

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

Visualizing and exploring patterns of large mutational events with SigProfilerMatrixGenerator

Azhar Khandekar¹⁻³, Raviteja Vangara¹⁻³, Mark Barnes¹⁻³, Marcos Díaz-Gay¹⁻³, Ammal Abbasi¹⁻³, Erik N. Bergstrom¹⁻³, Christopher D. Steele¹⁻³, Nischalan Pillay⁴⁻⁵, and Ludmil B. Alexandrov^{1-3*}

Affiliations:

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

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

³Moore's Cancer Center, UC San Diego, La Jolla, CA, 92037, USA

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

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

*Correspondence should be addressed to L2alexandrov@health.ucsd.edu.

Keywords: copy-number signatures, structural variant signatures, mutational patterns

28 **ABSTRACT**

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

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

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

50

51

52 **BACKGROUND**

53 Large-scale cancer genomics projects have comprehensively surveyed the molecular landscapes
54 of most types of human cancer [1, 2]. These studies have provided a compendium of somatic
55 mutations for each examined cancer genome and revealed both the mutations driving cancer
56 development and the processes generating most somatic mutations within each cancer [1-3]. One
57 commonly performed type of genomics analysis is the examination of mutational patterns within
58 a set of cancer genomes and the extraction of mutational signatures that have given rise to these
59 patterns [3, 4]. Historically, mutational patterns have been predominately examined in the
60 context of small mutational events, which include single base substitutions (SBS), doublet base
61 substitutions (DBS), and small insertions and deletions (IDs) [3, 5]. Recent studies have also
62 started exploring the patterns of large mutational events, including ones due to copy-number
63 alterations and/or structural variations [6, 7]. Previously, we developed a computational tool,
64 termed, SigProfilerMatrixGenerator, designed exclusively for examining the mutational patterns
65 of all types of small mutational events [8]. Here, we present a new version of
66 SigProfilerMatrixGenerator that now provides the capabilities for optimized exploration and
67 visualization of large mutational events.

68

69 Large mutational events, generally defined as genomic alterations greater than 50 base pairs, are
70 an important class of somatic aberrations in human cancer [6]. In principle, there are two
71 commonly examined and closely interrelated types of large mutational events: *(i)* a structural
72 variation (SV, also known as a genomic rearrangement), where a large-scale genomic segment
73 gets altered; and *(ii)* a copy number variation (CNV), where the number of DNA copies of a
74 genomic segment gets modified. Not all structural variations are related to CNVs, as SVs do not

75 necessarily alter the copy number of a genomic segment; examples include copy neutral events
76 such as inversions and reciprocal translocations. Similarly, not all changes in copy number
77 require prior SVs, as is the case of chromosomal duplications and whole-genome doubling.
78 Importantly, SVs and CNVs also differ in the types of genomics approaches that can detect them.
79 In most cases, comprehensive detection of SVs requires whole-genome sequencing (WGS) data
80 as it relies on either read alignment [9] or genome assembly methods [10]. In contrast, in
81 addition to WGS data, CNVs can be detected from whole-exome sequencing, RNA-sequencing,
82 single-cell sequencings approaches, and genotyping microarrays [11-13].
83
84 Deciphering mutational signatures from catalogues of somatic mutations, a process known as *de*
85 *novo* signature extraction, relies on a biologically meaningful classification of mutational events
86 [5]. We previously created the mathematical concept of mutational signatures and provided a set
87 of tools for deciphering signatures of small mutational [4, 8]. Mutational patterns of SBSs,
88 DBSs, IDs, have been extensively explored with more than 100 distinct mutational signatures
89 published in the literature [3, 14]. These signatures reflect the activities of endogenous and/or
90 exogenous mutational processes with an approximately half of all signatures being, at least
91 putatively, linked with a proposed etiology [15-18]. Recently, mutational signature analyses of
92 larger copy number alterations and structural alterations have emerged [6, 7, 19, 20]. A crucial
93 first step in extracting mutational signatures is the derivation of features according to a
94 predefined schema for mutational classification. This step involves transforming the mutational
95 catalogues of a set of cancer genomes into a matrix, which is then amenable to subsequent matrix
96 decomposition techniques [8]. Here, we present a computational package for classification of
97 large-scale alterations and the generation of mutational matrices for signature decomposition.

98 Two separate classification schemas are implemented: one for copy number variations and one
99 for structural variations. Both schemas were previously developed and applied to large cohorts of
100 cancer samples [7, 19, 21]. To the best of our knowledge, there is currently no tool that allows
101 matrix generation and visualization of SVs and CNVs classified under these schemas.
102 SigProfilerMatrixGenerator's capabilities for analyzing SVs and CNVs are implemented in
103 Python and the tool allows using multiple input formats, including segmentation and browser
104 extensible data paired-end (BEDPE) files generated by commonly used algorithms for detecting
105 copy number variations and structural variations, respectively. Additionally,
106 SigProfilerMatrixGenerator provides a comprehensive visualization of mutational patterns of
107 large mutational events and an R wrapper package for users that prefer working within the R
108 environment.
109

110 IMPLEMENTATION

111 Classification of Copy Number Variations

112 The schema for classifying copy number variations is based on Steele *et al.* [7] and it utilizes
113 allele-specific copy number, which quantifies the number of segments for each allele at each
114 variant loci rather than the total number of chromosome copies. In this schema, the copy-number
115 profile of a sample can be represented by a mutational vector with 48 dimensions. Specifically,
116 copy number segments are categorized into three heterozygosity states: heterozygous segments
117 with total copy number (TCN) of $A > 0, B > 0$ (numbers reflect the counts for major allele A and
118 minor allele B ; **Figure 1a**), segments with loss of heterozygosity (LOH) with total copy number
119 of $A > 0, B = 0$ (**Figure 1b**), and segments with homozygous deletions and TCN of $A = 0, B = 0$
120 (**Figure 1c**). Segments are further subclassified into 5 categories based on total copy number,
121 which reflects the sum of the copies on the major allele A and the copies on the minor allele B :
122 $TCN = 0, TCN = 1, TCN = 2, TCN = 3$ or $4, TCN = 5$ to $8, TCN \geq 9$. Each of these total copy
123 number states accounts for the phenomenon of whole-genome duplication, for example a diploid
124 ($TCN = 2$) state transitioning to a doubled state ($TCN = 4$), and a subsequent doubling of this state
125 to $TCN = 8$ is accounted for by the $TCN = 5-8$ category (**Figure 1a**). The categories for total copy
126 number have been chosen for biological relevance (**Figure 1**): $TCN = 0$ reflects homozygous
127 deletions, $TCN = 1$ represents a genomic deletion resulting in an LOH, $TCN = 2$ is equivalent to a
128 diploid state including copy neutral LOH (a phenomenon whereby one of two homologous
129 chromosomal regions is lost, but two identical copies of this region still remain; **Figure 1b**),
130 $TCN = 3$ or 4 reflect a gained state of tri- to tetra-ploidy, $TCN = 5$ to 8 represent a penta- to octo-
131 ploidy state, and $TCN \geq 9$ represents high-level amplifications such as ones found in samples
132 containing extrachromosomal DNA (ecDNA) [22]. Each of the heterozygous and LOH total

133 copy number categories are additionally subclassified into five additional categories based on the
134 size of their segments: 0 – 100kb, 100kb – 1Mb, 1Mb – 10Mb, 10Mb – 40Mb, and >40Mb.
135 Three size bins are used for the additional subcategorization of homozygous deletions: 0 –
136 100kb, 100kb – 1Mb, and >1Mb. The partitioning by segment sizes was chosen to ensure that a
137 sufficient proportion of segments are classified within each category [7]. This classification
138 allows summarizing copy number profiles using 48 distinct channels and can be represented
139 using a vector with 48 components. For example, a sample harboring multiple focal
140 amplifications, either contained on linear or extrachromosomal DNA, will have many events in
141 the 9+ total copy number category and the first 3 size bins (0 – 100kb, 100kb – 1Mb, 1Mb –
142 10Mb; **Figure 2a-b**). Conversely, a sample containing a large number of focal deletions or losses
143 of entire chromosomes or chromosome arms will have numerous events in the LOH category,
144 spanning all size bins (**Figure 2c-d**). Another example will be a sample with a whole-genome
145 doubling where copy number changes will primarily encompass segments with large genomic
146 sizes (10Mb – 40Mb; 40Mb) and total copy number between 3 and 4 (**Figure 2e-f**). Overall, this
147 48-channel classification schema can effectively summarize a diverse array of copy number
148 states seen across tumor types [7], whether they contain broad or focal events that result in
149 amplifications or deletions.

150

151 **Input Data for Classifying Copy Number Variations**

152 SigProfilerMatrixGenerator allows examining allele specific CNV data that, at a minimum,
153 include the following information for each CNV segment: chromosome, start coordinate, end
154 coordinate, and copy number of both the minor and major allele. Output files from the following
155 tools for detecting CNVs are automatically supported: ASCAT [23], ABSOLUTE [24],

156 Sequenza [25], FACETS [12], Battenberg [23], and PURPLE [26]. Additionally, custom
157 segmentation files from other CNV detection tools can be used if these files contain the
158 aforementioned information.

159

160 **Classification of Structural Variations**

161 A classification schema consisting of 32 features, based on Nik-Zainal *et al.* [21], is used to
162 construct a mutational vector with 32 dimensions for each sample. In principle, each structural
163 variant consists of two breakpoints which are at single-base resolution, where a breakpoint is
164 defined as a junction that indicates a structurally variable genomic segment greater than 50 base
165 pairs [10]. Breakpoints are typically detected using three signals from aligned sequencing reads:
166 depth of sequence coverage, discordant read-pairs, and split read-pairs [27-29]. Breakpoints can
167 also be detected via genome assembly, where reads are assembled into contigs, the contigs are
168 aligned to the reference genome, and these alignments are analyzed for structural variants [10].
169 The previously developed classification of structural variants considers the following canonical
170 SVs: tandem duplications, deletions, inversions, and translocations (**Figure 3**). A tandem
171 duplication refers to a segment of genomic material that has been duplicated and inserted on the
172 same chromosome adjacent to the original segment (**Figure 3a**). It should be noted that a tandem
173 duplication is not necessarily the same as a copy-number amplification. For example, ecDNA
174 copy-number amplifications are not tandem duplications as they are not inserted adjacent to the
175 original chromosome segment. A somatic deletion is an event that has removed a set of existing
176 base-pairs from a given location of a chromosome (**Figure 3b**). An inversion is when a segment
177 of the chromosome breaks off and reattaches at the same locus but in a reverse orientation
178 (**Figure 3c**). A translocation event occurs when a piece of one chromosome breaks off and some

179 (or all) fragments from that piece re-attach to either another chromosome or to a different locus
180 of the same chromosome (**Figure 3d**). The classification schema bins all SVs, apart from
181 translocations, according to the size of the event in base pairs: 0–10kb, 10kb–100kb, 100kb–
182 1Mb, 1Mb–10Mb, and >10Mb [21]. Translocations, which may involve more than one
183 chromosome, are not binned by size because they can be either balanced (where there is no net
184 loss of genetic material on the chromosomes involved and thus the size can be described by one
185 number) or unbalanced (where there is a net loss or gain of genetic material on the chromosomes
186 involved and thus the sizes of the segments cannot be described by just one number). Note that
187 whether a translocation is balanced or unbalanced is not considered in this classification schema.
188 The different types of SVs are then further divided into *clustered* and *non-clustered* events to
189 account for the non-random distribution of these events along the genome. Clustered events are
190 defined as events that occur closer to each other on a chromosome than purely expected by
191 chance. These clusters often arise as a result of complex events, such as chromothripsis [30] or
192 chromoplexy [31], generating many breakpoints in a single instantaneous event as opposed to the
193 gradual accumulation of events over many cell cycles which results in more dispersed non-
194 clustered events. Clusters of breakpoints can also form as a result of other mechanisms,
195 including, for example, rearrangement hotspots in the genome [32]. Clustering of SVs is
196 determined based on a previously developed algorithm that utilizes the Potts’ filter method [33].
197 This method segments a chromosome based on inter-mutational distance of SV breakpoints, and
198 if the average distance in a particular segment is less than 10 times the average inter-mutational
199 distance in the sample, all breakpoints in the segment are considered clustered. A minimum of 10
200 breakpoints must be present for a given segment to be considered clustered, otherwise all
201 breakpoints in that segment are considered non-clustered.

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

210

211 **Input Data for Classifying Structural Variants**

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

225 **DISCUSSION**

226 The newly developed version of SigProfilerMatrixGenerator allows transforming a set of
227 mutational catalogues of copy-number changes and structural rearrangements into matrices
228 amenable to decomposition, including, subsequent mutational signature analysis. The tool
229 provides support for two previously developed [7, 21] classification schemas for large mutational
230 events. Further, the tool also delivers an extensive plotting functionality that seamlessly
231 integrates with matrix generation to visualize the majority of output in a single analysis.
232 SigProfilerMatrixGenerator is the first tool to provide support for the 48 channel CNV schema
233 across a wide variety of popular tools for detecting CNV. Importantly, this schema can be
234 applied across several data modalities, including whole-genome sequencing, whole-exome
235 sequencing, RNA-sequencing, single-cell sequencing approaches, and genotyping microarrays.
236 In addition, SigProfilerMatrixGenerator is the first Python package that provides support for the
237 32 channel SV schema in a fast and intuitive manner with minimal preprocessing.

238

239 **CONCLUSION**

240 A breadth of computational tools exists for exploring the patterns for small mutational events,
241 including our initial implementation of SigProfilerMatrixGenerator [8]. However, to the best of
242 our knowledge, there are currently no tool for exploration and visualization of large mutational
243 events. We recently demonstrated that a classification of CNVs into 48 channels provides the
244 means to better elucidate and understand the mutational processes operative in human cancer [7].
245 Similarly, we and others have previously demonstrated that the classification of SVs into 32
246 channels can be used to understand the mutational processes giving rise to SVs across multiple
247 cancer types [19]. Our newly developed version of SigProfilerMatrixGenerator provides the
248 capability to examine these classification schemas from cancer genomics sequencing data. The

249 tool can scale to large datasets and will serve as foundation to future analysis of both mutational
250 patterns and mutational signatures of large mutational events.

251

252 AVAILABILITY AND REQUIREMENTS

253 **Project name:** SigProfilerMatrixGenerator

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

255 <https://github.com/AlexandrovLab/SigProfilerMatrixGeneratorR>

256 **Operating system(s):** Unix, Linux, and Windows

257 **Programming language:** Python 3 and R

258 **Other requirements:** None

259 **License:** BSD 2-Clause "Simplified" License

260 **Any restrictions to use by non-academics:** None

261

262 ABBREVIATIONS

263 **BEDPE:** browser extensible data paired-end

264 **CNV:** copy number variation

265 **DBS:** doublet base substitution

266 **ecDNA:** extrachromosomal DNA

267 **ID:** small insertions and deletions

268 **LOH:** loss of heterozygosity

269 **SBS:** single base substitution

270 **SV:** structural variation

271 **TCN:** total copy-number

272 **WGS:** whole-genome sequencing

273

274 DECLARATIONS

275 **Ethics approval and consent to participate:** Not applicable.

276 **Consent for publication:** Not applicable.

277 **Availability of data and materials:** Data sharing is not applicable to this article as no datasets were generated or analyzed during the current study.

279 **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.

284 **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 Cancer Research UK Grand Challenge Award C98/A24032. This work was also supported 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.

289 **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.

294 **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.

296

297 **FIGURE LEGENDS**

298 **Figure 1. Description of the Copy Number Classification Schema.** The copy number
299 classification schema consists of 48 mutually exclusive channels, divided by heterozygosity
300 status, segment size, and total copy number (TCN). **a)** In the heterozygous state, both alleles are
301 retained and either one or both alleles can be amplified. This amplification can be focal (top
302 panel) or it can encompass a chromosome or even the whole genome (bottom panel). The
303 heterozygous category is further subdivided based on TCN (TCN=1, TCN=2, TCN=3 or 4,
304 TCN=5 to 8, and TCN>=9). **b)** In a state of loss of heterozygosity (LOH), one of the alleles is
305 lost. The remaining allele can then be duplicated (i.e., copy neutral LOH), and undergo more
306 amplification resulting in higher total copy number states. The LOH category is further
307 subdivided based on TCN (TCN=0, TCN=1, TCN=2, TCN=3 or 4, TCN=5 to 8, and TCN>=9).
308 The heterozygous and LOH categories are further divided on the basis of the size of the segment:
309 0 – 100kb, 100kb – 1Mb, 1Mb – 10Mb, 10Mb – 40Mb, >40Mb. High-level LOH or
310 heterozygous amplifications (e.g., TCN=5 to 8 or TCN>= 9) can be carried on
311 extrachromosomal DNA (depicted as red circles) as well as on linear chromosomes. **c)**
312 Homozygous deletions result in the loss of both alleles, and are divided on the basis of the size of
313 the deleted segment: 0 – 100kb, 100kb – 1Mb, and >1Mb.

314

315 **Figure 2. Converting Copy Number Segmentation Profiles into Copy Number Mutational**
316 **Vectors.** The CNV classification schema converts a sample's segmentation profile (**a, c, e**) into a
317 count vector of 48 mutually exclusive components (**b, d, f**). These components are based on
318 segment size, heterozygosity status, and total copy number. A breast cancer sample with many
319 highly amplified segments, possibly due to the presence of extrachromosomal DNA, is shown in

320 (a, b). This sample's count vector is characterized by peaks in the 5-8 and 9+ total copy number
321 categories. A gastric cancer sample with extensive loss of heterozygosity is shown in (c, d). This
322 sample's count vector is characterized by peaks in the LOH category, specifically with a total
323 copy number of 1 indicating a loss of an allele. A sarcoma sample with a whole-genome
324 duplication event, characterized by peaks in the 3-4 total copy number category and the 40+ Mb
325 size bin, is shown in (e, f).

326

327 **Figure 3. Description of the Structural Variant Classification Schema.** Structural variants
328 (SVs) are categorized as tandem-duplications, deletions, inversions, or translocations. **a)** Tandem
329 duplication of a segment containing the A allele. A tandem duplication occurs when a segment is
330 duplicated and inserted adjacent to the original chromosomal segment. **b)** Deletion of the
331 segment containing the A allele. A deletion occurs when there is a loss of genetic material from a
332 chromosome. **c)** An inversion of the segment containing the B allele. An inversion occurs when a
333 segment breaks off and reattaches in a reverse orientation within the same chromosome. **d)** A
334 translocation of a chromosomal segment. A translocation event occurs when a piece of one
335 chromosome breaks off and some (or all) fragments from that piece re-attach to either another
336 chromosome or to a different locus of the same chromosome.

337

338 **Figure 4. Classifying Structural Variants into Mutational Vectors.** **a)** An example of a bone
339 cancer sample from PCAWG with a highly rearranged genome consisting of both clustered and
340 non-clustered structural variants (SVs) is shown as a Circos plot representation. **b)** Zooming into
341 SVs specifically found on chromosome 12 in the bone cancer sample. SVs are shown as a linear
342 representation (top) and as a rainfall plot (bottom). The rainfall plot depicts all breakpoints on

343 chromosome 12 according to their genomic coordinate (x-axis) and the \log_{10} inter-mutational
344 distance (y-axis), which is the distance to the breakpoint immediately preceding it. The tendency
345 of breakpoints to cluster in a specific genomic region on chromosome 12 due to a chromothripsis
346 event is evident in all representations. **c)** Zooming into SVs specifically found on chromosome 8
347 in the bone cancer sample. SVs are shown as a linear representation (top) and as a rainfall plot
348 (bottom). The rainfall plot depicts all breakpoints on chromosome 8 according to their genomic
349 coordinate (x-axis) and the \log_{10} inter-mutational distance (y-axis), which is the distance to the
350 breakpoint immediately preceding it. There are no clustered SVs on chromosome 8 as, per the
351 SV classification schema, clustering requires a minimum of 10 breakpoints in a segment of a
352 chromosome. **d)** The SV classification schema is applied to the SVs found on chromosome 12 in
353 the bone cancer sample. SVs are classified by the event type (denoted by color) and are binned
354 according to the size of the event (0 – 10kb, 10kb – 100kb, 100kb – 1Mb, 1Mb – 10Mb, and
355 >10Mb). **e)** The SV classification schema is applied to the SVs found on chromosome 8 in the
356 bone cancer sample. SVs are classified by the event type (denoted by color) and are binned
357 according to the size of the event (0 – 10kb, 10kb – 100kb, 100kb – 1Mb, 1Mb – 10Mb, and
358 >10Mb).
359

360 REFERENCES

- 361 1. Cancer Genome Atlas Research N, Weinstein JN, Collisson EA, Mills GB, Shaw KR,
362 Ozenberger BA, Ellrott K, Shmulevich I, Sander C, Stuart JM: **The Cancer Genome Atlas**
363 **Pan-Cancer analysis project**. *Nat Genet* 2013, **45**(10):1113-1120.
- 364 2. Consortium ITP-CAoWG: **Pan-cancer analysis of whole genomes**. *Nature* 2020,
365 **578**(7793):82-93.
- 366 3. Alexandrov LB, Kim J, Haradhvala NJ, Huang MN, Tian Ng AW, Wu Y, Boot A, Covington
367 KR, Gordenin DA, Bergstrom EN *et al*: **The repertoire of mutational signatures in human**
368 **cancer**. *Nature* 2020, **578**(7793):94-101.
- 369 4. Islam SMA, Diaz-Gay M, Wu Y, Barnes M, Vangara R, Bergstrom EN, He Y, Vella M, Wang
370 J, Teague JW *et al*: **Uncovering novel mutational signatures by de novo extraction with**
371 **SigProfilerExtractor**. *Cell Genom* 2022, **2**(11):None.
- 372 5. Alexandrov LB, Nik-Zainal S, Wedge DC, Campbell PJ, Stratton MR: **Deciphering**
373 **signatures of mutational processes operative in human cancer**. *Cell Rep* 2013, **3**(1):246-
374 259.
- 375 6. Li Y, Roberts ND, Wala JA, Shapira O, Schumacher SE, Kumar K, Khurana E, Waszak S,
376 Korbel JO, Haber JE *et al*: **Patterns of somatic structural variation in human cancer**
377 **genomes**. *Nature* 2020, **578**(7793):112-121.
- 378 7. Steele CD, Abbasi A, Islam SMA, Bowes AL, Khandekar A, Haase K, Hames-Fathi S, Ajayi
379 D, Verfaillie A, Dhimi P *et al*: **Signatures of copy number alterations in human cancer**.
380 *Nature* 2022, **606**(7916):984-991.
- 381 8. Bergstrom EN, Huang MN, Mahto U, Barnes M, Stratton MR, Rozen SG, Alexandrov LB:
382 **SigProfilerMatrixGenerator: a tool for visualizing and exploring patterns of small**
383 **mutational events**. *BMC Genomics* 2019, **20**(1):685.
- 384 9. Cameron DL, Di Stefano L, Papenfuss AT: **Comprehensive evaluation and**
385 **characterisation of short read general-purpose structural variant calling software**. *Nat*
386 *Commun* 2019, **10**(1):3240.
- 387 10. Cosenza MR, Rodriguez-Martin B, Korbel JO: **Structural Variation in Cancer: Role,**
388 **Prevalence, and Mechanisms**. *Annu Rev Genomics Hum Genet* 2022, **23**:123-152.
- 389 11. Talevich E, Shain AH, Botton T, Bastian BC: **CNVkit: Genome-Wide Copy Number**
390 **Detection and Visualization from Targeted DNA Sequencing**. *PLoS Comput Biol* 2016,
391 **12**(4):e1004873.
- 392 12. Shen R, Seshan VE: **FACETS: allele-specific copy number and clonal heterogeneity**
393 **analysis tool for high-throughput DNA sequencing**. *Nucleic Acids Res* 2016,
394 **44**(16):e131.
- 395 13. Serin Harmanci A, Harmanci AO, Zhou X: **CaSpER identifies and visualizes CNV events by**
396 **integrative analysis of single-cell or bulk RNA-sequencing data**. *Nat Commun* 2020,
397 **11**(1):89.
- 398 14. Degasperi A, Zou X, Amarante TD, Martinez-Martinez A, Koh GCC, Dias JML, Heskin L,
399 Chmelova L, Rinaldi G, Wang VYW *et al*: **Substitution mutational signatures in whole-**
400 **genome-sequenced cancers in the UK population**. *Science* 2022, **376**(6591).

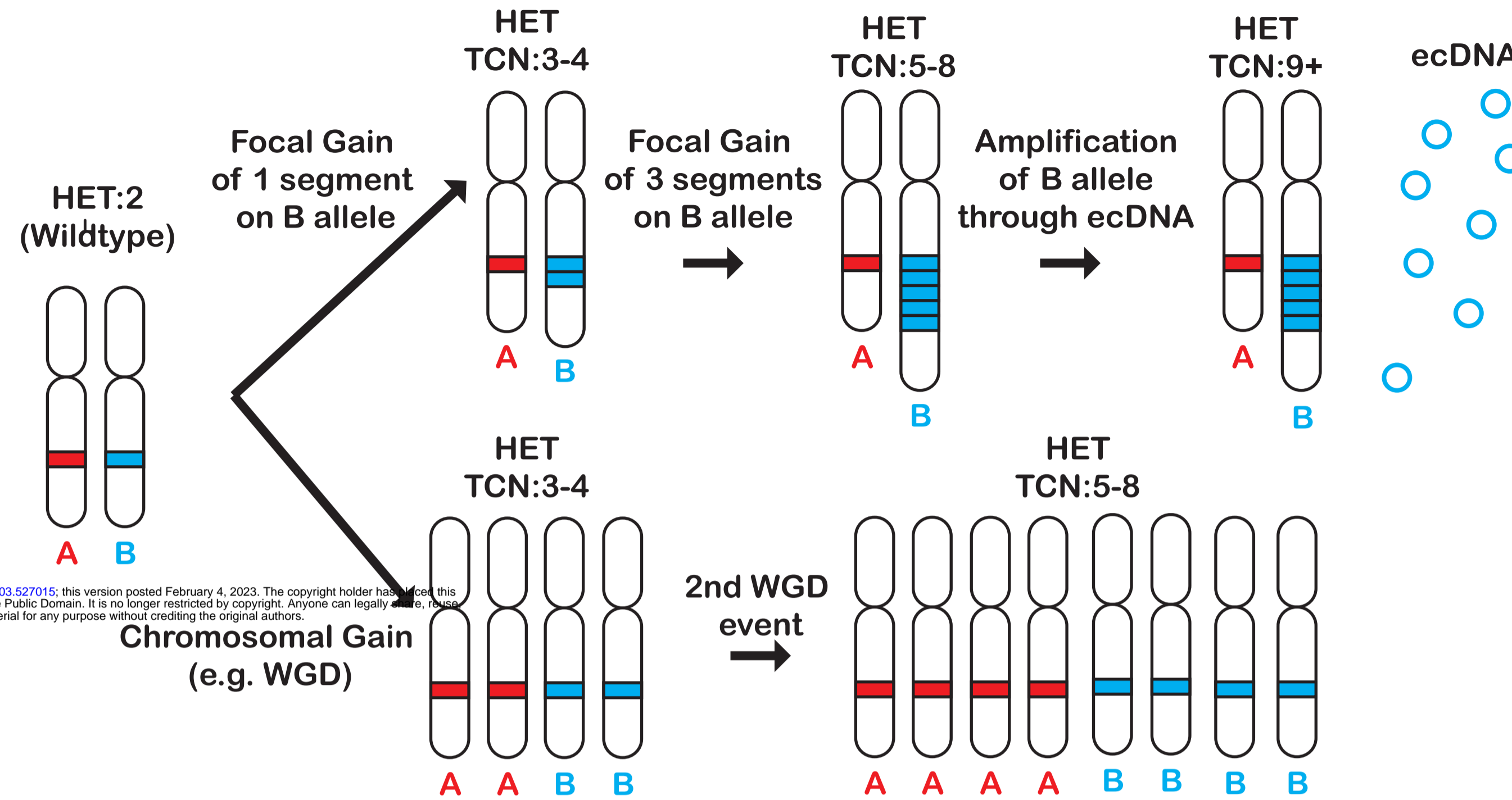
- 401 15. Alexandrov LB, Ju YS, Haase K, Van Loo P, Martincorena I, Nik-Zainal S, Totoki Y,
402 Fujimoto A, Nakagawa H, Shibata T *et al*: **Mutational signatures associated with**
403 **tobacco smoking in human cancer**. *Science* 2016, **354**(6312):618-622.
- 404 16. Petljak M, Alexandrov LB, Brammell JS, Price S, Wedge DC, Grossmann S, Dawson KJ, Ju
405 YS, Iorio F, Tubio JMC *et al*: **Characterizing Mutational Signatures in Human Cancer Cell**
406 **Lines Reveals Episodic APOBEC Mutagenesis**. *Cell* 2019, **176**(6):1282-1294 e1220.
- 407 17. Riva L, Pandiri AR, Li YR, Droop A, Hewinson J, Quail MA, Iyer V, Shepherd R, Herbert RA,
408 Campbell PJ *et al*: **The mutational signature profile of known and suspected human**
409 **carcinogens in mice**. *Nat Genet* 2020, **52**(11):1189-1197.
- 410 18. Alexandrov LB, Jones PH, Wedge DC, Sale JE, Campbell PJ, Nik-Zainal S, Stratton MR:
411 **Clock-like mutational processes in human somatic cells**. *Nat Genet* 2015, **47**(12):1402-
412 1407.
- 413 19. Degasperi A, Amarante TD, Czarnecki J, Shooter S, Zou X, Glodzik D, Morganella S, Nanda
414 AS, Badja C, Koh G: **A practical framework and online tool for mutational signature**
415 **analyses show intertissue variation and driver dependencies**. *Nature cancer* 2020,
416 **1**(2):249-263.
- 417 20. Drews RM, Hernando B, Tarabichi M, Haase K, Lesluyes T, Smith PS, Morrill Gavarro L,
418 Couturier DL, Liu L, Schneider M *et al*: **A pan-cancer compendium of chromosomal**
419 **instability**. *Nature* 2022, **606**(7916):976-983.
- 420 21. Nik-Zainal S, Davies H, Staaf J, Ramakrishna M, Glodzik D, Zou X, Martincorena I,
421 Alexandrov LB, Martin S, Wedge DC *et al*: **Landscape of somatic mutations in 560 breast**
422 **cancer whole-genome sequences**. *Nature* 2016, **534**(7605):47-54.
- 423 22. Kim H, Nguyen NP, Turner K, Wu S, Gujar AD, Luebeck J, Liu J, Deshpande V, Rajkumar U,
424 Namburi S *et al*: **Extrachromosomal DNA is associated with oncogene amplification**
425 **and poor outcome across multiple cancers**. *Nat Genet* 2020, **52**(9):891-897.
- 426 23. Van Loo P, Nordgard SH, Lingjaerde OC, Russnes HG, Rye IH, Sun W, Weigman VJ,
427 Marynen P, Zetterberg A, Naume B *et al*: **Allele-specific copy number analysis of**
428 **tumors**. *Proc Natl Acad Sci U S A* 2010, **107**(39):16910-16915.
- 429 24. Carter SL, Cibulskis K, Helman E, McKenna A, Shen H, Zack T, Laird PW, Onofrio RC,
430 Winckler W, Weir BA *et al*: **Absolute quantification of somatic DNA alterations in**
431 **human cancer**. *Nat Biotechnol* 2012, **30**(5):413-421.
- 432 25. Favero F, Joshi T, Marquard AM, Birkbak NJ, Krzystanek M, Li Q, Szallasi Z, Eklund AC:
433 **Sequenza: allele-specific copy number and mutation profiles from tumor sequencing**
434 **data**. *Ann Oncol* 2015, **26**(1):64-70.
- 435 26. Shale C, Cameron DL, Baber J, Wong M, Cowley MJ, Papenfuss AT, Cuppen E, Priestley P:
436 **Unscrambling cancer genomes via integrated analysis of structural variation and copy**
437 **number**. *Cell Genomics* 2022, **2**(4):100112.
- 438 27. Wala JA, Bandopadhyay P, Greenwald NF, O'Rourke R, Sharpe T, Stewart C,
439 Schumacher S, Li Y, Weischenfeldt J, Yao X *et al*: **SvABA: genome-wide detection of**
440 **structural variants and indels by local assembly**. *Genome Res* 2018, **28**(4):581-591.
- 441 28. Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Kallberg M, Cox AJ, Kruglyak
442 S, Saunders CT: **Manta: rapid detection of structural variants and indels for germline**
443 **and cancer sequencing applications**. *Bioinformatics* 2016, **32**(8):1220-1222.

- 444 29. Rausch T, Zichner T, Schlattl A, Stutz AM, Benes V, Korbel JO: **DELLY: structural variant**
445 **discovery by integrated paired-end and split-read analysis.** *Bioinformatics* 2012,
446 **28(18):i333-i339.**
- 447 30. Stephens PJ, Greenman CD, Fu B, Yang F, Bignell GR, Mudie LJ, Pleasance ED, Lau KW,
448 Beare D, Stebbings LA *et al*: **Massive genomic rearrangement acquired in a single**
449 **catastrophic event during cancer development.** *Cell* 2011, **144(1):27-40.**
- 450 31. Shen MM: **Chromoplexy: a new category of complex rearrangements in the cancer**
451 **genome.** *Cancer Cell* 2013, **23(5):567-569.**
- 452 32. Glodzik D, Morganella S, Davies H, Simpson PT, Li Y, Zou X, Diez-Perez J, Staaf J,
453 Alexandrov LB, Smid M *et al*: **A somatic-mutational process recurrently duplicates**
454 **germline susceptibility loci and tissue-specific super-enhancers in breast cancers.** *Nat*
455 *Genet* 2017, **49(3):341-348.**
- 456 33. Winkler G, Liebscher V: **Smoothers for discontinuous signals.** *Journal of Nonparametric*
457 *Statistics* 2002, **14(1-2):203-222.**
458

Figure 1.

a)

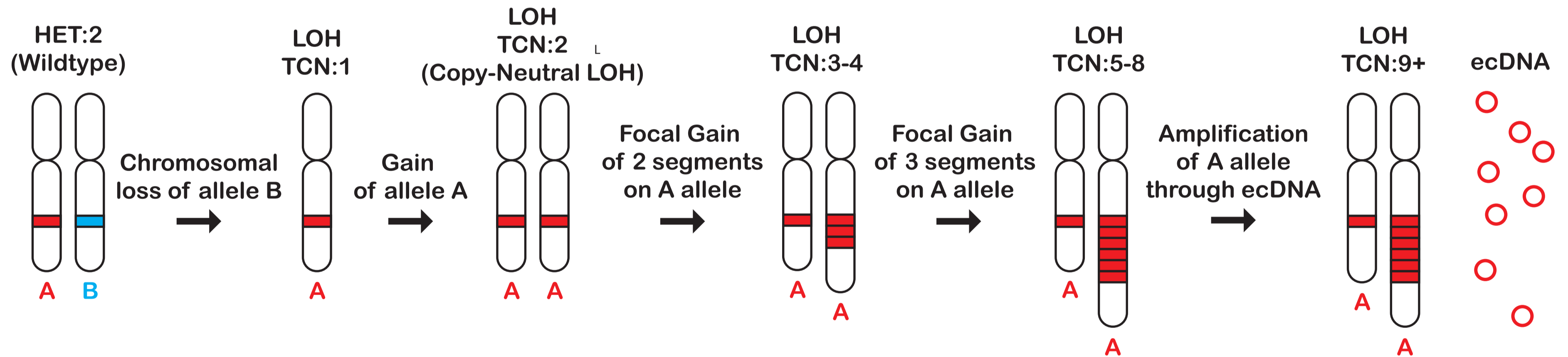
Heterozygous CN
{A>0, B>0}



bioRxiv preprint doi: <https://doi.org/10.1101/2023.02.03.527015>; this version posted February 4, 2023. The copyright holder for this preprint (which was not certified by peer review) in the Public Domain. It is no longer restricted by copyright. Anyone can legally share, reuse, remix, or adapt this material for any purpose without crediting the original authors.

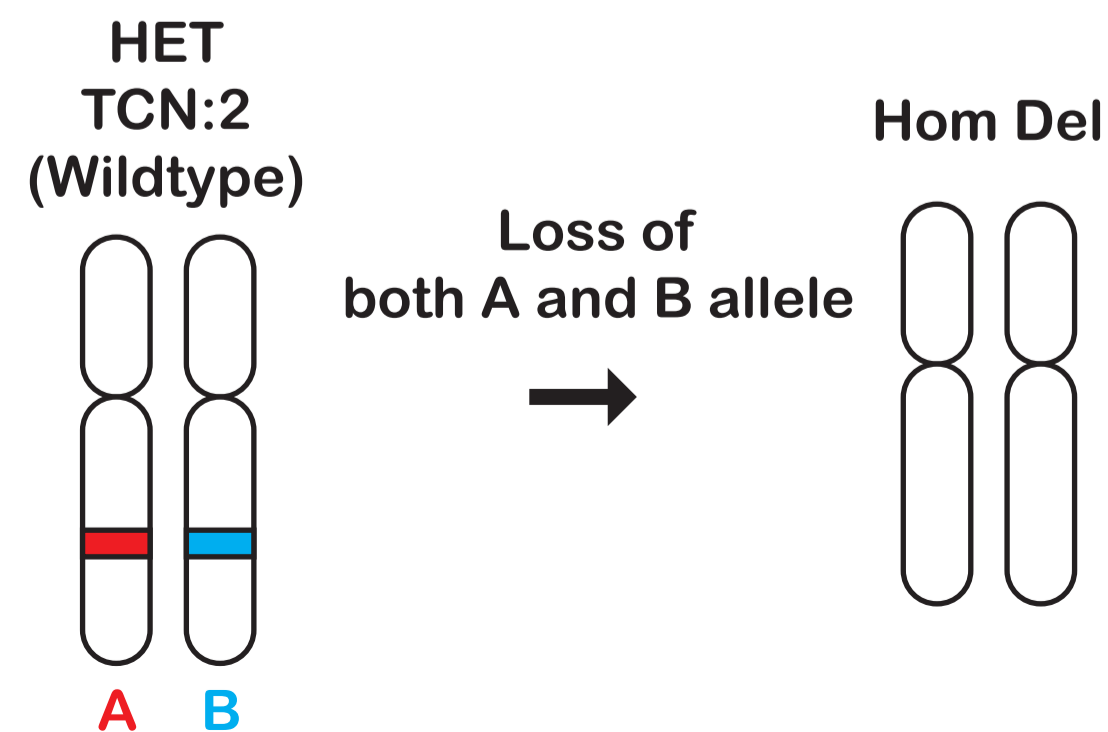
b)

Loss of Heterozygosity
{A>0, B=0}



c)

Homozygous Deletion
{A=0, B=0}



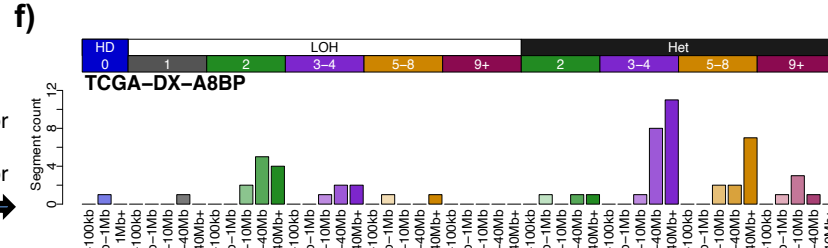
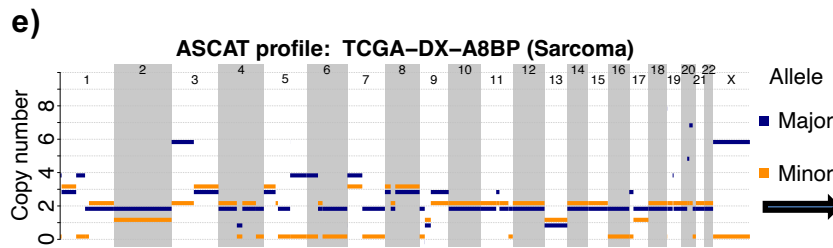
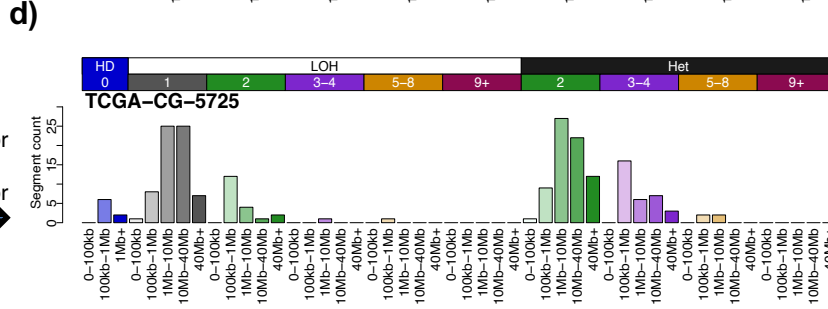
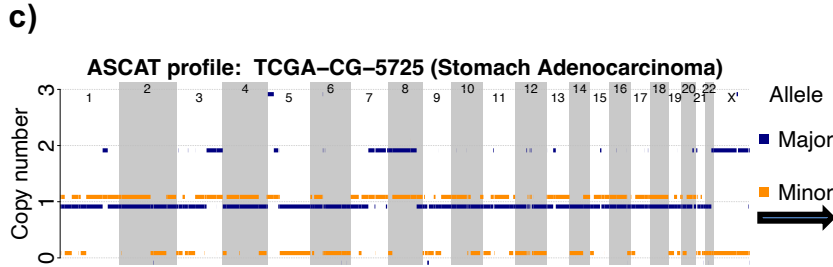
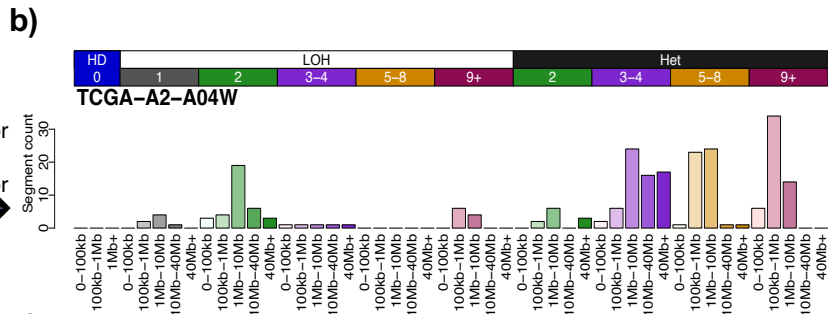
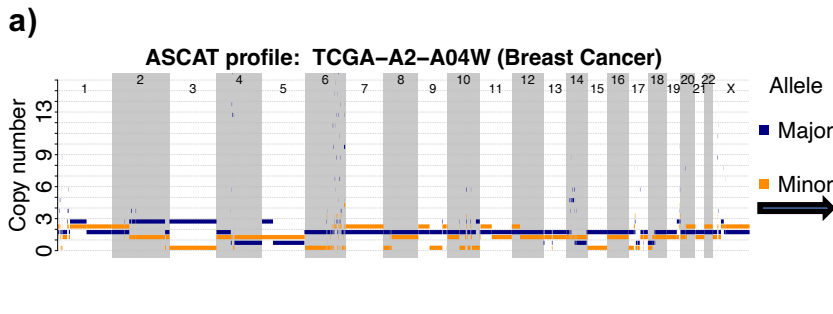
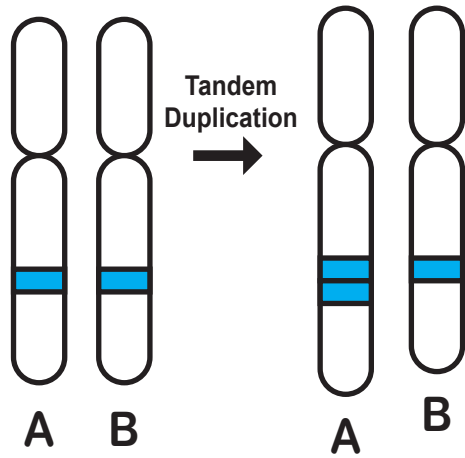
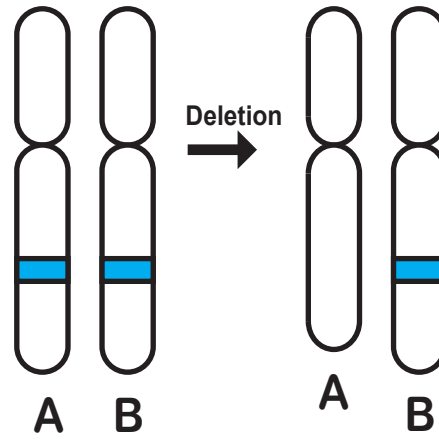


Figure 3.

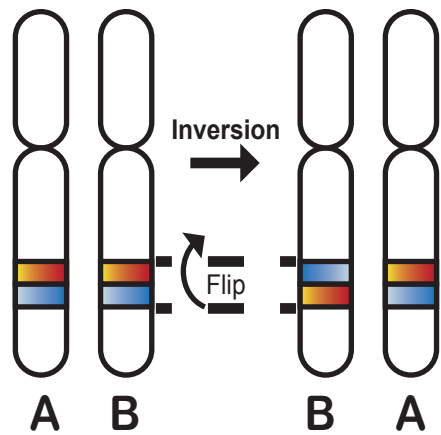
a)



b)



c)



d)

