

Whole exome-seq and RNA-seq data reveal unique neoantigen profiles in Kenyan breast cancer patients

1 **Godfrey Wagutu¹, John Gitau^{1,2,3}, Kennedy Mwangi⁴, Mary Murithi⁵, Elias Melly⁸, Alexandra**
2 **R. Harris⁷, Shahin Sayed⁶, Stefan Ambs⁷, Francis Makokha^{1*}**

3 ¹Directorate of Research and Innovation, Mount Kenya University, Thika, Kenya

4 ²African Institute for Mathematical Science, Rwanda

5 ³Center for Epidemiological Modeling and analysis, Nairobi, Kenya

6 ⁴International Livestock Research Institute, Nairobi, Kenya

7 ⁵Kabarak University, Nakuru, Kenya

8 ⁶Aga Khan University Hospital, Nairobi, Kenya

9 ⁷Laboratory of Human Carcinogenesis, National Cancer Institute, Bethesda, USA

10 ⁸National Cancer Institute, Kenya

11 *** Correspondence:**

12 Francis Makokha

13 fmakokha@mku.ac.ke

14 **Keywords: Neoantigen, breast cancer, exome-seq, RNA-seq, somatic mutations**

15

16 ABSTRACT

17 **Background:** The immune response against tumors relies on distinguishing between self and non-
18 self, the basis of cancer immunotherapy. Neoantigens from somatic mutations are central to many
19 immunotherapeutic strategies and understanding their landscape in breast cancer is crucial for
20 targeted interventions. We aimed to profile neoantigens in Kenyan breast cancer patients using
21 genomic DNA and total RNA from paired tumor and adjacent non-cancerous tissue samples of 23
22 patients.

23 **Methods:** We sequenced the genome-wide exome (WES) and RNA, from which somatic mutations
24 were identified and their expression quantified, respectively. Neoantigen prediction focused on
25 human leukocyte antigens (HLA) crucial to cancer, HLA type I. HLA alleles were predicted from
26 WES data covering the adjacent non-cancerous tissue samples, identifying four alleles that were
27 present in at least 50% of the patients. Neoantigens were deemed potentially immunogenic if their
28 predicted median IC50 binding scores were $\leq 500\text{nM}$ and were expressed [transcripts per million
29 (TPM) >1] in tumor samples.

30 **Results:** An average of 1465 neoantigens covering 10260 genes had $\leq 500\text{nM}$ median IC50 binding
31 score and >1 TPM in the 23 patients and their presence significantly correlated with the somatic
32 mutations ($R^2 = 0.570$, $P = 0.001$). Assessing 58 genes reported in the catalog of somatic mutations in
33 cancer (COSMIC, v99) to be commonly mutated in breast cancer, 44 (76%) produced >2 neoantigens
34 among the 23 patients, with a mean of 10.5 ranging from 2 to 93. For the 44 genes, a total of 477
35 putative neoantigens were identified, predominantly derived from missense mutations (88%), indels
36 (6%), and frameshift mutations (6%). Notably, 78% of the putative breast cancer neoantigens were
37 patient-specific. HLA-C*06:01 allele was associated with the majority of neoantigens (194),
38 followed by HLA-A*30:01 (131), HLA-A*02:01 (103), and HLA-B*58:01 (49). Among the genes of

39 interest that produced putative neoantigens were *MUC17*, *TTN*, *MUC16*, *AKAP9*, *NEB*, *RP111*,
40 *CDH23*, *PCDHB10*, *BRCA2*, *TP53*, *TG*, and *RBI*.

41 **Conclusions:** The unique neoantigen profiles in our patient group highlight the potential of
42 immunotherapy in personalized breast cancer treatment as well as potential biomarkers for prognosis.
43 The unique mutations producing these neoantigens, compared to other populations, provide an
44 opportunity for validation in a much larger sample cohort.

45 INTRODUCTION

46 Breast cancer is among the most frequent causes of cancer-related mortality in women. Disease
47 heterogeneity and limited immunogenicity contribute to the lethality of breast cancer (Benvenuto et al.,
48 2019). Immune evasion, an important hallmark of cancer, adds to the complexity of cancer burden
49 through induction of immunosuppression (Bates et al., 2018). Immune checkpoint blockade (CKB)
50 therapy has been developed to target and block immune regulatory molecules (PD-1/PD-L1 and
51 CTLA-4) and in the process reactivate T cell immunity (Touchaei & Vahidi, 2024). This approach has
52 been reported to improve clinical responses and survival, especially in tumors with high mutational
53 burdens, such as lung cancer and melanoma (Shiravand et al., 2022). However, CKB therapy is not
54 universally successful among all patients and shows increased efficacy with higher mutational burden
55 tumors (Brahmer et al., 2015). Another immunotherapy approach that has been tested in clinical studies
56 is the targeting of tumor-associated antigens (TAAs) that are expressed in tumors at abnormally high
57 levels and rarely detectable in normal tissues (Valilou & Rezaei, 2019). One of the limitations of this
58 therapy approach is that many TAAs represent normal self-antigens and thus can be tolerated by T-
59 cells, resulting in poor immune response (Benvenuto et al., 2019). This poses challenges for
60 applicability in breast cancer because it generally has a lower mutational burden. Thus, CKB and TAAs
61 immunotherapy have had limited success in breast cancer patients (Narang et al., 2019).

62 Tumor neoantigens are tumor-specific antigens derived from somatic mutations in expressed
63 genes and are presentable to the major histocompatibility complex (MHC) by both class I human
64 leukocyte antigen (HLA-I) molecules present on surface of cancer cell, as well as class II HLA
65 molecules present on professional antigen-presenting cells (Blass & Ott, 2021). This elicits anti-tumor
66 immune responses that have the potential of eliminating the tumor cells with minimal off-target effects
67 (Pan et al., 2018). Neoantigens are encoded in various mutational types, including single nucleotide
68 substitution, insertion and deletions (INDELs), splice sites, stop codons gains and silent change, which

69 can result in translational frameshifts or novel open reading frames (Benvenuto et al., 2019). As such,
70 these neoantigens offer an advantage over TAAs in that they are only expressed by cancer cells and
71 not by normal cells, which enables specific recognition by the immune system (Benvenuto et al., 2019).
72 Although some neoantigens are shared among patients, most of them are patient-specific and are not
73 subject to immune tolerance mechanisms (Yarchoan et al., 2017). The specificity of neoantigens could
74 provide an opportunity for future personalized therapy in a cancer with a low tumor mutational burden
75 and a high disease heterogeneity, such as breast cancer. Moreover, neoantigens can potentially be used
76 as biomarkers in cancer immunotherapy to assess or predict the response of a patient to treatment
77 (Benvenuto et al., 2019).

78 Despite advancements in next generation sequencing and high-performance computing that has
79 resulted in improved cancer immunotherapy research and neoantigen-based treatments, there remains
80 a scarcity of information regarding neoantigens in specific populations from sub-Saharan African
81 countries such as Kenya. This lack of data poses a significant challenge in tailoring immunotherapeutic
82 strategies for breast cancer patients in such regions that have a high cancer burden, especially when
83 compounded by germline ancestral factors and a distinct mutational spectrum that may influence tumor
84 biology and immune response. Thus, it is critical to profile the neoantigen burden in this population to
85 contribute to the global collection of breast cancer immunogenic antigens for future drug development.
86 To this end, we sought to profile neoantigens in Kenyan women diagnosed with breast cancer *in silico*
87 through analysis of the whole exome and RNA sequencing data from 23 patients. We characterized the
88 mutation burden for each patient using WES, identified gene expression patterns in tumor tissue, and
89 predicted the putative neoantigens incorporating these datasets.

90

91 **MATERIALS AND METHODS**

92 **Patients and samples**

93 Tumor and adjacent normal tissue pairs were obtained from 23 breast cancer patients at the Aga
94 Khan Hospital, Nairobi, Kenya and AIC Kijabe Hospital, Kijabe, Kenya between 2019 and 2021.
95 Samples were collected through surgical excision, after which tissues were snap frozen in liquid
96 nitrogen and temporarily stored at Aga Khan Hospital. Frozen tissue samples were shipped to the
97 National Cancer Institute, Bethesda, MD, USA, for sequencing. Prior to tissue collection, all patients
98 provided written informed consent and the study was approved by Research and Ethics Committees at
99 Aga Khan University Hospital, Nairobi (Ref: 2018/REC-80) and AIC Kijabe Hospital (KH IERC-
100 02718/2019).

101 **Whole-exome sequencing (WES) and RNA-sequencing**

102 Genomic DNA was extracted from the samples using the DNeasy Blood and Tissue Kit (Qiagen,
103 Hilden, Germany), following manufacturer's instructions. Total RNA was extracted from the frozen
104 tissues using TRIzol reagent (Invitrogen). WES was performed by the company, Psomagen
105 (<https://www.psomagen.com/>). This service provider is Clinical Laboratory Improvement
106 Amendments-certified and College of American Pathologists (CAP)-accredited, achieving a sequence
107 depth of 250x for tumor tissues and 150x for adjacent non-cancerous tissues, as previously described
108 by us (Tang et al., 2023). Total RNA from the 23 sample pairs was processed by a NCI Leidos core
109 facility, where library preparation was performed using the TruSeq Poly A kit (Illumina, San Diego,
110 USA). Samples were sequenced on a Novaseq system with 150 bp paired-end reads and a depth of 30
111 million reads.

112 **Reads mapping and variant calling**

113 For WES, raw reads were quality checked using FASTQC (Andrews, 2010) and results
114 summarized using MultiQC (Ewels et al., 2016). The reads were trimmed for low quality reads and

115 adapter sequences using Trimmomatic (Bolger et al., 2014) and quality-checked again using FASTQC
116 and MultiQC. All samples passed the QC test after trimming and the reads were aligned using BWA-
117 MEM (Li, 2013) to the hg38 human reference genome, where >95% of the reads aligned properly to
118 the genome. The aligned reads were deduplicated and read groups added to the deduplicated bam files
119 using Picard. This was followed by base quality recalibration in GATK (McKenna et al., 2010).
120 Somatic variant calling was performed using MuTect2 (McKenna et al., 2010) in paired tumor-normal
121 mode utilizing the panel of normal option that was derived from normal reads. Variants were
122 normalized using a variant tool set (vt; Tan et al., 2015), filtered using GATK and
123 functional/consequence-annotated using a variant effect predictor (VEP; McLaren et al., 2016).
124 Annotated variants were converted to MAF files using *vcf2maf* (Kandath et al., 2020) and concatenated
125 into a single file. The MAF files were imported into R package *mafTools* (Mayakonda et al., 2018) for
126 further processing.

127 For RNA-seq, a quality check was performed using FASTQC and MultiQC after which the
128 reads were trimmed and quality checked again. All samples passed the quality check and the reads
129 were pseudo-aligned to the hg38 reference genome using Kallisto aligner (Bray et al., 2016) with
130 default settings to obtain count matrix. Alignment statistics showed that over >50% reads mapped
131 uniquely to the genome. The raw counts were normalized into estimated Transcripts Per Million (TPM),
132 and scaled using the average transcript length over samples and the library size by *tximport* (Soneson
133 et al., 2016).

134 **Variant expression annotation**

135 VCF files containing the variants were annotated for expression using the vcf-expression-
136 annotator (<https://github.com/griffithlab/VAtools>) with default setting except for choosing the use of
137 gene names instead of transcripts and thereby ignoring the Ensembl id version. The tool takes the
138 output of Kallisto and adds the data contained in the file to the VEP annotated VCF's INFO column.

139 Each of the variant annotated gets its expression value (TPM) added to the annotation information and
140 this is used to determine the level of variant expression during neoantigen filtering.

141 **Neoantigen prediction**

142 Human leukocyte antigen (HLA) class I alleles (HLA a, b and c) were predicted from each
143 patient's normal sample exome-seq data using HLA-HD v.1.2.1 (Kawaguchi et al., 2017). Here, the
144 putative HLA reads are aligned to an imputed library of full-length HLA alleles. Neoantigens were
145 then predicted using pVACseq (Hundal et al, 2016) with MHCflurry, MHCnuggetsI, SMM, and
146 SMMPMBEC algorithms and keeping the default parameters, except for turning off the VAF and
147 coverage filters. Here, the neopeptides that could bind to the patient-specific HLA alleles were
148 predicted from the Immune Epitope Database (IEDB; Vita et al., 2019). This involved matching patient
149 HLA type to the existing IEDB list keeping all amino acids with lengths for 9, 10 and 11-mers.
150 Predicted epitopes were filtered to retain only those with high affinity ($IC_{50} \leq 500nM$) and were
151 expressed (transcripts per million, $TPM > 1$) in tumor samples. The bioinformatic analysis workflow is
152 outlined in Figure 1.

153 Sample summary statistics and the pairwise tests for differences among mutations and neoan-
154 tigen abundance among the BC subtypes using Wilcoxon test and visualization of the results were
155 performed in R software (R Core Team, 2023).

156 **RESULTS**

157 **Patients and sample characteristics**

158 The demographic and clinical characteristics of the 23 breast cancer patients are summarized
159 in supplementary Table S1. We grouped the tumors into 3 subtypes based on expression of either the
160 hormone receptors (HR) or human epidermal growth factor receptor 2 (HER2) (Narang et al., 2019):
161 those that were HER2+ regardless of the HR status, those that were negative for all hormone receptors

162 (triple negative breast cancer; TNBC) and those that were HR+ and HER2-. Majority of the samples
163 were HR+/HER2- constituting 52.2%, followed by HER2+ at 34.8% and TNBC at 13.0%. Most of the
164 patients had invasive carcinoma (invasive ductal carcinoma, 78.26% and invasive carcinoma; 4.35%).
165 For tumor grade, 65.22% of the patients had grade 3 tumors (65.22%), while the rest had grade 2 tumors
166 (34.78%). Clinically, 39.13% of the patients were in stage II, 30.44% in stage III, and 8.7% in stage I
167 (Table S1).

168 **Mutation profiles for the 23 patients**

169 Across all genes, the average number of detected mutations in the 23 patients was 2809
170 mutations. Considering the different subtypes, TNBC had the highest average number of mutations at
171 3202, followed by HR+/HER2- at 2757, and HER2+ at 2740 mutations (Figure S1). From the catalog
172 of somatic mutations in cancer (COSMIC, v99), we identified 73 genes reported to be mutated in breast
173 cancer and among those, 62 (84.9%) had at least one mutation in our samples. The mutation frequency
174 among the 62 genes ranged from 1 to 55 mutations per individual. The majority of the mutations were
175 of the missense type, most of which were substitutions of C>T (Figure 2). The top 10 mutated genes
176 among the 62 are shown in Figure 3. Four genes (*MUC16*, *MUC17*, *TTN*, *RP1L1*) were altered in more
177 than 95% of the patients (Figure 3). Moreover, mutations in genes *TP53-ERBB3*, *PTEN-CFAP46* were
178 found to significantly co-occur, while *BRCA1-MUC17* mutations were significantly mutually exclusive
179 ($P<0.05$) (Figure 4). Furthermore, the majority of the single nucleotide mutations were substitution of
180 C to T, whereas T to A substitutions were most uncommon. Transitions occurred more frequently than
181 transversion in these substitutions (Figure 5).

182 **Neoantigen burden**

183 In an analysis that included all the genes (10260), an average of 1465 neoantigens had a
184 ≤ 500 nM median IC50 binding score and >1 TPM expression level in any of the 23 patients and their
185 presence significantly correlated with the somatic mutations ($R^2=0.570$, $P=0.001$) (Figure 6). Out of

186 the 62 COSMIC genes that were mutated in the tumor tissue, 58 genes produced at least one neoantigen.
187 After filtering for genes that produced at least two neoantigens, 44 genes had a mean of 10.5
188 neoantigens ranging from 2 to 93. A total of 477 putative neoantigens were identified in these 44 genes
189 across the 23 patients (Figure 7) predominantly derived from missense mutations (88%), indels (6%)
190 and frameshift mutations (6%) (Figure 8). Most of the neoantigens were produced in the TNBC subtype
191 with an average of 25 neoantigens, followed by HR+/HER2- at 20 neoantigens and HER2+ with an
192 average of 19 neoantigens (Figure S1). Notably, 78% of the putative breast cancer neoantigens were
193 patient-specific (Table S2). HLA-C*06:01 allele was associated with majority of neoantigens (194),
194 followed by HLA-A*30:01 (131), HLA-A*02:01 (103), and HLA-B*58:01 (49). Among the genes of
195 interest that produced putative neoantigens include *MUC17*, *TTN*, *MUC16*, *AKAP9*, *NEB*, *RP1L1*,
196 *CDH23*, *PCDHB10*, *BRCA2*, *TP53*, *TG*, *RBI* among others (Figure 7, Table S3).

197

198 DISCUSSION

199 We analyzed the mutational burden and predicted the neoantigen repertoire in 23 Kenyan breast
200 cancer patients using WES and RNA sequencing data. Among the different breast cancer subtypes, we
201 found that the TNBC molecular subtype had the highest mutational and neoantigen burden although
202 there was no significant difference among the subtypes (Figure S1, Table S4). This is consistent with
203 other studies (Narang et al., 2019). TNBC origin is not well understood although it is reported to be
204 heterogeneous in nature relying on different signaling pathways such as JAK/STAT,
205 PI3K/AKT/mTOR or NOTCH, cell cycle regulators (*TP53*) and genome integrity genes (*BRCA1/2*)
206 (Benvenuto et al., 2019). This makes it a disease that is difficult to manage because we do not have a
207 clear understanding of the molecular mechanisms driving it. Yet, the high mutational and neoantigen
208 burden combined with the patient specificity may provide an untapped opportunity to design and
209 optimize personalized immunotherapy for this subtype.

210 In contrast to most populations where *TP53*, *PIK3CA* and *GATA3* are the most mutated genes
211 (Pan et al., 2020; Pipek et al., 2023; Tang et al., 2023), in our study population, three genes *MUC16*,
212 *MUC17* and *TTN* were highly mutated in over 50% of the samples and produced the highest number
213 of neoantigens. *MUC16* has been reported to take part in breast cancer progression and metastasis when
214 overexpressed due to its influence on cell cycle and survival through the JAK2/STAT3 pathway
215 (Lakshmanan et al., 2012). It has been reported as one of the highly mutated genes in breast cancer
216 (Wang & Guda, 2016). *MUC16* has also been described as a marker for disease progression, recurrence,
217 and chemotherapy response (Felder et al., 2014). A high mutation frequency for *MUC17* and *TTN* have
218 recently been reported as an unexpected finding in a study of early onset breast cancer (EOBC) in
219 Taiwanese women (Midha et al., 2020). *MUC17* may influence chemoresistance and has recently been
220 reported as a driver gene in adult gliomas (Al Amri et al., 2020; Machado & Ferrer, 2023). For *TTN*,
221 Oh et al. (2020) found that mutations in *TTN* correlate with tumor mutational burden and high

222 microsatellite instability, which is associated with poor breast cancer prognosis. Thus, the role of
223 *MUC17* and *TTN* should further be investigated on how mutations in them may relate to early onset of
224 breast cancer in Kenyan patients (Tang et al., 2023).

225 We found that *TP53* gene mutations significantly co-occurred with *ERBB3* mutations and so
226 did mutations in *PTEN* and *CFAP46*, whereas *BRCA1* and *MUC17* mutations never co-occurred. *TP53*
227 mutations are associated with tumor aggression and are found in about half of HER2-amplified tumors
228 (Marvalim et al., 2023). The *TP53* mutations have been implicated in poor prognosis of HER2+
229 subtypes compared to other subtypes (Dumay et al., 2013). *PTEN* is a tumor suppressor gene, whose
230 mutation has been associated with initiation, progression, and metastasis of breast cancer (Chen et al.,
231 2022). On the other hand, although *CFAP46* role in breast cancer is not yet clear, gene fusion involving
232 various other genes such as *VTIIA* (reported to cause the initiation of glioma and other cancers) has
233 been reported to play a role in breast cancer (Tsuge et al., 2019).

234 Breast tumors with either germline or somatic *BRCA1* mutations show no difference in their
235 cancer biology, but inherited mutations in this gene confers a very high lifetime risk of developing
236 breast cancer (Milne & Antoniou, 2011; den Brok et al., 2017; Bodily et al., 2020). This could be the
237 reason such mutations do not necessarily need to co-occur with other gene mutations to initiate or
238 promote breast cancer progression. In our study, *BRCA1* was not among the highly mutated genes
239 considering all mutations but was among the genes with high number of missense mutations (Figure
240 4). In contrast, *MUC17* mutations were among the most prevalent. Given the role of *MUC17* mutations
241 in chemoresistance and in early onset breast cancer (Al Amri et al., 2020; Machado & Ferrer, 2023),
242 its high prevalence and exclusive occurrence in the Kenyan samples that are prone to early onset of
243 breast cancer should be investigated further.

244 Similar to most studies on neoantigen prediction in breast cancer, we have found that
245 neoantigen burden is positively correlated with tumor mutational burden and that neoantigens were
246 patient-specific (Narang et al., 2019; Animesh et al., 2022). Although most of the top 10 mutated genes
247 (80%) were also the top 10 in the number of neoantigens generated, genes like *TP53* and *PIK3CA* that
248 are reported to be highly mutated in most patient cohorts were not among the top 10 mutated genes in
249 this study, but generated among the highest number of neoantigens (Figure 6; Figure 7). *ARID1A* gene,
250 which showed unique mutational profile in Kenyan population using exome data compared to African
251 American and Asian population (Tang et al., 2023), was not among the highly mutated, but produced
252 neoantigens. We found that most neoantigens were derived predominantly from missense mutations
253 (88%), compared to indels and frameshift mutations (12%). This is consistent with other studies
254 although the majority do not predict neoantigens from indels and frameshift mutations (Morisaki et al.,
255 2021). Similar to other studies, the TNBC subtype had more neoantigens, compared to HR+/HER2-
256 and HER2+ subtypes (Narang et al., 2019; Morisaki et al., 2021).

257 In our small sample cohort, we have been able to identify putative neoantigens that show
258 patient-specificity and thus are important in tailored treatment. Interestingly, the mutations and
259 neoantigens in this population are predominantly derived from a unique set of genes (*MUC16*, *MUC17*,
260 *TNT*) compared to other populations, which provide an opportunity for validation in a much larger
261 sample cohort. We predicted neoantigens based on binding affinity to HLA class I only as it is the most
262 important class of antigen binding proteins in cancer immunity. However, HLA class II-based
263 neoantigens may also have a role in tumor immune response (Alspach et al., 2019). Moreover, we did
264 not investigate the expression of the predicted neoantigens on tumor cells alongside the MHC class I
265 molecules and their ability to activate T cells. This being a discovery study, validation of the findings
266 need to be done in a larger cohort while addressing the highlighted limitations of this study.

267 Taken together, our findings corroborate the neoantigen profile in breast cancer, highlighting
268 the patient specificity in Kenyan population breast cancer mutational and neoantigens signatures. We
269 also describe putative neoantigens that could be used as markers for breast cancer diagnosis, treatment
270 monitoring, and development of novel immunotherapy.

271 **Acknowledgment**

272 We would like to thank the patients for their consent to provide samples and Aga Khan University
273 Hospital (Nairobi) and AIC Kijabe Hospital (Kijabe) for granting access to patient samples.

274 **Funding**

275 This work was funded by the National Research Fund – Kenya that supported sample collection, and
276 by the Center for Cancer Research, National Cancer Institute, USA, that supported the sequencing
277 work.

278 **Data accessibility**

279 WES data is accessible at SRA database Accession number: PRJNA913947, while RNA-seq data is
280 accessible at the GEO database under Accession number: GSE225846. All other datasets for this study
281 are included in the article’s Supplementary Material.

282 **Authors’ contribution**

283 FM conceived the idea and designed the research project, collected the samples and assembled
284 experiment materials. GW, JG, KM and MM performed the data analysis. FM and GW wrote the
285 manuscript. ARH, SS, SA helped with drafting and reviewing the manuscript. All authors contributed
286 to the revision and final editing of the manuscript prior to submission.

287

288

289

290

291

292 References

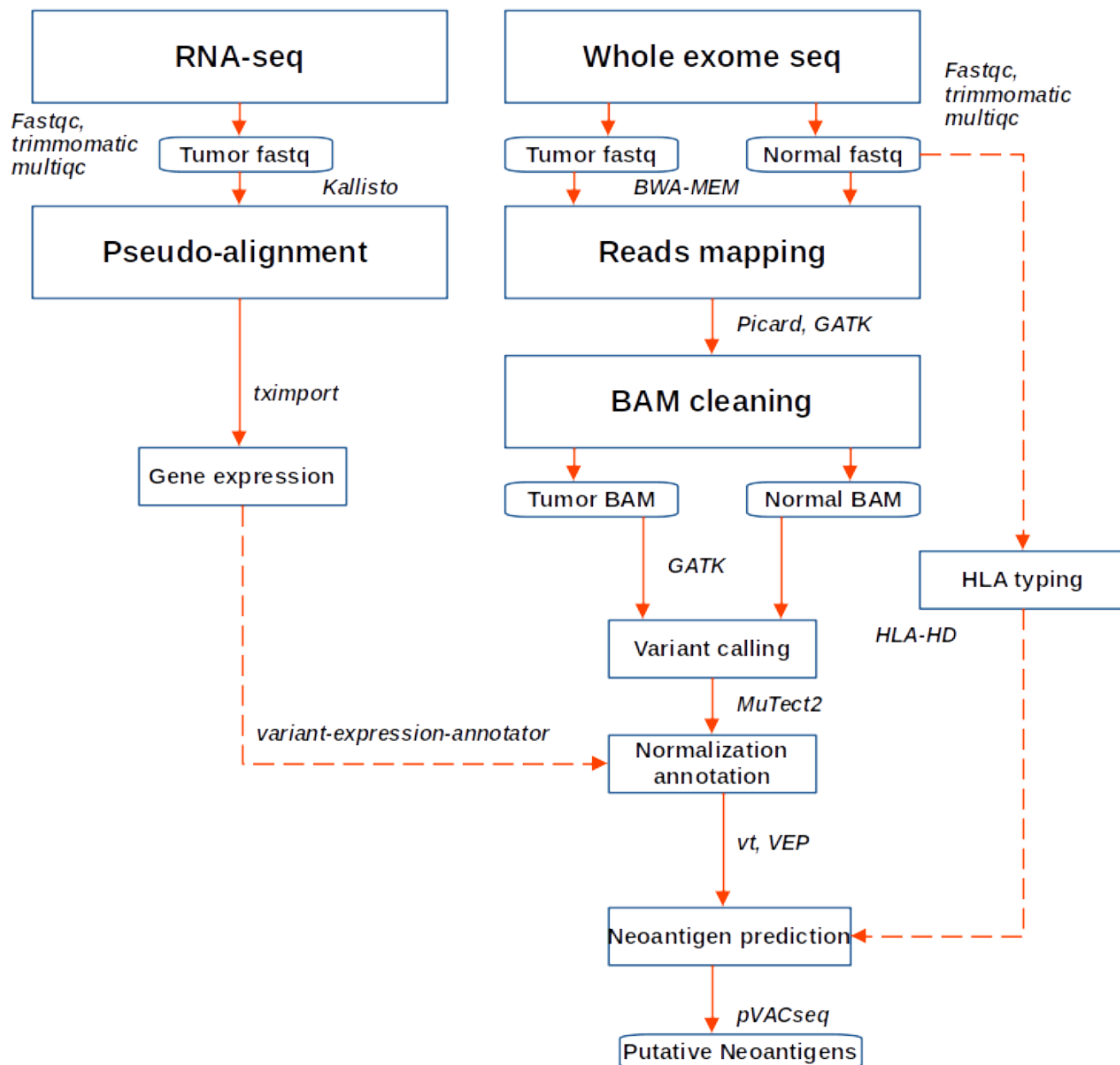
- 293 Andrews S. (2010). FastQC: a quality control tool for high throughput sequence data. Available
294 online at: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>
- 295 Al Amri, W. S., Allinson, L. M., Baxter, D. E., Bell, S. M., Hanby, A. M., Jones, S. J., Shaaban, A.
296 M., Stead, L. F., Verghese, E. T., & Hughes, T. A. (2020). Genomic and Expression Analyses Define
297 MUC17 and PCNX1 as Predictors of Chemotherapy Response in Breast Cancer. *Molecular Cancer*
298 *Therapeutics*, 19(3), 945–955. <https://doi.org/10.1158/1535-7163.MCT-19-0940>
- 299 Alspach, E., Lussier, D. M., Miceli, A. P., Kizhvatov, I., DuPage, M., Luoma, A. M., Meng, W.,
300 Lichti, C. F., Esaulova, E., Vomund, A. N., Runci, D., Ward, J. P., Gubin, M. M., Medrano, R. F. V.,
301 Arthur, C. D., White, J. M., Sheehan, K. C. F., Chen, A., Wucherpennig, K. W., ... Schreiber, R. D.
302 (2019). MHC-II neoantigens shape tumour immunity and response to immunotherapy. *Nature*,
303 574(7780), 696–701. <https://doi.org/10.1038/s41586-019-1671-8>
- 304 Animesh, S., Ren, X., An, O., Chen, K., Lee, S. C., Yang, H., & Fullwood, M. J. (2022). *Exploring*
305 *the Neoantigen burden in Breast Carcinoma Patients*. <https://doi.org/10.1101/2022.03.03.482669>
- 306 Bates, J. P., Derakhshandeh, R., Jones, L., & Webb, T. J. (2018). Mechanisms of immune evasion in
307 breast cancer. *BMC cancer*, 18(1), 556. <https://doi.org/10.1186/s12885-018-4441-3>
- 308 Benvenuto, M., Focaccetti, C., Izzi, V., Masuelli, L., Modesti, A., & Bei, R. (2021). Tumor antigens
309 heterogeneity and immune response-targeting neoantigens in breast cancer. *Seminars in cancer*
310 *biology*, 72, 65–75. <https://doi.org/10.1016/j.semcancer.2019.10.023>
- 311 Benvenuto, M., Focaccetti, C., Izzi, V., Masuelli, L., Modesti, A., & Bei, R. (2021). Tumor antigens
312 heterogeneity and immune response-targeting neoantigens in breast cancer. *Seminars in Cancer*
313 *Biology*, 72, 65–75. <https://doi.org/10.1016/j.semcancer.2019.10.023>
- 314 Blass, E., & Ott, P. A. (2021). Advances in the development of personalized neoantigen-based
315 therapeutic cancer vaccines. *Nature Reviews Clinical Oncology*, 18(4), 215–229.
316 <https://doi.org/10.1038/s41571-020-00460-2>
- 317 Bodily W.R., Shirts B.H., Walsh T., Gulsuner S., King M.-C., Parker A., Roosan M., Piccolo S.R.
318 Effects of Germline and Somatic Events in Candidate BRCA-like Genes on Breast-Tumor
319 Signatures. *PLoS ONE*. 2020;15:e0239197. doi: 10.1371/journal.pone.0239197
- 320 Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: A flexible trimmer for Illumina
321 sequence data. *Bioinformatics*, 30(15), 2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>
- 322 Brahmer, J., Reckamp, K. L., Baas, P., Crinò, L., Eberhardt, W. E. E., Poddubskaya, E., Antonia, S.,
323 Pluzanski, A., Vokes, E. E., Holgado, E., Waterhouse, D., Ready, N., Gainor, J., Arén Frontera, O.,
324 Havel, L., Steins, M., Garassino, M. C., Aerts, J. G., Domine, M., ... Spigel, D. R. (2015).
325 Nivolumab versus Docetaxel in Advanced Squamous-Cell Non–Small-Cell Lung Cancer. *New*
326 *England Journal of Medicine*, 373(2), 123–135. <https://doi.org/10.1056/NEJMoa1504627>

- 327 Bray, N. L., Pimentel, H., Melsted, P., & Pachter, L. (2016). Near-optimal probabilistic RNA-seq
328 quantification. *Nature Biotechnology*, 34(5), 525–527. <https://doi.org/10.1038/nbt.3519>
- 329 Chen, J., Sun, J., Wang, Q., Du, Y., Cheng, J., Yi, J., Xie, B., Jin, S., Chen, G., Wang, L., Wang, X.,
330 & Wei, H. (2022). Systemic Deficiency of PTEN Accelerates Breast Cancer Growth and Metastasis.
331 *Frontiers in Oncology*, 12, 825484. <https://doi.org/10.3389/fonc.2022.825484>
- 332 den Brok W.D., Schrader K.A., Sun S., Tinker A.V., Zhao E.Y., Aparicio S., Gelmon K.A.
333 Homologous Recombination Deficiency in Breast Cancer: A Clinical Review. *JCO Precis. Oncol.*
334 2017;1:1–13. doi: 10.1200/PO.16.00031.
- 335 Dumay, A., Feugeas, J., Wittmer, E., Lehmann-Che, J., Bertheau, P., Espié, M., Plassa, L., Cottu, P.,
336 Marty, M., André, F., Sotiriou, C., Pusztai, L., & De Thé, H. (2013). Distinct *tumor protein p53*
337 mutants in breast cancer subgroups. *International Journal of Cancer*, 132(5), 1227–1231.
338 <https://doi.org/10.1002/ijc.27767>
- 339 Ewels, P., Magnusson, M., Lundin, S., & Källér, M. (2016). MultiQC: Summarize analysis results for
340 multiple tools and samples in a single report. *Bioinformatics*, 32(19), 3047–3048.
341 <https://doi.org/10.1093/bioinformatics/btw354>
- 342 Felder, M., Kapur, A., Gonzalez-Bosquet, J., Horibata, S., Heintz, J., Albrecht, R., Fass, L., Kaur, J.,
343 Hu, K., Shojaei, H., Whelan, R. J., & Patankar, M. S. (2014). MUC16 (CA125): Tumor biomarker to
344 cancer therapy, a work in progress. *Molecular Cancer*, 13(1), 129. <https://doi.org/10.1186/1476-4598-13-129>
- 346 Hundal, J., Carreno, B. M., Petti, A. A., Linette, G. P., Griffith, O. L., Mardis, E. R., & Griffith, M.
347 (2016). pVAC-Seq: A genome-guided in silico approach to identifying tumor neoantigens. *Genome*
348 *Medicine*, 8(1), 11. <https://doi.org/10.1186/s13073-016-0264-5>
- 349 Kawaguchi, S., Higasa, K., Shimizu, M., Yamada, R., & Matsuda, F. (2017). HLA-HD: An accurate
350 HLA typing algorithm for next-generation sequencing data. *Human Mutation*, 38(7), 788–797.
351 <https://doi.org/10.1002/humu.23230>
- 352 Kandoth, C., Gao, J., Qwangmsk, Mattioni, M., Struck, A., Boursin, Y., Penson, A., & Chavan, S.
353 (2018). *mskcc/vcf2maf: Vcf2maf v1.6.16* (v1.6.16) [Computer software].
354 <https://doi.org/10.5281/ZENODO.593251>
- 355 Lakshmanan, I., Ponnusamy, M. P., Das, S., Chakraborty, S., Haridas, D., Mukhopadhyay, P., Lele,
356 S. M., & Batra, S. K. (2012). MUC16 induced rapid G2/M transition via interactions with JAK2 for
357 increased proliferation and anti-apoptosis in breast cancer cells. *Oncogene*, 31(7), 805–817.
358 <https://doi.org/10.1038/onc.2011.297>
- 359 Li, H. (2013). *Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.*
360 <https://doi.org/10.48550/ARXIV.1303.3997>
- 361 Machado, G. C., & Ferrer, V. P. (2023). MUC17 mutations and methylation are associated with poor
362 prognosis in adult-type diffuse glioma patients. *Journal of the Neurological Sciences*, 452, 120762.
363 <https://doi.org/10.1016/j.jns.2023.120762>

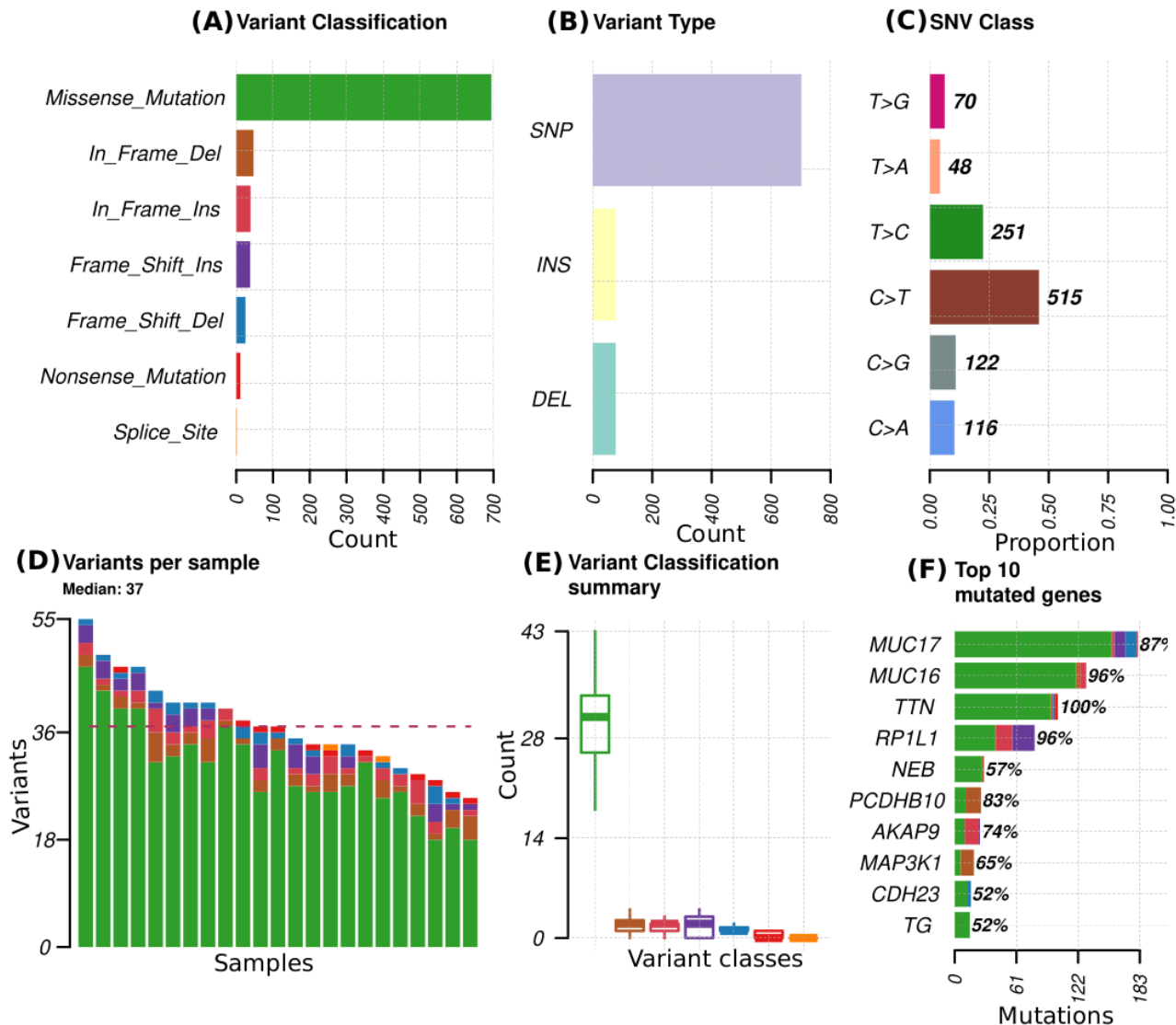
- 364 Marvalim, C., Datta, A., & Lee, S. C. (2023). Role of p53 in breast cancer progression: An insight
365 into p53 targeted therapy. *Theranostics*, 13(4), 1421–1442. <https://doi.org/10.7150/thno.81847>
- 366 Mayakonda, A., Lin, D.-C., Assenov, Y., Plass, C., & Koeffler, H. P. (2018). Maftools: Efficient and
367 comprehensive analysis of somatic variants in cancer. *Genome Research*, 28(11), 1747–1756.
368 <https://doi.org/10.1101/gr.239244.118>
- 369 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K.,
370 Altshuler, D., Gabriel, S., Daly, M., & DePristo, M. A. (2010). The Genome Analysis Toolkit: A
371 MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research*,
372 20(9), 1297–1303. <https://doi.org/10.1101/gr.107524.110>
- 373 McLaren, W., Gil, L., Hunt, S. E., Riat, H. S., Ritchie, G. R. S., Thormann, A., Flicek, P., &
374 Cunningham, F. (2016). The Ensembl Variant Effect Predictor. *Genome Biology*, 17(1), 122.
375 <https://doi.org/10.1186/s13059-016-0974-4>
- 376 Midha, M. K., Huang, Y.-F., Yang, H.-H., Fan, T.-C., Chang, N.-C., Chen, T.-H., Wang, Y.-T., Kuo,
377 W.-H., Chang, K.-J., Shen, C.-Y., Yu, A. L., Chiu, K.-P., & Chen, C.-J. (2020). Comprehensive
378 Cohort Analysis of Mutational Spectrum in Early Onset Breast Cancer Patients. *Cancers*, 12(8),
379 2089. <https://doi.org/10.3390/cancers12082089>
- 380 Milne, R. L. & Antoniou, A. C. (2011). Genetic modifiers of cancer risk for BRCA1 and BRCA2
381 mutation carriers. *Annals of Oncology*, 22 (1): i11-i17.
- 382 Morisaki, T., Kubo, M., Umebayashi, M., Yew, P. Y., Yoshimura, S., Park, J.-H., Kiyotani, K., Kai,
383 M., Yamada, M., Oda, Y., Nakamura, Y., Morisaki, T., & Nakamura, M. (2021). Neoantigens elicit T
384 cell responses in breast cancer. *Scientific Reports*, 11(1), 13590. <https://doi.org/10.1038/s41598-021-91358-1>
- 386 Narang, P., Chen, M., Sharma, A. A., Anderson, K. S., & Wilson, M. A. (2019). The neoepitope
387 landscape of breast cancer: Implications for immunotherapy. *BMC Cancer*, 19(1), 200.
388 <https://doi.org/10.1186/s12885-019-5402-1>
- 389 Oh, J.-H., Jang, S. J., Kim, J., Sohn, I., Lee, J.-Y., Cho, E. J., Chun, S.-M., & Sung, C. O. (2020).
390 Spontaneous mutations in the single TTN gene represent high tumor mutation burden. *Npj Genomic
391 Medicine*, 5(1), 33. <https://doi.org/10.1038/s41525-019-0107-6>
- 392 Pan, J.-W., Zabidi, M. M. A., Ng, P.-S., Meng, M.-Y., Hasan, S. N., Sandey, B., Sammut, S.-J., Yip,
393 C.-H., Rajadurai, P., Rueda, O. M., Caldas, C., Chin, S.-F., & Teo, S.-H. (2020). The molecular
394 landscape of Asian breast cancers reveals clinically relevant population-specific differences. *Nature
395 Communications*, 11(1), 6433. <https://doi.org/10.1038/s41467-020-20173-5>
- 396 Pan, R.-Y., Chung, W.-H., Chu, M.-T., Chen, S.-J., Chen, H.-C., Zheng, L., & Hung, S.-I. (2018).
397 Recent Development and Clinical Application of Cancer Vaccine: Targeting Neoantigens. *Journal of
398 Immunology Research*, 2018, 1–9. <https://doi.org/10.1155/2018/4325874>
- 399 Pipek, O., Alpár, D., Ruzs, O., Bődör, C., Udvarnoki, Z., Medgyes-Horváth, A., Csabai, I., Szállási,
400 Z., Madaras, L., Kahán, Z., Cserni, G., Kővári, B., Kulka, J., & Tőkés, A. M. (2023). Genomic

- 401 Landscape of Normal and Breast Cancer Tissues in a Hungarian Pilot Cohort. *International Journal*
402 *of Molecular Sciences*, 24(10), 8553. <https://doi.org/10.3390/ijms24108553>
- 403 R Core Team (2023). R: A language and environment for statistical computing. R Foundation for
404 Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>
- 405 Shiravand, Y., Khodadadi, F., Kashani, S. M. A., Hosseini-Fard, S. R., Hosseini, S., Sadeghirad, H.,
406 Ladwa, R., O'Byrne, K., & Kulasinghe, A. (2022). Immune Checkpoint Inhibitors in Cancer
407 Therapy. *Current Oncology*, 29(5), 3044–3060. <https://doi.org/10.3390/curroncol29050247>
- 408 Sonesson, C., Love, M. I., & Robinson, M. D. (2016). Differential analyses for RNA-seq: Transcript-
409 level estimates improve gene-level inferences. *F1000Research*, 4, 1521.
410 <https://doi.org/10.12688/f1000research.7563.2>
- 411 Tan, A., Abecasis, G. R., & Kang, H. M. (2015). Unified representation of genetic variants.
412 *Bioinformatics*, 31(13), 2202–2204. <https://doi.org/10.1093/bioinformatics/btv112>
- 413 Tang, W., Zhang, F., Byun, J. S., Dorsey, T. H., Yfantis, H. G., Ajao, A., Liu, H., Pichardo, M. S.,
414 Pichardo, C. M., Harris, A. R., Yang, X. R., Figueroa, J. D., Sayed, S., Makokha, F. W., & Ambs, S.
415 (2023). Population-specific Mutation Patterns in Breast Tumors from African American, European
416 American, and Kenyan Patients. *Cancer Research Communications*, 3(11), 2244–2255.
417 <https://doi.org/10.1158/2767-9764.CRC-23-0165>
- 418 Touchaei, Z. A., & Vahidi, S. (2024). MicroRNAs as regulators of immune checkpoints in cancer
419 immunotherapy: Targeting PD-1/PD-L1 and CTLA-4 pathways. *Cancer Cell International*, 24(1),
420 102. <https://doi.org/10.1186/s12935-024-03293-6>
- 421 Tsuge, S., Saberi, B., Cheng, Y., Wang, Z., Kim, A., Luu, H., Abraham, J. M., Ybanez, M. D.,
422 Hamilton, J. P., Selaru, F. M., Villacorta-Martin, C., Schlesinger, F., Philosophe, B., Cameron, A.
423 M., Zhu, Q., Anders, R., Gurakar, A., & Meltzer, S. J. (2019). Detection of Novel Fusion Transcript
424 VTI1A-CFAP46 in Hepatocellular Carcinoma. *Gastrointestinal Tumors*, 6(1–2), 11–27.
425 <https://doi.org/10.1159/000496795>
- 426 Valilou, S. F., & Rezaei, N. (2019). Tumor Antigens. In *Vaccines for Cancer Immunotherapy* (pp.
427 61–74). Elsevier. <https://doi.org/10.1016/B978-0-12-814039-0.00004-7>
- 428 Vita, R., Mahajan, S., Overton, J. A., Dhanda, S. K., Martini, S., Cantrell, J. R., Wheeler, D. K.,
429 Sette, A., & Peters, B. (2019). The Immune Epitope Database (IEDB): 2018 update. *Nucleic Acids*
430 *Research*, 47(D1), D339–D343. <https://doi.org/10.1093/nar/gky1006>
- 431 Wang, X., & Guda, C. (2016). Integrative exploration of genomic profiles for triple negative breast
432 cancer identifies potential drug targets. *Medicine*, 95(30), e4321.
433 <https://doi.org/10.1097/MD.0000000000004321>
- 434 Yarchoan, M., Johnson, B. A., 3rd, Lutz, E. R., Laheru, D. A., & Jaffee, E. M. (2017). Targeting
435 neoantigens to augment antitumour immunity. *Nature reviews. Cancer*, 17(4), 209–222.
436 <https://doi.org/10.1038/nrc.2016.154>

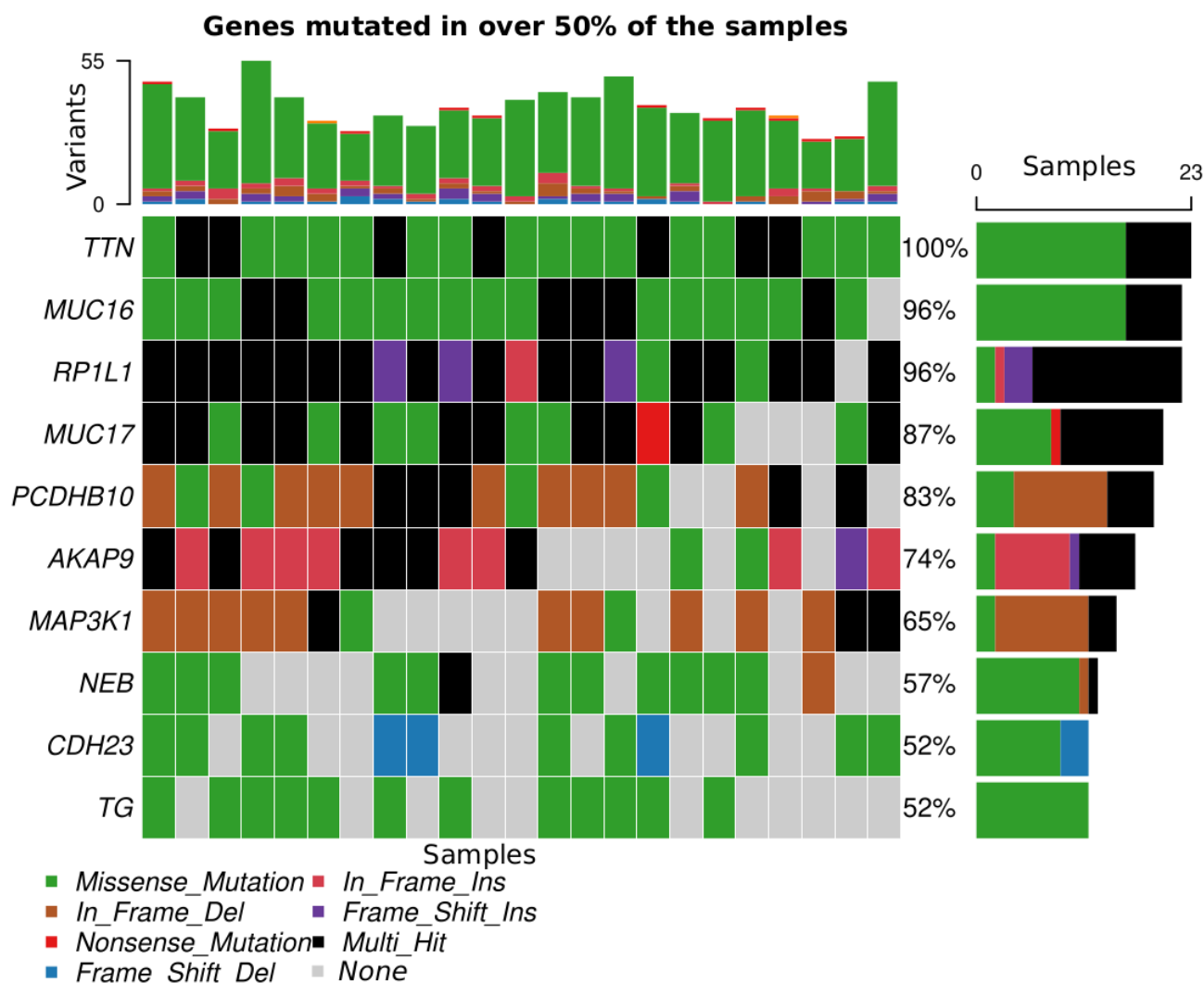
437 **Figures**
438



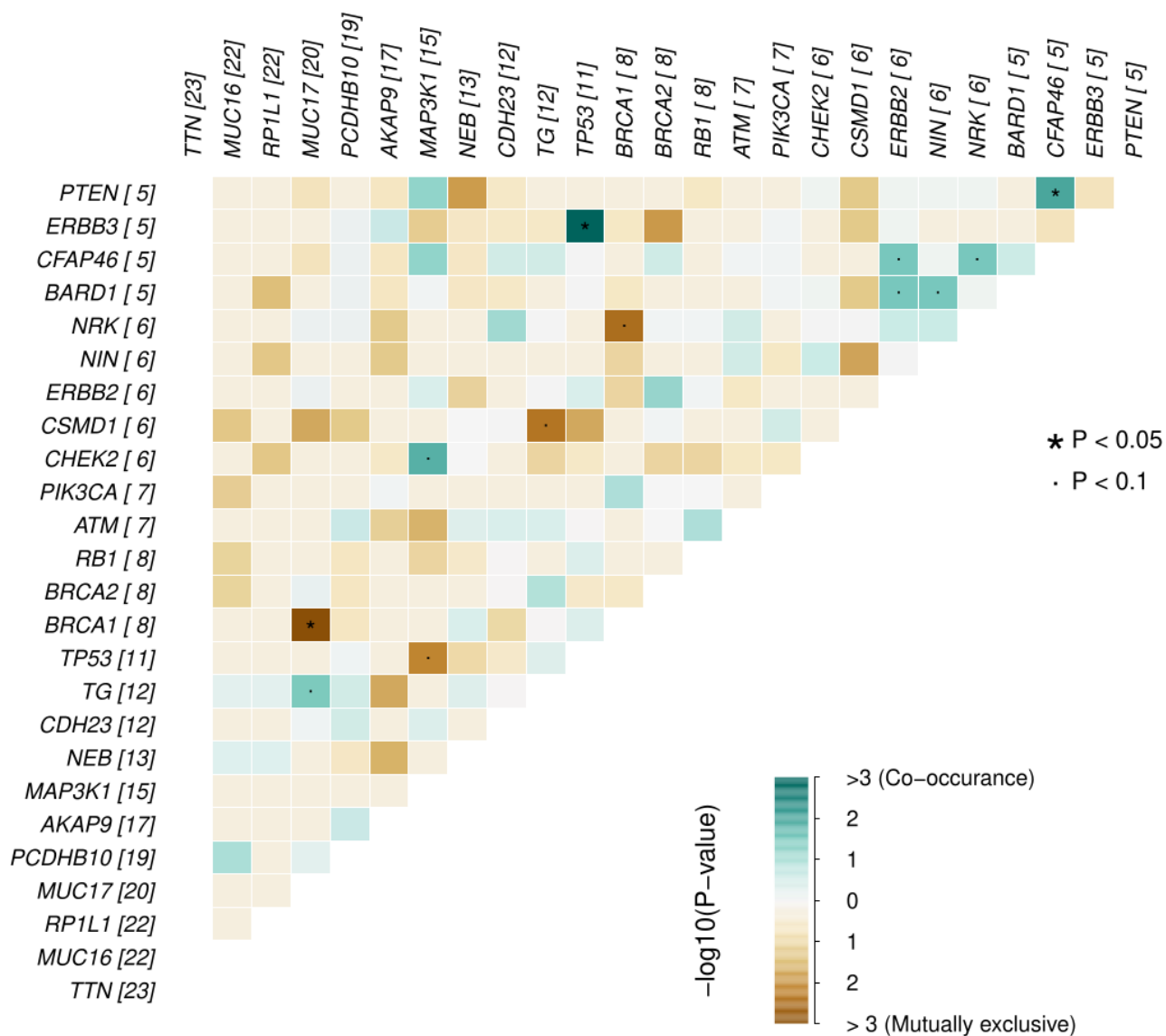
439 **Figure 1:** Workflow for neoantigen prediction from WES and RNA sequencing data. Fastq files
440 were quality checked, trimmed and aligned to the hg38 genome. Variant calling was performed
441 following GATK best practice, while gene expression was quantified using Kallisto. Variants were
442 annotated and expression data added, after which neoantigen prediction was performed in PVACseq
443 pipeline
444
445



446
 447 **Figure 2:** Mutational profiles in 23 patients for 73 genes reported to be mutated in breast cancer. **A)**
 448 variant classes abundance in the total mutations, **B)** variant types that include single nucleotide
 449 polymorphism (SNP), insertions (INS) and deletions (DEL), **C)** proportion of different single
 450 nucleotide variant (SNV), **D)** distribution of variants per sample with colors representing the different
 451 variant classes denoted in A, **E)** summary of the variant classes distribution and numbers in all
 452 samples, **F)** Top 10 mutated genes, with colors representing different variant classes and the
 453 percentages indicating the proportion of samples in which the genes mutations are present.
 454



455
 456 **Figure 3:** Top 10 genes mutated in >50% of the samples. Each color corresponds to a variant class
 457 listed at the bottom of the figure apart from gray, which indicates absence of mutation.
 458



459

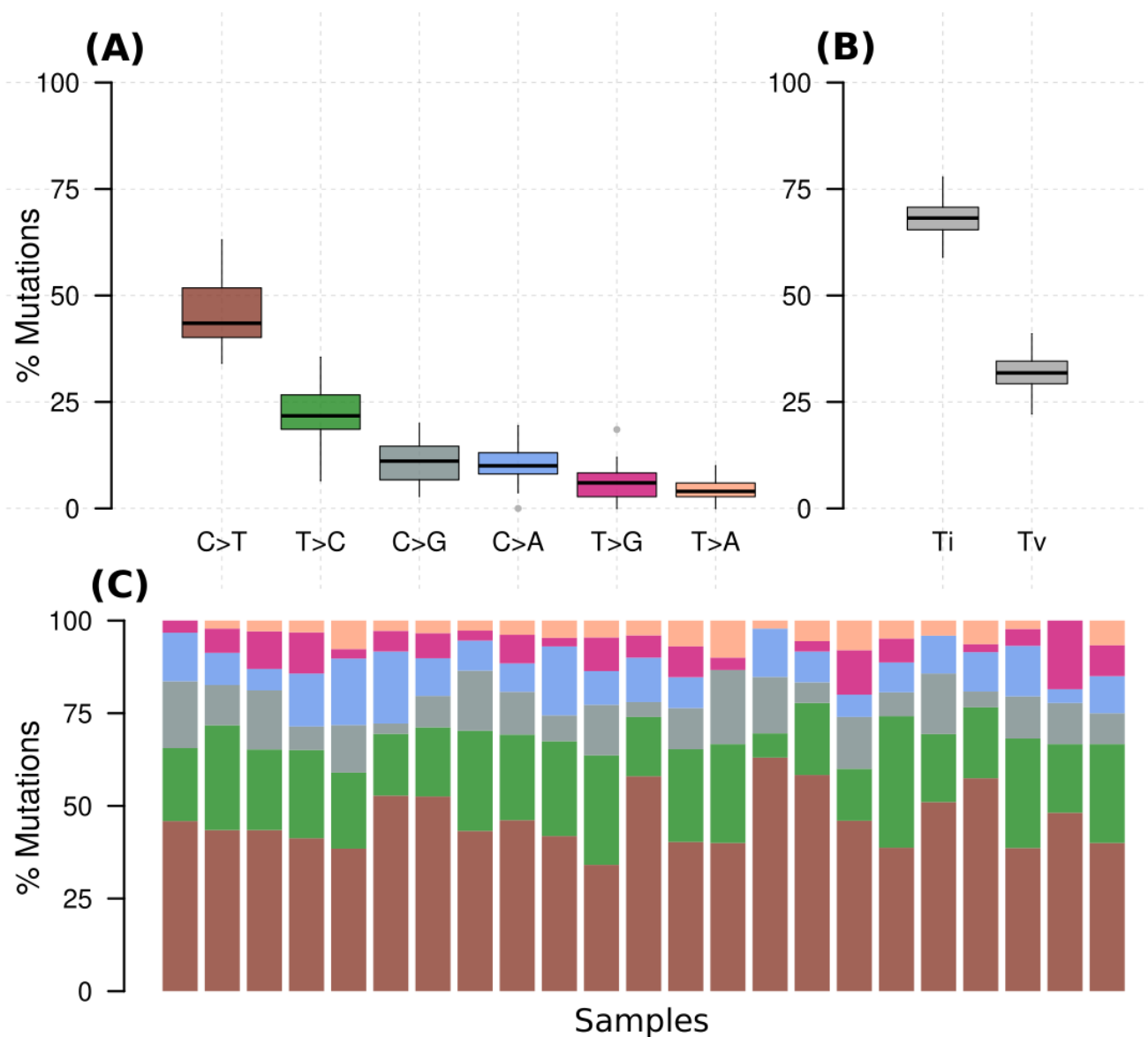
460

461

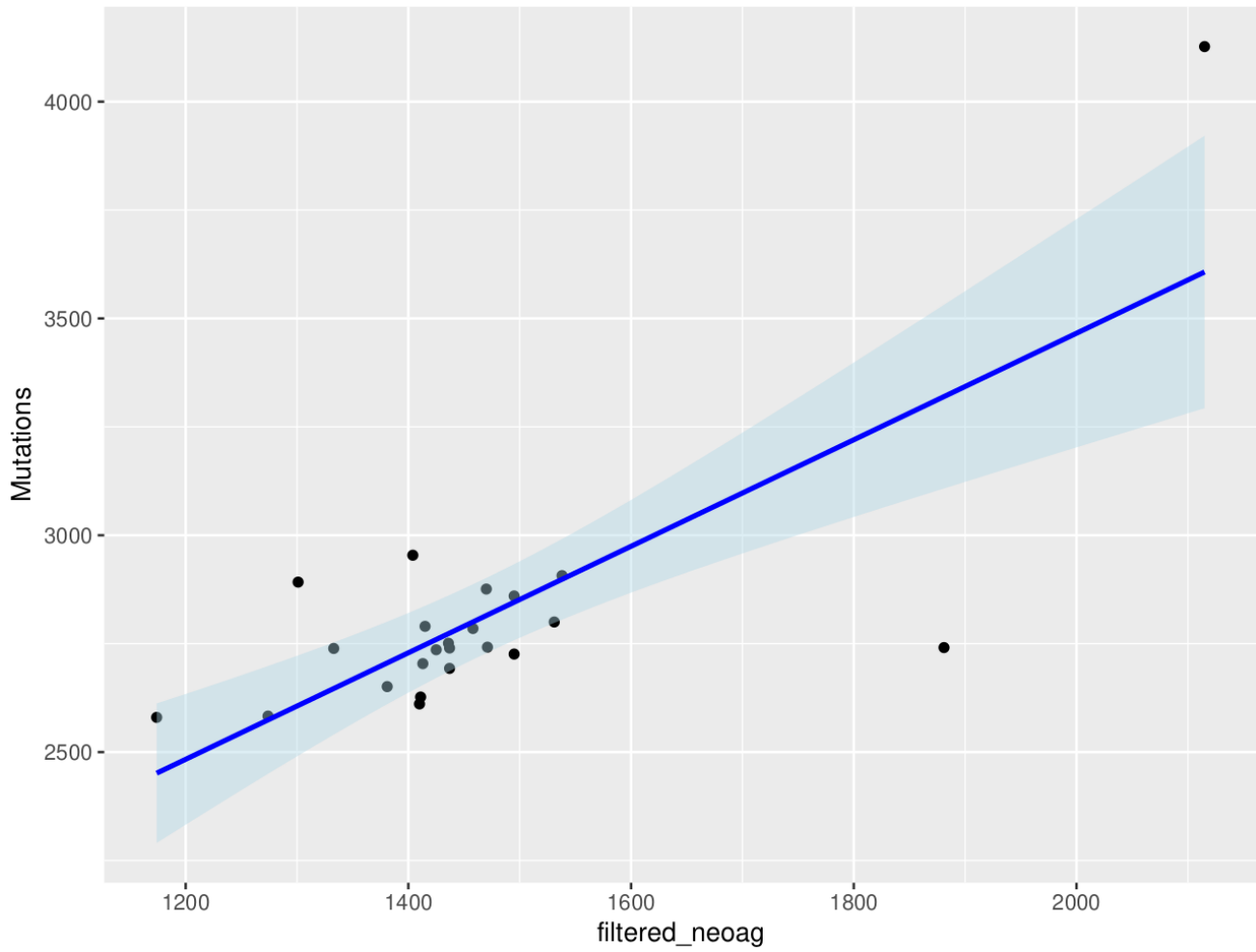
462

463

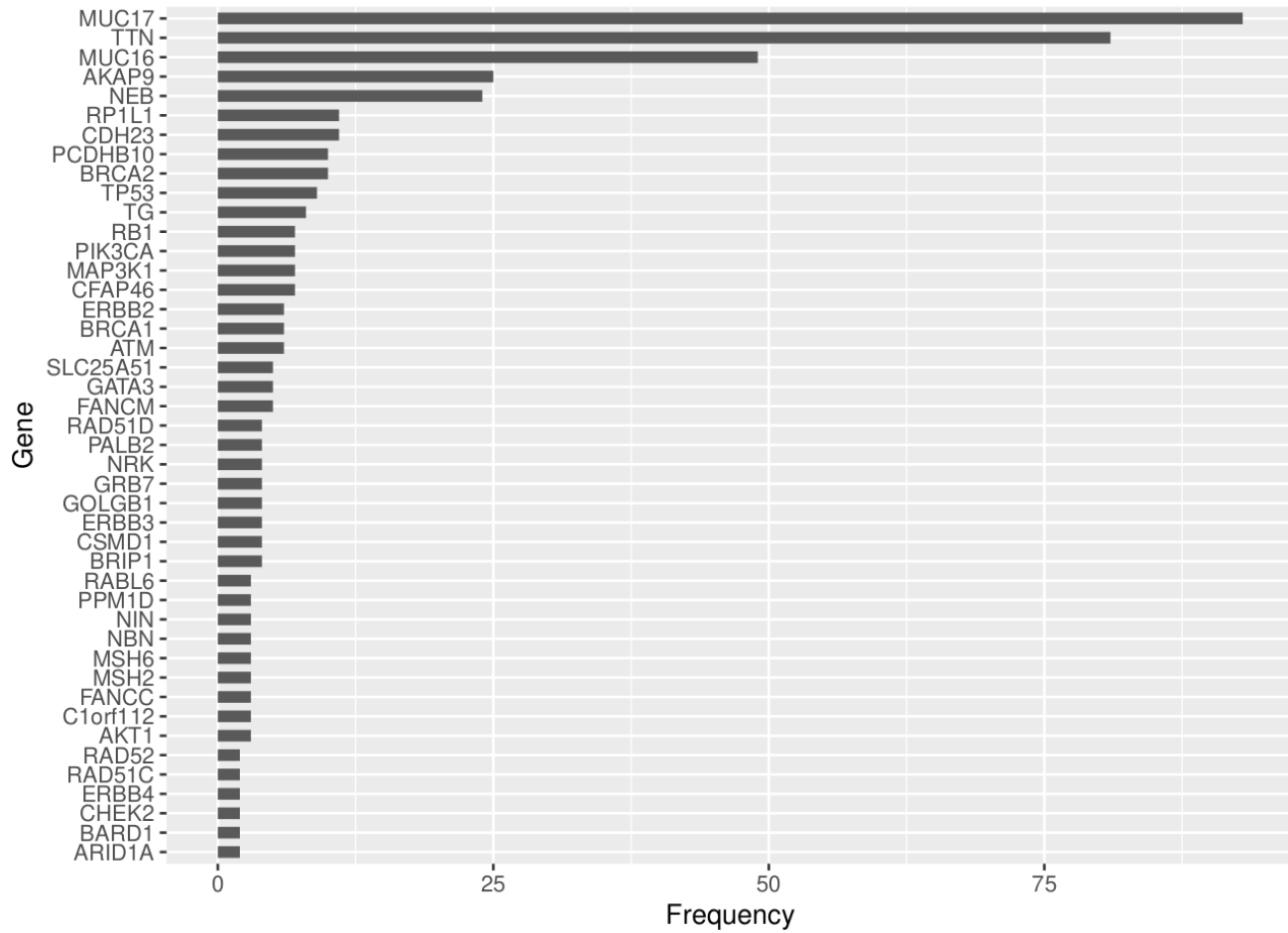
Figure 4: Probability of mutations in any two genes co-occurrence or being mutually exclusive in the breast cancer genes for the 23 Kenyan patients. The numbers in parenthesis alongside each gene represents the number of missense mutations for that gene in the samples.



464
465 **Figure 5:** A) Percentage of various substitution types in all samples, B) percentage of transversions
466 (interchange of purines for pyrimidine) and transition (interchange of either purines or pyrimidines)
467 for all samples, C) percentage of the substitutions in each of the samples with colors denoting the
468 various types in A.
469



470
471 **Figure 6:** Correlation between tumor mutational burden and neoantigen burden for all the genes in
472 the 23 patients. The neoantigens are filtered for high affinity ($IC_{50} \leq 500nM$) and expression
473 (transcripts per million, $TPM > 1$) in tumor samples.
474

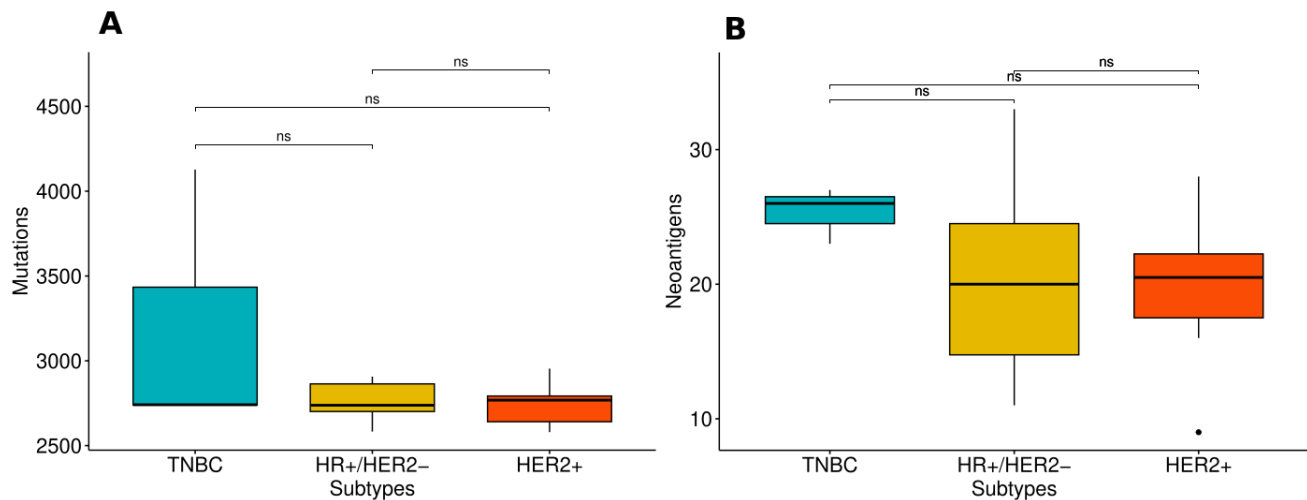


475
476
477
478

Figure 7: Frequency of neoantigens derived from the COSMIC genes that were mutated in the tumor tissue and produced >1 neoantigens for the 23 patients.

499 **Supplementary Materials**

500



501

502 **Figure S1:** Statistical pairwise test (Wilcoxon's test) for differences in mutational burden (A) and
503 neoantigen counts (B) for the 23 samples.

504

505 **Table S1:** Sample characteristics of the 23 Kenyan patients used in this study.

506

507 **Table S2:** Putative neoantigens for each of the 23 Kenyan patients. Cells in red indicate that the
508 neoantigen is shared by at least 2 patients.

509

510 **Table S3:** Summary of the roles of the top ten genes that generated a high number of neoantigens.

511

512 **Table S4:** Summary of total mutations and proportion of mutation types, total neoantigens, filtered
513 total neoantigens (filtered for high affinity ($IC_{50} \leq 500nM$) and expression [transcripts per million,
514 $TPM > 1$] in tumor samples) and filtered putative neoantigens from COSMIC 44 genes mutated in
515 tumor tissue, and proportion of mutation types that generated them per breast cancer subtype for the
516 23 samples