1    **MutationalPatterns: The one stop shop for the analysis of mutational processes**

2    Freek Manders[1,2], Arianne M. Brandsma[1,2], Jurrian de Kanter[1,2], Mark Verheul[1,2], Rurika

3    Oka[1,2], Markus J. van Roosmalen[1,2], Bastiaan van der Roest[2,3,4], Arne van Hoeck[2,3], Edwin

4    Cuppen[2,3] and Ruben van Boxtel[1,2,*]

5    [1]Princess Máxima Center for Pediatric Oncology, Heidelberglaan 25, 3584CS Utrecht, The

6    Netherlands

7    [2]Oncode Institute, Jaarbeursplein 6, 3521 AL Utrecht, The Netherlands

8    [3]Center for Molecular Medicine, University Medical Center Utrecht, Utrecht University,

9    Universiteitsweg 100, 3584, CG, Utrecht, The Netherlands

10   [4]Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht,

11   Utrecht University, Universiteitsweg 100, 3584, CG, Utrecht, The Netherlands

12   [*]Corresponding author: R.vanBoxtel@prinsesmaximacentrum.nl

13

14

## Abstract

### Background

The collective of somatic mutations in a genome represents a record of mutational processes that have been operative in a cell. These processes can be investigated by extracting relevant mutational patterns from sequencing data.

### Results

Here, we present the next version of MutationalPatterns, an R/Bioconductor package, which allows in-depth mutational analysis of catalogues of single and double base substitutions as well as small insertions and deletions. Major features of the package include the possibility to perform regional mutation spectra analyses and the possibility to detect strand asymmetry phenomena, such as lesion segregation. On top of this, the package also contains functions to determine how likely it is that a signature can cause damaging mutations (i.e., mutations that affect protein function). This updated package supports stricter signature refitting on known signatures in order to prevent overfitting. Using simulated mutation matrices containing varied signature contributions, we showed that reliable refitting can be achieved even when only 50 mutations are present per signature. Additionally, we incorporated bootstrapped signature refitting to assess the robustness of the signature analyses. Finally, we applied the package on genome mutation data of cell lines in which we deleted specific DNA repair processes and on large cancer datasets, to show how the package can be used to generate novel biological insights.

**Conclusions**

This novel version of MutationalPatterns allows for more comprehensive analyses and visualization of mutational patterns in order to study the underlying processes. Ultimately, in-depth mutational analyses may contribute to improved biological insights in mechanisms of mutation accumulation as well as aid cancer diagnostics. MutationalPatterns is freely available at http://bioconductor.org/packages/MutationalPatterns.

**Keywords**

R, regional mutation patterns, mutagenic processes, mutational signatures, indels, base substitutions, somatic mutations

**Background**

Mutational landscapes in the genomes of cells are the result of a balance between mutagenic and DNA-repair processes (1). The somatic mutations that shape these landscapes gradually accumulate throughout life in both healthy and malignant cells (2,3). As a result, the complete collection of somatic mutations in the genome of a cell forms a record of the mutational processes that have been active throughout the life of that cell. In-depth analyses of somatic mutations can allow us to better understand the mutational processes that caused them (4).

3

57    First, such analyses can provide insight into the etiology of cancer by identifying mutagenic

58    exposures, which ultimately contribute to the accumulation of cancer driving mutations. For

59    example, we recently identified a mutational pattern caused by a carcinogenic strain

60    of *Escherichia coli* found in the gut of ~20% of healthy individuals (5). This pattern matched

61    mutations found in colorectal cancer driver genes, indicating a direct role in tumorigenesis.

62    Mutational patterns have been systematically determined *in vitro* for many environmental

63    mutagenic agents, which can be used to deduce cancer causes (6). The effects of such

64    agents can also be found *in vivo*. For example, we recently found mutations caused by

65    exposure to the antiviral drug ganciclovir, which patients received to treat a viral infection

66    after a hematopoietic stem cell transplant (7). Second, studying mutational processes can

67    be useful for improved cancer diagnostics. For example, the presence of certain mutational

68    signatures can be used as a functional readout for deficiency of  homologous recombination

69    (HR)-mediated double strand break repair (8,9). Cancers with a defect in this repair pathway

70    are selectively sensitive to poly(ADP-ribose) polymerase (PARP) inhibitors, providing a

71    targeted therapy for the patients (10,11).

72         One of the most popular tools to analyze somatic mutation profiles is the

73    R/Bioconductor package MutationalPatterns, which can be used to easily investigate

74    mutation spectra (12–19). It can also be used to identify new signatures in mutation data

75    using Nonnegative Matrix Factorization (NMF) and to determine the contribution of

76    previously defined signatures to a sample using a method known as "signature refitting" (4).

77    However, the original version of this package has several limitations. First, the package is

78    limited to single base substitutions (SBSs) and cannot be used for small insertions and

79    deletions (indels) or double base substitutions (DBSs) even though signatures for these

80    mutation types have recently been identified in large pan-cancer sequencing efforts (13).

81    The package also suffers from signature overfitting when determining the contribution of

82    known patterns to a sample, which can result in too many signatures being attributed (20).

83    Additionally, the package only allows for analyzing spectra for mutations in the entire

84    genome, making it difficult to study the involvement of specific genomic elements, such as

85    enhancers or secondary hairpin structures. The ability to investigate the role of such

86    elements in mutation accumulation is important, because this allows for identifying the

87    molecular mechanisms by which certain processes induce mutagenesis (21–23).

88    Here we present a novel, almost completely rewritten version of MutationalPatterns for the

89    analysis of mutational processes, which is easy-to-use and contains many new features,

90    such as DNA lesion segregation (24). Existing features have also been improved, resulting in

91    a very comprehensive package that can be used for both basic and more advanced

92    mutational pattern analyses. MutationalPatterns supports DBSs, multi base substitutions

93    (MBSs) and indels, and can automatically extract all these mutation types from a single

94    variant call format (VCF) file. The package can generate region specific spectra and signature

95    contributions to study the varying activities of mutational processes across the genome. The

96    package also generates more accurate results by supporting stricter signature refitting. This

97    refitting can also be bootstrapped to determine the confidence of the results. Additionally, a

98    process known as lesion segregation can be investigated.

99        The MutationalPatterns package can be used to generate novel biological insights,

100   which we demonstrate by applying it to whole genome sequencing (WGS) data obtained

101   from a lymphoblastoid cell line, in which specific DNA repair processes were deleted using

102   CRISPR-Cas9 genome editing, as well as by applying the package on large cancer datasets.

5

103    Additionally, we demonstrate that the package scales well on these large datasets. Finally,

104    we show the improved accuracy of the stricter signature refitting using simulated data.

105

106    **Implementation**

107    *Mutation profiles*

108    MutationalPatterns supports SBSs, DBSs, MBSs and indels. Multiple mutation types are

109    allowed to be present in a single VCF file so that users do not have to split them beforehand.

110    A specific mutation type can be selected as an argument with the "read_vcfs_as_granges"

111    function when reading in the VCF files. Alternatively, the "get_mut_type" function can be

112    used on data that is already loaded in the memory.

113    DBS and MBS variants can be called by various variant callers, such as the Genome

114    Analysis ToolKit (GATK) Mutect2, in two different ways (25). The variants can be called

115    explicitly as DBS and MBS variants or as neighboring SBSs. A downside of the first approach

116    is that neighboring germline and somatic mutations can be called as a single combined DBS

117    or MBS, because the variants are compared to the reference instead of the control sample.

118    MutationalPatterns supports both approaches. When the second approach is used,

119    neighboring SBSs will be merged into somatic DBS or MBS variants.

120    Because they get merged, DBS and MBS variants are no longer incorrectly identified

121    as separate SBSs by MutationalPatterns. This improves the quality of the SBS profiles, as

122    DBS and MBS mutations often have a very different context on account of them being

123    caused by different processes (13) (Additional file 1: Figure S1).

124

125    The COSMIC contexts of SBS, indel and DBS variants can be retrieved with fast vectorized

126    functions, namely "mut_context", "get_indel_context" and "get_dbs_context". The context

127    of SBS variants consisted of its direct 5' and 3' bases in the original package. These contexts

128    were chosen because they are generally the most informative and adding more bases

129    drastically increases the feature space, leading to sparsity (4). Indeed, adding only one extra

130    base to both the upstream and downstream context increases the number of features from

131    96 to 1536. However, with the increasing availability of large sequencing cohorts such large

132    feature spaces have become more manageable, making it easier to examine nucleotide

133    preference more upstream or downstream of the mutated base. Therefore,

134    MutationalPatterns' users can now choose any context size for SBSs.

135

136    The mutation contexts can be used for custom analyses. Alternatively, the number of

137    mutations per context can be counted, resulting in a count matrix, where each row is a

138    context and each column a sample. These matrices are created with the "mut_matrix",

139    "mut_matrix_stranded", "count_indel_contexts", "count_dbs_contexts" and

140    "count_mbs_contexts" functions. The "count_mbs_contexts" function uses the length of

141    the MBSs, because to date no COSMIC consensus has been defined.

142         The count matrices can be plotted as spectra or profiles for all the mutation types

143    (Fig. 1a, b, c). The SBS spectra can be displayed in the individual samples. Additionally, the

144    error bars can be displayed as standard deviation, 95% confidence interval (CI) and the

145    standard error of the mean. A count matrix with a larger context can be visualized using the

146    new "plot_profile_heatmap" or "plot_river" functions (Fig. 1d, Additional file 1: Figure S2).

147    This last function can be especially helpful to provide a quick overview of a mutation

7

148    spectrum with a wider context. Next to visualizing them, a count matrix can also be used for

149    downstream analyses, such as a *de novo* extraction of mutational signatures. In some cases,

150    it can be useful to pool multiple samples within a count matrix to increase statistical power.

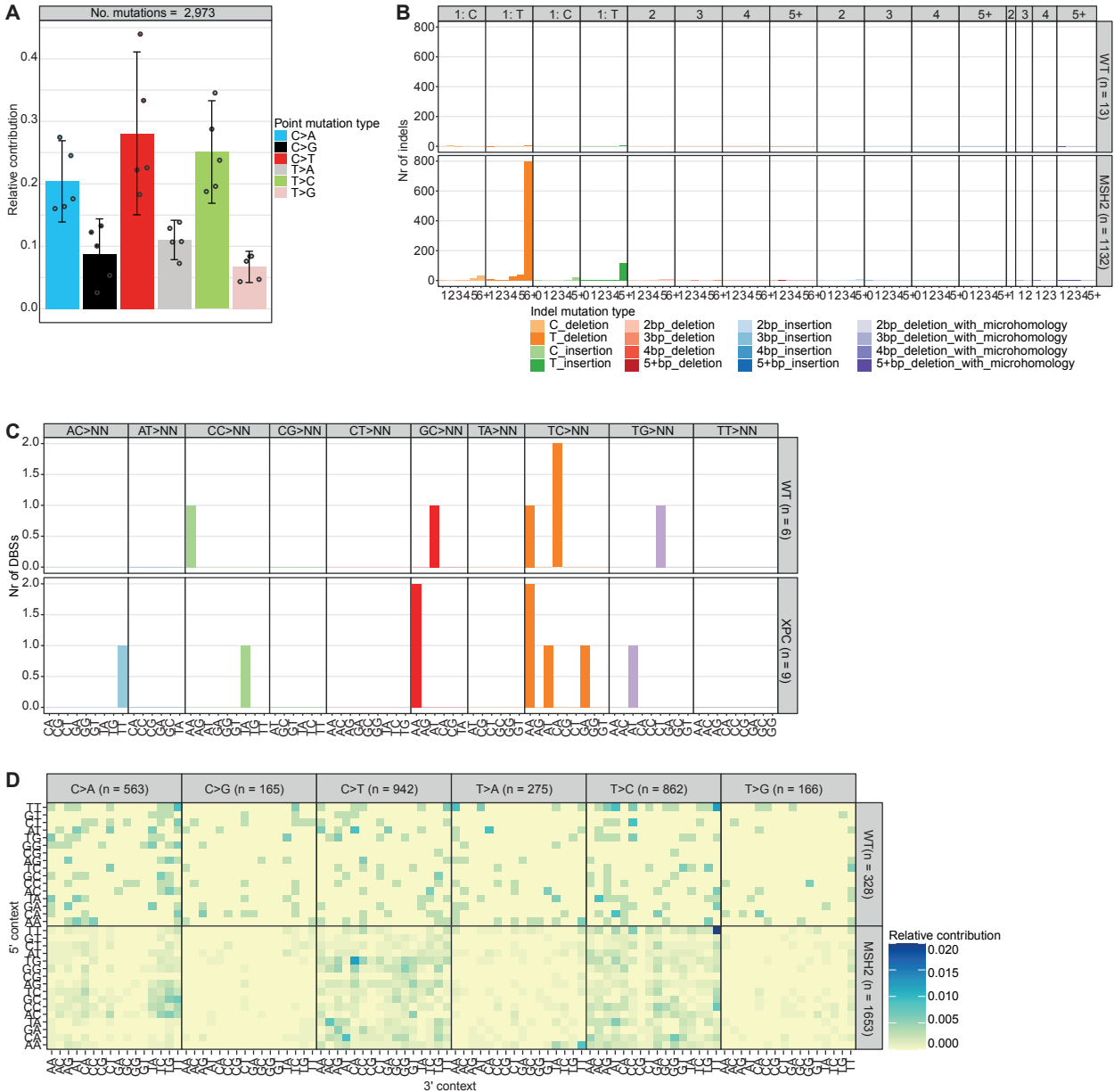151    This can be done using the new "pool_mut_mat" function.



152

153    Fig. 1 Mutation profiles can be made for multiple mutation types

154    **a** Relative contribution of the indicated mutation types to the point mutation spectrum.

155    Bars depict the mean relative contribution of each mutation type over all the samples and

8

156    error bars indicate the 95% confidence interval. The dots show the relative contributions of

157    the individual samples. The total number of somatic point mutations per tissue is indicated.

158    **b** Absolute contribution of the indicated mutation types to the indel spectrum for the wild-

159    type (WT) and *MSH2* knockout. The total number of indels per sample is indicated. **c**

160    Absolute contribution of the indicated mutation types to the DBS spectrum for the wild-type

161    (WT) and *XPC* knockout. The total number of DBSs per sample is indicated. **d** Heatmap

162    depicting the relative contribution of the indicated mutation types and the surrounding

163    bases to the point mutation spectrum for the WT and *MSH2* knockout. The total number of

164    somatic point mutations per tissue is indicated.

165

166    *Region specific analyses*

167    Mutational processes can be influenced by regional genomic features at multiple scales,

168    such as chromatin landscape, secondary hairpin structures as well as the major and minor

169    groove of the DNA (21–23). With the previous version of MutationalPatterns, it was possible

170    to test for enrichment and/or depletion of the mutation load in such regions. However, the

171    package lacked the possibility to automatically correct for multiple testing. In addition,

172    mutational profiles in genomic regions could not be easily assessed. In MutationalPatterns,

173    multiple testing correction is now automatically performed when testing for enrichment and

174    depletion. In addition, multiple significance levels are now supported, which can be

175    visualized using one or multiple asterisks. Furthermore, regional mutation profiles can be

176    determined in detail. This is done by first splitting mutations based on pre-defined genomic

177    regions, with the new "split_muts_region" function, which requires a GRanges or

178    GRangesList object containing chromosome coordinates as its input. These coordinates can

179    be read into R from file types like ".txt" or ".bed" files or they can be directly read from

180    databases, such as Ensembl (26). This analysis can be performed for multiple samples and

181    multiple types of regions at once. A user could, for example, split a set of mutations into

182    "promoter", "enhancer" and "other" mutations.

183        Splitting the mutations according to different genomic regions results in a

184    GRangesList containing sample/region combinations. These combinations can be treated as

185    separate samples by, for example, performing *de novo* signature analysis to identify

186    processes that are specifically active in certain genomic regions. Knowing in which regions a

187    signature is predominantly present, can lead to a better understanding of its etiology.

188    Instead of treating the sample/region combinations as separate samples, the genomic

189    regions can also be incorporated into the mutational contexts, using the new

190    "lengthen_mut_matrix" function. This means that a mutational context like "A[C>A]A" could

191    be split into "A[C>A]A-promoter" and "A[C>A]A-enhancer". This analysis allows users to

192    generate signatures that contain different mutation contexts in different genomic regions.

193    Such signatures could be more specific than the regular COSMIC signatures.

194    Region-specific mutation spectra can be visualized with the "plot_spectrum_region"

195    function, which contains the same arguments as the "plot_spectrum" function (Fig. 2a, b). In

196    addition, region-specific 96-channel mutation profiles can be visualized with the new

197    "plot_profile_region" function, which contains the same arguments as the "plot_96_profile"

198    function (Fig. 2c). Both the "plot_spectrum_region" and "plot_profile_region" functions

199    contain a "mode" argument, which allows users to normalize for the occurrence of the

200    different mutation types per sample/region combination, per sample, or not at all.

201    Instead of using pre-determined genomic regions, it is also possible to compare the

202    mutation spectra of regions with different mutation densities. These regions can be

203    identified using the new "bin_mutation_density" function.

204    Regional mutational patterns can also be investigated using an unsupervised

205    approach, which is unique to MutationalPatterns, with the new

206    "determine_regional_similarity" function. This function uses a sliding window approach to

207    calculate the cosine similarity between the global mutation profile and the mutation profile

208    of smaller genomic windows, allowing for the unbiased identification of regions with a

209    mutation profile, that differs from the rest of the genome. Users can correct for the

210    oligonucleotide frequency of the genomic windows using the "oligo_correction" argument.

211    The function returns an S4 object, containing the genomic windows with their associated

212    cosine similarities and the settings used to run the function. Because of the unbiased

213    approach of this function, it works best on a large dataset containing at least 100,000

214    substitutions. The result of this analysis can be visualized using the new

215    "plot_regional_similarity" function.

216

217    *Lesion segregation*

218    Mutation spectra sometimes contain Watson versus Crick strand asymmetries (24). These

219    asymmetries can be the result of many DNA lesions occurring during a single cell cycle. If

220    these lesions are not properly repaired before the next genome duplication, then the

221    resulting sister chromatids will segregate into different daughter cells, which will each

222    inherit the lesions on opposite strands. This process is known as lesion segregation (24). The

223    presence of lesion segregation in mutation data can be calculated with the new

224    "calculate_lesion_segregation" function. This calculation can be done for all mutations

225    together or separately for the different mutation contexts. The results can be visualized

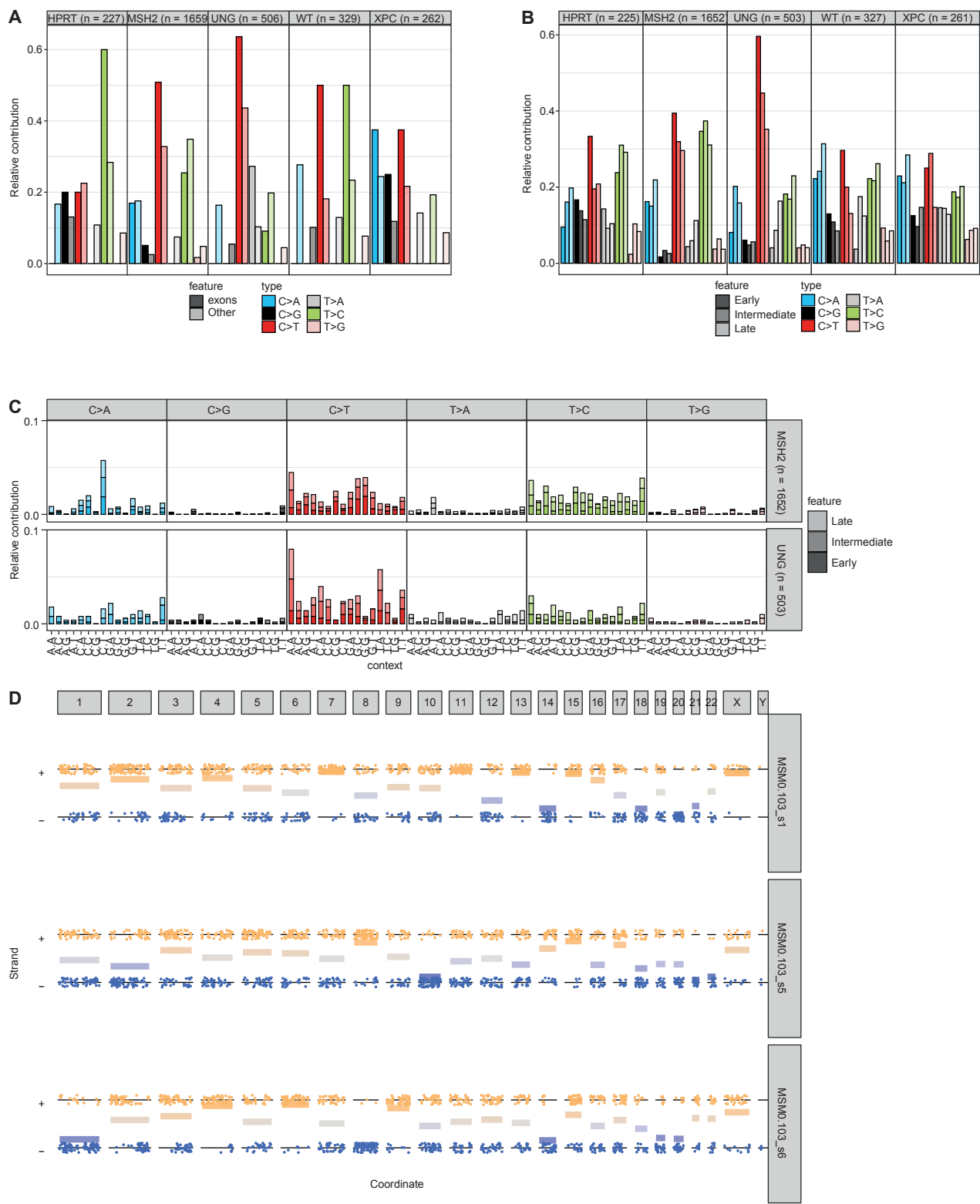226    using the "plot_lesion_segregation" function (Fig. 2d, Additional file 1: Figure S3).

227

228

Fig. 2 Regional spectra show differences between genomic regions

a Relative contribution of the indicated mutation types to the point mutation spectrum split

between exons and the rest of the genome for each sample. b Relative contribution of the

232    indicated mutation types to the point mutation spectrum split between early-,

233    intermediate-, and late-replicating DNA for each sample. **c** Relative contribution of each

234    trinucleotide change to the point mutation spectrum split between early- intermediate and

235    late-replicating DNA for each sample. **d** A jitter plot depicting the presence of lesion

236    segregation for each sample per chromosome. Each dot depicts a single base substitution.

237    Any C>N or T>N is shown as a "+" strand mutation, while G>N and A>N mutations are shown

238    on the "-" strand. The x-axis shows the position of the mutations. The horizontal lines are

239    calculated as the mean of the "+" and "-" strand, where "+" equals 1 and "-" equals 0. They

240    indicate per chromosome on which strand most of the mutations are located. The

241    mutations were downsampled to 33% to reduce the file size.

242

243    *Mutational signature analysis*

244    When performing signature analyses, it is possible to either extract novel signatures using

245    NMF or to fit previously defined signatures to a mutation count matrix (signature refitting).

246    Both approaches can be applied for all mutation types. By combining count matrices of

247    different types, it is even possible to create a composite signature.

248         MutationalPatterns now supports a variational Bayesian (Bayes) NMF algorithm from

249    the ccfindR package to help choose the optimal number of signatures, in addition to the

250    regular NMF algorithm (27) (Additional file 1: Figure S4). One challenge with *de novo*

251    signature extraction is that extracted signatures can be very similar to previously defined

252    signatures with known etiology. With the new "rename_nmf_signatures" function, these

253    extracted signatures can be identified using cosine similarity scores and their names can be

14

254    changed from an arbitrary naming to a custom naming that reflects their similarity to these

255    previously defined signatures.

256         The original MutationalPatterns package already contained the "fit_to_signatures"

257    function, which finds the optimal combination of signatures to reconstruct a profile and

258    calculates a reconstructed profile based on this combination of signatures.  However, this

259    approach could lead to too many signatures being used to explain the data (20). One simple

260    method to reduce this overfitting, which was used in the vignette of the previous version of

261    MutationalPatterns, is to remove all signatures with less than 10 mutations. However, this

262    method, which we will call "regular_10+", only reduced overfitting slightly. To reduce

263    overfitting, we introduce the new "fit_to_signatures_strict" function. The default backwards

264    selection method of this function iteratively refits a set of signatures to the data, each time

265    removing the signature with the lowest contribution. During each iteration the cosine

266    similarity between the original and reconstructed profile is calculated. The iteration process

267    stops when the change in cosine similarity between two iterations is bigger than the user-

268    specified "max_delta" cutoff (Additional file 1: Figure S5). Users can set the "max_delta"

269    cutoff based on their desired sensitivity and specificity. Stricter refitting, with this method, is

270    comparable to a previously described approach and results in less signatures being chosen

271    when tested on mutation data obtained from cell lines that lack specific DNA repair

272    pathways (Fig. 3a, b; see Additional file 2) (13). The "fit_to_signatures_strict" function also

273    has a best subset selection approach. This method works similarly to the backwards

274    selection approach. However, instead of removing the signature with the lowest

275    contribution, each combination of x signatures is tried. This includes signatures that were

276    not included in a previous iteration. Here, x is the number of signatures used during

277    refitting, which is reduced by one in each iteration step. By default,

15

278    "fit_to_signatures_strict" uses the backwards selection method, because the best subset

279    method becomes very slow when fitting against more than 10-15 signatures. Therefore, we

280    used the backwards selection method for all "strict" signature refitting analyses in the rest

281    of this manuscript. Another way to reduce overfitting is to only use signatures that are

282    known to be potentially active in your tissue/cells of interest. We recommend using this

283    method in combination with "fit_to_signatures_strict" for optimal results.

284    In addition to estimating contributions of signatures to mutation spectra, it is also

285    vital to know how confident these contributions are. The confidence of signature

286    contributions can be determined using a bootstrapping approach with the new

287    "fit_to_signatures_bootstrapped" function, which can use both the strict and the regular

288    refitting methods. Its output can be visualized in multiple ways using the

289    "plot_bootstrapped_contribution" function (Fig. 3c, Additional file 1: Figure S6). The

290    signature contributions can be correlated between signatures across the different bootstrap

291    iterations. This correlation can be visualized using the "plot_correlation_bootstrap" function

292    (Fig. 3d). A negative correlation between two signatures means that each signature had a

293    high contribution in iterations in which the other had a low contribution, which can occur

294    when the refitting process has difficulty distinguishing between two similar signatures. One

295    simple way to deal with highly similar signatures is to merge them. This can be done using

296    the new "merge_signatures" function.

297    To test the accuracy of signature analysis, the cosine similarity between the

298    reconstructed and original mutation profile needs to be determined. A high cosine similarity

299    between the reconstructed and original profile indicates that the used signatures can

300    explain the original spectrum well. This comparison between reconstructed and original

16

301    mutation profiles can be visualized with the new "plot_original_vs_reconstructed" function

302    (Fig. 3e).

303    In order to perform refitting, a matrix is required of the predefined signatures.

304    Signature matrices of the Catalogue of Somatic Mutations in Cancer (COSMIC) (v3.1 + v3.2),

305    SIGNAL (v1) and SparseSignatures (v1) are now included in MutationalPatterns (6,13,15,28).

306    These matrices include general, tissue-specific and drug exposure signatures. The COSMIC

307    matrices also include DBS and indel signatures, next to the standard SBS signatures.

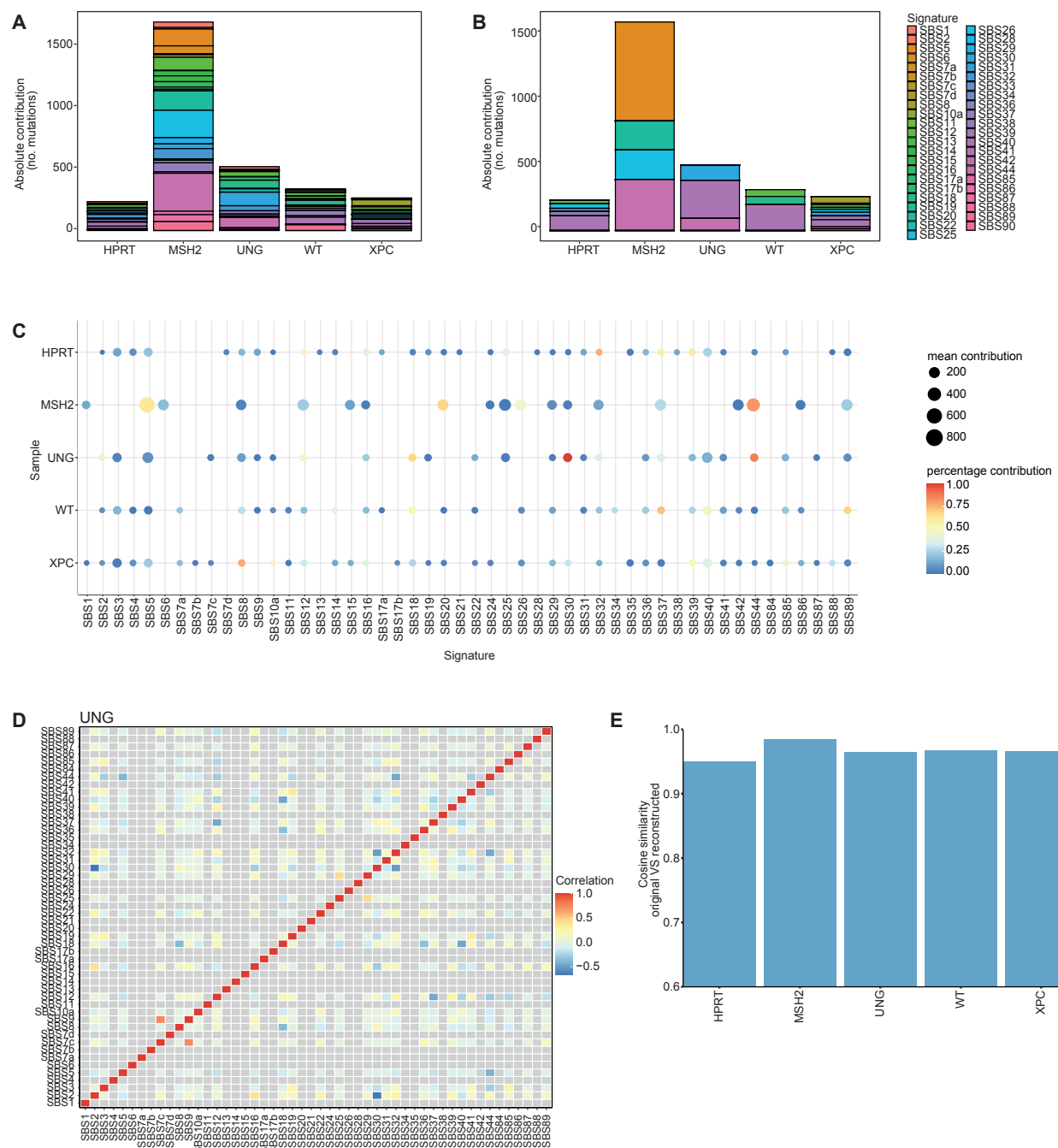308    Signature matrices can be easily loaded using the new "get_known_signatures" function.

309

310

311    Fig. 3 Signature refitting is improved

312    **a** Absolute contribution of each mutational signature for each sample using "regular"

313    signature refitting and **b** "strict" signature refitting. **c** Dot plot showing the contribution of

314    each mutational signature for each sample using bootstrapped signature refitting. The

315    colour of a dot indicates the fraction of bootstrap iterations in which a signature

316    contributed to a sample. The size indicates the mean number of contributing mutations

317   across bootstrap iterations in which the contribution was not zero. **d** Heatmap depicting the

318   Pearson correlation between signature contributions across the bootstrap iterations. **e** Bar

319   graph depicting the cosine similarity between the original and reconstructed profiles of each

320   sample based on signature refitting.

321

322   *Signature-specific damaging potential analysis*

323   Some signatures are more likely than others to have functional effects by causing premature

324   stop codons ("stop gain"), splice site mutations or missense mutations, because of sequence

325   specificity underlying these changes. With MutationalPatterns it is now possible to analyze

326   how likely it is for a signature to either cause "stop gain", "missense", "synonymous" or

327   "splice site" mutations for a set of genes of interest. For this analysis to be performed, the

328   potential damage first needs to be calculated per mutational context, with the

329   "context_potential_damage_analysis" function. Next, the potential damage per context is

330   combined using a weighted sum to calculate the potential damage per signature using the

331   "signature_potential_damage_analysis" function. The potential damage per signature is also

332   normalized using a "hypothetical" flat signature, which contains the same weight for each

333   mutation context.

334   This analysis will only take mutational contexts into account. Other features, such as

335   open/closed chromatin, are not considered, because they vary per tissue type. However,

336   this analysis can still give an indication of how damaging a signature might be, which could

337   be supplemented by further custom analyses.

338

19

339    This new version of MutationalPatterns also comes with many smaller updates and bugfixes.

340    A comprehensive list can be found in Additional file 3: Table S1.

341

342    **Results**

343

344    *Extended mutation context analysis and regional mutational patterns*

345    To demonstrate the importance of analyzing extended mutation contexts, regional

346    mutational patterns and lesion segregation for characterizing the underlying mutagenic

347    processes, we applied MutationalPatterns to three published mutation datasets. First, we

348    ran MutationalPatterns on 276 melanoma samples from the HMF database. After pooling

349    these samples, we observed that TT[C>T]CT mutations are the most common type of

350    substitution (Fig. 4a). This substitution type is more common than other T[C>T]C

351    substitutions, showing that the extended context has a large effect. Next, we compared the

352    mutation patterns of the melanoma samples between the different genomic regions

353    classified by the Ensembl regulatory build (30). While the patterns look similar, they are

354    significantly different (Fig. 4b) (p = 0.0005, chi-squared test). One subtle difference is the

355    low contribution of T[C>T]A in promoters compared to "Other" regions of the genome, not

356    present in the regulatory build.

357        Next, to show how MutationalPatterns can be used to identify regional activity of

358    specific mutation processes in an unsupervised manner, we applied the package on 217

359    pooled pediatric B-ALL WGS samples (31). These B-cell-derived leukemias have undergone

360    VDJ recombination, which is associated with somatic hypermutation at loci encoding for

361    immunoglobulin (32,33). As somatic hypermutation is associated with a specific signature,

362    these sites were expected to have a mutation spectrum that is different from the rest of the

363    genome. Indeed, MutationalPatterns was able to detect this for the two VDJ regions,

364    located on chromosomes 2 and 14 (Fig. 4c).  Some other regions also seem to have a

365    different mutational pattern, several of which contain PCDH genes. However, further

366    research is needed to explain these results. This example shows how MutationalPatterns

367    can identify region-specific mutational processes in an unsupervised manner.

368          Finally, to show how MutationalPatterns can identify lesion segregation, we applied

369    it on a dataset known to contain this phenomenon. We found significant lesion segregation

370    in data obtained from induced pluripotent stem cells treated with 0.109 uM of

371    dibenz[a,h]anthracene diol-epoxide (6,24), using the "plot_lesion_segregation" function of

372    MutationalPatterns (Fig. 2d). It was even possible to spot sister-chromatid-exchange events,

373    such as on chromosome 2 of sample MSM0.103_s6 (Fig. 2d, lower panel). To reduce the file

374    size of the figure, 66% of the mutations of each sample were removed using the

375    "downsample" argument of this function.  Using MutationalPatterns, we also found lesion

376    segregation in patients that received the antiviral drug ganciclovir (7).
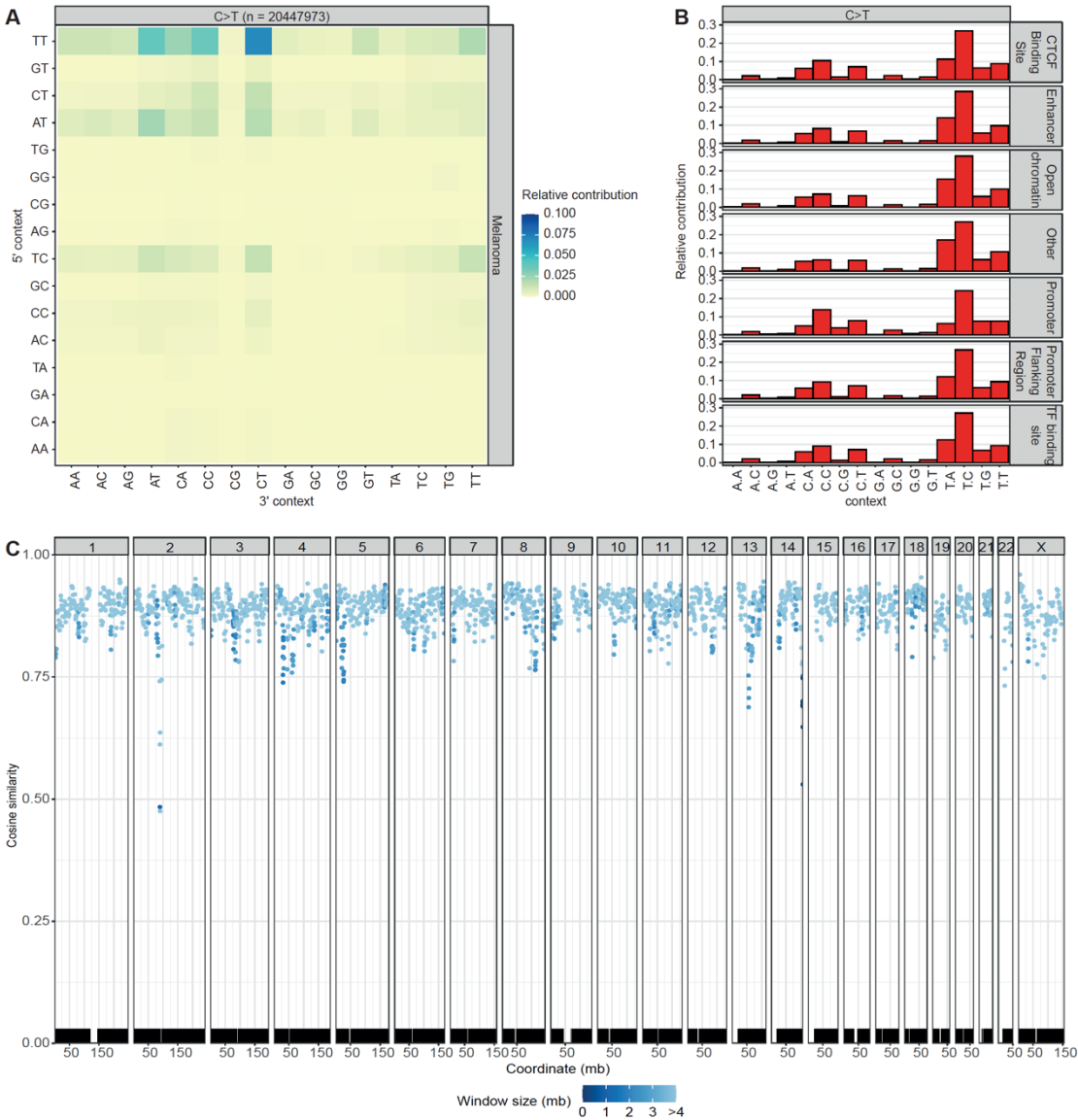
377

Fig. 4 Large cancer datasets show extended and regional mutation patterns

**a** Heatmap depicting the relative contribution of the indicated mutation types and the surrounding bases to the point mutation spectrum for metastatic melanomas. The total number of somatic point mutations is indicated. **b** Relative contribution of each C>T trinucleotide change to the point mutation spectrum split between different genomic regions. **c** Graph depicting the similarity in the mutation profile between genomic windows and the rest of the genome. Each dot shows the cosine similarity between the mutation profiles of a single window and the rest of the genome. The dots are colored based on the

386    sizes in mega bases of the windows. The locations of the mutations are plotted on the

387    bottom of the figure.

388

389    *MutationalPatterns offers more functionality than other mutation analysis tools*

390    An overview of the functions of MutationalPatterns and related tools is shown in Table 1.

391    The original version of MutationalPatterns is also included in this table. An important

392    advantage of the original package was that it combined many mutational analyses into a

393    single package. This new version improves many of these features and adds many new and

394    unique features.

395

396    *Mutation matrices can be generated faster*

397    To make MutationalPatterns scalable to large cancer datasets and suitable for interactive

398    analysis we improved the runtime of the "mut_matrix" and "mut_matrix_stranded"

399    functions by vectorizing them. The new functions for retrieving the mutation contexts and

400    generating the mutation matrices have also been written in a vectorized way. As a result,

401    these functions have O(n) or better scaling as tested on a large WGS database from the

402    Hartwig Medical Foundation (HMF) (Additional file 1: Figure S7) (29).

403        To test their improved performance, we benchmarked the "mut_matrix" and

404    "mut_matrix_stranded" functions on the example data provided in the previous version of

405    MutationalPatterns (Additional file 1: Figure S8). These functions are now respectively 3.4

406    and 2.6 times as fast on average. In other words, a mutation matrix for 1 million SBSs can

407    now be made in only 135 seconds on a laptop, which makes these functions suitable for

408    large cancer datasets.

409

**Table 1:** Feature comparison with other packages

| Group | Feature | Mutational Patterns | Mutational Patterns original (12) | Sigprofiler (13) | SignatureAnalyzer (13) | deconstruct Sigs (14) | sparseSignatures (15) | signeR (16) | somaticSignatures (17) | Maftools (18) | decompTumor2Sig (19) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Language | Language/platform | R (bioconductor) | R (bioconductor) | Python (+ R wrapper) | Python | R (cran) | R (bioconductor) | R (bioconductor) | R (bioconductor) | R (bioconductor) | R (bioconductor) |
| Genome | Supported genomes | Genome agnostic | Genome agnostic | Human, Mice, Rat, Yeast | - | Human | Genome agnostic | Genome agnostic | Genome agnostic | Genome agnostic | Genome agnostic |
| Mutation profile | 96 SNV profile | X | X | X | - | X | - | X | X | X | X |
|  | extended SNV profile | X | - | X | - | - | - | - | X | - | X |
|  | Indel profile | X | - | X | - | - | - | - | - | - | - |
|  | DBS profile | X | - | X | - | - | - | - | - | - | - |
|  | MBS profile | X | - | - | - | - | - | - | - | - | - |
|  | Transcriptional strand bias profile | X | X | X | - | - | - | - | - | - | - |
|  | Replicative strand bias profile | X | X | X | - | - | - | - | - | - | - |
|  | Pool samples | X | - | - | - | - | - | - | - | - | - |
| Signature extraction | Signature extraction (NMF) | X | X | X | - | - | - | - | X | X | - |
|  | Signature extraction (Bayes NMF) | X | - | - | X | - | - | X | - | - | - |
|  | Signature extraction (Lasso NMF) | - | - | - | - | - | X | - | - | - | - |
|  | Update signature names | X | - | - | - | - | - | - | - | - | - |
| Signature refitting | Signature refitting | X | X | X | X | X | - | - | - | - | X |
|  | Strict signature refitting | X | - | X | X | X | - | - | - | - | X |
|  | Strict signature refitting (best subset) | X | - | - | - | - | - | - | - | - | X |
| - | Bootstrapped signature refitting | X | - | - | - | - | - | - | - | - | - |
|  | Correlation bootstrapped refitting | X | - | - | - | - | - | - | - | - | - |
| Signature damage analysis | Signature potential damage analysis | X | - | - | - | - | - | - | - | - | - |
| Signature other | Plot supported profiles / signatures | X | X | X | X | X | X | X | X | X | X |
|  | Plot and compare supported profiles | X | X | - | - | - | - | - | - | - | - |
|  | Signature contribution heatmap | X | X | - | - | - | - | X | X | - | - |
|  | Signature contribution barplot | X | X | - | - | - | - | X | X | - | - |
|  | Signature/profile similarity heatmap | X | X | - | - | - | - | - | - | X | - |
|  | Similarity with reconstructed profile barplot | X | - | - | - | - | - | - | - | - | - |
| Genomic distribution | Rainfall plot | X | X | - | - | - | - | - | X | X | - |
|  | Enrichment/depletion in genomic region | X | X | - | - | - | - | - | - | - | - |
|  | Region specific profiles | X | - | - | - | - | - | - | - | - | - |
|  | Region specific signatures | X | - | - | - | - | - | - | - | - | - |
|  | Unsupervised regional similarity | X | - | - | - | - | - | - | - | - | - |
| Lesion segregation | Lesion segregation | X | - | - | - | - | - | - | - | - | - |

410

411

412     *Strict signature refitting improves performance*

413     To determine how well the strict refitting method of MutationalPatterns performs as

414     compared to the regular method, we used simulated mutation matrices. These matrices

415     were generated by sampling trinucleotide changes of 4 different randomly selected

416     signatures. This process was repeated 300 times per matrix, to generate 300 "samples".

417     Each of the samples in a matrix contained the same number of mutations per signature but

418     was composed of different signatures. The signatures were selected from the first 30

419     signatures of the COSMIC signature matrix. We limited our analysis to the first 30, because

420     these are the signatures that are most often observed in cancers and therefore more

421     accurately resemble real-life scenarios. In addition, this approach better resembles how the

422     package is used, because users will often fit against a limited number of signatures

423     associated with a specific tissue. By limiting ourselves to the first 30 COSMIC signatures we

424     also reduced overfitting. Any overfitting we observed was thus not caused by us using an

425     unusually large signature matrix. In total we generated 4 matrices, each containing 300

426     samples. The number of mutations per sample was respectively 200, 400, 2000 and 4000 for

427     the 4 different matrices.

428         The fraction of correctly attributed mutations to the specific signatures was

429     increased with the strict refitting approach of MutationalPatterns as compared to "regular"

430     or "regular_10+" refitting (Additional file 1: Figure S9a). All the tested refitting methods

431     work better when there are more mutations per signature. Instead of using the number of

432     correctly attributed mutations as a readout for performance, we determined whether the

433     presence and absence of specific signatures was correctly classified. This readout might be

25

434    more informative for mutational signature analysis because the presence of a signature can

435    be a clinically relevant finding. The strict refitting method achieved a much higher precision

436    than the original methods, while retaining a high correct recall rate (sensitivity) (Additional

437    file 1: Figure S9b). The strict method obtained an area under the curve (AUC) of 0.925, even

438    when only 50 mutations were present per signature, indicating that refitting can be

439    performed on relatively small amounts of mutations.

440

441    *SBS10a and SBS18 have a high damage potential*

442    We applied the "signature_potential_damage_analysis" function on the COSMIC signatures.

443    This analysis showed that SBS10a and SBS18 are respectively 3.6 and 2.0 times as likely to

444    cause a "stop gain" mutation compared to a completely flat signature, containing the same

445    weight for each mutation context, on a set of genes associated with cancer (Additional file

446    3: Table S2, Table S3). SBS18 is related to oxidative stress, suggesting that this type of stress

447    has a high potency of generating premature stop codons in genes that are recurrently

448    associated with tumorigenesis (13). In contrast, the clock-like signature SBS1, which also

449    occurs in healthy cells, was 0.81 and 0.40 times as likely to cause "stop gain" and "splice

450    site" mutations, respectively, as compared to a completely flat hypothetical signature (2,34)

451    (Additional file 3: Table S2). The damaging potential of this ageing-related mutational

452    process is thus relatively low. Overall, C>A heavy signatures, like the recently identified

453    ganciclovir signature, have more damage potential, because they are most likely to

454    introduce a premature stop codon in an open reading frame (7). Being able to quickly assess

455    the damage potential of existing and novel signatures can be very useful to prioritize

456    samples and mutagenic exposures for further investigation.

457

458    *Applying MutationalPatterns on mutation data of DNA repair-deficiencies*

459    To illustrate the functionality of MutationalPatterns on real-life data and to obtain novel

460    biological insights, we applied it to mutation data obtained from cell lines in which we

461    deleted specific DNA repair pathways using CRISPR-Cas9 genome editing technology

462    (Additional file 1: Figure S10, Additional file 2). In AHH-1 cells, a lymphoblastoid cell line, we

463    generated bi-allelic knockout lines of *MSH2, UNG* and *XPC* by transfecting the cells with a

464    plasmid containing Cas9 and a single gRNA against the gene of interest. By co-transfection

465    with a *HPRT*-targeting plasmid, we were able to select the transfected cells using 6-

466    thioguanine, to which only HPRT-sufficient cells are sensitive. Using this protocol, no

467    targeting vectors for each gene of interest were required. We analyzed somatic mutations in

468    *HPRT*-only knockout lines as well as the combination of *HPRT* with *MSH2, UNG* and *XPC*

469    (Additional file 2). To catalogue mutations that were acquired specifically in the absence of

470    the targeted DNA repair gene, we used a previously developed method (35). In brief, whole

471    genome sequencing was performed on generated clones and subclones. By subtracting

472    variants present in the clones from those in the subclones, the somatic mutations, that

473    accumulated in between the clonal steps, were determined.

474

475    The SBS profiles are shown in Additional file 1: Figure S11. Interestingly, the profile observed

476    in the *MSH2* knockout cell line displayed a large C[C>A]T peak. When extending the

477    sequence context surrounding the mutated base, the *MSH2* deficiency profile showed a

478    large TT[T>C]TT peak, suggesting that this extended context surrounding mutated thymine

479    residues is important for the underlying mutagenic process (Fig. 1d).

480

481 Next, we examined regional mutation patterns. The spectra of the *MSH2-* and *UNG-*

482 deficient cells varied between the exonic regions and the rest of the genome (Fig. 2a)(fdr =

483 0.0012, fdr = 0.0012; chi-squared test). Their exons contained more C>T and less T>C

484 mutations. The other samples did not show a significant difference in regional mutation

485 spectra. However, when we downsampled all the samples to 227 mutations, which is the

486 number of mutations in the *HPRT* only knockout, no significant regional mutation patterns

487 were observed in *MSH2* and *UNG* knockout cells. This suggests that with this number of

488 mutations insufficient statistical power was obtained for these analyses. Next to examining

489 mutation profiles in exonic regions, we also analyzed regions with different replication

490 timing dynamics, using the median replication timing data from 5 B-lymphocyte cell lines

491 from ENCODE (Fig. 2b, Additional file 3: Table S4) (40). The spectra of *MSH2* and *UNG*

492 knockouts were different between early-, intermediate- and late-replicating DNA (fdr =

493 0.0012, fdr = 0.0012; chi-squared test). Early replicating DNA has more C>T and less C>A

494 than late replicating DNA. These differences were still present when downsampling was

495 applied (fdr = 0.0025, fdr = 0.010; chi-squared test). Based on these region-specific analyses,

496 we can conclude that the mutational processes active in the *MSH2* and *UNG* knockouts

497 show varying activities in different regions of the genome, a result that cannot easily be

498 obtained with other tools.

499　　　We also tested if any of the DNA repair knockout cells displayed lesion segregation,

500 which would indicate that most of the mutations occurred during a single cell-cycle;

501 however, this was not the case (Additional file 1: Figure S6).

502

503    Finally, we looked at the mutational signatures in the knockout samples. Based on signature

504    refitting, the *MSH2* knockout contained contributions of SBS5, SBS20, SBS26 and SBS44 (Fig.

505    3b, c). Because of the bootstrapping we can be more confident in these results. SBS5 is a

506    clock-like signature, with unknown etiology. SBS20, SBS26 and SBS44 are all associated with

507    defective DNA mismatch repair in cancer mutation data (13). The UNG knockout contained

508    contributions from SBS30, which has previously been attributed to deficiency of the base

509    excision repair gene *NTHL1* (13). The glycosylase encoded by *NTHL1* is involved in the

510    removal of oxidized pyrimidines from the DNA and therefore SBS30 likely reflects an

511    alternative consequence of oxidative stress-induced mutagenesis as compared to SBS18.

512    However, *UNG* is a glycosylase that is believed to remove uracil residues from the DNA

513    (36,37). Therefore, our data suggests that SBS30 can be caused, besides oxidized

514    pyrimidines, by unremoved uracil residues. Alternatively, *UNG* may also, to a certain extent,

515    be involved in the removal of oxidized pyrimidines from the DNA. Even though the

516    contribution of SBS30 was relatively modest in the *UNG* knockout, it was consistently picked

517    up by the bootstrapping algorithm. This observation indicated that the number of mutations

518    attributed to a signature is not necessarily related to the confidence of its presence, which

519    further demonstrates the importance of our bootstrapping approach. Unexpectedly, the

520    contribution of SBS30 in *UNG* knockout cells was negatively correlated with SBS2, even

521    though their cosine similarity is only 0.46 (Fig. 3d). This indicates that the refitting algorithm

522    has difficulty choosing between SBS2 and SBS30. Such difficulties in signature selection

523    could lead to different and possibly incorrect signatures being attributed to similar sample

524    types. Understanding the correlation of estimated signature contributions between

525    different signatures, which can be achieved with bootstrapping, is important to prevent

526    incorrect interpretation of the data. The *XPC* knockout contained contributions from SBS8.

29

527    The etiology of this signature is not yet known. However, this finding further confirms the

528    association of SBS8 with nucleotide excision repair deficiency (38,39). Overall, the COSMIC

529    signatures could explain the mutation profiles of most samples quite well, even when strict

530    refitting was used (Fig. 3e).

531        Next, we studied the indel signatures in these knockout lines. Deletion of *MSH2*

532    resulted in an increased number of indels as compared to wild-type cells (Fig. 1b). Most of

533    these indels were single thymine deletions in thymine mononucleotide repeat regions.

534    Signature analysis indicated that ID1, ID2 and ID7 contributed to the indel pattern in the

535    MSH2-deficient cells (Fig. 5a, b). Of these, ID1 and ID2 are associated with polymerase

536    slippage during DNA replication and found in large numbers in cancers with mismatch repair

537    deficiency. ID7 is also associated with defective DNA mismatch repair, but not attributed to

538    polymerase slippage (13). Together these signatures could explain the mutational indel

539    profile of *MSH2* knockout cells very well (Fig. 5c), showing that MutationalPatterns can

540    perform indel signature refitting. None of the knockout cells displayed a strongly increased

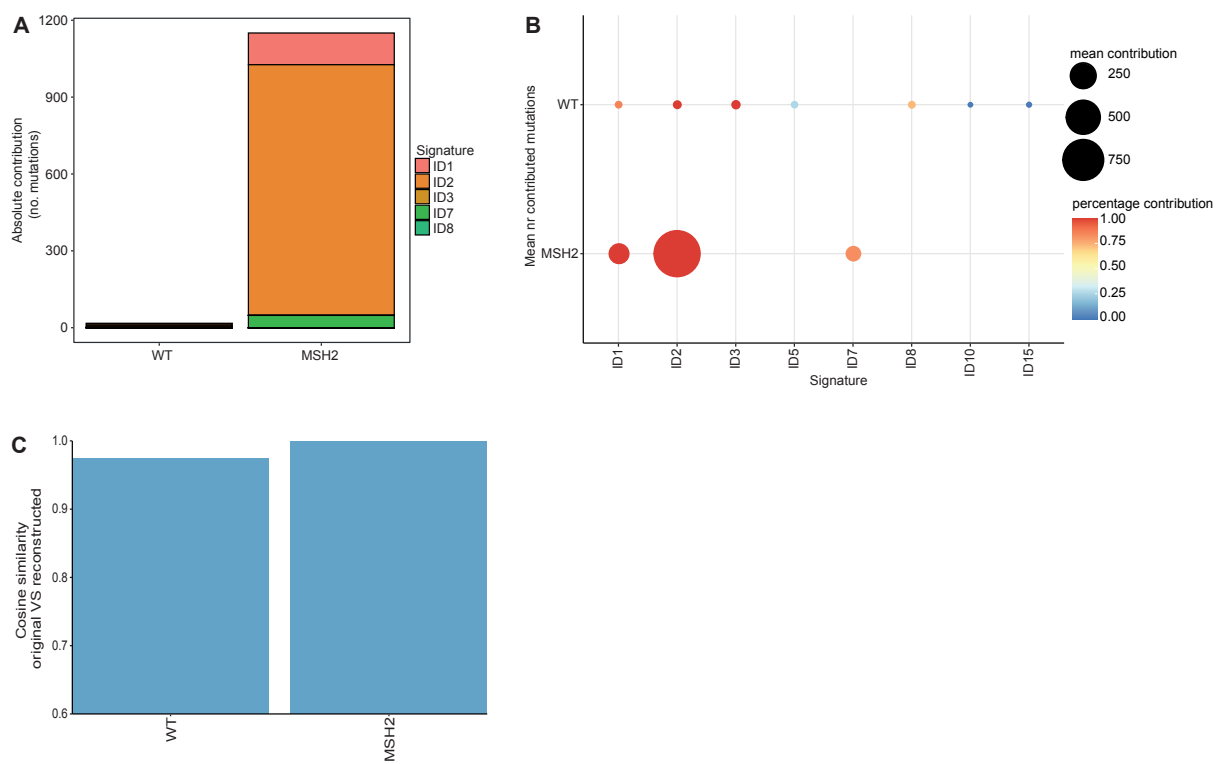541    number of DBSs as compared to the wild-type cells (Fig. 1c).

542

543    Fig. 5 Indel signatures can explain the MSH2 profile

544    **a** Relative contribution of each mutational signature for the wild-type (WT) and *MSH2*

545    samples using strict signature refitting. **b** Dot plot showing the contribution of each

546    mutational signature for the WT and *MSH2* samples using bootstrapped signature refitting.

547    The color of a dot indicates the fraction of bootstrap iterations in which a signature

548    contributed to a sample. The size indicates the mean number of contributing mutations

549    across bootstrap iterations in which the contribution was not zero. **c** Bar graph depicting the

550    cosine similarity between the original and reconstructed profiles of the WT and *MSH2*

551    samples based on signature refitting.

552

553    **Discussion**

554   The novel version of MutationalPatterns has been designed to be easy-to-use in such a way

555   that both experienced bioinformaticians and wet-lab scientists with a limited computational

556   background can use it. The code is written in the tidyverse style, which makes it more

557   similar to natural English and therefore easier to understand for non-programmers.

558   MutationalPatterns gives clear error messages with tips on how to solve them, in contrast to

559   the default error messages in R, which can sometimes be cryptic. The updated vignette,

560   accompanying the package, not only explains how the functions in the package can be used,

561   but also informs users on the pros and cons of the different analysis strategies.

562       Similar to the previous version of the package, plots are all generated using ggplot2

563   (41). This allows users to visualize their data in highly customizable plots that can be easily

564   modified. Because this feature was not readily apparent for many users of the original

565   MutationalPatterns package, we have now explicitly showed how to modify the elements of

566   a plot, such as the axis and theme, in the vignette.

567       We have adopted unit testing for this version of the package, resulting in more than

568   90% code coverage. This will improve the stability of the package and makes it easier to

569   maintain.

570       The novel version of MutationalPatterns is already available on Bioconductor as an

571   update of the previous version. MutationalPatterns does not break existing scripts and

572   pipelines, because backwards incompatible changes have been kept to a minimum.

573

574   **Conclusions**

575     MutationalPatterns is an easy-to-use R/Bioconductor package that allows in-depth analysis

576     of a broad range of patterns in somatic mutation catalogues, supporting single and double

577     base substitutions as well as small insertions and deletions. Here, we have described the

578     new and improved features of the package and shown how the package performs on

579     existing cancer data sets and on mutation data obtained from cell lines in which specific

580     DNA repair genes are deleted. These analyses demonstrate how the package can be used to

581     generate novel biological insights.

582

583     Mutational pattern analyses have proven to be a powerful approach to dissect mutational

584     processes that have operated in cancer and to support treatment decision making in

585     personalized medicine. Therefore, mutational patterns hold a great promise for improved

586     future cancer diagnosis. The MutationalPatterns package can be used to fulfill this promise

587     and we are confident that it will be embraced by the community.

588

589     **Availability and requirements**

590     The availability and requirements are listed as follows:

591     Project name: MutationalPatterns

592     Project home page: https://github.com/ToolsVanBox/MutationalPatterns

593     Archived version:

594     https://bioconductor.org/packages/3.14/bioc/html/MutationalPatterns.html

595     Operating system(s): Linux, Windows or MacOS

596     Programming language: R (version > = 4.1.0)

597     License: MIT

598

599     **List of abbreviations**

600     HR: homologous recombination

601     Indels: Insertions and deletions

602     DBS: double base substitutions

603     VCF: variant call format

604     MBS: Multi base substitutions

605     COSMIC: Catalogue of Somatic Mutations in Cancer

606     NMF: non-negative matrix factorization

607     Bayes: Bayesian

608     AUC: Area under the curve

609     PCA: Principal component analysis

610     CI: Confidence interval

611     WT: wild-type

612     Mb: mega bases

613

614     **Declarations**

615    *Ethics approval and consent to participate*

616    Not applicable

617

618    *Consent for publication*

619    Not applicable

620

621    *Availability of data and materials*

622    The datasets supporting this article are available on EGA under accession number (Study ID

623    EGAS00001004789).

624    Additionally, the VCF files and scripts that can be used to reproduce all figures in this paper

625    can be found at

626    https://github.com/ToolsVanBox/MutationalPatterns_manuscript2_data_scripts/

627

628    *Competing interests*

629    The authors declare that they have no competing interests.

630

631    *Funding*

632    This work was financially supported by a NWO VIDI grant project 016.Vidi.171.023 to R.v.B.

633

634    Authors' contributions

635    F.M., R.v.B. and A.M.B wrote the manuscript. F.M. and J.d.K. developed and implemented

636    the package. F.M. and R.O. maintain the package. A.M.B. and M.V. generated the data. F.M.

637    and M.J.v.R. analyzed the data. A.v.H., B.v.d.R. and E.C. tested the package and provided

638    feedback. All authors read and approved the final manuscript.

639

640    *Acknowledgements*

641    We would like to thank Francis Blokzijl and Roel Janssen for developing and maintaining the

642    first version of this package. We also want to thank Roel Janssen for his support during the

643    handover of the package. Finally, we would like to thank anybody who tested the package

644    for their feedback.

645

646    **References**

647    1.    Helleday T, Eshtad S, Nik-Zainal S. Mechanisms underlying mutational signatures in

648          human cancers. Nat Rev Genet. 2014;15:585–98.

649    2.    Blokzijl F, de Ligt J, Jager M, Sasselli V, Roerink S, Sasaki N, et al. Tissue-specific

650          mutation accumulation in human adult stem cells during life. Nature. 2016;538:260–

651          4.

652    3.    Campbell PJ, Getz G, Korbel JO, Stuart JM, Jennings JL, Stein LD, et al. Pan-cancer

653          analysis of whole genomes. Nature [Internet]. 2020;578(7793):82–93. Available from:

654          https://doi.org/10.1038/s41586-020-1969-6

655    4.    Alexandrov LB, Nik-Zainal S, Wedge DC, Campbell PJ, Stratton MR. Deciphering

656          Signatures of Mutational Processes Operative in Human Cancer. Cell Rep.

657          2013;3:246–59.

658    5.    Pleguezuelos-Manzano C, Puschhof J, Rosendahl Huber A, van Hoeck A, Wood HM,

659          Nomburg J, et al. Mutational signature in colorectal cancer caused by genotoxic pks+

660          E. coli. Nature. 2020;580:269–73.

661    6.    Kucab JE, Zou X, Morganella S, Joel M, Nanda AS, Nagy E, et al. A Compendium of

662          Mutational Signatures of Environmental Agents. Cell. 2019;177:821-836.e16.

663    7.    de Kanter JK, Peci F, Bertrums E, Rosendahl Huber A, van Leeuwen A, van Roosmalen

664          MJ, et al. Antiviral treatment causes a unique mutational signature in cancers of

665          transplantation recipients. Cell Stem Cell [Internet]. 2021; Available from:

666          https://www.sciencedirect.com/science/article/pii/S1934590921003374

667    8.    Davies H, Glodzik D, Morganella S, Yates LR, Staaf J, Zou X, et al. HRDetect is a

668          predictor of BRCA1 and BRCA2 deficiency based on mutational signatures. Nat Med.

669          2017;23:517–25.

670    9.    Nguyen L, W. M. Martens J, Van Hoeck A, Cuppen E. Pan-cancer landscape of

671          homologous recombination deficiency. Nat Commun. 2020;11:5584.

672    10.   Bryant HE, Schultz N, Thomas HD, Parker KM, Flower D, Lopez E, et al. Specific killing

673          of BRCA2-deficient tumours with inhibitors of poly(ADP-ribose) polymerase. Nature.

674          2005;434:913–7.

675    11.   Fong PC, Boss DS, Yap TA, Tutt A, Wu P, Mergui-Roelvink M, et al. Inhibition of

676          Poly(ADP-Ribose) Polymerase in Tumors from BRCA Mutation Carriers. N Engl J Med.

677    2009;361:557–68.

678    12.    Blokzijl F, Janssen R, van Boxtel R, Cuppen E. MutationalPatterns: comprehensive

679          genome-wide analysis of mutational processes. Genome Med. 2018;

680    13.    Alexandrov LB, Kim J, Haradhvala NJ, Huang MN, Tian Ng AW, Wu Y, et al. The

681          repertoire of mutational signatures in human cancer. Nature. 2020;578:94–101.

682    14.    Rosenthal R, McGranahan N, Herrero J, Taylor BS, Swanton C. deconstructSigs:

683          delineating mutational processes in single tumors distinguishes DNA repair

684          deficiencies and patterns of carcinoma evolution. Genome Biol. 2016;17:31.

685    15.    Ramazzotti D, Lal A, Liu K, Tibshirani R, Sidow A. De Novo Mutational Signature

686          Discovery in Tumor Genomes using SparseSignatures. bioRxiv. 2019;384834.

687    16.    Rosales RA, Drummond RD, Valieris R, Dias-Neto E, Da Silva IT. signeR: An empirical

688          Bayesian approach to mutational signature discovery. Bioinformatics. 2017;33:8–16.

689    17.    Gehring JS, Fischer B, Lawrence M, Huber W. SomaticSignatures: Inferring mutational

690          signatures from single-nucleotide variants. Bioinformatics. 2015;31:3673–5.

691    18.    Mayakonda A, Lin D-C, Assenov Y, Plass C, Koeffler HP. Maftools: efficient and

692          comprehensive analysis of somatic variants in cancer. Genome Res. 2018/10/19. 2018

693          Nov;28:1747–56.

694    19.    Krüger S, Piro RM. decompTumor2Sig: identification of mutational signatures active

695          in individual tumors. BMC Bioinformatics [Internet]. 2019;20(4):152. Available from:

696          https://doi.org/10.1186/s12859-019-2688-6

697    20.    Maura F, Degasperi A, Nadeu F, Leongamornlert D, Davies H, Moore L, et al. A

698     practical guide for mutational signature analysis in hematological malignancies. Nat

699     Commun. 2019;10:2969.

700  21.  Polak P, Karlic R, Koren A, Thurman R, Sandstrom R, Lawrence MS, et al. Cell-of-origin

701     chromatin organization shapes the mutational landscape of cancer. Nature.

702     2015;518:360–4.

703  22.  Buisson R, Langenbucher A, Bowen D, Kwan EE, Benes CH, Zou L, et al. Passenger

704     hotspot mutations in cancer driven by APOBEC3A and mesoscale genomic features.

705     Science (80- ). 2019;364:eaaw2872.

706  23.  Gonzalez-Perez A, Sabarinathan R, Lopez-Bigas N. Local Determinants of the

707     Mutational Landscape of the Human Genome. Cell. 2019;177:101–14.

708  24.  Aitken SJ, Anderson CJ, Connor F, Pich O, Sundaram V, Feig C, et al. Pervasive lesion

709     segregation shapes cancer genome evolution. Nature. 2020;583:265–70.

710  25.  Benjamin D, Sato T, Cibulskis K, Getz G, Stewart C, Lichtenstein L. Calling Somatic

711     SNVs and Indels with Mutect2. bioRxiv. 2019;861054.

712  26.  Yates AD, Achuthan P, Akanni W, Allen J, Allen J, Alvarez-Jarreta J, et al. Ensembl

713     2020. Nucleic Acids Res. 2020;48:D682–8.

714  27.  Woo J, Winterhoff BJ, Starr TK, Aliferis C, Wang J. De novo prediction of cell-type

715     complexity in single-cell RNA-seq and tumor microenvironments. Life Sci Alliance.

716     2019;2:e201900443.

717  28.  Degasperi A, Amarante TD, Czarnecki J, Shooter S, Zou X, Glodzik D, et al. A practical

718     framework and online tool for mutational signature analyses show inter-tissue

719     variation and driver dependencies. Nat cancer. 2020;1:249–63.

720    29.    Priestley P, Baber J, Lolkema MP, Steeghs N, de Bruijn E, Shale C, et al. Pan-cancer

721            whole-genome analyses of metastatic solid tumours. Nature. 2019;575:210–6.

722    30.    Zerbino DR, Wilder SP, Johnson N, Juettemann T, Flicek PR. The ensembl regulatory

723            build. Genome Biol [Internet]. 2015 Mar;16:56. Available from:

724            http://europepmc.org/articles/PMC4407537

725    31.    Ma X, Liu Y, Liu Y, Alexandrov LB, Edmonson MN, Gawad C, et al. Pan-cancer genome

726            and transcriptome analyses of 1,699 paediatric leukaemias and solid tumours.

727            Nature. 2018 Feb;

728    32.    Chi X, Li Y, Qiu X. V(D)J recombination, somatic hypermutation and class switch

729            recombination of immunoglobulins: mechanism and regulation. Immunology

730            [Internet]. 2020/02/27. 2020 Jul;160(3):233–47. Available from:

731            https://pubmed.ncbi.nlm.nih.gov/32031242

732    33.    Di Noia JM, Neuberger MS. Molecular Mechanisms of Antibody Somatic

733            Hypermutation. Annu Rev Biochem [Internet]. 2007 Jun 7;76(1):1–22. Available from:

734            https://doi.org/10.1146/annurev.biochem.76.061705.090740

735    34.    Alexandrov LB, Jones PH, Wedge DC, Sale JE, Campbell PJ, Nik-Zainal S, et al. Clock-

736            like mutational processes in human somatic cells. Nat Genet. 2015;47:1402–7.

737    35.    Drost J, van Boxtel R, Blokzijl F, Mizutani T, Sasaki N, Sasselli V, et al. Use of CRISPR-

738            modified human stem cell organoids to study the origin of mutational signatures in

739            cancer. Science (80- ). 2017;238:eaao3130.

740    36.    Prasad A, Wallace SS, Pederson DS. Initiation of Base Excision Repair of Oxidative

741            Lesions in Nucleosomes by the Human, Bifunctional DNA Glycosylase NTH1. Mol Cell

742         Biol. 2007;27:8442 LP – 8453.

743    37.   Li J, Braganza A, Sobol RW. Base Excision Repair Facilitates a Functional Relationship

744         Between Guanine Oxidation and Histone Demethylation. Antioxid Redox Signal.

745         2013;18:2429–43.

746    38.   Jager M, Blokzijl F, Kuijk E, Bertl J, Vougioukalaki M, Janssen R, et al. Deficiency of

747         nucleotide excision repair is associated with mutational signature observed in cancer.

748         Genome Res. 2019;29:1067–77.

749    39.   Yurchenko AA, Padioleau I, Matkarimov BT, Soulier J, Sarasin A, Nikolaev S. XPC

750         deficiency increases risk of hematologic malignancies through mutator phenotype

751         and characteristic mutational signature. Nat Commun [Internet]. 2020;11(1):5834.

752         Available from: https://doi.org/10.1038/s41467-020-19633-9

753    40.   Moore JE, Purcaro MJ, Pratt HE, Epstein CB, Shoresh N, Adrian J, et al. Expanded

754         encyclopaedias of DNA elements in the human and mouse genomes. Nature.

755         2020;583:699–710.

756    41.   Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York;

757         2016.

758

759    Additional file 1:

760    PDF (.pdf)

761    Additional figures

762    A PDF file containing the additional figures.

763

764     Additional file 2:

765     PDF (.pdf)

766     Additional methods

767     A PDF file describing the generation and sequencing analysis of the knockout lines.

768

769     Additional file 3:

770     Excel (.xlsx)

771     Additional tables

772     An Excel file containing the additional tables.