved  
547 the final version of the manuscript.

548

549

## 550 References

- 551 1. Lambert IB, Singer TM, Boucher SE, Douglas GR. Detailed review of transgenic rodent  
552 mutation assays. *Mutat Res* **590**, 1-280 (2005).  
553
- 554 2. OECD. *Detailed Review Paper on Transgenic Rodent Mutations Assay*, Paris (2009).  
555
- 556 3. Meier MJ, Beal MA, Schoenrock A, Yauk CL, Marchetti F. Whole genome sequencing of  
557 the Mutamouse model reveals strain- and colony-level variation, and genomic features  
558 of the transgene integration site. *Sci Rep* **9**, 13775 (2019).  
559
- 560 4. Shwed PS, Crosthwait J, Douglas GR, Seligy VL. Characterisation of MutaMouse  
561 lambda<sub>dagt10</sub>-lacZ transgene: evidence for in vivo rearrangements. *Mutagenesis* **25**, 609-  
562 616 (2010).  
563
- 564 5. Gingerich JD, Soper L, Lemieux CL, Marchetti F, Douglas GR. *Transgenic Rodent*  
565 *Gene Mutation Assay in Somatic Tissues*. Springer Science+Business Media (2014).  
566
- 567 6. O'Brien JM, *et al.* Transgenic rodent assay for quantifying male germ cell mutant  
568 frequency. *J Vis Exp*, e51576 (2014).  
569
- 570 7. Besaratinia A, Li H, Yoon JI, Zheng A, Gao H, Tommasi S. A high-throughput next-  
571 generation sequencing-based method for detecting the mutational fingerprint of  
572 carcinogens. *Nucleic Acids Res* **40**, e116 (2012).  
573
- 574 8. Beal MA, Gagne R, Williams A, Marchetti F, Yauk CL. Characterizing benzo[a]pyrene-  
575 induced lacZ mutation spectrum in transgenic mice using next-generation sequencing.  
576 *BMC Genomics* **16**, 812 (2015).  
577
- 578 9. Meier MJ, O'Brien JM, Beal MA, Allan B, Yauk CL, Marchetti F. In utero exposure to  
579 benzo[a]pyrene increases mutation burden in the soma and sperm of adult mice.  
580 *Environ Health Perspect* **125**, 82-88 (2017).  
581
- 582 10. Alexandrov LB, Nik-Zainal S, Wedge DC, Campbell PJ, Stratton MR. Deciphering  
583 signatures of mutational processes operative in human cancer. *Cell Rep* **3**, 246-259  
584 (2013).  
585
- 586 11. Nik-Zainal S, *et al.* Mutational processes molding the genomes of 21 breast cancers.  
587 *Cell* **149**, 979-993 (2012).  
588
- 589 12. Nik-Zainal S, *et al.* The life history of 21 breast cancers. *Cell* **149**, 994-1007 (2012).  
590
- 591 13. Forbes SA, *et al.* COSMIC: somatic cancer genetics at high-resolution. *Nucleic Acids*  
592 *Res* **45**, D777-D783 (2017).  
593
- 594 14. Helleday T, Eshtad S, Nik-Zainal S. Mechanisms underlying mutational signatures in  
595 human cancers. *Nat Rev Genet* **15**, 585-598 (2014).  
596
- 597 15. Alexandrov LB, *et al.* The Repertoire of Mutational Signatures in Human Cancer.  
598 *BioRxiv*, (2018).  
599



- 600 16. Alexandrov LB, *et al.* Mutational signatures associated with tobacco smoking in human  
601 cancer. *Science* **354**, 618-622 (2016).  
602
- 603 17. Alexandrov LB, *et al.* Signatures of mutational processes in human cancer. *Nature* **500**,  
604 415-421 (2013).  
605
- 606 18. Nik-Zainal S, *et al.* The genome as a record of environmental exposure. *Mutagenesis*  
607 **30**, 763-770 (2015).  
608
- 609 19. Olivier M, *et al.* Modelling mutational landscapes of human cancers in vitro. *Sci Rep* **4**,  
610 4482 (2014).  
611
- 612 20. Phillips DH. Mutational spectra and mutational signatures: Insights into cancer aetiology  
613 and mechanisms of DNA damage and repair. *DNA Repair (Amst)* **71**, 6-11 (2018).  
614
- 615 21. Pfeifer GP, Denissenko MF, Olivier M, Tretyakova N, Hecht SS, Hainaut P. Tobacco  
616 smoke carcinogens, DNA damage and p53 mutations in smoking-associated cancers.  
617 *Oncogene* **21**, 7435-7451 (2002).  
618
- 619 22. Alexandrov LB, *et al.* Clock-like mutational processes in human somatic cells. *Nat Genet*  
620 **47**, 1402-1407 (2015).  
621
- 622 23. Hollstein M, Alexandrov LB, Wild CP, Ardin M, Zavadil J. Base changes in tumour DNA  
623 have the power to reveal the causes and evolution of cancer. *Oncogene* **36**, 158-167  
624 (2017).  
625
- 626 24. Zhivagui M, Korenjak M, Zavadil J. Modelling mutation spectra of human carcinogens  
627 using experimental systems. *Basic Clin Pharmacol Toxicol* **121 Suppl 3**, 16-22 (2017).  
628
- 629 25. Chawanthayatham S, *et al.* Mutational spectra of aflatoxin B1 in vivo establish  
630 biomarkers of exposure for human hepatocellular carcinoma. *Proc Natl Acad Sci U S A*  
631 **114**, E3101-E3109 (2017).  
632
- 633 26. Huang MN, *et al.* Genome-scale mutational signatures of aflatoxin in cells, mice, and  
634 human tumors. *Genome Res* **27**, 1475-1486 (2017).  
635
- 636 27. Matsumura S, Sato H, Otsubo Y, Tasaki J, Ikeda N, Morita O. Genome-wide somatic  
637 mutation analysis via Hawk-Seq reveals mutation profiles associated with chemical  
638 mutagens. *Arch Toxicol* **93**, 2689-2701 (2019).  
639
- 640 28. Kucab JE, *et al.* A Compendium of Mutational Signatures of Environmental Agents. *Cell*  
641 **177**, 821-836 e816 (2019).  
642
- 643 29. Ikehata H, Masuda T, Sakata H, Ono T. Analysis of mutation spectra in UVB-exposed  
644 mouse skin epidermis and dermis: frequent occurrence of C-->T transition at methylated  
645 CpG-associated dipyrimidine sites. *Environ Mol Mutagen* **41**, 280-292 (2003).  
646
- 647 30. Ikehata H, Nakamura S, Asamura T, Ono T. Mutation spectrum in sunlight-exposed  
648 mouse skin epidermis: small but appreciable contribution of oxidative stress-mediated  
649 mutagenesis. *Mutat Res* **556**, 11-24 (2004).  
650

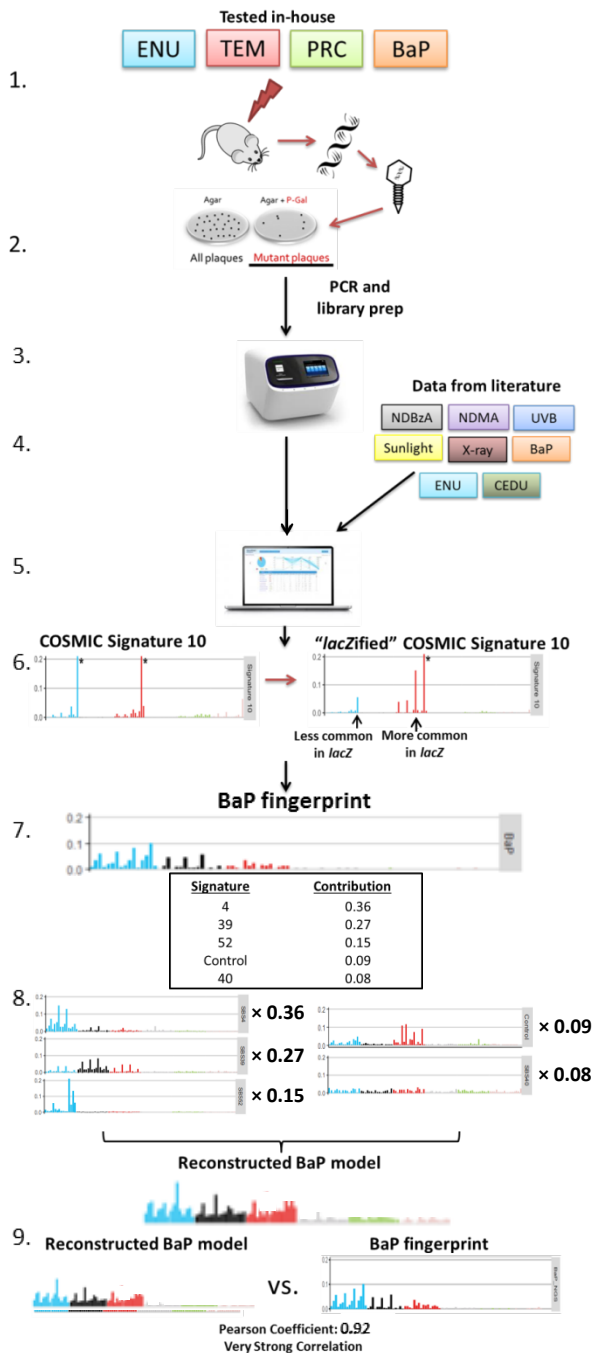
- 651 31. Frijhoff AF, *et al.* UVB-induced mutagenesis in hairless lambda lacZ-transgenic mice.  
652 *Environ Mol Mutagen* **29**, 136-142 (1997).  
653
- 654 32. Ono T, Ikehata H, Vishnu Priya P, Uehara Y. Molecular nature of mutations induced by  
655 irradiation with repeated low doses of X-rays in spleen, liver, brain and testis of lacZ-  
656 transgenic mice. *Int J Radiat Biol* **79**, 635-641 (2003).  
657
- 658 33. Ono T, *et al.* Molecular nature of mutations induced by a high dose of x-rays in spleen,  
659 liver, and brain of the lacZ-transgenic mouse. *Environ Mol Mutagen* **34**, 97-105 (1999).  
660
- 661 34. Souliotis VL, van Delft JH, Steenwinkel MJ, Baan RA, Kyrtopoulos SA. DNA adducts,  
662 mutant frequencies and mutation spectra in lambda lacZ transgenic mice treated with N-  
663 nitrosodimethylamine. *Carcinogenesis* **19**, 731-739 (1998).  
664
- 665 35. Suzuki T, *et al.* A comparison of the genotoxicity of ethylnitrosourea and ethyl  
666 methanesulfonate in lacZ transgenic mice (Muta Mouse). *Mutat Res* **395**, 75-82 (1997).  
667
- 668 36. Mientjes EJ, *et al.* DNA adducts, mutant frequencies, and mutation spectra in various  
669 organs of lambda lacZ mice exposed to ethylating agents. *Environ Mol Mutagen* **31**, 18-  
670 31 (1998).  
671
- 672 37. Jiao J, Douglas GR, Gingerich JD, Soper LM. Analysis of tissue-specific lacZ mutations  
673 induced by N-nitrosodibenzylamine in transgenic mice. *Carcinogenesis* **18**, 2239-2245  
674 (1997).  
675
- 676 38. Hakura A, Tsutsui Y, Sonoda J, Tsukidate K, Mikami T, Sagami F. Comparison of the  
677 mutational spectra of the lacZ transgene in four organs of the MutaMouse treated with  
678 benzo[a]pyrene: target organ specificity. *Mutat Res* **447**, 239-247 (2000).  
679
- 680 39. Douglas GR, Jiao J, Gingerich JD, Gossen JA, Soper LM. Temporal and molecular  
681 characteristics of mutations induced by ethylnitrosourea in germ cells isolated from  
682 seminiferous tubules and in spermatozoa of lacZ transgenic mice. *Proc Natl Acad Sci U*  
683 *S A* **92**, 7485-7489 (1995).  
684
- 685 40. Douglas GR, Jiao J, Gingerich JD, Soper LM, Gossen JA. Temporal and molecular  
686 characteristics of lacZ mutations in somatic tissues of transgenic mice. *Environ Mol*  
687 *Mutagen* **28**, 317-324 (1996).  
688
- 689 41. Staedtler F, Suter W, Martus HJ. Induction of A:T to G:C transition mutations by 5-(2-  
690 chloroethyl)-2'-deoxyuridine (CEDU), an antiviral pyrimidine nucleoside analogue, in the  
691 bone marrow of Muta Mouse. *Mutat Res* **568**, 211-220 (2004).  
692
- 693 42. Douglas GR, Gingerich JD, Gossen JA, Bartlett SA. Sequence spectra of spontaneous  
694 lacZ gene mutations in transgenic mouse somatic and germline tissues. *Mutagenesis* **9**,  
695 451-458 (1994).  
696
- 697 43. Dolle ME, Martus HJ, Novak M, van Orsouw NJ, Vijg J. Characterization of color  
698 mutants in lacZ plasmid-based transgenic mice, as detected by positive selection.  
699 *Mutagenesis* **14**, 287-293 (1999).  
700

- 701 44. Dolle ME, Snyder WK, Dunson DB, Vijg J. Mutational fingerprints of aging. *Nucleic Acids Res* **30**, 545-549 (2002).  
702  
703
- 704 45. Dolle ME, *et al.* Increased genomic instability is not a prerequisite for shortened lifespan  
705 in DNA repair deficient mice. *Mutat Res* **596**, 22-35 (2006).  
706
- 707 46. Ikehata H, Takatsu M, Saito Y, Ono T. Distribution of spontaneous CpG-associated G:C  
708 --> A:T mutations in the lacZ gene of Muta mice: effects of CpG methylation, the  
709 sequence context of CpG sites, and severity of mutations on the activity of the lacZ gene  
710 product. *Environ Mol Mutagen* **36**, 301-311 (2000).  
711
- 712 47. Maurice C, Dertinger SD, Yauk CL, Marchetti F. Integrated In Vivo Genotoxicity  
713 Assessment of Procarbazine Hydrochloride Demonstrates Induction of Pig-a and LacZ  
714 Mutations, and Micronuclei, in MutaMouse Hematopoietic Cells. *Environ Mol Mutagen*  
715 **60**, 505-512 (2019).  
716
- 717 48. Maurice C, O'Brien JM, Yauk CL, Marchetti F. Integration of sperm DNA damage  
718 assessment into OECD test guidelines for genotoxicity testing using the MutaMouse  
719 model. *Toxicol Appl Pharmacol* **357**, 10-18 (2018).  
720
- 721 49. Duret L. Mutation patterns in the human genome: more variable than expected. *PLoS Biol* **7**, e1000028 (2009).  
722  
723
- 724 50. Shelby MD, Tindall KR. Mammalian germ cell mutagenicity of ENU, IPMS and MMS,  
725 chemicals selected for a transgenic mouse collaborative study. *Mutat Res* **388**, 99-109  
726 (1997).  
727
- 728 51. Beranek DT. Distribution of methyl and ethyl adducts following alkylation with  
729 monofunctional alkylating agents. *Mutat Res* **231**, 11-30 (1990).  
730
- 731 52. Revollo J, *et al.* Spectrum of Pig-a mutations in T lymphocytes of rats treated with  
732 procarbazine. *Mutagenesis* **32**, 571-579 (2017).  
733
- 734 53. de Jong PJ, Grosovsky AJ, Glickman BW. Spectrum of spontaneous mutation at the  
735 APRT locus of Chinese hamster ovary cells: an analysis at the DNA sequence level.  
736 *Proc Natl Acad Sci U S A* **85**, 3499-3503 (1988).  
737
- 738 54. Peng W, Shaw BR. Accelerated deamination of cytosine residues in UV-induced  
739 cyclobutane pyrimidine dimers leads to CC-->TT transitions. *Biochemistry* **35**, 10172-  
740 10181 (1996).  
741
- 742 55. Kunkel TA. DNA-mismatch repair. The intricacies of eukaryotic spell-checking. *Curr Biol*  
743 **5**, 1091-1094 (1995).  
744
- 745 56. Meier B, *et al.* C. elegans whole-genome sequencing reveals mutational signatures  
746 related to carcinogens and DNA repair deficiency. *Genome Res* **24**, 1624-1636 (2014).  
747
- 748 57. OECD. *Test 488: Transgenic Rodent Somatic and Germ Cells Gene Mutation Assays*.  
749 OECD Publishing (2013).  
750

- 751 58. Gossen JA, Molijn AC, Douglas GR, Vijg J. Application of galactose-sensitive E. coli  
752 strains as selective hosts for LacZ- plasmids. *Nucleic Acids Res* **20**, 3254 (1992).  
753
- 754 59. Vijg J, Douglas GR. Bacteriophage lambda and plasmid lacZ transgenic mice for  
755 studying mutations in vivo. In: *Technologies for detection of DNA damage and mutations*  
756 (ed<sup>^</sup>(eds Pfeifer GP). Plenum Press (1996).  
757
- 758 60. Vijg J, Dolle ME, Martus HJ, Boerrigter ME. Transgenic mouse models for studying  
759 mutations in vivo: applications in aging research. *Mech Ageing Dev* **99**, 257-271 (1997).  
760
- 761 61. Kalnins A, Otto K, Ruther U, Muller-Hill B. Sequence of the lacZ gene of Escherichia coli.  
762 *EMBO J* **2**, 593-597 (1983).  
763
- 764 62. R Core Team. *R: a language and environment for statistical computing* (2016).  
765
- 766 63. Obenchain V, Lawrence M, Carey V, Gogarten S, Shannon P, Morgan M.  
767 VariantAnnotation: a Bioconductor package for exploration and annotation of genetic  
768 variants. *Bioinformatics* **30**, 2076-2078 (2014).  
769
- 770 64. Gehring JS, Fischer B, Lawrence M, Huber W. SomaticSignatures: inferring mutational  
771 signatures from single-nucleotide variants. *Bioinformatics* **31**, 3673-3675 (2015).  
772
- 773 65. Rosenthal R, McGranahan N, Herrero J, Taylor BS, Swanton C. DeconstructSigs:  
774 delineating mutational processes in single tumors distinguishes DNA repair deficiencies  
775 and patterns of carcinoma evolution. *Genome Biol* **17**, 31 (2016).  
776
- 777 66. Wickham H. *ggplot2: elegant graphics for data analysis*. Springer-Verlag (2016).  
778
- 779 67. Halekoh U, Højsgaard S, Yan J. The R package geepack for generalized estimating  
780 equations. *Journal of Statistical Software* **15**, 1-11 (2006).  
781  
782
- 783
- 784

785 **FIGURES**

786



1. Four chemicals were tested in-house against solvent controls using the TGR *in vivo* mutagenicity assay
2. Mutant plaques from controls and chemical-exposed mice were collected and pooled per individual
3. Mutant plaques were PCR amplified as 2 technical replicates, library prepped and sequenced on the Ion Proton Platform. SNVs were called and corrected for clonal expansion
4. Published Sanger sequencing data were compiled for 8 additional chemicals, plus controls, tested using the *lacZ* plasmid or MutaMouse mice
5. All sequencing data (Sanger and Ion Proton) were imported into the R console and trinucleotide mutation context were obtained using the "mutationContext" function
6. To compare human COSMIC signatures and *lacZ* mutation data, the COSMIC signatures were normalized to *lacZ* trinucleotide frequencies and each of the 96 trinucleotide substitutions were represented as relative frequency
7. The "deconstructSigs" package was used to identify COSMIC signatures that best describe the mutational fingerprint of chemical exposure
8. Identified signatures were multiplied by their respective contributions and used to reconstruct the mutational fingerprint of chemical exposure
9. Pearson correlation was used to determine how well reconstructed signature models compared to empirically-derived mutational fingerprint

787

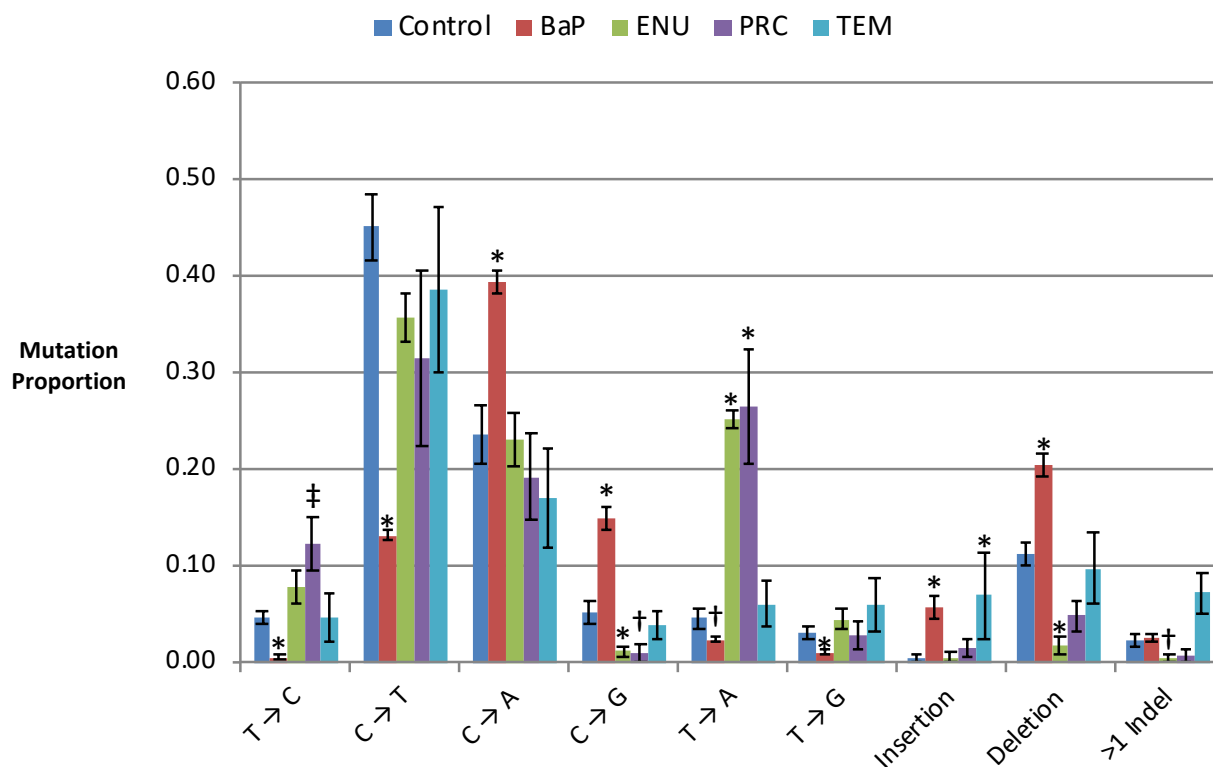
788

789 **Figure 1. Experimental design.** The experimental workflow included: animal exposure and

790 determination of mutant frequencies (Steps 1-2); sequencing of collected plaques and collection

791 of published *lacZ* sequenced data (Steps 3-4); generation of mutation profiles (steps 5-6); and  
792 query of the COSMIC database to identify mutational signatures that contributed to the mutation  
793 profile of tested agents (Steps 7-9).  
794

795



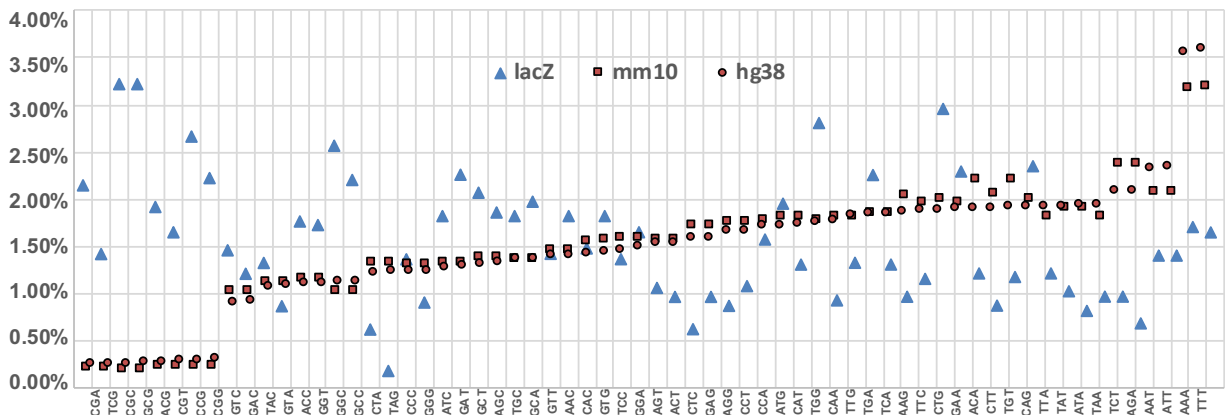
796

797

798 **Figure 2. Spontaneous and chemical-induced mutation proportions in bone marrow as**  
 799 **characterized by NGS.** BaP, shown in red has significantly higher proportions of C>A, C>G,  
 800 insertions, and deletions compared to control. In contrast, there is a lower proportion of T>C,  
 801 C>T, T>A, and T>G mutations than control. ENU, shown in green, has a higher proportion of  
 802 T>A mutations, while C>T, C>G, and deletions are lower. PRC, shown in purple, has a higher  
 803 proportion of T>A compared to control, and a marginally significant increase in T>C mutations  
 804 compared to control ( $P = 0.055$ ). The mutation spectrum for TEM, shown in turquoise, is most  
 805 similar to that of the control, with the exception of a significant increase in the proportion of  
 806 insertions. ‡  $P < 0.1$ , †  $P < 0.05$ , \*  $P < 0.0001$ .

807

808



809

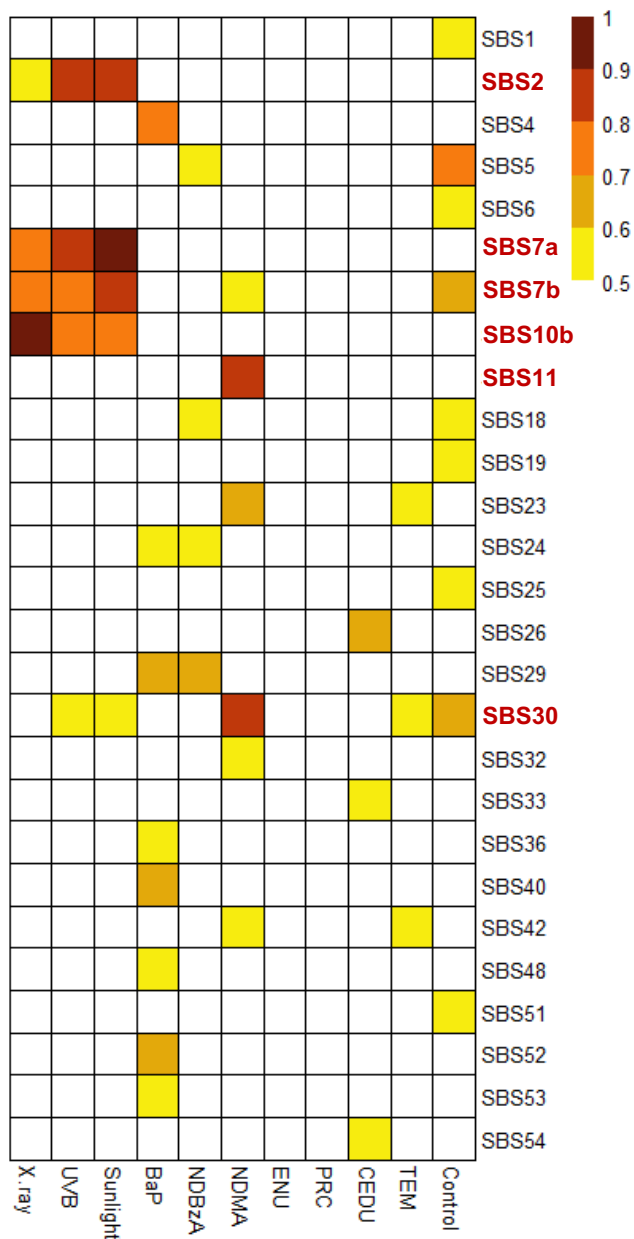
810

811 **Figure 3. Trinucleotide context differences between the *lacZ* transgene, mouse genome,**  
812 **and human genome.** Comparison of the frequencies of the 64 possible trinucleotides among  
813 the *lacZ* transgene (*lacZ*), mouse genome (mm10), and human genome (hg38) show that  
814 mouse and human genome frequencies are comparable with each other, while *lacZ* is more  
815 variable and biased towards some GC rich trinucleotides.

816



817



818

819

820 **Figure 4. Heatmap of similarities between obtained mutational profiles of tested agents**

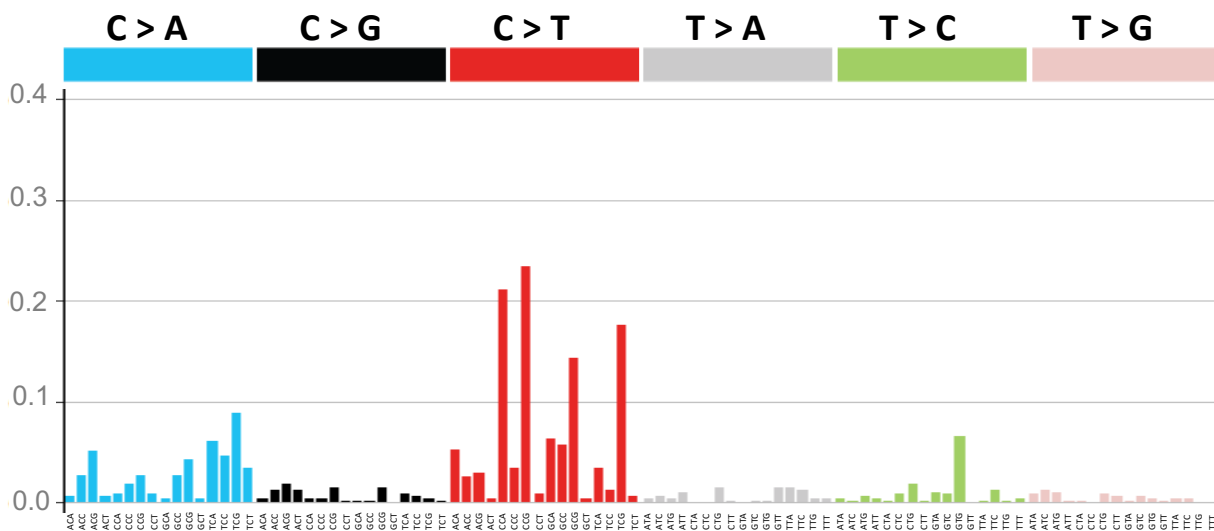
821 **and COSMIC SBS signatures.** All correlations that had a Pearson's correlation above 0.5 are

822 shown. The six SBS signatures that had a Pearson's correlation greater than 0.7 are indicated

823 in bold on the right of the heatmap.

824

## Control Signature



825

826

827 **Figure 5. The *lacZ* control signature.** The control signature is based on empirical mutation  
828 data from control animals in NGS and Sanger studies.

829

830

831

Signature	Electromagnetic radiation			Bulky adducts		Alkylating agents			Base analog	Clastogen
	X rays (35)	UVB (109)	Sunlight (62)	BaP (1165)	NDBzA (76)	NDMA (30)	ENU (613)	PRC (110)	CEDU (14)	TEM (115)
SBS 2		33	27							
SBS 4				36						
SBS 7a		7	27							
SBS 7b		15	14							
SBS 8					12		16	14		
SBS 10b	49	6								
SBS 11						37	8	17		
SBS 21	12									
SBS 24					14					
SBS 26									80	
SBS 30						50				
SBS 31										7
SBS 39				26						
SBS 40				8			12			32
SBS 42			7							19
SBS 51	6									
SBS 52				15			7			
SBS 54									19	
SBS 85							19	19		
Control	25	33	13	9	53		36	42		20
Residual	7	6	12	6	20	13	3	8	1	21
Pearson coef	0.97	0.93	0.98	0.92	0.84	0.89	0.66	0.63	0.67	0.74
Improvement	1.0	0	0	0	21.7	0	11.9	26.0	0	5.7

**Figure 6. The association of mutational signatures across different exposure groups.** The number below each agent indicates the number of unique mutants sequenced, while the number in each box represents the percent contribution of each signature to the

mutation profile of each tested agent. Brown rectangles indicate that the signature was present in the characterized mutations following the mutagenic exposure and that this association was moderate to strong (Pearson's coefficient  $> 0.5$ ; other signatures i.e., with a coefficient  $< 0.5$ ). The last two rows report the Pearson's coefficients between reconstructed signatures and observed mutation profiles and the percent increase due to inclusion of the control signature.