

The Repertoire of Mutational Signatures in Human Cancer

Ludmil B Alexandrov^{1*}, Jaegil Kim^{2*}, Nicholas J Haradhvala^{2,3*}, Mi Ni Huang^{4*}, Alvin WT Ng⁴, Arnoud Boot⁴, Kyle R Covington⁵, Dmitry A Gordenin⁶, Erik Bergstrom¹, Nuria Lopez-Bigas^{7,8,9}, Leszek J Klimczak¹⁰, John R McPherson⁴, Sandro Morganello¹¹, Radhakrishnan Sabarinathan^{7,8}, David A Wheeler⁵, Ville Mustonen¹², Gad Getz^{2,3,13,14**}, Steven G Rozen^{4**§}, Michael R Stratton^{11**§}, on behalf of the PCAWG Mutational Signatures Working Group and the ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Network

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

² Broad Institute, Cambridge, MA, USA

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

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

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

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

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

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

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

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

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

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

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

¹⁴ Harvard Medical School, Boston, MA 02215, USA

* Denotes equal contribution

** Denotes equal supervisory contribution

§Correspondence and requests for materials should be addressed to S.G.R. (steve.rozen@duke-nus.edu.sg) and M.R.S. (mrs@sanger.ac.uk)

Draft 2018 05 12

ABSTRACT

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

INTRODUCTION

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

Previous studies of multiple cancer types identified >30 single base substitution signatures, some of known but many of unknown aetiologies, some ubiquitous and others rare, some part of normal cell biology and others associated with abnormal exposures or operative during neoplastic progression^{7,16-27}. Six genome rearrangement signatures have also been identified in breast cancer¹⁸ and further patterns of rearrangements have been described^{15,28-30}. However, analysis of other mutation classes has been limited^{31,32}.

Thus far, mutational signature analysis has predominantly used cancer exome sequences. However, the many fold greater numbers of somatic mutations in whole-genome sequences provide substantially increased power for decomposition, enabling better separation of partially correlated signatures and extraction of signatures which contribute relatively small numbers of mutations. Furthermore, technical artefacts and differences in sequencing technologies and mutation calling algorithms can themselves generate mutational signatures. Therefore, the uniformly processed and highly curated sets of all classes of somatic mutations from the 2,780 cancer genome sequences of the Pan Cancer Analysis of Whole-Genomes (PCAWG) project, combined with almost all other cancer genomes and exomes for which suitable mutational catalogues are publicly available, presents a notable opportunity to establish the repertoire of mutational signatures and to determine their activities across the range of cancer types.

RESULTS

Cancer genomes and somatic mutations

Somatic mutational catalogues from 23,829 samples of most cancer types, including the 2,780 highly curated PCAWG whole-genomes, 1,865 additional whole-genomes and 19,184 exomes were studied. From these, 79,793,266 somatic single base substitutions, 814,191 doublet base substitutions and 4,122,233 small insertions and deletions (indels) were analysed for mutational signatures, ~10-fold more mutations than any previous study (<https://www.synapse.org/#!Synapse:syn11801889>)^{4,16}.

To enable mutational signature analysis classifications were developed for each type of mutation. For single base substitutions, the primary classification comprised 96 classes

constituted by the six base substitutions C>A, C>G, C>T, T>A, T>C, T>G (in which the mutated base is represented by the pyrimidine of the Watson-Crick base pair) plus the flanking 5' and 3' bases. In some analyses, two flanking bases 5' and 3' to the mutated base were considered (generating 1,536 classes) or mutations within transcribed genome regions were selected and classified according to whether the pyrimidine of the mutated base pair fell on the transcribed or untranscribed strand (192 classes). A classification was also derived for doublet base substitutions (78 classes). Indels were classified as deletions or insertions and, when of a single base, as C or T and according to the length of the mononucleotide repeat tract in which they occurred. Longer indels were classified as occurring at repeats or with overlapping microhomology at deletion boundaries, and according to the size of indel, repeat, and microhomology (83 classes, <https://www.synapse.org/#!Synapse:syn11726616>)^{3,16}.

Mutational signature analysis

The mutational catalogues from the 2,780 PCAWG whole-genome, 1,865 additional whole-genome, and 19,184 exome sequences of cancer were analysed separately. For each of these catalogue sets, signature extraction was conducted using methods based on nonnegative matrix factorisation (NMF)^{3,7} on each cancer type individually and also on all cancer types together. Analyses were carried out separately for single base substitutions (SBS signatures), doublet base substitutions (DBS signatures) and indels (ID signatures) and also for the three mutation types together (257 mutation classes, or 1697 if the 1536 SBS classification was employed) generating composite signatures.

Mutational signatures were extracted using two independently developed NMF-based methods: (i) SigProfiler, a further elaborated version of the framework used to generate the signatures shown in COSMIC^{3,16,18,33,34}, and (ii) SignatureAnalyzer, based on a Bayesian variant of NMF used in several previous publications^{6,7,35,36}. NMF determines both the signature profiles and the contributions of each signature to each cancer genome as part of its factorization of the input matrix of mutation spectra. However, given a substantial number of signatures and/or heterogeneous mutation burdens across samples, it is possible to reconstruct the mutations observed in a particular sample in multiple ways, often with very small and/or biologically implausible contributions from many signatures. Therefore, each method developed a separate procedure to estimate the contributions of signatures to each sample (Methods).

The results of the two methods exhibited many similarities. However, there were also noteworthy differences. The number of SBS signatures found in low mutation burden tumours in the PCAWG set (94.4% of cases that harbour 47% of mutations) was similar: 31 by SigProfiler and 35 by SignatureAnalyzer. The number of additional SBS signatures extracted from hyper-mutated PCAWG samples (5.6% of cases and 53% of mutations, <https://www.synapse.org/#!Synapse:syn12016215>), however, was different: 13 by SigProfiler and 25 by SignatureAnalyzer. There were also differences in SBS signature profiles, including among signatures found in low mutation burden cases. The latter primarily involved “flat”, relatively featureless signatures, which are mathematically challenging to deconvolute. Finally, there were differences in signature attributions to individual samples. In general, SignatureAnalyzer used more signatures to reconstruct the mutational profiles (Extended Data Figure 1, <https://www.synapse.org/#!Synapse:syn12169204>, <https://www.synapse.org/#!Synapse:syn12177011>) and the attribution to flat signatures was

different, with SigProfiler assigning mutations to SBS5 and SBS40 and SignatureAnalyzer using combinations of multiple signatures (Extended Data Figure 2ab, <https://www.synapse.org/#!Synapse:syn12169204>). The DBS and ID signatures were generally similar between the two methods (Extended Data Figure 2cd). These comparisons provide a useful perspective on both the consistency and variability of signature extraction and attribution depending on the methodology used.

The final sets of reference mutational signatures were determined from the PCAWG analysis supplemented by additional signatures from the other datasets. Signatures were supported by the outcomes of analyses using the 192 and 1536 mutation classifications, the existence of individual cancer samples dominated by a particular signature, and, where available, prior experimental evidence for certain mutational signatures (Methods and <https://www.synapse.org/#!Synapse:syn12009767>). Each signature was allocated a number consistent with, and extending, the COSMIC annotation³³. Some previous signatures split into multiple constituent signatures and these were numbered as before but with additional letter suffixes (eg, single SBS17 split into signatures SBS17a and SBS17b). DNA sequencing and analysis artefacts also generate mutational signatures, and we indicate which signatures are possible artefacts (<https://www.synapse.org/#!Synapse:syn12009767>) but do not present them below. However, future studies employing this signature set as a reference may consider utilizing artefact signatures for data quality control. The results of both SignatureAnalyzer and SigProfiler were used throughout the research reported here. However, for brevity and for continuity with the signature set previously displayed in COSMIC³³, which has been widely used as a reference, SigProfiler results are outlined below and SignatureAnalyzer results are provided at (Extended Data Figures 3,4, <https://www.synapse.org/#!Synapse:syn11738307>).

Single base substitution (SBS) mutational signatures

There were substantial differences in numbers of SBSs between samples (ranging from hundreds to millions) and between cancer types, as previously observed^{16,37} (Figure 1). In total, 67 SBS mutational signatures were extracted, of which 49 were considered to be likely real (Figure 2, Methods, <https://www.synapse.org/#!Synapse:syn12009783>). Except for SBS25, all mutational signatures previously reported on COSMIC^{4,23,33} were confirmed in the new set of analyses (median cosine similarity between the newly derived signatures and those on COSMIC: 0.95, excluding "split" signatures which are discussed below; <https://www.synapse.org/#!Synapse:syn12016215>). SBS25 was previously found only in cell lines derived from Hodgkin lymphomas, some of which had been previously treated with chemotherapy, and, to our knowledge, no data from primary cancers of this type are currently available. The newly derived signatures showed much improved separation from each other and hence more distinct signature profiles, presumably due to the substantially increased statistical power of this analysis (<https://www.synapse.org/#!Synapse:syn12009783>).

Thirteen new likely real SBS signatures compared to the set previously described in COSMIC³³ were extracted (excluding those that are the consequence of signature splitting). Some were in cancers with a previously unanalysed exogenous exposure (SBS42), some were in chemotherapy treated samples which have often been excluded from previous studies (SBS31, SBS32, SBS35) and some were rare and hence absent by chance from previous analyses (SBS36, SBS44). Others were more common, but contributed relatively few mutations to individual cancer genomes, or were similar to previously discovered signatures

and thus not isolated from datasets based predominantly on cancer exome sequences (eg SBS38, SBS39, SBS40). Notably, SBS40 was extracted from kidney cancer in which it appears to be required for optimal reconstruction of mutational catalogues. It is a relatively featureless (“flat”) signature, with similarity to SBS5 and other flat signatures, and this may account for it only clearly emerging now with the availability of whole cancer genomes. SBS40 may contribute to other cancer types but its similarity to SBS5 renders this uncertain and larger datasets will be required to clarify the extent of its activity. For some new signatures there were plausible underlying aetiologies (Figure 3, Extended Data Figures 4,5): SBS31 and SBS35, prior platinum compound chemotherapy³⁸; SBS32, prior azathioprine therapy; SBS36, inactivating germline or somatic mutations in *MUTYH* which encodes a component of the base excision repair machinery^{39,40}; SBS38, additional effects of ultraviolet light (UV) exposure; SBS42, occupational exposure to haloalkanes²⁷; SBS44, defective DNA mismatch repair due to MLH1 inactivation⁴¹. SBS33, SBS34, SBS37, SBS39, SBS40, and SBS41 are of unknown cause.

Three previously characterised base substitution signatures (SBS7, SBS10, SBS17) split into multiple constituent signatures (Figure 2). We previously regarded SBS7 as a single signature composed predominantly of C>T at CCN and TCN trinucleotides (the mutated base is underlined) together with many fewer T>N mutations. It was found in malignant melanomas and squamous skin carcinomas and is likely due to UV induced pyrimidine dimer formation followed by translesion DNA synthesis by error-prone polymerases which predominantly insert adenine opposite damaged bases. With the larger dataset now available, SBS7 has decomposed into four constituent signatures: SBS7a consisting mainly of C>T at TCN; SBS7b consisting of C>T mainly at CCN and to a lesser extent at TCN; SBS7c and SBS7d, which constituted relatively minor components of the previous SBS7 and consist predominantly of T>A at NTT and T>C at NTT respectively⁴². Splitting of a mutational signature likely reflects the existence of multiple distinct mutational processes, initiated by the same exposure, which have closely, but not perfectly, correlated activities. For example, the constituent signatures of SBS7 are probably all initiated by UV-induced DNA damage. SBS7a and SBS7b may reflect different dipyrimidine photoproducts whereas SBS7c and SBS7d may be due to low frequencies of misincorporation by translesion polymerases of T and G opposite thymines in pyrimidine dimers rather than the more frequent and non-mutagenic A. Splitting of SBS10 and SBS17 is described at <https://www.synapse.org/#!Synapse:syn12009783>.

Several base substitution signatures showed transcriptional strand bias (<https://www.synapse.org/#!Synapse:syn12009767>). Transcriptional strand bias is often attributable to transcription coupled nucleotide excision repair (TC-NER) acting on DNA damaged by exogenous exposures which cause covalently bound bulky adducts or crosslinking to other bases and consequent distortion of the helical structure. This results in stalling of RNA polymerase and hence recruitment of the TC-NER machinery. An excess of DNA damage on untranscribed compared to transcribed strands of genes may also contribute to transcriptional strand bias⁴³. Both mechanisms, however, result in more mutations of a damaged base on the untranscribed compared to the transcribed strands of genes. Assuming that either or both are responsible for the observed transcriptional strand biases (which may not always be the case), DNA damage to cytosine (SBS7a, SBS7b), guanine (SBS4, SBS8, SBS19, SBS23, SBS24, SBS31, SBS32, SBS35, SBS42), thymine (SBS7c, SBS7d, SBS21, SBS26, SBS33) and adenine (SBS5, SBS12, SBS16, SBS22, SBS25) may underlie these mutational signatures (see

<https://www.synapse.org/#!Synapse:syn12009783> for plots of strand bias). Although the likely underlying DNA damaging agents are known for SBS4 (tobacco mutagens), SBS7a, SBS7b, SBS7c, SBS7d (UV), SBS22 (aristolochic acid), SBS24 (aflatoxin), SBS25 (prior chemotherapy), SBS31 and SBS35 (platinum compounds), SBS32 (azathioprine), and SBS42 (haloalkanes), the causes of the remainder are unknown. Indeed, some signatures showing transcriptional strand bias are associated with defective DNA mismatch repair (SBS21 and SBS26) and it is conceivable that, for these, exogenous DNA damage is not involved. The extent of transcriptional strand bias appears to differ in different sectors of the genome. For example, consideration of the whole transcribed genome showed absent or minimal transcriptional strand bias in the APOBEC related SBS2 and SBS13 and in the defective polymerase epsilon proof-reading related SBS10a. However, consideration of exons alone showed clear evidence of transcriptional strand bias in these signatures (<https://www.synapse.org/#!Synapse:syn12009783>). The mechanism(s) underlying this amplification of transcriptional strand bias in exons is unknown and appears to be signature specific, since there is minimal difference in the extent of transcriptional strand bias between exons and other transcribed regions for other signatures (for example, SBS4).

Employing the single base substitution classification of 1536 mutation types, which uses the pentanucleotide sequence context two bases 5' and two bases 3' to each mutated base, yielded a set of signatures largely consistent with that based on substitutions in trinucleotide context alone. Notably, however, the pentanucleotide context enabled the extraction of two forms of both SBS2 and SBS13, one with mainly a pyrimidine (C or T) and the other with a purine (A or G) at the -2 base (the second base 5' to the mutated cytosine). These may represent the activities of the cytidine deaminases APOBEC3A and APOBEC3B, respectively⁴⁴. If so, APOBEC3A accounts for many more mutations than APOBEC3B in cancers with high APOBEC activity. Several other signatures showed non-random sequence contexts at +2 and -2 positions. In particular, the -2 bases in SBS17a and SBS17b and the -2 and +2 bases in SBS9 were predominantly A and T. In general, however, sequence context effects were much stronger for bases immediately 5' and 3' to the mutated bases.

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

Clustered single base substitution mutational signatures

Some mutational processes generate mutations that cluster in small regions of the genome. The relatively limited number of mutations generated by such processes, compared to those acting genome-wide, may result in failure to detect their signatures by standard methods. To obviate this problem, we first identified clustered mutations in each genome and analysed

these separately (Methods). Four main signatures associated with clustered mutations were identified (Figure 2). Two found in multiple cancer types were similar to single base substitution SBS2 and SBS13, which have been attributed to APOBEC enzyme activity (mostly APOBEC3B) and represent foci of *kataegis*^{17,45} (Methods). Two additional clustered mutational signatures, one characterised by C>T and C>G mutations at (A|G)C(C|T) trinucleotides⁴⁶ and the other T>A and T>C mutations at (A|T)I(A|T) were found in lymphoid neoplasms and likely represent direct and indirect consequences of activation induced cytidine deaminase (AID) mutagenesis (SBS84 and SBS85 respectively)⁶. The possibility that further processes may generate clustered mutations is not excluded.

Doublet base substitution (DBS) mutational signatures

Tandem doublet, triplet, quadruplet, quintuplet, and sextuplet base substitutions (<https://www.synapse.org/#!Synapse:syn11801938>, <https://www.synapse.org/#!Synapse:syn11726620>) at immediately adjacent bases were observed at ~1% the prevalence of single base substitutions. In most cancer genomes, the observed number of DBSs was considerably higher than expected from random adjacency of SBSs (<https://www.synapse.org/#!Synapse:syn12177057>) indicating the existence of commonly occurring, single mutagenic events that cause substitutions at neighbouring bases. There was substantial variation in the number of DBSs, ranging from zero to 20,818 in a sample. Across cancer types, the numbers of DBSs were generally proportional to the numbers of SBSs in that cancer type (Figure 1). However, colorectal adenocarcinomas had significantly fewer DBS than expected, and lung cancers and melanomas had more (Extended Data Table 1). The large dataset analysed here allowed, for the first time, systematic analysis of DBS and indel signatures (described below). Eleven DBS signatures were extracted (Figure 2).

DBS1 was characterised almost exclusively by CC>TT mutations (Figure 2), contributed 100s-10,000s of mutations in malignant melanomas (Figure 3) with SBS7a and SBS7b. DBS1 exhibited transcriptional strand bias consistent with damage to cytosines (<https://www.synapse.org/#!Synapse:syn12177063>). CC>TT mutations associated with UV induced DNA damage are well established in the literature and are thought to be due to generation of pyrimidine dimers and subsequent error-prone translesion DNA synthesis by polymerases that introduce adenines opposite the damaged bases⁴⁷.

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

DBS2 was composed predominantly of CC>AA mutations, with smaller numbers of CC>AG and CC>AT mutations, and contributed 100s-1000s of mutations in lung adenocarcinoma, lung squamous and head and neck squamous carcinomas, which are often caused by tobacco smoking (Figures 2 and 3). DBS2 showed transcriptional strand bias indicative of guanine damage (<https://www.synapse.org/#!Synapse:syn12177064>) and was associated with SBS4

which is caused by tobacco smoke exposure. It is likely, therefore, that DBS2 can be a consequence of DNA damage by tobacco smoke mutagens.

Analysis of each cancer type separately, however, revealed a signature very similar to DBS2 contributing 100s of mutations to liver cancers and 10s of mutations to cancers of other types without evidence of tobacco smoke exposure. A pattern closely resembling DBS2 and characterised predominantly by CC>AA mutations, together with smaller contributions of CC>AG and CC>AT, dominates DBSs in normal mouse cells and is particularly frequent in the liver⁴⁸. The nature of the mutational processes underlying these doublet signatures in smoking-unrelated human cancers and in normal mice is unknown. However, acetaldehyde exposure in experimental systems generates a mutational signature characterised primarily by CC>AA and lower burdens of CC>AG and CC>AT mutations together with single base substitution C>A mutations⁴⁹. Acetaldehyde is an oxidation product of alcohol and a constituent of cigarette smoke. The role of acetaldehyde, and perhaps other aldehydes, in generating DBS2, whether associated with tobacco smoking, alcohol consumption or in non-exposed cells, merits further investigation⁵⁰.

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

Small insertion and deletion (ID) mutational signatures

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

Indel signature 1 (ID1) was composed predominantly of insertions of thymine and ID2 of deletions of thymine, both at long (≥ 5) thymine mononucleotide repeats (Figure 2). 10s to 100s of mutations of both signatures were found in the large majority of most cancer types but were particularly common in colorectal, stomach, endometrial and oesophageal cancers and in diffuse large B cell lymphoma (Figure 3). Most of these cancers are likely to be DNA mismatch repair proficient on the basis of the relatively limited numbers of indels and absence of the SBS signatures (SBS6, SBS14, SBS15, SBS20, SBS21, SBS26, and SBS44) associated with DNA mismatch repair deficiency. Together, ID1 and ID2 accounted for 97% and 45% of indels in hypermutated and non-hypermutated cancer genomes, respectively (Extended Data Table 2), and both signatures have also been found in non-neoplastic cells⁵¹. They are likely due to the intrinsic tendency to slippage during DNA replication of long mononucleotide tracts. However, the mechanistic basis for separation into two signatures, one presumably due to slippage of the nascent strand (ID1) and the other the template strand (ID2) is unclear. Similarly, the substantial differences in their mutation frequencies between cancer types are not well understood.

ID3 was characterised predominantly by deletions of cytosine at short (≤ 5 bp long) mononucleotide cytosine repeats and exhibited 100s of mutations in tobacco smoking associated cancers of the lung and head and neck (Figures 2 and 3). There was transcriptional strand bias of mutations, with more guanine deletions than cytosine deletions on the untranscribed strands of genes, compatible with TC-NER of adducted guanine (<https://www.synapse.org/#!/Synapse:syn12177065>, <https://www.synapse.org/#!/Synapse:syn12177066>). The numbers of ID3 mutations in cancer samples positively correlated with the numbers of SBS4 and DBS2 mutations, both of which have been associated with tobacco smoking (Extended Data Figure 6). It is therefore likely that DNA damage by components of tobacco smoke underlie ID3 but the mechanism(s) by which indels are generated is unclear.

ID13 was characterised predominantly by deletions of thymine at thymine-thymine dinucleotides and exhibited large numbers of mutations in malignant melanomas of the skin (Figures 2 and 3). The numbers of ID13 mutations correlated with the numbers of SBS7a, SBS7b and DBS1 mutations, which have been attributed to DNA damage induced by UV (Extended Data Figure 6). It is, however, notable that a similar mutation of the other pyrimidine, ie deletion of cytosine at cytosine-cytosine dinucleotides, does not feature strongly in ID13, perhaps reflecting the predominance of thymine compared to cytosine dimers induced by UV⁵². The mechanism(s) underlying thymine deletion is unclear.

ID6 and ID8 were both characterised predominantly by deletions ≥ 5 bp (Figure 2). ID6 exhibited overlapping microhomology at deletion boundaries with a mode of 2bp and often longer stretches. This signature was correlated with SBS3 which has been attributed to defective homologous recombination based repair (Extended Data Figure 6). By contrast, ID8 deletions showed shorter or no microhomology at deletion boundaries, with a mode of 1bp, and did not strongly correlate with SBS3 mutations (Figures 2 and 3). These indel patterns are characteristic of DNA double strand break repair by non-homologous recombination based end-joining mechanisms and indicate that at least two distinct forms of end-joining mechanism are operative in human cancer⁵³.

A small fraction of cancers exhibited very large numbers of ID1 and ID2 mutations ($>10,000$) (Figure 3, <https://www.synapse.org/#!/Synapse:syn12009775>). These were usually accompanied by SBS6, SBS14, SBS15, SBS20, SBS21, SBS26 and/or SBS44 which are associated with DNA mismatch repair deficiency, sometimes combined with POLE or POLD1 proofreading deficiency (SBS14, SBS20)³⁶. Occasional cases with these signatures additionally showed large numbers of ID7 indels (<https://www.synapse.org/#!/Synapse:syn11738668>). In addition, rare samples showed large numbers of either ID4, ID11, ID14, ID15, ID16 or ID17 mutations but did not show ID1 and ID2 mutations or the single base substitution signatures usually associated with DNA mismatch repair deficiency. The mechanisms underlying these signatures are unknown.

Composite mutational signatures

In the analyses described above mutational signatures were extracted for each mutation type separately. However, mutational processes in nature generate composite signatures that may include SBSs, DBSs, IDs, genome rearrangements and chromosome number changes. We

therefore also extracted signatures using combined catalogues of SBSs, DBSs, and IDs (257 mutation subclasses or 1697 if the 1536 classification of single base substitutions was used). Fifty two composite signatures were extracted.

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

Correlations with age

A positive correlation between age of cancer diagnosis and the number of mutations attributable to a signature suggests that the mutational process underlying the signature has been operative, at a more or less constant rate, throughout the cell lineage from fertilized egg to cancer cell, and thus in normal cells from which that cancer type develops^{4,54}. Confirming previous reports, the numbers of SBS1 and SBS5 mutations correlated with age, exhibiting different rates in different tissue types (Q values in <https://www.synapse.org/#!Synapse:syn12030687>, <https://www.synapse.org/#!Synapse:syn12217988>). In addition, SBS40 correlated with age in multiple cancer types. However, given the similarity in signature profile between SBS5 and SBS40 the possibility of misattribution between these signatures cannot currently be excluded. The numbers of DBSs and IDs were much lower than the numbers of SBSs and the numbers of samples in which DBS and ID signatures could be attributed were also lower. Nevertheless, DBS2 and DBS4 correlated with age and, consistent with the interpretation of activity in normal cells, the profiles of DBS2 and DBS4 together closely resemble the spectrum of DBS mutations found in normal mouse cells⁴⁸. Neither DBS2 nor DBS4, however, was clearly correlated with an SBS or ID signature that correlates with age. ID1, ID2, ID5 and ID8 showed correlations with age in multiple tissues. ID1 and ID2 indels are likely due to slippage at poly T repeats during DNA replication and correlated with the number of SBS1 substitutions. SBS1 has previously been proposed to reflect the number of mitoses a cell has experienced and thus SBS1, ID1 and ID2 may all be generated during DNA replication at mitosis⁴. The number of ID5 mutations correlated with the number of SBS40 mutations and thus the mutational processes underlying these two age correlated signatures may also harbour common components. ID8 is predominantly composed of deletions >5bp with no or 1bp of microhomology at their boundaries. These are likely due to DNA double strand breaks which have not been repaired by homologous recombination based mechanisms, but instead by a non-homologous-end joining mechanism. The features of ID8 resemble those of some

ionising radiation associated mutations and this may, therefore, be an underlying aetiological factor⁵⁵. Taken together, the results indicate that multiple mutational processes operate in normal cells.

DISCUSSION

Cancers arise as a result of somatic mutations. Mutational signature analysis therefore provides important insights into cancer development through comprehensive characterisation of the underlying mutational processes. There are, however, important constraints, limitations and assumptions in the analytic frameworks we have used that should be recognised. Although designed to reflect the mutational consequences of recurrent mutational processes, mutational signatures extracted from sample sets in which multiple mutational processes are operative remain mathematical approximations, with profiles that can be influenced by the mathematical approach used and by additional factors, such as the other mutational processes present. For conceptual and practical simplicity we have assumed that there is a single signature associated with each mutational process and have provided an average, reference signature to represent it. However, we do not discount the possibility that further nuances and variations of signature profiles exist, for example between different tissues. Moreover, although the extent of separation between partially correlated signatures has been improved in this analysis, some signatures may still represent combinations of constituent signatures. Contributions from each signature to the burden of mutations in each sample have been estimated. However, with increasing numbers of signatures and multiple orders of magnitude differences in mutation burdens from certain signatures, prior knowledge can help to avoid biologically implausible results. Thus further development of methods for deciphering mutational signatures and attribution of mutations is warranted and this needs to be supplemented by signatures derived from experimental systems in which the causes of the mutations are known. The numbers of DBSs, clustered substitutions, IDs and genome rearrangements (reported in ³⁰) are small compared to single base substitutions. Thus, larger datasets may be required to robustly characterise their mutational signatures. Nevertheless, the results outlined here indicate that signatures with many similarities and some differences can be found by different mathematical approaches, and that these are confirmed in many different ways, including experimentally elucidated signatures^{22,31,38,41,42,54,56-62} and the observation of tumors dominated by a single signature (<https://www.synapse.org/#!/Synapse:syn12016215>)

This analysis includes almost all publicly available exome and whole-genome cancer sequences, amounting in aggregate to 23,829 cancers of most cancer types. Some rare or geographically restricted signatures may not have been captured and signatures of therapeutic mutagenic exposures have not been exhaustively explored. Nevertheless, it is likely that a substantial proportion of the naturally-occurring mutational signatures found in human cancer have now been described. This comprehensive repertoire provides a foundation for future research into (i) geographical and temporal differences in cancer incidence to elucidate underlying differences in aetiology, (ii) the mutational processes and signatures present in normal tissues and caused by non-neoplastic disease states, (iii) clinical and public health applications of signatures as indicators of sensitivity to therapeutics and past exposure to mutagens, and (iv) mechanistic understanding of the mutational processes underlying carcinogenesis.

Acknowledgements

The results here are partly based on data generated by the TCGA Research Network (<http://cancergenome.nih.gov/>).

This work was supported by Wellcome grant reference 206194 (M.R.S., S.M., L.B.A.), <G.G. funding>, Singapore National Medical Research Council grant NMRC/CIRG/1422/2015 and the Singapore Ministry of Health via the Duke-NUS Signature Research Programmes (M.N.H., A.W.T.N., A.B., S.G.R.), US National Institute of Health Intramural Research Program Project Z1AES103266 (D.A.G.), the European Research Council Consolidator Grant 682398 (N.L.-B.), US National Cancer Institute U24CA143843 (D.A.W.), and. Cancer Research UK Grand Challenge Award C98/A24032 (L.B.A., M.R.S.). We thank XXX

Figure legends.

Figure 1. Mutation burdens of single base substitutions, doublet base substitutions and small insertions and deletions for the 2,780 PCAWG tumors. Each sample is displayed according to its tumor type. Tumor types are ordered according to the median number of single base substitutions. The numbers of cases of each tumor type are shown. The proportions of each mutation subclass in each sample are shown as coloured bar charts.

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

Figure 3. The number of mutations contributed by each mutational signature to the 2,780 PCAWG tumors. The numbers of mutations attributed are shown by cancer type. The size of each dot represents the proportion of samples of each tumor type that show the mutational signature. The colour of each dot represents the median mutation burden of the signature in samples which show the signature. Contributions are shown for single base substitution, doublet base substitution and indel mutational signatures separately. Contributions of composite signatures to the PCAWG cancers and single base substitution signatures to the complete set of cancer samples analysed are shown in Supplementary information.

Figure 4. Illustrative examples of mutational spectra of individual cancer samples (a breast cancer, a lung cancer and a malignant melanoma) and their contributory single base substitution, doublet base substitution and small insertion and deletion mutational signatures.

Methods (will appear on-line)

Principles and strategy of mutational signature analysis adopted in this report

Conceptual principles.

- Multiple mutational processes generate the somatic mutations present in human cancer.
- Each mutational process generates a particular pattern of somatic mutations known as a mutational signature.
- 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.
- 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.

Aims of the study.

- 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.

Approach used.

- Several mathematical approaches have been used to deconvolute (extract) the mutational signatures present in a set of mutational catalogues. 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.
- 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>).
- In brief: SigProfiler employs an elaboration of previously presented approaches for signature extraction and for attribution of mutation counts to mutational signatures in individual tumours^{3,4,16,18}.
- In brief: SignatureAnalyzer employs a Bayesian variant of NMF^{6,7,35}. 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.
- 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.

Role of NMF in extraction and attribution of mutational signatures.

- 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 tumor samples present in V . For example, this may be especially evident when some tumors in V have high numbers of mutations (eg 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 tumor, 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 tumors 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 tumors that is separate from the process of extracting (discovering) the signatures.
- The use of multiple extractions to support confidence in results:
 - For the bulk of the signatures reported here, for both SigProfiler and SignatureAnalyzer, the main extraction procedure was carried out on (1) the majority of the PCAWG tumors excluding certain highly mutated tumours and (2) the corresponding highly mutated tumours. For SigProfiler this latter set consisted of the melanomas, and for SignatureAnalyzer it consisted of melanomas, microsatellite-unstable tumours, and a single temozolomide-exposed tumour (<https://www.synapse.org/#!Synapse:syn11738314>). The decision to partition the PCAWG tumours was made after exploratory extractions, which included all the PCAWG tumours together, each PCAWG tumor type separately, and other partitions.
 - In addition, SigProfiler extracted signatures from
 - non-PCAWG whole genomes, on each tumor type separately and also on all samples from all tumor types together.

- TCGA exomes, on each tumor type separately and also on all samples from all tumors together.
- non-TCGA exomes, on each tumor type separately and also on all samples from all tumors together.

This allowed the extraction of signatures that were not present in the PCAWG tumors (eg SBS42, which is due to haloalkane exposure and seen only in whole exome data). It also served as an important validation, as extraction of similar signatures from single tumour types and other sample sets supports the correctness of the signature extracted from the PCAWG samples (<https://www.synapse.org/#!Synapse:syn12016215>).

- Signature extraction from each tumour type (or from some other subset of cancers) separately has the advantages of:
 - usually including fewer (and different) mutational signatures in each tumour type sample set than in the set of all cancers together and thus fewer (and different) opportunities for inter-signature interference.
 - allowing multiple independent opportunities for extraction of a signature that is present in multiple tumour types, and thus of obtaining validation/confirmation of the signature's existence and profile.
 - allowing extraction of a signature that may (for a number of reasons) fail to be extracted in analysis of all tumour types together.
 - providing primary evidence for the existence of the signature in each tumour type.
 - allowing separation of highly mutated cancer types/samples from cancer types/samples with low mutation burdens.
- Signature extraction from multiple tumour types together has the advantages of:
 - usually including more samples with a particular signature than in each individual cancer type and thus being better powered to separate a signature from other partially correlated signatures and/or from signatures with similar profiles.
 - providing a single profile for a signature rather than the multiple slightly different profiles which emerge from extraction of each tumour type separately.
- The profiles of the mutational signatures extracted from cancer are highly variable. They range from some that have contributions from mutations of all subtypes in the mutation classification ("flat" or "featureless" signatures eg SBS5 and SBS40) to others that are essentially defined by mutations at only one (or a small number) of the mutation subtypes (eg signatures SBS2, SBS13, SBS10a and SBS10b). There appears to be less concordance between the results of SigProfiler and SignatureAnalyzer for flat signatures than for signatures with distinct features indicating that generally, these may be more difficult to accurately extract and distinguish from each other. However, there is experimental support for the existence of SBS5 and SBS3^{54,61}.
- We represented each signature as a single reference. This selection of a single reference signature does not exclude the possibility that signature profiles may show nuances and further complexity and may vary in different contexts (eg in different tissues). The rationale for selecting a single reference signature was the view that this

would be a level of granularity useful to most researchers. For those with specialised interests in particular mutational processes and their components, we also provided the signatures extracted from individual tumor types, comprising PCAWG and non-PCAWG genomes and exomes (<https://www.synapse.org/#!Synapse:syn12025142>).

- Attribution of signatures to cancer samples:
 - The reference signatures from SigProfiler and SignatureAnalyzer were used to estimate the number of mutations due to each signature in each tumour (<https://www.synapse.org/#!Synapse:syn11804065>).
 - SigProfiler and SignatureAnalyzer differ in their approaches for attributing signatures. However, both incorporate a set of rules based on prior knowledge and biological plausibility, and incorporate techniques to encourage sparsity in the number of signatures attributed to a given tumor.
 - Sparsity (limiting the numbers of signatures and limiting the numbers of signatures attributed to each cancer sample) is an important concept and feature of both SigProfiler and SignatureAnalyzer (both in signature extraction and attribution). Our prior beliefs are that i) there is a limited set of significantly contributing mutational processes (and hence a limited set of mutational signatures) operating to generate somatic mutations across all cancers and ii) that a limited set of mutational processes contribute to individual cancer genomes (as opposed to all mutational signatures contributing to all samples). Our aim in discovering mutational signatures is to reflect the underlying biological processes and to attribute them appropriately. It is not a mathematical exercise in which the main objective and priority is to minimize the difference between $W \times H$ and the original spectra in V . Indeed, if the latter was the main aim, for 96 mutation classes a set of 96 signatures each constituted entirely of mutations in just one class (and therefore ignoring sparsity), will always provide error free reconstruction but will provide absolutely no information about underlying mutational processes.

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

- The results (signatures and attributions) of the two methods have been presented separately. We have done this in preference to combining them. We have handled the two outputs in this way because we believe that this provides a simpler conceptual and technical basis on which the research community can understand the results, can employ the methods in future and can compare results with those shown in this paper. We also do not have a basis for believing that a combined/averaged/overlapping single result set is a better representation of the natural truth than either of the two result sets individually and do not have a well-founded and simple technical approach for combining them. We have, however, provided comparisons of the outputs.
- For brevity and for continuity with previous publications, the results from SigProfiler, a further elaborated version of previously described approaches^{3,4,16,18} that generated the 30 signatures previously shown in COSMIC³³, are shown in the main manuscript, and the results from SignatureAnalyzer in supplementary data (<https://www.synapse.org/#!Synapse:syn11738307>).
- Nomenclature of signatures is based on and extends the nomenclature previously used in COSMIC³³.

- Both methods analysed each mutation type (SBSs, DBSs and IDs) separately and also together as a composite signature. In future, however, SigProfiler will usually use the separately extracted single base substitution, indel and doublet base substitution signatures as its standard. This generally facilitates portability, and comparison of signature profiles with those from a variety of sample sets including targeted sequences, exomes etc.
- SBS signatures reported in Supplementary Data include possible artefacts (<https://www.synapse.org/#!Synapse:syn12009783> and see below).

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

- Sequencing artefacts and differences in analysis pipelines can also generate mutational signatures. We have annotated which signatures are likely real or “possible artefact”.
- There are multiple reasons for believing a signature reflects a biological mutational signature rather than an artefact.
 - The input data supporting the signature seem correct: key mutational features of the putative signature look real in a mapped-read browser such as IGV, or characteristic mutations are experimentally confirmed in the tumor and normal samples. Inspection in a mapped read browser is especially important in checking for possible problems in potentially new signatures arising in datasets other than the highly scrutinized and checked PCAWG and TCGA sets. Features associated with experimental, mapping, or other computational artifacts include strong preference for the first read, very low variant allele fractions, variants in regions of low germ-line sequencing coverage, variants found near indels in low-complexity regions, variants from a signature only found in one sequencing center etc.
 - The 96-mutation profile and additional features (eg strand asymmetry, association with replication timing), are known to result from a particular process in experimental systems. Examples: UV, polymerase epsilon proofreading deficiency, aristolochic acid and cisplatin exposure.
 - The putative signature is broadly consistent with previous biochemical knowledge of mutational processes (eg preference for G adducts in aflatoxin).
 - The putative signature dominates the spectra of some tumors (column J of <https://www.synapse.org/#!Synapse:syn12016215>).
 - The putative mutational signature is consistently deciphered from multiple independent datasets; if so it is either a common sequencing artefact or something real.
 - The putative signature correlates with known or suspected mutational exposures, endogenous processes, or repair defects, especially if some of those exposures/processes/repair defects result in overwhelming mutational spectra. Examples: melanoma / fair skin / UV exposure, POLE mutations, MMR deficiency and APOBEC germ line variants.
 - The putative signature correlates with other clinical characteristics, such as age at diagnosis (examples SBS1 and SBS5) or tobacco smoking (SBS4).
 - The mutational signature exhibits a strong transcriptional strand bias; it is hard to imagine an artefact with transcriptional strand bias.

- The putative signature shows association with other genomic features, such as microindels in homopolymers, replication strand, replication timing, or nucleosome occupancy.

Cancer sample sets on which different analyses have been conducted.

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

Splitting of mutational signatures.

- Certain previously existing single signatures have split into multiple constituent signatures in this analysis. This is likely due to the existence of multiple, partially correlated mutational processes with the same initiating factor (for example, UV exposure) but subsequent differences in underlying mechanisms which differ in intensity in different tissues or other contexts. A previous example of this for which we have allocated different signature numbers is the split of the usually co-occurring but independently varying consequences of APOBEC mutagenesis into signatures SBS2 and SBS13 (<https://www.synapse.org/#!Synapse:syn12009783>).
- Depending on the extent of correlation of the two signatures, and the available dataset/statistical power such signatures may manifest as a single signature, overlapping partially separated signatures or as two separate signatures.
- We are aware that splitting of signatures can also be a mathematical artefact. However, we have used multiple extractions to confirm and validate signature splits and applied the principle of sparsity to limit artefactual splits (<https://www.synapse.org/#!Synapse:syn12009783>).

SigProfiler overview

SigProfiler incorporates two distinct steps for identification of mutational signatures based on the previously described methodology^{3,4,16,18}. The first step, SigProfilerExtraction, encompasses a hierarchical *de novo* extraction of mutational signatures based on somatic mutations and their immediate sequence context, while the second step,

SigProfilerAttribution, focuses on accurately estimating the number of somatic mutations associated with each extracted mutational signature in each sample.

SigProfilerExtraction

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

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

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

Transcriptional strand bias associated with mutational signatures was assessed by applying SigProfilerExtraction to catalogs of in-transcript mutations that capture strand information (192 mutations classes, <https://www.synapse.org/#!Synapse:syn12026195>). These 192-class signatures were collapsed to strand-invariant 96-class signatures and compared to the signatures extracted from the 96-class data, revealing very high cosine similarities (median 0.9, column F in <https://www.synapse.org/#!Synapse:syn12016215>).

SigProfilerAttribution (single sample attribution)

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

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

In equation (1), \vec{C} and each \vec{S}_r are vectors of K nonnegative components reflecting, respectively, the mutational spectrum of a sample and the r -th reference mutational signature. All mutational signatures, \vec{S}_r , were identified in the SigProfilerExtraction step. Each E_r is unknown scalar reflecting the number of mutations contributed by signature \vec{S}_r in the mutational spectrum \vec{C} . The minimization of equation (1) is always performed under two additional constraints: (i) $E_r \geq 0$ and (ii) $\|\vec{C}\|_1 \geq E_r$; The constrained minimization of equation (1) is performed using a nonlinear convex optimization programming solver using the interior-point algorithm⁶³.

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

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

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

SignatureAnalyzer overview

SignatureAnalyzer employs a Bayesian variant of NMF that infers the number of signatures through the automatic relevance determination technique and delivers highly interpretable and sparse representations for both signature profiles and attributions that strike a balance between data fitting and model complexity. Please see references^{6,7,35} for details.

SignatureAnalyzer signature extraction

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

In step 1 of the two-step extraction process, global signature extraction was performed for the low mutation burden samples ($n = 2,624$). These excluded hyper-mutated tumors: those with putative polymerase epsilon (POLE) defects or mismatch repair defects (microsatellite

instable tumors - MSI), skin tumours (which had intense UV mutagenesis), and one tumour with temozolomide (TMZ) exposure. In step 2, additional signatures unique to hyper-mutated samples were extracted while allowing all signatures found in the low mutation burden samples to explain some of the spectra of hyper-mutated samples. This approach was designed to minimize a well-known "signature bleeding" effect or a bias of hyper- or ultra-mutated samples on the signature extraction. In addition, this approach provided information about which signatures are unique to the hyper-mutated samples which is later used when attributing signatures to samples.

SignatureAnalyzer signature attribution

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

Data Availability

Data are available at <https://www.synapse.org/#!Synapse:syn11726601/wiki/513478>.

Code Availability

SigProfiler code is available at

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

SignatureAnalyzer code is available at <https://www.synapse.org/#!Synapse:syn11801492>.

References

- 1 Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. *Nature* **458**, 719-724, doi:10.1038/nature07943 (2009).
- 2 Alexandrov, L. B. & Stratton, M. R. Mutational signatures: the patterns of somatic mutations hidden in cancer genomes. *Curr Opin Genet Dev* **24**, 52-60, doi:10.1016/j.gde.2013.11.014 (2014).
- 3 Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M. R. Deciphering signatures of mutational processes operative in human cancer. *Cell Rep* **3**, 246-259, doi:10.1016/j.celrep.2012.12.008 (2013).
- 4 Alexandrov, L. B. *et al.* Clock-like mutational processes in human somatic cells. *Nat Genet* **47**, 1402-1407, doi:10.1038/ng.3441 (2015).
- 5 Morganella, S. *et al.* The topography of mutational processes in breast cancer genomes. *Nat Commun* **7**, 11383, doi:10.1038/ncomms11383 (2016).
- 6 Kasar, S. *et al.* Whole-genome sequencing reveals activation-induced cytidine deaminase signatures during indolent chronic lymphocytic leukaemia evolution. *Nature communications* **6**, 8866 (2015).
- 7 Kim, J. *et al.* Somatic ERCC2 mutations are associated with a distinct genomic signature in urothelial tumors. *Nat Genet* **48**, 600-606, doi:10.1038/ng.3557 (2016).
- 8 Fischer, A., Illingworth, C. J., Campbell, P. J. & Mustonen, V. EMu: probabilistic inference of mutational processes and their localization in the cancer genome. *Genome Biol* **14**, R39, doi:10.1186/gb-2013-14-4-r39 (2013).
- 9 Roberts, N. hdp (hierarchical Dirichlet process) R package. <https://github.com/nicolaroberts/hdp> (2015).
- 10 Gehring, J. S., Fischer, B., Lawrence, M. & Huber, W. SomaticSignatures: inferring mutational signatures from single-nucleotide variants. *Bioinformatics* **31**, 3673-3675 (2015).
- 11 Blokzijl, F., Janssen, R., Van Boxtel, R. & Cuppen, E. MutationalPatterns: an integrative R package for studying patterns in base substitution catalogues. *bioRxiv*, doi:10.1101/071761 (2016).
- 12 Shiraishi, Y., Tremmel, G., Miyano, S. & Stephens, M. A simple model-based approach to inferring and visualizing cancer mutation signatures. *PLoS genetics* **11**, e1005657 (2015).
- 13 Rosales, R. A., Drummond, R. D., Valieris, R., Dias-Neto, E. & da Silva, I. T. signeR: an empirical Bayesian approach to mutational signature discovery. *Bioinformatics* **33**, 8-16, doi:10.1093/bioinformatics/btw572 (2017).
- 14 Ardin, M. *et al.* MutSpec: a Galaxy toolbox for streamlined analyses of somatic mutation spectra in human and mouse cancer genomes. *BMC bioinformatics* **17**, 170 (2016).
- 15 Funnell, T. *et al.* Integrated single-nucleotide and structural variation signatures of DNA-repair deficient human cancers. *bioRxiv*, doi:10.1101/267500 (2018).
- 16 Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415-421, doi:10.1038/nature12477 (2013).
- 17 Nik-Zainal, S. *et al.* Mutational processes molding the genomes of 21 breast cancers. *Cell* **149**, 979-993, doi:10.1016/j.cell.2012.04.024 (2012).
- 18 Nik-Zainal, S. *et al.* Landscape of somatic mutations in 560 breast cancer whole-genome sequences. *Nature* **534**, 47-54, doi:10.1038/nature17676 (2016).

- 19 Petljak, M. & Alexandrov, L. B. Understanding mutagenesis through delineation of mutational signatures in human cancer. *Carcinogenesis* **37**, 531-540, doi:10.1093/carcin/bgw055 (2016).
- 20 Schulze, K. *et al.* Exome sequencing of hepatocellular carcinomas identifies new mutational signatures and potential therapeutic targets. *Nat Genet* **47**, 505-511, doi:10.1038/ng.3252 (2015).
- 21 Poon, S. L. *et al.* Mutation signatures implicate aristolochic acid in bladder cancer development. *Genome Med* **7**, 38, doi:10.1186/s13073-015-0161-3 (2015).
- 22 Poon, S. L. *et al.* Genome-wide mutational signatures of aristolochic acid and its application as a screening tool. *Sci Transl Med* **5**, 197ra101, doi:10.1126/scitranslmed.3006086 (2013).
- 23 Alexandrov, L. B. Understanding the origins of human cancer. *Science* **350**, 1175, doi:10.1126/science.aad7363 (2015).
- 24 Hayward, N. K. *et al.* Whole-genome landscapes of major melanoma subtypes. *Nature* **545**, 175-180, doi:10.1038/nature22071 (2017).
- 25 Polak, P. *et al.* A mutational signature reveals alterations underlying deficient homologous recombination repair in breast cancer. *Nat Genet* **49**, 1476-1486, doi:10.1038/ng.3934 (2017).
- 26 Merlevede, J. *et al.* Mutation allele burden remains unchanged in chronic myelomonocytic leukaemia responding to hypomethylating agents. *Nat Commun* **7**, 10767, doi:10.1038/ncomms10767 (2016).
- 27 Mimaki, S. *et al.* Hypermutation and unique mutational signatures of occupational cholangiocarcinoma in printing workers exposed to haloalkanes. *Carcinogenesis* **37**, 817-826, doi:10.1093/carcin/bgw066 (2016).
- 28 Drier, Y. *et al.* Somatic rearrangements across cancer reveal classes of samples with distinct patterns of DNA breakage and rearrangement-induced hypermutability. *Genome Res* **23**, 228-235, doi:10.1101/gr.141382.112 (2013).
- 29 Stephens, P. J. *et al.* Massive genomic rearrangement acquired in a single catastrophic event during cancer development. *Cell* **144**, 27-40, doi:10.1016/j.cell.2010.11.055 (2011).
- 30 Li, Y. *et al.* Patterns of structural variation in human cancer. *bioRxiv*, doi:10.1101/181339 (2017).
- 31 Meier, B. *et al.* *C. elegans* whole-genome sequencing reveals mutational signatures related to carcinogens and DNA repair deficiency. *Genome Res* **24**, 1624-1636, doi:10.1101/gr.175547.114 (2014).
- 32 Supek, F. & Lehner, B. Clustered Mutation Signatures Reveal that Error-Prone DNA Repair Targets Mutations to Active Genes. *Cell* **170**, 534-547 e523, doi:10.1016/j.cell.2017.07.003 (2017).
- 33 Wellcome Trust Sanger Institute. *COSMIC, Catalog of Somatic Mutations in Cancer - Signatures of Mutational Processes in Human Cancer*, <<http://cancer.sanger.ac.uk/cosmic/signatures>> (2017).
- 34 Forbes, S. A. *et al.* The Catalogue of Somatic Mutations in Cancer (COSMIC). *Curr Protoc Hum Genet* **Chapter 10**, Unit 10 11, doi:10.1002/0471142905.hg1011s57 (2008).
- 35 Tan, V. Y. & Févotte, C. Automatic relevance determination in nonnegative matrix factorization with the/spl beta/-divergence. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **35**, 1592-1605 (2013).

- 36 Haradhvala, N. J. *et al.* Distinct mutational signatures characterize concurrent loss of polymerase proofreading and mismatch repair. *Nature Communications* **9**, 1746, doi:10.1038/s41467-018-04002-4 (2018).
- 37 Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* **499**, 214, doi:10.1038/nature12213 (2013).
- 38 Boot, A. *et al.* In-depth characterization of the cisplatin mutational signature in human cell lines and in esophageal and liver tumors. *Genome Res* **In press** (2018).
- 39 Viel, A. *et al.* A Specific Mutational Signature Associated with DNA 8-Oxoguanine Persistence in MUTYH-defective Colorectal Cancer. *EBioMedicine* **20**, 39-49, doi:10.1016/j.ebiom.2017.04.022 (2017).
- 40 Pilati, C. *et al.* Mutational signature analysis identifies MUTYH deficiency in colorectal cancers and adrenocortical carcinomas. *J Pathol* **242**, 10-15, doi:10.1002/path.4880 (2017).
- 41 Drost, J. *et al.* Use of CRISPR-modified human stem cell organoids to study the origin of mutational signatures in cancer. *Science* **358**, 234-238 (2017).
- 42 Saini, N. *et al.* The Impact of Environmental and Endogenous Damage on Somatic Mutation Load in Human Skin Fibroblasts. *PLoS Genetics* **12**, e1006385, doi:10.1371/journal.pgen.1006385 (2016).
- 43 Haradhvala, N. J. *et al.* Mutational Strand Asymmetries in Cancer Genomes Reveal Mechanisms of DNA Damage and Repair. *Cell* **164**, 538-549, doi:10.1016/j.cell.2015.12.050 (2016).
- 44 Chan, K. *et al.* An APOBEC3A hypermutation signature is distinguishable from the signature of background mutagenesis by APOBEC3B in human cancers. *Nat Genet* **47**, 1067-1072, doi:10.1038/ng.3378 (2015).
- 45 Roberts, S. A. *et al.* Clustered mutations in yeast and in human cancers can arise from damaged long single-strand DNA regions. *Molecular cell* **46**, 424-435 (2012).
- 46 Kasar, S. & Brown, J. R. Mutational landscape and underlying mutational processes in chronic lymphocytic leukemia. *Mol Cell Oncol* **3**, e1157667, doi:10.1080/23723556.2016.1157667 (2016).
- 47 Brash, D. E. UV Signature Mutations. *Photochemistry and Photobiology* **91**, 15-26, doi:doi:10.1111/php.12377 (2015).
- 48 Hill, K. A., Wang, J., Farwell, K. D. & Sommer, S. S. Spontaneous tandem-base mutations (TBM) show dramatic tissue, age, pattern and spectrum specificity. *Mutat Res* **534**, 173-186 (2003).
- 49 Matsuda, T., Kawanishi, M., Yagi, T., Matsui, S. & Takebe, H. Specific tandem GG to TT base substitutions induced by acetaldehyde are due to intra-strand crosslinks between adjacent guanine bases. *Nucleic Acids Res* **26**, 1769-1774 (1998).
- 50 Garaycochea, J. I. *et al.* Alcohol and endogenous aldehydes damage chromosomes and mutate stem cells. *Nature* (2018).
- 51 Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat Genet* **48**, 126-133, doi:10.1038/ng.3469 (2016).
- 52 Pfeifer, G. P. Formation and processing of UV photoproducts: effects of DNA sequence and chromatin environment. *Photochemistry and photobiology* **65**, 270-283 (1997).
- 53 Ceccaldi, R., Rondinelli, B. & D'Andrea, A. D. Repair pathway choices and consequences at the double-strand break. *Trends in cell biology* **26**, 52-64 (2016).
- 54 Blokzijl, F. *et al.* Tissue-specific mutation accumulation in human adult stem cells during life. *Nature* **538**, 260-264, doi:10.1038/nature19768 (2016).

- 55 Behjati, S. *et al.* Mutational signatures of ionizing radiation in second malignancies. *Nature communications* **7**, 12605 (2016).
- 56 Huang, M. N. *et al.* Genome-Scale Mutational Signatures of Aflatoxin In Cells, Mice And Human Tumors. *Genome Research* **27**, 1475-1486, doi:10.1101/gr.220038.116 (2017).
- 57 Nik-Zainal, S. *et al.* The genome as a record of environmental exposure. *Mutagenesis* **30**, 763-770, doi:10.1093/mutage/gev073 (2015).
- 58 Olivier, M. *et al.* Modelling mutational landscapes of human cancers in vitro. *Sci Rep* **4**, 4482, doi:10.1038/srep04482 (2014).
- 59 Szikriszt, B. *et al.* A comprehensive survey of the mutagenic impact of common cancer cytotoxics. *Genome biology* **17**, 99 (2016).
- 60 Zhivagui, M. *et al.* Experimental analysis of exome-scale mutational signature of glycidamide, the reactive metabolite of acrylamide. *bioRxiv*, 254664 (2018).
- 61 Zámboorszky, J. *et al.* Loss of BRCA1 or BRCA2 markedly increases the rate of base substitution mutagenesis and has distinct effects on genomic deletions. *Oncogene* **36**, 746 (2017).
- 62 Zou, X. *et al.* Validating the concept of mutational signatures with isogenic cell models. *Nature communications* **9**, 1744 (2018).
- 63 Byrd, R. H., Hribar, M. E. & Nocedal, J. An interior point algorithm for large-scale nonlinear programming. *SIAM Journal on Optimization* **9**, 877-900 (1999).

Extended Data Figure and Table Legends

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

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

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

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

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

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

Extended Data Figure 7. Mutational signatures extracted from the composite feature set consisting of SBSs in pentanucleotide context, DBSs, and IDs. For each of the four composite

mutational signatures shown, the top panel is the SBS signature collapsed to 96 SBS classes, the middle panel is the co-extracted DBS signature, and the lower panel is the co-extracted ID signature. Note the similarities between the DBS portion of Composite 4 and DBS2, between the ID portion of Composite 4 and ID3, and other similarities noted in the figure.

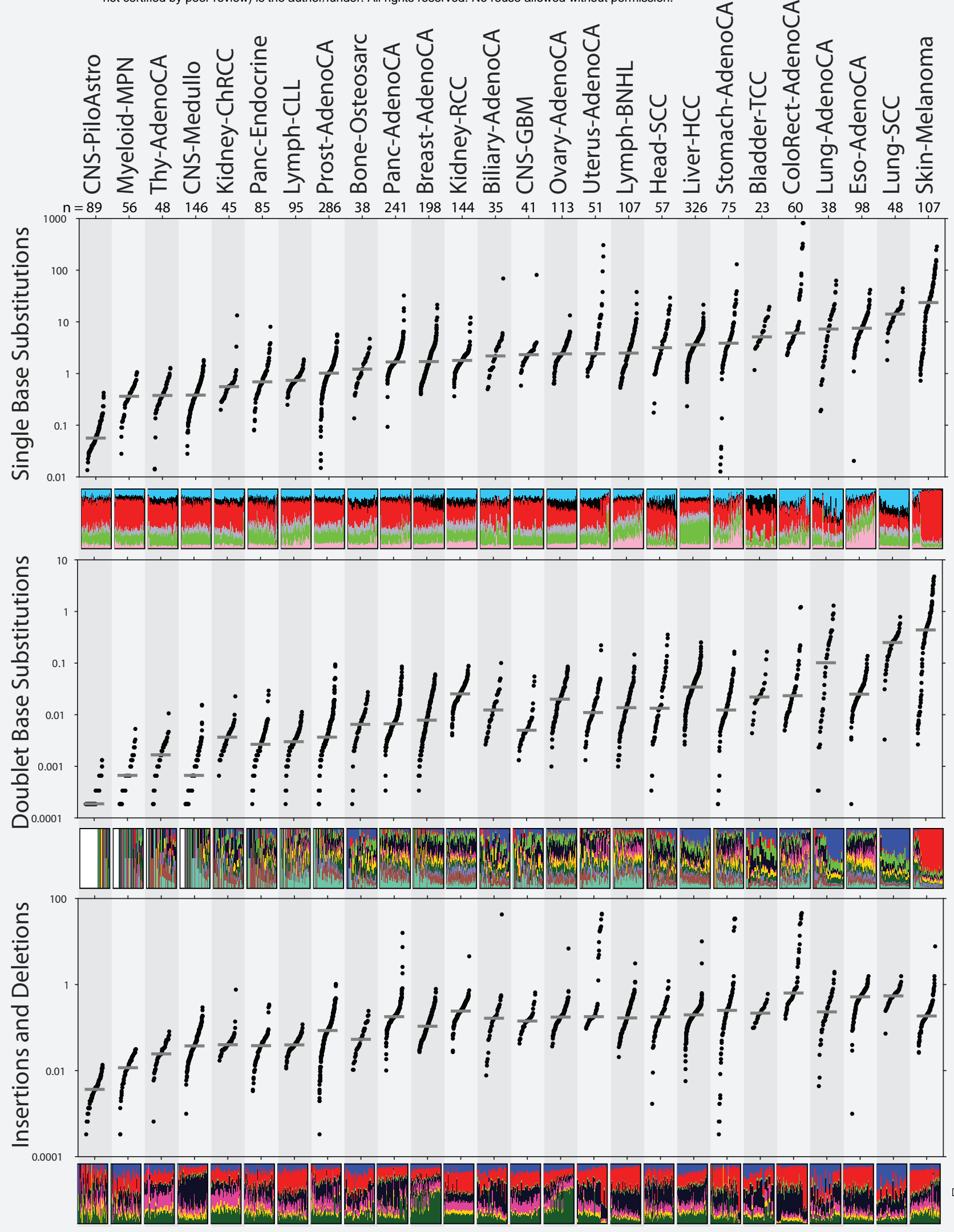
Extended Data Figure 8. SigProfiler signature extraction **(a)** and attribution **(b)**. See Methods for description.

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), R function call:

```
glm(DBS.counts ~ SBS.counts + Cancer.Types)
```

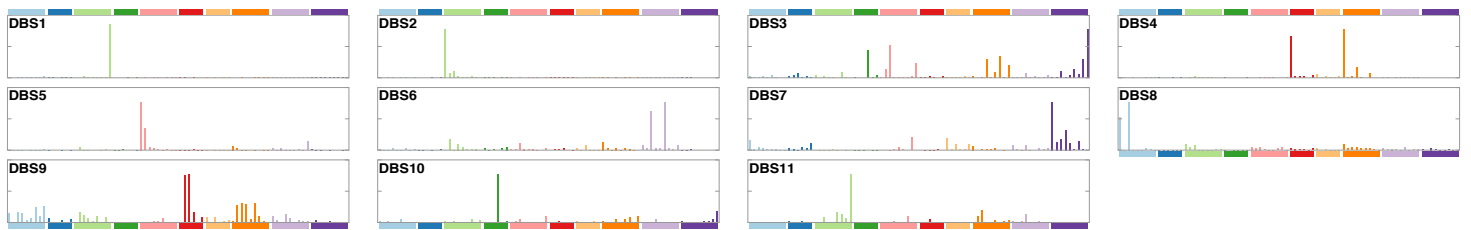
Extended Data Table 2. Numbers of insertion/deletion mutations due to ID1, ID2, and all other ID signatures in hypermutators and non-hypermutators.

Mutations per Megabase





Doublet Base Substitution



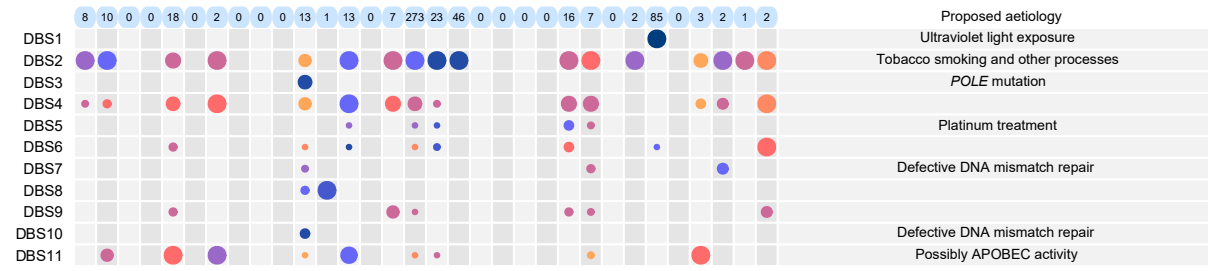
Insertion and Deletion



Single Base Substitution



Doublet Base Substitution



Insertion and Deletion

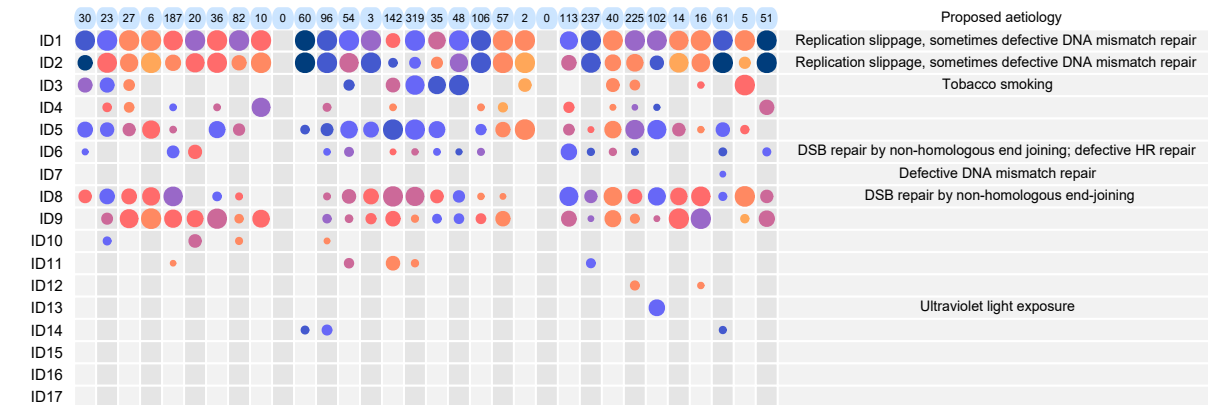
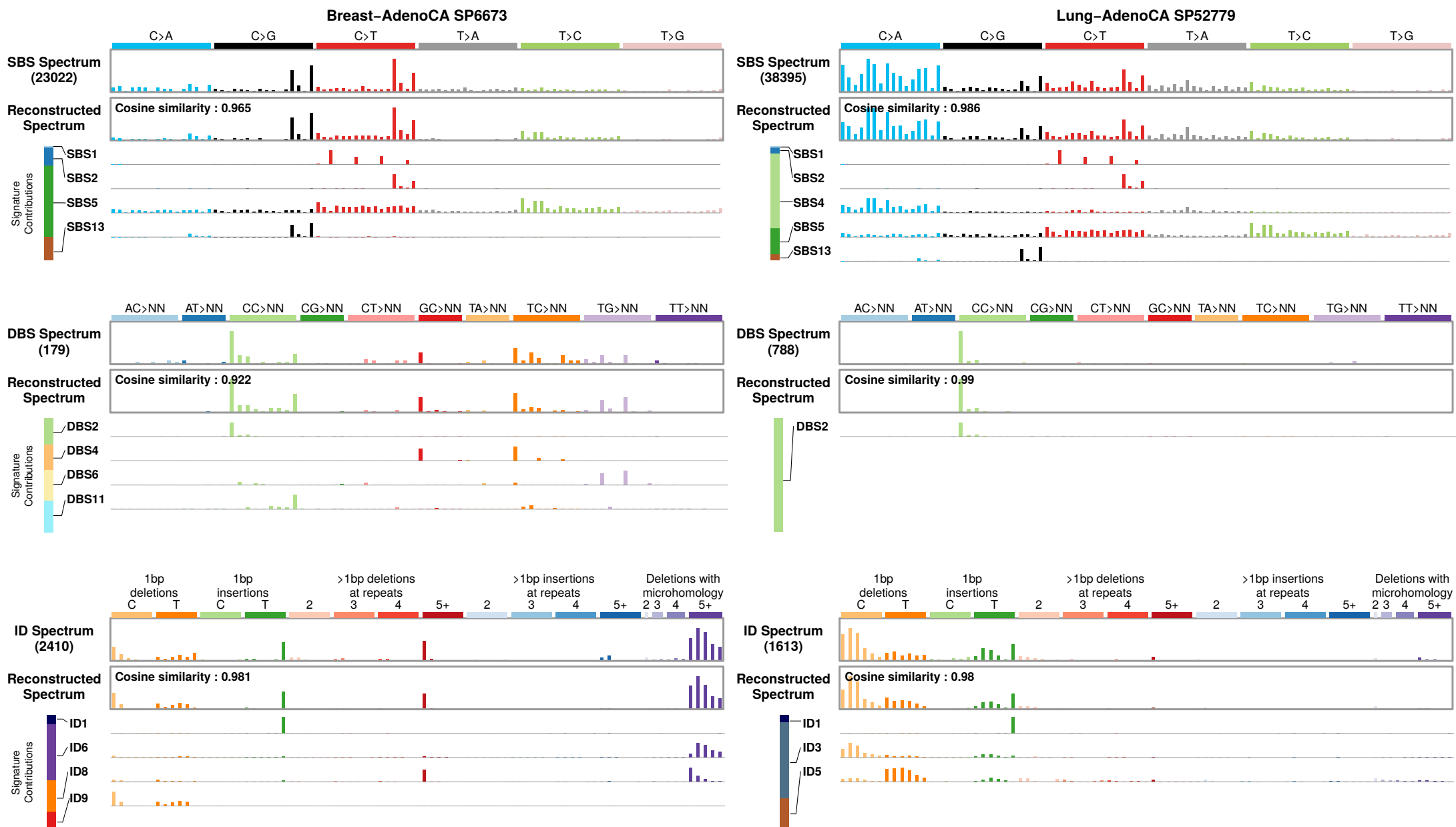
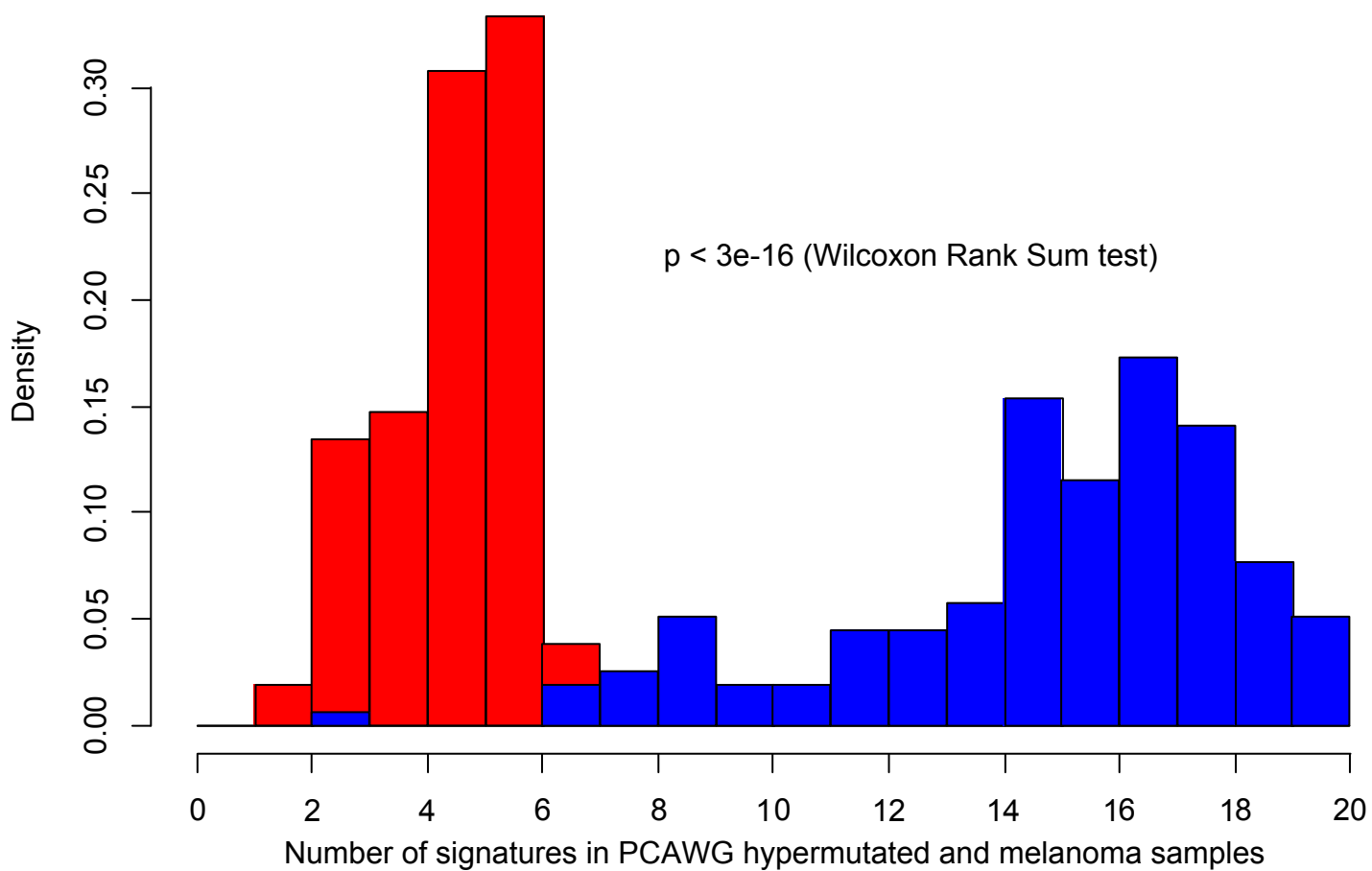
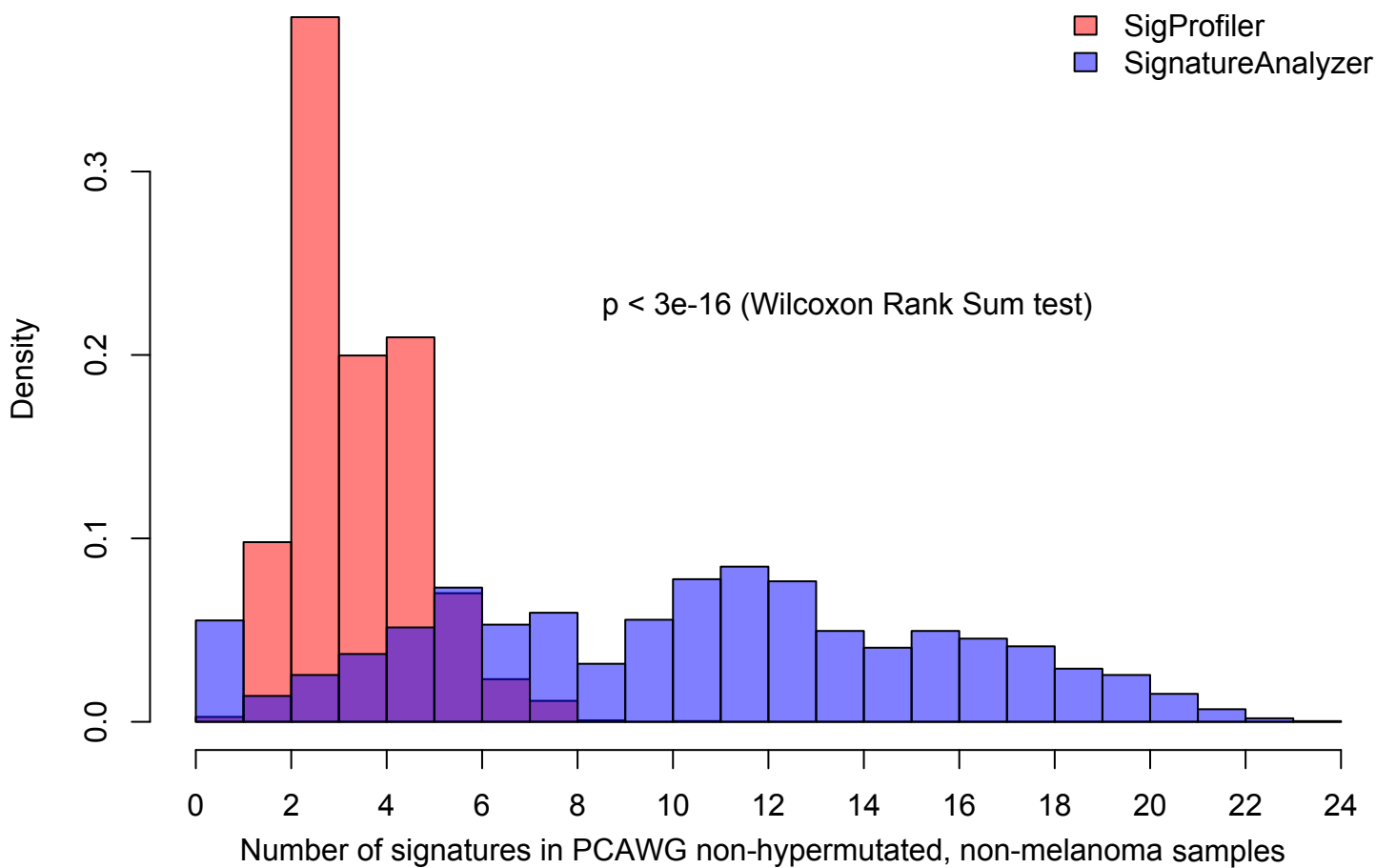


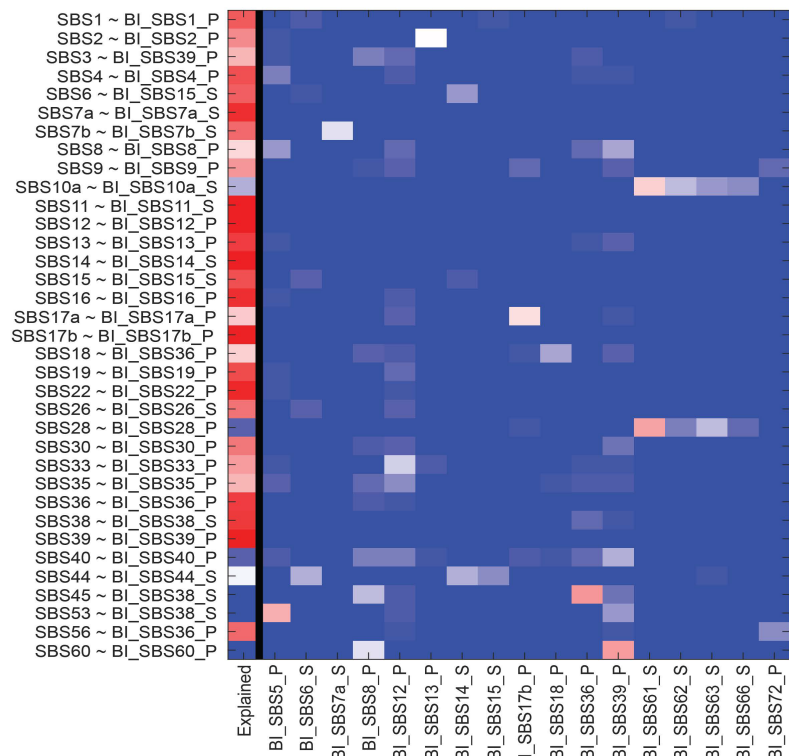
Fig 4



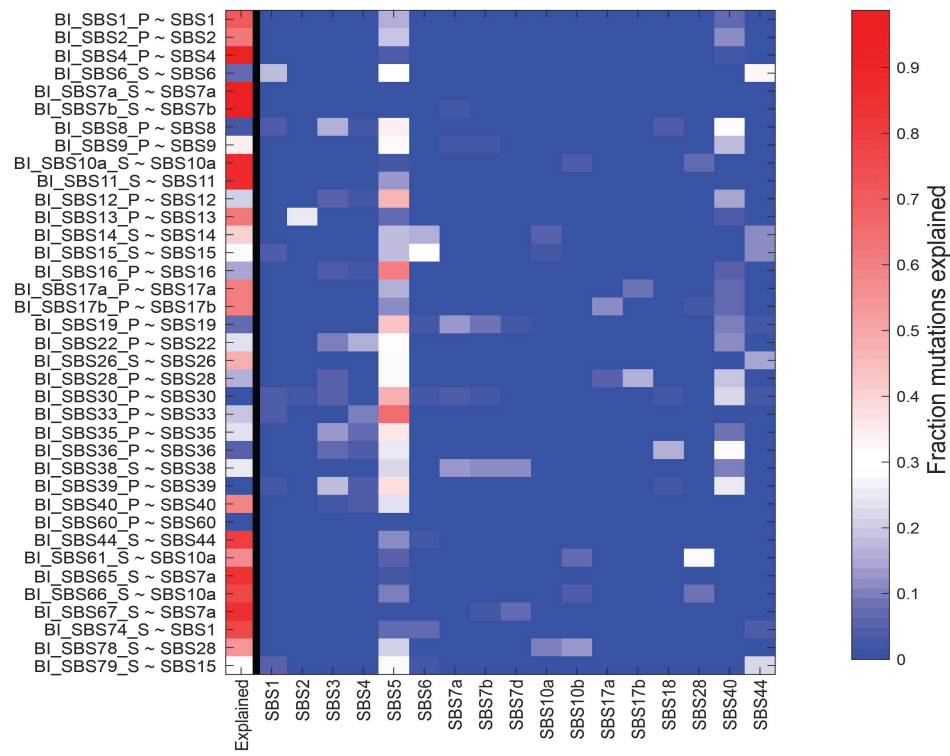
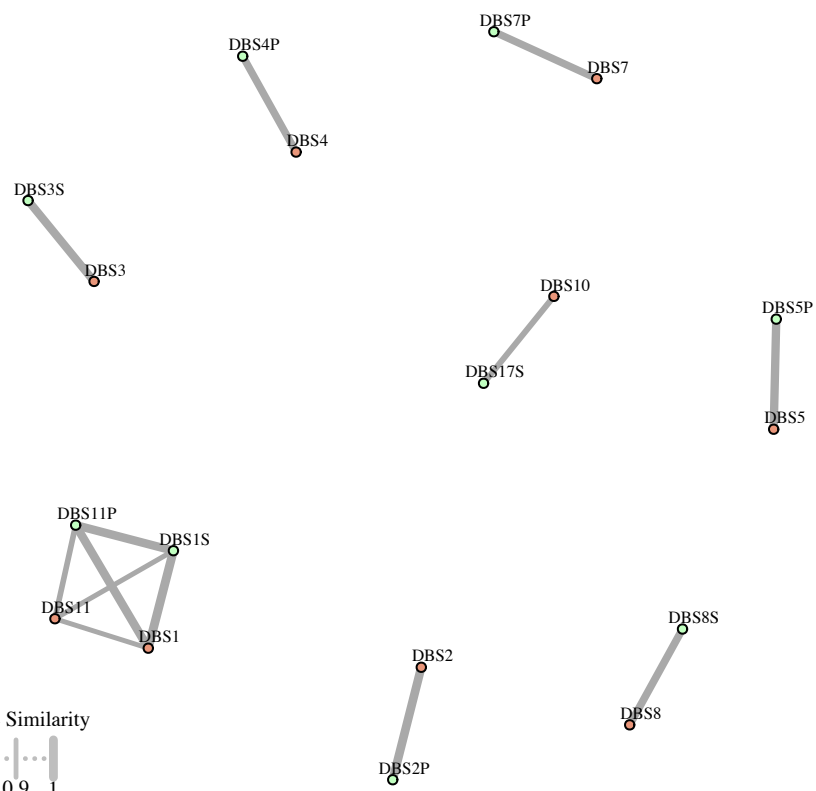
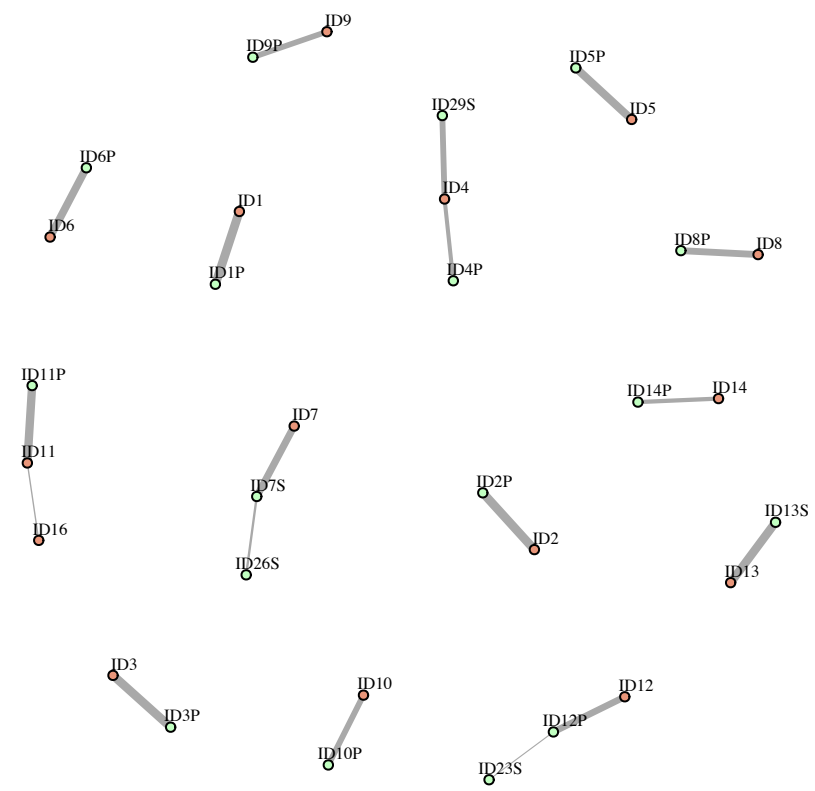


a

SigProfiler signatures explained by SignatureAnalyzer

**b**

SignatureAnalyzer signatures explained by sigprofiler

**c****d**

Isolated Signatures

SigProfiler DBS

● DBS6 ● DBS9

SignatureAnalyzer DBS

● DBS12 ● DBS13S ● DBS14S ● DBS15S ● DBS16S ● DBS18P

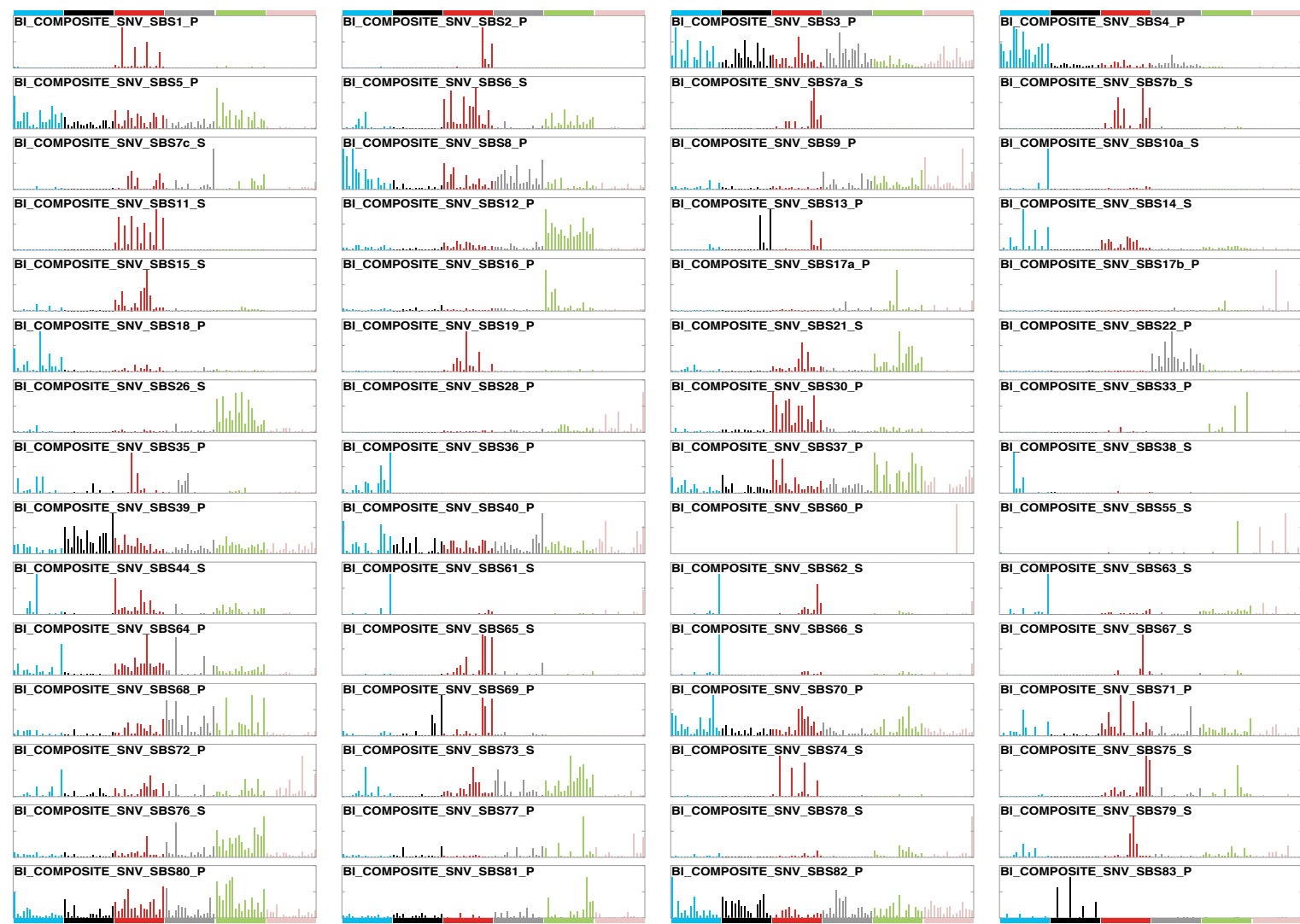
SigProfiler ID

● ID15 ● ID17

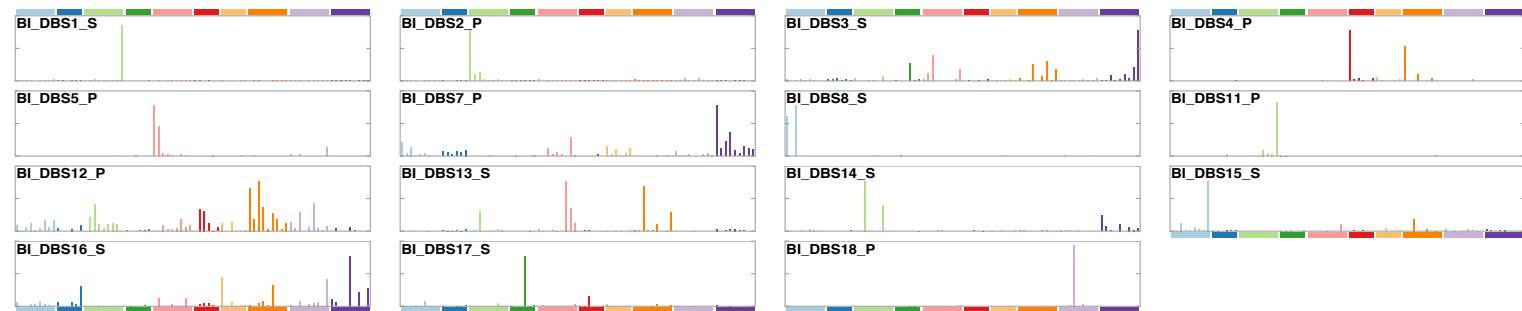
SignatureAnalyzer ID

● ID18S ● ID19S ● ID20S ● ID21S ● ID22S ● ID24S ● ID25S ● ID27S ● ID28S ● ID30S ● ID31P ● ID32S

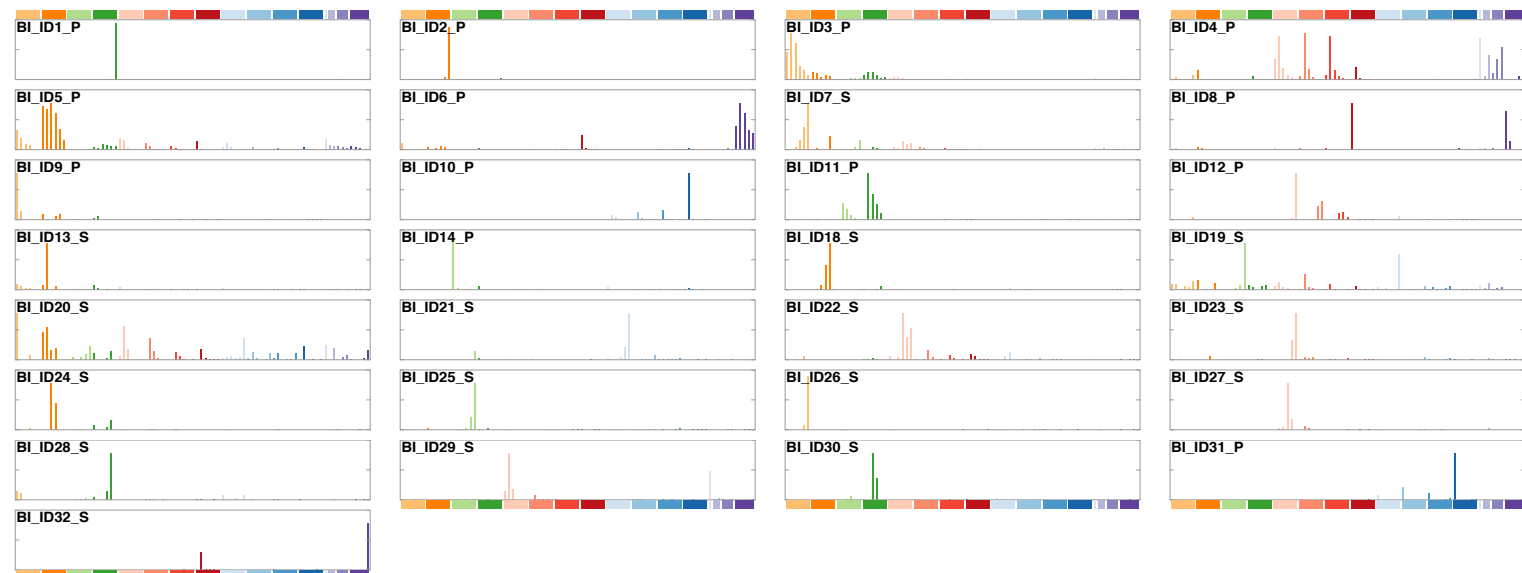
SignatureAnalyzer reference SBS signatures



SignatureAnalyzer reference DBS signatures



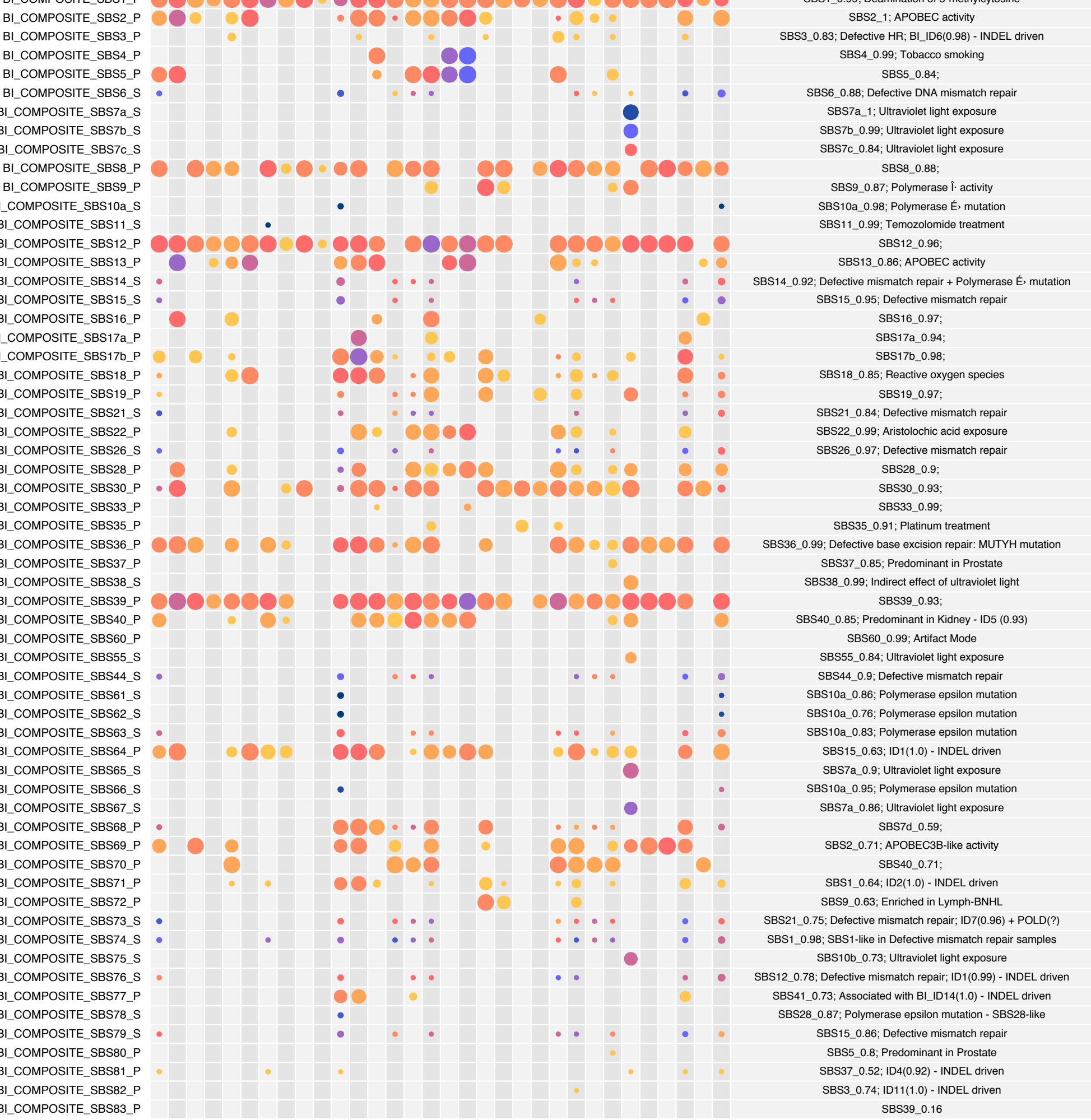
SignatureAnalyzer reference ID signatures



Extended Data Fig 4

Biliary-AdenoCA
Bladder-TCC
Bone-Osteosarc
Bone-Other
Breast
Cervix
CNS-GBM
CNS-Medullo
CNS-Oligo
CNS-PiloAstro
ColoRect-AdenoCA
Eso-AdenoCA
Head-SCC
Kidney-ChrRCC
Kidney-RCC
Liver-HCC
Lung-AdenoCA
Lung-SCC
Lymph-BNHL
Lymph-CLL
Myeloid-AML
Myeloid-MDSMPN
Ovary-AdenoCA
Panc-AdenoCA
Panc-AdenoCA
Prost-AdenoCA
Skin-Melanoma
SoftTissue-Leiomyo
SoftTissue-Liposarc
Stomach-AdenoCA
Thy-AdenoCA
Uterus-AdenoCA

bioRxiv preprint doi: <https://doi.org/10.1101/328550>; this version posted May 15, 2018. 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.



Proposed etiology

SBS1_0.99; Deamination of 5-methylcytosine

SBS2_1; APOBEC activity

SBS3_0.83; Defective HR; BI_ID6(0.98) - INDEL driven

SBS4_0.99; Tobacco smoking

SBS5_0.84;

SBS6_0.88; Defective DNA mismatch repair

SBS7a_1; Ultraviolet light exposure

SBS7b_0.99; Ultraviolet light exposure

SBS7c_0.84; Ultraviolet light exposure

SBS8_0.88;

SBS9_0.87; Polymerase δ activity

SBS10a_0.98; Polymerase ϵ mutation

SBS11_0.99; Temozolomide treatment

SBS12_0.96;

SBS13_0.86; APOBEC activity

SBS14_0.92; Defective mismatch repair + Polymerase ϵ mutation

SBS15_0.95; Defective mismatch repair

SBS16_0.97;

SBS17a_0.94;

SBS17b_0.98;

SBS18_0.85; Reactive oxygen species

SBS19_0.97;

SBS21_0.84; Defective mismatch repair

SBS22_0.99; Aristolochic acid exposure

SBS26_0.97; Defective mismatch repair

SBS28_0.9;

SBS30_0.93;

SBS33_0.99;

SBS35_0.91; Platinum treatment

SBS36_0.99; Defective base excision repair: MUTYH mutation

SBS37_0.85; Predominant in Prostate

SBS38_0.99; Indirect effect of ultraviolet light

SBS39_0.93;

SBS40_0.85; Predominant in Kidney - ID5 (0.93)

SBS60_0.99; Artifact Mode

SBS55_0.84; Ultraviolet light exposure

SBS44_0.9; Defective mismatch repair

SBS10a_0.86; Polymerase epsilon mutation

SBS10a_0.76; Polymerase epsilon mutation

SBS10a_0.83; Polymerase epsilon mutation

SBS15_0.63; ID1(1.0) - INDEL driven

SBS7a_0.9; Ultraviolet light exposure

SBS10a_0.95; Polymerase epsilon mutation

SBS7a_0.86; Ultraviolet light exposure

SBS7d_0.59;

SBS2_0.71; APOBEC3B-like activity

SBS40_0.71;

SBS1_0.64; ID2(1.0) - INDEL driven

SBS9_0.63; Enriched in Lymph-BNHL

SBS21_0.75; Defective mismatch repair; ID7(0.96) + POLD(?)

SBS1_0.98; SBS1-like in Defective mismatch repair samples

SBS10b_0.73; Ultraviolet light exposure

SBS12_0.78; Defective mismatch repair; ID1(0.99) - INDEL driven

SBS41_0.73; Associated with BI_ID14(1.0) - INDEL driven

SBS28_0.87; Polymerase epsilon mutation - SBS28-like

SBS15_0.86; Defective mismatch repair

SBS5_0.8; Predominant in Prostate

SBS37_0.52; ID4(0.92) - INDEL driven

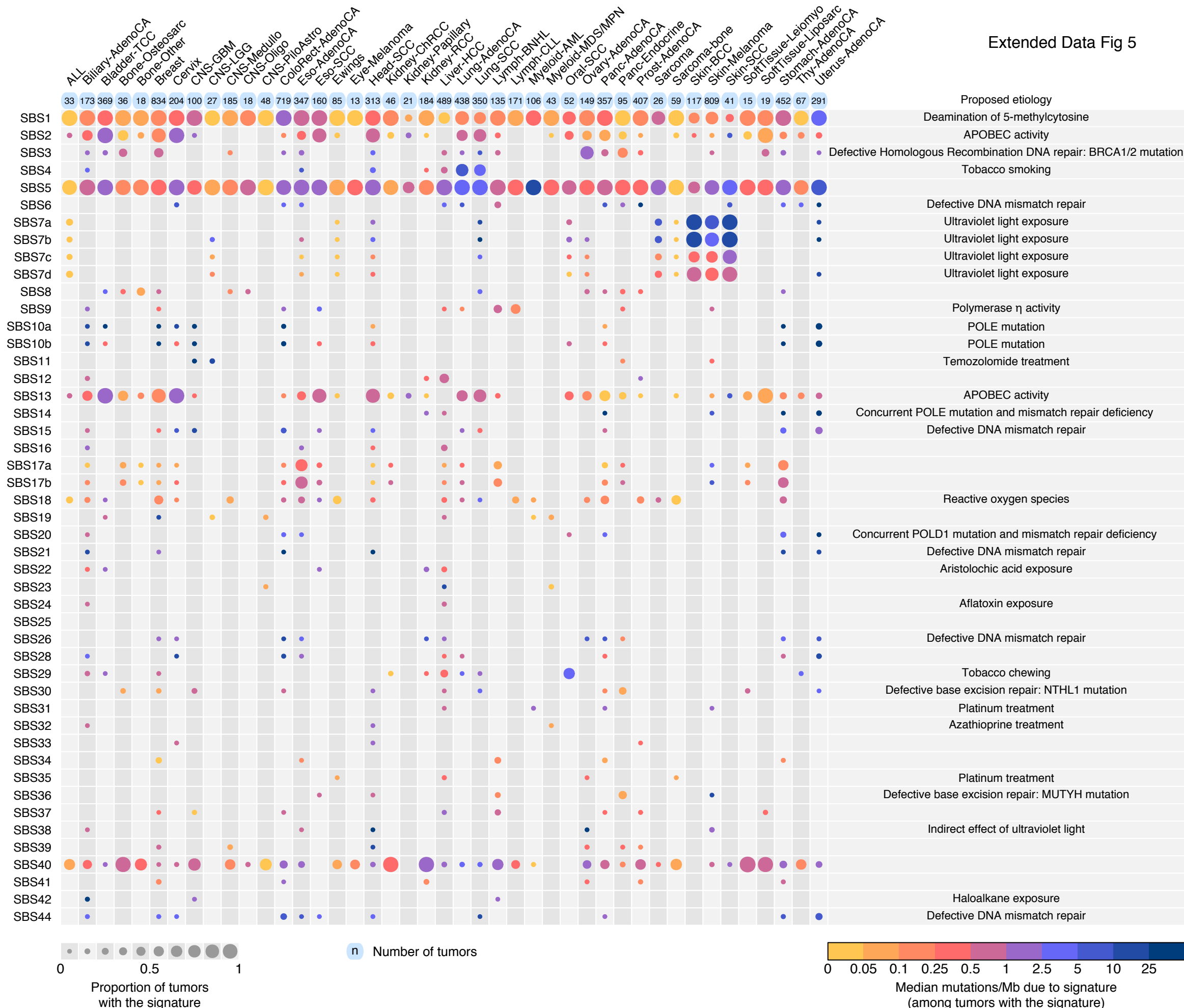
SBS3_0.74; ID11(1.0) - INDEL driven

SBS39_0.16



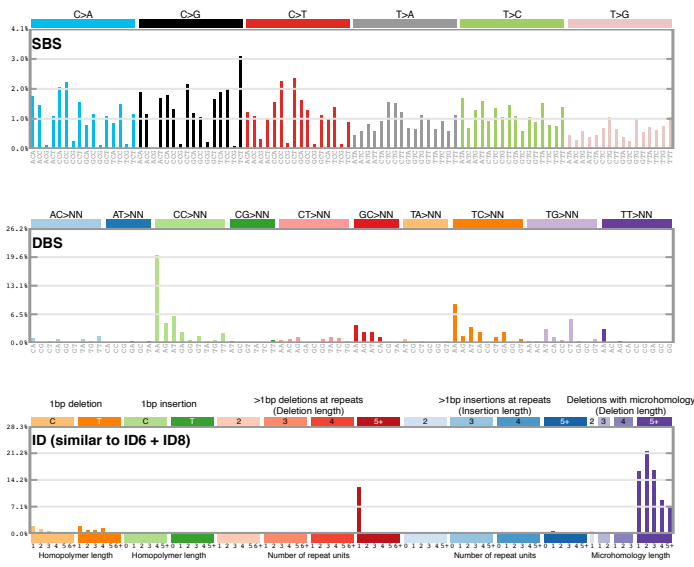
Median mutations/Mb due to signature
(among tumors with the signature)

Extended Data Fig 5



Extended Data Fig 7

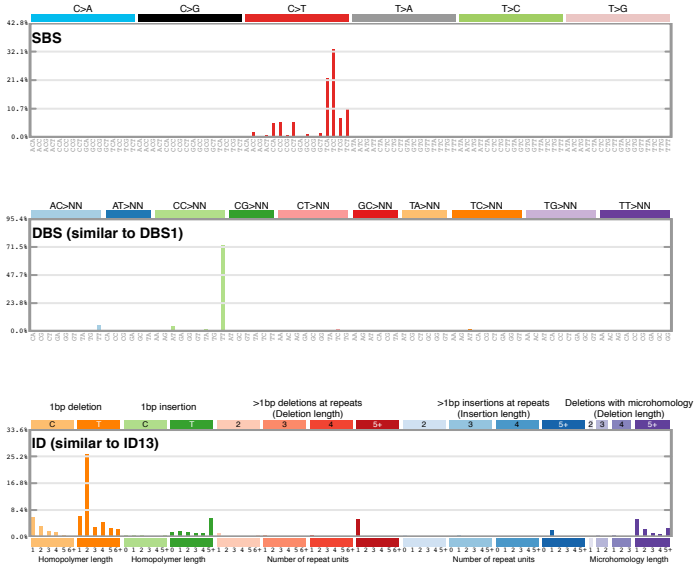
COMP-3



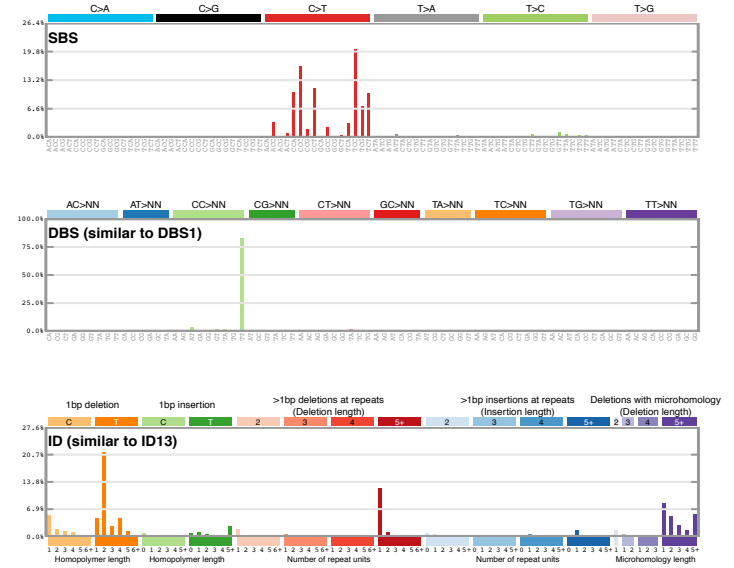
COMP-4



COMP-7a



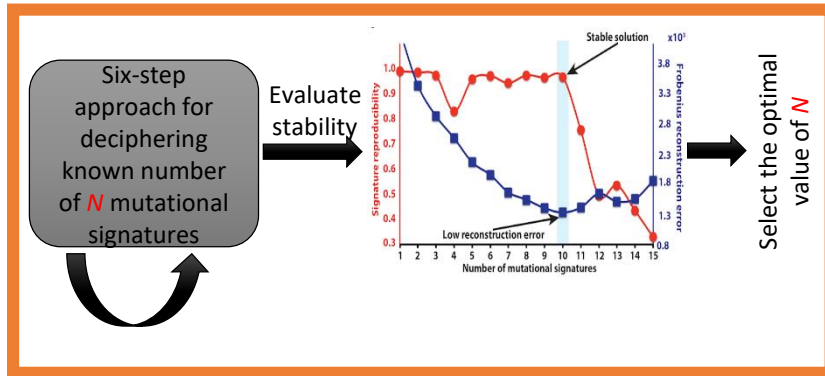
COMP-7b



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

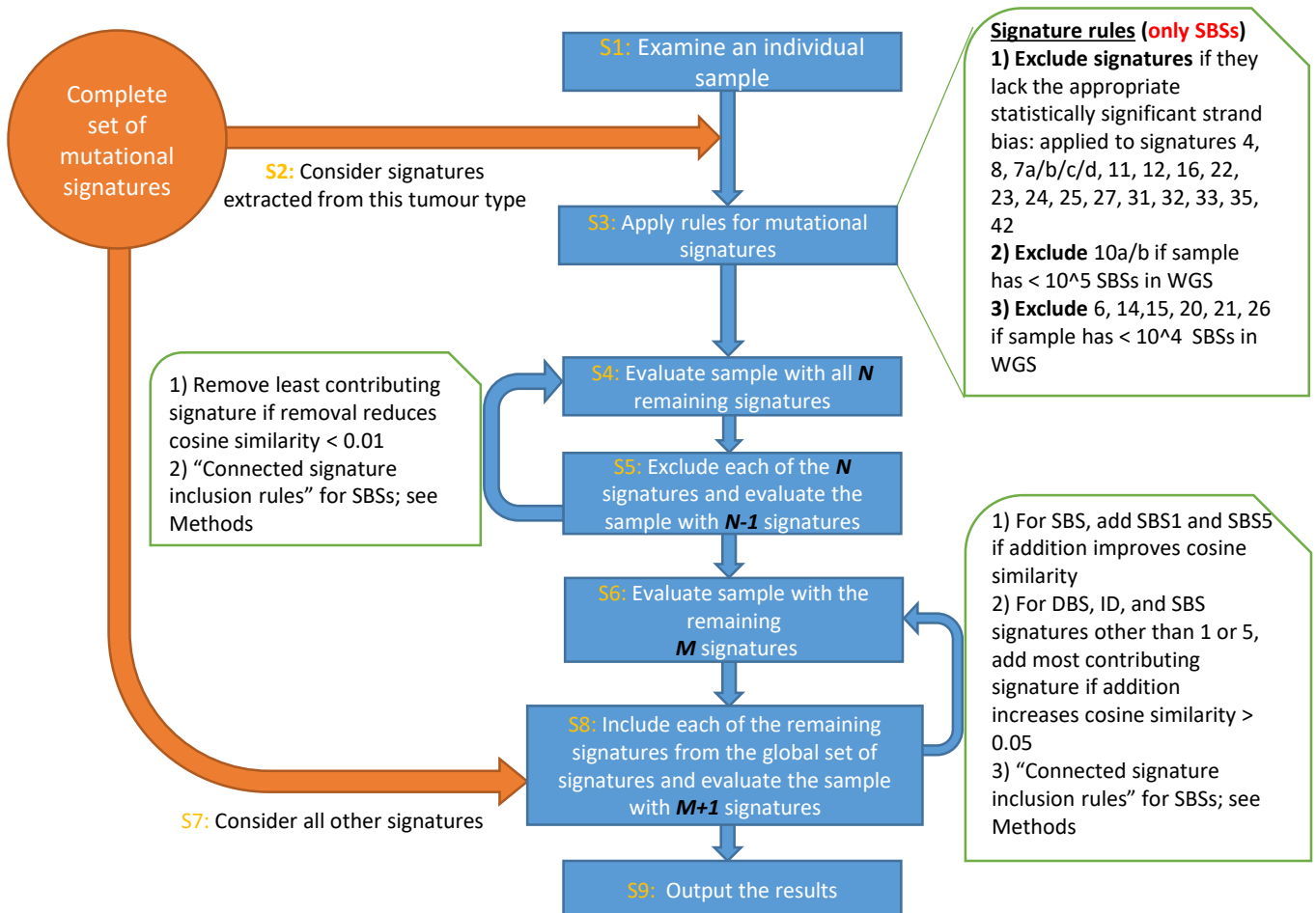
Step B (Solution evaluation and re-iteration)

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), R function call:

`glm(DBS.counts ~ SBS.counts + 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