

1

2

The Repertoire of Mutational Signatures in Human Cancer

3

Ludmil B Alexandrov^{1*}, Jaegil Kim^{2*}, Nicholas J Haradhvala^{2,3*}, Mi Ni Huang^{4*}, Alvin WT Ng⁴,
4 Yang Wu⁴, Arnoud Boot⁴, Kyle R Covington⁵, Dmitry A Gordenin⁶, Erik N Bergstrom¹, S M
5 Ashiqul Islam¹, Nuria Lopez-Bigas^{7,8,9}, Leszek J Klimczak¹⁰, John R McPherson⁴, Sandro
6 Morganella¹¹, Radhakrishnan Sabarinathan^{7,8}, David A Wheeler⁵, Ville Mustonen¹², the
7 PCAWG Mutational Signatures Working Group, Gad Getz^{2,3,13,14**}, Steven G Rozen^{4**[§]},
8 Michael R Stratton^{11**[§]}

9

10

11 ¹Department of Cellular and Molecular Medicine and Department of Bioengineering and
12 Moores Cancer Center, University of California, La Jolla, San Diego, CA, USA

13 ²Broad Institute, Cambridge, MA, USA

14 ³Center for Cancer Research, Massachusetts General Hospital, Boston, MA 02129, USA

15 ⁴Centre for Computational Biology and Programme in Cancer & Stem Cell Biology, Duke NUS
16 Medical School, Singapore

17 ⁵Human Genome Sequencing Center, and Dan L. Duncan Cancer Center, Baylor College of
18 Medicine, Houston, TX, USA

19 ⁶Genome Integrity and Structural Biology Laboratory, National Institute of Environmental
20 Health Sciences, US National Institutes of Health, Research Triangle Park, NC, USA

21 ⁷Institute for Research in Biomedicine (IRB Barcelona), The Barcelona Institute of Science
22 and Technology, Baldori Reixac, 10, 08028 Barcelona, Spain

23 ⁸Research Program on Biomedical Informatics, Universitat Pompeu Fabra, Barcelona, Spain

24 ⁹Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain

25 ¹⁰Integrative Bioinformatics Support Group, National Institute of Environmental Health
26 Sciences, US National Institutes of Health, Research Triangle Park, NC, USA

27 ¹¹Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SA, UK

28 ¹²Organismal and Evolutionary Biology Research Programme, Department of Computer
29 Science, Institute of Biotechnology, University of Helsinki, Helsinki, Finland

30 ¹³Department of Pathology, Massachusetts General Hospital, Boston, MA 02114, USA

31 ¹⁴Harvard Medical School, Boston, MA 02215, USA

32

33

34 *Denotes equal contribution

35 **Denotes equal supervisory contribution

36

37 [§]Correspondence and requests for materials should be addressed to Steven G Rozen
38 (steverozen@gmail.com) and Michael R Stratton (mrs@sanger.ac.uk).

39

40 **ABSTRACT**

41 Somatic mutations in cancer genomes are caused by multiple mutational processes each of
42 which generates a characteristic mutational signature. Using 84,729,690 somatic mutations
43 from 4,645 whole cancer genome and 19,184 exome sequences encompassing most cancer
44 types we characterised 49 single base substitution, 11 doublet base substitution, four
45 clustered base substitution, and 17 small insertion and deletion mutational signatures. The
46 substantial dataset size compared to previous analyses enabled discovery of new signatures,
47 separation of overlapping signatures and decomposition of signatures into components that
48 may represent associated, but distinct, DNA damage, repair and/or replication mechanisms.
49 Estimation of the contribution of each signature to the mutational catalogues of individual
50 cancer genomes revealed associations with exogenous and endogenous exposures and
51 defective DNA maintenance processes. However, many signatures are of unknown cause.
52 This analysis provides a systematic perspective on the repertoire of mutational processes
53 contributing to the development of human cancer including a comprehensive reference set
54 of mutational signatures in human cancer.

55

56

57 INTRODUCTION

58 Somatic mutations in cancer genomes are caused by mutational processes of both
59 exogenous and endogenous origins that have operated during the cell lineage between the
60 fertilised egg and the cancer cell¹. Each mutational process may involve components of DNA
61 damage/modification, DNA repair and DNA replication, any of which may be normal or
62 abnormal, and generates a characteristic mutational signature that may incorporate base
63 substitutions, small insertions and deletions, genome rearrangements, and chromosome
64 copy number changes². The catalogue of mutations from an individual cancer genome may
65 have been generated by multiple mutational processes and thus incorporates multiple
66 superimposed mutational signatures. Therefore, in order to systematically characterise the
67 mutational processes contributing to cancer, mathematical methods have been developed
68 that can be used to (i) decipher mutational signatures from a set of somatic mutational
69 catalogues, (ii) estimate the numbers of mutations attributable to each signature in each
70 sample, and (iii) annotate each mutation class in each tumour with the probability of arising
71 from each signature³⁻¹⁵.

72

73 Previous studies of multiple cancer types identified >30 single base substitution signatures,
74 some of known but many of unknown aetiologies, some ubiquitous and others rare, some
75 part of normal cell biology and others associated with abnormal exposures or operative
76 during neoplastic progression^{6,16-27}. Six genome rearrangement signatures have also been
77 identified in breast cancer¹⁸ and further patterns of rearrangements have been
78 described^{13,28-30}. However, analysis of other mutation classes has been relatively
79 limited^{31,17,18,32,33}.

80

81 Thus far, mutational signature analysis has predominantly used cancer exome sequences.
82 However, the many fold greater numbers of somatic mutations in whole-genome sequences
83 provide substantially increased power for signature decomposition, enabling better
84 separation of partially correlated signatures and extraction of signatures that contribute
85 relatively small numbers of mutations. Furthermore, technical artefacts and differences in
86 sequencing technologies and mutation calling algorithms can themselves generate
87 mutational signatures. Therefore, the uniformly processed and highly curated sets of all
88 classes of somatic mutations from the 2,780 cancer genome sequences of the Pan Cancer
89 Analysis of Whole-Genomes (PCAWG) project^{34,35}, combined with almost all other cancer
90 genomes and exomes for which suitable mutational catalogues are publicly available,
91 <https://www.synapse.org/#!Synapse:syn11801788>, presents a notable opportunity to
92 establish the repertoire of mutational signatures and to determine their activities across the
93 range of cancer types.

94

95

96 RESULTS

97 Cancer genomes and somatic mutations

98 Somatic mutational catalogues from 23,829 samples of most cancer types, including the
99 2,780 highly curated PCAWG whole-genomes^{34,35}, 1,865 additional whole-genomes and
100 19,184 exomes were studied. From these, 79,793,266 somatic single base substitutions,
101 814,191 doublet base substitutions and 4,122,233 small insertions and deletions (indels)
102 were analysed for mutational signatures, ~10-fold more mutations than any previous study
103 (<https://www.synapse.org/#!Synapse:syn11801889>)^{4,36}.

104

105 To enable mutational signature analysis classifications were developed for each type of
106 mutation. For single base substitutions, the primary classification comprised 96 classes
107 constituted by the six base substitutions C>A, C>G, C>T, T>A, T>C, T>G (in which the
108 mutated base is represented by the pyrimidine of the Watson-Crick base pair) plus the
109 flanking 5' and 3' bases (<https://cancer.sanger.ac.uk/cosmic/signatures/SBS>). In some
110 analyses, two flanking bases 5' and 3' to the mutated base were considered (generating
111 1,536 classes) or mutations within transcribed genome regions were selected and classified
112 according to whether the pyrimidine of the mutated base pair fell on the transcribed or
113 untranscribed strand (192 classes). A classification was also derived for doublet base
114 substitutions (78 classes, <https://cancer.sanger.ac.uk/cosmic/signatures/DBS>). Indels were
115 classified as deletions or insertions and, when of a single base, as C or T and according to the
116 length of the mononucleotide repeat tract in which they occurred. Longer indels were
117 classified as occurring at repeats or with overlapping microhomology at deletion
118 boundaries, and according to the size of indel, repeat, and microhomology (83 classes,
119 <https://cancer.sanger.ac.uk/cosmic/signatures/ID>).

120

121 **Mutational signature analysis**

122 The mutational catalogues from the 2,780 PCAWG whole-genome, 1,865 additional whole-
123 genome, and 19,184 exome sequences of cancer were analysed separately
124 (<https://doi.org/10.7303/syn11801889>)^{34,35}. For each of these catalogue sets, signature
125 extraction was conducted using methods based on nonnegative matrix factorisation
126 (NMF)^{3,6} on each cancer type individually and also on all cancer types together. Analyses
127 were carried out separately for single base substitutions (SBS signatures), doublet base
128 substitutions (DBS signatures) and indels (ID signatures) and also for the three mutation
129 types together (1697 mutation classes if the 1536 classes of SBS in pentanucleotide context
130 was employed) generating composite signatures.

131

132 Mutational signatures were extracted using two independently developed NMF-based
133 methods: (i) SigProfiler, a further elaborated version of the framework used to generate the
134 signatures shown in the previous version of the COSMIC compendium of mutational
135 signatures (COSMICv2)^{3,18,36-38}, and (ii) SignatureAnalyzer, based on a Bayesian variant of
136 NMF used in several previous publications^{6,15,39,40}. NMF determines both the signature
137 profiles and the contributions of each signature to each cancer genome as part of its
138 factorization of the input matrix of mutation spectra. However, given a substantial number
139 of signatures and/or heterogeneous mutation burdens across samples, it is possible to
140 reconstruct the mutations observed in a particular sample in multiple ways, often with very
141 small and/or biologically implausible contributions from many signatures. Therefore, each
142 method developed a separate procedure to estimate the contributions of signatures to each
143 sample (Methods).

144

145 We tested SignatureAnalyzer and SigProfiler on 11 sets of synthetic data, encompassing a
146 total of 64,400 synthetic samples, in which known signature profiles were used to generate
147 catalogues of synthetic mutational spectra. Both approaches performed well in re-extracting
148 the known signatures in realistically complex data. The tests highlighted the importance of,
149 and challenges in, selecting the number of signatures, because extracted signatures
150 discordant from the known input usually arose from difficulty in selecting the correct

151 number of signatures. Thus, these tests confirmed that use of NMF-based approaches to
152 extract signatures is not a purely algorithmic process. Instead, signature extraction requires
153 human judgement that considers all of the available data, including evidence from
154 experimental delineation of mutational signatures and the literature on DNA damage and
155 repair, and prior evidence of biological plausibility. In addition, signature extraction requires
156 human-guided sensitivity analysis to confirm that extractions from different groupings of
157 tumours yield essentially the same signatures. These types of evidence and techniques were
158 used in the determination of the signature profiles reported here. The findings we report
159 from tests on synthetic data are consistent with results regarding NMF, and the related
160 areas of probabilistic topic modelling and latent Dirichlet allocation, in multiple problem
161 domains⁴¹⁻⁴³. It is widely understood that the choice of the number of latent variables (for
162 our purposes, the number of mutational signatures) is rarely amenable to complete
163 automation. In further simulations, we also found that mutation catalogues from whole
164 genomes allowed substantially better signature extraction than the much smaller catalogues
165 from whole exomes and that signature extraction on whole genome data from half as many
166 tumours would have supported inferior signature extraction. See Methods for further
167 details; all results are at <https://doi.org/10.7303/syn18497223> and a summary can be found
168 at <https://doi.org/10.7303/syn18511087.1>.

169
170 The results of SigProfiler and SignatureAnalyzer exhibited many similarities, and we assigned
171 the same identifiers to similar signatures extracted by the two methods
172 <https://www.synapse.org/#!Synapse:syn12016215>. However, there were also noteworthy
173 differences. The number of SBS signatures found in low mutation burden tumours in the
174 PCAWG set (94.4% of cases that harbour 47% of mutations) was similar: 31 by SigProfiler
175 and 35 by SignatureAnalyzer. The number of additional SBS signatures extracted from
176 hyper-mutated PCAWG samples (5.6% of cases and 53% of mutations), however, was
177 different: 13 by SigProfiler and 25 by SignatureAnalyzer. There were also differences in SBS
178 signature profiles, including among signatures found in low mutation burden cases. The
179 latter primarily involved “flat”, relatively featureless signatures, which are mathematically
180 challenging to deconvolute. Finally, there were differences in signature attributions to
181 individual samples. In general, SignatureAnalyzer used more signatures to reconstruct the
182 mutational profiles (Extended Data Figure 1,
183 <https://www.synapse.org/#!Synapse:syn12169204>,
184 <https://www.synapse.org/#!Synapse:syn12177011>) and the attribution to flat signatures
185 was different, with SigProfiler assigning mutations to SBS5 and SBS40 and SignatureAnalyzer
186 using combinations of multiple signatures (Extended Data Figure 2ab,
187 <https://www.synapse.org/#!Synapse:syn12169204>). The DBS and ID signatures were
188 generally similar between the two methods (Extended Data Figure 2cd). These comparisons
189 provide a useful perspective on both the consistency and variability of signature extraction
190 and attribution depending on the methodology used.

191
192 The final sets of reference mutational signatures were determined from the PCAWG analysis
193 supplemented by additional signatures from the other datasets. SBS signatures using the 96
194 mutation classification were supported by the outcomes of analyses using the 192 and 1536
195 mutation classifications, the existence of individual cancer samples dominated by a
196 particular signature, and, where available, prior experimental evidence for certain
197 mutational signatures (Methods,

198 <https://doi.org/10.7303/syn12025148>, <https://doi.org/10.7303/syn12009645>, COSMIC at
199 <https://cancer.sanger.ac.uk/cosmic/signatures>). Each signature was allocated a number
200 consistent with, and extending, the COSMICv2 annotation³⁷. Some previous signatures split
201 into multiple constituent signatures and these were numbered as before but with additional
202 letter suffixes (e.g., single SBS17 split into signatures SBS17a and SBS17b). DNA sequencing
203 and analysis artefacts also generate mutational signatures, and we indicate which signatures
204 are possible artefacts (<https://www.synapse.org/#!Synapse:syn12009767>) but do not
205 present them below. However, future studies employing this signature set as a reference
206 may consider utilizing artefact signatures for data quality control. The results of both
207 SignatureAnalyzer and SigProfiler were used throughout the research reported here.
208 However, for brevity and for continuity with the signature set previously displayed in
209 COSMIC³⁷, which has been widely used as a reference, SigProfiler results are outlined below
210 and SignatureAnalyzer results are provided at (Extended Data Figures 3,4,
211 <https://www.synapse.org/#!Synapse:syn11738307>).

212

213 **Single base substitution (SBS) mutational signatures**

214 There were substantial differences in numbers of SBSs between samples (ranging from
215 hundreds to millions) and between cancer types, as previously observed⁴⁴ (Figure 1). In total,
216 67 SBS mutational signatures were extracted, of which 49 were considered to be likely of
217 biological origin (Figure 2, Methods, <https://cancer.sanger.ac.uk/cosmic/signatures/SBS/>).
218 Except for SBS25, all mutational signatures reported in COSMICv2 (i.e.,
219 https://cancer.sanger.ac.uk/cosmic/signatures_v2)^{4,23,37} were confirmed in the new set of
220 analyses (median cosine similarity between the newly derived signatures and those on
221 COSMICv2: 0.95, excluding "split" signatures which are discussed below; range 0.74 to
222 0.9996 <https://www.synapse.org/#!Synapse:syn12016215>). SBS14, SBS16, and SBS20
223 changed the most; for explanation, see <https://cancer.sanger.ac.uk/cosmic/signatures/SBS/>.
224 SBS25 was previously found only in cell lines derived from Hodgkin lymphomas, at least one
225 of which had been previously treated with chemotherapy, and, to our knowledge, no data
226 from primary cancers of this type are currently available. The newly derived signatures
227 showed much improved separation from each other and hence more distinct signature
228 profiles, presumably due to the substantially increased statistical power of this analysis
229 (online Methods section *Better separation compared to COSMICv2 signatures*).

230

231 Thirteen new likely real SBS signatures compared to the set previously described in
232 COSMICv2³⁷ were extracted (excluding those that are the consequence of signature
233 splitting). Some were in cancers with a previously unanalysed exogenous exposure (SBS42),
234 some were in chemotherapy treated samples which have often been excluded from
235 previous studies (SBS31, SBS32, SBS35) and some were rare and hence absent by chance
236 from previous analyses (SBS36, SBS44). Others were more common, but contributed
237 relatively few mutations to individual cancer genomes, or were similar to previously
238 discovered signatures and thus not isolated from datasets based predominantly on cancer
239 exome sequences (e.g., SBS38, SBS39, SBS40). Notably, SBS40 was extracted from kidney
240 cancer in which it appears to be required for optimal reconstruction of mutational
241 catalogues. It is a relatively featureless ("flat") signature, with similarity to SBS5 and other
242 flat signatures, and this may account for it only clearly emerging now with the availability of
243 whole cancer genomes. SBS40 may contribute to other cancer types but its similarity to
244 SBS5 renders this uncertain and larger datasets will be required to clarify the extent of its

245 activity. For some new signatures there were plausible underlying aetiologies (Figure 3,
246 Extended Data Figures 4,5): SBS31 and SBS35, prior platinum compound chemotherapy⁴⁵;
247 SBS32, prior azathioprine therapy; SBS36, inactivating germline or somatic mutations in
248 *MUTYH* which encodes a component of the base excision repair machinery^{46,47}; SBS38,
249 additional effects of ultraviolet light (UV) exposure; SBS42, occupational exposure to
250 haloalkanes²⁷; SBS44, defective DNA mismatch repair due to MLH1 inactivation⁴⁸. SBS33,
251 SBS34, SBS37, SBS39, SBS40, and SBS41 are of unknown cause.

252

253 Three previously characterised base substitution signatures (SBS7, SBS10, SBS17) split into
254 multiple constituent signatures (Figure 2). We previously regarded SBS7 as a single signature
255 composed predominantly of C>T at CCN and TCN trinucleotides (the mutated base is
256 underlined) together with many fewer T>N mutations. It was found in malignant melanomas
257 and squamous skin carcinomas and is likely due to UV induced pyrimidine dimer formation
258 followed by translesion DNA synthesis by error-prone polymerases which predominantly
259 insert adenine opposite damaged bases. With the larger dataset now available, SBS7 has
260 decomposed into four constituent signatures: SBS7a consisting mainly of C>T at TCN; SBS7b
261 consisting of C>T mainly at CCN and to a lesser extent at TCN; SBS7c and SBS7d, which
262 constituted relatively minor components of the previous SBS7 and consist predominantly of
263 T>A at NTT and T>C at NTT respectively⁴⁹. Splitting of a mutational signature likely reflects
264 the existence of multiple distinct mutational processes, initiated by the same exposure,
265 which have closely, but not perfectly, correlated activities. For example, the constituent
266 signatures of SBS7 are probably all initiated by UV-induced DNA damage. SBS7a and SBS7b
267 may reflect different dipyrimidine photoproducts whereas SBS7c and SBS7d may be due to
268 low frequencies of misincorporation by translesion polymerases of T and G opposite
269 thymines in pyrimidine dimers rather than the more frequent and non-mutagenic A.
270 Splitting of SBS10 and SBS17 is described at
271 <https://cancer.sanger.ac.uk/cosmic/signatures/SBS/>.

272

273 Several base substitution signatures showed transcriptional strand bias
274 (<https://www.synapse.org/#!Synapse:syn12009767>). Transcriptional strand bias is often
275 attributable to transcription coupled nucleotide excision repair (TC-NER) acting on DNA
276 damaged by exogenous exposures which cause covalently bound bulky adducts or
277 crosslinking to other bases and consequent distortion of the helical structure. This results in
278 stalling of RNA polymerase and hence recruitment of the TC-NER machinery. An excess of
279 DNA damage on untranscribed compared to transcribed strands of genes may also
280 contribute to transcriptional strand bias⁵⁰. Both mechanisms, however, result in more
281 mutations of a damaged base on the untranscribed compared to the transcribed strands of
282 genes. Assuming that either or both are responsible for the observed transcriptional strand
283 biases (which may not always be the case), DNA damage to cytosine (SBS7a, SBS7b),
284 guanine (SBS4, SBS8, SBS19, SBS23, SBS24, SBS31, SBS32, SBS35, SBS42), thymine (SBS7c,
285 SBS7d, SBS21, SBS26, SBS33) and adenine (SBS5, SBS12, SBS16, SBS22, SBS25) may underlie
286 these mutational signatures (see <https://cancer.sanger.ac.uk/cosmic/signatures/SBS/>
287 for plots of strand bias). Although the likely underlying DNA damaging agents are known for
288 SBS4 (tobacco mutagens), SBS7a, SBS7b, SBS7c, SBS7d (UV), SBS22 (aristolochic acid), SBS24
289 (aflatoxin), SBS25 (prior chemotherapy), SBS31 and SBS35 (platinum compounds), SBS32
290 (azathioprine), and SBS42 (haloalkanes), the causes of the remainder are unknown. Indeed,
291 some signatures showing transcriptional strand bias are associated with defective DNA

292 mismatch repair (SBS21 and SBS26) and it is conceivable that, for these, exogenous DNA
293 damage is not involved. The extent of transcriptional strand bias appears to differ in
294 different sectors of the genome. For example, consideration of the whole transcribed
295 genome showed absent or minimal transcriptional strand bias in the APOBEC related SBS2
296 and SBS13 and in the defective polymerase epsilon proof-reading related SBS10a. However,
297 consideration of exons alone showed clear evidence of transcriptional strand bias in these
298 signatures (<https://cancer.sanger.ac.uk/cosmic/signatures/SBS/>). The mechanism(s)
299 underlying this amplification of transcriptional strand bias in exons is unknown and appears
300 to be signature specific, since there is minimal difference in the extent of transcriptional
301 strand bias between exons and other transcribed regions for other signatures (for example,
302 SBS4 and SBS22).

303
304 Employing the single base substitution classification of 1536 mutation types, which uses the
305 pentanucleotide sequence context two bases 5' and two bases 3' to each mutated base,
306 yielded a set of signatures largely consistent with that based on substitutions in
307 trinucleotide context alone. Notably, however, the pentanucleotide context enabled the
308 extraction of two forms of both SBS2 and SBS13, one with mainly a pyrimidine (C or T) and
309 the other with a purine (A or G) at the -2 base (the second base 5' to the mutated cytosine).
310 These may represent the activities of the cytidine deaminases APOBEC3A and APOBEC3B,
311 respectively⁵¹. If so, APOBEC3A accounts for many more mutations than APOBEC3B in
312 cancers with high APOBEC activity. Several other signatures showed non-random sequence
313 contexts at +2 and -2 positions. In particular, the -2 bases in SBS17a and SBS17b and the -2
314 and +2 bases in SBS9 were predominantly A and T. In general, however, sequence context
315 effects were much stronger for bases immediately 5' and 3' to the mutated bases.

316
317 SBS signatures showed substantial variation in the numbers of cancer types and cancer
318 samples in which they were found, ranging from SBS1 and SBS5 which were present in
319 almost every cancer type and almost every cancer sample, to SBS23 which was only
320 observed in a small subset of liver cancers (Figure 3). The numbers of mutations per cancer
321 sample attributed to each signature also varied greatly, from a few tens of mutations for
322 SBS1 to millions of mutations for SBS10b. Almost all individual cancer samples exhibited
323 multiple signatures, with a mode of three signatures per sample in the PCAWG set
324 (<https://www.synapse.org/#!Synapse:syn12169204>). The assigned signatures reconstruct
325 well the mutational spectra of the cancer samples (in PCAWG samples, median cosine
326 similarity 0.97; 96.3% of samples with cosine similarity >0.90) (illustrative examples are
327 shown in Figure 4).

328 329 **Clustered single base substitution mutational signatures**

330 Some mutational processes generate mutations that cluster in small regions of the genome.
331 The relatively limited number of mutations generated by such processes, compared to those
332 acting genome-wide, may result in failure to detect their signatures by standard methods.
333 To obviate this problem, we first identified clustered mutations in each genome and
334 analysed these separately (Methods). Four main signatures associated with clustered
335 mutations were identified (Figure 2) and were consistent with previous reports^{15,16,32}. Two
336 found in multiple cancer types were similar to single base substitution SBS2 and SBS13,
337 which have been attributed to APOBEC enzyme activity (mostly APOBEC3B) and represent
338 foci of *kataegis*^{17,32,52}. Two additional clustered mutational signatures, one characterised by

339 C>T and C>G mutations at (A|G)C(C|T) trinucleotides⁵³ and the other T>A and T>C
340 mutations at (A|T)I(A|T) were found in lymphoid neoplasms and likely represent direct and
341 indirect consequences of activation induced cytidine deaminase (AID) mutagenesis and
342 translesion DNA synthesis by error-prone polymerases (SBS84 and SBS85 respectively)¹⁵.
343 The possibility that further processes may generate clustered mutations is not excluded.

344
345

346 **Doublet base substitution (DBS) mutational signatures**

347 Tandem doublet, triplet, quadruplet, quintuplet, and sextuplet base substitutions
348 (<https://www.synapse.org/#!Synapse:syn11801938>,
349 <https://www.synapse.org/#!Synapse:syn11726620>) at immediately adjacent bases were
350 observed at ~1% the prevalence of single base substitutions. In most cancer genomes, the
351 observed number of DBSs was considerably higher than expected from random adjacency of
352 SBSs (<https://www.synapse.org/#!Synapse:syn12177057>) indicating the existence of
353 commonly occurring, single mutagenic events that cause substitutions at neighbouring
354 bases. There was substantial variation in the number of DBSs, ranging from zero to 20,818 in
355 a sample. Across cancer types, the numbers of DBSs were generally proportional to the
356 numbers of SBSs in that cancer type (Figure 1). However, colorectal adenocarcinomas had
357 significantly fewer DBSs than expected, and lung cancers and melanomas had more
358 (Extended Data Table 1). The large dataset analysed here allowed, for the first time,
359 systematic analysis of DBS and indel signatures (described below). Eleven DBS signatures
360 were extracted (Figure 2). Of these, to our knowledge, only two have been previously
361 reported³³ evidencing further the value of the large numbers of mutations from whole
362 genome data.

363

364 DBS1 was characterised almost exclusively by CC>TT mutations (Figure 2), contributed 100s-
365 10,000s of mutations in malignant melanomas (Figure 3) with SBS7a and SBS7b. DBS1
366 exhibited transcriptional strand bias consistent with damage to cytosines
367 (<https://www.synapse.org/#!Synapse:syn12177063>). CC>TT mutations associated with UV
368 induced DNA damage are well established in the literature, were previously reported in
369 melanomas, and are thought to be due to generation of pyrimidine dimers and subsequent
370 error-prone translesion DNA synthesis by polymerases that introduce adenines opposite the
371 damaged bases^{33,54}.

372

373 Reanalysis after exclusion of malignant melanomas and other cancers with evidence of UV
374 exposure still yielded a signature (termed DBS11) characterised predominantly by CC>TT
375 mutations and smaller numbers of other doublet base substitutions at CC and TC which
376 contributed 10s of mutations to many samples of multiple cancer types (Figures 2 and 3).
377 DBS11 was associated with SBS2 which is due to APOBEC activity. Thus, APOBEC activity may
378 also generate DBS11, although the mechanism by which it induces doublet base
379 substitutions is not well understood.

380

381 DBS2 was composed predominantly of CC>AA mutations, with smaller numbers of CC>AG
382 and CC>AT mutations, and contributed 100s-1000s of mutations in lung adenocarcinoma,
383 lung squamous and head and neck squamous carcinomas, which are often caused by
384 tobacco smoking, which has been reported previously (Figures 2 and 3)³³. DBS2 showed
385 transcriptional strand bias indicative of guanine damage

386 (<https://www.synapse.org/#!Synapse:syn12177064>) and was associated with SBS4 which is
387 caused by tobacco smoke exposure. It is likely, therefore, that DBS2 can be a consequence
388 of DNA damage by tobacco smoke mutagens.

389

390 Analysis of each cancer type separately, however, revealed a signature very similar to DBS2
391 contributing 100s of mutations to liver cancers and 10s of mutations to cancers of other
392 types without evidence of tobacco smoke exposure. A pattern closely resembling DBS2 and
393 characterised predominantly by CC>AA mutations, together with smaller contributions of
394 CC>AG and CC>AT, dominates DBSs in normal mouse cells and is particularly frequent in the
395 liver⁵⁵. The nature of the mutational processes underlying these doublet signatures in
396 smoking-unrelated human cancers and in normal mice is unknown. However, acetaldehyde
397 exposure in experimental systems generates a mutational signature characterised primarily
398 by CC>AA and lower burdens of CC>AG and CC>AT mutations together with C>A single base
399 substitutions⁵⁶. Acetaldehyde is an oxidation product of alcohol and a constituent of
400 cigarette smoke. The role of acetaldehyde, and perhaps other aldehydes, in generating
401 DBS2, whether associated with tobacco smoking, alcohol consumption or in non-exposed
402 cells, merits further investigation⁵⁷.

403

404 DBS3, DBS7, DBS8 and DBS10 showed 100s-1000s of mutations in rare colorectal, stomach
405 and oesophageal cancers some of which showed evidence of defective DNA mismatch
406 repair (DBS7, DBS10) or polymerase epsilon exonuclease domain mutations (DBS3)
407 generating hypermutator phenotypes (Figures 2, 3). DBS5 was found in cancers previously
408 exposed to platinum chemotherapy and is associated with SBS31 and SBS35. The remaining
409 DBS signatures are of uncertain cause.

410

411 **Small insertion and deletion (ID) mutational signatures**

412 Indels were usually present at ~10% the frequency of base substitutions (Figure 1). There
413 was substantial variation between cancer genomes in numbers of indels, even when cancers
414 with evidence of defective DNA mismatch repair were excluded. Overall, the numbers of
415 deletions and insertions were similar, but there was variation between cancer types with
416 some showing more deletions and others more insertions of various subtypes (Figure 1).
417 Seventeen indel mutational signatures were extracted (Figure 2).

418

419 Indel signature 1 (ID1) was composed predominantly of insertions of thymine and ID2 of
420 deletions of thymine, both at long (≥ 5) thymine mononucleotide repeats (Figure 2). 10s to
421 100s of mutations of both signatures were found in the large majority of most cancer types
422 but were particularly common in colorectal, stomach, endometrial and oesophageal cancers
423 and in diffuse large B cell lymphoma (Figure 3). Most of these cancers are likely to be DNA
424 mismatch repair proficient on the basis of the relatively limited numbers of indels and
425 absence of the SBS signatures (SBS6, SBS14, SBS15, SBS20, SBS21, SBS26, and SBS44)
426 associated with DNA mismatch repair deficiency. Together, ID1 and ID2 accounted for 97%
427 and 45% of indels in hypermutated and non-hypermutated cancer genomes, respectively
428 (Extended Data Table 2), and both signatures have also been found in non-neoplastic cells⁵⁸.
429 They are likely due to the intrinsic tendency to slippage during DNA replication of long
430 mononucleotide tracts. However, the mechanistic basis for separation into two signatures,
431 one presumably due to slippage of the nascent strand (ID1) and the other the template

432 strand (ID2) is unclear. Similarly, the substantial differences in their mutation frequencies
433 between cancer types are not well understood.

434

435 ID3 was characterised predominantly by deletions of cytosine at short (≤ 5 bp long)
436 mononucleotide cytosine repeats and exhibited 100s of mutations in tobacco smoking
437 associated cancers of the lung and head and neck (Figures 2 and 3). There was
438 transcriptional strand bias of mutations, with more guanine deletions than cytosine
439 deletions on the untranscribed strands of genes, compatible with TC-NER of adducted
440 guanine

(<https://www.synapse.org/#!Synapse:syn12177065>,
441 <https://www.synapse.org/#!Synapse:syn12177066>). The numbers of ID3 mutations in

442 cancer samples positively correlated with the numbers of SBS4 and DBS2 mutations, both of
443 which have been associated with tobacco smoking (Extended Data Figure 6). It is therefore
444 likely that DNA damage by components of tobacco smoke underlie ID3 but the
445 mechanism(s) by which indels are generated is unclear.

446

447 ID13 was characterised predominantly by deletions of thymine at thymine-thymine
448 dinucleotides and exhibited large numbers of mutations in malignant melanomas of the skin
449 (Figures 2 and 3). The numbers of ID13 mutations correlated with the numbers of SBS7a,
450 SBS7b and DBS1 mutations, which have been attributed to DNA damage induced by UV
451 (Extended Data Figure 6). It is, however, notable that a similar mutation of the other
452 pyrimidine, i.e., deletion of cytosine at cytosine-cytosine dinucleotides, does not feature
453 strongly in ID13, perhaps reflecting the predominance of thymine compared to cytosine
454 dimers induced by UV⁵⁹. The mechanism(s) underlying thymine deletion is unclear.

455

456 ID6 and ID8 were both characterised predominantly by deletions ≥ 5 bp (Figure 2). ID6
457 exhibited overlapping microhomology at deletion boundaries with a mode of 2bp and often
458 longer stretches. This signature was correlated with SBS3 which has been attributed to
459 defective homologous recombination based repair (Extended Data Figure 6). By contrast,
460 ID8 deletions showed shorter or no microhomology at deletion boundaries, with a mode of
461 1bp, and did not strongly correlate with SBS3 mutations (Figures 2 and 3). These patterns of
462 deletion may be characteristic of DNA double strand break repair by non-homologous
463 recombination based end-joining mechanisms, and if so, suggest that at least two distinct
464 forms of end-joining mechanism are operative in human cancer⁶⁰.

465

466 A small fraction of cancers exhibited very large numbers of ID1 and ID2 mutations ($>10,000$)
467 (Figure 3, <https://cancer.sanger.ac.uk/cosmic/signatures/ID>). These were usually
468 accompanied by SBS6, SBS14, SBS15, SBS20, SBS21, SBS26 and/or SBS44 which are
469 associated with DNA mismatch repair deficiency, sometimes combined with POLE or POLD1
470 proofreading deficiency (SBS14, SBS20)⁴⁰. Occasional cases with these signatures
471 additionally showed large numbers of ID7 indels
472 (<https://www.synapse.org/#!Synapse:syn11738668>). In addition, rare samples showed large
473 numbers of either ID4, ID11, ID14, ID15, ID16 or ID17 mutations but did not show ID1 and
474 ID2 mutations or the single base substitution signatures usually associated with DNA
475 mismatch repair deficiency. The mechanisms underlying these signatures are unknown.

476

477 **Composite mutational signatures**

478 In the analyses described above mutational signatures were extracted for each mutation
479 type separately. However, mutational processes in nature generate composite signatures
480 that may include SBSs, DBSs, IDs, genome rearrangements and chromosome number
481 changes. We therefore also extracted signatures using combined catalogues of SBSs, DBSs,
482 and IDs (257 mutation subclasses or 1697 if the 1536 classification of single base
483 substitutions was used). Fifty-two composite signatures were extracted.

484
485 A composite signature with components similar to SBS4, DBS2 (characterised predominantly
486 by CC>AA mutations) and ID3 (characterised predominantly by deletion of cytosine at short
487 runs of cytosines) was found mainly in lung cancers, suggesting that it is the consequence of
488 tobacco smoke exposure (Extended Data Figure 7). Similarly, composite signatures with
489 components similar to SBS7a, SBS7b, DBS1 (characterised predominantly by CC>TT
490 mutations) and ID13 (characterised predominantly by deletion of thymine at thymine–
491 thymine dinucleotides) were found in skin cancers and are thus likely due to UV induced
492 DNA damage (Extended Data Figure 7). A further composite signature in breast and ovarian
493 cancers included features of SBS3 and ID6 combined with ID8 (deletions >5bp with varying
494 degrees of overlapping microhomology) and is likely associated with defective homologous
495 recombination based repair (Extended Data Figure 7). In these composite signatures
496 attributions of the constituent SBS, DBS and ID signatures extracted independently in the
497 main analyses were correlated with each other, adding support to the existence of the
498 composite signatures (Extended Data Figure 6). Various forms of defective DNA mismatch
499 repair were also associated with multiple SBS, DBS and ID signatures.

500

501 **Correlations with age**

502 A positive correlation between age of cancer diagnosis and the number of mutations
503 attributable to a signature suggests that the mutational process underlying the signature
504 has been operative, at a more or less constant rate, throughout the cell lineage from
505 fertilized egg to cancer cell, and thus in normal cells from which that cancer type
506 develops^{4,61}. Confirming previous reports, the numbers of SBS1 and SBS5 mutations
507 correlated with age, exhibiting different rates in different tissue types (q-values in
508 <https://www.synapse.org/#!Synapse:syn12030687>,
509 <https://www.synapse.org/#!Synapse:syn20317940>
510 <https://www.synapse.org/#!Synapse:syn12217988>). In addition, SBS40 correlated with age
511 in multiple cancer types. However, given the similarity in signature profile between SBS5
512 and SBS40 the possibility of misattribution between these signatures cannot currently be
513 excluded. The numbers of DBSs and IDs were much lower than the numbers of SBSs and the
514 numbers of samples in which DBS and ID signatures could be attributed were also lower.
515 Nevertheless, DBS2 and DBS4 correlated with age and, consistent with the interpretation of
516 activity in normal cells, the profiles of DBS2 and DBS4 together closely resemble the
517 spectrum of DBS mutations found in normal mouse cells⁵⁵. Neither DBS2 nor DBS4,
518 however, was clearly correlated with an SBS or ID signature that correlates with age. ID1,
519 ID2, ID5 and ID8 showed correlations with age in multiple tissues. ID1 and ID2 indels are
520 likely due to slippage at poly T repeats during DNA replication and correlated with the
521 number of SBS1 substitutions. SBS1 has previously been proposed to reflect the number of
522 mitoses a cell has experienced and thus SBS1, ID1 and ID2 may all be generated during DNA
523 replication at mitosis⁴. The number of ID5 mutations correlated with the number of SBS40
524 mutations and thus the mutational processes underlying these two age-correlated

525 signatures may also harbour common components. ID8 is predominantly composed of
526 deletions >5bp with no or 1bp of microhomology at their boundaries. These are likely due to
527 DNA double strand breaks which have not been repaired by homologous recombination
528 based mechanisms, but instead by a non-homologous-end joining mechanism. The features
529 of ID8 resemble those of some ionising radiation associated mutations and this may,
530 therefore, be an underlying aetiological factor⁶². Taken together, the results indicate that
531 multiple mutational processes operate in normal cells.

532

533

534 **DISCUSSION**

535 Cancers arise as a result of somatic mutations. Mutational signature analysis therefore
536 provides important insights into cancer development through comprehensive
537 characterisation of the underlying mutational processes. There are, however, important
538 constraints, limitations and assumptions in the analytic frameworks we have used that
539 should be recognised. Although designed to reflect the mutational consequences of
540 recurrent mutational processes, mutational signatures extracted from sample sets in which
541 multiple mutational processes are operative remain mathematical approximations, with
542 profiles that can be influenced by the mathematical approach used and by additional
543 factors, such as the other mutational processes present. For conceptual and practical
544 simplicity, we have assumed that there is a single signature associated with each mutational
545 process and have provided an average reference signature to represent it. However, we do
546 not discount the possibility that further nuances and variations of signature profiles exist,
547 for example between different tissues. Moreover, although the extent of separation
548 between partially correlated signatures has been improved in this analysis, some signatures
549 may still represent combinations of constituent signatures. Contributions from each
550 signature to the burden of mutations in each sample have been estimated. However, with
551 increasing numbers of signatures and multiple orders of magnitude differences in mutation
552 burdens from certain signatures, prior knowledge can help to avoid biologically implausible
553 results. Thus, further development of methods for deciphering mutational signatures and
554 attribution of mutations is warranted and this needs to be supplemented by signatures
555 derived from experimental systems in which the causes of the mutations are known. The
556 numbers of DBSs, clustered substitutions, IDs and genome rearrangements (reported in ref.
557 ³⁰) are small compared to single base substitutions. Thus, larger datasets may be required to
558 robustly characterise their mutational signatures. Nevertheless, the results outlined here
559 indicate that signatures with many similarities and some differences can be found by
560 different mathematical approaches, and that these are confirmed in many different ways,
561 including experimentally elucidated signatures^{22,31,45,48,49,61,63-69} and the observation of
562 tumours dominated by a single signature
563 (<https://www.synapse.org/#!Synapse:syn12016215>).

564

565 Prior reports have provided only a relatively limited examination of doublet and indel
566 mutational spectra and, to the best of our knowledge, no previous comprehensive analysis
567 of doublet and indel mutational signatures has been performed. Here, we provide the first
568 systematic analysis of these mutation types by considering 83 mutational subtypes for
569 indels and 78 mutational subtypes for doublets. This analysis also includes almost all publicly
570 available exome and whole-genome cancer sequences, amounting in aggregate to 23,829
571 cancers of most cancer types. Some rare or geographically restricted signatures may not

572 have been captured and signatures of therapeutic mutagenic exposures have not been
573 exhaustively explored. Nevertheless, it is likely that a substantial proportion of the naturally-
574 occurring mutational signatures found in human cancer have now been described. This
575 comprehensive repertoire provides a foundation for future research into *(i)* geographical
576 and temporal differences in cancer incidence to elucidate underlying differences in
577 aetiology, *(ii)* the mutational processes and signatures present in normal tissues and caused
578 by non-neoplastic disease states, *(iii)* clinical and public health applications of signatures as
579 indicators of sensitivity to therapeutics and past exposure to mutagens, and *(iv)* mechanistic
580 understanding of the mutational processes underlying carcinogenesis.
581

582 **ACKNOWLEDGEMENTS**

583 The results here are partly based on data generated by the TCGA Research Network
584 (<http://cancergenome.nih.gov/>) and the ICGC/TCGA Pan-cancer Analysis of Whole Genomes
585 Network. This work was supported by Wellcome grant reference 206194 (M.R.S.), Singapore
586 National Medical Research Council grants NMRC/CIRG/1422/2015 and MOH-000032/MOH-
587 CIRG18may-0004 and the Singapore Ministry of Health via the Duke-NUS Signature Research
588 Programmes (M.N.H., A.W.T.N., Y.W., A.B., S.G.R.), US National Institute of Health
589 Intramural Research Program Project Z1AES103266 (D.A.G.), the European Research Council
590 Consolidator Grant 682398 (N.L.-B.), US National Cancer Institute U24CA143843 (D.A.W.),
591 and Cancer Research UK Grand Challenge Award C98/A24032 (E.N.B., S.M.A.I., L.B.A.,
592 M.R.S.). G.G and J.K were partially supported by the National Cancer Institute grants
593 U24CA210999 and U24CA143845. G.G. was partially supported by the Paul C. Zamecnick,
594 MD, Chair in Oncology at the Massachusetts General Hospital Cancer Center. N.J.H and G.G
595 were partially supported by G.G's start up funds at MGH.

596

597

598

599 **Figure legends.**

600 **Figure 1. Mutation burdens of single base substitutions, doublet base substitutions and**
601 **small insertions and deletions for the 2,780 PCAWG tumours.** Each sample is displayed
602 according to its tumour type. Tumour types are ordered according to the median number of
603 single base substitutions. The numbers of cases of each tumour type are shown. The
604 proportions of each mutation subclass in each sample are shown as coloured bar charts.

605
606 **Figure 2. Profiles of single base substitution, doublet base substitution and small insertion**
607 **and deletion mutational signatures.** The subclassifications of each mutation type (single
608 base substitutions, 96 subtypes; doublet base substitutions, 78 subtypes; indels, 83
609 subtypes) are described in the main text. Magnified versions of signatures SBS4, DBS2 and
610 ID3 (which are all associated with tobacco smoking) are shown to illustrate the positions of
611 each mutation subtype on each plot.

612
613 **Figure 3. The number of mutations contributed by each mutational signature to the 2,780**
614 **PCAWG tumours.** The numbers of mutations attributed are shown by cancer type. The size
615 of each dot represents the proportion of samples of each tumour type that show the
616 mutational signature. The colour of each dot represents the median mutation burden of the
617 signature in samples which show the signature. Contributions are shown for single base
618 substitution, doublet base substitution and indel mutational signatures separately.
619 Contributions of composite signatures to the PCAWG cancers and single base substitution
620 signatures to the complete set of cancer samples analysed are shown in Supplementary
621 information.

622
623 **Figure 4. Illustrative examples of mutational spectra of individual cancer samples.** A breast
624 cancer, a lung cancer, and a malignant melanoma and their contributory single base
625 substitution, doublet base substitution, and small insertion and deletion mutational
626 signatures.

627

628 **Online Methods**

629

630 ***Principles and strategy of mutational signature analysis adopted in this report***

631 *Conceptual principles.*

- 632
- 633 • Multiple mutational processes generate the somatic mutations present in each individual human cancer.
 - 634 • Each mutational process generates a particular pattern of somatic mutations known as a mutational signature.
 - 635
 - 636 • Each mutational process may incorporate a component of DNA damage/modification, DNA repair and DNA replication, each of which may be part of normal or abnormal cell biology. Differences in any of the three components may result in a different mutational signature, thus, by definition, constituting a distinct mutational process.
 - 637
 - 638
 - 639
 - 640
 - 641 • Multiple mutational processes operating continuously or intermittently during the cell lineage from the fertilised egg to the cancer cell may contribute to the aggregate set of mutations found in the cancer cell. Thus, the catalogue of somatic mutations from a single cancer sample often includes mutations of many different mutational signatures.
 - 642
 - 643
 - 644
 - 645
 - 646

647 *Aims of the study.*

- 648
- 649 • To decipher the mutational signatures present in essentially the full set of whole genome and exome sequenced human cancers from which data is currently available and subsequently to estimate the contributions of each signature to each cancer genome.
 - 650
 - 651
 - 652

653 *Approach used.*

- 654
- 655 • Several mathematical approaches have been used to deconvolute/extract the mutational signatures present in a set of mutational catalogues^{3,6,7,9,14-16,39,70-72}. They are all based on the premise that different mutational processes (and thus their signatures) contribute to different extents to different samples within the set.
 - 656
 - 657
 - 658 • Two independently developed methods based on NMF (SigProfiler and SignatureAnalyzer) were applied separately to the sets of mutational catalogues. By using two methods we aimed to provide perspective on the impact different methodologies can have on numbers of signatures generated, signature profiles and attributions. The two methods are described in detail below and the code for both is available (<https://www.synapse.org/#!Synapse:syn11801488>). Results from the two methods have been compared (<https://www.synapse.org/#!Synapse:syn12177006>).
 - 659
 - 660
 - 661
 - 662
 - 663
 - 664
 - 665 • Briefly, SigProfiler employs an elaboration of previously presented approaches for signature extraction and for attribution of mutation counts to mutational signatures in individual tumours^{3,4,18,36}.
 - 666
 - 667
 - 668 • Briefly, SignatureAnalyzer employs a Bayesian variant of NMF^{6,15,39}. This method enables inferences for the number of signatures through the automatic relevance determination technique and delivers highly interpretable and sparse representations for both signature profiles and attributions at a balance between data fitting and model complexity.
 - 669
 - 670
 - 671
 - 672

- 673
- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- The methods that SigProfiler and SignatureAnalyzer use for determining the number of extracted signatures are presented in the detailed descriptions of each of these methods, below.
 - Both methods assume that the spectra of individual tumours can be represented as linear combinations of signatures. Thus, if the combination of two simultaneously operating mutational processes were to create a signature profile that is not a linear combination of the two, both SigProfiler and SignatureAnalyzer would extract this as a separate signature. We believe this is the case for SBS20, which appears to be due to the simultaneous operation of *POLD1* mutation and mismatch repair deficiency.

683 *Role of NMF in extraction and attribution of mutational signatures.*

- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707
- 708
- 709
- 710
- 711
- 712
- 713
- 714
- 715
- 716
- 717
- NMF is the approximate representation of a nonnegative matrix V , in this case the observed mutational spectra (or profiles) of a set of tumors, as the product of two usually smaller nonnegative matrices, W and H , which are the signatures and the attributions respectively.
 - In our experience, however, calculating a single NMF is rarely sufficient to allow confident extraction and attribution of signatures that reflect the underlying biological mutational processes. There are two main reasons for this:
 - The profiles of extracted signatures can vary substantially depending on the tumour samples present in V . For example, this may be especially evident when some tumors in V have high numbers of mutations (e.g., samples due to UV exposure or DNA mismatch repair deficiency), while others have low numbers. In situations such as this, signatures due to highly mutagenic processes sometimes capture mutations from other processes and also "bleed" into other signatures.
 - With multiple potentially similar signatures operating, there are multiple possible and reasonably accurate reconstruction solutions for each tumour, often with many small and/or biologically implausible contributions.
 - To address these challenges two key additional analytic features have been incorporated into our analyses:
 - Both SigProfiler and SignatureAnalyzer carried out multiple NMFs on different subsets of tumours for signature extraction, and indeed, each signature extraction by SigProfiler entails 1024 NMFs with different random initial conditions. We describe below how we selected representative mutational signature profiles.
 - Both SigProfiler and SignatureAnalyzer developed a process of attributing signature activities to tumours that is separate from the process of extracting (discovering) the signatures.
 - The use of multiple extractions to support confidence in results:
 - SignatureAnalyzer, carried out the main extraction procedure on (1) the majority of the PCAWG tumours excluding certain highly mutated tumours and (2) the melanomas, microsatellite-unstable tumours, and a single temozolomide-exposed tumour (<https://www.synapse.org/#!Synapse:syn11738314>).
 - SigProfiler extracted signatures from

- 718 ▪ Separate extraction of SBS, DBS, and ID signatures from all PCAWG
719 whole-genomes together (the main source of the reference
720 mutational signature).
- 721 ▪ Separate extraction of SBS, DBS, and ID signatures from PCAWG
722 whole-genomes with each tumour type examined by itself.
- 723 ▪ Extraction of SBS signatures from all non-PCAWG whole-genomes
724 together.
- 725 ▪ Extraction of SBS signatures from non-PCAWG whole-genomes with
726 each tumour type examined by itself.
- 727 ▪ Separate extraction of SBS and ID signatures from all TCGA exomes
728 together.
- 729 ▪ Separate extraction of SBS and ID signatures from TCGA exomes with
730 each tumour type examined by itself.
- 731 ▪ Separate extraction of SBS and ID signatures from all non-TCGA
732 exomes together.
- 733 ▪ Separate extraction of SBS and ID signatures from non-TCGA exomes
734 with each tumour type examined by itself.
- 735 This allowed the extraction of signatures that were not present in the PCAWG
736 tumours (e.g., SBS42, which has been attributed to haloalkane exposure and
737 seen only in whole exome sequencing data). It also served as an important
738 validation, as extraction of similar signatures from single tumour types and
739 other sample sets supports the correctness of the signature extracted from
740 the PCAWG samples (<https://www.synapse.org/#!Synapse:syn12016215>).
- 741 ○ Signature extraction from each tumour type (or from some other subset of
742 cancers) separately has the advantages of:
- 743 ▪ Usually including fewer (and different) mutational signatures in each
744 tumour type sample set than in the set of all cancers together and
745 thus fewer (and different) opportunities for inter-signature
746 interference.
- 747 ▪ Allowing multiple independent opportunities for extraction of a
748 signature that is present in multiple tumour types, and thus of
749 obtaining validation/confirmation of the signature's existence and
750 profile.
- 751 ▪ Allowing extraction of a signature that may (for a number of reasons)
752 fail to be extracted in analysis of all tumour types together.
- 753 ▪ Providing primary evidence for the existence of the signature in each
754 tumour type.
- 755 ▪ Allowing separation of highly mutated cancer types/samples from
756 cancer types/samples with low mutation burdens.
- 757 ○ Signature extraction from multiple tumour types together has the
758 advantages of:
- 759 ▪ Usually including more samples with a particular signature than in
760 each individual cancer type and thus being better powered to
761 separate a signature from other partially correlated signatures and/or
762 from signatures with similar profiles.

- 763 ▪ Providing a single profile for a signature rather than the multiple
764 slightly different profiles which emerge from extraction of each
765 tumour type separately.
- 766 • The profiles of the mutational signatures extracted from cancer are highly variable.
767 They range from some that have contributions from mutations of all subtypes in the
768 mutation classification (“flat” or “featureless” signatures, e.g., SBS5 and SBS40) to
769 others that are essentially defined by mutations at only one (or a small number) of
770 the mutation subtypes (e.g., signatures SBS2, SBS13, SBS10a and SBS10b). There
771 appears to be less concordance between the results of SigProfiler and
772 SignatureAnalyzer for flat signatures than for signatures with distinct features
773 indicating that generally, these may be more difficult to accurately extract and
774 distinguish from each other. However, there is experimental support for the
775 existence of SBS5 and SBS3^{61,68}.
 - 776 • We represented each signature as a single reference. This selection of a single
777 reference signature does not exclude the possibility that signature profiles may show
778 nuances and further complexity and may vary in different contexts (e.g., in different
779 tissues). The rationale for selecting a single reference signature was the view that
780 this would be a level of granularity useful to most researchers. For those with
781 specialised interests in particular mutational processes and their components, we
782 also provided the signatures extracted from individual tumour types, comprising
783 PCAWG and non-PCAWG genomes and exomes
784 (<https://www.synapse.org/#!/Synapse:syn12025142>).
 - 785 • Attribution of signatures to cancer samples:
 - 786 ○ The reference signatures from SigProfiler and SignatureAnalyzer were used to
787 estimate the number of mutations due to each signature in each tumour
788 (<https://www.synapse.org/#!/Synapse:syn11804065>).
 - 789 ○ SigProfiler and SignatureAnalyzer differ in their approaches for attributing
790 signatures. However, both incorporate a set of rules based on prior
791 knowledge and biological plausibility, and incorporate techniques to
792 encourage sparsity in the number of signatures attributed to a given tumour.
 - 793 ○ Sparsity (limiting the numbers of signatures and limiting the numbers of
794 signatures attributed to each cancer sample) is an important concept and
795 feature of both SigProfiler and SignatureAnalyzer (both in signature
796 extraction and attribution). Our prior beliefs are that (i) there is a limited set
797 of significantly contributing mutational processes (and hence a limited set of
798 mutational signatures) operating to generate somatic mutations across all
799 cancers and (ii) that a limited set of mutational processes contribute to
800 individual cancer genomes (as opposed to all mutational signatures
801 contributing to all samples). Our aim in discovering mutational signatures is
802 to reflect the underlying biological processes and to attribute them
803 appropriately. It is not a mathematical exercise in which the main objective
804 and priority is to minimize the difference between $W \times H$ and the original
805 spectra in V . Indeed, if the latter was the main aim, for 96 mutation classes a
806 set of 96 signatures each constituted entirely of mutations in just one class
807 (and therefore ignoring sparsity), will always provide error free
808 reconstruction but will provide absolutely no information about underlying
809 mutational processes.

810

811 *Presentation of the results of signature extraction and attribution from SigProfiler and*
812 *SignatureAnalyzer.*

- 813 • The results (signatures and attributions) of the two methods have been presented
814 separately. We have done this in preference to combining them. We have handled
815 the two outputs in this way because we believe that this provides a simpler
816 conceptual and technical basis on which the research community can understand
817 the results, can employ the methods in future and can compare results with those
818 shown in this paper. We also do not have a basis for believing that a
819 combined/averaged/overlapping single result set is a better representation of the
820 natural truth than either of the two result sets individually and do not have a well-
821 founded and simple technical approach for combining them. We have, however,
822 provided comparisons of the outputs.
- 823 • For brevity and for continuity with previous publications, the results from SigProfiler,
824 a further elaborated version of previously described approaches^{3,4,18,36} that
825 generated the 30 signatures previously shown in COSMICv2³⁷, are shown in the main
826 manuscript, and the results from SignatureAnalyzer in supplementary data
827 (<https://www.synapse.org/#!/Synapse:syn11738307>).
- 828 • Nomenclature of signatures is based on and extends the nomenclature previously
829 used in COSMIC (COSMICv2, https://cancer.sanger.ac.uk/cosmic/signatures_v2)³⁷.
- 830 • Both methods analysed each mutation type (SBSs, DBSs and IDs) separately and also
831 together as a composite signature. In future, however, SigProfiler will usually use the
832 separately extracted single base substitution, indel and doublet base substitution
833 signatures as its standard. This generally facilitates portability, and comparison of
834 signature profiles with those from a variety of sample sets including targeted
835 sequences, exomes etc.
- 836 • SBS signatures reported in Supplementary Data include possible artefacts
837 (<https://cancer.sanger.ac.uk/cosmic/signatures/SBS/> and see below).

838

839 *Quality control: annotating signatures as likely real or a possible artefact*

- 840 • Sequencing artefacts and differences in analysis pipelines can also generate
841 mutational signatures. We have annotated which signatures are likely real or
842 “possible artefact”.
- 843 • There are multiple reasons for believing a signature reflects a biological mutational
844 signature rather than an artefact.
 - 845 ○ The input data supporting the signature seem correct: key mutational
846 features of the putative signature look real in a mapped-read browser such as
847 Integrative Genomics Viewer (IGV,
848 <https://software.broadinstitute.org/software/igv/>), or characteristic mutations
849 are experimentally confirmed in the tumour and normal samples. Inspection
850 in a mapped read browser is especially important in checking for possible
851 problems in potentially new signatures arising in datasets other than the
852 highly scrutinized and checked PCAWG and TCGA sets. Features associated
853 with experimental, mapping, or other computational artefacts include strong
854 preference for the first read, very low variant allele fractions, variants in
855 regions of low germ-line sequencing coverage, variants found near indels in

- 856 low-complexity regions, variants from a signature only found in one
857 sequencing centre etc.
- 858 ○ The 96-mutation profile and additional features (e.g., strand asymmetry,
859 association with replication timing), are known to result from a particular
860 process in experimental systems. Examples: UV, polymerase epsilon
861 proofreading deficiency, aristolochic acid and cisplatin exposure.
 - 862 ○ The putative signature is broadly consistent with previous biochemical
863 knowledge of mutational processes (e.g., preference for G adducts in
864 aflatoxin).
 - 865 ○ The putative signature dominates the spectra of some tumours (column J of
866 <https://www.synapse.org/#!/Synapse:syn12016215>).
 - 867 ○ The putative mutational signature is consistently deciphered from multiple
868 independent datasets; this indicates that the signatures is either a common
869 sequencing artefact or something real.
 - 870 ○ The putative signature correlates with known or suspected mutational
871 exposures, endogenous processes, or repair defects, especially if some of
872 those exposures/processes/repair defects result in overwhelming mutational
873 spectra. Examples: melanoma / fair skin / UV exposure, POLE mutations,
874 MMR deficiency and APOBEC germ line variants.
 - 875 ○ The putative signature correlates with other clinical characteristics, such as
876 age at diagnosis (examples SBS1 and SBS5) or tobacco smoking (SBS4).
 - 877 ○ The mutational signature exhibits a strong transcriptional strand bias; it is
878 hard to imagine an artefact with transcriptional strand bias.
 - 879 ○ The putative signature shows association with other genomic features, such
880 as microindels in homopolymers, replication strand, replication timing, or
881 nucleosome occupancy.

882

883 *Cancer sample sets on which different analyses have been conducted.*

884

- 884 • Because PCAWG genomes are of high quality with respect to the calling of all
885 mutation types, all our analyses (all types of signature extraction and all types of
886 signature attribution) have been conducted on the 2,780 PCAWG genomes.
- 887 • SigProfiler also extracted SBS signatures from the non-PCAWG whole genomes,
888 TCGA exomes, and non-TCGA exomes and attributed SBS signatures to them.
- 889 • ID signatures have been extracted and attributed to PCAWG genomes and to a
890 subset of TCGA exomes with large numbers of indels (the latter SigProfiler only). We
891 have not done this for indels in non-PCAWG whole genome sequences and non-
892 TCGA exomes (*i*) because of the unknown and variable accuracy and standardisation
893 of indel mutation calls from different groups generating the data, (*ii*) because in
894 some cases no indel calls were provided by the data generator and (*iii*) because for
895 exomes in most cases there would be very few mutations.
- 896 • DBS signatures have been extracted and attributed to PCAWG genomes only. We
897 have not done this for the other categories of samples because of the unknown and
898 variable quality of the mutation calls, the possibility that filters introduced for quality
899 control might deliberately exclude doublet mutations, and the small numbers of
900 doublet mutations in exomes.
- 901 • Consistent with the above, composite mutational signatures have only been
902 extracted and attributed for PCAWG genomes.

903

904 *Splitting of mutational signatures.*

- 905 • Certain previously existing single signatures have split into multiple constituent
906 signatures in this analysis. This is likely due to the existence of multiple, partially
907 correlated mutational processes with the same initiating factor (for example, UV
908 exposure) but subsequent differences in underlying mechanisms which differ in
909 intensity in different tissues or other contexts. A previous example of this for which
910 we have allocated different signature numbers is the split of the usually co-occurring
911 but independently varying consequences of APOBEC mutagenesis into signatures
912 SBS2 and SBS13 (<https://cancer.sanger.ac.uk/cosmic/signatures/SBS/>).
- 913 • Depending on the extent of correlation of the two signatures, and the available
914 dataset/statistical power such signatures may manifest as a single signature,
915 overlapping partially separated signatures or as two separate signatures.
- 916 • We are aware that splitting of signatures can also be a mathematical artefact.
917 However, we have used multiple extractions to confirm and validate signature splits
918 and applied the principle of sparsity to limit artefactual splits
919 (<https://cancer.sanger.ac.uk/cosmic/signatures/SBS/>).

920

921 ***Better separation compared to COSMICv2 signatures***

922 As described in the manuscript, all mutational signatures previously reported on COSMIC
923 were confirmed in the new set of analyses with median cosine similarity of 0.95. However,
924 the separation between the COSMICv2 mutational signatures
925 (https://cancer.sanger.ac.uk/cosmic/signatures_v2) is much worse compared to the
926 separation between the PCAWG mutational signatures. One can easily discern this by visual
927 examination of signature profiles. For example, in COSMICv2, signatures 5 and 16 have a
928 cosine similarity of 0.90, thus making them hard to distinguish from one another. In
929 contrast, in the current PCAWG analysis, SBS5 and SBS16 have a cosine similarity of 0.65.
930 This allows unambiguously assigning SBS5 and SBS16 to different samples. In the PCAWG
931 analysis, the larger number of samples has allowed reducing the bleeding between
932 signatures and has given more unique and easily distinguishable signatures. One can
933 evaluate the overall separation of a set of mutational signatures by examining the
934 distribution of cosine similarities between the signatures in the set. The COSMICv2
935 signatures have a median cosine similarity between the signatures in COSMICv2 of 0.238. In
936 contrast, the PCAWG signatures have a much lower median cosine similarity between the
937 signatures in PCAWG of 0.098. This 2-fold reduction in similarity is highly statistically
938 significant (p -value: 9.1×10^{-25}) and indicates a better separation between the signatures in
939 the current PCAWG analysis.

940

941 ***Correlations of mutational signature activity with age***

942 Prior to evaluating the association between age and the activity of a mutational signatures,
943 all outliers for both age and numbers of mutations attributed to a signature in a cancer type
944 were removed from the data. Outlier was defined as any value outside three standard
945 deviations from the mean value. A robust linear regression model that estimates the slope
946 of the line and whether this slope is significantly different from zero (F-test; p -value<0.05)
947 was performed using the MATLAB function `robustfit`
948 (<https://www.mathworks.com/help/stats/robustfit.html>) with default parameters. The p -
949 values yielded from the F-tests were corrected using the Benjamini-Hochberg procedure for

950 false discovery rate. Results are at <https://www.synapse.org/#!Synapse:syn12030687> and
951 <https://www.synapse.org/#!Synapse:syn20317940>.

952

953 **SigProfiler overview**

954 SigProfiler incorporates two distinct steps for identification of mutational signatures based
955 on the previously described methodology^{3,4,18,36}. The first step, SigProfilerExtraction,
956 encompasses a hierarchical *de novo* extraction of mutational signatures based on somatic
957 mutations and their immediate sequence context, while the second step,
958 SigProfilerAttribution, focuses on accurately estimating the number of somatic mutations
959 associated with each extracted mutational signature in each sample.

960

961 **SigProfilerExtraction**

962 (Note: This phase is termed SigProfiler in the MATLAB code and SigProfilerExtractor in
963 Python). The hierarchical *de novo* extraction approach is an extension of our previous
964 framework for analysis of mutational signatures (Extended Data Figure 8a)^{3,18}. Briefly, for a
965 given set of mutational catalogues, the previously developed algorithm was hierarchically
966 applied to an input matrix $M \in \mathbb{R}_+^{K \times G}$ of non-negative integers with dimension $K \times G$,
967 where K is the number of mutation types and G is the number of samples. This previously
968 described algorithm deciphers a minimal set of mutational signatures that optimally
969 explains the proportion of each mutation type and estimates the contribution of each
970 signature to each sample. The algorithm uses multiple NMFs to identify the matrix of
971 mutational signatures, $P \in \mathbb{R}_+^{K \times N}$, and the matrix of the activities of these signatures, $E \in$
972 $\mathbb{R}_+^{N \times G}$, as previously described³. The unknown number of signatures, N , is determined by
973 human assessment of the stability and accuracy of solutions for a range of values for N , as
974 described³. The identification of M and P is done by minimizing the generalized Kullback-
975 Leibler divergence:

976

$$\min_{P \in \mathbb{R}_+^{(K,N)}, E \in \mathbb{R}_+^{(N,G)}} \sum_{ij} (M_{ij} \log \frac{M_{ij}}{\widehat{M}_{ij}} - M_{ij} + \widehat{M}_{ij}),$$

977

978 where $\widehat{M} \in \mathbb{R}_+^{K \times G}$ is the unnormalized approximation of M , i.e., $\widehat{M} = P \times E$. The
979 framework is applied hierarchically to increase its ability to find mutational signatures
980 generating few mutations or present in few samples. In detail, after application to a
981 matrix M containing the original samples, the accuracy of reconstructing the mutational
982 spectrum of each sample with the extracted mutational signatures is evaluated. Samples
983 that are well-reconstructed are removed, after which the framework is applied to the
984 remaining sub-matrix of M .

985

986 Transcriptional strand bias associated with mutational signatures was assessed by applying
987 SigProfilerExtraction to catalogues of in-transcript mutations that capture strand
988 information (192 mutations classes, <https://www.synapse.org/#!Synapse:syn12026195>).
989 These 192-class signatures were collapsed to strand-invariant 96-class signatures and
990 compared to the signatures extracted from the 96-class data, revealing very high cosine
991 similarities (median 0.90, column F in <https://www.synapse.org/#!Synapse:syn12016215>).
992

993 **SigProfilerAttribution (single sample attribution)**

994 (Note: This phase is termed SigProfilerSingleSample in both the MATLAB and Python code).
995 After signatures are discovered by SigProfilerExtraction, another procedure,
996 SigProfilerAttribution, estimates their contributions to individual samples. For each
997 examined sample, $C \in \mathbb{R}_+^{K \times 1}$, the estimation algorithm involves finding the minimum of the
998 Frobenius norm of a constrained function (see below for constraints) for a set of vectors
999 $S_{i=1..q} \in Q$, where Q is a (not necessarily proper) subset of the set of mutational signatures,
1000 P , ie, $Q \subseteq P$.
1001

$$\min \left\| \vec{C} - \sum_{r=1}^q (\vec{S}_r \times E_r) \right\|_F^2 \quad (1)$$

1002
1003 In equation (1), \vec{C} and each \vec{S}_r are vectors of K nonnegative components reflecting,
1004 respectively, the mutational spectrum of a sample and the r -th reference mutational
1005 signature. All mutational signatures, \vec{S}_r , were identified in the SigProfilerExtraction step.
1006 Each E_r is unknown scalar reflecting the number of mutations contributed by signature \vec{S}_r in
1007 the mutational spectrum \vec{C} . The minimization of equation (1) is always performed under
1008 two additional constraints: (i) $E_r \geq 0$ and (ii) $\|\vec{C}\|_1 \geq E_r$; The constrained minimization of
1009 equation (1) is performed using a nonlinear convex optimization programming solver using
1010 the interior-point algorithm⁷³.
1011

1012 SigProfilerAttribution follows a multistep process, wherein equation (1) is minimized
1013 multiple times with additional constraints (Extended Data Figure 8b).
1014

1015 In the first phase, the subset Q contains all signatures that were found by
1016 SigProfilerExtraction in the same cancer type as the examined sample. Furthermore,
1017 signatures violating biologically meaningful constraints based on transcriptional strand bias
1018 and/or total number of somatic mutations are excluded from the set Q
1019 (<https://www.synapse.org/#!Synapse:syn12177009>). Further, any $\vec{S}_r \times E_r$ for which the
1020 cosine similarity between \hat{C} and \vec{C} is ≤ 0.01 are sequentially removed, where $\hat{C} =$
1021 $\sum_{r=1}^q (\vec{S}_r \times E_r)$. Let T be the final set of signatures attributed to the sample at the end of
1022 the first phase.
1023

1024 In the second phase, equation (1) is minimized by sequentially allowing each signature,
1025 $S_r \in P \setminus Q$, to be added provided that it increases the cosine similarity between \hat{C} and \vec{C}
1026 by >0.05 . During this second phase, several additional biological conditions are enforced: (i)
1027 signatures SBS1 and SBS5 are allowed in all samples, (ii) if one connected SBS signature is
1028 found in a sample than another one is also allowed in the sample (e.g., if SBS17a is found in
1029 a sample then SBS17b is allowed in the sample).
1030

1031

1032 **SignatureAnalyzer overview**

1033 SignatureAnalyzer employs a Bayesian variant of NMF that infers the number of signatures
1034 through the automatic relevance determination technique and delivers highly interpretable

1035 and sparse representations for both signature profiles and attributions that strike a balance
1036 between data fitting and model complexity. Please see references^{6,15,39} for more details.

1037

1038 ***SignatureAnalyzer signature extraction***

1039 In 2,780 PCAWG samples, we applied a two-step signature extraction strategy using 1536
1040 penta-nucleotide contexts for SBSs, 83 ID features, and 78 DBS features. In addition to
1041 separate extraction of SBS, ID, and DBS signatures, we performed a "COMPOSITE" signature
1042 extraction based on all 1697 features (1536 SBS + 78 DBS + 83 ID). For SBSs, the 1536 SBS
1043 COMPOSITE signatures are preferred, and for DBSs and IDs, the separately extracted
1044 signatures are preferred.

1045 In step 1 of the two-step extraction process, global signature extraction was performed for
1046 the low mutation burden samples ($n = 2,624$). These excluded hyper-mutated tumours:
1047 those with putative polymerase epsilon (POLE) defects or mismatch repair defects
1048 (microsatellite instable tumours - MSI), skin tumours (which had intense UV mutagenesis),
1049 and one tumour with temozolomide (TMZ) exposure. Because SignatureAnalyzer's
1050 underlying algorithm performs a stochastic search, different runs can produce different
1051 results. In step 1 we ran SignatureAnalyzer 10 times and selected the solution with the
1052 highest posterior probability. In step 2, additional signatures unique to hyper-mutated
1053 samples were extracted (again selecting the highest posterior probability over 10 runs),
1054 while allowing all signatures found in the low mutation burden-samples to explain some of
1055 the spectra of hyper-mutated samples. This approach was designed to minimize a well-
1056 known "signature bleeding" effect or a bias of hyper- or ultra-mutated samples on the
1057 signature extraction. In addition, this approach provided information about which
1058 signatures are unique to the hyper-mutated samples which is later used when attributing
1059 signatures to samples.

1060

1061 ***SignatureAnalyzer signature attribution***

1062 A similar strategy was used for signature attribution; we performed a separate attribution
1063 process for low- and hyper-mutated samples in all COMPOSITE, SBS, DBS, and ID signatures.
1064 For downstream analyses, we preferred to use the COMPOSITE attributions for SBSs and the
1065 separately calculated attributions for DBSs and IDs. Signature attribution in low-mutation
1066 burden samples was performed separately in each tumour type (e.g., Biliary-AdenoCA,
1067 Bladder-TCC, Bone-Osteosarc, etc.). Attribution was also performed separately in the
1068 combined MSI ($n=39$), POLE ($n=9$), skin melanoma ($n=107$), and TMZ-exposed samples
1069 (<https://www.synapse.org/#!Synapse:syn11738314>). In both groups, signature availability
1070 (i.e., which signatures were active or not) was primarily inferred through the automatic
1071 relevance determination process applied to the activity matrix H only, while fixing the
1072 signature matrix, W . The attribution in low-mutation burden samples was performed using
1073 only signatures found in the step 1 of the signature extraction. Two additional rules were
1074 applied in SBS signature attribution to enforce biological plausibility and minimize a
1075 signature bleeding: (i) allow signature SBS4 (smoking signature) only in lung and head and
1076 neck cases; (ii) allow signature SBS11 (TMZ signature) in a single GBM sample. This was
1077 enforced by introducing a binary, signature-by-sample, signature indicator matrix Z (1 -
1078 allowed and 0 - not allowed), which was multiplied by the H matrix in every multiplication
1079 update of H . No additional rules were applied to ID or DBS signature attributions, except
1080 that signatures found in hyper-mutated samples were not allowed in low-mutation burden
1081 samples.

1082

1083 **Tests on Synthetic Data**

1084 Our goal was to evaluate SignatureAnalyzer (SA) and SigProfiler (SP) on realistic synthetic
1085 data. We operationally defined "realistic" as corresponding to either SA's or SP's analysis of
1086 the PCAWG genome data. SA's reference signature profiles were based on "COMPOSITE"
1087 signatures, consisting of 1536 strand-agnostic single base substitutions (SBSs) in
1088 pentanucleotide context, 78 doublet base substitutions and 83 types of small insertions and
1089 deletions, for a total of 1,697 mutation types. SP's reference analysis was based on strand-
1090 agnostic single base substitutions in the context of one 5' and one 3' base; we term this
1091 "SBS96" data. For each test, we generated two sets of "realistic" data: *SP-realistic*, based on
1092 SP's reference signatures and attributions, and *SA-realistic*, based on SA's reference
1093 signatures and attributions, as well as two other types of data that involved using SA profiles
1094 with SP attributions and vice versa.

1095

1096 **Generating synthetic data – overview.** For tests (i) through (x) below, Synthetic data for
1097 sets of synthetic tumours of a given cancer type, t , were generated based on three
1098 parameters that were in turn based on the observed statistics for each signature, s , in
1099 cancer type t :

1100

1101 π , the proportion of tumours of cancer type t with signature s

1102

1103 μ , the mean of \log_{10} of the number of s mutations across those tumours of type t that have
1104 signature s

1105

1106 σ , the standard deviation of \log_{10} of the numbers of s mutations across those t tumours that
1107 have s

1108

1109 To generate synthetic data,

1110 (i) the proportion of tumours affected by s was drawn from the binomial distribution based
1111 on π ,

1112 (ii) the number of mutations due to s in an affected tumour was drawn from a normal
1113 distribution based on μ and σ .

1114 The code used to generate the synthetic data and summarize SignatureAnalyzer and
1115 SigProfiler results is open-source and freely available as the SynSig package:
1116 <https://github.com/steverozen/SynSig/tree/v0.2.0>.

1117

1117 **Description of each suite of synthetic data sets**

1118

1119 **i. Synthetic pancreatic adenocarcinoma (1,000 spectra).**

1120 <https://doi.org/10.7303/syn18500212.1>

1121

1122 **ii. 2,700 synthetic whole-genome mutational spectra – 300 spectra from each of 9 cancer**
1123 **types.** These spectra consist of 300 synthetic spectra from each of the following cancer
1124 types: bladder transitional cell carcinoma, oesophageal adenocarcinoma, breast
1125 adenocarcinoma, lung squamous cell carcinoma, renal cell carcinoma, ovarian
1126 adenocarcinoma, osteosarcoma, cervical adenocarcinoma, and stomach adenocarcinoma.

1127 <https://doi.org/10.7303/syn18500213.1>

1128

1129 **iii. Mutational spectra generated from combinations of flat, relatively featureless**
1130 **mutational signatures -- version 1**, 1000 synthetic tumours comprised of 500 synthetic
1131 Kidney-RCCs (high prevalence and mutation load from SBS5 and SBS40 signatures) and 500
1132 synthetic ovarian adenocarcinomas (high prevalence of and mutation load from SBS3). This
1133 data set embodies tumours with high prevalence of the main flat signatures, SBS3, SBS5,
1134 and SBS40, in a realistic context.

1135 <https://doi.org/10.7303/syn18500214.1>

1136

1137 **iv. Mutational spectra generated from combinations of flat, relatively featureless**
1138 **mutational signatures -- version 2**, 1000 synthetic spectra all constructed entirely from
1139 SBS3, SBS5, and SBS40, using mutational loads modelled on kidney-RCC (SBS5 and SBS40)
1140 and ovarian adenocarcinoma (SBS3). Most synthetic spectra have contributions from all
1141 three signatures.

1142 <https://doi.org/10.7303/syn18500215.1>

1143

1144 **v. Mutational spectra generated from signatures with overlapping and potentially**
1145 **interfering profiles - version 1**. 500 synthetic bladder transitional cell carcinomas (high in
1146 SBS2 and SBS13) and 500 synthetic skin melanomas (high in SBS7a,b,c,d). The potential
1147 interference is between SBS2 (mainly C > T) and SBS7a,b (mainly C > T).

1148 <https://doi.org/10.7303/syn18500217.1>

1149

1150 **vi. Mutational spectra generated from signatures with overlapping and potentially**
1151 **interfering profiles - version 2**. 1000 synthetic tumours composed from SBS2 and
1152 SBS7a,b. Mutational load distributions were drawn from bladder transitional cell carcinoma
1153 (SBS2) and skin melanoma (SBS7a,b). Most spectra contain both signatures. The potential
1154 interference is between SBS2 (mainly C > T) and SBS7a,b (mainly C > T).

1155 <https://doi.org/10.7303/syn18500216.1>

1156

1157 **vii. Mutational spectra generated from combinations of signatures conferring high and**
1158 **low mutation burdens**. Based on 500 synthetic non-hypermutated tumours (parameters for
1159 SBS1 and SBS5 estimated from colorectal and uterine adenocarcinomas) and 500 hyper-
1160 mutated tumours (parameters for SBS26 and SBS44 estimated from hypermutated
1161 colorectal and uterine adenocarcinomas). High and low mutation burden tumours are
1162 segregated for SignatureAnalyzer (which analyses low mutation burden tumours first, then
1163 high-burden tumours). SigProfiler analyses all tumours together.

1164 <https://doi.org/10.7303/syn18500218.1>

1165 <https://doi.org/10.7303/syn18500219.1>

1166 <https://doi.org/10.7303/syn18500216.1>

1167

1168 **viii. A set of 30 random 96-feature mutational signature profiles and a set of 30 random**
1169 **1697-feature signature profiles (mimicking COMPOSITE signatures, which have 1697**
1170 **mutation types)**. Each of these are used in two types of exposures, one with more (mean
1171 ~15.6) signatures per tumour and one with fewer (mean ~4) signatures per tumour.

1172 <https://doi.org/10.7303/syn18500221.1>

1173

1174 **ix. 2,700 whole-exome mutational spectra consisting of 300 synthetic spectra from each of**
1175 **9 different cancer types.** This test data set was generated from *test ii* by reducing the
1176 number of mutations of each type by 0.013 (approximately ratio of mutation counts
1177 between whole exome and whole genome mutational spectra).

1178 <https://doi.org/10.7303/syn18909829.4>

1179

1180 *Summary of findings:* Both SA and SP extracted substantially fewer signatures in this test
1181 than in *test ii*. In particular:

1182

1183 **SA:** SA extracted only 18 signatures from the SA-realistic whole-exome data in this suite,
1184 compared to the 40 signatures it extracted from the corresponding whole-genome synthetic
1185 data in *test ii* and compared to the 39 ground-truth signatures in the synthetic spectra. The
1186 average cosine similarity between ground-truth and extracted signatures for the synthetic
1187 exome data was 0.863, compared to 0.968 for the signatures extracted from the whole-
1188 genome spectra in *test ii*.

1189

1190 **SP:** SP extracted only 8 signatures from the SP-realistic whole-exome data in this suite,
1191 compared to the 19 it extracted from the whole-genome data in *test ii* and the 21 ground-
1192 truth signatures in the synthetic spectra. The average cosine similarity between ground-
1193 truth and extracted signatures for the synthetic exome data was 0.825, compared to 0.965
1194 for the signatures extracted from the whole-genome spectra in *test ii*.

1195

1196 **x. 1,350 synthetic whole-genome mutational spectra: 150 spectra from each of 9 cancer**
1197 **types.** This test data set consisted of every other tumour from *test ii*.

1198

1199 *Summary of findings:* On the SA-realistic synthetic data, SA extracted fewer signatures in
1200 this data set than in *test ii*, and in fact the number of signatures extracted was closer to the
1201 ground truth and the cosine similarities were there higher. SA over-split in the
1202 corresponding set of 2,700 tumours, and we speculate that SA's tendency to over-split
1203 signatures is partly dependent on the number of input spectra, with the result that
1204 extraction on 1,350 led to less over-splitting. SP extracted fewer signatures on this data set
1205 than on *test ii*. In particular:

1206

1207 **SA:** SA extracted 38 signatures from the SA-realistic data in this suite, compared to the 40
1208 signatures it extracted from the 2,700 whole-genome spectra in *test ii* and compared to the
1209 39 ground-truth signatures. The average cosine similarity between ground-truth and
1210 extracted signatures for 1,350 genomes was 0.979 compared to 0.968 for the signatures
1211 extracted from the 2,700 whole-genome spectra in *test ii*.

1212

1213 **SP:** SP extracted 16 signatures from the SP-realistic data in this suite, compared to the 19
1214 signatures it extracted from the 2,700 whole-genome spectra in *test ii* and the 21 ground-
1215 truth signatures. The average cosine similarity between ground-truth and extracted
1216 signatures for the 1,350 spectra was 0.939 compared to 0.965 for the signatures extracted
1217 from the 2,700 spectra in *test ii*.

1218

1219 **xi. Extraction of signatures from exome subsets of PCAWG mutational spectra.** Our
1220 objective was to further test whether availability of mutations from whole-genome

1221 mutational spectra, as opposed to whole-exome spectra, enabled us to extract larger
1222 numbers of more accurate mutational signature profiles. In this test, we extracted
1223 signatures from mutational spectra that were based on only the exome regions of the actual
1224 PCAWG tumours (rather than on the purely synthetic data in *test ix*). The input data and
1225 extraction results are at <https://doi.org/10.7303/syn18818766>. We next summarize our
1226 findings for each of the SBS, DBS, and ID mutational signatures.

1227
1228 *xi-1 SBS signatures.* SignatureAnalyzer on COMPOSITE mutational classes (1536 SBS in
1229 pentanucleotide context plus DBS and ID) extracted 12 mutational signature profiles from
1230 the whole-exome data, none of which strongly resembled any of the 58 signatures it
1231 extracted from the whole-genome data. However, some signatures were unions or splits of
1232 the signatures extracted from the whole genome data. For example, WI was a union of the
1233 APOBEC signatures BI_COMPOSITE_SBS2_P and BI_COMPOSITE_SBS13_P. More broadly,
1234 somewhat recognizable SBS portions of the signatures were combined with the DBS and ID
1235 portions of the signatures in difficult-to-interpret combinations. We believe that SBS
1236 mutation counts were too low when spread across 1536 mutational classes to support
1237 robust mutational signature extraction.

1238
1239 SigProfiler on 96 SBS mutational classes extracted 17 mutational signature profiles from the
1240 exome data, compared to 48 that it extracted from the whole-genome data. The median
1241 cosine similarity of the exome-extracted signature profiles to the mutational signature
1242 profiles extracted from the whole genome data was 0.94. An outlier was SBS-E-2, which was
1243 a union of SBS2 and SBS13 (which tend to co-occur).

1244
1245 *xi-2 DBS signatures.* SignatureAnalyzer extracted 2 DBS signatures from the whole-exome
1246 data, compared to 15 DBS signatures that it extracted from the full whole genome data. One
1247 exome-extracted signature was essentially identical to BI_DBS1 (consisting almost entirely
1248 of CC > TT mutations), and one somewhat similar to BI_DBS2 (mostly CC > AA) but with
1249 many other mutational classes in addition.

1250
1251 SigProfiler extracted 3 DBS signatures from the whole-exome data, compared to the 11 DBS
1252 signatures that it extracted from the whole genome data. The exome-extracted signatures
1253 were good approximations of DBS1, DBS2, and DBS10 (cosine similarities 1, 0.93, and 0.98).

1254
1255 *xi-3 ID signatures.* SignatureAnalyzer extracted 4 ID signatures from the whole-exome data,
1256 compared to 29 ID signatures extracted from the whole-genome data. It extracted close
1257 approximations of BI_ID1_P and BI_ID2_P with cosine similarities 0.97 and 0.94. These are
1258 insertions (signature W.3) and deletions (signature W.1) of T:A in poly T:A.
1259 SignatureAnalyzer extracted 2 additional signatures. One of these (W.4) was a version of
1260 BI_ID4_P with several mutational classes absent. The other (W.2) appeared to be a union of
1261 many of the remaining ID signatures.

1262
1263 SigProfiler extracted 6 ID signatures from the whole-exome data, compared to the 17 ID
1264 signatures that it extracted from the whole genome data. Signatures ID-E-1, ID-E-2, ID-E-3,
1265 and ID-E-4 were good approximations of ID1, ID2, ID3, and ID4, respectively. An additional
1266 signature, ID-E-5, was approximately a union of ID6 and ID8. The remaining signature, ID-E-6
1267 was a partial version (deletions in C homopolymers only) of ID7.

1268

1269 **Detailed Summary of Results (including links to input synthetic data sets and the signature**
1270 **profiles extracted)**; <https://doi.org/10.7303/syn18497223> provides a table with the number
1271 of signatures extracted by SigProfiler and SignatureAnalyzer for each synthetic data set and
1272 the cosine similarities to the input ground-truth signatures. See above for overall
1273 interpretation of the results.

1274

1275 **Data Availability**

1276 Data are available at <https://www.synapse.org/#!Synapse:syn11726601/wiki/513478>. All
1277 figures and extended data figures have associated raw data.

1278

1279 **Code Availability**

1280 SigProfiler is available both as a MATLAB framework and as a Python package. In both cases,
1281 SigProfiler is fully functional, free, and open-source tool distributed under the permissive 2-
1282 Clause BSD License. SigProfiler in MATLAB can be downloaded from:

1283 <https://www.mathworks.com/matlabcentral/fileexchange/38724-sigprofiler>

1284 SigProfiler in Python can be downloaded from:

1285 <https://github.com/AlexandrovLab/SigProfilerExtractor>. SignatureAnalyzer code is available at

1286 <https://www.synapse.org/#!Synapse:syn11801492>. The code used to generate the synthetic data

1287 and summarize SignatureAnalyzer and SigProfiler results is open-source and freely available as the

1288 SynSig package: <https://github.com/steverozen/SynSig/tree/v0.2.0> under the GPL3 license.

1289

1290

References

1291

- 1292 1 Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. *Nature* **458**, 719-
1293 724, doi:10.1038/nature07943 (2009).
- 1294 2 Alexandrov, L. B. & Stratton, M. R. Mutational signatures: the patterns of somatic
1295 mutations hidden in cancer genomes. *Current opinion in genetics & development* **24**,
1296 52-60 (2014).
- 1297 3 Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M. R.
1298 Deciphering signatures of mutational processes operative in human cancer. *Cell Rep*
1299 **3**, 246-259 (2013).
- 1300 4 Alexandrov, L. B. *et al.* Clock-like mutational processes in human somatic cells. *Nat*
1301 *Genet* **47**, 1402-1407, doi:10.1038/ng.3441 (2015).
- 1302 5 Morganella, S. *et al.* The topography of mutational processes in breast cancer
1303 genomes. *Nat Commun* **7**, 11383, doi:10.1038/ncomms11383 (2016).
- 1304 6 Kim, J. *et al.* Somatic ERCC2 mutations are associated with a distinct genomic
1305 signature in urothelial tumors. *Nat Genet* **48**, 600-606, doi:10.1038/ng.3557 (2016).
- 1306 7 Fischer, A., Illingworth, C. J., Campbell, P. J. & Mustonen, V. EMu: probabilistic
1307 inference of mutational processes and their localization in the cancer genome.
1308 *Genome Biol* **14**, R39, doi:10.1186/gb-2013-14-4-r39 (2013).
- 1309 8 Roberts, N. hdp (hierarchical Dirichlet process) R package.
1310 <https://github.com/nicolaroberts/hdp> (2015).
- 1311 9 Gehring, J. S., Fischer, B., Lawrence, M. & Huber, W. SomaticSignatures: inferring
1312 mutational signatures from single-nucleotide variants. *Bioinformatics* **31**, 3673-3675
1313 (2015).
- 1314 10 Shiraishi, Y., Tremmel, G., Miyano, S. & Stephens, M. A simple model-based
1315 approach to inferring and visualizing cancer mutation signatures. *PLoS genetics* **11**,
1316 e1005657 (2015).
- 1317 11 Rosales, R. A., Drummond, R. D., Valieris, R., Dias-Neto, E. & da Silva, I. T. signeR: an
1318 empirical Bayesian approach to mutational signature discovery. *Bioinformatics* **33**, 8-
1319 16, doi:10.1093/bioinformatics/btw572 (2017).
- 1320 12 Ardin, M. *et al.* MutSpec: a Galaxy toolbox for streamlined analyses of somatic
1321 mutation spectra in human and mouse cancer genomes. *BMC bioinformatics* **17**, 170
1322 (2016).
- 1323 13 Funnell, T. *et al.* Integrated single-nucleotide and structural variation signatures of
1324 DNA-repair deficient human cancers. *bioRxiv*, doi:10.1101/267500 (2018).
- 1325 14 Blokzijl, F., Janssen, R., van Boxtel, R. & Cuppen, E. MutationalPatterns:
1326 comprehensive genome-wide analysis of mutational processes. *Genome medicine*
1327 **10**, 33 (2018).
- 1328 15 Kasar, S. *et al.* Whole-genome sequencing reveals activation-induced cytidine
1329 deaminase signatures during indolent chronic lymphocytic leukaemia evolution. *Nat*
1330 *Commun* **6**, 8866, doi:10.1038/ncomms9866 (2015).
- 1331 16 Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature*
1332 **500**, 415-421 (2013).
- 1333 17 Nik-Zainal, S. *et al.* Mutational processes molding the genomes of 21 breast cancers.
1334 *Cell* **149**, 979-993 (2012).
- 1335 18 Nik-Zainal, S. *et al.* Landscape of somatic mutations in 560 breast cancer whole-
1336 genome sequences. *Nature* **534**, 47-54, doi:10.1038/nature17676 (2016).

- 1337 19 Petljak, M. & Alexandrov, L. B. Understanding mutagenesis through delineation of
1338 mutational signatures in human cancer. *Carcinogenesis* **37**, 531-540,
1339 doi:10.1093/carcin/bgw055 (2016).
- 1340 20 Schulze, K. *et al.* Exome sequencing of hepatocellular carcinomas identifies new
1341 mutational signatures and potential therapeutic targets. *Nat Genet* **47**, 505-511,
1342 doi:10.1038/ng.3252 (2015).
- 1343 21 Poon, S. L. *et al.* Mutation signatures implicate aristolochic acid in bladder cancer
1344 development. *Genome Med* **7**, 38, doi:10.1186/s13073-015-0161-3 (2015).
- 1345 22 Poon, S. L. *et al.* Genome-wide mutational signatures of aristolochic acid and its
1346 application as a screening tool. *Sci Transl Med* **5**, 197ra101,
1347 doi:10.1126/scitranslmed.3006086 (2013).
- 1348 23 Alexandrov, L. B. Understanding the origins of human cancer. *Science* **350**, 1175,
1349 doi:10.1126/science.aad7363 (2015).
- 1350 24 Hayward, N. K. *et al.* Whole-genome landscapes of major melanoma subtypes.
1351 *Nature* **545**, 175-180, doi:10.1038/nature22071 (2017).
- 1352 25 Polak, P. *et al.* A mutational signature reveals alterations underlying deficient
1353 homologous recombination repair in breast cancer. *Nat Genet* **49**, 1476-1486,
1354 doi:10.1038/ng.3934 (2017).
- 1355 26 Merlevede, J. *et al.* Mutation allele burden remains unchanged in chronic
1356 myelomonocytic leukaemia responding to hypomethylating agents. *Nat Commun* **7**,
1357 10767, doi:10.1038/ncomms10767 (2016).
- 1358 27 Mimaki, S. *et al.* Hypermutation and unique mutational signatures of occupational
1359 cholangiocarcinoma in printing workers exposed to haloalkanes. *Carcinogenesis* **37**,
1360 817-826, doi:10.1093/carcin/bgw066 (2016).
- 1361 28 Drier, Y. *et al.* Somatic rearrangements across cancer reveal classes of samples with
1362 distinct patterns of DNA breakage and rearrangement-induced hypermutability.
1363 *Genome Res* **23**, 228-235, doi:10.1101/gr.141382.112 (2013).
- 1364 29 Stephens, P. J. *et al.* Massive genomic rearrangement acquired in a single
1365 catastrophic event during cancer development. *Cell* **144**, 27-40,
1366 doi:10.1016/j.cell.2010.11.055 (2011).
- 1367 30 Li, Y. *et al.* Patterns of structural variation in human cancer. *bioRxiv*,
1368 doi:10.1101/181339 (2017).
- 1369 31 Meier, B. *et al.* C. elegans whole-genome sequencing reveals mutational signatures
1370 related to carcinogens and DNA repair deficiency. *Genome research* **24**, 1624-1636
1371 (2014).
- 1372 32 Supek, F. & Lehner, B. Clustered Mutation Signatures Reveal that Error-Prone DNA
1373 Repair Targets Mutations to Active Genes. *Cell* **170**, 534-547 e523,
1374 doi:10.1016/j.cell.2017.07.003 (2017).
- 1375 33 Chen, J. M., Férec, C. & Cooper, D. N. Patterns and mutational signatures of tandem
1376 base substitutions causing human inherited disease. *Human mutation* **34**, 1119-1130
1377 (2013).
- 1378 34 Campbell, P. J. *et al.* Pan-cancer analysis of whole genomes. *Preprint on bioRxiv*,
1379 <https://doi.org/10.1101/162784> (2017).
- 1380 35 Yung, C. K. *et al.* Large-scale uniform analysis of cancer whole genomes in multiple
1381 computing environments. *BioRxiv*, 161638 (2017).
- 1382 36 WTSI Mutational Signature Framework,
1383 <http://www.mathworks.com/matlabcentral/fileexchange/38724> (2013).

- 1384 37 Wellcome Trust Sanger Institute. *COSMIC, Catalog of Somatic Mutations in Cancer -*
1385 *Signatures of Mutational Processes in Human Cancer,*
1386 <http://cancer.sanger.ac.uk/cosmic/signatures> (2017).
- 1387 38 Forbes, S. A. *et al.* The Catalogue of Somatic Mutations in Cancer (COSMIC). *Curr*
1388 *Protoc Hum Genet* **Chapter 10**, Unit 10 11, doi:10.1002/0471142905.hg1011s57
1389 (2008).
- 1390 39 Tan, V. Y. & Févotte, C. Automatic relevance determination in nonnegative matrix
1391 factorization with the/spl beta/-divergence. *IEEE Transactions on Pattern Analysis*
1392 *and Machine Intelligence* **35**, 1592-1605 (2013).
- 1393 40 Haradhvala, N. J. *et al.* Distinct mutational signatures characterize concurrent loss of
1394 polymerase proofreading and mismatch repair. *Nature Communications* **9**, 1746,
1395 doi:10.1038/s41467-018-04002-4 (2018).
- 1396 41 Cichocki, A., Zdunek, R., Phan, A. H. & Amari, S.-i. *Nonnegative matrix and tensor*
1397 *factorizations: applications to exploratory multi-way data analysis and blind source*
1398 *separation.* (John Wiley & Sons, 2009).
- 1399 42 Devarajan, K. Nonnegative matrix factorization: an analytical and interpretive tool in
1400 computational biology. *PLoS computational biology* **4**, e1000029 (2008).
- 1401 43 Blei, D., Carin, L. & Dunson, D. Probabilistic Topic Models: A focus on graphical model
1402 design and applications to document and image analysis. *IEEE signal processing*
1403 *magazine* **27**, 55 (2010).
- 1404 44 Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new
1405 cancer-associated genes. *Nature* **499**, 214-218, doi:10.1038/nature12213 (2013).
- 1406 45 Boot, A. *et al.* In-depth characterization of the cisplatin mutational signature in
1407 human cell lines and in esophageal and liver tumors. *Genome Res* **28**, 654-665
1408 doi:10.1101/gr.230219.117 (2018).
- 1409 46 Viel, A. *et al.* A Specific Mutational Signature Associated with DNA 8-Oxoguanine
1410 Persistence in MUTYH-defective Colorectal Cancer. *EBioMedicine* **20**, 39-49,
1411 doi:10.1016/j.ebiom.2017.04.022 (2017).
- 1412 47 Pilati, C. *et al.* Mutational signature analysis identifies MUTYH deficiency in colorectal
1413 cancers and adrenocortical carcinomas. *J Pathol* **242**, 10-15, doi:10.1002/path.4880
1414 (2017).
- 1415 48 Drost, J. *et al.* Use of CRISPR-modified human stem cell organoids to study the origin
1416 of mutational signatures in cancer. *Science* **358**, 234-238 (2017).
- 1417 49 Saini, N. *et al.* The Impact of Environmental and Endogenous Damage on Somatic
1418 Mutation Load in Human Skin Fibroblasts. *PLoS Genetics* **12**, e1006385,
1419 doi:10.1371/journal.pgen.1006385 (2016).
- 1420 50 Haradhvala, N. J. *et al.* Mutational Strand Asymmetries in Cancer Genomes Reveal
1421 Mechanisms of DNA Damage and Repair. *Cell* **164**, 538-549,
1422 doi:10.1016/j.cell.2015.12.050 (2016).
- 1423 51 Chan, K. *et al.* An APOBEC3A hypermutation signature is distinguishable from the
1424 signature of background mutagenesis by APOBEC3B in human cancers. *Nat Genet* **47**,
1425 1067-1072, doi:10.1038/ng.3378 (2015).
- 1426 52 Roberts, S. A. *et al.* Clustered mutations in yeast and in human cancers can arise
1427 from damaged long single-strand DNA regions. *Molecular cell* **46**, 424-435 (2012).
- 1428 53 Kasar, S. & Brown, J. R. Mutational landscape and underlying mutational processes in
1429 chronic lymphocytic leukemia. *Mol Cell Oncol* **3**, e1157667,
1430 doi:10.1080/23723556.2016.1157667 (2016).

- 1431 54 Brash, D. E. UV Signature Mutations. *Photochemistry and Photobiology* **91**, 15-26,
1432 doi:doi:10.1111/php.12377 (2015).
- 1433 55 Hill, K. A., Wang, J., Farwell, K. D. & Sommer, S. S. Spontaneous tandem-base
1434 mutations (TBM) show dramatic tissue, age, pattern and spectrum specificity. *Mutat*
1435 *Res* **534**, 173-186 (2003).
- 1436 56 Matsuda, T., Kawanishi, M., Yagi, T., Matsui, S. & Takebe, H. Specific tandem GG to
1437 TT base substitutions induced by acetaldehyde are due to intra-strand crosslinks
1438 between adjacent guanine bases. *Nucleic Acids Res* **26**, 1769-1774 (1998).
- 1439 57 Garaycochea, J. I. *et al.* Alcohol and endogenous aldehydes damage chromosomes
1440 and mutate stem cells. *Nature* (2018).
- 1441 58 Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat Genet*
1442 **48**, 126-133, doi:10.1038/ng.3469 (2016).
- 1443 59 Pfeifer, G. P. Formation and processing of UV photoproducts: effects of DNA
1444 sequence and chromatin environment. *Photochemistry and photobiology* **65**, 270-
1445 283 (1997).
- 1446 60 Ceccaldi, R., Rondinelli, B. & D'Andrea, A. D. Repair pathway choices and
1447 consequences at the double-strand break. *Trends in cell biology* **26**, 52-64 (2016).
- 1448 61 Blokzijl, F. *et al.* Tissue-specific mutation accumulation in human adult stem cells
1449 during life. *Nature* **538**, 260-264, doi:10.1038/nature19768 (2016).
- 1450 62 Behjati, S. *et al.* Mutational signatures of ionizing radiation in second malignancies.
1451 *Nature communications* **7**, 12605 (2016).
- 1452 63 Huang, M. N. *et al.* Genome-Scale Mutational Signatures of Aflatoxin In Cells, Mice
1453 And Human Tumors. *Genome Research* **27**, 1475-1486, doi:10.1101/gr.220038.116
1454 (2017).
- 1455 64 Nik-Zainal, S. *et al.* The genome as a record of environmental exposure. *Mutagenesis*
1456 **30**, 763-770, doi:10.1093/mutage/gev073 (2015).
- 1457 65 Olivier, M. *et al.* Modelling mutational landscapes of human cancers in vitro. *Sci Rep*
1458 **4**, 4482, doi:10.1038/srep04482 (2014).
- 1459 66 Szikriszt, B. *et al.* A comprehensive survey of the mutagenic impact of common
1460 cancer cytotoxics. *Genome biology* **17**, 99 (2016).
- 1461 67 Zhivagui, M. *et al.* Experimental analysis of exome-scale mutational signature of
1462 glycidamide, the reactive metabolite of acrylamide. *bioRxiv*, 254664 (2018).
- 1463 68 Záborszky, J. *et al.* Loss of BRCA1 or BRCA2 markedly increases the rate of base
1464 substitution mutagenesis and has distinct effects on genomic deletions. *Oncogene*
1465 **36**, 746 (2017).
- 1466 69 Zou, X. *et al.* Validating the concept of mutational signatures with isogenic cell
1467 models. *Nature communications* **9**, 1744 (2018).
- 1468 70 Roberts, N. D. *Patterns of somatic genome rearrangement in human cancer*,
1469 University of Cambridge, (2018).
- 1470 71 Rosales, R. A., Drummond, R. D., Valieris, R., Dias-Neto, E. & da Silva, I. T. signeR: an
1471 empirical Bayesian approach to mutational signature discovery. *Bioinformatics* **33**, 8-
1472 16 (2016).
- 1473 72 Gori, K. & Baez-Ortega, A. sigfit: flexible Bayesian inference of mutational signatures.
1474 *bioRxiv*, 372896 (2018).
- 1475 73 Byrd, R. H., Hribar, M. E. & Nocedal, J. An interior point algorithm for large-scale
1476 nonlinear programming. *SIAM Journal on Optimization* **9**, 877-900 (1999).

1477

1478 **Extended Data Figure and Table Legends**

1479

1480 **Extended Data Figure 1. Histogram of number of signatures attributed in each of 2,780**
1481 **PCAWG samples by SigProfiler and SignatureAnalyzer.** Hypermutated tumours and
1482 melanomas (156) are listed at <https://www.synapse.org/#!Synapse:syn11738314>.

1483

1484 **Extended Data Figure 2. Comparisons between SigProfiler and SignatureAnalyzer results.**
1485 Comparison of the attributions for corresponding SigProfiler (**a**) and SignatureAnalyzer (**b**)
1486 signatures. Each of the SBS signatures extracted by SigProfiler and SignatureAnalyzer was
1487 paired with the signature of highest cosine similarity in the extraction by the other method
1488 (if one with >0.85 cosine similarity exists). The first column of the plot corresponds to the
1489 fraction of mutations assigned by one method (summed across samples and mutation
1490 types) that were also assigned by the other method. The remaining mutations were then re-
1491 distributed to the other signatures in the extraction, weighted by their relative probabilities
1492 of having been generated by each signature, and the resulting fraction of mutations is
1493 plotted. Signatures on the x-axis are only shown if they contribute at least 0.1 fraction of
1494 mutations to at least one signature on the y-axis. Cosine similarities between SigProfiler and
1495 SignatureAnalyzer DBS (**c**) and ID (**d**) signatures. Brown nodes represent SigProfiler
1496 signatures; green nodes represent SignatureAnalyzer signatures. Matches with cosine
1497 similarities > 0.8 are show as edges, with the width of the edge indicate the strength of the
1498 similarity. The locations of the nodes have no significance. Signatures with no matches of $>$
1499 0.8 cosine similarity are show below. Note that SigProfiler ID15 and ID17 were extracted
1500 from data that were not analysed by SignatureAnalyzer. Suffixes 'P' and 'S' on
1501 SignatureAnalyzer signature names indicate (1) signatures extracted from non-
1502 hypermutated, non-melanoma tumours and (2) hypermutated and melanoma tumours,
1503 respectively.

1504

1505 **Extended Data Figure 3. SignatureAnalyzer reference signatures.** See legend of main text
1506 Figure 2.

1507

1508 **Extended Data Figure 4. The number of SBS mutations attributed to each mutational**
1509 **signature for each cancer type over the 2,780 PCAWG tumours by SignatureAnalyzer.** See
1510 main text Figure 3 for explanation.

1511

1512 **Extended Data Figure 5. The number of SBS mutations attributed to each mutational**
1513 **signature to each cancer type over the complete set of 23,829 cancer samples analysed by**
1514 **SigProfiler.** See main text Figure 3 for explanation.

1515

1516 **Extended Data Figure 6. Associations of between SBS, DBS, and ID signature activities for**
1517 **SigProfiler (a) and SignatureAnalyzer (b).** Each node represents an SBS (light green), DBS
1518 (dark green) or ID (black) signature. Any two signatures with sample attributions that
1519 significantly correlated with $R^2 > 0.3$ (SigProfiler) or > 0.5 (SignatureAnalyzer) are connected
1520 by edges. Edge widths are proportional to the strength of the correlation. Signatures with
1521 no significant correlation to any other signature above the relevant threshold are not
1522 shown. Signature locations are fit for display purposes only and do not indicate similarity.

1523

1524 **Extended Data Figure 7. Mutational signatures extracted from the composite feature set**
1525 **consisting of SBSs in pentanucleotide context, DBSs, and IDs.** For each of the four
1526 composite mutational signatures shown, the top panel is the SBS signature collapsed to 96
1527 SBS classes, the middle panel is the co-extracted DBS signature, and the lower panel is the
1528 co-extracted ID signature. Note the similarities between the DBS portion of Composite 4 and
1529 DBS2, between the ID portion of Composite 4 and ID3, and other similarities noted in the
1530 figure.

1531
1532 **Extended Data Figure 8. SigProfiler signature extraction (a) and attribution (b).** See
1533 Methods for description.

1534
1535 **Extended Data Table 1. The number of DBSs is proportional to the number of SBSs with**
1536 **the exception of a few cancer types (ColoRect-AdenoCA, Lung-AdenoCA, Lung-SCC, Skin-**
1537 **Melanoma) analysed by the following linear regression (computed by an R function call):**
1538
$$\text{glm}(\text{DBS.counts} \sim \text{SBS.counts} + \text{Cancer.Types}).$$

1539

1540 **Extended Data Table 2. Numbers of insertion/deletion mutations due to ID1, ID2, and all**
1541 **other ID signatures in hypermutators and non-hypermutators.**

bioRxiv preprint doi: <https://doi.org/10.1101/322859>; this version posted July 3, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

Mutations per Megabase

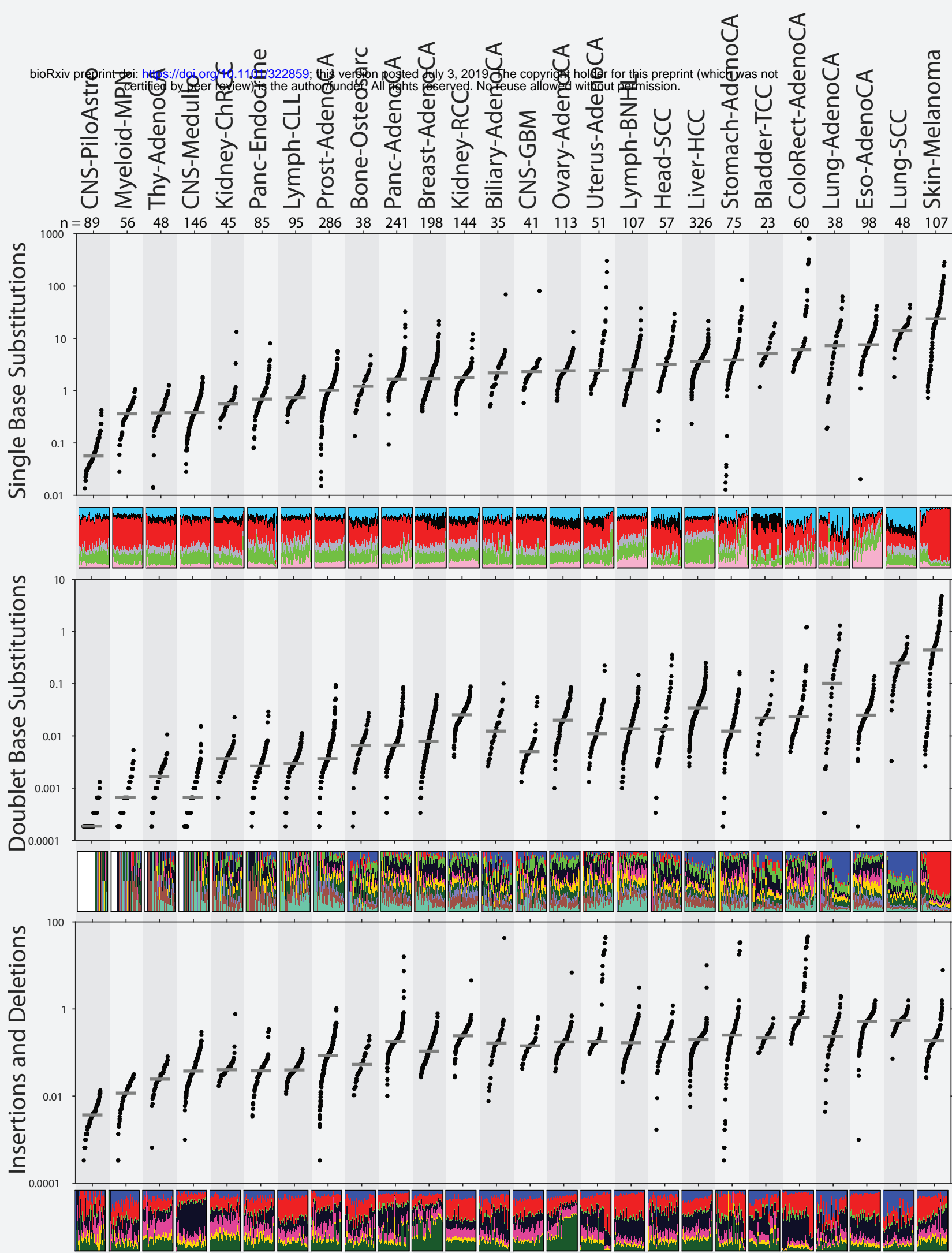
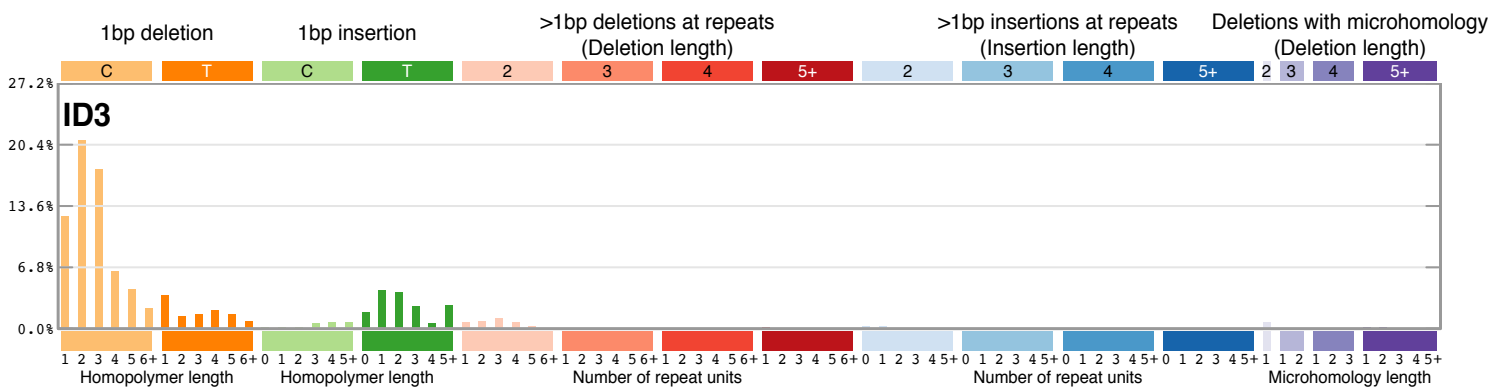
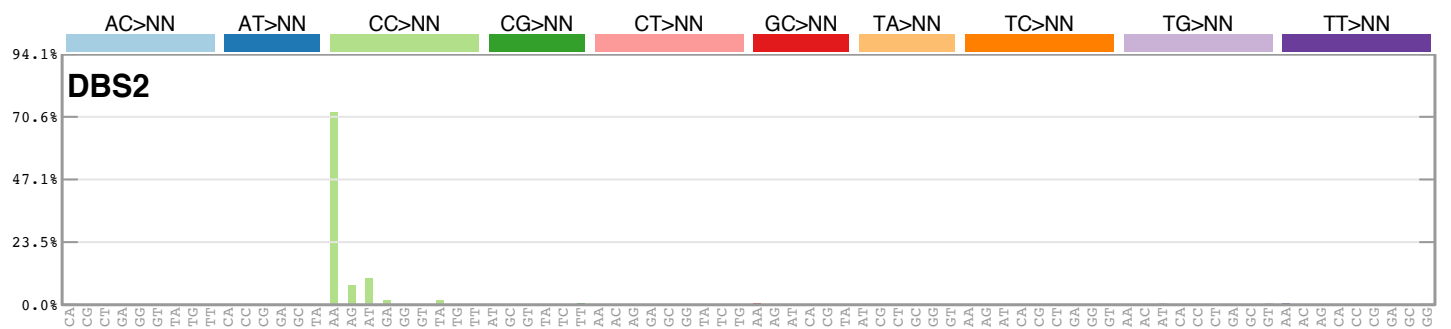
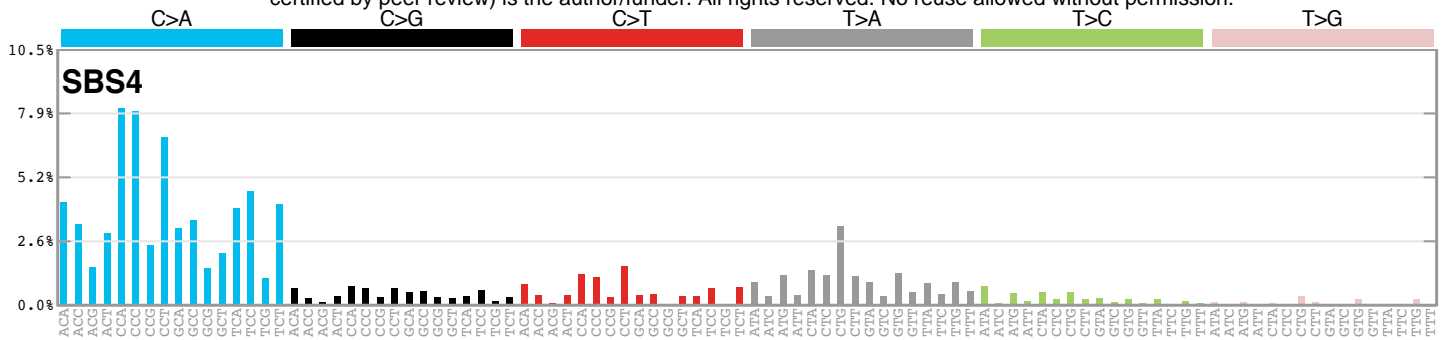
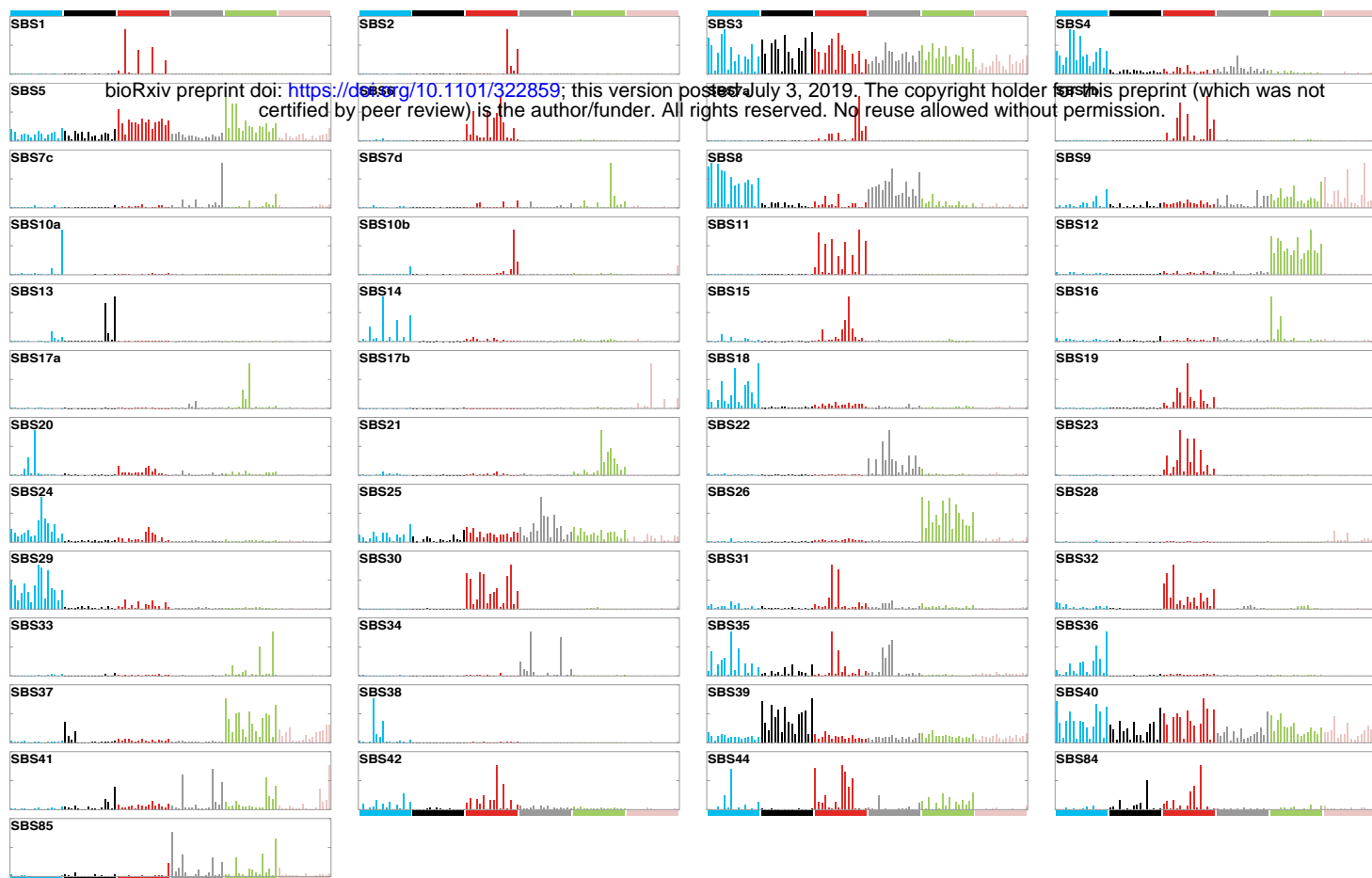


Fig 2 part 1

bioRxiv preprint doi: <https://doi.org/10.1101/322859>; this version posted July 3, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



Single Base Substitution



Doublet Base Substitution



Insertion and Deletion

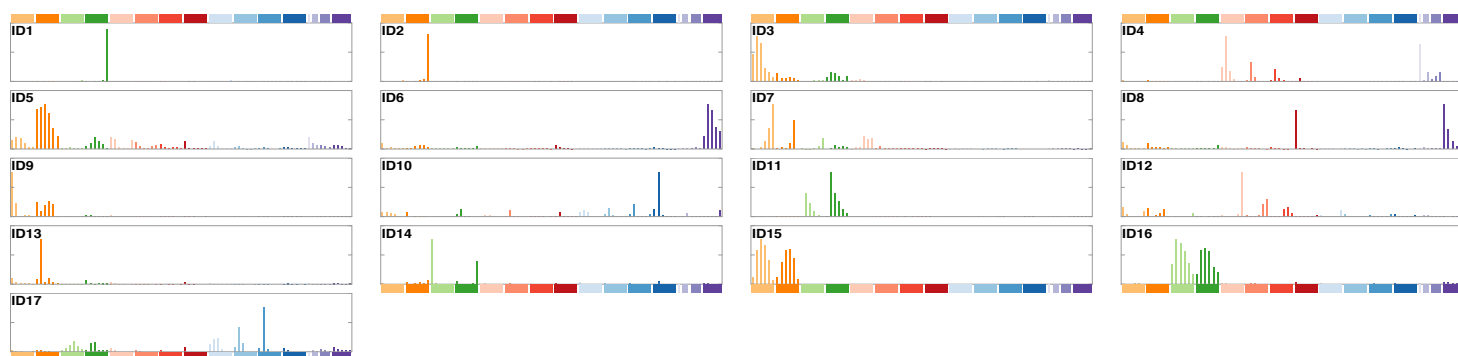
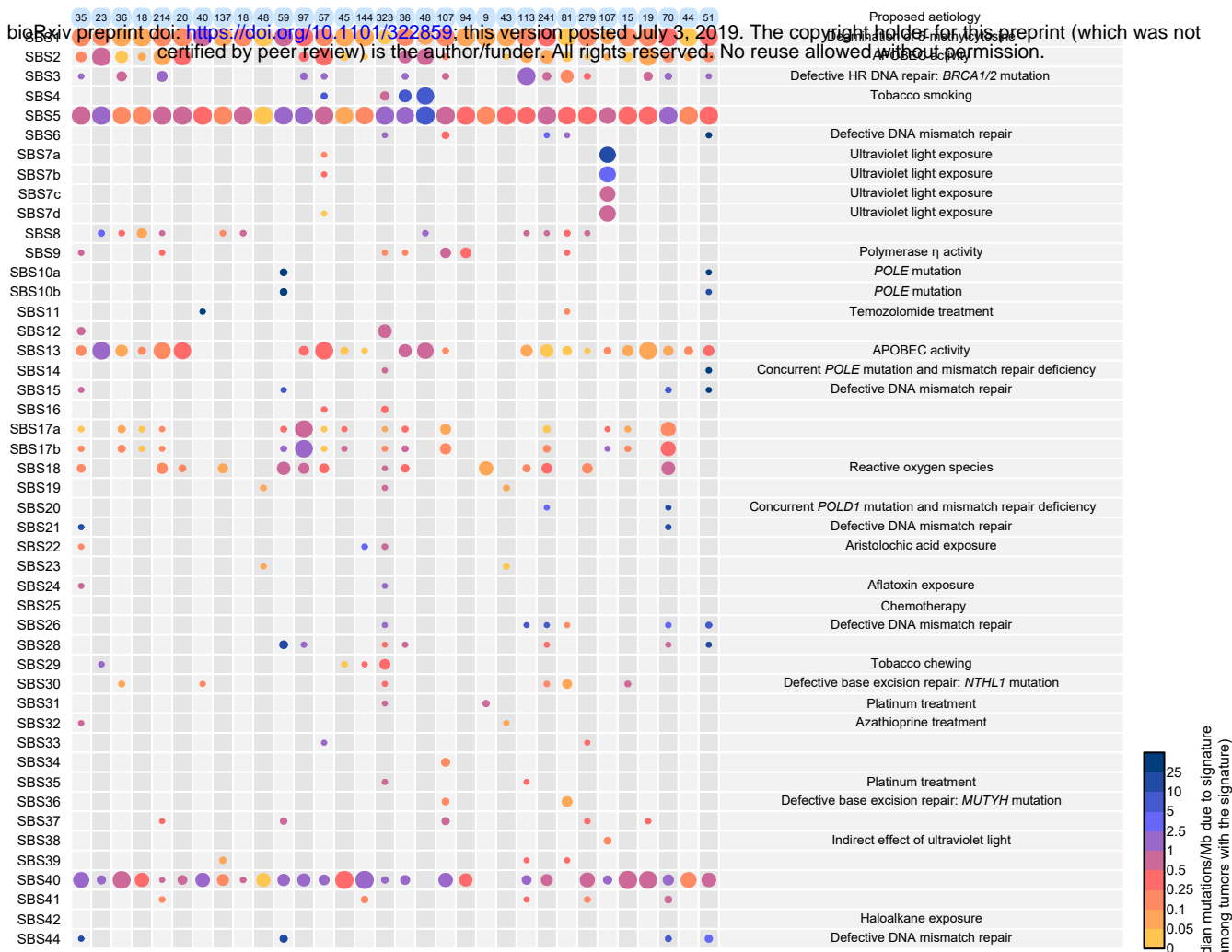
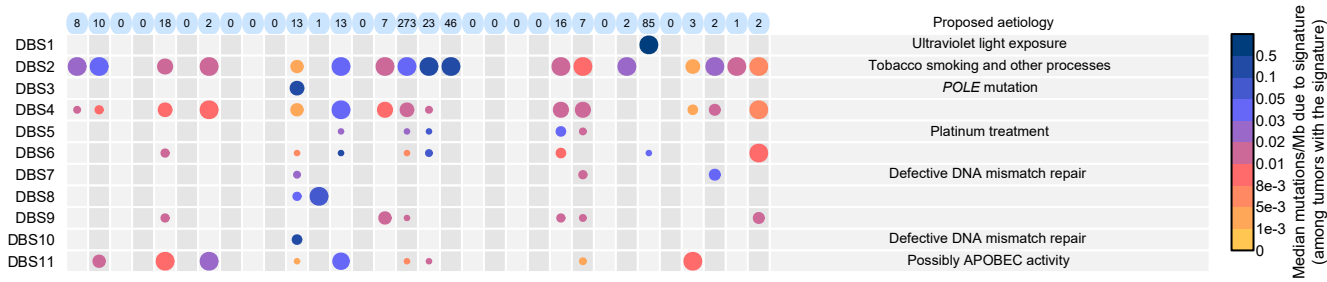


Fig 3

Single Base Substitution



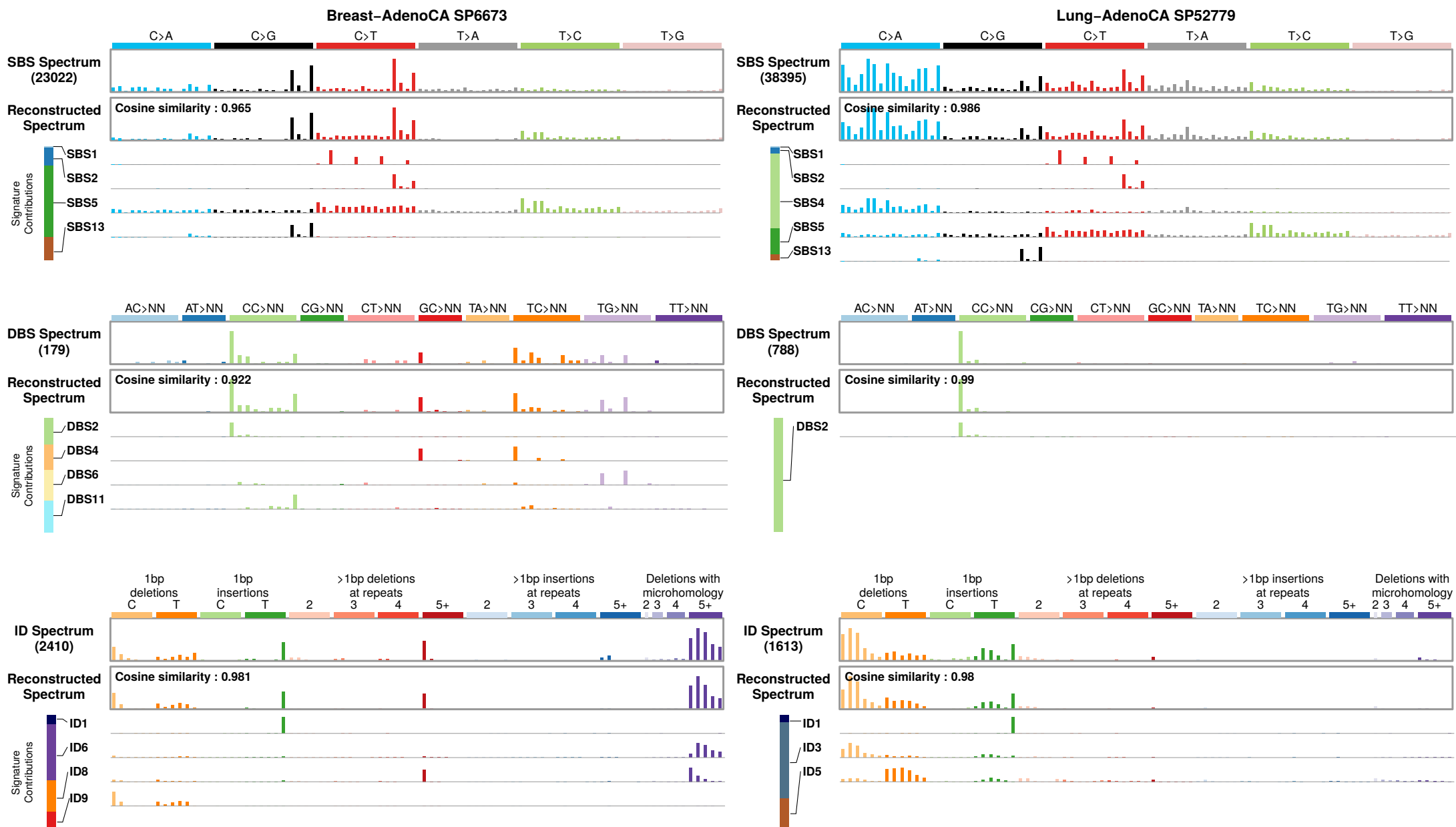
Doublet Base Substitution



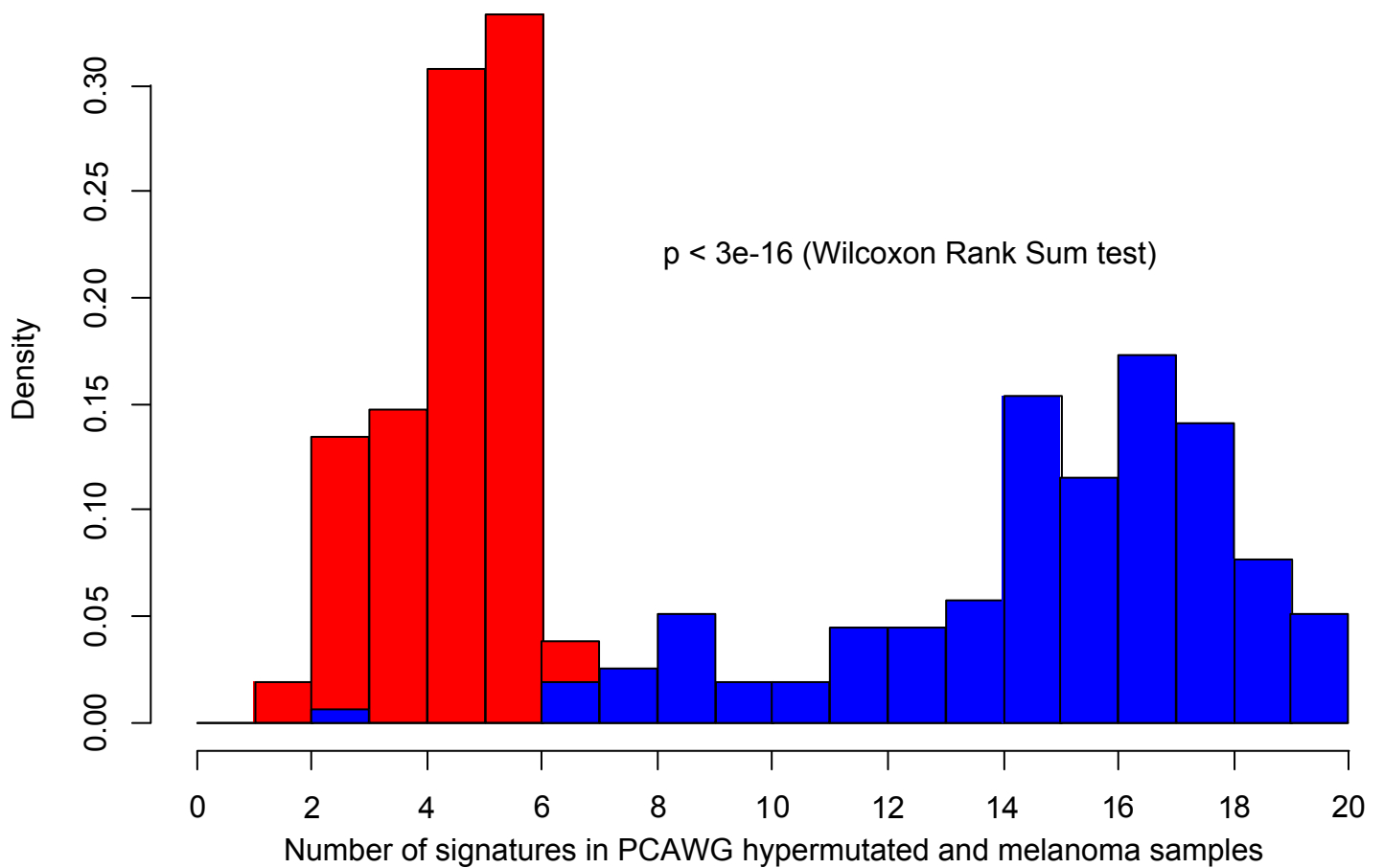
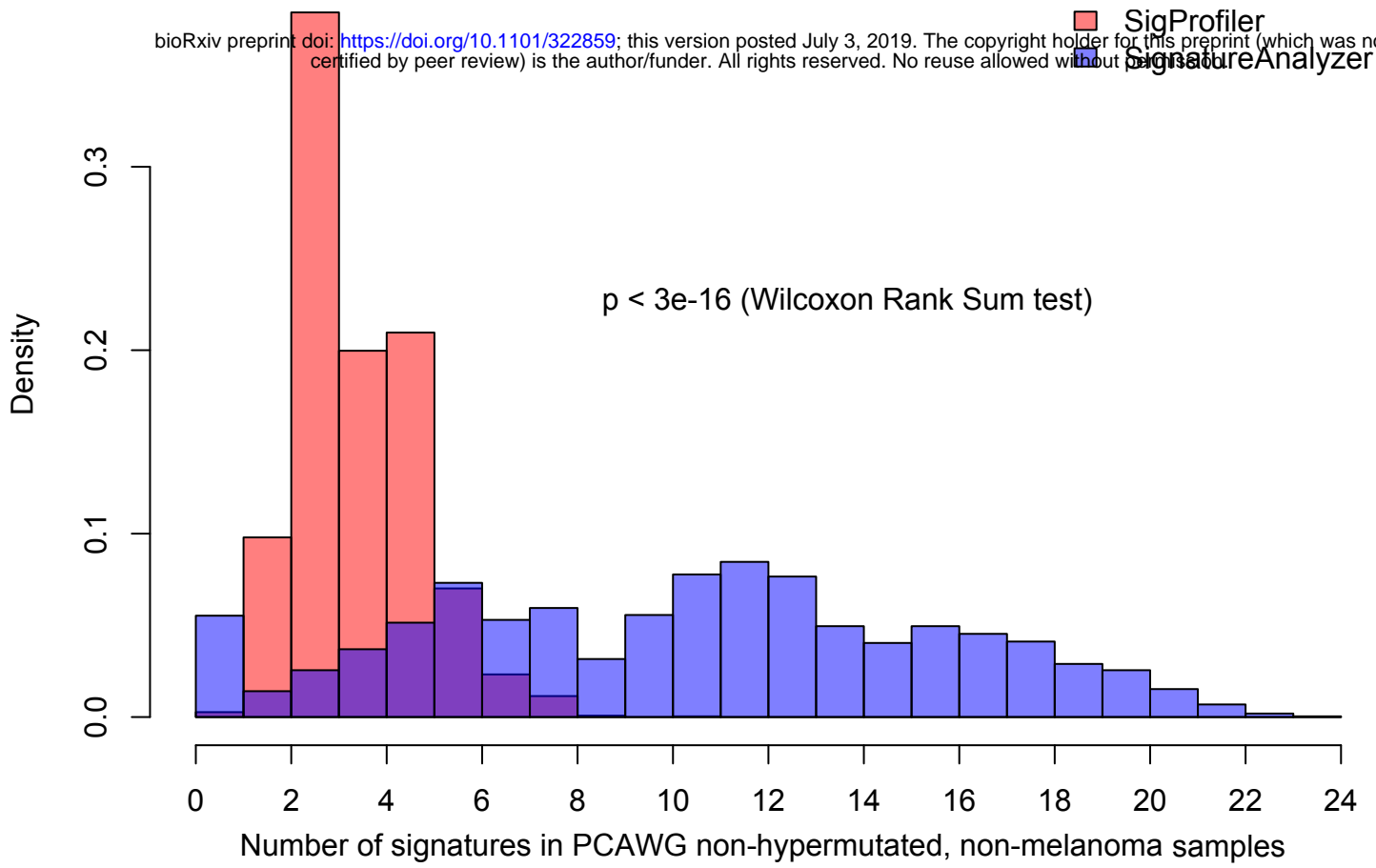
Insertion and Deletion



Fig 4

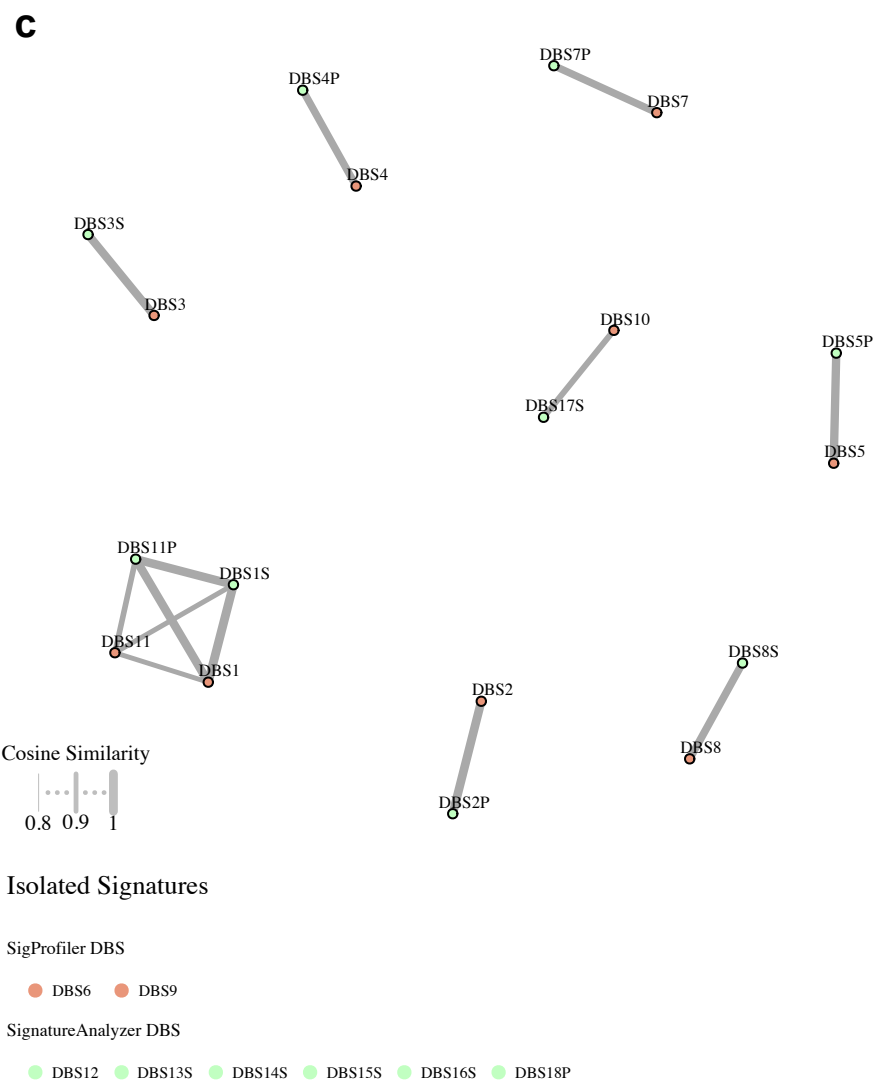
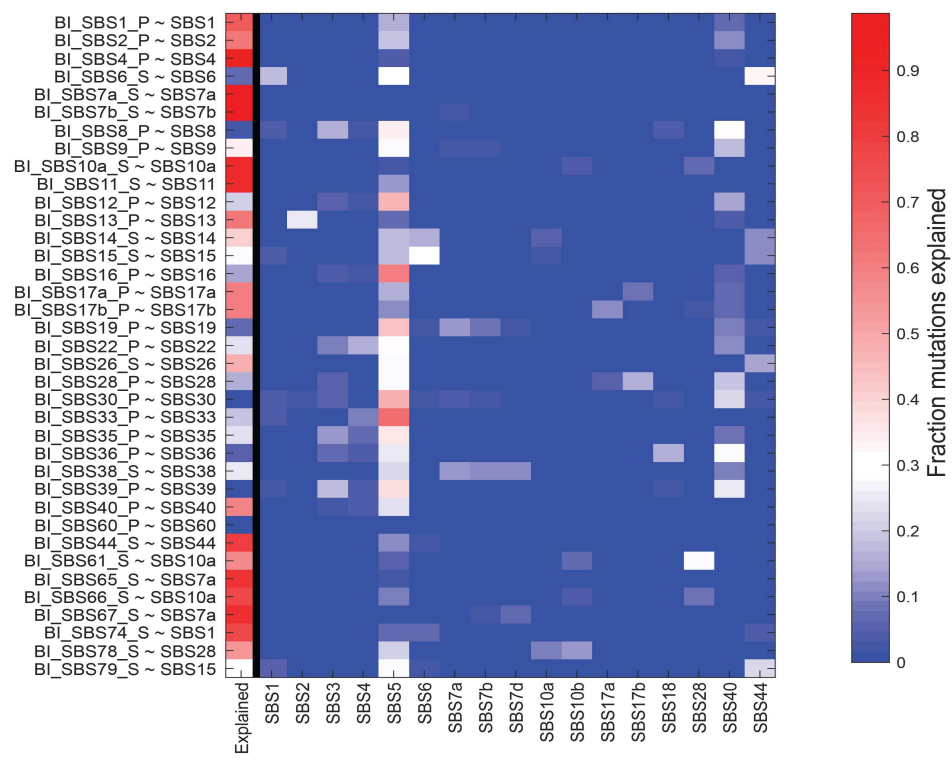
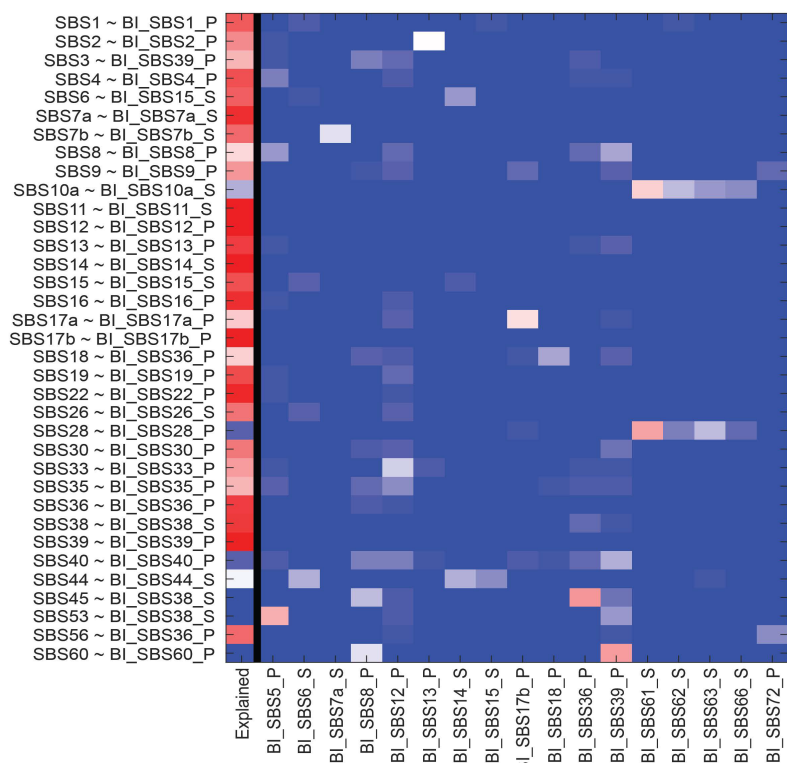


bioRxiv preprint doi: <https://doi.org/10.1101/322859>; this version posted July 3, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



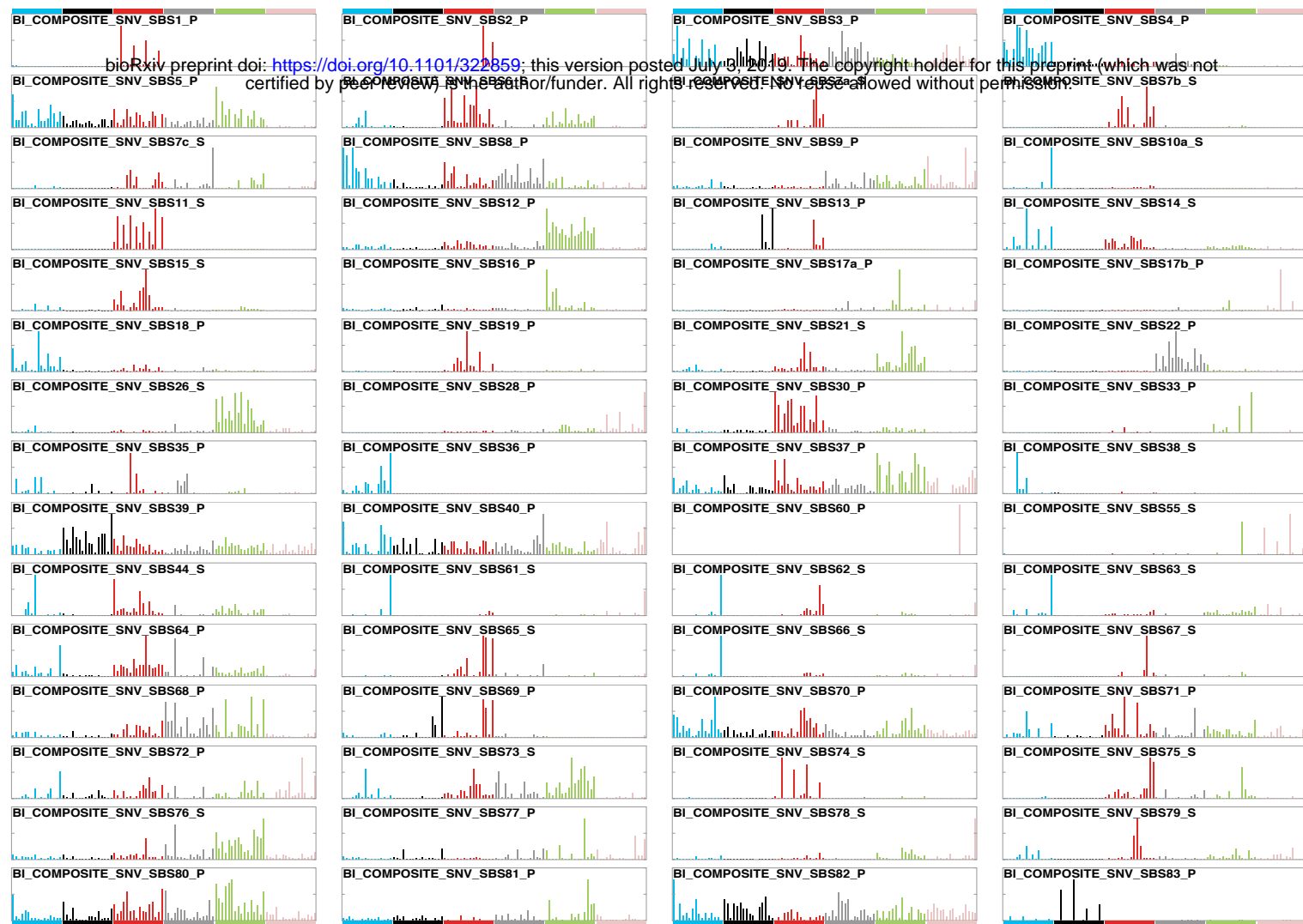
bioRxiv preprint doi: <https://doi.org/10.1101/322859>; this version posted July 3, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

SignatureAnalyzer signatures explained by sigprofiler

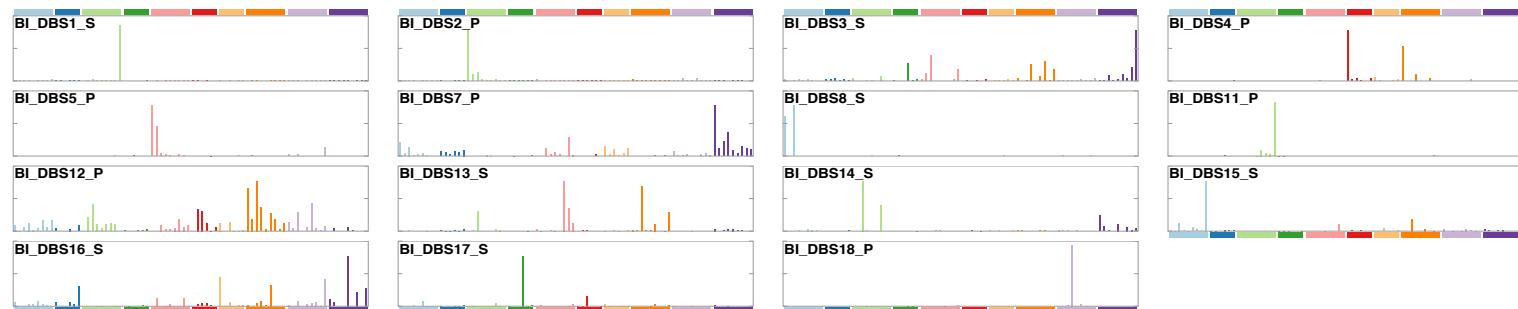


SignatureAnalyzer reference SBS signatures

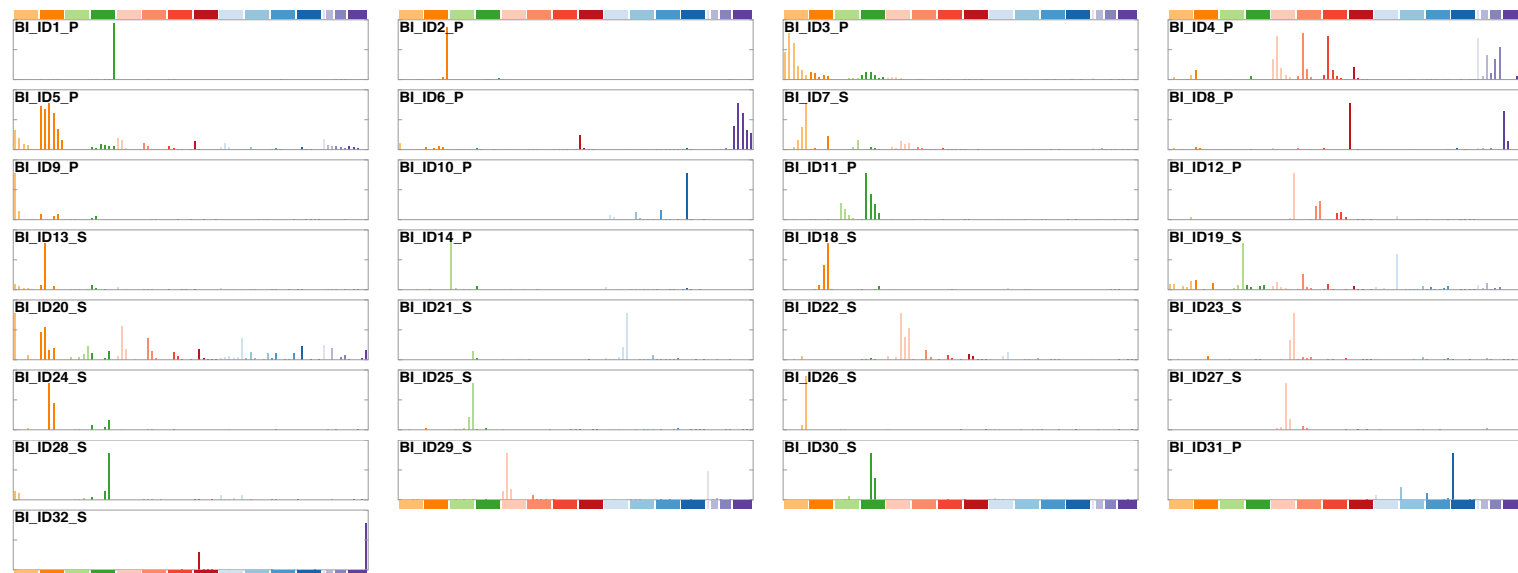
Extended data Fig 3



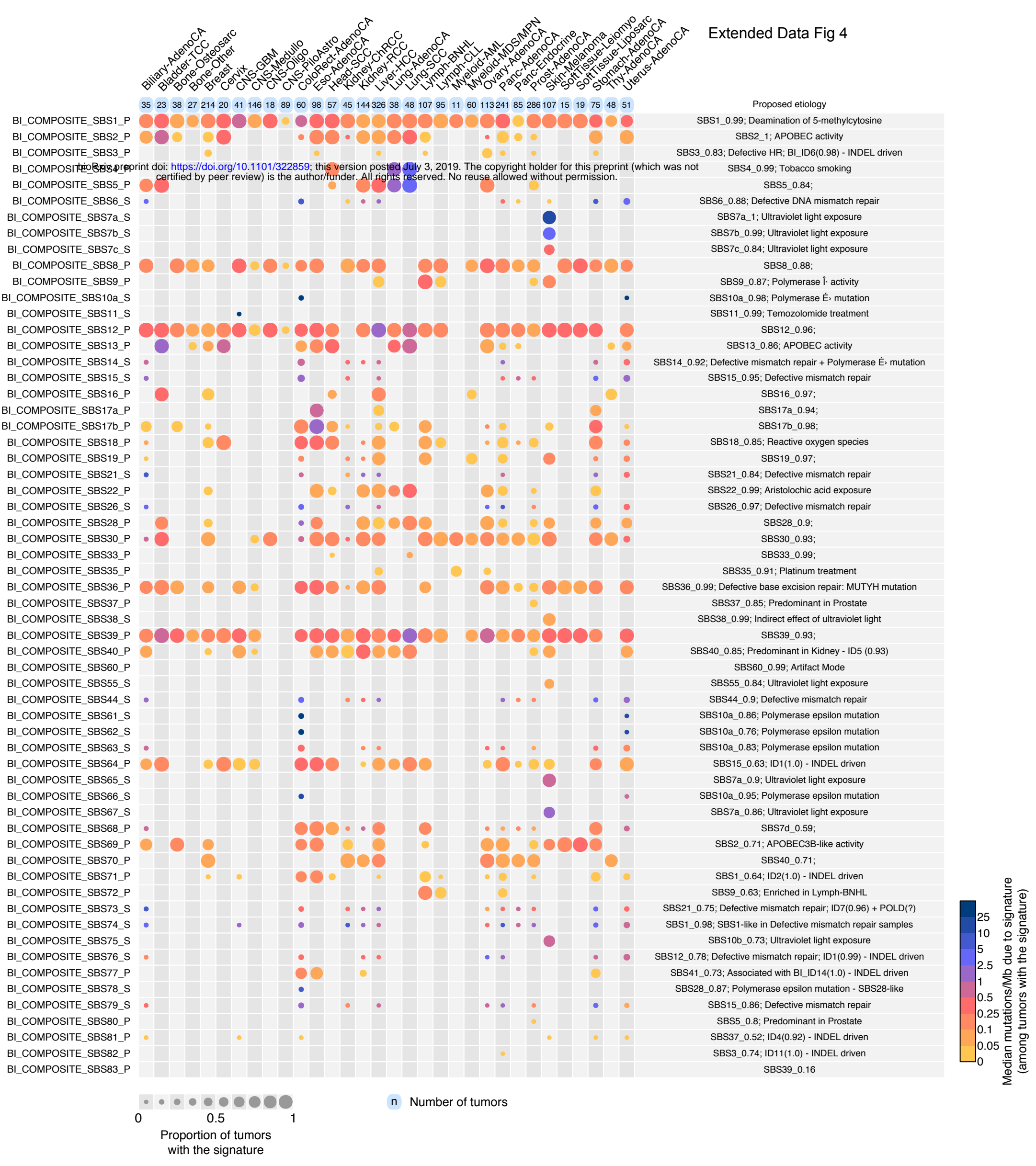
SignatureAnalyzer reference DBS signatures



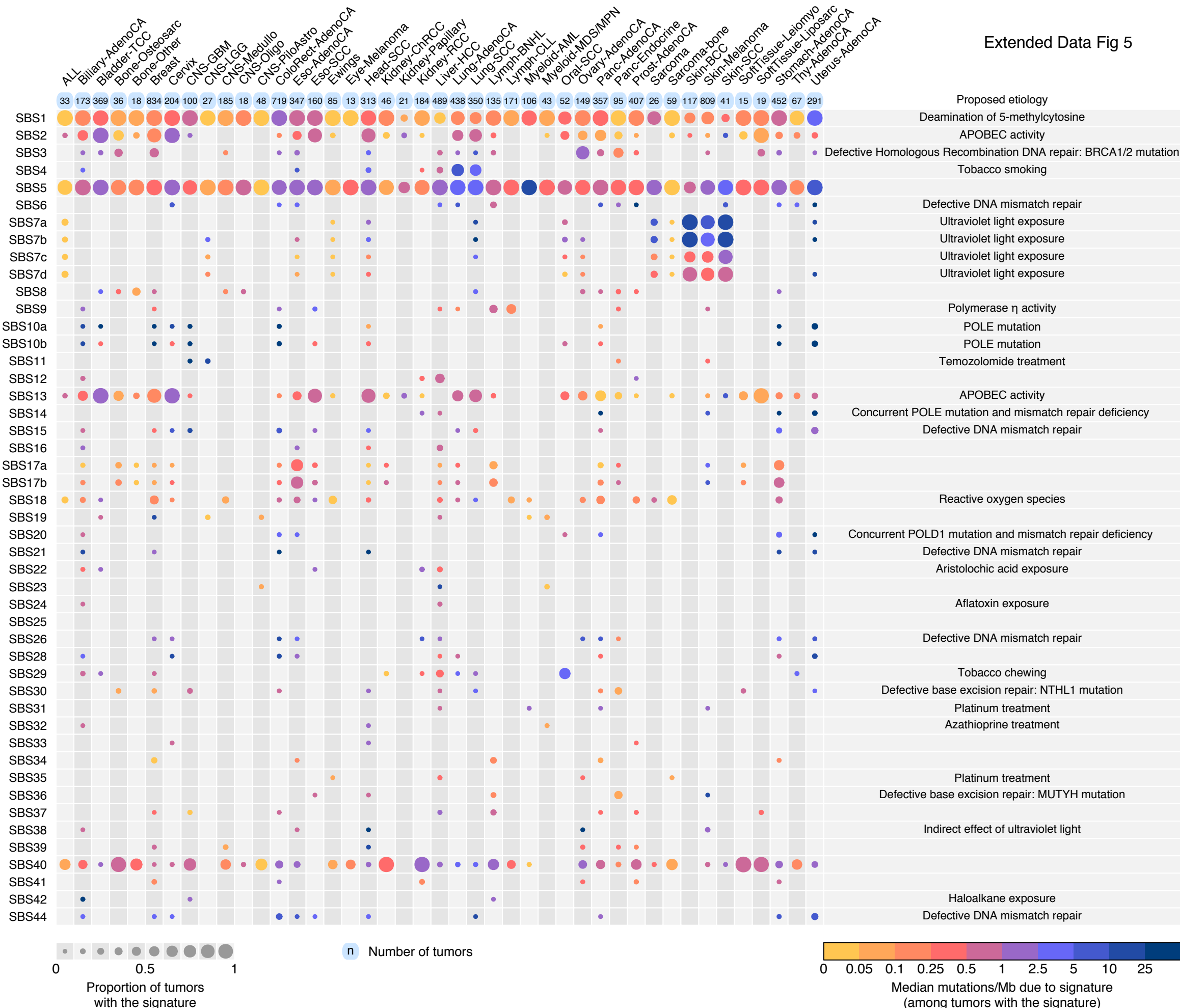
SignatureAnalyzer reference ID signatures



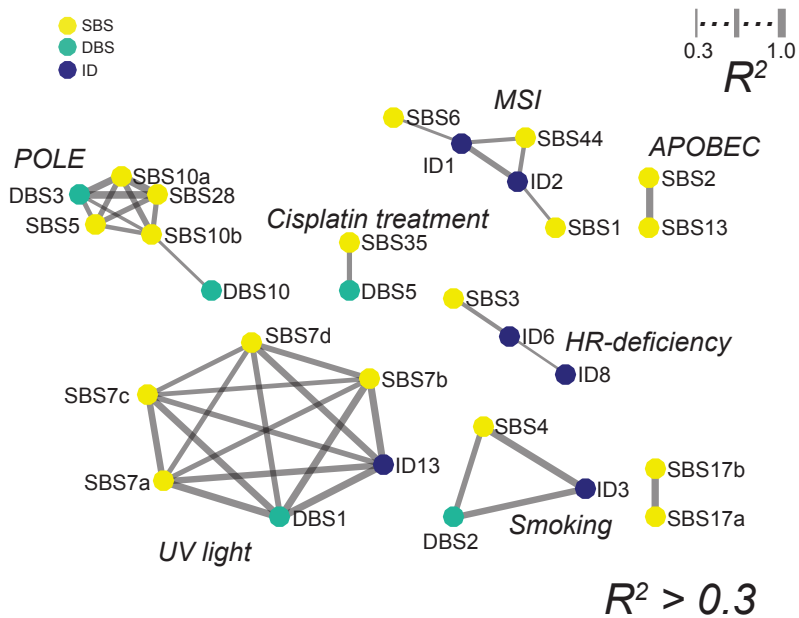
Extended Data Fig 4



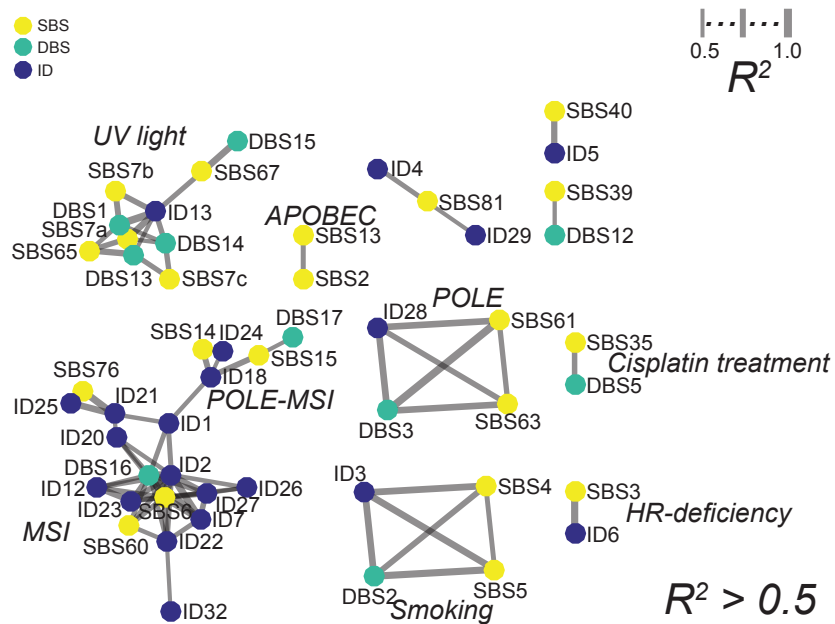
Extended Data Fig 5



a SigProfiler

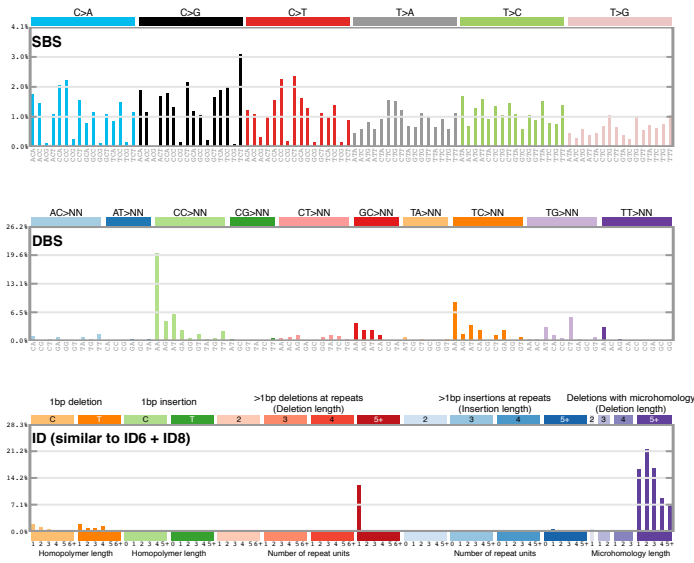


b SignatureAnalyzer



Extended Data Fig 7

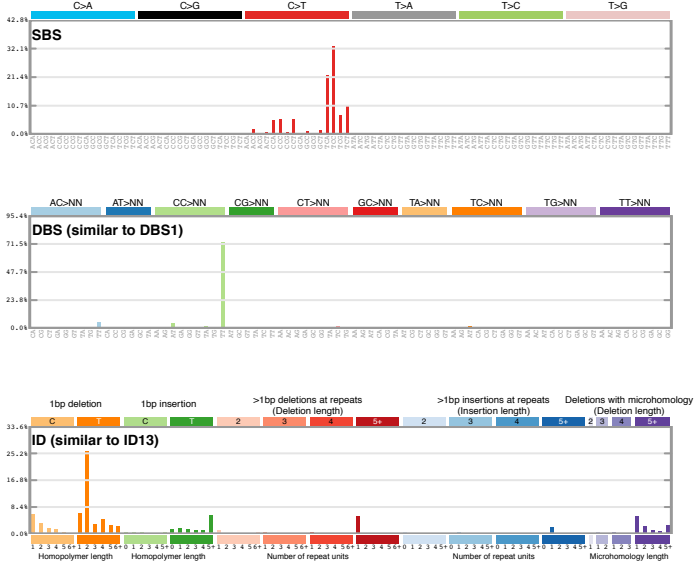
COMP-3



COMP-4



COMP-7a



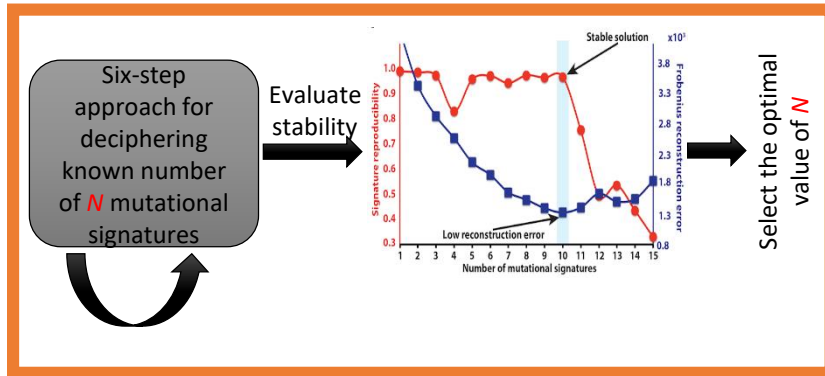
COMP-7b



a Extraction of mutational signatures

Step A (Apply the approach to a set of samples D ; initially D contains all samples, i.e., $D=M$)

Described in detail in (Alexandrov et al., Cell Rep. 2013;3(1):246-59).



Repeat $N = 1 \dots (G - 1)$

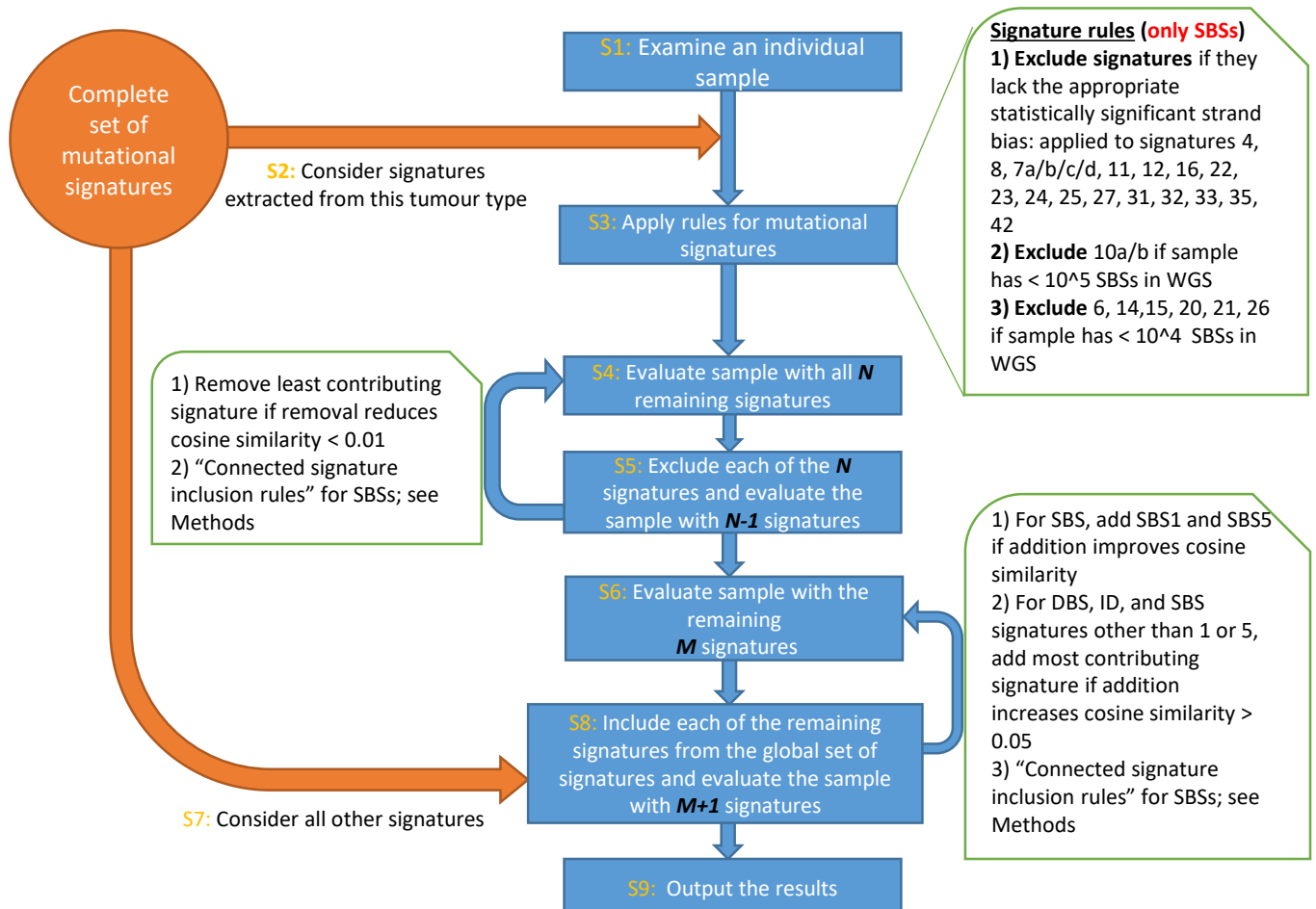
Step B (Solution evaluation and re-iteration)

Extracted mutational signatures and their activities to individual samples are saved into a set S . The activity of any signature that does not increase the cosine similarity of a sample with more 0.01 was removed from the sample (i.e., assigned a value of zero). **Step A** is repeated for all samples for which the identified signatures do not explain their patterns (cosine similarity < 0.95). The algorithm continues to the **step C** when **step A** cannot find any stable signatures.

Step C (Clustering of mutational signatures)

Hierarchical consensus clustering was applied to the set S to derive the consensus mutational signatures across the set of samples M .

b Attribution of activities of mutational signatures in samples



Extended Data Table 1. The number of DBSs is proportional to the number of SBSs with the exception of a few cancer types (ColoRect-AdenoCA, Lung-AdenoCA, Lung-SCC, Skin-Melanoma), analysed by the following linear regression (computed by an R function call): $\text{glm}(\text{DBS.counts} \sim \text{SBS.counts} + \text{Cancer.Types})$

| | Estimate | Std.Error | t value | Pr(> t) | |
|---------------------|-----------|-----------|---------|-----------|-----|
| (Intercept) | 5.61E+00 | 8.76E+01 | 0.064 | 0.9489 | |
| SBS.counts | 3.74E-03 | 1.25E-04 | 29.841 | <2.00E-16 | *** |
| Bladder-TCC | 1.32E+01 | 1.39E+02 | 0.095 | 0.92432 | |
| Bone-Osteosarc | 2.18E+00 | 1.21E+02 | 0.018 | 0.98567 | |
| Bone-Other | -2.81E+00 | 1.33E+02 | -0.021 | 0.9831 | |
| Breast | 5.32E+00 | 9.44E+01 | 0.056 | 0.95511 | |
| Cervix | -1.06E+01 | 1.45E+02 | -0.073 | 0.94185 | |
| CNS-GBM | -2.81E+01 | 1.19E+02 | -0.236 | 0.81352 | |
| CNS-Medullo | -7.04E+00 | 9.75E+01 | -0.072 | 0.94239 | |
| CNS-Oligo | -1.03E+01 | 1.50E+02 | -0.069 | 0.94539 | |
| CNS-PiloAstro | -5.87E+00 | 1.03E+02 | -0.057 | 0.95467 | |
| ColoRect-AdenoCA | -4.11E+02 | 1.12E+02 | -3.667 | 0.00025 | *** |
| Eso-AdenoCA | -1.56E+01 | 1.02E+02 | -0.153 | 0.87838 | |
| Head-SCC | 5.27E+01 | 1.11E+02 | 0.474 | 0.63541 | |
| Kidney-ChRCC | -3.14E+00 | 1.17E+02 | -0.027 | 0.97857 | |
| Kidney-RCC | 5.61E+01 | 9.76E+01 | 0.574 | 0.56584 | |
| Liver-HCC | 7.82E+01 | 9.21E+01 | 0.849 | 0.39575 | |
| Lung-AdenoCA | 5.02E+02 | 1.21E+02 | 4.136 | 3.63E-05 | *** |
| Lung-SCC | 5.85E+02 | 1.15E+02 | 5.078 | 4.08E-07 | *** |
| Lymph-BNHL | 1.04E+01 | 1.01E+02 | 0.103 | 0.91765 | |
| Lymph-CLL | -4.30E+00 | 1.02E+02 | -0.042 | 0.96655 | |
| Myeloid-AML | -1.89E+00 | 1.79E+02 | -0.011 | 0.99156 | |
| Myeloid-MDS/MPN | -7.43E+00 | 1.10E+02 | -0.067 | 0.94622 | |
| Ovary-AdenoCA | 3.59E+01 | 1.00E+02 | 0.358 | 0.72023 | |
| Panc-AdenoCA | -8.34E-01 | 9.37E+01 | -0.009 | 0.99289 | |
| Panc-Endocrine | -5.70E+00 | 1.04E+02 | -0.055 | 0.95628 | |
| Prost-AdenoCA | 2.52E+00 | 9.27E+01 | 0.027 | 0.97831 | |
| Skin-Melanoma | 1.67E+03 | 1.02E+02 | 16.47 | <2.00E-16 | *** |
| SoftTissue-Leiomyo | 5.98E+00 | 1.60E+02 | 0.037 | 0.97016 | |
| SoftTissue-Liposarc | 7.77E+00 | 1.48E+02 | 0.053 | 0.95804 | |
| Stomach-AdenoCA | -3.04E+01 | 1.06E+02 | -0.287 | 0.77417 | |
| Thy-AdenoCA | -4.80E+00 | 1.15E+02 | -0.042 | 0.96676 | |
| Uterus-AdenoCA | -1.25E+02 | 1.14E+02 | -1.096 | 0.27304 | |

Extended Data Table 2. Numbers of insertion/deletion mutations due to ID1, ID2, and all other ID signatures combined, in hypermutators and non-hypermutators

| Signature | Hypermutators | | Non-hypermutators | | All Tumours | |
|---------------------|------------------|----------|-------------------|----------|------------------|----------|
| | Count | Fraction | Count | Fraction | Count | Fraction |
| ID1 | 593,935 | 0.236 | 399,633 | 0.276 | 993,568 | 0.250 |
| ID2 | 1,838,867 | 0.730 | 252,893 | 0.174 | 2,091,760 | 0.527 |
| ID1+ID2 | 2,432,802 | 0.966 | 652,526 | 0.450 | 3,085,328 | 0.777 |
| Other ID signatures | 85,038 | 0.034 | 797,964 | 0.550 | 883,002 | 0.223 |
| Total | 2,517,840 | 1 | 1,450,490 | 1 | 3,968,330 | 1 |