

# 1 PanelCAT: an Open-Source Comparative Analysis Tool for 2 Next-Generation Sequencing Panel Target Regions

3 André Oszwald<sup>1</sup>, Lucia Zisser<sup>2</sup>, Eva Compérat<sup>1</sup>, Leonhard Müllauer<sup>1</sup>

4 <sup>1</sup>Department of Pathology, Medical University of Vienna, Währinger Gürtel 18-20, 1090 Vienna,  
5 Austria

6 <sup>2</sup>Department of Biomedical Imaging and Image-guided Therapy, Division of Nuclear  
7 Medicine., Medical University of Vienna, Währinger Gürtel 18-20, 1090 Vienna, Austria

## 8 Abstract

9 Multi-gene next-generation sequencing (NGS) panels have become a routine diagnostic method in  
10 the contemporary practice of personalised medicine. To avoid inadequate test choice or  
11 interpretation, a detailed understanding of the precise panel target regions is required. However, the  
12 necessary bioinformatic expertise is not always available, and publicly accessible and easily  
13 interpretable analyses of target regions are scarce. To address this critical knowledge gap, we present  
14 the Panel Comparative Analysis Tool (PanelCAT) an open-source application to analyze, visualize and  
15 compare NGS panel DNA target regions. PanelCat uses RefSeq, ClinVar and COSMIC cancer mutation  
16 census databases to quantify the exon and mutation coverage of target regions and provides  
17 interactive graphical representations and search functions to inspect the results. We demonstrate  
18 the utility of PanelCAT by analyzing two large NGS panels (Illumina TSO500 and Qiagen pan-cancer  
19 panel) to validate the advertised target genes, quantify targeted exons and mutations, and identify  
20 differences between panels. PanelCat will enable institutions and researchers to catalogue and  
21 visualize NGS panel target regions independent of the manufacturer, promote transparency of panel  
22 limitations, and share this information with employees and requisitioners.

## 23 Introduction

24 Precision oncology routinely involves next-generation sequencing (NGS) of tumor DNA to identify  
25 therapeutically actionable targets or diagnostically relevant mutations that critically direct patient  
26 management<sup>1</sup>. Most multi-gene sequencing panels do not cover entire genes, but only variable  
27 portions of genes that are considered most relevant, i.e., predominantly protein-coding sequences  
28 and tumor mutational hotspots. For this reason, both the choice of an adequate test and its  
29 interpretation, especially regarding the certainty of negative findings, crucially depend on detailed  
30 knowledge of the portions of genes and genetic alterations that may be assessed by a panel.

31 Target regions of commercial NGS panels are typically specified in a panel-specific BED file by a list of  
32 chromosome numbers, start and stop coordinates<sup>2</sup>. Although this information is an essential part of  
33 the test documentation, it is not useful to understand panel target regions in detail without further  
34 analysis for several reasons: it does not inform on the non-targeted portions of genes without  
35 comparison to a reference genome; the provided information on target genes, transcripts and exons  
36 is not updated alongside the transcript databases (e.g., RefSeq<sup>3</sup>); the target regions must be  
37 systematically compared to mutation/variant databases in order to determine pathogenic mutations  
38 that can be detected; lastly, genomic positions with known high rates of erroneous variant calls are  
39 often masked during secondary analysis, but these positions are defined in separate files.

40 Consequently, the lack of detailed publicly available data on precise panel targets, and the barriers to  
41 generate it due to the required bioinformatic expertise, portend the risk of inadequate test choice  
42 and test misinterpretation. To reduce this risk, we developed the “Panel Comparative Analysis Tool”  
43 (PanelCAT), an application that allows to analyze, visualize and compare DNA target regions of NGS  
44 panels within a user-friendly interface, and provides a platform to clearly communicate this  
45 information to others. We demonstrate the use of this tool by analyzing two large multi-gene NGS  
46 panels, the Illumina TrueSight Oncology 500 (TSO) and Qiagen Pan-Cancer (QPC) Panel, in order to

47 provide a more detailed documentation of their targeted genes, exons, known pathogenic mutations,  
48 and differences between the panels, than has been available to date.

## 49 **Materials and Methods**

### 50 **PanelCat code**

51 PanelCat code was created, and all analyses were performed, in R<sup>4</sup> v4.3.0 within RStudio v2023.03.0.

52 Analysis of genomic ranges (target regions and variant coordinates) were performed using the

53 GenomicFeatures<sup>5</sup> package. Graphs were drawn using ggplot2<sup>6</sup> and plotly<sup>7</sup> packages. A browser-

54 based implementation of the script was created using ShinyR<sup>8</sup>. PanelCat is provided under the open-

55 source license AGPLv3 and the source code, and R session information is available at

56 <https://github.com/aoszwald/panelcat>. The basic procedure of the panel analysis is outlined below.

57 PanelCat accepts target region files as input (containing columns for chromosome, start, and end

58 position of target regions), and optionally the mask region file (also containing chromosome, start

59 and end coordinates). The application first determines the intersection between the panel target

60 regions and RefSeq<sup>3</sup> exon coordinates in order to systematically identify target genes, and

61 subsequently all exon ranges of target genes. The exon ranges of each targeted gene are then

62 intersected with the panel target regions to quantify the targeted portion of protein coding bases per

63 gene. Targeted mutations are then identified by intersection of the panel target regions with the

64 coordinates of mutations in the ClinVar<sup>9</sup> and COSMIC<sup>10</sup> databases. Optionally, the mask file is

65 incorporated in the analysis to identify and determine the portion of masked bases and mutations.

66 The summarized output data are combined into lists of items and saved as R data objects. Panels that

67 were previously analyzed and saved in this form are pre-loaded the next time the application is

68 started. The panel output files can also be used for analysis outside of PanelCat; within R, the panel

69 data and listed sub-items can be accessed via the "\$" operator. In the course of this study, individual

70 data were accessed in this way and further processed in R independently of the PanelCat functions

71 to answer specific questions, including the cumulative percentage of non-targeted mutations, and  
72 discrepancies between the advertized gene list and the confirmed gene list.

### 73 Data sources

74 BED files indicating target regions of NGS panels and corresponding mask files, including Illumina  
75 TSO500 and Qiagen Pan-Cancer Panels, were obtained from the customer support of the  
76 manufacturers, or obtained in the course of using a product. The TSO500 mask file was provided by  
77 Illumina. These files are not provided as part of PanelCat.

78 The following databases are not provided as part of the software download, and need to be either  
79 downloaded manually (COSMIC) or automatically by PanelCat (ClinVar and ClinVar assembly report,  
80 RefSeq):

81 ClinVar data (last accessed 25.05.2023) was obtained from  
82 [https://ftp.ncbi.nlm.nih.gov/pub/clinvar/vcf\\_GRCh37/weekly/clinvar.vcf.gz](https://ftp.ncbi.nlm.nih.gov/pub/clinvar/vcf_GRCh37/weekly/clinvar.vcf.gz). The COSMIC cancer  
83 mutation census data (v98, last accessed 23.05.2023) was obtained from  
84 <https://cancer.sanger.ac.uk/cosmic/download>. RefSeq data (last accessed 25.05.2023) was obtained  
85 from  
86 [https://ftp.ncbi.nlm.nih.gov/refseq/H\\_sapiens/annotation/GRCh37\\_latest/refseq\\_identifiers/GRCh37\\_latest\\_genomic.gff.gz](https://ftp.ncbi.nlm.nih.gov/refseq/H_sapiens/annotation/GRCh37_latest/refseq_identifiers/GRCh37_latest_genomic.gff.gz), and the GrCh37. The NCBI assembly report was obtained (last accessed  
87 25.05.2023) from  
88 [https://ftp.ncbi.nlm.nih.gov/genomes/refseq/vertebrate\\_mammalian/Homo\\_sapiens/annotation\\_releases/105.20220307/GCF\\_000001405.25\\_GRCh37.p13/GCF\\_000001405.25\\_GRCh37.p13\\_assembly\\_report.txt](https://ftp.ncbi.nlm.nih.gov/genomes/refseq/vertebrate_mammalian/Homo_sapiens/annotation_releases/105.20220307/GCF_000001405.25_GRCh37.p13/GCF_000001405.25_GRCh37.p13_assembly_report.txt).

## 92 Results

### 93 General function of PanelCat

94 PanelCat (<https://github.com/aoszwald/panelcat> and <https://aoszwald.shinyapps.io/panelcat/>)  
95 provides functions to automatically distill descriptive information from panel target region files and

96 public databases, and to display this data to facilitate evaluation and comparison of panels. To  
97 analyze a panel, PanelCat is provided with the target region file (typically with a .bed file suffix, but  
98 others may be acceptable), and optionally a mask file (indicating regions where variant calls are  
99 unreliable and will be filtered out). The application then determines the overlap between target  
100 regions and protein-coding bases per gene in RefSeq<sup>3</sup>, known pathogenic and likely pathogenic  
101 mutations in ClinVar<sup>9</sup> and tier 1-3 oncogenic mutations in COSMIC<sup>10</sup> cancer mutation census (CMC)  
102 databases. The output data is saved in a compact form that can be used in PanelCat or explored  
103 independently in R Statistics<sup>4</sup>.

104 The application provides several visualization options to analyze and compare panels, briefly outlined  
105 here. The first option is a scatter plot, where users can contrast target coverage metrics of RefSeq,  
106 ClinVar and COSMIC databases from all analyzed panels. This function can be used to compare  
107 between panels, e.g., to determine differences in target gene coverage (Fig. 1A), or within panels,  
108 e.g., to evaluate the relationship between coverage of protein-coding bases and pathogenic  
109 mutations. The second option is a horizontal column plot showing per-gene coverage, and the  
110 masked portion, of protein-coding bases, ClinVar variants and COSMIC mutations, whereby users can  
111 select multiple panels for comparison (Fig. 1B). The third visualization is a column plot optimized to  
112 search for and display one or multiple specific genes of interest across several (or all) analyzed  
113 panels, with indication of a panel that targets all searched genes (Fig. 1C). The fourth visualization is  
114 a column plot of the estimated frequency of COSMIC CMC mutations that are not targeted by the  
115 panels, both in target genes only and in all genes, with and without considering variant masking (Fig.  
116 1D). Lastly, PanelCat provides a graph to compare base coverage of individual exons per transcript  
117 between two panels in paired violin plots (Fig. 1E).

118 Besides graphs, PanelCat provides table views to compare raw numerical data, or investigate the  
119 coverage of specific exons and mutations. In the first option, users can inspect the PanelCat output  
120 data (representing the raw data from which the graphs are generated), whereby multiple panels and

121 metrics can be simultaneously displayed and searched for specific genes. The second option displays  
122 the coverage of each exon of every transcript of all target genes of a panel (Fig. 1F). Similarly, the  
123 third and fourth options provide for each panel a complete list of targeted ClinVar variants and  
124 COSMIC mutations, including the ability to display masked variants. The exon, ClinVar and COSMIC  
125 tables can be searched/filtered for each column independently, e.g., for specific genes, transcripts,  
126 exons, coding or amino acid changes, or genome coordinates.

127 When run locally, PanelCat automatically obtains the current ClinVar (released weekly) and RefSeq  
128 databases upon first use. These, along with all previously processed panels, can be updated in a  
129 single step; previous versions of databases and panel analyses are stored for later reference and  
130 documentation. The COSMIC CMC database (updated every several months) requires manual  
131 download and replacement of the local file.

132 In summary, PanelCat offers multiple useful and intuitive functions to substantially improve the  
133 transparency and accessibility of NGS panel target region documentation. The newly developed tool  
134 was next used to analyze two very large panels currently used in both clinical and research settings;  
135 the Illumina TrueSight Oncology 500 (TSO) and the Qiagen Pan-Cancer (QPC) panel. Although the  
136 target regions can be requested as part of the panel documentation, the explicit coverage of genes  
137 and mutations are not provided. However, informed clinical use requires detailed information, so we  
138 used PanelCat to characterize their target regions in detail and explore their subtle differences.

### 139 Mutation coverage in Qiagen PanCancer panel is similar or greater than TSO500, 140 despite lower exon coverage

141 In their respective product documentation, the TSO and QPC panels advertise the same 523 gene  
142 targets for analysis of small variants (e.g, SNV, insertions and deletions). We first compared the  
143 advertized genes to the target genes identified using PanelCat. The QPC target regions overlapped  
144 with exons of 603 genes, including all advertized genes. By contrast, the TSO target regions  
145 overlapped with exons of 625 genes, but these included only 521 of the 523 advertised targets. The

146 two remaining target genes (HLA-B and HLA-C) do not overlap with the TSO target regions;  
147 accordingly, we did not find any variant calls in HLA-B or HLA-C in a representative set of 400 samples  
148 analyzed with the TSO panel (including unfiltered variant call files in 10 samples). Alterations in HLA-I  
149 genes have been postulated to promote tumor evasion of immune surveillance, e.g., by restricting  
150 neoantigen presentation<sup>11,12</sup>, although no guidelines recommendations to test HLA genes exist  
151 momentarily.

152 We next identified the targeted exon-coding bases of each gene. We first searched for genes with the  
153 greatest coverage and found that in the TSO, 20 genes had exon coverage > 95% (including NAB2,  
154 TERC, CD74, TFE3, KIF5B, EML4, EWSR1, FLI1, ETV1, ETV5, PAX3, all over 99%), whereas in the QPC, it  
155 was only six (TERC, ZRSR2, ATR, POLD1, KMT2B, RECQL4). We found a strong direct correlation  
156 between the base coverage of TSO and QPC panels (Pearson's  $r = 0.81$ ,  $p < 2e-16$ ), and no significant  
157 difference in mean exon base coverage per gene (TSO 50.3% vs. QPC 48.4%,  $p = 0.23$ ) (Fig 2A). We  
158 identified target genes where relative coverage was considerably greater in the TSO than in the QPC,  
159 including NTRK2, ETV1, AKT3, ERG, and PAX7, but only few genes with greater coverage in the QPC  
160 panel, notably PMS2, TERT, HLA-B and HLA-C (Fig. 2B). Importantly, four targeted genes (HLA-A,  
161 KMT2B, KMT2C, KMT2D) showed total masking of all target regions in the TSO panel (but not in the  
162 QPC, which does not use a mask file). Accordingly, we did not find any variants calls in these genes in  
163 a representative set of 400 samples analyzed with the TSO panel.

164 ClinVar lists approximately 50,000 known variants labelled "pathogenic" or "likely pathogenic" in the  
165 advertized TSO and QPC target genes. While the TSO targeted 92.5% of these variants (94.8%  
166 without variant masking), the QPC targeted 97.4%, despite lower exon coverage. Consequently,  
167 targeting of all pathogenic variants was achieved for 182 genes in the TSO panel (200 without  
168 masking), and 223 genes in the QPC panel. In the TSO panel, no pathogenic variants were targeted in  
169 8 genes (of which 6 were due to variant masking), whereas in the QPC it was only three. In both  
170 panels, the majority (QPC: 51%, TSO: 64%) of non-targeted variants occurred in two similar sets of

171 only 10 genes, in both cases including NF1 and the DNA repair genes MLH1, MSH2, BRCA1, BRCA2  
172 and ATM (Supplemental Table 1). Most differences between panels were attributable to either  
173 greater coverage in the QPC (PMS2, TERT) or extensive masking in the TSO (KMT2B, KMT2C, KMT2D)  
174 (Fig 2C, 2D).

175 The COSMIC cancer mutation census database (CMC, v98) lists ~43,000 unique mutations occurring  
176 in genes targeted by the TSO and QPC panels, of which the majority (93.4%, or 99.5% without variant  
177 masking) are targeted by the TSO, and all by the QPC (100%) (Fig. 2E, 2F). Due to masking, no  
178 mutations are targeted by the TSO500 in HLA-A, KMT2C and KMT2D. Independent of masking, the  
179 QPC panel more extensively targeted mutations in NF1, TERT and PMS2 than the TSO. We estimated  
180 the frequency of samples to harbor non-targeted mutations by calculating the positive sample  
181 proportion of unique mutations in the CMC dataset, and cumulating the frequency all non-targeted  
182 mutations. The rate of non-covered CMC tier 1-3 mutations in targeted genes per sample was lower  
183 in the QPC (0.005) than in the TSO (0.14, 0.02 without masking), suggesting that one in eight (TSO, or  
184 one in 50 without masking) or one in 200 (QPC) samples would harbor oncogenic mutations that  
185 cannot be detected with the panels, in one of the panel target genes.

## 186 Discussion

187 Detailed knowledge of the target regions of NGS panels is important for the correct choice and  
188 interpretation of molecular tests, but is not typically well illustrated by the test manufacturer, and  
189 usually requires bioinformatic analysis to acquire. In this study, we present PanelCat, a novel open-  
190 source tool that can be used by laboratories or NGS panel distributors to analyze NGS target regions  
191 and share this information to enable more informed decisions. A limitation of PanelCat is that it does  
192 not assess fusion or copy number events detected by panels, but respective features may be  
193 implemented in future.

194 PanelCat enables rapid assessment and rich visualisation of the designed target regions of NGS  
195 panels without bioinformatic expertise. As an example, we expanded in detail on the existing and

196 incomplete documentation of two large NGS tests (Illumina TSO500 and Qiagen Pan-Cancer).  
197 PanelCat quantified precise exon coverage and identified of genes with poor coverage, extensive  
198 variant masking, differences between panels, and even discrepancies to the advertised gene list.  
199 Thus, we found that unlike the QPC, the TSO500 does not target HLA-B and HLA-C; that KMT2B,  
200 KMT2C and KMT2D are extensively masked in the TSO500 and will not yield variant calls after  
201 filtering, and that PMS2 and TERT are substantially better covered in the QPC panel independent of  
202 variant masking. In addition, we used PanelCat to describe the different coverage of individual exons  
203 of PMS2 in the panels.

204 PanelCat is different from the Panel Informativity Optimizer (PIO) method, previously demonstrated  
205 to assist in optimizing NGS panel design<sup>13</sup>. Because PanelCat does not generate new target regions, it  
206 only indirectly assists in panel design by highlighting deficits in exon or mutation coverage in  
207 particular genes of interest. Compared to PIO, PanelCat provides superior functions to analyze and  
208 compare existing (or proposed) panels. Crucially, PIO cannot process complex target regions or  
209 conventional BED-format files, only lists of complete genes or exons. Most panels target incomplete  
210 genes or exons, and would thus be inaccurately represented using PIO. Consequently, the limited  
211 panel benchmarking functions of PIO cannot inform on precise exons coverage, whereas PanelCat  
212 provides detailed information on the level of genes, exons and individual mutations. In contrast to  
213 PIO, which provides a linear data pipeline from input to output, PanelCat is a platform to collect  
214 panel analyses and visualize them for frequent inspection in a routine clinical setting. In summary,  
215 PanelCat provides opportunities that have not yet been demonstrated with previous methods, albeit  
216 with features designed more for panel end-users than panel developers.

217 In contrast to PIO, PanelCat does not use a variety of mutation databases to account for the  
218 heterogeneity of mutation frequencies across different disease entities. Although panels are often  
219 designed for specific disease entities or groups thereof, many widely used panels (e.g., Thermo  
220 Fisher Oncomine Focus, or Illumina TSO500) were designed to cover a wide range of disease entities,

221 and we therefore also chose a disease-agnostic approach for PanelCat. However, the variant  
222 databases used by PanelCat can be pre-processed by users to focus the analysis entirely on  
223 mutations that are relevant in a specific disease context.

224 PanelCat was designed for users with limited or no IT support. For this reason, the software is  
225 designed to function on a local device (without installation of software besides R statistics); however,  
226 a slightly modified script can be hosted using ShinyServer to provide a network service. PanelCat  
227 reference databases can be easily updated, and stored panel analyses can be managed within the  
228 operating system's file system.

229 Although multi-gene NGS panels are currently the standard procedure in many institutions, routine  
230 whole exome sequencing of tumor specimens is being increasingly performed. Due to the high  
231 performance of the underlying packages<sup>5</sup>, PanelCat could be used to analyze target regions of whole  
232 exome panels. However, the increased rendering time of some of the implemented visualization  
233 methods could be impractical. Nevertheless, the PanelCat output data, saved as R objects, could be  
234 used outside of PanelCat to plot custom graphs demanding less computation.

235 In conclusion, we present PanelCat as a powerful solution to current shortcomings in the  
236 presentation, analysis and awareness of NGS panel target regions. We believe this software will  
237 improve the transparency of NGS panels and facilitate more informed decisions in test choice and  
238 interpretation, thus constituting a valuable addition to the expanding repertoire of available tools.

## 239 **References**

- 240 1. Mosele, F. *et al.* Recommendations for the use of next-generation sequencing (NGS) for patients  
241 with metastatic cancers: a report from the ESMO Precision Medicine Working Group. *Ann. Oncol.*  
242 **31**, 1491–1505 (2020).
- 243 2. Niu, J., Denisko, D. & Hoffman, M. M. The Browser Extensible Data (BED) format.
- 244 3. O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status, taxonomic  
245 expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733-745 (2016).

- 246 4. R Core Team. *R: A Language and Environment for Statistical Computing*. (R Foundation for  
247 Statistical Computing, 2020).
- 248 5. Lawrence, M. *et al.* Software for Computing and Annotating Genomic Ranges. *PLoS Comput. Biol.*  
249 **9**, (2013).
- 250 6. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag New York, 2016).
- 251 7. Sievert, C. *Interactive Web-Based Data Visualization with R, plotly, and shiny*. (Chapman and  
252 Hall/CRC, 2020).
- 253 8. Chang, W. *et al.* *shiny: Web Application Framework for R*. (2022).
- 254 9. Landrum, M. J. *et al.* ClinVar: improvements to accessing data. *Nucleic Acids Res.* **48**, D835–D844  
255 (2020).
- 256 10. Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res.* **47**,  
257 D941–D947 (2019).
- 258 11. Hazini, A., Fisher, K. & Seymour, L. Deregulation of HLA-I in cancer and its central importance for  
259 immunotherapy. *J. Immunother. Cancer* **9**, e002899 (2021).
- 260 12. Fangazio, M. *et al.* Genetic mechanisms of HLA-I loss and immune escape in diffuse large B cell  
261 lymphoma. *Proc. Natl. Acad. Sci.* **118**, e2104504118 (2021).
- 262 13. Alcazer, V. & Sujobert, P. Panel Informativity Optimizer: An R Package to Improve Cancer Next-  
263 Generation Sequencing Panel Informativity. *J. Mol. Diagn.* **24**, 697–709 (2022).
- 264

## 265 **Figure Legends**

266 **Figure 1: PanelCat user interface and visualisation options.** Users can select tabs with several graphs  
267 to investigate panel DNA target region coverage. Options include an X/Y scatterplot to compare any  
268 coverage metric from any previous analysed panel (A), barplots to view or compare absolute or  
269 relative RefSeq exon base/ClinVar variant/COSMIC mutation coverage per gene (B), a search function  
270 to compare specific genes of interest across multiple panels (C), a graph to visualise the estimated  
271 average rate of non-targeted COSMIC CMC tier 1-3 mutations (D), violin plots to compare coverage  
272 per exon for a transcript of interest between two panels (E), multiple table views to inspect Exons (F),  
273 ClinVar variants or COSMIC CMC mutations.

274 **Figure 2: Exon base and mutation coverage in Qiagen PanCancer and TSO500 panels.** Coverage of  
275 exon bases (A,B), ClinVar variants (C,D) and COSMIC CMC mutations (E,F). Note that the TSO uses a  
276 mask file to filter variant calls at positions with high error rates (shaded purple bars), but the QPC  
277 does not. Figures B, D, and F show only genes with large relative differences (filtered by fold changes  
278 of 2, 1.5, and 1.1, respectively)



