

SOFTWARE

Open Access



Visualizing and exploring patterns of large mutational events with SigProfilerMatrixGenerator

Azhar Khandekar^{1,2,3}, Raviteja Vangara^{1,2,3}, Mark Barnes^{1,2,3}, Marcos Díaz-Gay^{1,2,3}, Ammal Abbasi^{1,2,3}, Erik N. Bergstrom^{1,2,3}, Christopher D. Steele^{1,2,3}, Nischalan Pillay^{4,5} and Ludmil B. Alexandrov^{1,2,3*}

Abstract

Background All cancers harbor somatic mutations in their genomes. In principle, mutations affecting between one and fifty base pairs are generally classified as small mutational events. Conversely, large mutational events affect more than fifty base pairs, and, in most cases, they encompass copy-number and structural variants affecting many thousands of base pairs. Prior studies have demonstrated that examining patterns of somatic mutations can be leveraged to provide both biological and clinical insights, thus, resulting in an extensive repertoire of tools for evaluating small mutational events. Recently, classification schemas for examining large-scale mutational events have emerged and shown their utility across the spectrum of human cancers. However, there has been no computationally efficient bioinformatics tool that allows visualizing and exploring these large-scale mutational events.

Results Here, we present a new version of SigProfilerMatrixGenerator that now delivers integrated capabilities for examining large mutational events. The tool provides support for examining copy-number variants and structural variants under two previously developed classification schemas and it supports data from numerous algorithms and data modalities. SigProfilerMatrixGenerator is written in Python with an R wrapper package provided for users that prefer working in an R environment.

Conclusions The new version of SigProfilerMatrixGenerator provides the first standardized bioinformatics tool for optimized exploration and visualization of two previously developed classification schemas for copy number and structural variants. The tool is freely available at <https://github.com/AlexandrovLab/SigProfilerMatrixGenerator> with an extensive documentation at <https://osf.io/s93d5/wiki/home/>.

Keywords Copy-number signatures, Structural variant signatures, Mutational patterns

*Correspondence:

Ludmil B. Alexandrov
L2alexandrov@health.ucsd.edu

¹ Department of Cellular and Molecular Medicine, UC San Diego, La Jolla, CA 92093, USA

² Department of Bioengineering, UC San Diego, La Jolla, CA 92093, USA

³ Moores Cancer Center, UC San Diego, La Jolla, CA 92037, USA

⁴ Research Department of Pathology, Cancer Institute, University College London, London WC1E 6BT, UK

⁵ Department of Cellular and Molecular Pathology, Royal National Orthopaedic Hospital NHS Trust, Stanmore HA7 4LP, Middlesex, UK

Background

Large-scale cancer genomics projects, including, The Cancer Genome Atlas (TCGA) and the Pan-cancer Analysis of Whole Genomes (PCAWG) initiatives, have comprehensively surveyed the molecular landscapes of most types of human cancer [1, 2]. These studies have provided a compendium of somatic mutations for each examined cancer genome and revealed both the mutations driving cancer development and the processes generating most somatic mutations within each cancer [1–3]. One commonly performed type of genomics analysis is the



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

examination of mutational patterns within a set of cancer genomes and the extraction of mutational signatures that have generated these patterns [3, 4]. Historically, mutational patterns have been predominately examined in the context of small mutational events, which include single base substitutions (SBS), doublet base substitutions (DBS), and small insertions and deletions (IDs) [3, 5]. Recent studies have also started exploring the patterns of large mutational events, including ones due to copy-number alterations and/or structural variations [6, 7]. Previously, we developed a computational tool, termed, SigProfilerMatrixGenerator, designed exclusively for examining the mutational patterns of all types of small mutational events [8]. Here, we present a new version of SigProfilerMatrixGenerator that now provides the capabilities for optimized exploration and visualization of large mutational events.

Large mutational events, generally defined as genomic alterations greater than 50 base pairs, are an important class of somatic aberrations in human cancer [6]. In principle, there are two commonly examined and closely interrelated types of large mutational events: (i) a structural variation (SV, also known as a genomic rearrangement), where a large-scale genomic segment gets altered; and (ii) a copy number variation (CNV), where the number of DNA copies of a genomic segment gets modified. Not all structural variations are related to CNVs, as SVs do not necessarily alter the copy number of a genomic segment; examples include copy neutral events such as inversions and reciprocal translocations. Similarly, not all changes in copy number require prior SVs, as is the case of chromosomal duplications and whole-genome doubling. Importantly, SVs and CNVs also differ in the types of genomics approaches that can detect them. In most cases, comprehensive detection of SVs requires whole-genome sequencing (WGS) data as it relies on either read alignment [9] or genome assembly methods [10]. In contrast, in addition to WGS data, CNVs can be detected from whole-exome sequencing, RNA-sequencing, single-cell sequencings approaches, and genotyping microarrays [11–13].

Deciphering mutational signatures from catalogues of somatic mutations, a process known as *de novo* signature extraction, relies on a biologically meaningful classification of mutational events [5]. We previously created the mathematical concept of mutational signatures and provided a set of tools for deciphering signatures of small mutational [4, 8]. Mutational patterns of SBSs, DBSs, IDs, have been extensively explored with more than 100 distinct mutational signatures published in the literature [3, 14]. These signatures reflect the activities of endogenous and/or exogenous mutational processes with approximately half of all signatures being, at least putatively,

linked with a proposed etiology [15–18]. Recently, mutational signature analyses of larger copy number alterations and structural alterations have emerged [6, 7, 19, 20]. A crucial first step in extracting mutational signatures is the derivation of features according to a predefined schema for mutational classification. This step involves transforming the mutational catalogues of a set of cancer genomes into a matrix, which is then amenable to subsequent matrix decomposition techniques [8]. Here, we present a computational package for classification of large-scale alterations and the generation of mutational matrices for signature decomposition. Two separate classification schemas are implemented: one for copy number variations and one for structural variations. Both schemas were previously developed and applied to large cohorts of cancer samples [7, 19, 21]. To the best of our knowledge, SigProfilerMatrixGenerator is the first tool that allows matrix generation and visualization of the CNV scheme used for generating the global reference set of Catalogue of Somatic Mutations in Cancer (COSMIC) copy-number signatures [7]. SigProfilerMatrixGenerator's capabilities for analyzing SVs and CNVs are implemented in Python and R, and the tool allows using multiple input formats, including segmentation and browser extensible data paired-end (BEDPE) files generated by commonly used algorithms for detecting copy number variations and structural variations, respectively. Additionally, SigProfilerMatrixGenerator provides a comprehensive visualization of mutational patterns of large mutational events and an R wrapper package for users that prefer working within the R environment.

Implementation

Classification of copy number variations

The schema for classifying copy number variations is based on Steele et al. [7] and it utilizes allele-specific copy number, which quantifies the number of segments for each allele at each variant loci rather than the total number of chromosome copies. In this schema, the copy-number profile of a sample can be represented by a mutational vector with 48 dimensions. Specifically, copy number segments are categorized into three heterozygosity states: heterozygous segments with total copy number (TCN) of $A > 0, B > 0$ (numbers reflect the counts for major allele A and minor allele B ; Fig. 1a), segments with loss of heterozygosity (LOH) with total copy number of $A > 0, B = 0$ (Fig. 1b), and segments with homozygous deletions and TCN of $A = 0, B = 0$ (Fig. 1c). Segments are further subclassified into 5 categories based on total copy number, which reflects the sum of the copies on the major allele A and the copies on the minor allele B : $TCN = 0, TCN = 1, TCN = 2, TCN = 3$ or $4, TCN = 5$ to 8 , and $TCN \geq 9$. Each of these total copy number

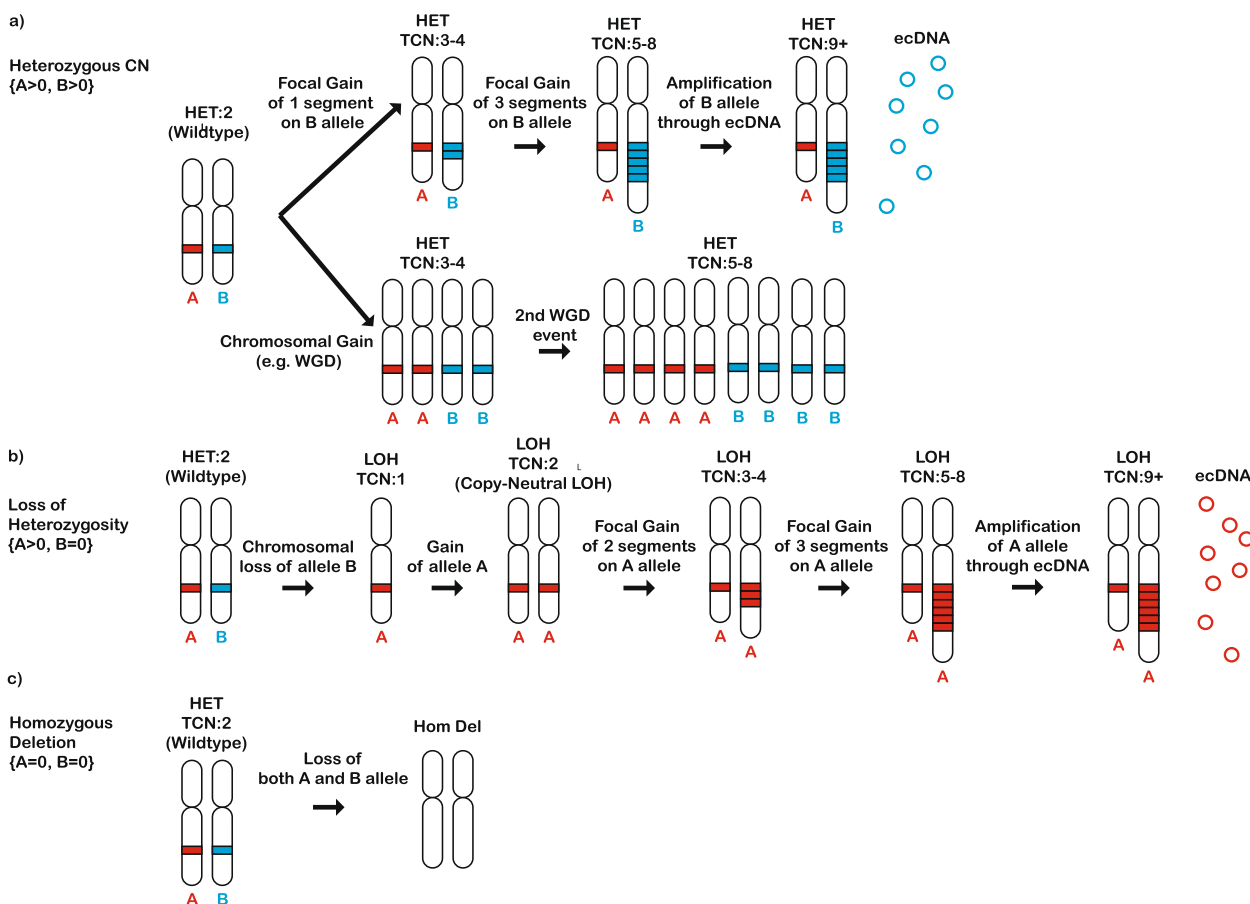


Fig. 1 Description of the copy number classification schema. The copy number classification schema consists of 48 mutually exclusive channels, divided by heterozygosity status, segment size, and total copy number (TCN). **a** In the heterozygous state, both alleles are retained and either one or both alleles can be amplified. This amplification can be focal (top panel) or it can encompass a chromosome or even the whole genome (bottom panel). The heterozygous category is further subdivided based on TCN (TCN = 1, TCN = 2, TCN = 3 or 4, TCN = 5 to 8, and TCN >= 9). **b** In a state of loss of heterozygosity (LOH), one of the alleles is lost. The remaining allele can then be duplicated (i.e., copy neutral LOH), and undergo more amplification resulting in higher total copy number states. The LOH category is further subdivided based on TCN (TCN = 1, TCN = 2, TCN = 3 or 4, TCN = 5 to 8, and TCN >= 9). The heterozygous and LOH categories are further divided on the basis of the size of the segment: 0 – 100 kb, 100 kb – 1 Mb, 1 Mb – 10 Mb, 10 Mb – 40 Mb, > 40 Mb. High-level LOH or heterozygous amplifications (e.g., TCN = 5 to 8 or TCN >= 9) can be carried on extrachromosomal DNA (depicted as red circles) as well as on linear chromosomes. **c** Homozygous deletions result in the loss of both alleles, and are divided on the basis of the size of the deleted segment: 0 – 100 kb, 100 kb – 1 Mb, and > 1 Mb

states accounts for the phenomenon of whole-genome duplication, for example a diploid (TCN=2) state transitioning to a doubled state (TCN=4), and a subsequent doubling of this state to TCN=8 is accounted for by the TCN=5–8 category (Fig. 1a). The categories for total copy number have been chosen for biological relevance (Fig. 1): TCN=0 reflects homozygous deletions, TCN=1 represents a genomic deletion resulting in an LOH, TCN=2 is equivalent to a diploid state including copy neutral LOH (a phenomenon whereby one of two homologous chromosomal regions is lost, but two identical copies of this region still remain; Fig. 1b), TCN=3 or 4 reflect a gained state of tri- to tetra-ploidy, TCN=5 to 8 represent a penta- to octo-ploidy state, and TCN >= 9

represents high-level amplifications such as ones found in samples containing extrachromosomal DNA (ecDNA) [22]. Each of the heterozygous and LOH total copy number categories are additionally subclassified into five additional categories based on the size of their segments: 0 – 100 kb, 100 kb – 1 Mb, 1 Mb – 10 Mb, 10 Mb – 40 Mb, and > 40 Mb. Three size bins are used for the additional subcategorization of homozygous deletions: 0 – 100 kb, 100 kb – 1 Mb, and > 1 Mb. The partitioning by segment sizes was chosen to ensure that a sufficient proportion of segments are classified within each category [7]. This classification allows summarizing copy number profiles using 48 distinct channels and can be represented using a vector with 48 components. For example,

a sample harboring multiple focal amplifications, either contained on linear or extrachromosomal DNA, will have many events in the 9+ total copy number category and the first 3 size bins (0 – 100 kb, 100 kb – 1 Mb, 1 Mb – 10 Mb; Fig. 2a, b). Conversely, a sample containing a large number of focal deletions or losses of entire chromosomes or chromosome arms will have numerous events in the LOH category, spanning all size bins (Fig. 2c, d). Another example will be a sample with a whole-genome doubling where copy number changes will primarily encompass segments with large genomic sizes (10 Mb – 40 Mb; 40 Mb) and total copy number between 3 and 4 (Fig. 2e, f). Overall, this 48-channel classification schema can effectively summarize a diverse array of copy number states seen across tumor types [7], whether they contain broad or focal events that result in amplifications or deletions.

Input data for classifying copy number variations

SigProfilerMatrixGenerator allows examining allele specific CNV data that, at a minimum, include the following information for each CNV segment: chromosome, start coordinate, end coordinate, and copy number of both the

minor and major alleles. Output files from the following tools for detecting CNVs are automatically supported: ASCAT [23], ABSOLUTE [24], Sequenza [25], FACETS [12], Battenberg [23], and PURPLE [26]. Additionally, custom segmentation files from other CNV detection tools can be used if these files contain the aforementioned information.

Classification of structural variants

A classification schema consisting of 32 features, based on Nik-Zainal et al. [21], is used to construct a mutational vector with 32 dimensions for each sample. In principle, each structural variant consists of two breakpoints which are at single-base resolution, where a breakpoint is defined as a junction that indicates a structurally variable genomic segment greater than 50 base pairs [10]. Breakpoints are typically detected using three signals from aligned sequencing reads: depth of sequence coverage, discordant read-pairs, and split read-pairs [27–29]. Breakpoints can also be detected via genome assembly, where reads are assembled into contigs, the contigs are aligned to the reference genome, and these alignments are analyzed for structural

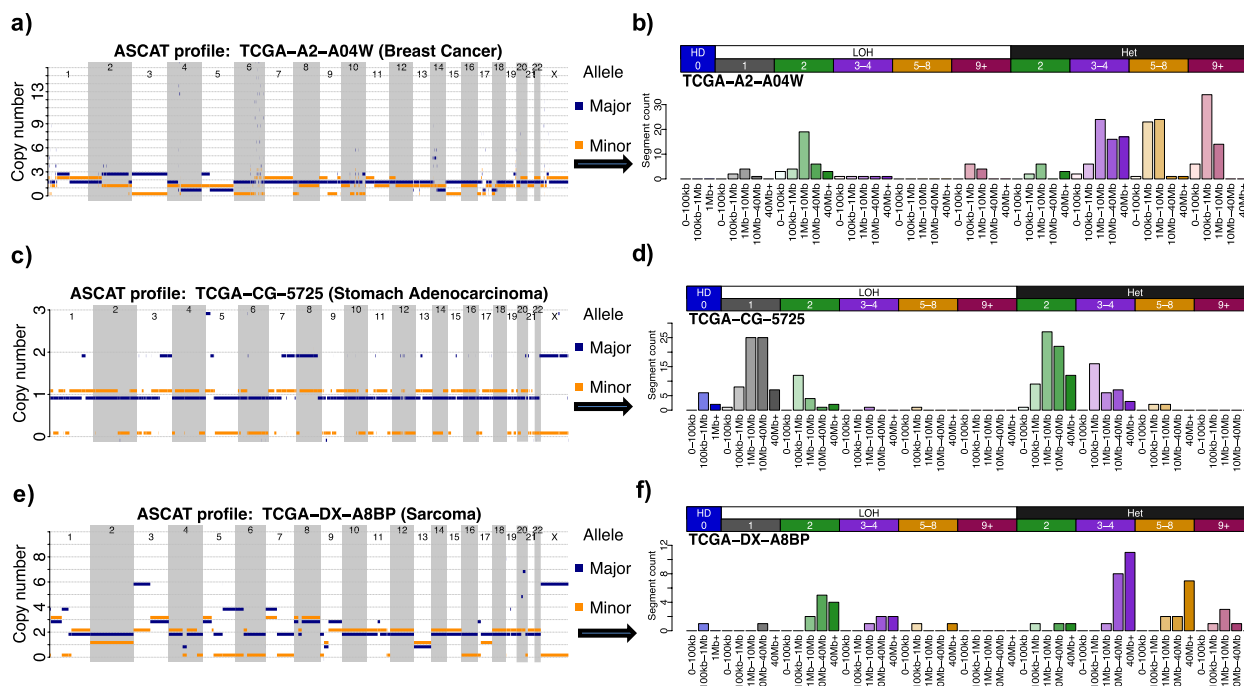


Fig. 2 Converting copy number segmentation profiles into copy number mutational vectors. The CNV classification schema converts a sample’s segmentation profile (a, c, e) into a count vector of 48 mutually exclusive components (b, d, f). These components are based on segment size, heterozygosity status, and total copy number. A breast cancer sample with many highly amplified segments, possibly due to the presence of extrachromosomal DNA, is shown in (a, b). This sample’s count vector is characterized by peaks in the 5–8 and 9+ total copy number categories. A gastric cancer sample with extensive loss of heterozygosity is shown in (c, d). This sample’s count vector is characterized by peaks in the LOH category, specifically with a total copy number of 1 indicating a loss of an allele. A sarcoma sample with a whole-genome duplication event, characterized by peaks in the 3–4 total copy number category and the 40+ Mb size bin, is shown in (e, f)

variants [10]. The previously developed classification of structural variants considers the following canonical SVs: tandem duplications, deletions, inversions, and translocations (Fig. 3). A tandem duplication refers to a segment of genomic material that has been duplicated and inserted on the same chromosome adjacent to the original segment (Fig. 3a). It should be noted that a tandem duplication is not necessarily the same as a copy-number amplification. For example, ecDNA copy-number amplifications are not tandem duplications as they are not inserted adjacent to the original chromosome segment. A somatic deletion is an event that has removed a set of existing base-pairs from a given location of a chromosome (Fig. 3b). An inversion is when a segment of the chromosome breaks off and reattaches at the same locus but in a reverse orientation (Fig. 3c). A translocation event occurs when a piece of one chromosome breaks off and some (or all) fragments from that piece re-attach to either another chromosome or to a different locus of the same chromosome (Fig. 3d). The classification schema bins all SVs, apart from translocations, according to the size of the event in base pairs: 0–10 kb, 10 kb–100 kb, 100 kb–1 Mb, 1 Mb–10 Mb, and > 10 Mb [21]. Translocations, which may involve

more than one chromosome, are not binned by size because they can be either balanced (where there is no net loss of genetic material on the chromosomes involved and thus the size can be described by one number) or unbalanced (where there is a net loss or gain of genetic material on the chromosomes involved and thus the sizes of the segments cannot be described by just one number). Note that whether a translocation is balanced or unbalanced is not considered in this classification schema. The different types of SVs are then further divided into *clustered* and *non-clustered* events to account for the non-random distribution of these events along the genome. Clustered events are defined as events that occur closer to each other on a chromosome than purely expected by chance. These clusters often arise as a result of complex events, such as chromothripsis [30] or chromoplexy [31], generating many breakpoints in a single instantaneous event as opposed to the gradual accumulation of events over many cell cycles which results in more dispersed non-clustered events. Clusters of breakpoints can also form as a result of other mechanisms, including, for example, rearrangement hotspots in the genome [32]. Clustering of SVs is determined based on a previously developed algorithm

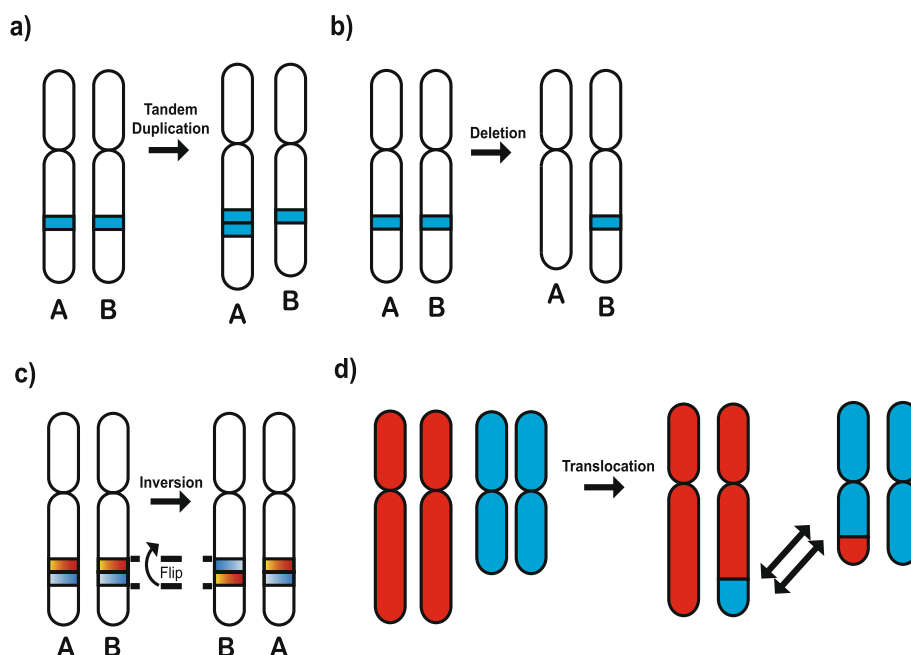


Fig. 3 Description of the structural variant classification schema. Structural variants (SVs) are categorized as tandem-duplications, deletions, inversions, or translocations. **a** Tandem duplication of a segment containing the A allele. A tandem duplication occurs when a segment is duplicated and inserted adjacent to the original chromosomal segment. **b** Deletion of the segment containing the A allele. A deletion occurs when there is a loss of genetic material from a chromosome. **c** An inversion of the segment containing the B allele. An inversion occurs when a segment breaks off and reattaches in a reverse orientation within the same chromosome. **d** A translocation of a chromosomal segment. A translocation event occurs when a piece of one chromosome breaks off and some (or all) fragments from that piece re-attach to either another chromosome or to a different locus of the same chromosome

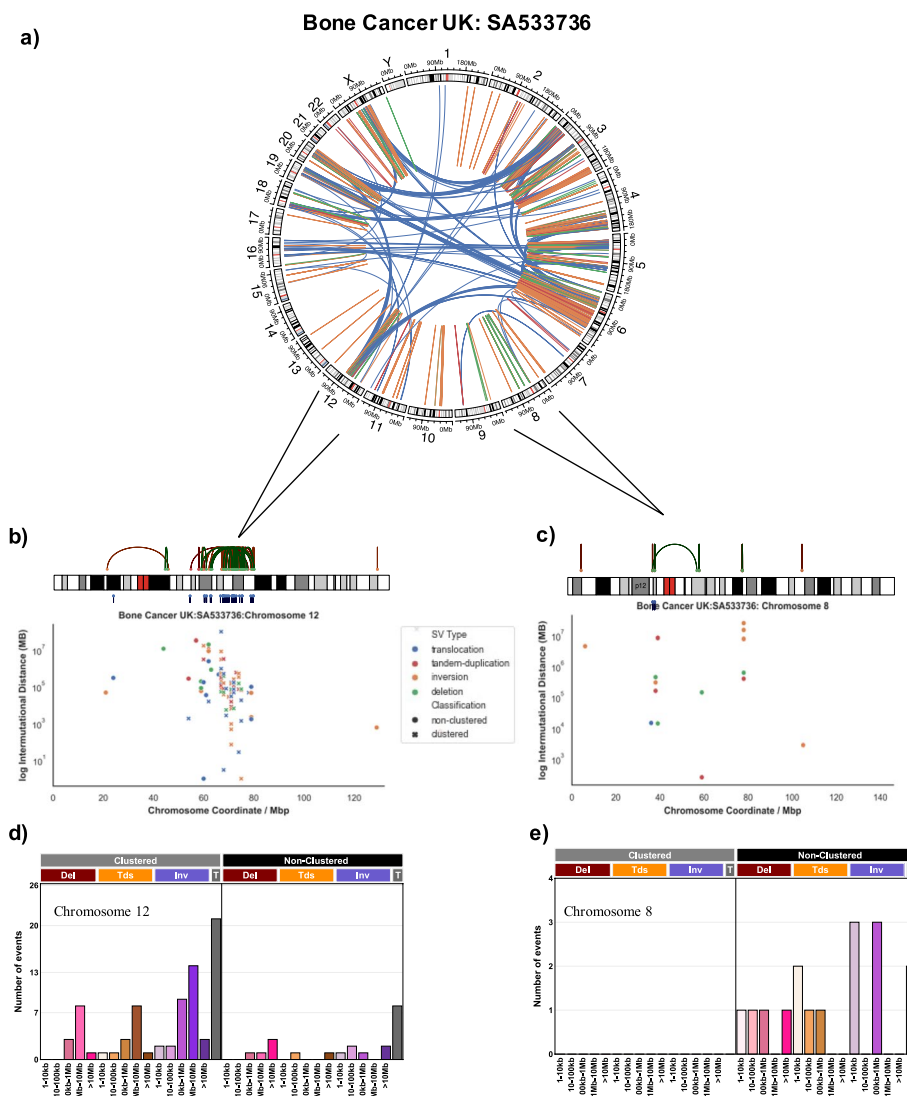


Fig. 4 Classifying Structural Variants into Mutational Vectors. **a** An example of a bone cancer sample from PCAWG with a highly rearranged genome consisting of both clustered and non-clustered structural variants (SVs) is shown as a Circos plot representation. **b** Zooming into SVs specifically found on chromosome 12 in the bone cancer sample. SVs are shown as a linear representation (top) and as a rainfall plot (bottom). The rainfall plot depicts all breakpoints on chromosome 12 according to their genomic coordinate (x-axis) and the \log_{10} inter-mutational distance (y-axis), which is the distance to the breakpoint immediately preceding it. The tendency of breakpoints to cluster in a specific genomic region on chromosome 12 due to a chromothripsis event is evident in all representations. **c** Zooming into SVs specifically found on chromosome 8 in the bone cancer sample. SVs are shown as a linear representation (top) and as a rainfall plot (bottom). The rainfall plot depicts all breakpoints on chromosome 8 according to their genomic coordinate (x-axis) and the \log_{10} inter-mutational distance (y-axis), which is the distance to the breakpoint immediately preceding it. There are no clustered SVs on chromosome 8 as, per the SV classification schema, clustering requires a minimum of 10 breakpoints in a segment of a chromosome. **d** The SV classification schema is applied to the SVs found on chromosome 12 in the bone cancer sample. SVs are classified by the event type (denoted by color) and are binned according to the size of the event (0 – 10 kb, 10 kb – 100 kb, 100 kb – 1 Mb, 1 Mb – 10 Mb, and > 10 Mb). **e** The SV classification schema is applied to the SVs found on chromosome 8 in the bone cancer sample. SVs are classified by the event type (denoted by color) and are binned according to the size of the event (0 – 10 kb, 10 kb – 100 kb, 100 kb – 1 Mb, 1 Mb – 10 Mb, and > 10 Mb)

that utilizes the Potts’ filter method [33]. This method segments a chromosome based on inter-mutational distance of SV breakpoints, and if the average distance in a particular segment is less than 10 times the average inter-mutational distance in the sample, all breakpoints

in the segment are considered clustered. A minimum of 10 breakpoints must be present for a given segment to be considered clustered, otherwise all breakpoints in that segment are considered non-clustered.

An example of a whole-genome sequenced bone cancer with a highly rearranged genome that contains chromosomes with clustered events as well as chromosomes with only non-clustered events is shown in Fig. 4a. For instance, in this sample, chromosome 12 contains a high number of SV breakpoints in close proximity to one another (Fig. 4b) and the SV pattern of this chromosome can be summarized in a vector with 32 components containing a high number of clustered SVs (Fig. 4d). In contrast, chromosome 8 has SV breakpoints randomly scattered throughout the chromosome (Fig. 4c) and the SV pattern of chromosome 8 is exclusively one of non-clustered SVs (Fig. 4e).

Input data for classifying structural variants

SigProfilerMatrixGenerator allows examining SV data that contains genomics information for each of the two breakpoints of a structural variant. In principle, the tool can process files in browser extensible data paired-end (BEDPE) format that, at a minimum, contain the following six columns: *chrom1*, *start1*, *end1*, *chrom2*, *start2*, and *end2*. Here, the genomics coordinates of the first breakpoint are annotated as *chrom1*, *start1*, and *end1*, while the genomics coordinates of the second breakpoint are provided as *chrom2*, *start2*, and *end2*. If the type of SV has been predetermined, then its annotation can be provided using a column named *svclass*. Otherwise, the columns *strand1* and *strand2*, which indicate the strands of the read mate-pairs, are required. If the mates are on the same chromosome, the convention followed is inversion (+/- or -/+), deletion (+/+), and tandem-duplication (-/-). If mates are on different chromosomes, the SV is automatically classified as a translocation. SigProfilerMatrixGenerator supports SV in BEDPE format, which is utilized by most bioinformatics tools for detecting SVs, as well as being the native output files from BRASS [21].

Discussion

The newly developed version of SigProfilerMatrixGenerator allows transforming a set of mutational catalogues of copy-number changes and structural rearrangements into matrices amenable to decomposition, including, subsequent mutational signature analysis. The tool provides support for two previously developed [7, 21] classification schemas for large mutational events and seamlessly integrates with other components of the SigProfiler software suite, such as downstream signature extraction with SigProfilerExtractor [4] and visualization of both mutational patterns and signatures with SigProfilerPlotting [8]. Plots for CNV and SV patterns can now be generated for each cancer sample (as shown in Figs. 2 and 4), and plots for CNV and SV signatures are automatically generated following signature extraction from a cohort of samples.

This enables a streamlined workflow for end-to-end analysis of mutational signatures of large-scale events. Additionally, SigProfilerMatrixGenerator rapidly scales to large datasets. For example, the tool can generate an SV count matrix for all 2,658 PCAWG samples in 3.6 s and a CNV count matrix for the entire TCGA array data (9,875 samples) in 14.3 s. SigProfilerMatrixGenerator is also the first tool to provide support for the 48 channel CNV schema across a wide variety of popular tools for detecting CNV. Importantly, this schema can be applied across several data modalities, including whole-genome sequencing, whole-exome sequencing, RNA-sequencing, single-cell sequencing approaches, and genotyping microarrays. In addition, SigProfilerMatrixGenerator is the first Python package that provides support for the 32 channel SV schema in a fast and intuitive manner with minimal preprocessing, and the only package to provide support for SV and CNV schemas in both a Python and R environment.

Conclusion

A breadth of computational tools exists for exploring the patterns for small mutational events, including our initial implementation of SigProfilerMatrixGenerator [8]. We recently demonstrated that a classification of CNVs into 48 channels provides the means to better elucidate and understand the mutational processes operative in human cancer [7]. Similarly, we and others have previously demonstrated that the classification of SVs into 32 channels can be used to understand the mutational processes giving rise to SVs across multiple cancer types [19]. Our newly developed version of SigProfilerMatrixGenerator provides the capability to examine these classification schemas from cancer genomics sequencing data. The tool can scale to large datasets and will serve as foundation for future analysis of both mutational patterns and mutational signatures of large mutational events.

Availability and requirements

Project name: SigProfilerMatrixGenerator.

Project home page: <https://github.com/AlexandrovLab/SigProfilerMatrixGenerator>,

Operating system(s): Unix, Linux, and Windows.

Programming language: Python 3 and R.

Other requirements: None.

License: BSD 2-Clause "Simplified" License.

Any restrictions to use by non-academics: None.

Abbreviations

BEDPE	Browser extensible data paired-end
CNV	Copy number variation
DBS	Doublet base substitution
ecDNA	Extrachromosomal DNA
ID	Small insertions and deletions

LOH	Loss of heterozygosity
PCAWG	Pan-cancer Analysis of Whole Genomes
SBS	Single base substitution
SV	Structural variation
TCGA	The Cancer Genome Atlas
TCN	Total copy-number
WGS	Whole-genome sequencing

Acknowledgements

The computational development reported in this manuscript have utilized the Triton Shared Computing Cluster at the San Diego Supercomputer Center of UC San Diego.

Authors' contributions

AK developed the Python and R code with assistance from RV, MB, AA, ENB, and MDG. AA, CDS, and NP tested and evaluated the performance of the code. CDS, NP, and LBA developed the copy number classifications schema. AK wrote the manuscript with assistance from RV, MB, CDS, AA, and MDG. LBA supervised the overall development of the code and writing of the manuscript. All authors read and approved the final manuscript.

Funding

This work was supported by the US National Institute of Health grants R01ES030993-01A1, R01ES032547-01, and R01CA269919-01 to LBA as well as by Cancer Research UK Grand Challenge Award C98/A24032. This work was also supported by a Packard Fellowship for Science and Engineering. The funders had no roles in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Availability of data and materials

Data sharing is not applicable to this article as no datasets were generated or analyzed during the current study.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

LBA is a compensated consultant and has equity interest in io9, LLC. His spouse is an employee of Biotheranostics, Inc. LBA is also an inventor of a US Patent 10,776,718 for source identification by non-negative matrix factorization. LBA declares U.S. provisional applications with serial numbers: 63/289,601; 63/269,033; 63/366,392; 63/367,846; 63/412,835. All other authors declare that they have no competing interests.

Received: 3 February 2023 Accepted: 14 August 2023

Published online: 21 August 2023

References

- Cancer Genome Atlas Research N, Weinstein JN, Collisson EA, Mills GB, Shaw KR, Ozenberger BA, Ellrott K, Shmulevich I, Sander C, Stuart JM. The Cancer genome atlas pan-cancer analysis project. *Nat Genet*. 2013;45(10):1113–20.
- Consortium ITP-CAoWG. Pan-cancer analysis of whole genomes. *Nature*. 2020;578(7793):82–93.
- Alexandrov LB, Kim J, Haradhvala NJ, Huang MN, Tian Ng AW, Wu Y, Boot A, Covington KR, Gordenin DA, Bergstrom EN, et al. The repertoire of mutational signatures in human cancer. *Nature*. 2020;578(7793):94–101.
- Islam SMA, Diaz-Gay M, Wu Y, Barnes M, Vangara R, Bergstrom EN, He Y, Vella M, Wang J, Teague JW, et al. Uncovering novel mutational signatures by de novo extraction with SigProfilerExtractor. *Cell Genom*. 2022;2(11):None.
- Alexandrov LB, Nik-Zainal S, Wedge DC, Campbell PJ, Stratton MR. Deciphering signatures of mutational processes operative in human cancer. *Cell Rep*. 2013;3(1):246–59.
- Li Y, Roberts ND, Wala JA, Shapira O, Schumacher SE, Kumar K, Khurana E, Waszak S, Korbel JO, Haber JE, et al. Patterns of somatic structural variation in human cancer genomes. *Nature*. 2020;578(7793):112–21.
- Steele CD, Abbasi A, Islam SMA, Bowes AL, Khandekar A, Haase K, Hames-Fathi S, Ajayi D, Verfaillie A, Dhami P, et al. Signatures of copy number alterations in human cancer. *Nature*. 2022;606(7916):984–91.
- Bergstrom EN, Huang MN, Mahto U, Barnes M, Stratton MR, Rozen SG, Alexandrov LB. SigProfilerMatrixGenerator: a tool for visualizing and exploring patterns of small mutational events. *BMC Genomics*. 2019;20(1):685.
- Cameron DL, Di Stefano L, Papefuss AT. Comprehensive evaluation and characterisation of short read general-purpose structural variant calling software. *Nat Commun*. 2019;10(1):3240.
- Cosenza MR, Rodriguez-Martin B, Korbel JO. Structural variation in cancer: role, prevalence, and mechanisms. *Annu Rev Genomics Hum Genet*. 2022;23:123–52.
- Talevich E, Shain AH, Botton T, Bastian BC. CNVkit: genome-wide copy number detection and visualization from targeted DNA sequencing. *PLoS Comput Biol*. 2016;12(4):e1004873.
- Shen R, Seshan VE. FACETS: allele-specific copy number and clonal heterogeneity analysis tool for high-throughput DNA sequencing. *Nucleic Acids Res*. 2016;44(16):e131.
- Serin Harmanci A, Harmanci AO, Zhou X. CaSpER identifies and visualizes CNV events by integrative analysis of single-cell or bulk RNA-sequencing data. *Nat Commun*. 2020;11(1):89.
- Degasperi A, Zou X, Amarante TD, Martinez-Martinez A, Koh GCC, Dias JML, Heskin L, Chmelova L, Rinaldi G, Wang VYW, et al. Substitution mutational signatures in whole-genome-sequenced cancers in the UK population. *Science*. 2022;376(6591):science.abl9283.
- Alexandrov LB, Ju YS, Haase K, Van Loo P, Martincorena I, Nik-Zainal S, Totoki Y, Fujimoto A, Nakagawa H, Shibata T, et al. Mutational signatures associated with tobacco smoking in human cancer. *Science*. 2016;354(6312):618–22.
- Petljak M, Alexandrov LB, Brummel DS, Price S, Wedge DC, Grossmann S, Dawson KJ, Ju YS, Iorio F, Tubio JMC, et al. Characterizing mutational signatures in human cancer cell lines reveals episodic APOBEC mutagenesis. *Cell*. 2019;176(6):1282–1294 e1220.
- Riva L, Pandiri AR, Li YR, Droop A, Hewinson J, Quail MA, Iyer V, Shepherd R, Herbert RA, Campbell PJ, et al. The mutational signature profile of known and suspected human carcinogens in mice. *Nat Genet*. 2020;52(11):1189–97.
- Alexandrov LB, Jones PH, Wedge DC, Sale JE, Campbell PJ, Nik-Zainal S, Stratton MR. Clock-like mutational processes in human somatic cells. *Nat Genet*. 2015;47(12):1402–7.
- Degasperi A, Amarante TD, Czarnecki J, Shooter S, Zou X, Glodzik D, Morganello S, Nanda AS, Badja C, Koh G. A practical framework and online tool for mutational signature analyses show intertissue variation and driver dependencies. *Nature cancer*. 2020;1(2):249–63.
- Drews RM, Hernando B, Tarabichi M, Haase K, Leslyes T, Smith PS, Morrill Gavarro L, Couturier DL, Liu L, Schneider M, et al. A pan-cancer compendium of chromosomal instability. *Nature*. 2022;606(7916):976–83.
- Nik-Zainal S, Davies H, Staaf J, Ramakrishna M, Glodzik D, Zou X, Martincorena I, Alexandrov LB, Martin S, Wedge DC, et al. Landscape of somatic mutations in 560 breast cancer whole-genome sequences. *Nature*. 2016;534(7605):47–54.
- Kim H, Nguyen NP, Turner K, Wu S, Gujar AD, Luebeck J, Liu J, Deshpande V, Rajkumar U, Namburi S, et al. Extrachromosomal DNA is associated with oncogene amplification and poor outcome across multiple cancers. *Nat Genet*. 2020;52(9):891–7.
- Van Loo P, Nordgard SH, Lingjaerde OC, Russnes HG, Rye IH, Sun W, Weigman VJ, Marynen P, Zetterberg A, Naume B, et al. Allele-specific copy number analysis of tumors. *Proc Natl Acad Sci U S A*. 2010;107(39):16910–5.
- Carter SL, Cibulskis K, Helman E, McKenna A, Shen H, Zack T, Laird PW, Onofrio RC, Winckler W, Weir BA, et al. Absolute quantification of somatic DNA alterations in human cancer. *Nat Biotechnol*. 2012;30(5):413–21.

25. Favero F, Joshi T, Marquard AM, Birkbak NJ, Krzystanek M, Li Q, Szallasi Z, Eklund AC. Sequenza: allele-specific copy number and mutation profiles from tumor sequencing data. *Ann Oncol.* 2015;26(1):64–70.
26. Shale C, Cameron DL, Baber J, Wong M, Cowley MJ, Papenfuss AT, Cuppen E, Priestley P. Unscrambling cancer genomes via integrated analysis of structural variation and copy number. *Cell Genomics.* 2022;2(4): 1001–12.
27. Wala JA, Bandopadhyay P, Greenwald NF, O'Rourke R, Sharpe T, Stewart C, Schumacher S, Li Y, Weischenfeldt J, Yao X, et al. SvABA: genome-wide detection of structural variants and indels by local assembly. *Genome Res.* 2018;28(4):581–91.
28. Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Kallberg M, Cox AJ, Kruglyak S, Saunders CT. Manta: rapid detection of structural variants and indels for germline and cancer sequencing applications. *Bioinformatics.* 2016;32(8):1220–2.
29. Rausch T, Zichner T, Schlattl A, Stutz AM, Benes V, Korbel JO. DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics.* 2012;28(18):i333–9.
30. Stephens PJ, Greenman CD, Fu B, Yang F, Bignell GR, Mudie LJ, Pleasance ED, Lau KW, Beare D, Stebbings LA, et al. Massive genomic rearrangement acquired in a single catastrophic event during cancer development. *Cell.* 2011;144(1):27–40.
31. Shen MM. Chromoplexy: a new category of complex rearrangements in the cancer genome. *Cancer Cell.* 2013;23(5):567–9.
32. Glodzik D, Morganella S, Davies H, Simpson PT, Li Y, Zou X, Diez-Perez J, Staaf J, Alexandrov LB, Smid M, et al. A somatic-mutational process recurrently duplicates germline susceptibility loci and tissue-specific super-enhancers in breast cancers. *Nat Genet.* 2017;49(3):341–8.
33. Winkler G, Liebscher V. Smoothers for discontinuous signals. *Journal of Nonparametric Statistics.* 2002;14(1–2):203–22.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

