

1 **Katdetectr: utilising unsupervised changepoint analysis for robust kataegis detection**

2

3 Daan M. Hazelaar^{1,2,†}, Job van Riet^{1-3,†}, Youri Hoogstrate^{1,4}, Martijn P. Lolkema² and Harmen
4 J. G. van de Werken^{1,3,5}

5

6

7 ¹Cancer Computational Biology Centre, Erasmus MC Cancer Institute, Erasmus University
8 Medical Centre, Dr. Molewaterplein 40, 3015 GD, Rotterdam, the Netherlands, ²Department
9 of Medical Oncology, Erasmus MC Cancer Institute, Erasmus University Medical Centre, Dr.
10 Molewaterplein 40, 3015 GD, Rotterdam, the Netherlands, ³Department of Urology,
11 Erasmus MC Cancer Institute, Erasmus University Medical Centre, Dr. Molewaterplein 40,
12 3015 GD, Rotterdam, the Netherlands, ⁴Department of Neurology, Erasmus MC Cancer
13 Institute, Erasmus University Medical Center, Dr. Molewaterplein 40, 3015 GD, Rotterdam,
14 the Netherlands, ⁵Department of Immunology, Erasmus MC Cancer Institute, Erasmus
15 University Medical Center, Dr. Molewaterplein 40, 3015 GD, Rotterdam, the Netherlands.

16

17 [†]Shared first-author

18

19 All Authors have seen and approved this manuscript.

20 **Abstract**

21 **Motivation:**

22 Kataegis refers to the occurrence of regional hypermutation in cancer genomes and is a
23 phenomenon that has been observed in a wide range of malignancies. Robust detection of
24 kataegis is necessary to advance research regarding the origins and clinical impact of
25 kataegis. Multiple kataegis detection packages are publicly available; however, the
26 performance of their respective approaches have not been evaluated extensively. Here, we
27 introduce katdetectr, an R-based, open-source, computationally fast, and robust package
28 for the detection, characterisation and visualisation of kataegis.

29 **Results:**

30 The performance of katdetectr and five publicly available packages for kataegis detection
31 were evaluated using an in-house generated synthetic dataset and an *a priori* labelled pan-
32 cancer dataset of whole genome sequenced malignancies. The performance evaluation
33 revealed that katdetectr has the highest accuracy and normalized Matthews Correlation
34 Coefficient for kataegis classification on both the synthetic and the *a priori* labelled dataset.
35 Katdetectr is in particularly more robust for kataegis detection within samples with a high
36 tumour mutational burden.

37 **Availability and Implementation:**

38 Katdetectr imports standardised variant calling formats (MAF and VCF) as well as standard
39 Bioconductor classes (GRanges and VRanges). Katdetectr segments genomic variants
40 utilising unsupervised changepoint detection and the Pruned Exact Linear Time search
41 algorithm. Kataegis foci are flagged based on the historical definition, namely that a kataegis
42 foci is a continuous segment harbouring ≥ 6 variants and has a mean intermutation distance
43 ≤ 1000 bp. Additionally, the implementation of changepoint detection utilised by katdetectr
44 results in fast computation. Furthermore, katdetectr is available on Bioconductor which
45 ensures reliability, and operability on common operating systems (Windows, macOS and
46 Linux). Katdetectr is available on Bioconductor at
47 <https://www.bioconductor.org/packages/devel/bioc/html/katdetectr.html> and on GitHub at
48 <https://github.com/ErasmusMC-CCBC/katdetectr>. All code used for the performance
49 evaluation is available on GitHub at: [https://github.com/ErasmusMC-
50 CCBC/evaluation_katdetectr](https://github.com/ErasmusMC-CCBC/evaluation_katdetectr)

51 **Contact:** h.vandewerken@erasmusmc.nl

52

53 **Introduction**

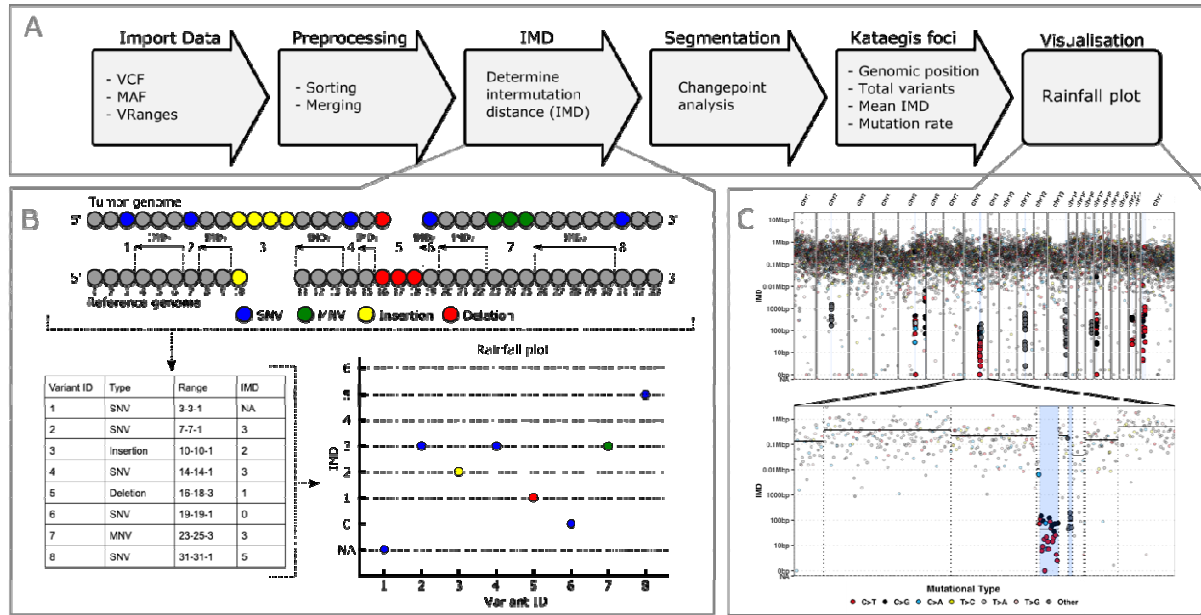
54 Next-generation sequencing of cancer genomes has revealed that mutations can cluster
55 together, i.e., the acquired mutations are found in proximity to one another, much closer
56 than would be expected if they had been dispersed uniformly throughout the genome
57 purely by chance (Alexandrov et al., 2013a; Nik-Zainal et al., 2012a). This phenomenon was
58 termed kataegis and its respective genomic location was termed a kataegis foci. Kataegis,
59 which is Greek for thunderstorm or shower, was first observed and visualised in whole
60 genome sequencing (WGS) data of 21 primary breast cancers (Nik-Zainal et al., 2012b).
61 Alexandrov et al. subsequently detected 873 kataegis foci in a pan-cancer dataset containing
62 507 WGS samples from primary malignancies (Alexandrov et al., 2013b).

63
64 Extensive exploration of the aetiology of kataegis revealed a significant positive correlation
65 between kataegis and two distinct mutational signatures both attributed to the APOBEC
66 enzyme-family (Alexandrov et al., 2020; Bergstrom, Luebeck, et al., 2022; Burns et al., 2013;
67 Taylor et al., 2013b).

68
69 Subsequently, multiple studies confirmed the importance of the APOBEC enzymes in cancer,
70 showing that APOBEC is a major cause of mutagenesis, both seen in clusters, dispersed
71 throughout the cancer genome and in extrachromosomal DNA (Bergstrom et al., 2021;
72 Bergstrom, Luebeck, et al., 2022; Langenbucher et al., 2021; Maciejowski et al., n.d.; Taylor
73 et al., 2013a).

74
75 Previous studies have shown that kataegis occurs within known cancer genes including
76 TP53, EGFR and BRAF which are associated with overall survival (Bergstrom, Luebeck, et al.,
77 2022). Still, the clinical significance of kataegis remains to be validated and therefore
78 obfuscates kataegis as a clinical biomarker for predicting prognosis. Nevertheless, any future
79 clinical application requires accurate and robust detection of kataegis.

80
81 Here, we introduce katdetectr, an R-based and Bioconductor package that contains a
82 complete suite for the detection, characterisation and visualisation of kataegis. Additionally,
83 we have evaluated the performance of katdetectr and five publicly available kataegis
84 detection packages (Bergstrom, Kundu, et al., 2022; Lin et al., 2021; Lora, 2016; Mayakonda
85 et al., 2018; Yousif et al., 2020).



86

87 **Figure 1, Overview of the katdetectr workflow, intermutation distance and rainfall plots.**

88 A) General workflow of katdetectr represented by arrows. B) The intermutation distance
 89 (IMD) is determined for each two subsequent genomic variants per chromosome and
 90 rainfall plots are used to visualise these IMDs and corresponding detected changepoint
 91 segments. C) Rainfall plot of PD7049a (breast cancer) from the Alexandrov dataset as
 92 interrogated by katdetectr (Alexandrov et al., 2013a). Y-axis: IMD, x-axis: variant ID ordered
 93 on genomic appearance, light blue rectangles: kataegis foci with genomic variants within
 94 kataegis foci shown in bold. The mutational type is depicted by the colour. The determined
 95 segmentation (as mean IMD per segment) is shown by black horizontal solid lines whilst
 96 vertical lines represent detected changepoints. Note that the first variant of a kataegis foci
 97 has a high IMD due to the usage of the upstream-oriented IMD.

98

99

100 Approach

101 Katdetectr was programmed in the R statistical programming language (v4.1.2) (R Core
 102 Team, 2022). Briefly, katdetectr can import standardised formats denoting genomic variants
 103 including: Variant Calling Format (VCF), Mutation Annotation Format (MAF) and VRanges
 104 objects. Per sample, the genomic variants are pre-processed and subsequently the
 105 upstream-oriented intermutation distance (IMD) is calculated (Nik-Zainal et al., 2012a). The
 106 distribution of IMDs is then segmented based on unsupervised detection of changepoints
 107 using the changepoint package (v2.2.3) and the Pruned Exact Linear Time (PELT) search
 108 method (Haynes et al., 2017; Haynes & Killick, 2021; Killick et al., 2012; Killick & Eckley,
 109 2014).

110

111 After segmentation, putative kataegis foci are called based on the following definition: 1) a
 112 continuous segment harbouring ≥ 6 variants and 2) the captured IMDs within the segment
 113 contain a mean IMD of ≤ 1000 bp (Alexandrov et al., 2013a). Moreover, katdetectr can

114 visualise the IMD, changepoints and their continuous segments and can highlight all
 115 putative kataegis foci within a sample using an intuitive rainfall plot (Figure 1).
 116 The output of katdetectr consists of an S4 object containing the putative kataegis foci
 117 (GRanges), the annotated genomic variants (VRanges) and the annotated segments
 118 (GRanges).

119

120 See supplementary note 1 for an extended description of the design of katdetectr and
 121 parameters settings.

122

123 Method

124 The performance of katdetectr (v1.0.0) was compared to alternative packages by utilising an
 125 in-house generated synthetic dataset containing 1024 samples and a publicly available pan-
 126 cancer dataset containing 507 WGS samples with a priori labelled kataegis foci as curated by
 127 Alexandrov et al. (2013) (Alexandrov et al., 2013a; Bergstrom, Kundu, et al., 2022; Lin et al.,
 128 2021; Lora, 2016; Mayakonda et al., 2018; Yousif et al., 2020).

129

130 In order to quantify and compare performances, the task of kataegis detection was reduced
 131 to a binary classification problem. The task of the kataegis detection packages was to
 132 correctly label each variant for kataegis, i.e., whether or not a genomic variant lies within a
 133 kataegis foci.

134

135 In order to assess performance related to sample-specific Tumour Mutational Burden
 136 (TMB), we binned samples based on TMB. The synthetic dataset contained eight TMB
 137 classes (0.1, 0.5, 1, 5, 10, 50, 100, 500) whilst the Alexandrov dataset was binned into three
 138 TMB classes (low: TMB < 0.1, middle: 0.1 ≥ TMB < 10, high: TMB ≥ 10).

139 Due to large class imbalance, we used the normalised Matthews Correlation Coefficient
 140 (nMCC) as the main performance metric for performance evaluation (Chicco & Jurman,
 141 2020).

142

143 See supplementary note 1 for an extended description of the datasets, synthetic data
 144 generation and confusion matrices.

145

Performance kataegis classification

| | Package | Reference | Language | Synthetic dataset | | | | | Dataset labelled by Alexandrov et al. | | | | |
|---|---------------------|---------------------------------|----------|-------------------|-------------|-------------|-------------|-------------|---------------------------------------|-------------|-------------|-------------|-------------|
| | | | | Accuracy | nMCC | F1 | TPR | TNR | Accuracy | nMCC | F1 | TPR | TNR |
| 1 | katdetectr | Hazelaar, van Riet et al., 2022 | R | <u>0.99</u> | <u>0.98</u> | <u>0.97</u> | 0.94 | 0.99 | <u>0.99</u> | <u>0.92</u> | <u>0.83</u> | 0.91 | <u>0.99</u> |
| 2 | SeqKat | Taylor et al., 2013 | R | 0.84 | 0.54 | 0.02 | 0.93 | 0.84 | 0.99 | 0.85 | 0.69 | 0.59 | 0.99 |
| 3 | MafTools | Mayakonda et al., 2018 | R | 0.74 | 0.53 | 0.01 | 0.96 | 0.74 | 0.99 | 0.85 | 0.66 | 0.93 | 0.99 |
| 4 | SigProfilerClusters | Bergstrom, Kundu, et al., 2022 | Python | 0.65 | 0.52 | 0.01 | 0.88 | 0.65 | 0.99 | 0.84 | 0.68 | 0.66 | <u>0.99</u> |
| 5 | ClusteredMutations | Lora, 2016 | R | 0.70 | 0.53 | 0.01 | <u>0.99</u> | 0.74 | 0.99 | 0.83 | 0.61 | <u>0.99</u> | 0.99 |
| 6 | kataegis | Lin et al., 2021 | R | 0.99 | 0.80 | 0.52 | 0.36 | <u>0.99</u> | 0.99 | 0.56 | 0.03 | 0.02 | 0.99 |

146

147 **Table 1, performance metrics of evaluated kataegis detection packages.** Accuracy,
 148 normalized Matthews Correlation Coefficient (nMCC), F1 score, True Positive Rate (TPR) and
 149 True Negative Rate (TNR) of the kataegis detection packages on 1024 synthetic samples and
 150 507 a priori labelled WGS samples (Alexandrov et al., 2013a). Rows were sorted in
 151 descending order based on nMCC score on the Alexandrov dataset (grey transparent
 152 background). For each performance metric, the highest score is underlined.

153

154 **Results**

155 Out of all evaluated packages, katdetectr revealed the highest overall accuracy and nMCC in
156 correctly labelling kataegis foci within both the synthetic and Alexandrov et al. dataset
157 (Table 1). The performance of all packages was found to be associated with the sample-
158 respective TMB (**Supplementary Figure 1**). Performance evaluation per TMB-binned
159 category revealed that katdetectr is on par with alternative packages for samples with TMB
160 ≤ 50 . However, in contrast to alternative packages, the nMCC of katdetectr remains high for
161 samples with high TMB (ranging between 50-500; **Supplementary Figures 2-3**).
162 Furthermore, katdetectr demonstrated the fastest computational runtimes of all evaluated
163 packages (**Supplementary Figures 4**).

165 **Conclusion**

166 Here, we described katdetectr; an R-based Bioconductor package capable of the detection,
167 characterization and visualization of putative kataegis foci within genomic variants.
168 Performance evaluation revealed that katdetectr robustly detects kataegis in a wide range
169 of malignancies, irrespectively of low or high TMB. Additionally, katdetectr is user-friendly
170 and computationally inexpensive with fast runtimes. In conclusion, the robust and
171 reproducible methodologies of katdetectr can help facilitate further research into the
172 clinical significance and underlying biological mechanism of kataegis.

174 **Acknowledgements**

175 We would like to thank John Martens, Marcel Smid and Guido Jenster for their discussions,
176 input and support. Additionally, we would like to thank Coen Berns and Yi Ping for their
177 initial efforts on detecting kataegis.

179 **Funding**

180 This research received funding from the Daniel den Hoed Fonds - Cancer Computational
181 Biology Center (DDHF-CCBC) grant.
182 Conflict of Interest: none declared.

184 **Data availability**

185 All data used in the performance evaluation can be found on Zenodo at:
186 <https://zenodo.org/record/6623289#.YqBxHi8RrOo>

189 **References**

- 190 Alexandrov, L. B., Kim, J., Haradhvala, N. J., Huang, M. N., Tian Ng, A. W., Wu, Y., ... &
191 Stratton, M. R. (2020). The repertoire of mutational signatures in human cancer. *Nature*,
192 578(7793), 94-101. <https://doi.org/10.1038/s41586-020-1943-3>
- 193
194 Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Aparicio, S. A., Behjati, S., Biankin, A. V., ... &
195 Stratton, M. R. (2013). Signatures of mutational processes in human cancer. *Nature*,
196 500(7463), 415-421. <https://doi.org/10.1038/nature12477>
- 197
198 Bergstrom, E. N., Kundu, M., Tbeileh, N., & Alexandrov, L. B. (2022). Examining clustered
199 somatic mutations with SigProfilerClusters. *bioRxiv*.
200 <https://doi.org/10.1093/BIOINFORMATICS/BTAC335>

201

202 Bergstrom, E. N., Luebeck, J., Petljak, M., Khandekar, A., Barnes, M., Zhang, T., ... &
203 Alexandrov, L. B. (2022). Mapping clustered mutations in cancer reveals APOBEC3
204 mutagenesis of ecDNA. *Nature*, 602(7897), 510-517.

205 <https://doi.org/10.1038/s41586-022-04398-6>

206

207 Bergstrom, E. N., Luebeck, J., Petljak, M., Bafna, V., Mischel, P. S., Harris, R. S., & Alexandrov,
208 L. B. (2021). Comprehensive analysis of clustered mutations in cancer reveals recurrent
209 APOBEC3 mutagenesis of ecDNA. *bioRxiv*. <https://doi.org/10.1101/2021.05.27.445689>

210

211 Burns, M. B., Lackey, L., Carpenter, M. A., Rathore, A., Land, A. M., Leonard, B., ... & Harris,
212 R. S. (2013). APOBEC3B is an enzymatic source of mutation in breast cancer. *Nature*,
213 494(7437), 366-370. <https://doi.org/10.1038/NATURE11881>

214

215 Chicco, D., & Jurman, G. (2020). The advantages of the Matthews correlation coefficient
216 (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics*, 21(1),
217 1–13. <https://doi.org/10.1186/s12864-019-6413-7>

218 Haynes, K., Fearnhead, P., & Eckley, I. A. (2017). A computationally efficient nonparametric
219 approach for changepoint detection. *Statistics and computing*, 27(5), 1293-1305.

220 <https://doi.org/10.1007/s11222-016-9687-5>

221

222 Haynes, K., & Killick, R. (2021). changepoint.np: Methods for Nonparametric Changepoint
223 Detection. <https://CRAN.R-project.org/package=changepoint.np>

224

225 Killick, R., & Eckley, I. (2014). changepoint: An R package for changepoint analysis. *Journal of*
226 *statistical software*, 58(3), 1-19.

227

228 Killick, R., Fearnhead, P., & Eckley, I. A. (2012). Optimal detection of changepoints with a
229 linear computational cost. *Journal of the American Statistical Association*, 107(500), 1590–
230 1598. <https://doi.org/10.1080/01621459.2012.737745>

231

232 Langenbucher, A., Bowen, D., Sakhtemani, R., Bournique, E., Wise, J. F., Zou, L., ... &
233 Lawrence, M. S. (2021). An extended APOBEC3A mutation signature in cancer. *Nature*
234 *communications*, 12(1), 1-11. <https://doi.org/10.1038/s41467-021-21891-0>

235

236 Lin, X., Hua, Y., Gu, S., Lv, L., Li, X., Chen, P., Dai, P., Hu, Y., Liu, A., & Li, J. (2021). kataegis: an
237 R package for identification and visualization of the genomic localized hypermutation
238 regions using high-throughput sequencing. *BMC Genomics*, 22(1), 1–6.

239 <https://doi.org/10.1186/s12864-021-07696-x>

240

241 Lora, D. (2016). ClusteredMutations: Location and Visualization of Clustered Somatic
242 Mutations. <https://CRAN.R-project.org/package=ClusteredMutations>

243 Maciejowski, J., Chatzipli, A., Dananberg, A., Chu, K., Toufektchan, E., Klimczak, L. J., ... & de

244

245 Lange, T. (2020). APOBEC3-dependent kataegis and TREX1-driven chromothripsis during
246 telomere crisis. *Nature genetics*, 52(9), 884-890.

247 <https://doi.org/10.1038/s41588-020-0667-5>

248

249 Mayakonda, A., Lin, D. C., Assenov, Y., Plass, C., & Koeffler, H. P. (2018). Maftools: efficient
250 and comprehensive analysis of somatic variants in cancer. *Genome research*, 28(11), 1747-
251 1756. <https://doi.org/10.1101/gr.239244.118>

252

253 Nik-Zainal, S., Alexandrov, L. B., Wedge, D. C., Van Loo, P., Greenman, C. D., Raine, K., ... &
254 Breast Cancer Working Group of the International Cancer Genome Consortium. (2012).
255 Mutational processes molding the genomes of 21 breast cancers. *Cell*, 149(5), 979-993.
256 <https://doi.org/10.1016/j.cell.2012.04.024>

257

258 R Core Team. (2022). R: A Language and Environment for Statistical Computing.
259 <https://www.R-project.org/>

260

261 Taylor, B. J., Nik-Zainal, S., Wu, Y. L., Stebbings, L. A., Raine, K., Campbell, P. J., ... &
262 Neuberger, M. S. (2013). DNA deaminases induce break-associated mutation showers with
263 implication of APOBEC3B and 3A in breast cancer kataegis. *elife*, 2.
264 doi: 10.7554/eLife.00534

265

266 Yousif, F., Lin, X., Fan, F., Lalansingh, C., & Macdonald, J. (2020). SeqKat: Detection of
267 Kataegis. <https://CRAN.R-project.org/package=SeqKat>