

Lawrence Berkeley National Laboratory

Recent Work

Title

Experimental and pan-cancer genome analyses reveal widespread contribution of acrylamide exposure to carcinogenesis in humans.

Permalink

<https://escholarship.org/uc/item/31b437c9>

Journal

Genome research, 29(4)

ISSN

1088-9051

Authors

Zhivagui, Maria
Ng, Alvin WT
Ardin, Maude
[et al.](#)

Publication Date

2019-04-01

DOI

10.1101/gr.242453.118

Peer reviewed

**Experimental analysis of exome-scale mutational signature
of glycidamide, the reactive metabolite of acrylamide**

Journal:	<i>Carcinogenesis</i>
Manuscript ID	CARCIN-2017-00290.R1
Manuscript Type:	Original Manuscript
Date Submitted by the Author:	n/a
Complete List of Authors:	<p>Zhivagui, Maria; International Agency for Research on Cancer, Molecular Mechanisms and Biomarkers Ardin, Maude; International Agency for Research on Cancer, Molecular Mechanisms and Biomarkers Ng, Alvin; Duke-NUS Graduate Medical School, Centre for Computational Biology Churchwell, Mona; National Center for Toxicological Research, Division of Biochemical Toxicology Pandey, Manuraj; International Agency for Research on Cancer, Molecular Mechanisms and Biomarkers Villar, Stephanie; International Agency for Research on Cancer, Molecular Mechanisms and Biomarkers Cahais, Vincent; International Agency for Research on Cancer, Epigenetics Robitaille, Alexis; International Agency for Research on Cancer, Infections and Cancer Biology Bouaoun, Liacine; International Agency for Research on Cancer, Environment and Radiation Heguy, Adriana; New York University-Langone Medical Center, Pathology and Genome Technology Center Guyton, Kate; International Agency for Research on Cancer, IARC Monographs Stampfer, Martha; Berkeley National Laboratory, Life Sciences Division McKay, James; International Agency for Research on Cancer, Genetic Cancer Susceptibility Hollstein, Monica; University of Leeds, Faculty of Medicine and Health, LIGHT Laboratories; Deutsches Krebsforschungszentrum, Department C016 Olivier, Magali; International Agency for Research on Cancer, Molecular Mechanisms and Biomarkers Rozen, Steve; Duke-NUS Graduate Medical School, Centre for Computational Biology Beland, Frederick; National Center for Toxicological Research, Division of Biochemical Toxicology Korenjak, Michael; International Agency for Research on Cancer, Molecular Mechanisms and Biomarkers Zavadil, Jiri; International Agency for Research on Cancer, Molecular Mechanisms and Biomarkers</p>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Keywords:	acrylamide, glycidamide, DNA adducts, massively parallel sequencing, mutational signatures

SCHOLARONE™
Manuscripts

For Peer Review

1
2
3
4
5 **1 Title**

6 **2 Experimental analysis of exome-scale mutational signature of glycidamide, the**
7 **3 reactive metabolite of acrylamide**

8
9
10 **4 Authors**

11 Maria Zhivagui¹, Maude Ardin¹, Alvin W. T. Ng^{2,3,4}, Mona I. Churchwell⁵, Manuraj Pandey¹,
12 Stephanie Villar¹, Vincent Cahais⁶, Alexis Robitaille⁷, Liacine Bouaoun⁸, Adriana Heguy⁹,
13 Kathryn Guyton¹⁰, Martha R. Stampfer¹¹, James McKay¹², Monica Hollstein^{1,13,14}, Magali
14 Olivier¹, Steven G. Rozen^{2,3}, Frederick A. Beland⁵, Michael Korenjak¹ and Jiri Zavadil¹

15
16
17 **9 Affiliations**

18 ¹ Molecular Mechanisms and Biomarkers Group, International Agency for Research on
19 Cancer, Lyon 69008, France

20 ² Centre for Computational Biology, Duke-NUS Medical School, Singapore 169857,
21 Singapore

22 ³ Program in Cancer and Stem Cell Biology, Duke-NUS Medical School, 169857, Singapore

23 ⁴ NUS Graduate School for Integrative Sciences and Engineering, 117456, Singapore

24 ⁵ Division of Biochemical Toxicology, National Center for Toxicological Research, Jefferson,
25 AR 72079, USA

26 ⁶ Epigenetics Group, International Agency for Research on Cancer, Lyon 69008, France

27 ⁷ Infections and Cancer Biology Group, International Agency for Research on Cancer, Lyon
28 69008, France

29 ⁸ Environment and Radiation Section, International Agency for Research on Cancer, Lyon
30 69008, France

31 ⁹ Department of Pathology and Genome Technology Center, New York University, Langone
32 Medical Center, New York, NY 10016, USA

33 ¹⁰ IARC Monographs Section, International Agency for Research on Cancer, Lyon 69008,
34 France

35 ¹¹ Biological Systems and Engineering Division, Lawrence Berkeley National Laboratory,
36 Berkeley, CA, 94720, USA

37 ¹² Genetic Cancer Susceptibility Group, International Agency for Research on Cancer, Lyon
38 69008, France

39 ¹³ Deutsches Krebsforschungszentrum, 69120 Heidelberg, Germany

40 ¹⁴ Faculty of Medicine and Health, University of Leeds, LIGHT Laboratories, Leeds LS2 9JT,
41 United Kingdom

42
43
44
45
46
47
48
49 **35 Keywords:** Acrylamide, glycidamide, DNA adducts, massively parallel sequencing,
50 **36 mutational signatures**

51
52
53 **37 Correspondence:**

54 **38 ZavadilJ@iarc.fr and/or KorenjakM@iarc.fr**

Abstract

Acrylamide, a probable human carcinogen, is ubiquitously present in the human environment, with sources including heated starchy foods, coffee and cigarette smoke. Humans are also exposed to acrylamide occupationally. Acrylamide is genotoxic, inducing gene mutations and chromosomal aberrations in various experimental settings. Covalent haemoglobin adducts were reported in acrylamide-exposed humans and DNA adducts in experimental systems. The carcinogenicity of acrylamide has been attributed to the effects of glycidamide, its reactive and mutagenic metabolite capable of inducing rodent tumors at various anatomical sites. In order to characterize the pre-mutagenic DNA lesions and global mutation spectra induced by acrylamide and glycidamide, we combined DNA-adduct and whole-exome sequencing analyses in an established exposure-clonal immortalization system based on mouse embryonic fibroblasts. Sequencing and computational analysis revealed a unique mutational signature of glycidamide, characterized by predominant T:A>A:T transversions, followed by T:A>C:G and C:G>A:T mutations exhibiting specific trinucleotide contexts and significant transcription strand bias. Computational interrogation of human cancer genome sequencing data indicated that a combination of the glycidamide signature and an experimental benzo[a]pyrene signature are nearly equivalent to the COSMIC tobacco-smoking related signature 4 in lung adenocarcinomas and squamous cell carcinomas. We found a more variable relationship between the glycidamide- and benzo[a]pyrene-signatures and COSMIC signature 4 in liver cancer, indicating more complex exposures in the liver. Our study demonstrates that the controlled experimental characterization of specific genetic damage associated with glycidamide exposure facilitates identifying corresponding patterns in cancer genome data, thereby underscoring how mutation signature laboratory experimentation contributes to the elucidation of cancer causation.

64

A 40-word summary

Innovative experimental approaches identify a novel mutational signature of glycidamide, a metabolite of the probable human carcinogen acrylamide. The results may elucidate the cancer risks associated with exposure to acrylamide, commonly found in tobacco smoke, thermally processed foods and beverages.

70 Introduction

71 Cancer can be caused by chemicals, complex mixtures, occupational exposures, physical
72 agents, and biological agents, as well as lifestyle factors. Many human carcinogens show a
73 number of characteristics that are shared among carcinogenic agents (1). Different human
74 carcinogens may exhibit a spectrum of these key characteristics, and operate through
75 separate mechanisms to generate patterns of genetic alterations. Recognizable patterns of
76 genetic alterations or mutational signatures characterize carcinogens that are genotoxic.
77 Recent work shows that these DNA sequence changes can be expressed in simple
78 mathematical terms that enable mutational signatures to be extracted from thousands of
79 cancer genome sequencing data sets (2). Several of the over 30 identified mutational
80 signatures have been attributed to specific external exposures or endogenous factors
81 through epidemiological and experimental studies (2). However, about 40% of the current
82 signatures remain of unknown origin, and additional, thus far unrecognized, signatures are
83 likely to be defined in rapidly accumulating cancer genome data. Well-controlled
84 experimental exposure systems can thus help identify the underlying causes of known
85 orphan mutational signatures as well as define new patterns generated by candidate
86 carcinogens (reviewed in (3,4)).

87 Various diet-related exposures contribute to the human cancer burden. Examples
88 include contaminants in food or alternative medicines, such as aflatoxin B1 (AFB1) or
89 aristolochic acid (AA). The mutagenicity of these compounds is well-documented; AFB1
90 induces predominantly C:G>A:T base substitutions and AA causes T:A>A:T transversions.
91 The characteristic mutations coupled with information on the preferred sequence contexts in
92 which they are likely to arise allowed unequivocal association of exposure to AFB1 or AA
93 with specific subtypes of hepatobiliary or urological cancers, respectively (5-13).

94 Among dietary compounds with carcinogenic potential, acrylamide is of special
95 interest due to extensive human exposure. Important sources of exposure to acrylamide
96 include tobacco smoke (14), coffee (15), and a broad spectrum of occupational settings (16).
97 Dietary sources of acrylamide comprise carbohydrate-rich food products that have been
98 subject to heating at high temperatures. This is due to Maillard reactions, which involve
99 reducing sugars and the amino acid asparagine, present in potatoes and cereals (17). There
100 is sufficient evidence that acrylamide is carcinogenic in experimental animals (18,19) and it
101 has been classified as a probable carcinogen (Group 2A) by the International Agency for
102 Research on Cancer in 1994 (16). The association of dietary acrylamide exposure with
103 renal, endometrial and ovarian cancers has been explored in recent epidemiological studies
104 (20,21). However, accurate acrylamide exposure assessment in epidemiological studies
105 based on questionnaires has been difficult, and more direct measures of molecular markers,
106 such as hemoglobin adduct levels, may not yield conclusive findings on past exposures (22-

1
2
3 107 27). An improved understanding of its mechanism of action using well-controlled
4 108 experimental systems is critical for understanding the potential carcinogenic risk associated
5 109 with exposure.

6
7 110 Acrylamide undergoes oxidation by cytochrome P450, producing the reactive
8 111 metabolite glycidamide that is highly efficient in DNA binding due to its electrophilic epoxide
9 112 structure (28-30). The *Hras* mutation load in neoplasms of mice exposed to acrylamide or
10 113 glycidamide was found to be considerably higher in mice treated with glycidamide (31). This
11 114 finding is corroborated by a considerably higher mutation frequency in the *cII* reporter gene
12 115 of Big Blue mouse embryonic fibroblasts treated with glycidamide in comparison to
13 116 acrylamide (32,33). Mutation analysis in different experimental *in vivo* and *in vitro* models
14 117 using reporter genes showed an increased association of acrylamide and glycidamide
15 118 exposure with T:A>C:G transitions, as well as T:A>A:T and C:G>G:C transversion mutations
16 119 (31-36), whereas glycidamide exposure was also characterized by C:G>A:T transversions
17 120 (33). However, these proposed acrylamide- and glycidamide-specific mutation patterns were
18 121 based on limited mutation counts in reporter genes and thus do not reflect the complexity of
19 122 genome-wide distributions and profiles. Based on the limited data available thus far, it is not
20 123 possible to translate adequately the reported mutation types (T:A>C:G, T:A>A:T, C:G>G:C,
21 124 C:G>A:T) to global alteration patterns.

22
23
24
25
26
27
28
29 125 The advent of massively parallel sequencing has created the opportunity to study a
30 126 large number of mutations in a single sample, thus significantly enhancing the power of
31 127 mutation analysis in experimental models and enabling reliable identification of specific
32 128 sequence contexts for the induced alterations. Analogously to human cancer genome
33 129 projects, genome-scale mutational signatures can be extracted from highly controlled
34 130 carcinogen exposure experiments using mammalian cell and animal models coupled with
35 131 advanced mathematical approaches (2,3,37,38).

36
37
38
39
40 132 Here we report the systematic assessment of acrylamide and glycidamide
41 133 mutagenicity based on DNA adduct formation and mutation profile analysis using massively
42 134 parallel sequencing in a cell model amenable to the analysis of carcinogen-induced mutation
43 135 patterns and their impact on the resulting cell phenotype (3,37-39). We identify a specific
44 136 and robust mutational signature attributable to glycidamide, and by computationally
45 137 interrogating human cancer genome-wide mutation data, we characterize glycidamide
46 138 signature-positive tumors, thereby highlighting a potential contribution of
47 139 acrylamide/glycidamide exposure to carcinogenesis in humans.

48
49
50
51
52 140

141 **Materials and methods**

142 **Source and authentication of primary cells**

143 Primary Human-p53 knock-in mouse embryonic fibroblasts (Hupki MEFs) were isolated from
144 13.5-day old *Trp53^{tm/Holl}* mouse embryos from the Central Animal Laboratory of the
145 Deutsches Krebsforschungszentrum, Heidelberg, as described previously (40). The mice
146 had been tested for Specific Pathogen-Free (SPF) status. The derived primary cells were
147 genotyped for the human *TP53* codon 72 polymorphism (Table 1) to authenticate the
148 embryo of origin. Cells from three different embryos (E210, E213 and E214) were used for
149 the exposure experiments (Table 1). All subsequent cell cultures were routinely tested at all
150 stages for the absence of mycoplasma.

152 **Cell culture, exposure and immortalization**

153 The primary MEF cells were expanded in Advanced DMEM supplemented with 15% fetal
154 calf serum, 1% penicillin/streptomycin, 1% pyruvate, 1% glutamine, and 0.1% β -mercapto-
155 ethanol. The cells were then seeded in six-well plates and, at passage 2, exposed for 24
156 hours to acrylamide (A4058, Sigma), glycidamide (04704, Sigma), or vehicle (PBS).
157 Acrylamide exposure was carried out in the absence or presence of 2% human S9 fraction
158 (Life Technologies) complemented with NADPH (Sigma). Exposed and control primary cells
159 were cultivated until they bypassed senescence and immortalized clonal cell populations
160 could be isolated (41). The human mammary epithelial cell (HMEC) cultures utilized in this
161 study for whole-genome sequencing (WGS) were generated from benzo[a]pyrene (B[a]P)
162 exposed HMEC described previously (42,43).

164 **MTT assay for cell metabolic activity and viability**

165 Cells were seeded in 96-well plates and treated as indicated. Cell viability was measured 48
166 hours after treatment cessation using CellTiter 96® Aqueous One solution Cell Proliferation
167 Assay (Promega). Plates were incubated for 4 hours at 37°C and absorbance was
168 measured at 492 nm using the APOLLO 11 LB913 plate reader. The MTT assay was
169 performed in triplicates for each experimental condition.

171 **γ H2Ax Immunofluorescence**

172 Immunofluorescence staining was carried out using an antibody specific for Ser139-
173 phosphorylated H2Ax (γ H2Ax) (9718, Cell Signaling Technology). Primary MEFs were
174 seeded on coverslips in 12 well-plates. The cells were incubated in with γ H2Ax-antibody
175 (1:500 in 1% BSA) at 4°C overnight. Subsequent incubation with a fluorochrome-conjugated
176 secondary antibody (4412, Cell Signaling Technology) was carried out for 60 minutes at

1
2
3 177 room temperature. Coverslips were mounted in Vectashield mounting medium with DAPI
4 178 (Eurobio). Immunofluorescence images were captured using a Nikon Eclipse Ti.

5
6 179

7 180 **DNA adduct analysis**

8
9 181 Glycidamide-DNA adducts (N7-(2-carbamoy-2-hydroxyethyl)-guanine (N7-GA-Gua) and N3-
10 182 (2-carbamoy-2-hydroxyethyl)-adenine (N3-GA-Ade)) were quantified by liquid
11 183 chromatography-mass spectrometry (LC-MS/MS) with stable isotope dilution as previously
12 184 described (44) (see Supplementary Materials and Methods for details). The LC-MS/MS used
13 185 for quantification consisted of an Acquity UPLC system (Waters) and a Xevo TQ-S triple
14 186 quadrupole mass spectrometer (Waters). The same MRM transitions as previously
15 187 described (44) were monitored with a cone voltage of 50V and collision energy of 20eV for
16 188 each adduct transition and its corresponding labeled isotope transition.
17
18
19
20
21

22 189

23 190 **TP53 genotyping**

24 191 Exons 4 to 8 of the knocked-in human *TP53* gene (NC_000017.11) were sequenced using
25 192 standard protocols. Sanger sequencing of PCR products was performed at Biofidal (Lyon,
26 193 France). *TP53* primer sequences are listed in Supplementary Materials and Methods.
27 194 Resulting sequences were analyzed using the CodonCode Aligner software.
28
29
30

31 195

32 196 **Library preparation and whole-exome sequencing (WES)**

33 197 Library preparation was carried out using the Kapa Hyper Plus library preparation kit (Kapa
34 198 Biosystems) according the manufacturer's instructions. Exome capture was performed using
35 199 the SureSelect XT Mouse All Exon Kit (Agilent Technologies). Eighteen exome-captured
36 200 libraries were sequenced in the paired-end 150 base-pair run mode using the Illumina
37 201 HiSeq4000 sequencer.
38
39
40

41 202

42 203 **Processing of WES data**

43 204 Fastq files were analyzed for data amount and quality using FastQC (0.11.3) and were
44 205 processed with an in-house pipeline for adapter trimming and alignment to the mm10
45 206 genome (release GRCm38). These components of the pipeline are publicly available at
46 207 <https://github.com/IARCBioinfo/alignment-nf>. The resulting alignment files had a mean depth-
47 208 of-coverage of 135 and 175 for acrylamide and glycidamide samples, respectively. All
48 209 alignment files can be accessed from the NCBI Sequence Read Archive (SRA) data portal
49 210 under the BioProject accession number PRJNA238303. Two somatic variant callers were
50 211 employed with default parameters in order to detect single base substitutions (SBS) and
51 212 small insertions/deletions (indels) (MuTect 1.1.6-4 and Strelka 1.015) in exposed clones,
52 213 using primary cells as normal samples. Each immortalized clone was compared to primary
53
54
55
56
57
58
59
60

1
2
3 214 MEFs from three different embryos (conditions Prim_1, Prim_2, and Prim_3). The overlap of
4 215 the variant calling outcome with respect to the different primary MEFs showed concordance
5 216 close to 80% (Suppl. Fig. S1) with MuTect exhibiting more stringent calling performance.
6 217 Thus, mutation data obtained from the MuTect variant caller were further processed with the
7 218 MutSpec suite ((45); <https://github.com/IARCbioinfo/mutspec>). For more details, see
8 219 Supplementary Materials and Methods and the summary of sequencing metrics (Suppl.
9 220 Table S1), the list of identified MuTect SBS variants (Suppl. Table S2) and indels (Suppl.
10 221 Table S3).
11
12
13
14
15

16 223 **Bioinformatics and statistical analyses**

17 224 The FactoMiner R package (R package version 3.3.2; [https://cran.r-](https://cran.r-project.org/web/packages/FactoMineR)
18 225 [project.org/web/packages/FactoMineR](https://cran.r-project.org/web/packages/FactoMineR)) was used to perform the principal component
19 226 analysis (PCA). To perform the transcription strand bias (SB) analyses, p -values were
20 227 calculated using Pearson's χ^2 test. As multiple comparisons were assessed, the p -value
21 228 was adjusted by applying a false discovery rate (FDR). Statistical analyses were carried out
22 229 using the stats R package. The SB was considered statistically significant at p -value ≤ 0.05 .
23 230 To analyze samples mutation spectra and treatment-specific mutational signatures, filtered
24 231 mutations were classified into 96 types corresponding to the six possible base substitutions
25 232 (C:G>A:T, C:G>G:C, C:G>T:A, T:A>A:T, T:A>C:G, T:A>G:C) and the 16 combinations of
26 233 flanking nucleotides immediately 5' and 3' of the mutated base. Mutation patterns were then
27 234 deconvoluted into mutational signatures using the non-negative matrix factorization (NMF)
28 235 algorithm (46,47). The reconstruction error calculation evaluated the accuracy with which the
29 236 deciphered mutational signatures describe the original mutation spectra of each sample by
30 237 applying Pearson correlation and cosine similarity.

31 238 In order to clean up the profile of the glycidamide mutational signature from the
32 239 residual signature 17 signal and to increase the stability of NMF decomposition, we supplied
33 240 the NMF input by adding samples with a high level of signature 17 (over 65% contribution as
34 241 determined by independent NMF analysis, see Supplementary Materials and Methods).

35 242 Cosine similarity analysis was used to evaluate the concordance of the newly
36 243 identified T:A>A:T-rich mutational signature of glycidamide with the previously reported
37 244 mutational signatures characterized by a predominant T:A>A:T content. These comprised
38 245 COSMIC signatures 22 (AA), 25 and 27 (both of unknown etiology(2)), the experimentally
39 246 derived mutational signature of AA (37,45), 7,12-dimethylbenz[*a*]anthracene (DMBA)
40 247 (48,49), and urethane (50).

41 248 We employed the mutational signature activity (mSigAct) software's sparse signature
42 249 assignment function (`sparse.assign.activity`) (13) to assess the presence of the experimental
43 250 mutational signatures of glycidamide and benzo[*a*]pyrene in whole-genome somatic mutation

1
2
3 251 data from 38 lung adenocarcinomas, 48 lung squamous carcinomas, and 320 liver cancers
4 252 from the ICGC Pan-Cancer Analysis of Whole Genomes (PCAWG) study. We excluded 244
5 253 hyper-mutated microsatellite unstable and aristolochic acid signature-containing liver tumors
6 254 as the presence of high numbers of T>A mutations adversely prevented assessment of the
7 255 possible presence of the glycidamide signature. A set of 11 active COSMIC mutational
8 256 signatures were identified in the remaining tumor samples (excluding COSMIC signature 4).

9
10
11 We defined a 'pure' experimental C>N benzo[a]pyrene signature by WGS (using
12 257 Illumina HiSeq4000 by Genewiz, NJ, USA) of finite lifespan post-stasis clones derived from
13 258 primary human mammary epithelial cells (HMEC) treated with B[a]P as previously described
14 259 (42,43,51). The read alignment to NCBI GRCh38 genome build, variant calling, filtering and
15 260 annotation were consistent with the MutSpec pipeline described above (45). Proportion
16 261 matrices of the experimental GA-signature, the GA-signature normalized to the human
17 262 genome trinucleotide frequency to allow for human PCAWG data screening, and the whole-
18 263 genome B[a]P signature are available in Suppl. Table S4.
19 264

25 265 **Results**

26 266 **Acrylamide and glycidamide induce cytotoxic and genotoxic responses in Hupki** 27 267 **MEFs**

28 268 Upon exposure of primary Hupki MEFs to a range of concentrations of acrylamide (ACR) (in
29 269 the absence or presence of the S9 fraction) and its metabolite, glycidamide (GA), we
30 270 observed a dose-dependent cytotoxic effect on the cells for either compound (Fig. 1A). This
31 271 analysis informed the selection of two conditions for the ACR exposure to be used in the
32 272 subsequent exposure/immortalization experiments, 10 mM ACR for 24 hours in the absence
33 273 of human S9 fraction, and 5 mM ACR for 24 hours in the presence of S9 fraction, which
34 274 elicited 50% (range 30-70%) decrease in cell viability. The IC50 condition for GA was used
35 275 for subsequent mutagenesis analysis, corresponding to a 24-hour treatment with 3 mM of
36 276 the compound. The genotoxic effects of either ACR or GA manifested by a marked increase
37 277 in γ H2Ax staining in the exposed cell populations, in comparison to the mock-treated control
38 278 cells (Fig. 1B).

39 279

40 280 **Immortalized MEF cells accumulate TP53 mutations following acrylamide or** 41 281 **glycidamide treatment**

42 282 Primary MEF cultures from three different embryos (Prim_1, Prim_2, and Prim_3) were
43 283 exposed to ACR or GA using the established conditions and multiple immortalized clones
44 284 were derived. MEF senescence and immortalization phases were evident from the growth
45 285 curves generated for each culture (Suppl. Fig. S2). Subsequently, the clones derived from
46 286 ACR exposure (ACR clones) and GA exposure (GA clones) and spontaneous

1
2
3 287 immortalization (Spont), were pre-screened for *TP53* mutations by Sanger sequencing, to
4 288 assess the mutagenic process prior to exome-scale analysis. In the context of ACR
5 289 treatment, clones obtained from the Prim_2 MEFs that were heterozygous for the
6 290 polymorphic site in codon 72 showed a loss of heterozygosity involving a loss of the proline
7 291 allele in the ACR_1 clone whereas the arginine allele was lost in ACR_2, giving rise to a
8 292 hemizygous clone (Table 1). No *TP53* mutations were observed in any of the three Spont
9 293 clones, whereas 3 out of 7 ACR clones and 1 of 5 GA clones carried non-synonymous *TP53*
10 294 mutations (Table 1). The detected mutations indicated specific selection for mutations in the
11 295 *TP53* gene during cell immortalization and confirmed the clonal nature of MEF
12 296 immortalization.
13
14
15
16
17
18

297

298 **Analysis of mutation spectra**

19
20 299 Whole-exome sequencing of all spontaneously immortalized and exposed clones and
21 300 subsequent extraction of acquired variants revealed that the total number of acquired SBS
22 301 did not differ markedly between the ACR and Spont clones. The Spont clones harbored on
23 302 average 190 (median = 151, range = 141-277) SBS, whereas the ACR clones had on
24 303 average 208 (median = 173, range = 151-262) SBS. In contrast, the total number of SBS
25 304 was considerably increased in the GA clones, with an average of 485 SBS (median = 448,
26 305 range = 370-592) (Suppl. Table S1 and S2). This finding suggests markedly stronger
27 306 mutagenic properties of GA in the MEFs. To estimate the extent of sequencing-related
28 307 damage in our samples, we determined the GIV score of each sample as described in
29 308 Materials and Methods and in (52). No detectable damage for any of the mutation types was
30 309 observed in our dataset (data not shown). The ACR exposed samples exhibited an overall
31 310 diffuse pattern across the six different SBS types (Suppl. Fig. S3). The Spont clones showed
32 311 an enrichment of C:G>G:C SBS in the 5'-GCC-3' context, which was also present at varying
33 312 levels in the exposed cultures. This particular mutation type appears to be related to the
34 313 culture conditions used for the immortalization assay, as its presence has previously been
35 314 noted upon spontaneous as well as exposure-driven MEF immortalization (37). No
36 315 significant transcription strand bias was observed for any of the mutation classes in the
37 316 Spont or ACR clones (Suppl. Fig. S4). In the five clones derived from the GA-treated primary
38 317 MEF cultures, we observed an enrichment of acquired T:A>A:T and C:G>A:T transversions
39 318 and T:A>C:G transitions (Suppl. Fig. S3B), marked by significant transcription strand bias
40 319 (Suppl. Fig. S4).

41
42
43 320 PCA performed on the resulting 6-class SBS spectra unambiguously separated the
44 321 GA clones from the remaining experimental conditions (Fig. 2A). The analysis of indels
45 322 (listed in Suppl. Table S3) showed lower numbers of these alterations in the GA-associated
46 323 clones compared to the ACR or Spont clones (Fig. 2B). This suggests that a higher

1
2
3 324 accumulation of SBS may selectively promote the senescence bypass and selection of the
4 325 GA clones, with a decreased functional contribution of indels, while an inverse scenario is
5
6 326 plausible in case of the Spont and ACR clones, reminiscent of a previous report based on
7 327 the Big Blue mouse embryonic fibroblasts and *c/* transgene (53).
8
9 328

10 329 **Variant allele frequency analysis**

11 330 Variant allele frequency (VAF) analysis was carried out for GA clones. Overall, a significant
12 331 proportion of acquired mutations was present at allelic frequencies between 25-75% (Suppl.
13 332 Fig. S5). Upon grouping of substitutions into bins of high (67-100%), medium (34-66%) and
14 333 low (0-33%) VAF, the predominant GA-specific mutation types (T:A>A:T, T:A>C:G and
15 334 C:G>A:T) started manifesting at high VAF, whereas the 5'-NIT-3' alterations, corresponding
16 335 to the COSMIC signature 17 previously reported to arise in cultured mouse cells including
17 336 MEFs (38,54,55) showed lower VAF, therefore a later appearance in the cultures (Suppl.
18 337 Fig. S6). This observation suggests the early effects of the GA exposure and the
19 338 reproducible contribution of the induced mutations to the senescence bypass and their clonal
20 339 propagation during the immortalization stage.
21
22
23
24
25
26
27

28 341 **Mutational signature analysis**

29 342 Using NMF, we extracted the mutational signatures from all the MEF clones. Using
30 343 computed statistics for estimating the number of signatures, three signatures were identified
31 344 as an optimal number, with signatures A and C enriched in the Spont and ACR clones, and
32 345 signature B selectively enriched in the GA clones (Fig. 2C,D). Reconstruction of the
33 346 observed mutation spectra supports the robustness of the signature analysis with strong
34 347 Pearson's correlation and cosine similarity in GA-derived clones (Fig. 2D). In signature C
35 348 and also to a lesser extent in signatures A and B, we observed an admixture of a pattern
36 349 identical to the orphan COSMIC signature 17 (T:A>G:C in a 5'-NIT-3' trinucleotide context),
37 350 described in various human cancers (most notably esophageal adenocarcinoma), but also
38 351 seen in aflatoxin B1-driven mouse liver cancers (11), as well as primary MEF-derived clones
39 352 (37,38). In *in vitro* contexts, this signature has been linked to cell culture conditions and
40 353 associated oxidative stress (54,55). To refine further the obtained experimental signatures,
41 354 we developed a signature 'baiting' approach that combined the MEF clones data with
42 355 signature 17-rich data from esophageal adenocarcinomas from the ICGC ESAD-UK study
43 356 for new NMF analysis (56). This resulted in considerable reduction (average = 47%, median
44 357 = 48%) of the signature 17-specific most prominent T>G peaks and a more refined pattern
45 358 for signature B, associated primarily with GA treatment (Fig. 3A and Suppl. Fig. S7). This
46 359 putative GA signature retains the predominant enrichment for the T:A>A:T transversions and
47 360 T:A>C:G transitions in the 5'-CTG-3' and 5'-CTT-3' trinucleotide contexts, and the C:G>A:T

1
2
3 361 component. Moreover, these mutation types were marked by significant transcription strand
4 362 bias (Fig. 3B and Suppl. Fig. S4), exhibiting higher accumulation of mutations on the non-
5 363 transcribed strand consistent with the decreased efficiency of the transcription-coupled
6 364 nucleotide excision repair due to adduct formation.
7
8
9

365

366 **DNA adduct analysis**

10
11 367 Following metabolic activation, acrylamide induces well-characterized glycidamide DNA
12 368 adducts at the N7- and N3-positions of guanine and adenine, respectively. LC-MS/MS-based
13 369 adduct quantification revealed the absence of these adducts in the spontaneously
14 370 immortalized control samples as well as in MEFs exposed to acrylamide in the absence of
15 371 S9 fraction (levels below the limit of detection). This suggests the lack of CYP2E1 activity,
16 372 which is required for the metabolism of acrylamide to glycidamide, in the MEFs. Upon
17 373 addition of human S9 fraction, N7-GA-Gua levels increased to 11 adducts/ 10^8 nucleotides,
18 374 suggesting limited metabolic activation of acrylamide due to the presence of enzymatic
19 375 activity in the S9 fraction (Fig. 3C and Suppl. Fig. S8). Glycidamide-exposed cells exhibited
20 376 significantly increased DNA adduct levels, with both N7-GA-Gua and N3-GA-Ade observed
21 377 at very high average levels, 49 000 adducts/ 10^8 nucleotides and 350 adducts/ 10^8
22 378 nucleotides, respectively, after subtracting the trace amount of contamination from the
23 379 internal standard (Fig. 3C and Suppl. Fig. S8).
24
25
26
27
28
29
30

380

381 **Comparison of the glycidamide signature to known signatures characterized by** 382 **prominent T:A>A:T profiles**

35
36 383 We next performed cosine similarity analysis of the putative GA signature and all known
37 384 T:A>A:T-rich signatures extracted from primary cancers as well as experimental systems
38 385 (Fig. 3D and Suppl. Fig. S9). The best match was 84% pattern similarity with COSMIC
39 386 signature 25 (derived from four Hodgkin lymphoma cell lines) (Fig. 3D). However, unlike the
40 387 GA signature, COSMIC signature 25 exhibits strand bias for only T:A>A:T mutations and no
41 388 transcription strand bias for the T:A>C:G mutations. Thus, the mutation patterns and strand
42 389 bias on all three main mutation types generated by GA treatment (Fig. 3A,B) appear specific
43 390 and novel.
44
45
46

391

392 **Glycidamide signature screening in human tumor data from the ICGC PCAWG**

47
48
49 393 The initial mSigAct test performed on PCAWG data from lung and liver tumors indicated a
50 394 marked presence of the GA signature. This observation was in keeping with the presence of
51 395 acrylamide in tobacco smoke and was further corroborated by a cosine similarity of 94%
52 396 between the adenine (T>N) components of COSMIC signature 4 (tobacco smoking) and the
53 397 GA signature (Fig. 4A). We thus hypothesized that COSMIC signature 4 reflects co-
54
55
56
57
58
59
60

1
2
3 398 exposure to B[a]P (generating C>N/guanine mutations with transcription strand bias) and to
4 399 GA (generating T>N/adenine mutations with transcription strand bias) (Fig. 4A,B). To
5
6 400 provide further experimental evidence, we generated a 'pure' B[a]P mutational signature by
7
8 401 whole-genome sequencing of cell clones derived from B[a]P-exposed normal human
9
10 402 mammary epithelial cells (HMEC). This yielded a robust signature characterized by
11
12 403 predominant strand biased guanine (mainly C>A) mutation levels and negligibly mutated
13
14 404 adenines (T>N) (Fig. 4A,B). Next, we used mSigAct to interrogate the PCAWG tumor
15
16 405 samples for the level of exposure to the experimentally defined GA and B[a]P signatures
17
18 406 (alongside other COSMIC mutational signatures) in 48 lung squamous carcinomas, 38 lung
19
20 407 adenocarcinomas, and 320 liver cancers. We compared these to estimated levels of
21
22 408 exposure to COSMIC signature 4, and found that in the lung cancers, a combination of the
23
24 409 GA and B[a]P signatures accounted for very similar numbers of mutations as COSMIC
25
26 410 signature 4, thus further supporting the hypothesis that COSMIC signature 4 represents
27
28 411 combined and highly correlated exposure to GA and B[a]P (Fig. 4C). Compared to lung
29
30 412 cancers, we found more variability in the assignment of mutation numbers to GA and B[a]P
31
32 413 versus COSMIC signature 4 in liver cancers (Fig. 4C), which may reflect a decreased
33
34 414 relationship between GA and B[a]P exposure due to generally more complex exposure
35
36 415 history in the liver. The successful reconstruction of COSMIC signature 4 by the
37
38 416 experimental GA- and B[a]P- signatures in the lung and liver human tumors enabled correct
39
40 417 assignment of the GA-signature in a subset of 29 lung adenocarcinomas, 46 lung SCC and
41
42 418 26 liver tumors (Fig. 4D). The SBS counts corresponding to GA-mutational signature ranged
43
44 419 between 300 up to 43,000 mutations/per sample in lung tumors, and between 190 to 23,000
45
46 420 mutations/per sample in liver tumors (Fig. 4D and Suppl. Table S5). These findings indicate
47
48 421 exposure to glycidamide linked to tobacco smoking – when concomitant with B[a]P-
49
50 422 signature, or through diet or occupation – in the absence of B[a]P signature (samples Liver-
51
52 423 HCC::SP112224; Liver-HCC::SP49551; Liver-HCC::SP50105; Liver-HCC::SP98861; Liver-
53
54 424 HCC::SP50183, see Suppl. Fig. S10 and Suppl. Table S5).

425 **Discussion**

426 In this study we report the identification of an exome-wide mutational signature for
47
48 427 glycidamide, a metabolite of the probable human carcinogen acrylamide. The newly
49
50 428 identified signature is based on massively parallel sequencing performed in a well-controlled
51
52 429 experimental carcinogen exposure-clonal immortalization model, revealing characteristic
53
54 430 mutagenic effects of glycidamide. The glycidamide mutational signature presented here and
55
56 431 the results of statistical assessment of its presence in multiple human tumor types may help
57
58 432 clarify the thus-far tenuous association of acrylamide with human cancer.

1
2
3 433 In concordance with its *in vivo* carcinogenicity in rodents (16,19,31,57), our findings
4 434 in the established MEF carcinogen exposure and immortalization system suggest that
5 435 characteristic mutagenic effects may play a role during acrylamide/glycidamide-driven tumor
6 436 development. In contrast to glycidamide, acrylamide exposure led neither to an increased
7 437 number of SBS nor did it induce characteristic mutation types in the MEF exposure system.
8
9 438 Despite the absence of a mutagenic effect of acrylamide in our experiments, acrylamide and
10 439 glycidamide exposures induce an almost identical set of tumors in both mice and rats,
11 440 providing a substantial argument for a glycidamide-mediated tumorigenic effect of
12 441 acrylamide (19). This is further supported by mechanistic studies showing that lung tissue
13 442 from mice exposed to acrylamide and glycidamide displays comparable DNA adduct
14 443 patterns as well as similar mutation frequencies in the *cII* transgene (36). Similar
15 444 observations had been made in the context of *in vitro* mutagenicity of acrylamide in human
16 445 and mouse cells, suggesting the key role for epoxide metabolite glycidamide to form pre-
17 446 mutagenic DNA adducts (33).

18
19 447 As shown by our adduct analysis, acrylamide is not efficiently metabolized by MEFs.
20 448 This finding is in keeping with the results from previous animal carcinogenicity studies. In
21 449 fact, glycidamide induces hepatocellular carcinomas in neonatal B6C3F1 mice, whereas
22 450 administration of acrylamide does not increase the tumor incidence. This has been attributed
23 451 to the inability of neonatal mice to efficiently metabolize acrylamide (31). Moreover, in
24 452 contrast to acrylamide treatment, glycidamide induces tumors of the small intestine in a
25 453 dose-dependent manner upon perinatal exposure (57) and similar observations were made
26 454 for glycidamide mutagenicity *in vitro* (33). We compensated for the lack of proper acrylamide
27 455 metabolic activation by the addition of human S9 fraction, and the assessment of DNA
28 456 adducts indeed suggests acrylamide metabolic activation upon addition of S9. However, the
29 457 adduct levels are substantially lower compared to glycidamide exposure, which may account
30 458 for the observed differences in mutagenicity. Interestingly, a consistent minor contribution of
31 459 the glycidamide mutational signature was detected in the majority of ACR clones, whereas it
32 460 was absent in the Spont clones. This raises the possibility that partial metabolic activation of
33 461 acrylamide in the MEF system resulted in low levels of glycidamide. However, a clear
34 462 mutational signature in the employed experimental setting was achieved only by exposing
35 463 the cells directly to glycidamide.

36 464 Single reporter gene studies had previously linked acrylamide and glycidamide
37 465 exposure to multiple different mutation types. Thanks to the larger number of mutations
38 466 captured by exome sequencing, we were able to attribute to the glycidamide exposure a
39 467 particular mutational signature characterized by strand-biased C:G>A:T and T:A>A:T
40 468 transversions, and T:A>C:A transitions towards the non-transcribed strand suggesting a
41 469 formation of DNA-adducts. The presence of N7-GA-Gua and N3-GA-Ade, two well-

1
2
3 470 characterized glycidamide DNA adducts originating from the metabolic conversion of
4 471 acrylamide (30,44,53), shows a remarkable relationship between DNA adduct profiles and
5 472 the putative mutational signature of glycidamide. N3-GA-Ade and N7-GA-Gua are
6 473 depurinating adducts. They can result in apurinic/apyrimidinic sites, which, during replication,
7 474 induce the mis-incorporation of deoxyadenine, leading to the observed T:A>A:T and
8 475 C:G>A:T transversions of the glycidamide signature, respectively. The third mutation type
9 476 specifically enriched in the glycidamide signature, T:A>C:G transitions, has been ascribed to
10 477 the N1-GA-Ade adduct, a miscoding adduct and the most commonly identified adenine
11 478 adduct *in vitro* (35,44,53,58). Levels of the guanine adduct were especially high in the
12 479 exposed MEF cells, whereas the associated C:G>A:T transversions in the resulting post-
13 480 senescence clones were less represented. This could reflect differences in DNA repair
14 481 efficiency concerning individual GA-DNA adduct species, or the fact that the resulting clones
15 482 are derived from single cells whereas the GA-DNA adducts were measured on average in
16 483 the bulk primary cell population. A mechanism of negative selection of cells with high N7-
17 484 GA-Gua adduct burden is also plausible.

18 485 We observed consistent presence of COSMIC signature 17 in the data generated
19 486 from the untreated and treated MEF clones. The etiology of signature 17 remains unknown.
20 487 While some candidate causal factors have been proposed in esophageal adenocarcinoma
21 488 and gastric cancers (e.g., inflammatory conditions due to acid reflux, *H. pylori*) (56) and in
22 489 cultured mouse cell systems (54,55), further studies are required to establish why signature
23 490 17 tends to arise *in vitro* in immortalized clones derived from mouse embryonic fibroblasts as
24 491 observed in our study and also previous work (38).

25 492 Genome-scale sequencing of tumor tissues will be needed to verify, *in vivo*, the
26 493 glycidamide mutational signature identified in this study. The established animal models
27 494 (18,19) of acrylamide- and glycidamide-mediated tumorigenesis provide a suitable starting
28 495 point, and it would be interesting to compare mutational signatures derived from these
29 496 models with the *in vitro* results. The identified glycidamide signature with its extended
30 497 features of transcription strand bias for the major mutation types differs from the currently
31 498 known COSMIC signatures (Fig. 3D). In addition, we show that in the cancer genome
32 499 sequencing data sets from the ICGC PCAWG effort, the putative glycidamide-mutational
33 500 signature can be identified in a subset of tumors of the lung and liver (sites of possible
34 501 acrylamide exposure due to tobacco smoking), based on combining experimentally derived
35 502 signatures with sophisticated computational signature reconstruction approaches (Fig. 4).

36 503 The continued interest in understanding the contribution of acrylamide and its
37 504 electrophilic metabolite glycidamide to cancer development reflects recent accumulation of
38 505 new mechanistic data on the animal carcinogenicity of the compounds. The possible
39 506 carcinogenic effects in humans have been recommended for re-evaluation by the Advisory

1
2
3 507 Group to the Monographs Program of the International Agency for Research on Cancer (59).
4 508 Our findings related to the reconstruction of COSMIC signature 4 using the experimental
5 509 GA-signature and B[a]P signature, together with the presence of the GA signature in the
6 510 lung and liver cancer data are relevant given the established high contents of acrylamide in
7 511 tobacco smoke. Despite the absence of prominent T>N (adenine) mutations in the
8 512 experimental B[a]P exposure setting, we cannot exclude a possibility that in the human lung
9 513 cells the adenine residues can be additionally targeted by other tobacco carcinogens such
10 514 as benzo[a]pyrene derivatives or nitrosamines. Importantly, five liver tumor samples
11 515 identified in this study harbored the GA signature but the major features of signature 4 as
12 516 represented by the experimental B[a]P signature were absent (Suppl. Fig. S10, Suppl. Table
13 517 S5). These tumors are thus of particular interest as they could reflect dietary or occupational
14 518 exposure to acrylamide.

15
16 519 The presented mutational signature of glycidamide and its potential use for screening
17 520 of cancer genome sequencing data may provide a basis for relevant assessment of cancer
18 521 risk through new carefully designed molecular cancer epidemiology studies. Future
19 522 validation analyses involving e.g. GA-DNA adduct monitoring in non-tumor tissue of cancer
20 523 patients or in animal exposure models are warranted to provide additional evidence that the
21 524 predominant T>N mutations in the cancers identified in this study indeed originate from
22 525 exposure to acrylamide and its reactive metabolite glycidamide.

23 24 25 26 27 28 29 30 31 32 33 526 **Acknowledgments**

34 527 The views expressed in this manuscript do not necessarily represent those of the U.S. Food
35 528 and Drug Administration. The study was supported by funding obtained from INCa-INSERM
36 529 (Plan Cancer 2015 grant to J.Z.), NIH/NIEHS (1R03ES025023-01A1 grant to M.O.), and the
37 530 Singapore National Medical Research Council (NMRC/CIRG/1422/2015 grant to S.G.R.) and
38 531 the Singapore Ministry of Health via the Duke-NUS Signature Research Programmes to
39 532 S.G.R.. M.R.S. was supported by the U.S. Department of Energy under Contract No. DE-
40 533 AC02-05CH11231. We thank the NYU Genome Technology Center, funded in part by the
41 534 NIH/NCI Cancer Center Support Grant P30CA016087, and GENEWIZ, South Plainfield, NJ,
42 535 USA, for expert assistance with Illumina sequencing.

43
44
45
46
47
48 536

49 50 51 537 **References**

- 52 538 1. Smith, M.T., *et al.* (2016) Key Characteristics of Carcinogens as a Basis for
53 539 Organizing Data on Mechanisms of Carcinogenesis. *Environ Health Perspect*, **124**,
54 540 713-21.
55 541 2. Alexandrov, L.B., *et al.* (2013) Signatures of mutational processes in human cancer.
56 542 *Nature*, **500**, 415-421.

3. Zhivagui, M., *et al.* (2017) Modelling Mutation Spectra of Human Carcinogens Using Experimental Systems. *Basic Clin Pharmacol Toxicol*, **121 Suppl 3**, 16-22.
4. Hollstein, M., *et al.* (2017) Base changes in tumour DNA have the power to reveal the causes and evolution of cancer. *Oncogene*, **36**, 158-167.
5. Poon, S.L., *et al.* (2013) Genome-Wide Mutational Signatures of Aristolochic Acid and Its Application as a Screening Tool. *Science Translational Medicine*, **5**, 197ra101-197ra101.
6. Meier, B., *et al.* (2014) *C. elegans* whole-genome sequencing reveals mutational signatures related to carcinogens and DNA repair deficiency. *Genome Res*, **24**, 1624-36.
7. Scelo, G., *et al.* (2014) Variation in genomic landscape of clear cell renal cell carcinoma across Europe. *Nat Commun*, **5**, 5135.
8. Jelakovic, B., *et al.* (2015) Renal cell carcinomas of chronic kidney disease patients harbor the mutational signature of carcinogenic aristolochic acid. *Int J Cancer*, **136**, 2967-72.
9. Hoang, M.L., *et al.* (2016) Aristolochic Acid in the Etiology of Renal Cell Carcinoma. *Cancer Epidemiology, Biomarkers & Prevention*, **25**, 1600-1608.
10. Chawanthayatham, S., *et al.* (2017) Mutational spectra of aflatoxin B1 in vivo establish biomarkers of exposure for human hepatocellular carcinoma. *Proc Natl Acad Sci U S A*, **114**, E3101-E3109.
11. Huang, M.N., *et al.* (2017) Genome-scale mutational signatures of aflatoxin in cells, mice, and human tumors. *Genome Res*, **27**, 1475-1486.
12. Zhang, W., *et al.* (2017) Genetic Features of Aflatoxin-Associated Hepatocellular Carcinoma. *Gastroenterology*, **153**, 249-262 e2.
13. Ng, A.W.T., *et al.* (2017) Aristolochic acids and their derivatives are widely implicated in liver cancers in Taiwan and throughout Asia. *Sci Transl Med*, **9**.
14. Mojska, H., *et al.* (2016) Acrylamide content in cigarette mainstream smoke and estimation of exposure to acrylamide from tobacco smoke in Poland. *Annals of agricultural and environmental medicine: AAEM*, **23**, 456-461.
15. Takatsuki, S., *et al.* (2003) Determination of acrylamide in processed foods by LC/MS using column switching. *Shokuhin Eiseigaku Zasshi. Journal of the Food Hygienic Society of Japan*, **44**, 89-95.
16. IARC Monograph vol. 60 (1994) *Some industrial chemicals. Lyon, 15 - 22 February 1994*, Lyon.
17. Tareke, E., *et al.* (2002) Analysis of Acrylamide, a Carcinogen Formed in Heated Foodstuffs. *Journal of Agricultural and Food Chemistry*, **50**, 4998-5006.
18. Beland, F.A., *et al.* (2013) Carcinogenicity of acrylamide in B6C3F(1) mice and F344/N rats from a 2-year drinking water exposure. *Food and Chemical Toxicology*, **51**, 149-159.
19. Beland, F.A., *et al.* (2015) Carcinogenicity of glycidamide in B6C3F1 mice and F344/N rats from a two-year drinking water exposure. *Food and Chemical Toxicology*, **86**, 104-115.
20. Hogervorst, J.G., *et al.* (2008) Dietary acrylamide intake and the risk of renal cell, bladder, and prostate cancer. *The American Journal of Clinical Nutrition*, **87**, 1428-1438.
21. Virk-Baker, M.K., *et al.* (2014) Dietary Acrylamide and Human Cancer: A Systematic Review of Literature. *Nutrition and Cancer*, **66**, 774-790.
22. Olesen, P.T., *et al.* (2008) Acrylamide exposure and incidence of breast cancer among postmenopausal women in the Danish Diet, Cancer and Health Study. *International Journal of Cancer*, **122**, 2094-2100.
23. Wilson, K.M., *et al.* (2009) Acrylamide exposure measured by food frequency questionnaire and hemoglobin adduct levels and prostate cancer risk in the Cancer of the Prostate in Sweden Study. *International Journal of Cancer*, **124**, 2384-2390.

- 1
2
3 596 24. Xie, J., *et al.* (2013) Acrylamide Hemoglobin Adduct Levels and Ovarian Cancer
4 597 Risk: A Nested Case-Control Study. *Cancer Epidemiology Biomarkers & Prevention*,
5 598 **22**, 653-660.
- 6 599 25. Obón-Santacana, M., *et al.* (2016) Acrylamide and glycidamide hemoglobin adduct
7 600 levels and endometrial cancer risk: A nested case-control study in nonsmoking
8 601 postmenopausal women from the EPIC cohort. *International Journal of Cancer*, **138**,
9 602 1129-1138.
- 10 603 26. Obón-Santacana, M., *et al.* (2016) Acrylamide and Glycidamide Hemoglobin Adducts
11 604 and Epithelial Ovarian Cancer: A Nested Case-Control Study in Nonsmoking
12 605 Postmenopausal Women from the EPIC Cohort. *Cancer Epidemiology, Biomarkers &
13 606 Prevention*, **25**, 127-134.
- 14 607 27. Obón-Santacana, M., *et al.* (2016) Dietary and lifestyle determinants of acrylamide
15 608 and glycidamide hemoglobin adducts in non-smoking postmenopausal women from
16 609 the EPIC cohort. *European Journal of Nutrition*.
- 17 610 28. Sumner, S.C., *et al.* (1999) Role of cytochrome P450 2E1 in the metabolism of
18 611 acrylamide and acrylonitrile in mice. *Chemical Research in Toxicology*, **12**, 1110-
19 612 1116.
- 20 613 29. Ghanayem, B.I., *et al.* (2005) Role of CYP2E1 in the epoxidation of acrylamide to
21 614 glycidamide and formation of DNA and hemoglobin adducts. *Toxicological Sciences*,
22 615 **88**, 311-318.
- 23 616 30. Segerbäck, D., *et al.* (1995) Formation of N-7-(2-carbamoyl-2-hydroxyethyl) guanine
24 617 in DNA of the mouse and the rat following intraperitoneal administration of [14C]
25 618 acrylamide. *Carcinogenesis*, **16**, 1161-1165.
- 26 619 31. Von Tungeln, L.S., *et al.* (2012) Tumorigenicity of acrylamide and its metabolite
27 620 glycidamide in the neonatal mouse bioassay. *International Journal of Cancer*, **131**,
28 621 2008-2015.
- 29 622 32. Besaratinia, A., *et al.* (2003) Weak yet distinct mutagenicity of acrylamide in
30 623 mammalian cells. *Journal of the National Cancer Institute*, **95**, 889-896.
- 31 624 33. Besaratinia, A., *et al.* (2004) Genotoxicity of acrylamide and glycidamide. *Journal of
32 625 the National Cancer Institute*, **96**, 1023-1029.
- 33 626 34. Von Tungeln, L.S., *et al.* (2009) DNA adduct formation and induction of micronuclei
34 627 and mutations in B6C3F1/Tk mice treated neonatally with acrylamide or glycidamide.
35 628 *International Journal of Cancer*, **124**, 2006-2015.
- 36 629 35. Ishii, Y., *et al.* (2015) Acrylamide induces specific DNA adduct formation and gene
37 630 mutations in a carcinogenic target site, the mouse lung. *Mutagenesis*, **30**, 227-235.
- 38 631 36. Manjanatha, M.G., *et al.* (2015) Acrylamide-induced carcinogenicity in mouse lung
39 632 involves mutagenicity: *cil* gene mutations in the lung of big blue mice exposed to
40 633 acrylamide and glycidamide for up to 4 weeks. *Environ Mol Mutagen*, **56**, 446-56.
- 41 634 37. Olivier, M., *et al.* (2014) Modelling mutational landscapes of human cancers in vitro.
42 635 *Scientific Reports*, **4**.
- 43 636 38. Nik-Zainal, S., *et al.* (2015) The genome as a record of environmental exposure.
44 637 *Mutagenesis*, **30**, 763-70.
- 45 638 39. Huskova, H., *et al.* (2017) Modeling cancer driver events in vitro using barrier
46 639 bypass-clonal expansion assays and massively parallel sequencing. *Oncogene*, **36**,
47 640 6041-6048.
- 48 641 40. Liu, Z., *et al.* (2004) Human tumor p53 mutations are selected for in mouse
49 642 embryonic fibroblasts harboring a humanized p53 gene. *Proceedings of the National
50 643 Academy of Sciences of the United States of America*, **101**, 2963-2968.
- 51 644 41. Todaro, G.J., *et al.* (1963) Quantitative studies of the growth of mouse embryo cells
52 645 in culture and their development into established lines. *The Journal of Cell Biology*,
53 646 **17**, 299-313.
- 54 647 42. Severson, P.L., *et al.* (2014) Exome-wide mutation profile in benzo[a]pyrene-derived
55 648 post-stasis and immortal human mammary epithelial cells. *Mutation
56 649 Research/Genetic Toxicology and Environmental Mutagenesis*, **775-776**, 48-54.

- 1
2
3 650 43. Stampfer, M.R., *et al.* (1985) Induction of transformation and continuous cell lines
4 651 from normal human mammary epithelial cells after exposure to benzo[a]pyrene. *Proc*
5 652 *Natl Acad Sci U S A*, **82**, 2394-8.
6 653 44. Gamboa da Costa, G., *et al.* (2003) DNA adduct formation from acrylamide via
7 654 conversion to glycidamide in adult and neonatal mice. *Chemical Research in*
8 655 *Toxicology*, **16**, 1328-1337.
9 656 45. Ardin, M., *et al.* (2016) MutSpec: a Galaxy toolbox for streamlined analyses of
10 657 somatic mutation spectra in human and mouse cancer genomes. *BMC*
11 658 *Bioinformatics*, **17**, 170.
12 659 46. Brunet, J.-P., *et al.* (2004) Metagenes and molecular pattern discovery using matrix
13 660 factorization. *Proceedings of the National Academy of Sciences of the United States*
14 661 *of America*, **101**, 4164-4169.
15 662 47. Alexandrov, Ludmil B., *et al.* (2013) Deciphering Signatures of Mutational Processes
16 663 Operative in Human Cancer. *Cell Reports*, **3**, 246-259.
17 664 48. McCreery, M.Q., *et al.* (2015) Evolution of metastasis revealed by mutational
18 665 landscapes of chemically induced skin cancers. *Nature Medicine*, **21**, 1514-1520.
19 666 49. Nassar, D., *et al.* (2015) Genomic landscape of carcinogen-induced and genetically
20 667 induced mouse skin squamous cell carcinoma. *Nature Medicine*, **21**, 946-954.
21 668 50. Westcott, P.M.K., *et al.* (2014) The mutational landscapes of genetic and chemical
22 669 models of Kras-driven lung cancer. *Nature*, **517**, 489-492.
23 670 51. Stampfer, M.R., *et al.* (1988) Human mammary epithelial cells in culture:
24 671 differentiation and transformation. *Cancer Treat Res*, **40**, 1-24.
25 672 52. Chen, L., *et al.* (2017) DNA damage is a pervasive cause of sequencing errors,
26 673 directly confounding variant identification. *Science*, **355**, 752-756.
27 674 53. Besaratinia, A., *et al.* (2005) DNA adduction and mutagenic properties of acrylamide.
28 675 *Mutation Research*, **580**, 31-40.
29 676 54. Behjati, S., *et al.* (2014) Genome sequencing of normal cells reveals developmental
30 677 lineages and mutational processes. *Nature*, **513**, 422-425.
31 678 55. Milholland, B., *et al.* (2017) Differences between germline and somatic mutation rates
32 679 in humans and mice. *Nat Commun*, **8**, 15183.
33 680 56. Secrier, M., *et al.* (2016) Mutational signatures in esophageal adenocarcinoma define
34 681 etiologically distinct subgroups with therapeutic relevance. *Nature Genetics*, **48**,
35 682 1131-1141.
36 683 57. Olstørn, H.B.A., *et al.* (2007) Effects of perinatal exposure to acrylamide and
37 684 glycidamide on intestinal tumorigenesis in Min/+ mice and their wild-type litter mates.
38 685 *Anticancer Research*, **27**, 3855-3864.
39 686 58. Randall, S.K., *et al.* (1987) Nucleotide insertion kinetics opposite abasic lesions in
40 687 DNA. *Journal of Biological Chemistry*, **262**, 6864-6870.
41 688 59. Straif, K., *et al.* (2014) Future priorities for the IARC Monographs. *The Lancet*
42 689 *Oncology*, **15**, 683-684.

43 690
44
45

46 691 **Figure legends**

47 692 **Figure 1:** Acrylamide- and glycidamide-induced cytotoxicity and genotoxicity *in vitro*. **(A)** Cell
48 693 viability, following 24-hour treatment of primary MEFs with the indicated concentrations of
49 694 acrylamide (top panel), in the absence (diamonds) and presence (circles) of human S9
50 695 fraction, and glycidamide (bottom panel), as determined by MTT assay. Absorbance was
51 696 measured 48 hours after treatment cessation and was normalized to untreated cells. The
52 697 results are expressed as mean percent \pm SD of three replicates. **(B)** DNA damage

698 assessment by immunofluorescence with an antibody specific for Ser139-phosphorylated
699 histone H2Ax (γ H2Ax). Primary MEFs were treated with acrylamide or glycidamide for 24
700 hours prior to immunofluorescence. Compound concentrations used were based on 20-70%
701 viability reduction in the MTT assay: 10 mM acrylamide, 5 mM acrylamide in the presence of
702 S9 fraction and 3 mM glycidamide. ACR: acrylamide; GA: glycidamide.

703 **Figure 2:** Analysis of the mutation patterns derived from exome sequencing data from
704 immortalized Hupki MEF clones. **(A)** Principle component analysis (PCA) of WES data. PCA
705 was computed using as input the mutation count matrix of the clones that immortalized
706 spontaneously (Spont) or were derived from exposure to acrylamide (ACR) or glycidamide
707 (GA). Each sample is plotted considering the value of the first and second principal
708 components (Dim1 and Dim2). The percentage of variance explained by each component is
709 indicated within brackets on each axis. Spont, ACR- and GA-exposed samples are
710 represented by differently colored symbols. **(B)** Representation of small insertions and
711 deletions (indels) counts within the immortalized clones as determined by the Strelka variant
712 caller. **(C)** Mutational signatures identified by non-negative matrix factorization (NMF) in the
713 15 Hupki MEF-derived clones (sig A, sig B, and sig C). X-axis represents the trinucleotide
714 sequence context. Y-axis represents the frequency distribution of the mutations. The
715 predominant trinucleotide context for T:A > A:T mutations is indicated in sig B (5'-CTG-3').
716 The trinucleotide contexts for C:G > G:C (5'-GCC-3') and T:A > G:C mutations (5'-NTT-3')
717 are highlighted in sig C. **(D)** Contribution of the identified signatures to each sample (X-axis),
718 assigned either by absolute SBS counts or by proportion (bar graphs). The reconstruction
719 accuracy of the identified mutational signatures in individual samples is shown in the bottom
720 scatter plot (Y-axis value of 1 = 100% accuracy).

721 **Figure 3:** **(A)** Refinement of GA signature. The contribution of signature 17 (T:A>G:C in 5'-
722 NTT-3' context), present in all clones, was decreased by performing NMF on Hupki samples
723 pooled with primary tumor samples with high levels of signature 17 (see Methods). **(B)**
724 Transcription strand bias analysis for the six mutation types in GA-exposed clones. For each
725 mutation type, the number of mutations occurring on the transcribed (T) and non-transcribed
726 (N) strand is shown on the Y-axis. *** $p < 10^{-8}$; * $p < 10^{-2}$. **(C)** DNA adducts analysis as
727 determined by LC-MS/MS. Levels of N7-GA-Gua adduct in ACR+S9 and GA treated MEFs
728 and N3-GA-Ade DNA adduct level in GA treated MEFs. The data are presented as the
729 number of adducts in 10^8 nucleotides. $n \geq 2$. **(D)** Cosine similarity matrix comparing the
730 putative glycidamide mutational signature with other A>T rich mutational signatures from
731 COSMIC (signatures 22, 25, and 27) and from experimental exposure assays using specific
732 carcinogens (7,12-dimethylbenz[a]anthracene (DMBA), urethane, and aristolochic acid
733 (AA)).

1
2
3 734 **Figure 4:** GA signature in human primary cancer genome PCAWG data. **(A)** Comparison of
4 735 COSMIC signature 4 with two experimentally derived signatures (B[a]P_Exp = signature in
5 736 clones from benzo[a]pyrene treated HMEC cells; GA_Exp = signature in clones from
6 737 glycidamide-treated MEF cells). Cosine similarity between the T>N (adenine) components of
7 738 signature 4 and GA signature is shown to the right. **(B)** Transcription strand bias analysis for
8 739 the six mutation types underlying the signatures in panel A). For each mutation type, the
9 740 number of mutations occurring on the transcribed (T) and non-transcribed (N) strand is
10 741 shown on the left Y-axis. The significance is expressed as $-\log_{10}(\text{p-value})$ indicated on the
11 742 right Y-axis. *** $p < 10^{-8}$; ** $p < 10^{-4}$; * $p < 10^{-2}$. **(C)** Scatter plots show reconstruction of
12 743 COSMIC signature 4 using B[a]P- and glycidamide- experimental mutational signatures in
13 744 lung adenocarcinoma, lung squamous cell carcinoma and hepatocellular carcinoma from the
14 745 PCAWG data set. **(D)** mSigAct analysis identifies the assignment and the contributions of
15 746 mutational signatures (including the experimental signature_GA_Exp (red) and
16 747 signature_B[a]P_Exp (blue)) to the mutation burden of a total of 101 PCAWG lung and liver
17 748 tumors identified as positive for the GA signature signal.

Table 1: Summary of cell lines, treatment conditions and *TP53*¹ mutation status.

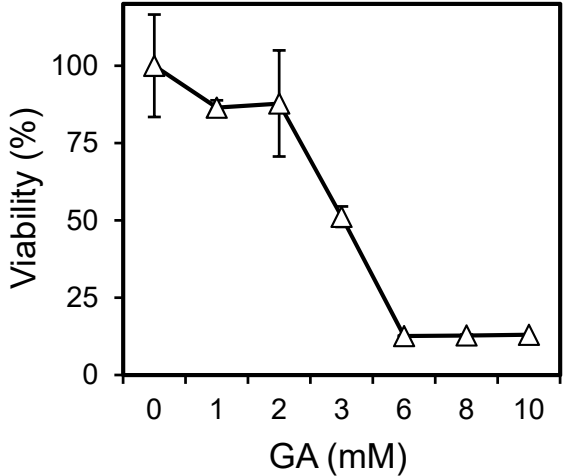
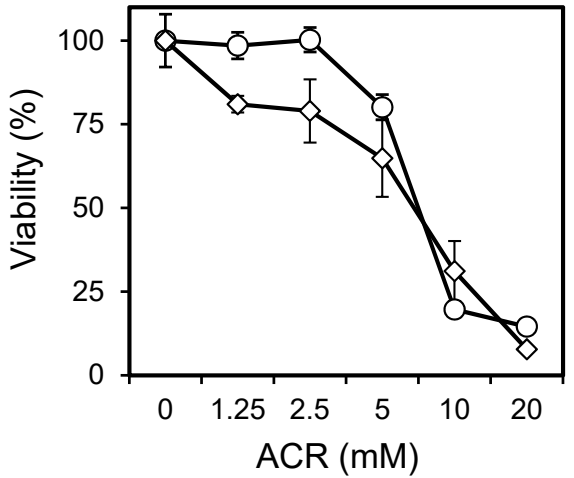
Sample ID	Embryo	Exposure	Conc. (mM)	Exposure duration (hrs)	coding DNA change ²	genomic DNA change ³	aa change	Codon 72 (rs1042522) ⁴
Prim_1	E210	-	-	-				Pro/Pro
Prim_2	E213	-	-	-				Arg/Pro
Prim_3	E214	-	-	-				Pro/Pro
Spont_1	E213	-	-	-				Arg/Pro
Spont_2	E214	-	-	-				Pro/Pro
Spont_3	E214	-	-	-				Pro/Pro
ACR_S9_1	E213	ACR	5	24				Arg/Pro
ACR_S9_2	E213	ACR	5	24				Arg/Pro
ACR_1	E213	ACR	10	24	c.881delA	g.7577057delT	p.E294fs	Arg/-
ACR_2	E213	ACR	10	24	c.818G>T	g.7577120C>A	p.R273L	Pro/-
ACR_3	E214	ACR	10	24	c.740A>T; c.839G>C	g.7577541T>A; g.7577099C>G	p.N247I; p.R280T	Pro/Pro
ACR_4	E214	ACR	10	24				Pro/Pro
ACR_5	E214	ACR	10	24				Pro/Pro
GA_1	E210	GA	3	24				Pro/Pro
GA_2	E210	GA	3	24				Pro/Pro
GA_3	E210	GA	3	24	c.309-310CC>TA	g.7579377-7579378GG>TA	[p.Y103Y; p.Q104K]	Pro/Pro
GA_4	E214	GA	3	24				Pro/Pro
GA_5	E214	GA	3	24				Pro/Pro

¹ human TP53 gene; ² NM_000546.4 coding sequence; ³ hg19 genomic coordinates; ⁴ human polymorphic site (rs1042522)

Prim = Primary cells; Spont = spontaneously immortalized clones; ACR = acrylamide-exposure derived clones; GA = glycidamide-exposure derived clones. Each exposure condition was carried out in two biological replicates (embryos). S9 = human S9 fraction; Pro = proline; Arg = arginine; Arg/- or Pro/- = loss of allele; fs = frameshift; aa = amino acid.

Figure 1

A



B

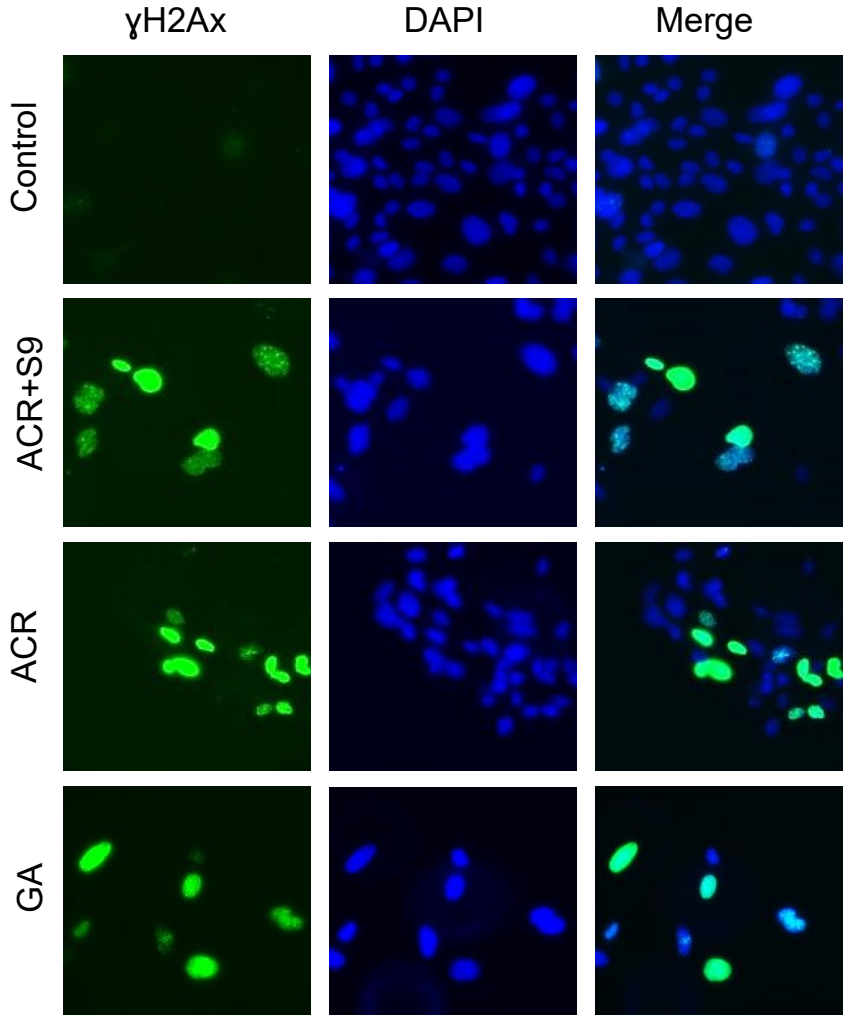
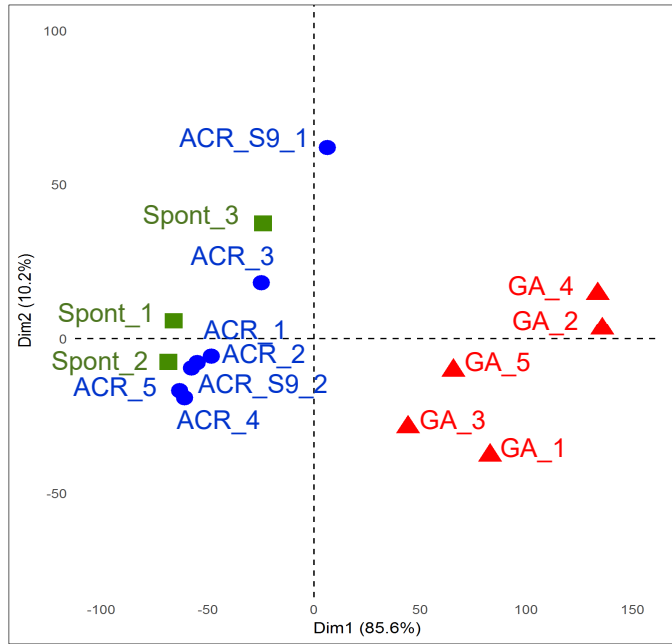
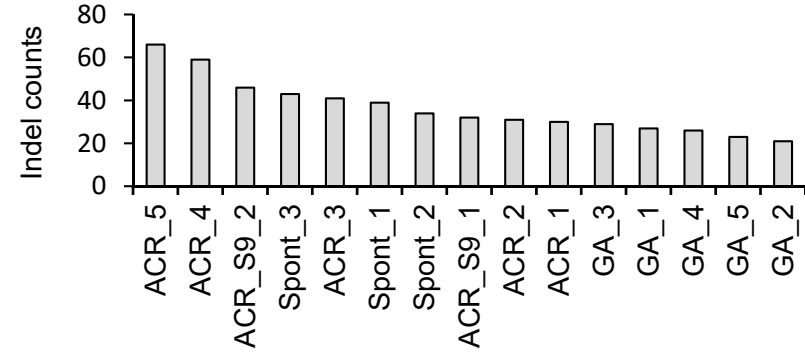


Figure 2

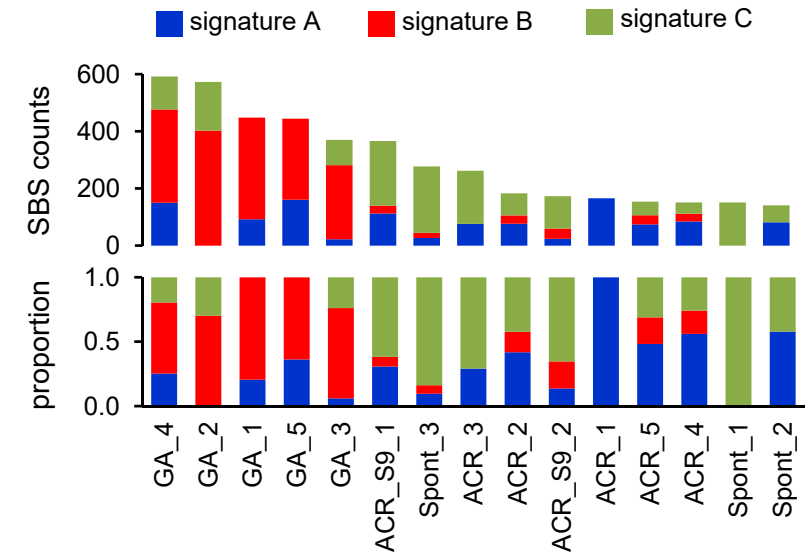
A



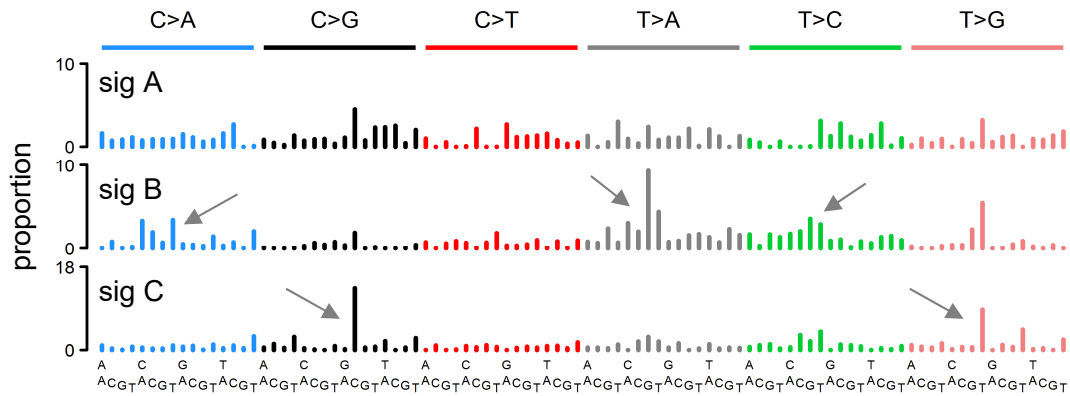
B



D



C



reconstruction accuracy

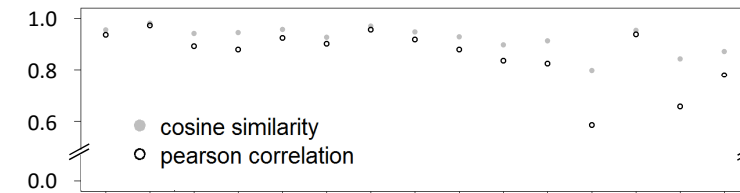


Figure 3

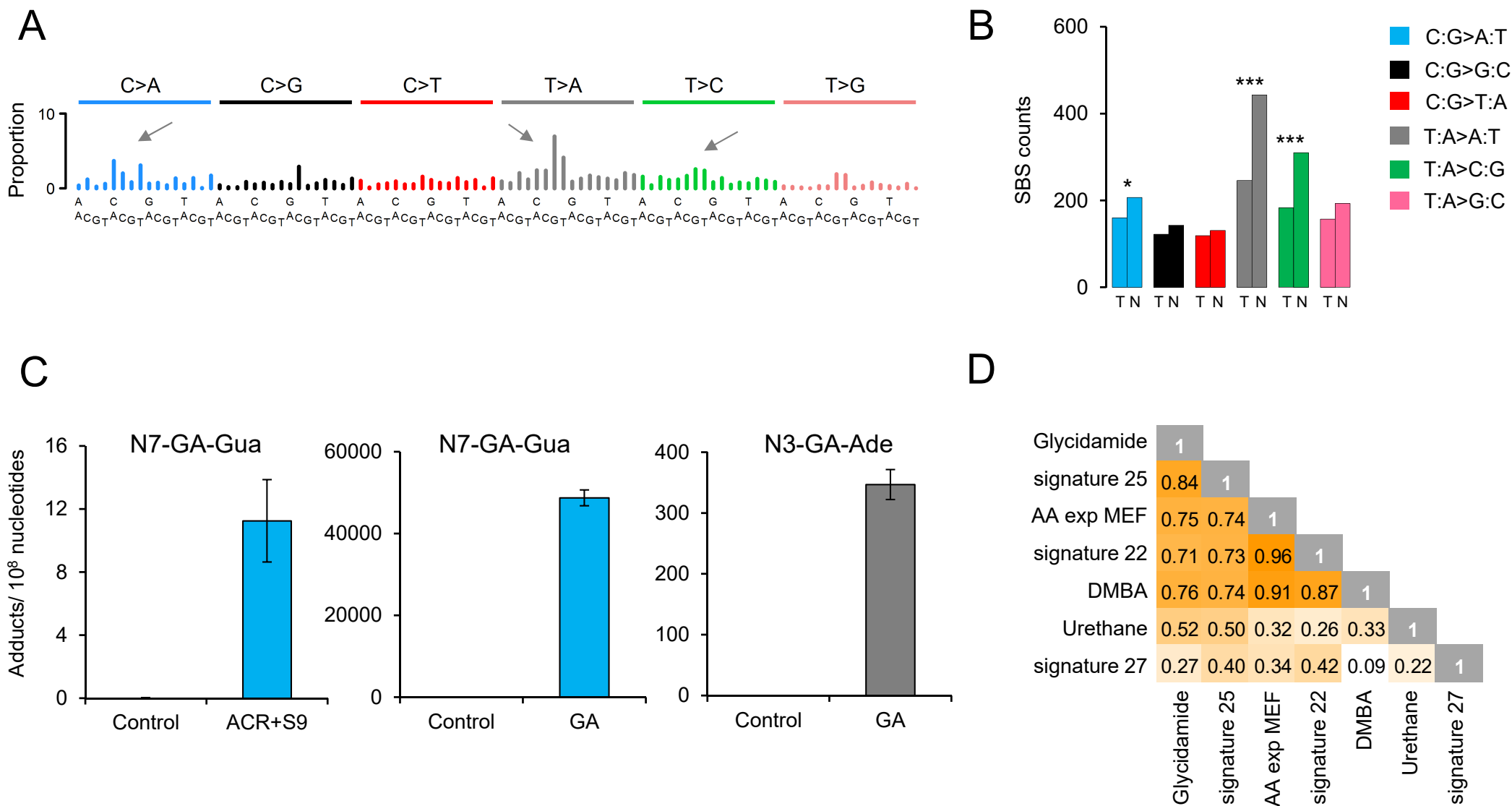
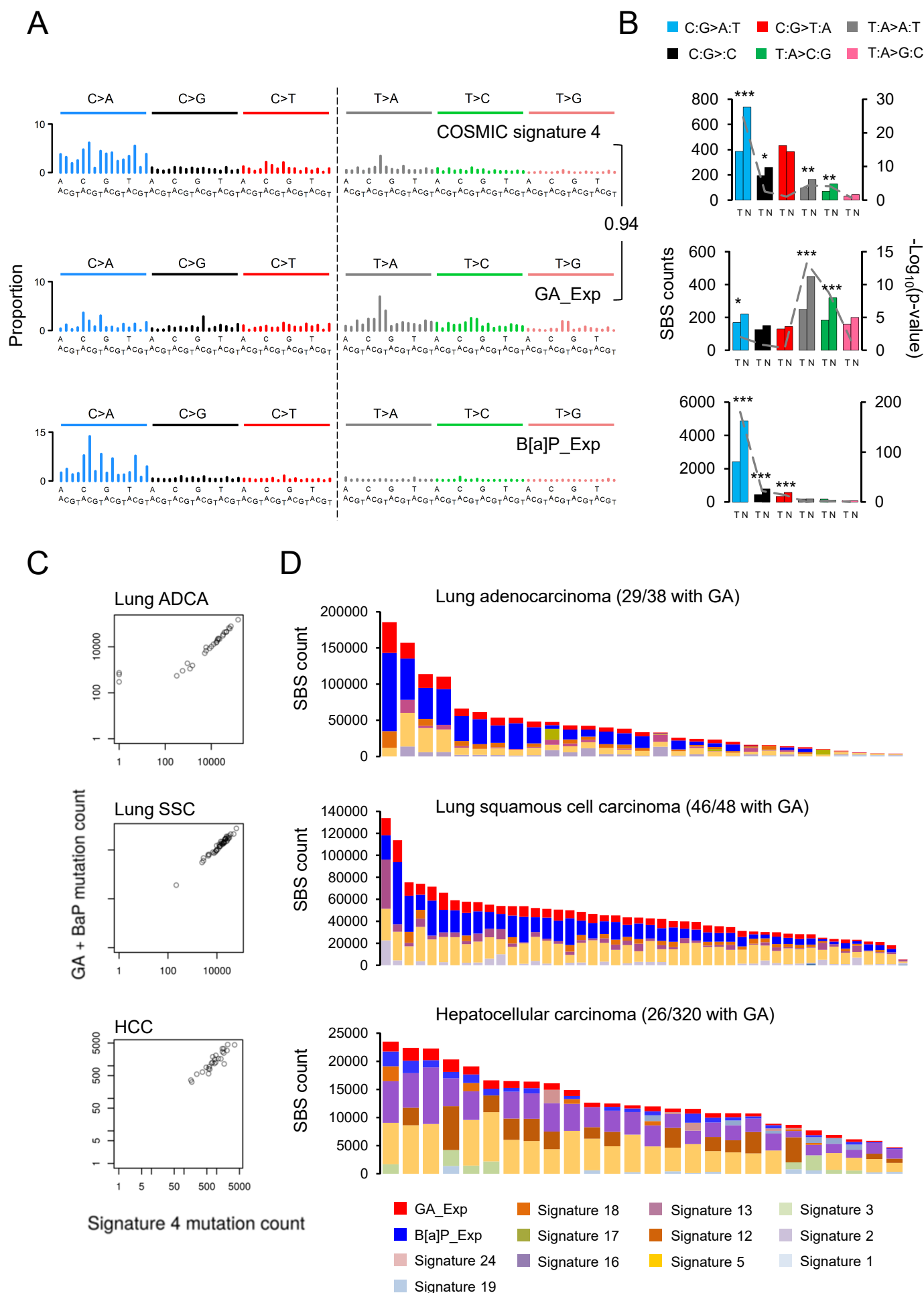


Figure 4



Supplementary Figure Legends

Supplementary Fig. S1: Comparison of different normalization and single-nucleotide variant calling strategies. Variant calling with respect to primary cell normalization. Venn diagrams show the overlap of variants called in glycidamide (GA)-derived clones after normalization to three different batches of primary cells (Prim_1, Prim_2, and Prim_3).

Supplementary Fig. S2: Growth curves of Hupki MEFs. Primary cells were either left untreated (Spont) or were exposed to acrylamide (ACR±S9) or glycidamide (GA). X-axis represents days in culture. Y-axis represents the cumulative doubling populations. The dashed vertical line represents the threshold of p-value < 0.05. Arrow: compound exposure; S*: senescence; SBI: senescence bypass/immortalization.

Supplementary Fig. S3: Mutation spectra derived from exome sequencing data from immortalized Hupki MEF clones derived from exposure to (A) acrylamide (ACR) or (B) glycidamide (GA), or (C) by spontaneous immortalization (Spont). X-axis represents the trinucleotide sequence context. Y-axis represents the frequency distribution of the mutations in each context.

Supplementary Fig. S4: Illustration of the transcription strand bias derived from the analysis of exome sequencing data from immortalized Hupki MEF cell lines. GA: glycidamide-derived clones; ACR: acrylamide-derived clones; Spont: spontaneously immortalized clones. The six mutation types are represented by different colors. For each mutation type, the number of mutations occurring on the transcribed (T) and non-transcribed (N) strand, as well as the p-values for strand bias is shown on the y-axes. The dashed grey line in each graph indicates the p-values for strand bias for each mutation type. The horizontal, dashed black line represents a significance threshold of p < 0.05.

Supplementary Fig. S5: Distribution of mutations based on their allelic frequencies in the five glycidamide (GA)-derived clones (left). Mutations in individual cell lines were ranked and plotted based on decreasing allelic frequency. Percentage of mutations with allelic frequency between 25% and 75% is indicated. Percentages of the six mutation types, color-coded, among all mutations identified in GA clones (right). The overall mutation number for each sample is indicated in the centre of the pie chart.

1
2
3 **Supplementary Fig. S6:** Mutation type and mutation spectra analysis with respect to variant
4 allele frequency (VAF). The analysis was carried out using exome sequencing data from
5 immortalized Hupki MEF clones derived from exposure to glycidamide. Top left: Mutation counts
6 were stratified into three VAF bins ([0-33% = low VAF]; [34-66% = medium VAF]; [67-100% =
7 high VAF]). Top right: The relative contribution of the six mutation types to the overall number of
8 mutations in each VAF bin is shown on the y-axis. Bottom panel: Mutation spectra (left) and
9 strand bias (right) analysis for the different VAF bins. Mutation spectra analysis: X-axis
10 represents the trinucleotide sequence context. Y-axis represents the frequency distribution of
11 the mutations. The counts for each mutation type are indicated in parentheses. Strand bias
12 analysis: For each mutation type, the number of mutations occurring on the transcribed and
13 non-transcribed strand is shown on the y-axis. T: transcribed strand; N: non-transcribed strand.

21 **Supplementary Fig. S7:** The 'baiting' clean-up of background signature 17 and the
22 quantification of its efficiency. COSMIC signature 17 (top track) marked by the arrows observed
23 in GA mutation spectra as well as in GA-mutational signature before and after baiting (clean).
24 The heat-map table on the right indicates the final proportionate reduction of signature 17-
25 specific peaks after re-running the NMF with signature 17-rich ICGC ESAD data sets listed in
26 the Supplementary Materials and Methods section.

31 **Supplementary Fig. S8:** (A) The structures of N7-GA-Gua and N3-GA-Ade adducts analyzed
32 by LC-MS/MS. (B) Representative multiple-reaction monitoring chromatograms (relative signal
33 intensity vs time) for N7-GA-Gua and N3-GA-Ade adducts in DNA from ACR treatment in the
34 presence of S9 fraction (ACR+S9) and GA-treated (GA) primary Hupki MEF. Internal standards
35 (IS) were added in amounts of 1000 fmol for N7-GA-Gua and 200 fmol for N3-GA-Ade.

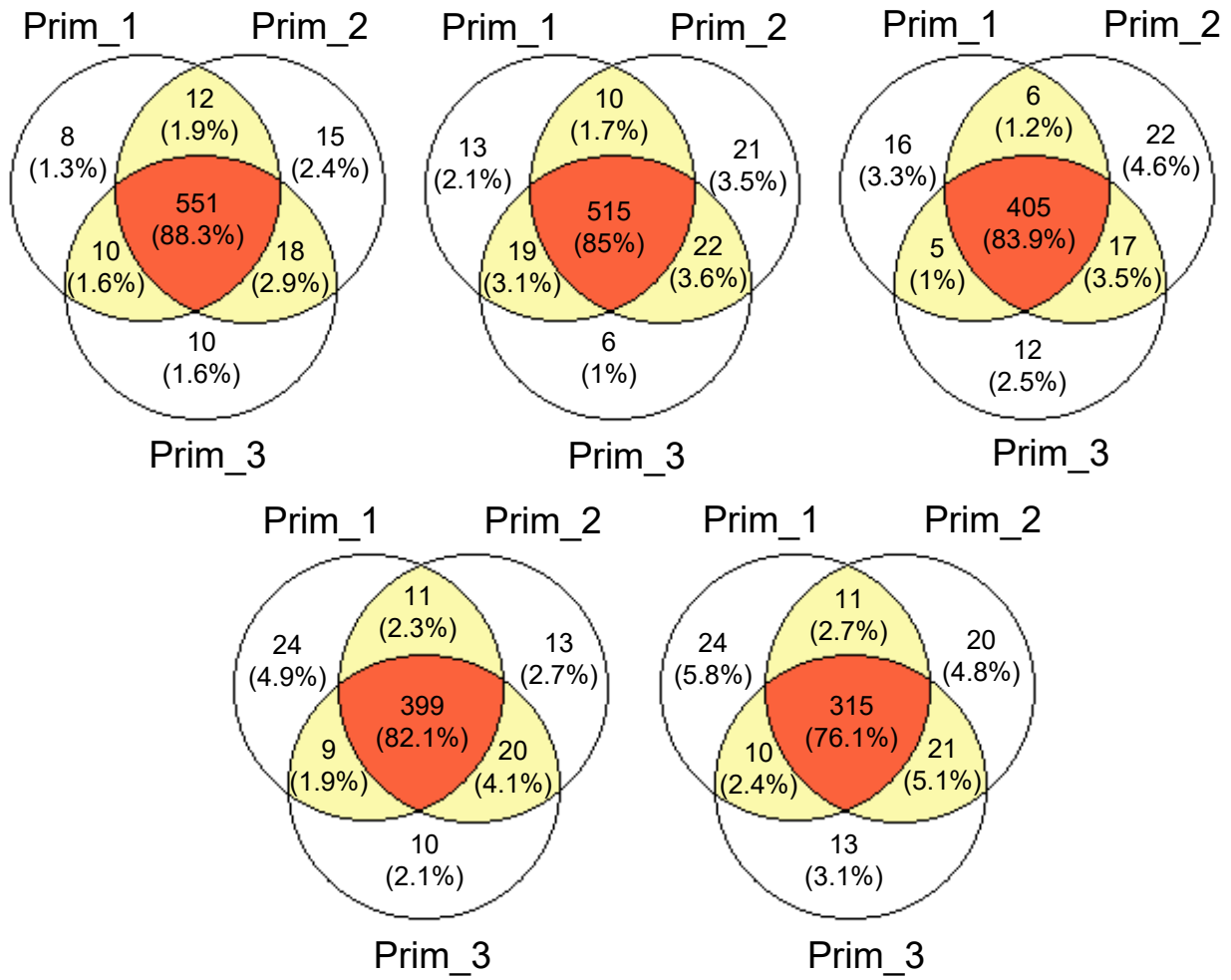
40 **Supplementary Fig. S9:** T:A>A:T enriched mutational signatures used for cosine similarity
41 analysis (see Fig. 3D). The individual signatures were originally derived from human cancer
42 sequencing data or experimental models (animal bioassays, cell lines) of carcinogen exposure.
43 X-axis represents the trinucleotide sequence context. Y-axis represents the frequency
44 distribution of the mutations. The predominant trinucleotide context for T:A>A:T mutations is
45 indicated by an arrow in the signature landscape. AA: aristolochic acid; DMBA: 7,12-
46 dimethylbenz[a]anthracene.

52 **Supplementary Fig. S10:** (A) Scatter plots show the measure of correlation of the GA-
53 signature versus B[a]P-signature (used to reconstruct COSMIC signature 4) in PCAWG lung
54 adenocarcinomas (ADCA), lung squamous cell carcinomas (SCC) and hepatocellular
55

1
2
3 carcinomas (HCC). **(B)** Bar-plots representing the proportion of the assignment of the
4 experimental GA_Exp and B[a]P_Exp signatures in lung adenocarcinomas, lung squamous cell
5 carcinomas and hepatocellular carcinomas from the PCAWG data set. The asterisk denotes
6 liver HCC samples harboring GA-signature only (no B[a]P-signature detected), indicating
7 possible dietary or occupational exposure. Full list of these samples is accessible from Suppl.
8 Table S5.
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

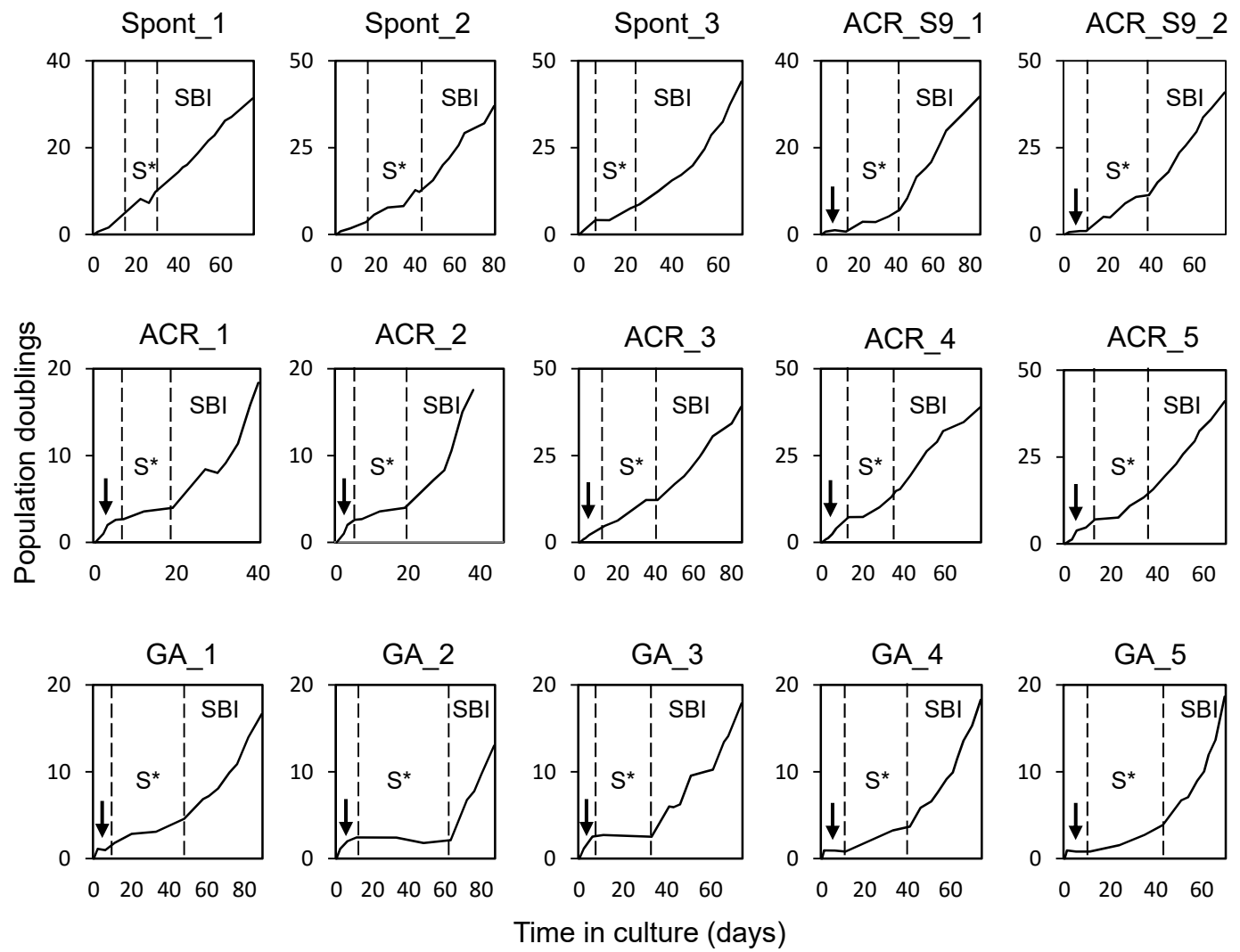
For Peer Review

Suppl. Fig. S1



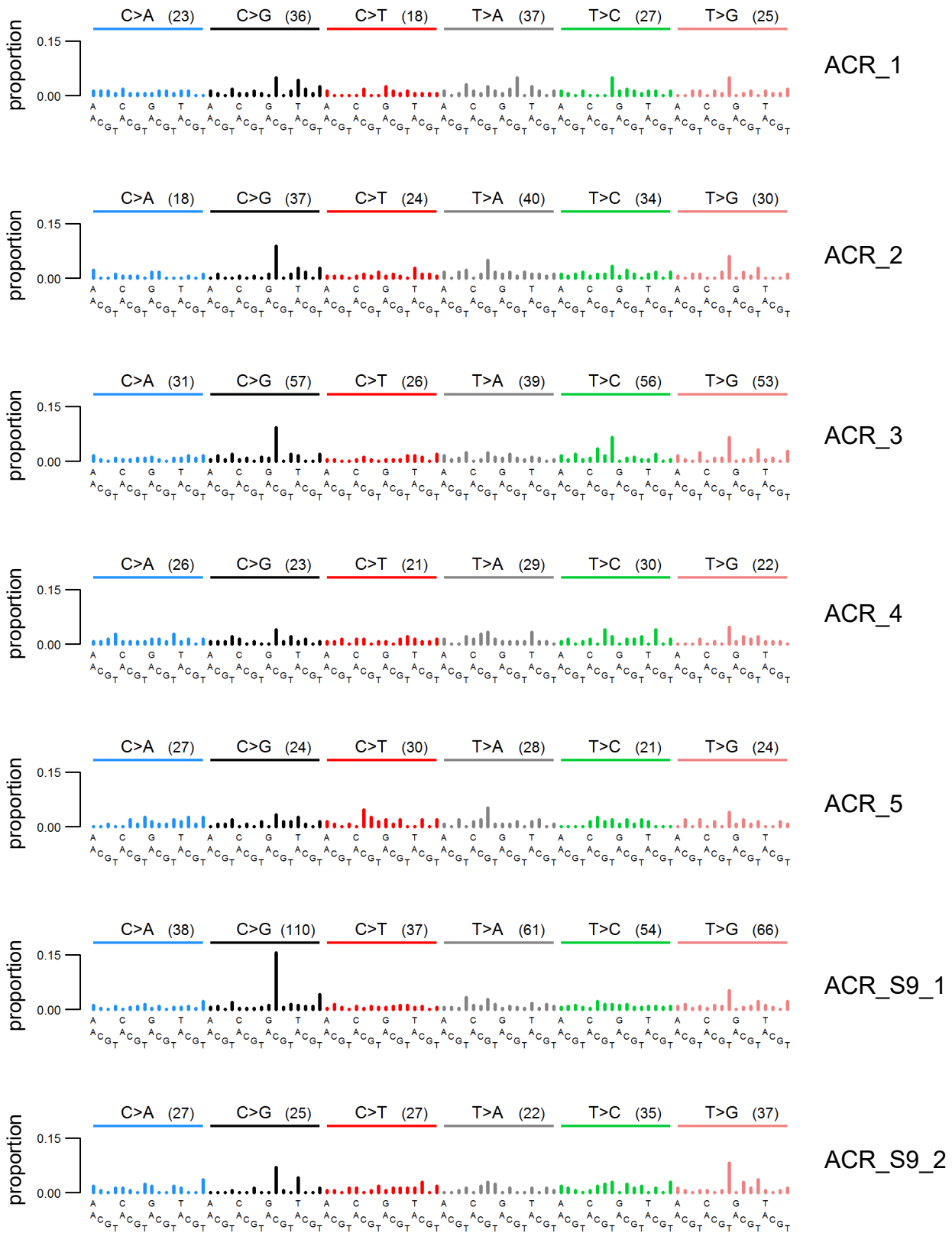
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Suppl. Fig. S2



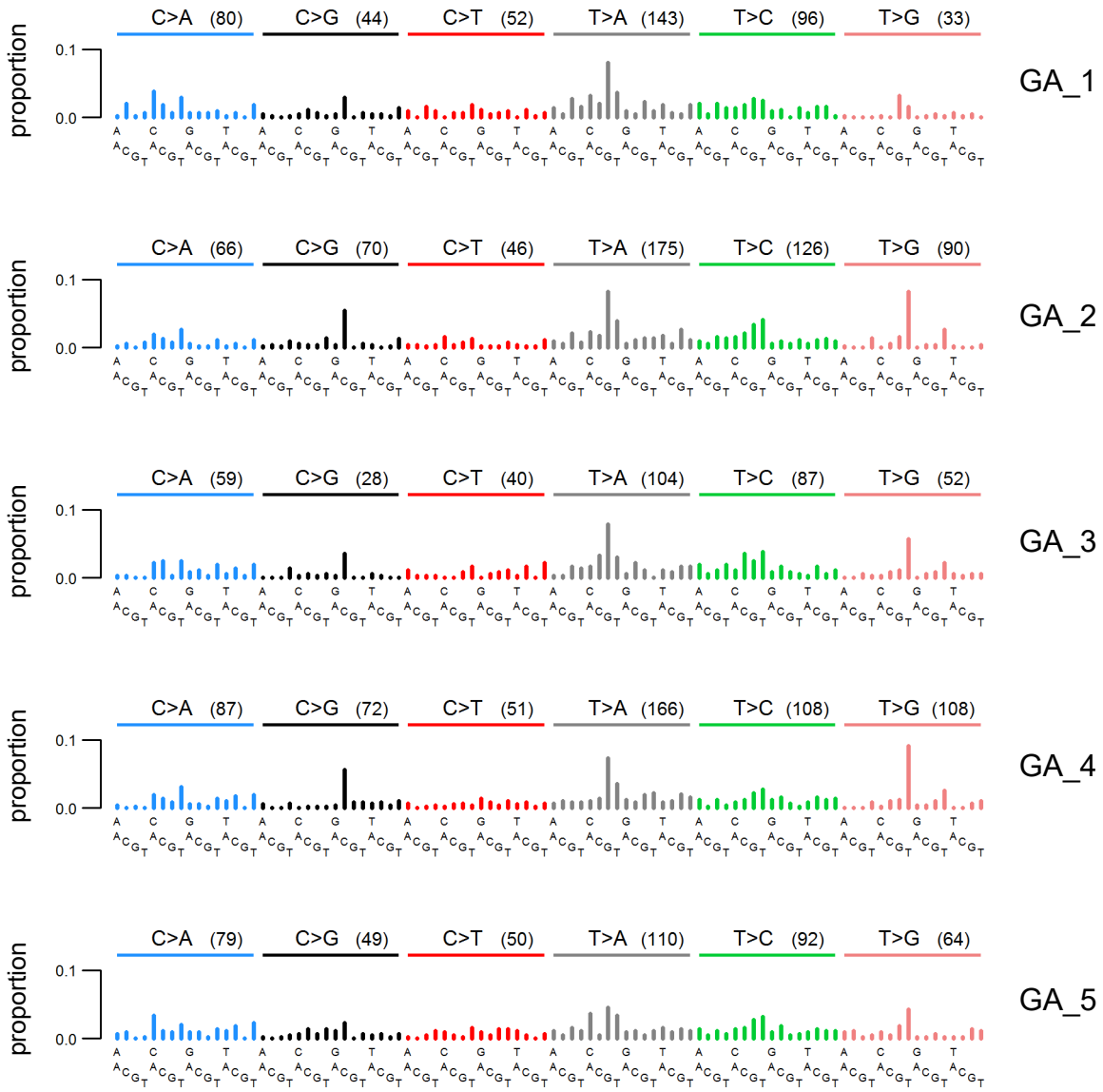
Suppl. Fig. S3

A



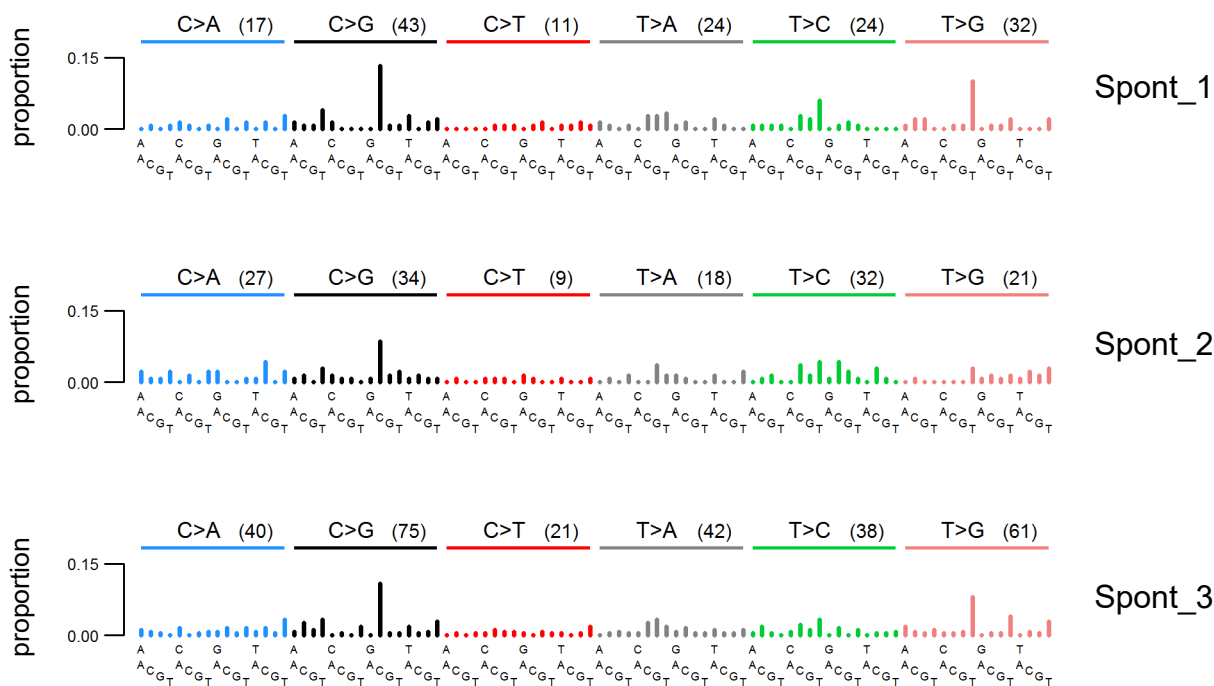
Suppl. Fig. S3

B

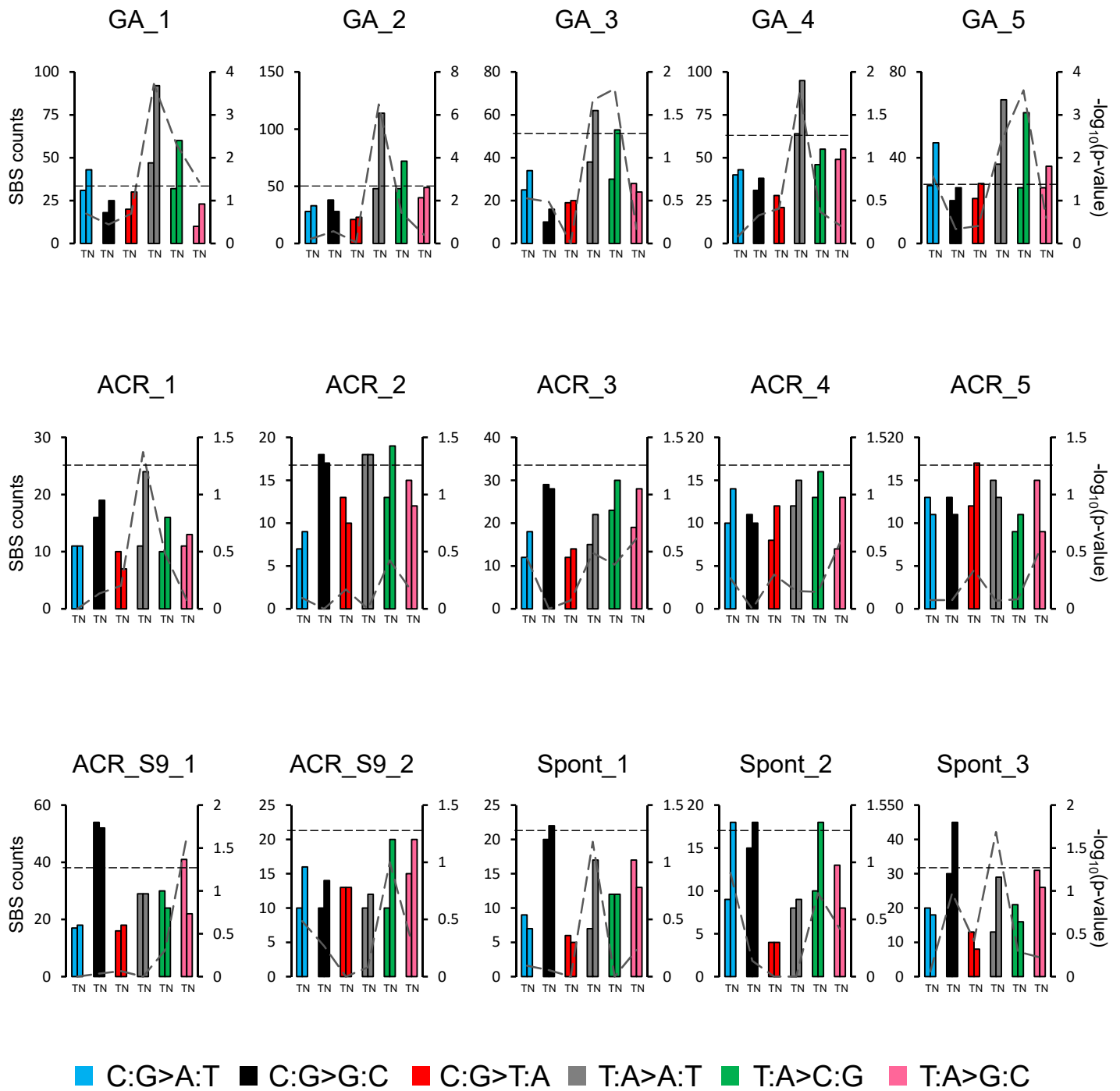


Suppl. Fig. S3

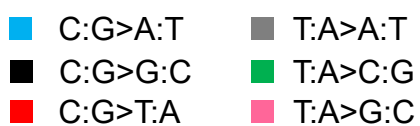
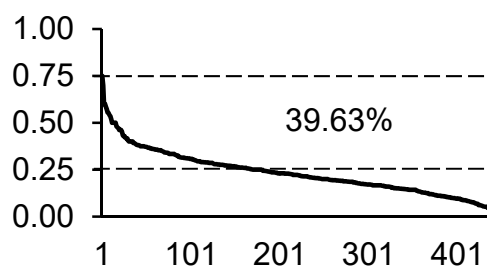
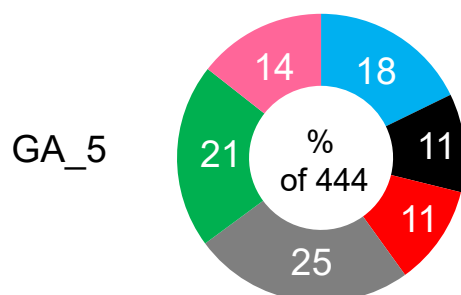
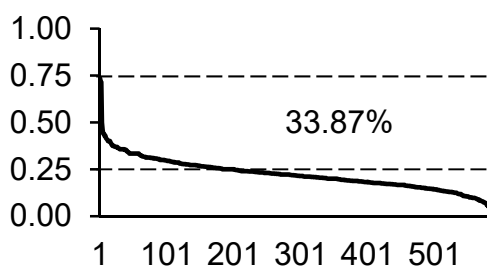
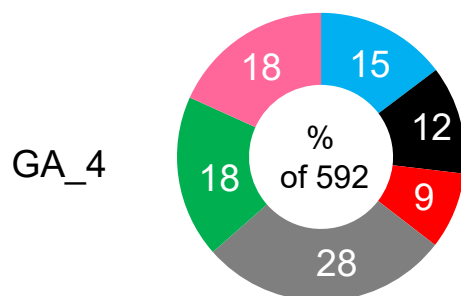
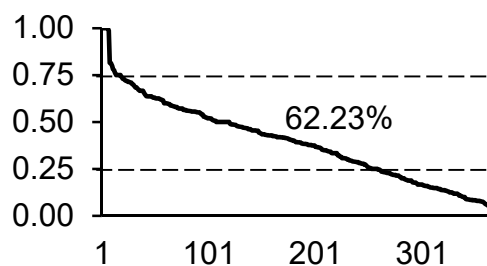
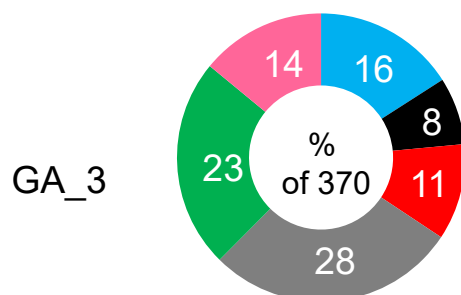
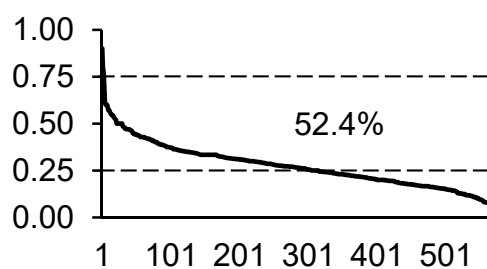
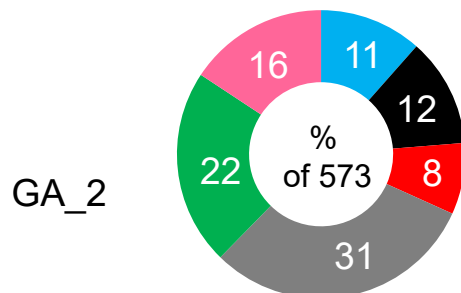
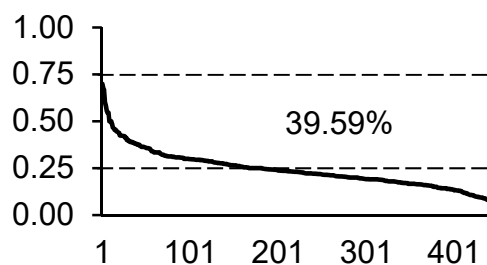
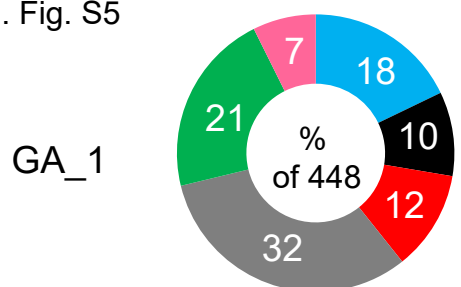
C



Suppl. Fig. S4



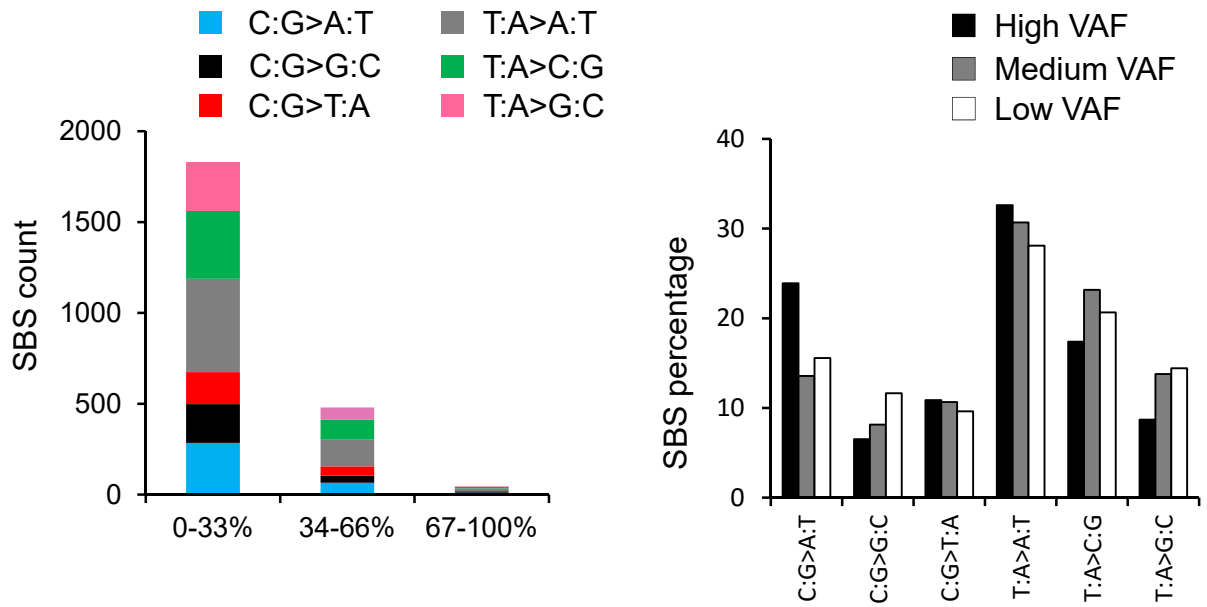
Suppl. Fig. S5



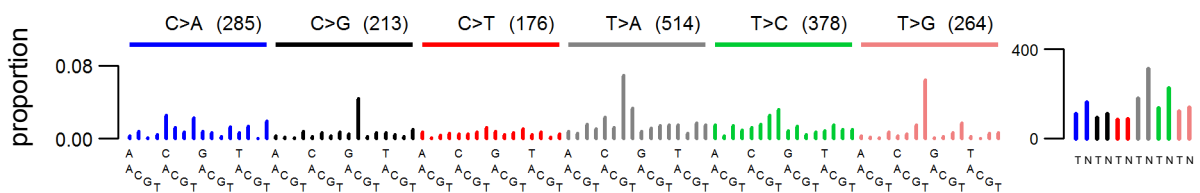
Rank

Allelic Frequency

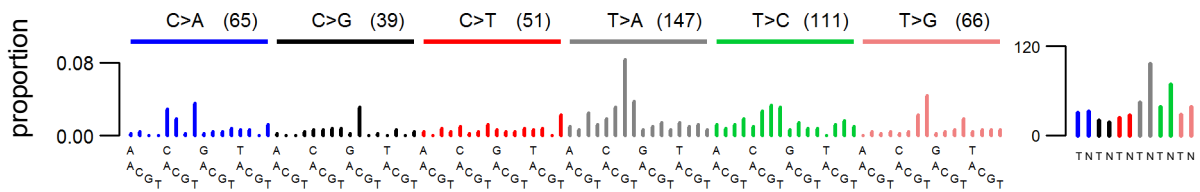
Suppl. Fig. S6



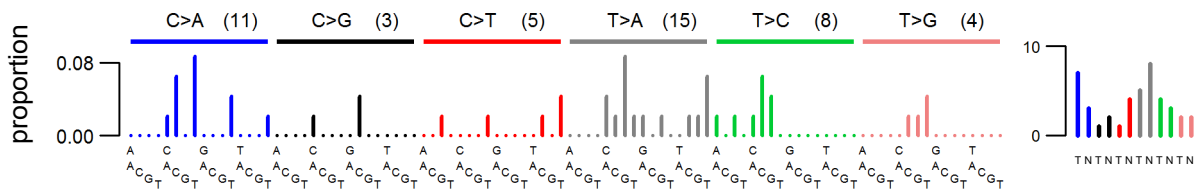
0-33%



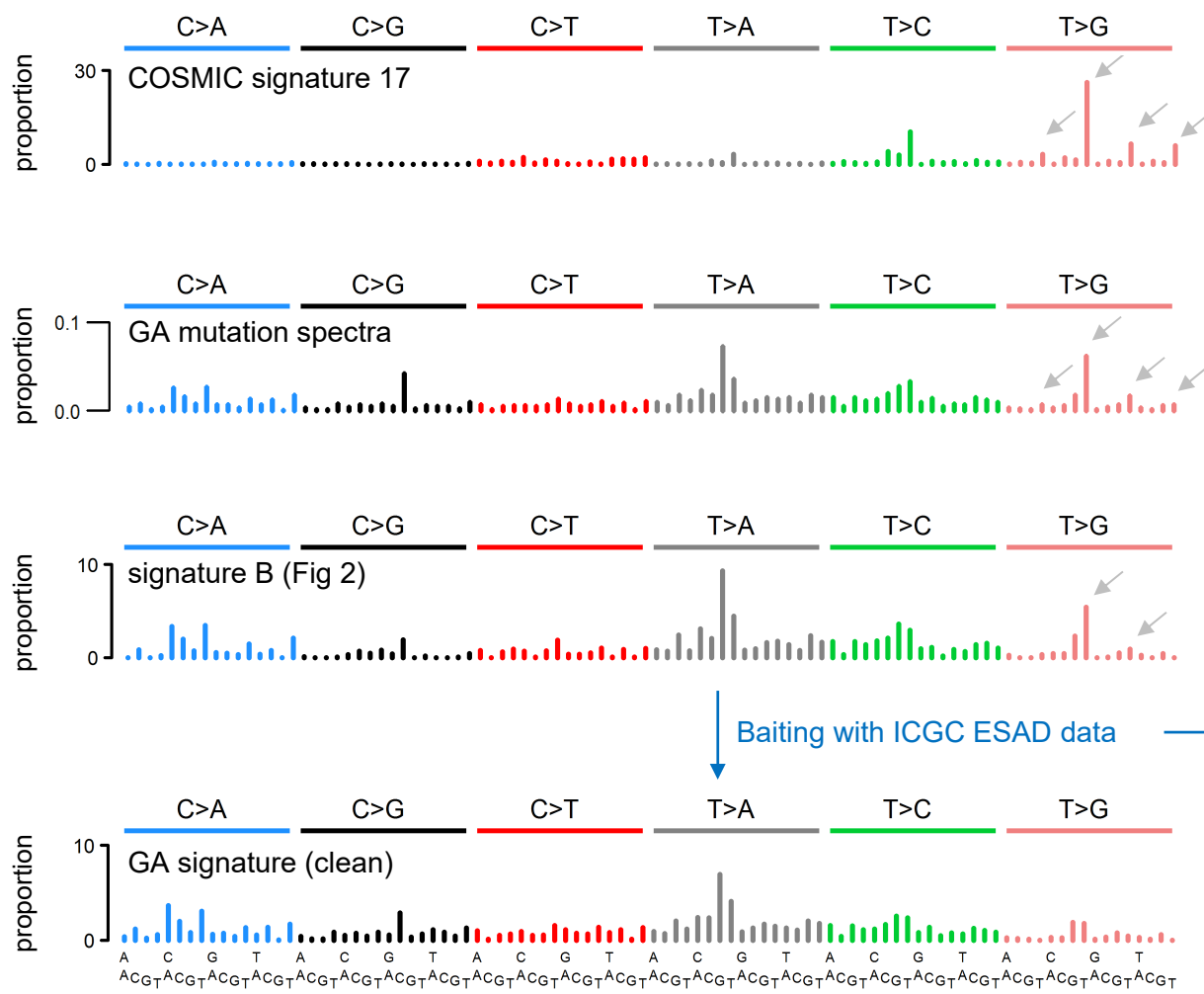
34-66%



67-100%



Suppl. Fig. S7

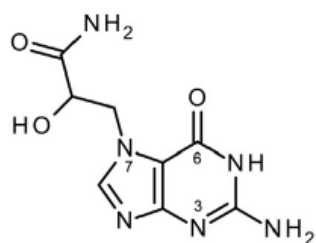


Sig. 17 peaks Decrease (%)	
C[T>A]C	0.0
C[T>A]G	25.7
C[T>A]T	8.0
C[T>C]C	22.4
C[T>C]G	29.7
C[T>C]T	20.0
A[T>G]T	100.0
C[T>G]C	42.9
C[T>G]G	19.1
C[T>G]T	68.2
G[T>G]T	52.7
T[T>G]T	0.0

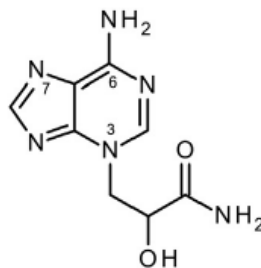
Suppl. Fig. S8

A

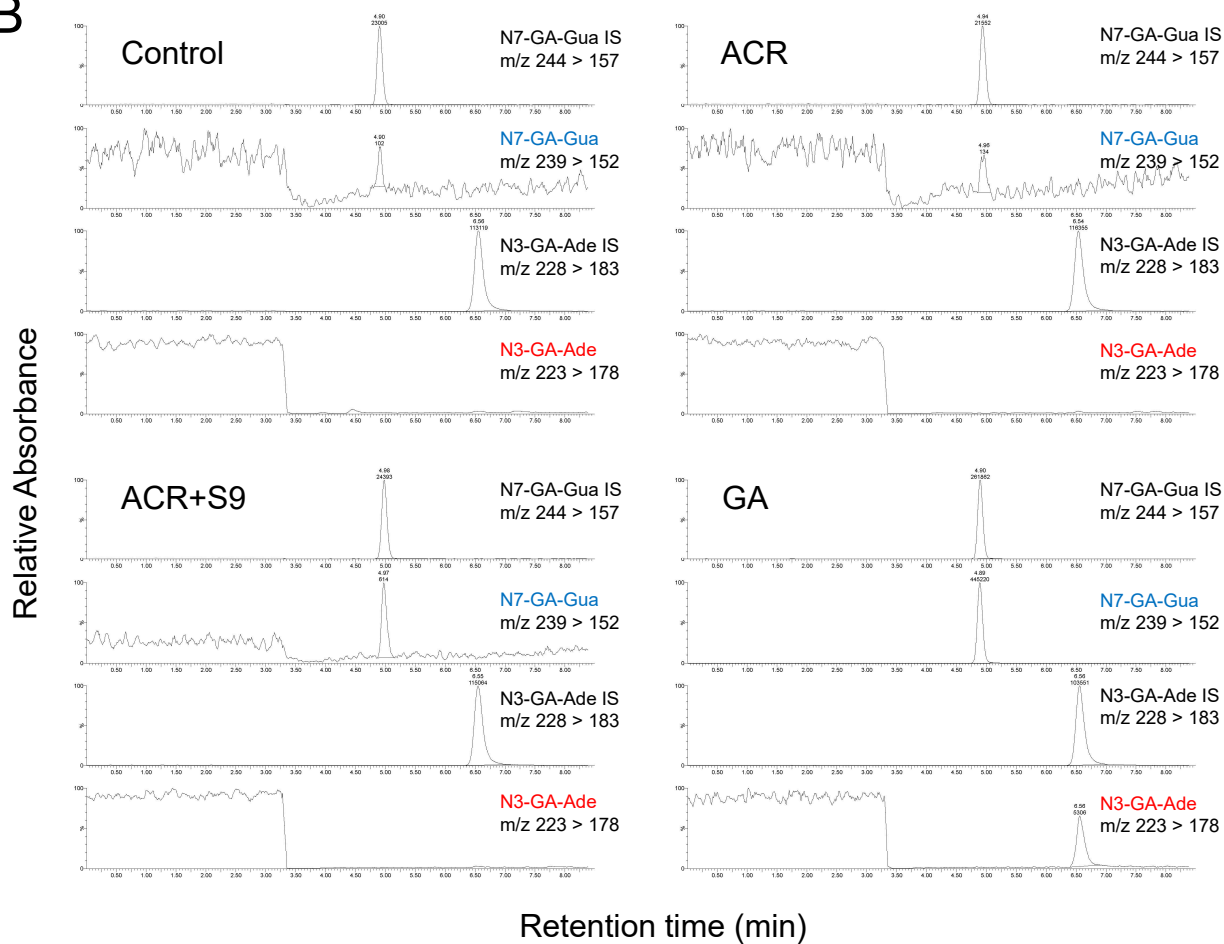
N7-GA-Gua



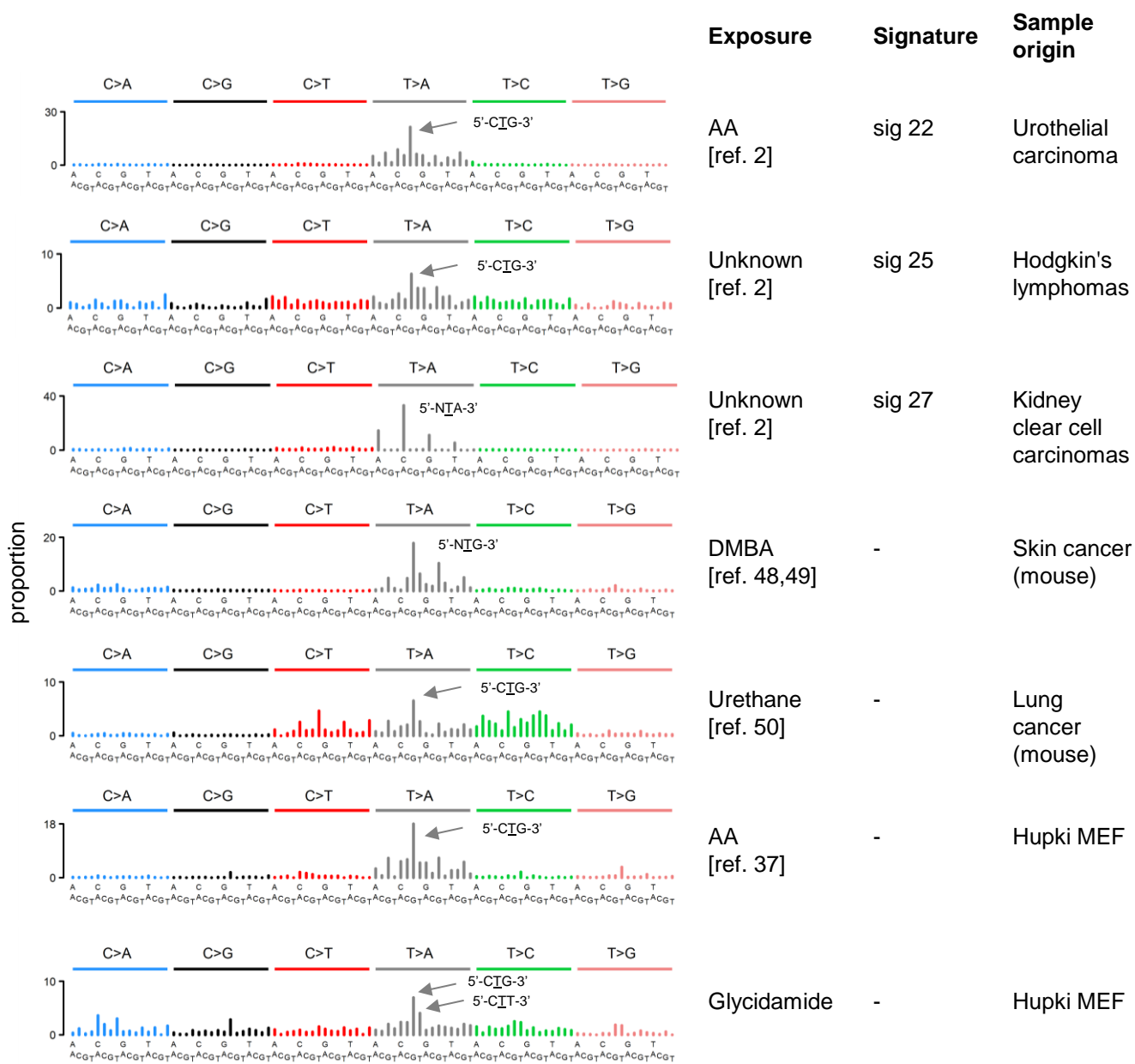
N3-GA-Ade



B



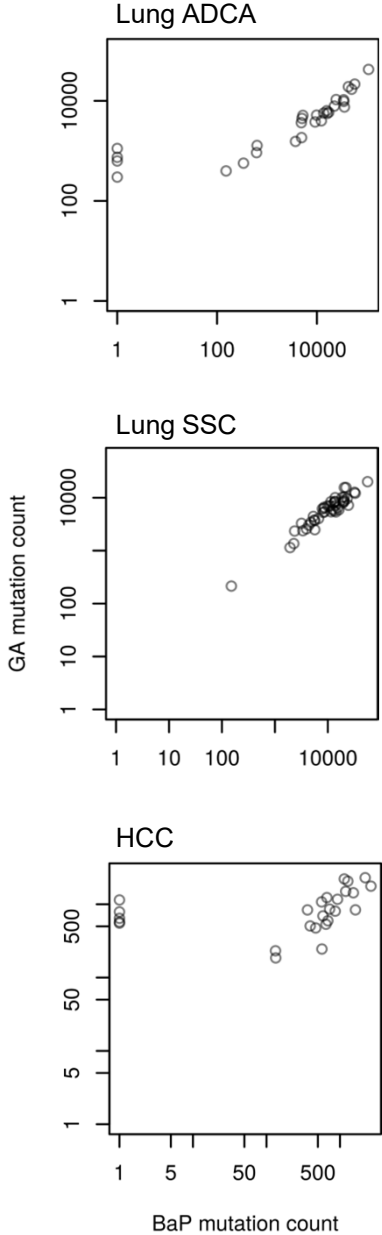
Suppl. Fig. S9



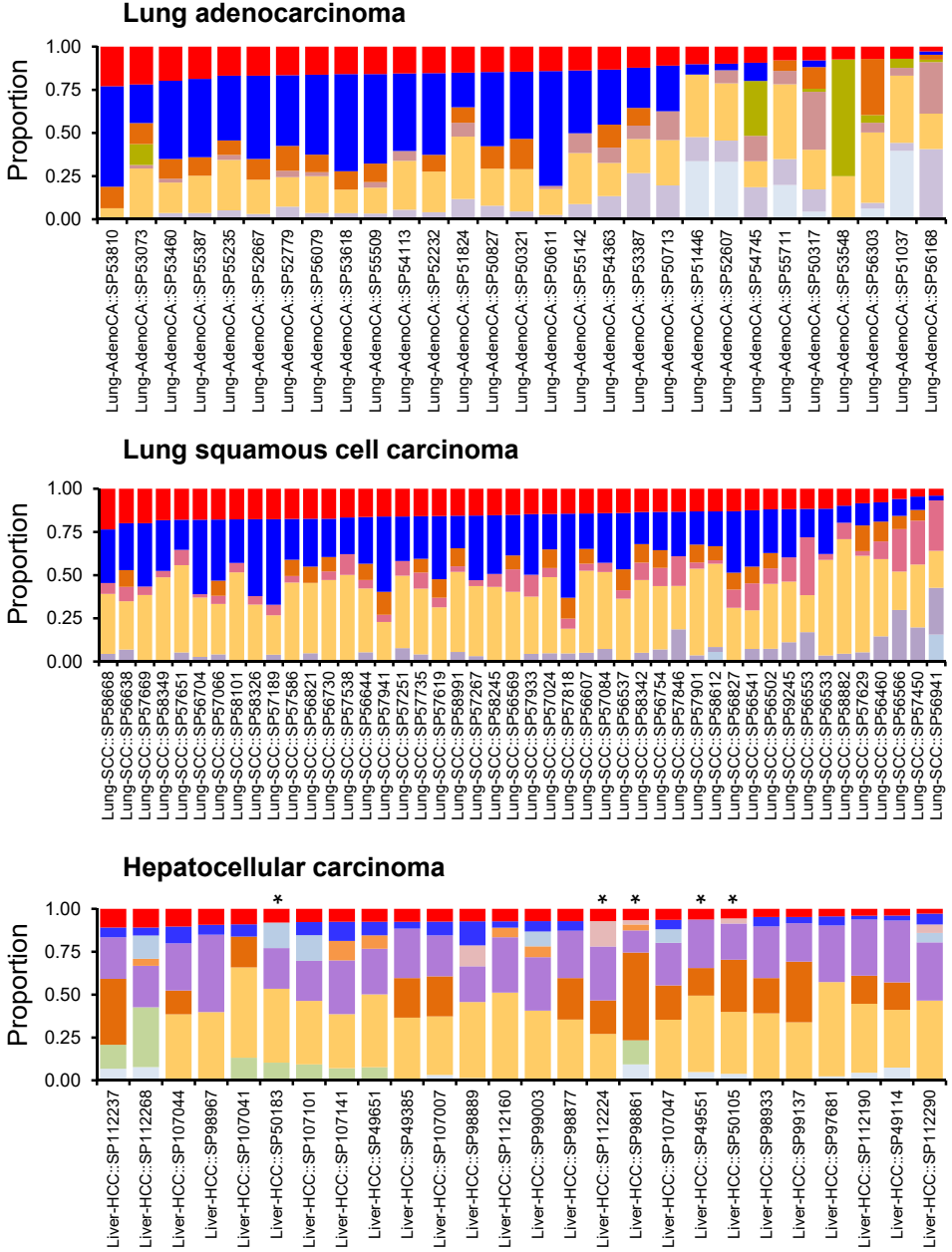
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56

Supplementary Figure S10

A



B



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Supplementary Materials and Methods

DNA adduct analysis

The DNA was isolated from the cells using standard digestion with proteinase K, followed by phenol-chloroform extraction and ethanol precipitation. The DNA was subsequently treated with RNase A and T1, extracted with phenol-chloroform, and reprecipitated with ethanol. N7 GA-Gua and N3 GA-Ade were released by neutral thermal hydrolysis for 15 minutes, using Eppendorf Thermomixer R (Eppendorf North America) set to 99 °C. The samples were filtered through Amicon 3K molecular weight cutoff filters (Merck Millipore) to separate the adducts from the intact DNA.

TP53 genotyping

The following are the *TP53* primers used for amplicon sequencing of mutations accumulated in human *TP53* of the Hupki MEFs. The sequences are presented in 5' to 3' orientation: Exon 4: fwd – TGCTCTTTTCACCCATCTAC, rev – ATACGGCCAGGCATTGAAGT; Exons 5-6: fwd – TGTTCACTTGTGCCCTGACT, rev – TTAACCCCTCCTCCCAGAGA; Exon 7: fwd – CTTGCCACAGGTCTCCCC, rev – CACTTGCCACCCTGCACA; Exon 8: fwd – TCCTTACTGCCTCTTGCTTCTCTT; rev – CCAAGGGTGCAGTTATGCCT. Sequences and their alterations were analyzed using the CodonCode Aligner software.

Processing of WES data

Prior to variant calling, recalibrated .bam files were interrogated for imbalanced base mismatch distribution between Read 1 and Read 2 sequences. We used the DNA damage estimator tool (as per (1); (<https://github.com/Ettwiller/Damage-estimator>)) to measure the Global Imbalance Value (GIV) score and to exclude sequencing-related DNA damage and artefacts due to oxidative damage that can confound the determination of treatment-specific variants. The MutSpec suite included tools for annotation of the vcf files with Annovar and variant filtering to remove dbSNP142 contents, segmental duplicates, repeats, and tandem repeat regions. Finally, to maximize the chance of robust variant calls and to exclude potential unfiltered single nucleotide polymorphisms (SNP), we considered only variants unique to each sample.

Bioinformatics and statistical analyses

The following are the International Cancer Genome Consortium (ICGC) esophageal carcinoma patient data (2,3) that were used in the step of cleaning the experimental signature from the COSMIC signature 17 signal: ESAD-UK-SP119768.hg19; ESAD-UK-SP191660.hg19; ESAD-UK-SP111113.hg19; ESAD-UK-SP111173.hg19; ESAD-UK-SP192267.hg19; ESAD-UK-SP111026.hg19; ESAD-UK-SP192494.hg19; ESAD-UK-SP111019.hg19; ESAD-UK-SP111058.hg19.

References

1. Chen, L., *et al.* (2017) DNA damage is a pervasive cause of sequencing errors, directly confounding variant identification. *Science*, **355**, 752-756.
2. Secrier, M., *et al.* (2016) Mutational signatures in esophageal adenocarcinoma define etiologically distinct subgroups with therapeutic relevance. *Nature Genetics*, **48**, 1131-1141.
3. Cancer Genome Atlas Research Network (2017) Integrated genomic characterization of oesophageal carcinoma. *Nature*, **541**, 169-175.

Title

Experimental analysis of exome-scale mutational signature of glycidamide, the reactive metabolite of acrylamide

Authors

Maria Zhivagui¹, Maude Ardin¹, Alvin W. T. Ng^{2,3,4}, Mona I. Churchwell⁵, Manuraj Pandey¹, Stephanie Villar¹, Vincent Cahais⁶, Alexis Robitaille⁷, Liacine Bouaoun⁸, Adriana Heguy⁹, Kathryn Guyton¹⁰, Martha R. Stampfer¹¹, James McKay¹², Monica Hollstein^{1,13,14}, Magali Olivier¹, Steven G. Rozen^{2,3}, Frederick A. Beland⁵, Michael Korenjak¹ and Jiri Zavadil¹

Affiliations

¹ Molecular Mechanisms and Biomarkers Group, International Agency for Research on Cancer, Lyon 69008, France

² Centre for Computational Biology, Duke-NUS Medical School, Singapore 169857, Singapore

³ Program in Cancer and Stem Cell Biology, Duke-NUS Medical School, 169857, Singapore

⁴ NUS Graduate School for Integrative Sciences and Engineering, 117456, Singapore

⁵ Division of Biochemical Toxicology, National Center for Toxicological Research, Jefferson, AR 72079, USA

⁶ Epigenetics Group, International Agency for Research on Cancer, Lyon 69008, France

⁷ Infections and Cancer Biology Group, International Agency for Research on Cancer, Lyon 69008, France

⁸ Environment and Radiation Section, International Agency for Research on Cancer, Lyon 69008, France

⁹ Department of Pathology and Genome Technology Center, New York University, Langone Medical Center, New York, NY 10016, USA

¹⁰ IARC Monographs Section, International Agency for Research on Cancer, Lyon 69008, France

¹¹ Biological Systems and Engineering Division, Lawrence Berkeley National Laboratory, Berkeley, CA, 94720, USA

¹² Genetic Cancer Susceptibility Group, International Agency for Research on Cancer, Lyon 69008, France

¹³ Deutsches Krebsforschungszentrum, 69120 Heidelberg, Germany

¹⁴ Faculty of Medicine and Health, University of Leeds, LIGHT Laboratories, Leeds LS2 9JT, United Kingdom

Keywords: Acrylamide, glycidamide, DNA adducts, massively parallel sequencing, mutational signatures

Correspondence:

ZavadilJ@iarc.fr and/or KorenjakM@iarc.fr

Abstract

Acrylamide, a probable human carcinogen, is ubiquitously present in the human environment, with sources including heated starchy foods, coffee and cigarette smoke. Humans are also exposed to acrylamide occupationally. Acrylamide is genotoxic, inducing gene mutations and chromosomal aberrations in various experimental settings. Covalent haemoglobin adducts were reported in acrylamide-exposed humans and DNA adducts in experimental systems. The carcinogenicity of acrylamide has been attributed to the effects of glycidamide, its reactive and mutagenic metabolite capable of inducing rodent tumors at various anatomical sites. In order to characterize the pre-mutagenic DNA lesions and global mutation spectra induced by acrylamide and glycidamide, we combined DNA-adduct and whole-exome sequencing analyses in an established exposure-clonal immortalization system based on mouse embryonic fibroblasts. Sequencing and computational analysis revealed a unique mutational signature of glycidamide, characterized by predominant T:A>A:T transversions, followed by T:A>C:G and C:G>A:T mutations exhibiting specific trinucleotide contexts and significant transcription strand bias. Computational interrogation of human cancer genome sequencing data indicated that a combination of the glycidamide signature and an experimental benzo[a]pyrene signature are nearly equivalent to the COSMIC tobacco-smoking related signature 4 in lung adenocarcinomas and squamous cell carcinomas. We found a more variable relationship between the glycidamide- and benzo[a]pyrene-signatures and COSMIC signature 4 in liver cancer, indicating more complex exposures in the liver. Our study demonstrates that the controlled experimental characterization of specific genetic damage associated with glycidamide exposure facilitates identifying corresponding patterns in cancer genome data, thereby underscoring how mutation signature laboratory experimentation contributes to the elucidation of cancer causation.

64

A 40-word summary

Innovative experimental approaches identify a novel mutational signature of glycidamide, a metabolite of the probable human carcinogen acrylamide. The results may elucidate the cancer risks associated with exposure to acrylamide, commonly found in tobacco smoke, thermally processed foods and beverages.

70 Introduction

71 Cancer can be caused by chemicals, complex mixtures, occupational exposures, physical
72 agents, and biological agents, as well as lifestyle factors. Many human carcinogens show a
73 number of characteristics that are shared among carcinogenic agents (1). Different human
74 carcinogens may exhibit a spectrum of these key characteristics, and operate through
75 separate mechanisms to generate patterns of genetic alterations. Recognizable patterns of
76 genetic alterations or mutational signatures characterize carcinogens that are genotoxic.
77 Recent work shows that these DNA sequence changes can be expressed in simple
78 mathematical terms that enable mutational signatures to be extracted from thousands of
79 cancer genome sequencing data sets (2). Several of the over 30 identified mutational
80 signatures have been attributed to specific external exposures or endogenous factors
81 through epidemiological and experimental studies (2). However, about 40% of the current
82 signatures remain of unknown origin, and additional, thus far unrecognized, signatures are
83 likely to be defined in rapidly accumulating cancer genome data. Well-controlled
84 experimental exposure systems can thus help identify the underlying causes of known
85 orphan mutational signatures as well as define new patterns generated by candidate
86 carcinogens (reviewed in (3,4)).

87 Various diet-related exposures contribute to the human cancer burden. Examples
88 include contaminants in food or alternative medicines, such as aflatoxin B1 (AFB1) or
89 aristolochic acid (AA). The mutagenicity of these compounds is well-documented; AFB1
90 induces predominantly C:G>A:T base substitutions and AA causes T:A>A:T transversions.
91 The characteristic mutations coupled with information on the preferred sequence contexts in
92 which they are likely to arise allowed unequivocal association of exposure to AFB1 or AA
93 with specific subtypes of hepatobiliary or urological cancers, respectively (5-13).

94 Among dietary compounds with carcinogenic potential, acrylamide is of special
95 interest due to extensive human exposure. Important sources of exposure to acrylamide
96 include tobacco smoke (14), coffee (15), and a broad spectrum of occupational settings (16).
97 Dietary sources of acrylamide comprise carbohydrate-rich food products that have been
98 subject to heating at high temperatures. This is due to Maillard reactions, which involve
99 reducing sugars and the amino acid asparagine, present in potatoes and cereals (17). There
100 is sufficient evidence that acrylamide is carcinogenic in experimental animals (18,19) and it
101 has been classified as a probable carcinogen (Group 2A) by the International Agency for
102 Research on Cancer in 1994 (16). The association of dietary acrylamide exposure with
103 renal, endometrial and ovarian cancers has been explored in recent epidemiological studies
104 (20,21). However, accurate acrylamide exposure assessment in epidemiological studies
105 based on questionnaires has been difficult, and more direct measures of molecular markers,
106 such as hemoglobin adduct levels, may not yield conclusive findings on past exposures (22-

1
2
3 107 27). An improved understanding of its mechanism of action using well-controlled
4 108 experimental systems is critical for understanding the potential carcinogenic risk associated
5 109 with exposure.

6
7 110 Acrylamide undergoes oxidation by cytochrome P450, producing the reactive
8 111 metabolite glycidamide that is highly efficient in DNA binding due to its electrophilic epoxide
9 112 structure (28-30). The *Hras* mutation load in neoplasms of mice exposed to acrylamide or
10 113 glycidamide was found to be considerably higher in mice treated with glycidamide (31). This
11 114 finding is corroborated by a considerably higher mutation frequency in the *cII* reporter gene
12 115 of Big Blue mouse embryonic fibroblasts treated with glycidamide in comparison to
13 116 acrylamide (32,33). Mutation analysis in different experimental *in vivo* and *in vitro* models
14 117 using reporter genes showed an increased association of acrylamide and glycidamide
15 118 exposure with T:A>C:G transitions, as well as T:A>A:T and C:G>G:C transversion mutations
16 119 (31-36), whereas glycidamide exposure was also characterized by C:G>A:T transversions
17 120 (33). However, these proposed acrylamide- and glycidamide-specific mutation patterns were
18 121 based on limited mutation counts in reporter genes and thus do not reflect the complexity of
19 122 genome-wide distributions and profiles. Based on the limited data available thus far, it is not
20 123 possible to translate adequately the reported mutation types (T:A>C:G, T:A>A:T, C:G>G:C,
21 124 C:G>A:T) to global alteration patterns.

22
23
24
25
26
27
28
29 125 The advent of massively parallel sequencing has created the opportunity to study a
30 126 large number of mutations in a single sample, thus significantly enhancing the power of
31 127 mutation analysis in experimental models and enabling reliable identification of specific
32 128 sequence contexts for the induced alterations. Analogously to human cancer genome
33 129 projects, genome-scale mutational signatures can be extracted from highly controlled
34 130 carcinogen exposure experiments using mammalian cell and animal models coupled with
35 131 advanced mathematical approaches (2,3,37,38).

36
37
38
39
40 132 Here we report the systematic assessment of acrylamide and glycidamide
41 133 mutagenicity based on DNA adduct formation and mutation profile analysis using massively
42 134 parallel sequencing in a cell model amenable to the analysis of carcinogen-induced mutation
43 135 patterns and their impact on the resulting cell phenotype (3,37-39). We identify a specific
44 136 and robust mutational signature attributable to glycidamide, and by computationally
45 137 interrogating human cancer genome-wide mutation data, we characterize glycidamide
46 138 signature-positive tumors, thereby highlighting a potential contribution of
47 139 acrylamide/glycidamide exposure to carcinogenesis in humans.

48
49
50
51
52 140

141 **Materials and methods**

142 **Source and authentication of primary cells**

143 Primary Human-p53 knock-in mouse embryonic fibroblasts (Hupki MEFs) were isolated from
144 13.5-day old *Trp53^{tm/Holl}* mouse embryos from the Central Animal Laboratory of the
145 Deutsches Krebsforschungszentrum, Heidelberg, as described previously (40). The mice
146 had been tested for Specific Pathogen-Free (SPF) status. The derived primary cells were
147 genotyped for the human *TP53* codon 72 polymorphism (Table 1) to authenticate the
148 embryo of origin. Cells from three different embryos (E210, E213 and E214) were used for
149 the exposure experiments (Table 1). All subsequent cell cultures were routinely tested at all
150 stages for the absence of mycoplasma.

151

152 **Cell culture, exposure and immortalization**

153 The primary MEF cells were expanded in Advanced DMEM supplemented with 15% fetal
154 calf serum, 1% penicillin/streptomycin, 1% pyruvate, 1% glutamine, and 0.1% β -mercapto-
155 ethanol. The cells were then seeded in six-well plates and, at passage 2, exposed for 24
156 hours to acrylamide (A4058, Sigma), glycidamide (04704, Sigma), or vehicle (PBS).
157 Acrylamide exposure was carried out in the absence or presence of 2% human S9 fraction
158 (Life Technologies) complemented with NADPH (Sigma). Exposed and control primary cells
159 were cultivated until they bypassed senescence and immortalized clonal cell populations
160 could be isolated (41). The human mammary epithelial cell (HMEC) cultures utilized in this
161 study for whole-genome sequencing (WGS) were generated from benzo[a]pyrene (B[a]P)
162 exposed HMEC described previously (42,43).

163

164 **MTT assay for cell metabolic activity and viability**

165 Cells were seeded in 96-well plates and treated as indicated. Cell viability was measured 48
166 hours after treatment cessation using CellTiter 96® Aqueous One solution Cell Proliferation
167 Assay (Promega). Plates were incubated for 4 hours at 37°C and absorbance was
168 measured at 492 nm using the APOLLO 11 LB913 plate reader. The MTT assay was
169 performed in triplicates for each experimental condition.

170

171 **γ H2Ax Immunofluorescence**

172 Immunofluorescence staining was carried out using an antibody specific for Ser139-
173 phosphorylated H2Ax (γ H2Ax) (9718, Cell Signaling Technology). Primary MEFs were
174 seeded on coverslips in 12 well-plates. The cells were incubated in with γ H2Ax-antibody
175 (1:500 in 1% BSA) at 4°C overnight. Subsequent incubation with a fluorochrome-conjugated
176 secondary antibody (4412, Cell Signaling Technology) was carried out for 60 minutes at

177 room temperature. Coverslips were mounted in Vectashield mounting medium with DAPI
178 (Eurobio). Immunofluorescence images were captured using a Nikon Eclipse Ti.

179

180 **DNA adduct analysis**

181 Glycidamide-DNA adducts (N7-(2-carbamoy-2-hydroxyethyl)-guanine (N7-GA-Gua) and N3-
182 (2-carbamoy-2-hydroxyethyl)-adenine (N3-GA-Ade)) were quantified by liquid
183 chromatography-mass spectrometry (LC-MS/MS) with stable isotope dilution as previously
184 described (44) (see Supplementary Materials and Methods for details). The LC-MS/MS used
185 for quantification consisted of an Acquity UPLC system (Waters) and a Xevo TQ-S triple
186 quadrupole mass spectrometer (Waters). The same MRM transitions as previously
187 described (44) were monitored with a cone voltage of 50V and collision energy of 20eV for
188 each adduct transition and its corresponding labeled isotope transition.

189

190 **TP53 genotyping**

191 Exons 4 to 8 of the knocked-in human *TP53* gene (NC_000017.11) were sequenced using
192 standard protocols. Sanger sequencing of PCR products was performed at Biofidal (Lyon,
193 France). *TP53* primer sequences are listed in Supplementary Materials and Methods.
194 Resulting sequences were analyzed using the CodonCode Aligner software.

195

196 **Library preparation and whole-exome sequencing (WES)**

197 Library preparation was carried out using the Kapa Hyper Plus library preparation kit (Kapa
198 Biosystems) according the manufacturer's instructions. Exome capture was performed using
199 the SureSelect XT Mouse All Exon Kit (Agilent Technologies). Eighteen exome-captured
200 libraries were sequenced in the paired-end 150 base-pair run mode using the Illumina
201 HiSeq4000 sequencer.

202

203 **Processing of WES data**

204 Fastq files were analyzed for data amount and quality using FastQC (0.11.3) and were
205 processed with an in-house pipeline for adapter trimming and alignment to the mm10
206 genome (release GRCm38). These components of the pipeline are publicly available at
207 <https://github.com/IARCBioinfo/alignment-nf>. The resulting alignment files had a mean depth-
208 of-coverage of 135 and 175 for acrylamide and glycidamide samples, respectively. All
209 alignment files can be accessed from the NCBI Sequence Read Archive (SRA) data portal
210 under the BioProject accession number PRJNA238303. Two somatic variant callers were
211 employed with default parameters in order to detect single base substitutions (SBS) and
212 small insertions/deletions (indels) (MuTect 1.1.6-4 and Strelka 1.015) in exposed clones,
213 using primary cells as normal samples. Each immortalized clone was compared to primary

1
2
3 214 MEFs from three different embryos (conditions Prim_1, Prim_2, and Prim_3). The overlap of
4 215 the variant calling outcome with respect to the different primary MEFs showed concordance
5 216 close to 80% (Suppl. Fig. S1) with MuTect exhibiting more stringent calling performance.
6 217 Thus, mutation data obtained from the MuTect variant caller were further processed with the
7 218 MutSpec suite ((45); <https://github.com/IARCbioinfo/mutspec>). For more details, see
8 219 Supplementary Materials and Methods and the summary of sequencing metrics (Suppl.
9 220 Table S1), the list of identified MuTect SBS variants (Suppl. Table S2) and indels (Suppl.
10 221 Table S3).

14 222

16 223 **Bioinformatics and statistical analyses**

17 224 The FactoMiner R package (R package version 3.3.2; [https://cran.r-](https://cran.r-project.org/web/packages/FactoMineR)
18 225 [project.org/web/packages/FactoMineR](https://cran.r-project.org/web/packages/FactoMineR)) was used to perform the principal component
19 226 analysis (PCA). To perform the transcription strand bias (SB) analyses, p -values were
20 227 calculated using Pearson's χ^2 test. As multiple comparisons were assessed, the p -value
21 228 was adjusted by applying a false discovery rate (FDR). Statistical analyses were carried out
22 229 using the stats R package. The SB was considered statistically significant at p -value ≤ 0.05 .
23 230 To analyze samples mutation spectra and treatment-specific mutational signatures, filtered
24 231 mutations were classified into 96 types corresponding to the six possible base substitutions
25 232 (C:G>A:T, C:G>G:C, C:G>T:A, T:A>A:T, T:A>C:G, T:A>G:C) and the 16 combinations of
26 233 flanking nucleotides immediately 5' and 3' of the mutated base. Mutation patterns were then
27 234 deconvoluted into mutational signatures using the non-negative matrix factorization (NMF)
28 235 algorithm (46,47). The reconstruction error calculation evaluated the accuracy with which the
29 236 deciphered mutational signatures describe the original mutation spectra of each sample by
30 237 applying Pearson correlation and cosine similarity.

31 238 In order to clean up the profile of the glycidamide mutational signature from the
32 239 residual signature 17 signal and to increase the stability of NMF decomposition, we supplied
33 240 the NMF input by adding samples with a high level of signature 17 (over 65% contribution as
34 241 determined by independent NMF analysis, see Supplementary Materials and Methods).

35 242 Cosine similarity analysis was used to evaluate the concordance of the newly
36 243 identified T:A>A:T-rich mutational signature of glycidamide with the previously reported
37 244 mutational signatures characterized by a predominant T:A>A:T content. These comprised
38 245 COSMIC signatures 22 (AA), 25 and 27 (both of unknown etiology(2)), the experimentally
39 246 derived mutational signature of AA (37,45), 7,12-dimethylbenz[*a*]anthracene (DMBA)
40 247 (48,49), and urethane (50).

41 248 We employed the mutational signature activity (mSigAct) software's sparse signature
42 249 assignment function (`sparse.assign.activity`) (13) to assess the presence of the experimental
43 250 mutational signatures of glycidamide and benzo[*a*]pyrene in whole-genome somatic mutation

1
2
3 251 data from 38 lung adenocarcinomas, 48 lung squamous carcinomas, and 320 liver cancers
4 252 from the ICGC Pan-Cancer Analysis of Whole Genomes (PCAWG) study. We excluded 244
5 253 hyper-mutated microsatellite unstable and aristolochic acid signature-containing liver tumors
6 254 as the presence of high numbers of T>A mutations adversely prevented assessment of the
7 255 possible presence of the glycidamide signature. A set of 11 active COSMIC mutational
8 256 signatures were identified in the remaining tumor samples (excluding COSMIC signature 4).

9
10
11 257 We defined a 'pure' experimental C>N benzo[a]pyrene signature by WGS (using
12 258 Illumina HiSeq4000 by Genewiz, NJ, USA) of finite lifespan post-stasis clones derived from
13 259 primary human mammary epithelial cells (HMEC) treated with B[a]P as previously described
14 260 (42,43,51). The read alignment to NCBI GRCh38 genome build, variant calling, filtering and
15 261 annotation were consistent with the MutSpec pipeline described above (45). Proportion
16 262 matrices of the experimental GA-signature, the GA-signature normalized to the human
17 263 genome trinucleotide frequency to allow for human PCAWG data screening, and the whole-
18 264 genome B[a]P signature are available in Suppl. Table S4.

25 265 **Results**

26 266 **Acrylamide and glycidamide induce cytotoxic and genotoxic responses in Hupki** 27 267 **MEFs**

28 268 Upon exposure of primary Hupki MEFs to a range of concentrations of acrylamide (ACR) (in
29 269 the absence or presence of the S9 fraction) and its metabolite, glycidamide (GA), we
30 270 observed a dose-dependent cytotoxic effect on the cells for either compound (Fig. 1A). This
31 271 analysis informed the selection of two conditions for the ACR exposure to be used in the
32 272 subsequent exposure/immortalization experiments, 10 mM ACR for 24 hours in the absence
33 273 of human S9 fraction, and 5 mM ACR for 24 hours in the presence of S9 fraction, which
34 274 elicited 50% (range 30-70%) decrease in cell viability. The IC50 condition for GA was used
35 275 for subsequent mutagenesis analysis, corresponding to a 24-hour treatment with 3 mM of
36 276 the compound. The genotoxic effects of either ACR or GA manifested by a marked increase
37 277 in γ H2Ax staining in the exposed cell populations, in comparison to the mock-treated control
38 278 cells (Fig. 1B).

39 279

40 280 **Immortalized MEF cells accumulate TP53 mutations following acrylamide or** 41 281 **glycidamide treatment**

42 282 Primary MEF cultures from three different embryos (Prim_1, Prim_2, and Prim_3) were
43 283 exposed to ACR or GA using the established conditions and multiple immortalized clones
44 284 were derived. MEF senescence and immortalization phases were evident from the growth
45 285 curves generated for each culture (Suppl. Fig. S2). Subsequently, the clones derived from
46 286 ACR exposure (ACR clones) and GA exposure (GA clones) and spontaneous

1
2
3 287 immortalization (Spont), were pre-screened for *TP53* mutations by Sanger sequencing, to
4 288 assess the mutagenic process prior to exome-scale analysis. In the context of ACR
5 289 treatment, clones obtained from the Prim_2 MEFs that were heterozygous for the
6 290 polymorphic site in codon 72 showed a loss of heterozygosity involving a loss of the proline
7 291 allele in the ACR_1 clone whereas the arginine allele was lost in ACR_2, giving rise to a
8 292 hemizygous clone (Table 1). No *TP53* mutations were observed in any of the three Spont
9 293 clones, whereas 3 out of 7 ACR clones and 1 of 5 GA clones carried non-synonymous *TP53*
10 294 mutations (Table 1). The detected mutations indicated specific selection for mutations in the
11 295 *TP53* gene during cell immortalization and confirmed the clonal nature of MEF
12 296 immortalization.
13
14
15
16
17
18

297

298 **Analysis of mutation spectra**

19
20 299 Whole-exome sequencing of all spontaneously immortalized and exposed clones and
21 300 subsequent extraction of acquired variants revealed that the total number of acquired SBS
22 301 did not differ markedly between the ACR and Spont clones. The Spont clones harbored on
23 302 average 190 (median = 151, range = 141-277) SBS, whereas the ACR clones had on
24 303 average 208 (median = 173, range = 151-262) SBS. In contrast, the total number of SBS
25 304 was considerably increased in the GA clones, with an average of 485 SBS (median = 448,
26 305 range = 370-592) (Suppl. Table S1 and S2). This finding suggests markedly stronger
27 306 mutagenic properties of GA in the MEFs. To estimate the extent of sequencing-related
28 307 damage in our samples, we determined the GIV score of each sample as described in
29 308 Materials and Methods and in (52). No detectable damage for any of the mutation types was
30 309 observed in our dataset (data not shown). The ACR exposed samples exhibited an overall
31 310 diffuse pattern across the six different SBS types (Suppl. Fig. S3). The Spont clones showed
32 311 an enrichment of C:G>G:C SBS in the 5'-GCC-3' context, which was also present at varying
33 312 levels in the exposed cultures. This particular mutation type appears to be related to the
34 313 culture conditions used for the immortalization assay, as its presence has previously been
35 314 noted upon spontaneous as well as exposure-driven MEF immortalization (37). No
36 315 significant transcription strand bias was observed for any of the mutation classes in the
37 316 Spont or ACR clones (Suppl. Fig. S4). In the five clones derived from the GA-treated primary
38 317 MEF cultures, we observed an enrichment of acquired T:A>A:T and C:G>A:T transversions
39 318 and T:A>C:G transitions (Suppl. Fig. S3B), marked by significant transcription strand bias
40 319 (Suppl. Fig. S4).

41
42
43
44
45
46
47
48
49
50
51
52 320 PCA performed on the resulting 6-class SBS spectra unambiguously separated the
53 321 GA clones from the remaining experimental conditions (Fig. 2A). The analysis of indels
54 322 (listed in Suppl. Table S3) showed lower numbers of these alterations in the GA-associated
55 323 clones compared to the ACR or Spont clones (Fig. 2B). This suggests that a higher
56
57
58
59
60

1
2
3 324 accumulation of SBS may selectively promote the senescence bypass and selection of the
4 325 GA clones, with a decreased functional contribution of indels, while an inverse scenario is
5
6 326 plausible in case of the Spont and ACR clones, reminiscent of a previous report based on
7
8 327 the Big Blue mouse embryonic fibroblasts and *c/* transgene (53).
9

328

329 **Variant allele frequency analysis**

330 Variant allele frequency (VAF) analysis was carried out for GA clones. Overall, a significant
331 proportion of acquired mutations was present at allelic frequencies between 25-75% (Suppl.
332 Fig. S5). Upon grouping of substitutions into bins of high (67-100%), medium (34-66%) and
333 low (0-33%) VAF, the predominant GA-specific mutation types (T:A>A:T, T:A>C:G and
334 C:G>A:T) started manifesting at high VAF, whereas the 5'-NIT-3' alterations, corresponding
335 to the COSMIC signature 17 previously reported to arise in cultured mouse cells including
336 MEFs (38,54,55) showed lower VAF, therefore a later appearance in the cultures (Suppl.
337 Fig. S6). This observation suggests the early effects of the GA exposure and the
338 reproducible contribution of the induced mutations to the senescence bypass and their clonal
339 propagation during the immortalization stage.

340

341 **Mutational signature analysis**

342 Using NMF, we extracted the mutational signatures from all the MEF clones. Using
343 computed statistics for estimating the number of signatures, three signatures were identified
344 as an optimal number, with signatures A and C enriched in the Spont and ACR clones, and
345 signature B selectively enriched in the GA clones (Fig. 2C,D). Reconstruction of the
346 observed mutation spectra supports the robustness of the signature analysis with strong
347 Pearson's correlation and cosine similarity in GA-derived clones (Fig. 2D). In signature C
348 and also to a lesser extent in signatures A and B, we observed an admixture of a pattern
349 identical to the orphan COSMIC signature 17 (T:A>G:C in a 5'-NIT-3' trinucleotide context),
350 described in various human cancers (most notably esophageal adenocarcinoma), but also
351 seen in aflatoxin B1-driven mouse liver cancers (11), as well as primary MEF-derived clones
352 (37,38). In *in vitro* contexts, this signature has been linked to cell culture conditions and
353 associated oxidative stress (54,55). To refine further the obtained experimental signatures,
354 we developed a signature 'baiting' approach that combined the MEF clones data with
355 signature 17-rich data from esophageal adenocarcinomas from the ICGC ESAD-UK study
356 for new NMF analysis (56). This resulted in considerable reduction (average = 47%, median
357 = 48%) of the signature 17-specific most prominent T>G peaks and a more refined pattern
358 for signature B, associated primarily with GA treatment (Fig. 3A and Suppl. Fig. S7). This
359 putative GA signature retains the predominant enrichment for the T:A>A:T transversions and
360 T:A>C:G transitions in the 5'-CTG-3' and 5'-CTT-3' trinucleotide contexts, and the C:G>A:T

1
2
3 361 component. Moreover, these mutation types were marked by significant transcription strand
4 362 bias (Fig. 3B and Suppl. Fig. S4), exhibiting higher accumulation of mutations on the non-
5 363 transcribed strand consistent with the decreased efficiency of the transcription-coupled
6 364 nucleotide excision repair due to adduct formation.
7
8
9

365

366 **DNA adduct analysis**

10 367 Following metabolic activation, acrylamide induces well-characterized glycidamide DNA
11 368 adducts at the N7- and N3-positions of guanine and adenine, respectively. LC-MS/MS-based
12 369 adduct quantification revealed the absence of these adducts in the spontaneously
13 370 immortalized control samples as well as in MEFs exposed to acrylamide in the absence of
14 371 S9 fraction (levels below the limit of detection). This suggests the lack of CYP2E1 activity,
15 372 which is required for the metabolism of acrylamide to glycidamide, in the MEFs. Upon
16 373 addition of human S9 fraction, N7-GA-Gua levels increased to 11 adducts/10⁸ nucleotides,
17 374 suggesting limited metabolic activation of acrylamide due to the presence of enzymatic
18 375 activity in the S9 fraction (Fig. 3C and Suppl. Fig. S8). Glycidamide-exposed cells exhibited
19 376 significantly increased DNA adduct levels, with both N7-GA-Gua and N3-GA-Ade observed
20 377 at very high average levels, 49 000 adducts/10⁸ nucleotides and 350 adducts/10⁸
21 378 nucleotides, respectively, after subtracting the trace amount of contamination from the
22 379 internal standard (Fig. 3C and Suppl. Fig. S8).
23
24
25
26
27
28
29
30

380

381 **Comparison of the glycidamide signature to known signatures characterized by** 382 **prominent T:A>A:T profiles**

35 383 We next performed cosine similarity analysis of the putative GA signature and all known
36 384 T:A>A:T-rich signatures extracted from primary cancers as well as experimental systems
37 385 (Fig. 3D and Suppl. Fig. S9). The best match was 84% pattern similarity with COSMIC
38 386 signature 25 (derived from four Hodgkin lymphoma cell lines) (Fig. 3D). However, unlike the
39 387 GA signature, COSMIC signature 25 exhibits strand bias for only T:A>A:T mutations and no
40 388 transcription strand bias for the T:A>C:G mutations. Thus, the mutation patterns and strand
41 389 bias on all three main mutation types generated by GA treatment (Fig. 3A,B) appear specific
42 390 and novel.
43
44
45
46
47

391

392 **Glycidamide signature screening in human tumor data from the ICGC PCAWG**

48 393 The initial mSigAct test performed on PCAWG data from lung and liver tumors indicated a
49 394 marked presence of the GA signature. This observation was in keeping with the presence of
50 395 acrylamide in tobacco smoke and was further corroborated by a cosine similarity of 94%
51 396 between the adenine (T>N) components of COSMIC signature 4 (tobacco smoking) and the
52 397 GA signature (Fig. 4A). We thus hypothesized that COSMIC signature 4 reflects co-
53
54
55
56
57
58
59
60

1
2
3 398 exposure to B[a]P (generating C>N/guanine mutations with transcription strand bias) and to
4 399 GA (generating T>N/adenine mutations with transcription strand bias) (Fig. 4A,B). To
5 400 provide further experimental evidence, we generated a 'pure' B[a]P mutational signature by
6 401 whole-genome sequencing of cell clones derived from B[a]P-exposed normal human
7 402 mammary epithelial cells (HMEC). This yielded a robust signature characterized by
8 403 predominant strand biased guanine (mainly C>A) mutation levels and negligibly mutated
9 404 adenines (T>N) (Fig. 4A,B). Next, we used mSigAct to interrogate the PCAWG tumor
10 405 samples for the level of exposure to the experimentally defined GA and B[a]P signatures
11 406 (alongside other COSMIC mutational signatures) in 48 lung squamous carcinomas, 38 lung
12 407 adenocarcinomas, and 320 liver cancers. We compared these to estimated levels of
13 408 exposure to COSMIC signature 4, and found that in the lung cancers, a combination of the
14 409 GA and B[a]P signatures accounted for very similar numbers of mutations as COSMIC
15 410 signature 4, thus further supporting the hypothesis that COSMIC signature 4 represents
16 411 combined and highly correlated exposure to GA and B[a]P (Fig. 4C). Compared to lung
17 412 cancers, we found more variability in the assignment of mutation numbers to GA and B[a]P
18 413 versus COSMIC signature 4 in liver cancers (Fig. 4C), which may reflect a decreased
19 414 relationship between GA and B[a]P exposure due to generally more complex exposure
20 415 history in the liver. The successful reconstruction of COSMIC signature 4 by the
21 416 experimental GA- and B[a]P- signatures in the lung and liver human tumors enabled correct
22 417 assignment of the GA-signature in a subset of 29 lung adenocarcinomas, 46 lung SCC and
23 418 26 liver tumors (Fig. 4D). The SBS counts corresponding to GA-mutational signature ranged
24 419 between 300 up to 43,000 mutations/per sample in lung tumors, and between 190 to 23,000
25 420 mutations/per sample in liver tumors (Fig. 4D and Suppl. Table S5). These findings indicate
26 421 exposure to glycidamide linked to tobacco smoking – when concomitant with B[a]P-
27 422 signature, or through diet or occupation – in the absence of B[a]P signature (samples Liver-
28 423 HCC::SP112224; Liver-HCC::SP49551; Liver-HCC::SP50105; Liver-HCC::SP98861; Liver-
29 424 HCC::SP50183, see Suppl. Fig. S10 and Suppl. Table S5).

425 Discussion

426 In this study we report the identification of an exome-wide mutational signature for
427 glycidamide, a metabolite of the probable human carcinogen acrylamide. The newly
428 identified signature is based on massively parallel sequencing performed in a well-controlled
429 experimental carcinogen exposure-clonal immortalization model, revealing characteristic
430 mutagenic effects of glycidamide. The glycidamide mutational signature presented here and
431 the results of statistical assessment of its presence in multiple human tumor types may help
432 clarify the thus-far tenuous association of acrylamide with human cancer.

1
2
3 433 In concordance with its *in vivo* carcinogenicity in rodents (16,19,31,57), our findings
4 434 in the established MEF carcinogen exposure and immortalization system suggest that
5
6 435 characteristic mutagenic effects may play a role during acrylamide/glycidamide-driven tumor
7
8 436 development. In contrast to glycidamide, acrylamide exposure led neither to an increased
9
10 437 number of SBS nor did it induce characteristic mutation types in the MEF exposure system.
11
12 438 Despite the absence of a mutagenic effect of acrylamide in our experiments, acrylamide and
13
14 439 glycidamide exposures induce an almost identical set of tumors in both mice and rats,
15
16 440 providing a substantial argument for a glycidamide-mediated tumorigenic effect of
17
18 441 acrylamide (19). This is further supported by mechanistic studies showing that lung tissue
19
20 442 from mice exposed to acrylamide and glycidamide displays comparable DNA adduct
21
22 443 patterns as well as similar mutation frequencies in the *cII* transgene (36). Similar
23
24 444 observations had been made in the context of *in vitro* mutagenicity of acrylamide in human
25
26 445 and mouse cells, suggesting the key role for epoxide metabolite glycidamide to form pre-
27
28 446 mutagenic DNA adducts (33).

29
30 447 As shown by our adduct analysis, acrylamide is not efficiently metabolized by MEFs.
31
32 448 This finding is in keeping with the results from previous animal carcinogenicity studies. In
33
34 449 fact, glycidamide induces hepatocellular carcinomas in neonatal B6C3F1 mice, whereas
35
36 450 administration of acrylamide does not increase the tumor incidence. This has been attributed
37
38 451 to the inability of neonatal mice to efficiently metabolize acrylamide (31). Moreover, in
39
40 452 contrast to acrylamide treatment, glycidamide induces tumors of the small intestine in a
41
42 453 dose-dependent manner upon perinatal exposure (57) and similar observations were made
43
44 454 for glycidamide mutagenicity *in vitro* (33). We compensated for the lack of proper acrylamide
45
46 455 metabolic activation by the addition of human S9 fraction, and the assessment of DNA
47
48 456 adducts indeed suggests acrylamide metabolic activation upon addition of S9. However, the
49
50 457 adduct levels are substantially lower compared to glycidamide exposure, which may account
51
52 458 for the observed differences in mutagenicity. Interestingly, a consistent minor contribution of
53
54 459 the glycidamide mutational signature was detected in the majority of ACR clones, whereas it
55
56 460 was absent in the Spont clones. This raises the possibility that partial metabolic activation of
57
58 461 acrylamide in the MEF system resulted in low levels of glycidamide. However, a clear
59
60 462 mutational signature in the employed experimental setting was achieved only by exposing
463 the cells directly to glycidamide.

464 Single reporter gene studies had previously linked acrylamide and glycidamide
465 exposure to multiple different mutation types. Thanks to the larger number of mutations
466 captured by exome sequencing, we were able to attribute to the glycidamide exposure a
467 particular mutational signature characterized by strand-biased C:G>A:T and T:A>A:T
468 transversions, and T:A>C:A transitions towards the non-transcribed strand suggesting a
469 formation of DNA-adducts. The presence of N7-GA-Gua and N3-GA-Ade, two well-

1
2
3 470 characterized glycidamide DNA adducts originating from the metabolic conversion of
4 471 acrylamide (30,44,53), shows a remarkable relationship between DNA adduct profiles and
5 472 the putative mutational signature of glycidamide. N3-GA-Ade and N7-GA-Gua are
6 473 depurinating adducts. They can result in apurinic/apyrimidinic sites, which, during replication,
7 474 induce the mis-incorporation of deoxyadenine, leading to the observed T:A>A:T and
8 475 C:G>A:T transversions of the glycidamide signature, respectively. The third mutation type
9 476 specifically enriched in the glycidamide signature, T:A>C:G transitions, has been ascribed to
10 477 the N1-GA-Ade adduct, a miscoding adduct and the most commonly identified adenine
11 478 adduct *in vitro* (35,44,53,58). Levels of the guanine adduct were especially high in the
12 479 exposed MEF cells, whereas the associated C:G>A:T transversions in the resulting post-
13 480 senescence clones were less represented. This could reflect differences in DNA repair
14 481 efficiency concerning individual GA-DNA adduct species, or the fact that the resulting clones
15 482 are derived from single cells whereas the GA-DNA adducts were measured on average in
16 483 the bulk primary cell population. A mechanism of negative selection of cells with high N7-
17 484 GA-Gua adduct burden is also plausible.

18
19 485 We observed consistent presence of COSMIC signature 17 in the data generated
20 486 from the untreated and treated MEF clones. The etiology of signature 17 remains unknown.
21 487 While some candidate causal factors have been proposed in esophageal adenocarcinoma
22 488 and gastric cancers (e.g., inflammatory conditions due to acid reflux, *H. pylori*) (56) and in
23 489 cultured mouse cell systems (54,55), further studies are required to establish why signature
24 490 17 tends to arise *in vitro* in immortalized clones derived from mouse embryonic fibroblasts as
25 491 observed in our study and also previous work (38).

26 492 Genome-scale sequencing of tumor tissues will be needed to verify, *in vivo*, the
27 493 glycidamide mutational signature identified in this study. The established animal models
28 494 (18,19) of acrylamide- and glycidamide-mediated tumorigenesis provide a suitable starting
29 495 point, and it would be interesting to compare mutational signatures derived from these
30 496 models with the *in vitro* results. The identified glycidamide signature with its extended
31 497 features of transcription strand bias for the major mutation types differs from the currently
32 498 known COSMIC signatures (Fig. 3D). In addition, we show that in the cancer genome
33 499 sequencing data sets from the ICGC PCAWG effort, the putative glycidamide-mutational
34 500 signature can be identified in a subset of tumors of the lung and liver (sites of possible
35 501 acrylamide exposure due to tobacco smoking), based on combining experimentally derived
36 502 signatures with sophisticated computational signature reconstruction approaches (Fig. 4).

37 503 The continued interest in understanding the contribution of acrylamide and its
38 504 electrophilic metabolite glycidamide to cancer development reflects recent accumulation of
39 505 new mechanistic data on the animal carcinogenicity of the compounds. The possible
40 506 carcinogenic effects in humans have been recommended for re-evaluation by the Advisory

1
2
3 507 Group to the Monographs Program of the International Agency for Research on Cancer (59).
4 508 Our findings related to the reconstruction of COSMIC signature 4 using the experimental
5 509 GA-signature and B[a]P signature, together with the presence of the GA signature in the
6 510 lung and liver cancer data are relevant given the established high contents of acrylamide in
7 511 tobacco smoke. Despite the absence of prominent T>N (adenine) mutations in the
8 512 experimental B[a]P exposure setting, we cannot exclude a possibility that in the human lung
9 513 cells the adenine residues can be additionally targeted by other tobacco carcinogens such
10 514 as benzo[a]pyrene derivatives or nitrosamines. Importantly, five liver tumor samples
11 515 identified in this study harbored the GA signature but the major features of signature 4 as
12 516 represented by the experimental B[a]P signature were absent (Suppl. Fig. S10, Suppl. Table
13 517 S5). These tumors are thus of particular interest as they could reflect dietary or occupational
14 518 exposure to acrylamide.

15
16 519 The presented mutational signature of glycidamide and its potential use for screening
17 520 of cancer genome sequencing data may provide a basis for relevant assessment of cancer
18 521 risk through new carefully designed molecular cancer epidemiology studies. Future
19 522 validation analyses involving e.g. GA-DNA adduct monitoring in non-tumor tissue of cancer
20 523 patients or in animal exposure models are warranted to provide additional evidence that the
21 524 predominant T>N mutations in the cancers identified in this study indeed originate from
22 525 exposure to acrylamide and its reactive metabolite glycidamide.

526 **Acknowledgments**

527 The views expressed in this manuscript do not necessarily represent those of the U.S. Food
528 and Drug Administration. The study was supported by funding obtained from INCa-INSERM
529 (Plan Cancer 2015 grant to J.Z.), NIH/NIEHS (1R03ES025023-01A1 grant to M.O.), and the
530 Singapore National Medical Research Council (NMRC/CIRG/1422/2015 grant to S.G.R.) and
531 the Singapore Ministry of Health via the Duke-NUS Signature Research Programmes to
532 S.G.R.. M.R.S. was supported by the U.S. Department of Energy under Contract No. DE-
533 AC02-05CH11231. We thank the NYU Genome Technology Center, funded in part by the
534 NIH/NCI Cancer Center Support Grant P30CA016087, and GENEWIZ, South Plainfield, NJ,
535 USA, for expert assistance with Illumina sequencing.

536

537 **References**

- 538 1. Smith, M.T., *et al.* (2016) Key Characteristics of Carcinogens as a Basis for
539 Organizing Data on Mechanisms of Carcinogenesis. *Environ Health Perspect*, **124**,
540 713-21.
- 541 2. Alexandrov, L.B., *et al.* (2013) Signatures of mutational processes in human cancer.
542 *Nature*, **500**, 415-421.

- 1
- 2
- 3 543 3. Zhivagui, M., *et al.* (2017) Modelling Mutation Spectra of Human Carcinogens Using
4 544 Experimental Systems. *Basic Clin Pharmacol Toxicol*, **121 Suppl 3**, 16-22.
- 5 545 4. Hollstein, M., *et al.* (2017) Base changes in tumour DNA have the power to reveal the
6 546 causes and evolution of cancer. *Oncogene*, **36**, 158-167.
- 7 547 5. Poon, S.L., *et al.* (2013) Genome-Wide Mutational Signatures of Aristolochic Acid
8 548 and Its Application as a Screening Tool. *Science Translational Medicine*, **5**,
9 549 197ra101-197ra101.
- 10 550 6. Meier, B., *et al.* (2014) *C. elegans* whole-genome sequencing reveals mutational
11 551 signatures related to carcinogens and DNA repair deficiency. *Genome Res*, **24**,
12 552 1624-36.
- 13 553 7. Scelo, G., *et al.* (2014) Variation in genomic landscape of clear cell renal cell
14 554 carcinoma across Europe. *Nat Commun*, **5**, 5135.
- 15 555 8. Jelakovic, B., *et al.* (2015) Renal cell carcinomas of chronic kidney disease patients
16 556 harbor the mutational signature of carcinogenic aristolochic acid. *Int J Cancer*, **136**,
17 557 2967-72.
- 18 558 9. Hoang, M.L., *et al.* (2016) Aristolochic Acid in the Etiology of Renal Cell Carcinoma.
19 559 *Cancer Epidemiology, Biomarkers & Prevention*, **25**, 1600-1608.
- 20 560 10. Chawanthayatham, S., *et al.* (2017) Mutational spectra of aflatoxin B1 in vivo
21 561 establish biomarkers of exposure for human hepatocellular carcinoma. *Proc Natl*
22 562 *Acad Sci U S A*, **114**, E3101-E3109.
- 23 563 11. Huang, M.N., *et al.* (2017) Genome-scale mutational signatures of aflatoxin in cells,
24 564 mice, and human tumors. *Genome Res*, **27**, 1475-1486.
- 25 565 12. Zhang, W., *et al.* (2017) Genetic Features of Aflatoxin-Associated Hepatocellular
26 566 Carcinoma. *Gastroenterology*, **153**, 249-262 e2.
- 27 567 13. Ng, A.W.T., *et al.* (2017) Aristolochic acids and their derivatives are widely implicated
28 568 in liver cancers in Taiwan and throughout Asia. *Sci Transl Med*, **9**.
- 29 569 14. Mojska, H., *et al.* (2016) Acrylamide content in cigarette mainstream smoke and
30 570 estimation of exposure to acrylamide from tobacco smoke in Poland. *Annals of*
31 571 *agricultural and environmental medicine: AAEM*, **23**, 456-461.
- 32 572 15. Takatsuki, S., *et al.* (2003) Determination of acrylamide in processed foods by LC/MS
33 573 using column switching. *Shokuhin Eiseigaku Zasshi. Journal of the Food Hygienic*
34 574 *Society of Japan*, **44**, 89-95.
- 35 575 16. IARC Monograph vol. 60 (1994) *Some industrial chemicals. Lyon, 15 - 22 February*
36 576 *1994, Lyon*.
- 37 577 17. Tareke, E., *et al.* (2002) Analysis of Acrylamide, a Carcinogen Formed in Heated
38 578 Foodstuffs. *Journal of Agricultural and Food Chemistry*, **50**, 4998-5006.
- 39 579 18. Beland, F.A., *et al.* (2013) Carcinogenicity of acrylamide in B6C3F(1) mice and
40 580 F344/N rats from a 2-year drinking water exposure. *Food and Chemical Toxicology*,
41 581 **51**, 149-159.
- 42 582 19. Beland, F.A., *et al.* (2015) Carcinogenicity of glycidamide in B6C3F1 mice and
43 583 F344/N rats from a two-year drinking water exposure. *Food and Chemical*
44 584 *Toxicology*, **86**, 104-115.
- 45 585 20. Hogervorst, J.G., *et al.* (2008) Dietary acrylamide intake and the risk of renal cell,
46 586 bladder, and prostate cancer. *The American Journal of Clinical Nutrition*, **87**, 1428-
47 587 1438.
- 48 588 21. Virk-Baker, M.K., *et al.* (2014) Dietary Acrylamide and Human Cancer: A Systematic
49 589 Review of Literature. *Nutrition and Cancer*, **66**, 774-790.
- 50 590 22. Olesen, P.T., *et al.* (2008) Acrylamide exposure and incidence of breast cancer
51 591 among postmenopausal women in the Danish Diet, Cancer and Health Study.
52 592 *International Journal of Cancer*, **122**, 2094-2100.
- 53 593 23. Wilson, K.M., *et al.* (2009) Acrylamide exposure measured by food frequency
54 594 questionnaire and hemoglobin adduct levels and prostate cancer risk in the Cancer
55 595 of the Prostate in Sweden Study. *International Journal of Cancer*, **124**, 2384-2390.
- 56
- 57
- 58
- 59
- 60

- 1
2
3 596 24. Xie, J., *et al.* (2013) Acrylamide Hemoglobin Adduct Levels and Ovarian Cancer
4 597 Risk: A Nested Case-Control Study. *Cancer Epidemiology Biomarkers & Prevention*,
5 598 **22**, 653-660.
- 6 599 25. Obón-Santacana, M., *et al.* (2016) Acrylamide and glycidamide hemoglobin adduct
7 600 levels and endometrial cancer risk: A nested case-control study in nonsmoking
8 601 postmenopausal women from the EPIC cohort. *International Journal of Cancer*, **138**,
9 602 1129-1138.
- 10 603 26. Obón-Santacana, M., *et al.* (2016) Acrylamide and Glycidamide Hemoglobin Adducts
11 604 and Epithelial Ovarian Cancer: A Nested Case-Control Study in Nonsmoking
12 605 Postmenopausal Women from the EPIC Cohort. *Cancer Epidemiology, Biomarkers &*
13 606 *Prevention*, **25**, 127-134.
- 14 607 27. Obón-Santacana, M., *et al.* (2016) Dietary and lifestyle determinants of acrylamide
15 608 and glycidamide hemoglobin adducts in non-smoking postmenopausal women from
16 609 the EPIC cohort. *European Journal of Nutrition*.
- 17 610 28. Sumner, S.C., *et al.* (1999) Role of cytochrome P450 2E1 in the metabolism of
18 611 acrylamide and acrylonitrile in mice. *Chemical Research in Toxicology*, **12**, 1110-
19 612 1116.
- 20 613 29. Ghanayem, B.I., *et al.* (2005) Role of CYP2E1 in the epoxidation of acrylamide to
21 614 glycidamide and formation of DNA and hemoglobin adducts. *Toxicological Sciences*,
22 615 **88**, 311-318.
- 23 616 30. Segerbäck, D., *et al.* (1995) Formation of N-7-(2-carbamoyl-2-hydroxyethyl) guanine
24 617 in DNA of the mouse and the rat following intraperitoneal administration of [14C]
25 618 acrylamide. *Carcinogenesis*, **16**, 1161-1165.
- 26 619 31. Von Tungeln, L.S., *et al.* (2012) Tumorigenicity of acrylamide and its metabolite
27 620 glycidamide in the neonatal mouse bioassay. *International Journal of Cancer*, **131**,
28 621 2008-2015.
- 29 622 32. Besaratinia, A., *et al.* (2003) Weak yet distinct mutagenicity of acrylamide in
30 623 mammalian cells. *Journal of the National Cancer Institute*, **95**, 889-896.
- 31 624 33. Besaratinia, A., *et al.* (2004) Genotoxicity of acrylamide and glycidamide. *Journal of*
32 625 *the National Cancer Institute*, **96**, 1023-1029.
- 33 626 34. Von Tungeln, L.S., *et al.* (2009) DNA adduct formation and induction of micronuclei
34 627 and mutations in B6C3F1/Tk mice treated neonatally with acrylamide or glycidamide.
35 628 *International Journal of Cancer*, **124**, 2006-2015.
- 36 629 35. Ishii, Y., *et al.* (2015) Acrylamide induces specific DNA adduct formation and gene
37 630 mutations in a carcinogenic target site, the mouse lung. *Mutagenesis*, **30**, 227-235.
- 38 631 36. Manjanatha, M.G., *et al.* (2015) Acrylamide-induced carcinogenicity in mouse lung
39 632 involves mutagenicity: cll gene mutations in the lung of big blue mice exposed to
40 633 acrylamide and glycidamide for up to 4 weeks. *Environ Mol Mutagen*, **56**, 446-56.
- 41 634 37. Olivier, M., *et al.* (2014) Modelling mutational landscapes of human cancers in vitro.
42 635 *Scientific Reports*, **4**.
- 43 636 38. Nik-Zainal, S., *et al.* (2015) The genome as a record of environmental exposure.
44 637 *Mutagenesis*, **30**, 763-70.
- 45 638 39. Huskova, H., *et al.* (2017) Modeling cancer driver events in vitro using barrier
46 639 bypass-clonal expansion assays and massively parallel sequencing. *Oncogene*, **36**,
47 640 6041-6048.
- 48 641 40. Liu, Z., *et al.* (2004) Human tumor p53 mutations are selected for in mouse
49 642 embryonic fibroblasts harboring a humanized p53 gene. *Proceedings of the National*
50 643 *Academy of Sciences of the United States of America*, **101**, 2963-2968.
- 51 644 41. Todaro, G.J., *et al.* (1963) Quantitative studies of the growth of mouse embryo cells
52 645 in culture and their development into established lines. *The Journal of Cell Biology*,
53 646 **17**, 299-313.
- 54 647 42. Severson, P.L., *et al.* (2014) Exome-wide mutation profile in benzo[a]pyrene-derived
55 648 post-stasis and immortal human mammary epithelial cells. *Mutation*
56 649 *Research/Genetic Toxicology and Environmental Mutagenesis*, **775-776**, 48-54.

- 1
2
3 650 43. Stampfer, M.R., *et al.* (1985) Induction of transformation and continuous cell lines
4 651 from normal human mammary epithelial cells after exposure to benzo[a]pyrene. *Proc*
5 652 *Natl Acad Sci U S A*, **82**, 2394-8.
6 653 44. Gamboa da Costa, G., *et al.* (2003) DNA adduct formation from acrylamide via
7 654 conversion to glycidamide in adult and neonatal mice. *Chemical Research in*
8 655 *Toxicology*, **16**, 1328-1337.
9 656 45. Ardin, M., *et al.* (2016) MutSpec: a Galaxy toolbox for streamlined analyses of
10 657 somatic mutation spectra in human and mouse cancer genomes. *BMC*
11 658 *Bioinformatics*, **17**, 170.
12 659 46. Brunet, J.-P., *et al.* (2004) Metagenes and molecular pattern discovery using matrix
13 660 factorization. *Proceedings of the National Academy of Sciences of the United States*
14 661 *of America*, **101**, 4164-4169.
15 662 47. Alexandrov, Ludmil B., *et al.* (2013) Deciphering Signatures of Mutational Processes
16 663 Operative in Human Cancer. *Cell Reports*, **3**, 246-259.
17 664 48. McCreery, M.Q., *et al.* (2015) Evolution of metastasis revealed by mutational
18 665 landscapes of chemically induced skin cancers. *Nature Medicine*, **21**, 1514-1520.
19 666 49. Nassar, D., *et al.* (2015) Genomic landscape of carcinogen-induced and genetically
20 667 induced mouse skin squamous cell carcinoma. *Nature Medicine*, **21**, 946-954.
21 668 50. Westcott, P.M.K., *et al.* (2014) The mutational landscapes of genetic and chemical
22 669 models of Kras-driven lung cancer. *Nature*, **517**, 489-492.
23 670 51. Stampfer, M.R., *et al.* (1988) Human mammary epithelial cells in culture:
24 671 differentiation and transformation. *Cancer Treat Res*, **40**, 1-24.
25 672 52. Chen, L., *et al.* (2017) DNA damage is a pervasive cause of sequencing errors,
26 673 directly confounding variant identification. *Science*, **355**, 752-756.
27 674 53. Besaratinia, A., *et al.* (2005) DNA adduction and mutagenic properties of acrylamide.
28 675 *Mutation Research*, **580**, 31-40.
29 676 54. Behjati, S., *et al.* (2014) Genome sequencing of normal cells reveals developmental
30 677 lineages and mutational processes. *Nature*, **513**, 422-425.
31 678 55. Milholland, B., *et al.* (2017) Differences between germline and somatic mutation rates
32 679 in humans and mice. *Nat Commun*, **8**, 15183.
33 680 56. Secrier, M., *et al.* (2016) Mutational signatures in esophageal adenocarcinoma define
34 681 etiologically distinct subgroups with therapeutic relevance. *Nature Genetics*, **48**,
35 682 1131-1141.
36 683 57. Olstørn, H.B.A., *et al.* (2007) Effects of perinatal exposure to acrylamide and
37 684 glycidamide on intestinal tumorigenesis in Min/+ mice and their wild-type litter mates.
38 685 *Anticancer Research*, **27**, 3855-3864.
39 686 58. Randall, S.K., *et al.* (1987) Nucleotide insertion kinetics opposite abasic lesions in
40 687 DNA. *Journal of Biological Chemistry*, **262**, 6864-6870.
41 688 59. Straif, K., *et al.* (2014) Future priorities for the IARC Monographs. *The Lancet*
42 689 *Oncology*, **15**, 683-684.

43 690
44
45

46 691 **Figure legends**

47 692 **Figure 1:** Acrylamide- and glycidamide-induced cytotoxicity and genotoxicity *in vitro*. **(A)** Cell
48 693 viability, following 24-hour treatment of primary MEFs with the indicated concentrations of
49 694 acrylamide (top panel), in the absence (diamonds) and presence (circles) of human S9
50 695 fraction, and glycidamide (bottom panel), as determined by MTT assay. Absorbance was
51 696 measured 48 hours after treatment cessation and was normalized to untreated cells. The
52 697 results are expressed as mean percent \pm SD of three replicates. **(B)** DNA damage

698 assessment by immunofluorescence with an antibody specific for Ser139-phosphorylated
699 histone H2Ax (γ H2Ax). Primary MEFs were treated with acrylamide or glycidamide for 24
700 hours prior to immunofluorescence. Compound concentrations used were based on 20-70%
701 viability reduction in the MTT assay: 10 mM acrylamide, 5 mM acrylamide in the presence of
702 S9 fraction and 3 mM glycidamide. ACR: acrylamide; GA: glycidamide.

703 **Figure 2:** Analysis of the mutation patterns derived from exome sequencing data from
704 immortalized Hupki MEF clones. **(A)** Principle component analysis (PCA) of WES data. PCA
705 was computed using as input the mutation count matrix of the clones that immortalized
706 spontaneously (Spont) or were derived from exposure to acrylamide (ACR) or glycidamide
707 (GA). Each sample is plotted considering the value of the first and second principal
708 components (Dim1 and Dim2). The percentage of variance explained by each component is
709 indicated within brackets on each axis. Spont, ACR- and GA-exposed samples are
710 represented by differently colored symbols. **(B)** Representation of small insertions and
711 deletions (indels) counts within the immortalized clones as determined by the Strelka variant
712 caller. **(C)** Mutational signatures identified by non-negative matrix factorization (NMF) in the
713 15 Hupki MEF-derived clones (sig A, sig B, and sig C). X-axis represents the trinucleotide
714 sequence context. Y-axis represents the frequency distribution of the mutations. The
715 predominant trinucleotide context for T:A > A:T mutations is indicated in sig B (5'-CTG-3').
716 The trinucleotide contexts for C:G > G:C (5'-GCC-3') and T:A > G:C mutations (5'-NTT-3')
717 are highlighted in sig C. **(D)** Contribution of the identified signatures to each sample (X-axis),
718 assigned either by absolute SBS counts or by proportion (bar graphs). The reconstruction
719 accuracy of the identified mutational signatures in individual samples is shown in the bottom
720 scatter plot (Y-axis value of 1 = 100% accuracy).

721 **Figure 3:** **(A)** Refinement of GA signature. The contribution of signature 17 (T:A>G:C in 5'-
722 NTT-3' context), present in all clones, was decreased by performing NMF on Hupki samples
723 pooled with primary tumor samples with high levels of signature 17 (see Methods). **(B)**
724 Transcription strand bias analysis for the six mutation types in GA-exposed clones. For each
725 mutation type, the number of mutations occurring on the transcribed (T) and non-transcribed
726 (N) strand is shown on the Y-axis. *** $p < 10^{-8}$; * $p < 10^{-2}$. **(C)** DNA adducts analysis as
727 determined by LC-MS/MS. Levels of N7-GA-Gua adduct in ACR+S9 and GA treated MEFs
728 and N3-GA-Ade DNA adduct level in GA treated MEFs. The data are presented as the
729 number of adducts in 10^8 nucleotides. $n \geq 2$. **(D)** Cosine similarity matrix comparing the
730 putative glycidamide mutational signature with other A>T rich mutational signatures from
731 COSMIC (signatures 22, 25, and 27) and from experimental exposure assays using specific
732 carcinogens (7,12-dimethylbenz[a]anthracene (DMBA), urethane, and aristolochic acid
733 (AA)).

1
2
3 734 **Figure 4:** GA signature in human primary cancer genome PCAWG data. **(A)** Comparison of
4 735 COSMIC signature 4 with two experimentally derived signatures (B[a]P_Exp = signature in
5 736 clones from benzo[a]pyrene treated HMEC cells; GA_Exp = signature in clones from
6 737 glycidamide-treated MEF cells). Cosine similarity between the T>N (adenine) components of
7 738 signature 4 and GA signature is shown to the right. **(B)** Transcription strand bias analysis for
8 739 the six mutation types underlying the signatures in panel A). For each mutation type, the
9 740 number of mutations occurring on the transcribed (T) and non-transcribed (N) strand is
10 741 shown on the left Y-axis. The significance is expressed as $-\log_{10}(\text{p-value})$ indicated on the
11 742 right Y-axis. *** $p < 10^{-8}$; ** $p < 10^{-4}$; * $p < 10^{-2}$. **(C)** Scatter plots show reconstruction of
12 743 COSMIC signature 4 using B[a]P- and glycidamide- experimental mutational signatures in
13 744 lung adenocarcinoma, lung squamous cell carcinoma and hepatocellular carcinoma from the
14 745 PCAWG data set. **(D)** mSigAct analysis identifies the assignment and the contributions of
15 746 mutational signatures (including the experimental signature_GA_Exp (red) and
16 747 signature_B[a]P_Exp (blue)) to the mutation burden of a total of 101 PCAWG lung and liver
17 748 tumors identified as positive for the GA signature signal.