

1 Profile of the somatic mutational landscape in breast tumors from Hispanic/Latina women

2 Yuan C Ding^{1*}, Hanbing Song^{2*}, Aaron Adamson^{1*}, Daniel Schmolze³, Donglei Hu⁴, Scott

3 Huntsman⁴, Linda Steele¹, Carmina Patrick¹, Natalie Hernandez⁵, Charleen D Adams⁶, Laura

4 Fejerman⁷, Kevin L. Gardner⁸, Anna María Nápoles⁹, Eliseo J. Pérez-Stable¹⁰, Jeffrey N.

5 Weitzel¹¹, Henrik. Bengtsson^{12, 13}, Franklin W. Huang^{2,13,14,15,16,17}, Susan L. Neuhausen¹⁺, Elad

6 Ziv^{4,13,15+}

7 ¹Department of Population Sciences, Beckman Research Institute of City of Hope, Duarte, CA

8 91010

9 ² Division of Hematology/Oncology, Department of Medicine, University of California, San

10 Francisco, San Francisco, CA, USA.

11 ³Department of Pathology, City of Hope Medical Center, Duarte, CA 91010

12 ⁴Division of General Internal Medicine, Department of Medicine, University of California, San

13 Francisco, San Francisco, CA, USA.

14 ⁵ Western University of Health Sciences College of Graduate Nursing, Pomona, CA

15

16 ⁶Harvard T.H. Chan School of Public Health

17

18 ⁷Department of Public Health Sciences and Comprehensive Cancer Center, University of

19 California Davis, Davis, CA 95616

20 ⁸Department of Pathology and Cell Biology, Columbia University Irvine Medical Center

21 ⁹Division of Intramural Research, National Institute on Minority and Health Disparities, National

22 Institutes of Health, Bethesda, MD

23 ¹⁰National Institute on Minority and Health Disparities, National Institutes of Health, Bethesda,
24 MD

25 ¹¹Latin American School of Oncology, Sierra Madre, CA. 91024

26 ¹²Department of Epidemiology and Biostatistics, University of California, San Francisco, San
27 Francisco, CA

28 ¹³Helen Diller Family Comprehensive Cancer Center, University of California, San Francisco,
29 San Francisco, California 94143

30 ¹⁴Bakar Computational Health Sciences Institute, University of California, San Francisco,
31 San Francisco, CA, USA.

32 ¹⁵Institute for Human Genetics, University of California, San Francisco, San Francisco, CA,
33 USA.

34 ¹⁶Chan Zuckerberg Biohub, San Francisco, CA, USA.

35 ¹⁷Department of Medicine, San Francisco Veterans Affairs Medical Center, San Francisco, CA,
36 USA.

37 *All authors contributed equally

38 +Co-contributing authors

39

40 **ABSTRACT**

41

42 Breast cancer causes the most cancer deaths among Hispanic/Latinas (H/L). However, limited
43 tumor-sequencing data from H/L are available to guide treatment. To address this gap, we
44 performed whole-exome sequencing of DNA from 140 HL germline and 146 matched breast
45 tumors and RNA-seq for the tumors. We generated somatic-mutation profiles, identified copy-
46 number alterations (CNAs), and compared results to non-Hispanic White (White) women in The
47 Cancer Genome Atlas. Similar to Whites, *PIK3CA* and *TP53* were the most commonly mutated
48 genes in breast tumors from H/L. We found 4 common COSMIC mutation signatures (1, 2, 3,
49 13) and signature 16 not previously reported in other breast-cancer datasets. We observed
50 recurrent amplifications in breast-cancer drivers including *MYC*, *FGFR1*, *CCND1*, and *ERBB2*,
51 and a recurrent amplification on 17q11.2 associated with high *KIAA0100* gene expression,
52 implicated in breast-cancer aggressiveness. Expanded research is required to determine how
53 these characteristics of H/L tumors impact treatment response and survival.

54

55 Key words: Hispanic/Latino (H/L); Non-Hispanic Whites (Whites); copy-number alterations;
56 breast cancer; somatic mutations; expression outlier; disparities

57

58

59 INTRODUCTION

60 Sequencing studies of breast cancer have identified recurrently mutated genes and somatic
61 copy number alterations (SCNAs) affecting tumor suppressors and oncogenes¹⁻³. Both somatic
62 mutations and CNAs may be useful in determining prognosis. Currently, therapies for breast
63 cancer can be selected based on particular somatic mutations (i.e., alpelisib for *PIK3CA*⁴),
64 SCNAs (i.e., Trastuzumab for *HER2*), and germline mutations in genes in the homologous
65 recombination repair (HRR) pathway (polyADPribose polymerase inhibitors - PARPi's).

66 Genetic ancestry is associated with specific somatic mutations in many cancer types. *EGFR*
67 mutations are approximately four-fold more common in lung cancer from women and men of
68 East-Asian ancestry compared with lung cancer from women and men of other populations⁵ with
69 self-reported Hispanic/Latinos (H/L) representing an intermediate group^{6,7}. *FOXA1* mutations in
70 prostate cancer also are substantially more common in East-Asian ancestry populations
71 compared to European and African ancestry populations⁸. Comprehensive analyses of The
72 Cancer Genome Atlas (TCGA) have demonstrated that many mutations and CNAs are more
73 common in specific ancestral populations^{9,10}. In breast cancer, previous studies have
74 demonstrated that women of African ancestry have higher rates of *TP53* mutations and lower
75 rates of *PIK3CA* mutations, likely related to a higher incidence of a basal-like breast-cancer
76 subtype in African-American women^{11,12}. However, the genomic landscape of breast cancer has
77 not been well-characterized in H/L groups.

78 H/L represent the largest minority population in the US and have diverse origins, with the largest
79 subpopulations including Mexican Americans and Puerto Ricans. Genetically, H/L are a
80 population of mixed European, Indigenous-American (IA) and African ancestries with those
81 ancestry proportions varying widely depending on country of origin and regions within a country.
82 Although breast cancer is less common overall among H/L compared to self-reported non-

83 Hispanic White (White) women due to both environmental¹³ and genetic factors¹⁴, there is a
84 higher proportion of breast cancers diagnosed under age 50 years than in Whites¹⁵. Moreover,
85 outcomes are usually worse among H/L compared to White women¹⁶. In some studies, IA
86 ancestry was associated with poorer outcomes among H/L with breast cancer¹⁷. Human
87 epidermal growth factor (HER2) amplifications are over-represented among H/L and are more
88 common among H/L with more IA ancestry compared to those with more European ancestry¹⁸.
89 Few studies have investigated the distribution of somatic mutations and SCNAs in breast
90 tumors from H/L. In TCGA, out of 1,096 breast-cancer cases, only 39 are self-reported H/L. A
91 recent study analyzed data including whole-exome sequencing (WES) and gene expression
92 data from 109 Mexican women living in Mexico¹⁹. However, no similar-size study has been
93 conducted in H/L in the United States (US). To investigate the somatic-mutational spectrum in
94 breast cancer among H/L, we generated whole-exome sequencing (WES) and RNA sequencing
95 (RNA-seq) data from 146 tumors from 140 H/L from Southern California and performed
96 analyses of somatic mutations, SCNAs, and gene expression.

97 **METHODS**

98 **Participants.** One hundred and forty breast-cancer patients seen at City of Hope (COH) in
99 Duarte, California were included in this study. All participants signed a written informed consent
100 approved by the COH Institutional Review Board. Inclusion criteria were: 1) self-identified as
101 H/L; 2) tumor tissue from surgery was available and the sample contained more than 40% tumor
102 based on examination by a single breast pathologist (D. Schmolze). The percentage tumor
103 ranged from 40% to 90% with an average of 64% and a median of 65% tumor. An exclusion
104 criterion was neo-adjuvant therapy as treatment could change the mutation profile. Clinical data
105 were abstracted from medical records including date at diagnosis, date at surgery, tumor stage,
106 grade, histological estrogen receptor (ER), progesterone receptor (PR) and human epidermal
107 growth factor (HER2) status, second cancers, breast-cancer recurrence, parity, history of breast

108 feeding, age at menarche, and cause of death, if applicable. Six of the 140 breast-cancer
109 patients had two primary contralateral breast cancers with tissue available for study for a total of
110 146 tumors.

111

112 **DNA and RNA sequencing**

113 DNA extraction. Germline DNA was extracted from peripheral blood cells or from formalin-fixed
114 paraffin-embedded (FFPE) normal breast tissue adjacent to tumor tissue from surgery.

115 Peripheral blood cell DNA was extracted using a standard phenol chloroform method. For FFPE
116 tissue, DNA and RNA were extracted from ten 30- μ m sections from each tumor using the

117 QIAmp DNA FFPE Tissue Kit (Qiagen) and miRNeasy Kit (Qiagen) according to manufacturer's
118 instructions. DNA was quantified with the Quant-iT PicoGreen dsDNA Assay Kit (Thermo

119 Fisher Scientific, MA). After extraction and quantification, DNA was sent to The National Cancer
120 Institute (NCI) Cancer Genomics Research Laboratory (CGR) for WES. For RNA sequencing,

121 500 ng total RNA was sent to the COH Integrative Genomics Core (IGC).

122

123 DNA library construction, hybridization, and massively parallel sequencing. Library production

124 and sequencing for 146 tumors and 140 matching normal samples was performed at CGR. The

125 KAPA HyperPlus Kit (Kapa Biosystems, Inc., Wilmington, MA) was used to generate libraries

126 from 300ng DNA according to the KAPA-provided protocol. Libraries were pooled and sequence

127 capture was performed with NimbleGen's SeqCap EZ exome v3 (Roche NimbleGen, Inc.,

128 Madison, WI, USA), according to the manufacturer's protocol. The resulting post-capture

129 enriched multiplexed sequencing libraries were used in cluster formation on an Illumina cBOT

130 (Illumina, San Diego, CA, USA) and paired-end sequencing was performed using an Illumina

131 HiSeq 4000 following Illumina-provided protocols for 2 \times 100 bp paired-end sequencing to an

132 average-fold coverage of 80X for the tumors and 30X for the germline samples. Paired-end

133 reads from each sample were aligned to human reference genome (hg19) using Novoalign

134 (v3.00.05), and the aligned binary format sequence (BAM) files were sorted and indexed using
135 SAMtools (1, 2). The sorted and indexed BAMs were processed by Picard (v1.126,
136 <https://broadinstitute.github.io/picard/>) to remove duplicate sequencing reads. Local realignment
137 around suspected sites of indels was performed using Genome Analysis Toolkit (GATK)
138 IndelRealigner (v3.3-0-g37228af). These mapped sequence reads were then base-recalibrated
139 before being used for somatic mutation calling by MuTect2 in GATK (v4.0.11.0).

140 RNA-seq. In the COH IGC, sequencing libraries were prepared with Kapa RNA HyperPrep kit
141 with RiboErase (Roche) and sequenced on a HiSeq 2500 (Illumina) with 40 million reads per
142 sample. The RNA-seq reads were aligned to hg19 genome assembly using Tophat2 (v2.0.8)
143 with default settings. The gene-expression levels were counted by obtaining raw counts with
144 HTSeq (v0.6.1p1) against Ensembl v86 annotation. The counts data were normalized using the
145 trimmed mean of M values (TMM) method implemented in R package edgeR²⁰. Log2-
146 transformed counts were used to assign PAM50 subtypes based on the subgroup-specific gene
147 centering method developed by Zhao, *et al.*²¹. We estimated Z-scores based on the corrected
148 median absolute deviation (MAD) implemented by the *robStandardize* R function in the
149 robustHD R package and defined expression outliers as gene-sample data points with robust Z-
150 scores greater than three. Raw counts of RNA-seq data for 1,189 TCGA samples (including
151 both tumor and matched normal samples) were downloaded from the Genomic Data Commons
152 (GDC) using the GDCRNATools²² R package. RNA-seq data for H/L tumor samples and TCGA
153 samples were processed and analyzed separately.

154 **Data analysis**

155 Germline variant calling. Germline variant calling from the BAM files was performed in the COH
156 IGC using GATK HaplotypeCaller (<https://software.broadinstitute.org/gatk>). Variants with a call
157 quality less than 20, read depth less than 10, or allele fraction ratio less than 20% were
158 removed. Variants in variant call format files were evaluated for pathogenicity using Ingenuity

159 Variant Analysis (IVA) version 4 (Qiagen Inc, Alameda, CA) and American College of Medical
160 Genetics and Genomics (ACMGG) guidelines were applied using the IVA ACMGG calling
161 algorithm²³. Pathogenic or likely pathogenic variants were individually evaluated by the research
162 team using the available literature and ClinVar to make a final determination²⁴.

163

164 Genetic ancestry analysis. We performed genetic ancestry estimation for each of the 140
165 women using the germline whole-exome sequencing data. We used 90 European (1000
166 Genomes), 90 African (1000 Genomes), 90 East-Asian (1000 Genomes) and 71 IA ancestry²⁵
167 reference samples. We identified the SNPs that overlap all data sets (N=9,935). We combined
168 all SNPs and dropped SNPs that did not match based on reference and alternate alleles. To
169 estimate the ancestry for each sample, we used ADMIXTURE 1.3.0 setting the K parameter to 4
170 and running the unsupervised algorithm²⁶. In addition, we used principal components analysis,
171 calculated using PLINK 1.9²⁷ as a complementary method to assess ancestry.

172

173 Somatic variant calling. We identified somatic single nucleotide variants (SNV) using MuTect2
174 in GATK4 (v4.0.11.0) suite with default parameters²⁸ and indels using GATK Indelocator. Using
175 the SNV and indel filtering method described in Pereira et al.³, we focused on frameshift, non-
176 synonymous, canonical splicing site, and stop gain mutations. Briefly, somatic mutations were
177 manually curated and considered true positives in a sample if the mutation was observed in
178 >10% of reads or with a frequency of 5-10% if in frequently mutated breast-cancer genes or
179 seen in COSMIC database²⁹. Because the tumors include both tumor and normal stromal cells,
180 it is expected that the proportion of reads will have less than the expected 50% if 100% tumor.
181 Mutations in <5% of reads, in segmental duplication regions, or indels that overlapped
182 homopolymer stretches of six or more bases were considered false positives. We did visual
183 checking using the Integrative Genomics Viewer (IGV) to assess the quality of all somatic
184 mutations. We performed Sanger sequencing on a subset of samples to confirm specific

185 mutations in *AKT1*, *BARD1*, *MAP3K1*, and *MET*. Using the filtered and annotated somatic
186 mutations, we performed a somatic-mutation significance analysis via MutSigCV³⁰ (version
187 1.3.5) on Genepattern (<https://www.genepattern.org/modules/docs/MutSigCV>). Genes with false
188 discovery rate (FDR) $q < 0.05$ are considered to be significantly mutated genes.

189

190 We compared the significant somatic mutations in our analysis with the mutations from the
191 Romero-Cordoba dataset¹⁹. Using the publicly available somatic-mutation data from the
192 Romero-Cordoba study of the Mexican patients, we combined our somatic-mutation data and
193 performed a MutSigCV analysis to identify the common significant genes. Similarly, to
194 investigate if these significantly mutated genes were associated with ancestry, we performed
195 the same analysis on breast tumors from Whites in TCGA. Using 2% as the mutation frequency
196 threshold, we performed Fisher's exact test for each frequently mutated gene for comparison.

197

198 Copy-number analysis using FACETS. We used FACETS implemented in R package facets
199 version 0.6.1³¹ to calculate CNAs. The counts of reads with the reference (ref) allele, alternate
200 (alt) allele, errors (neither ref nor alt), and deletions at a specific genomic position were
201 generated using BAM files from the 146 matched tumor-normal sample pairs using the
202 application snp-pileup in the facets package. The segmentation of each tumor sample was then
203 estimated with the critical value (cval) 150.

204 The segmentation files generated by facets served as input files for the GISTIC2.0³² on the
205 GenePattern server (<https://genepattern.broadinstitute.org/gp>) to identify significant SCNAs
206 using a q-value cutoff < 0.05 . A gene was considered as copy number altered with GISTIC2-
207 thresholded scores of -2 (deep loss), -1 (shallow loss), 1 (low-level gain) and 2 (high-level gain).
208 The GISTIC2 copy-number results and clinical data for 816 TCGA tumor samples were
209 downloaded from the cBioPortal database³³ (<https://www.cbioportal.org>). Expression outliers
210 (defined by Z-scores greater than 3.0) were considered as driven by copy-number changes if

211 greater than 90% expression outliers in a gene had a GISTIC2-thresholded copy-number score
212 of 2 (high-level gain) or 1 (low-level gain). Fisher's exact test was used to identify genes with
213 frequency difference in expression outliers, driven by copy-number alterations, between 146
214 tumor samples from H/L and 452 TCGA Whites (determined as having > 95% European
215 ancestry as described below).

216

217 Mutation-signature analysis

218 Using the previously called SNVs, we performed a mutational signature analysis via the
219 MutationalPatterns R package³⁴. Hg19 was used as the reference genome. SNVs were parsed
220 and classified into six mutation patterns (C>T, T>A, C>G, T>C, C>A and T>G) and 96
221 trinucleotide changes. Then a non-negative matrix factorization algorithm was implemented to
222 extract mutation signatures. And we compared the similarities of these mutation signatures with
223 the COSMIC mutation signatures and each mutation signature could be treated as a linear
224 combination of the 30 COSMIC mutation signatures. The 30 COSMIC mutation signature
225 percentage contribution was then computed for each tumor and a contribution heatmap was
226 generated. Within these tumor samples, we performed a signature contribution comparison
227 using the two-sided Wilcoxon rank-sum tests among the five tumor subtypes (Luminal A, luminal
228 B, basal-like, HER2-enriched and normal-like).

229 We also compared the mutation-signature analysis with the breast tumors in the Romero-
230 Cordoba dataset and the breast tumors from Whites in TCGA SNV dataset. For the significant
231 COSMIC mutation signatures identified in our dataset, we performed two-sided Wilcoxon rank-
232 sum tests among the three datasets to test if the signature was enriched in Mexican patients.

233

234

235 **RESULTS**

236 Clinical/demographic data and germline pathogenic variants. Characteristics of the 140
237 participants are shown in Table 1. The mean age at diagnosis was 48.7 years with a range
238 from ages 31 to 75 years. Nearly all of the 140 H/L were of mixed European (Eur) and IA
239 ancestry. The mean ancestry composition was 50.6% Eur, 40.8% IA, 5.9% African, and 2.7%
240 Asian although the range of ancestry proportion varied widely from <1% to 96% IA at the
241 extremes (Figure 1). As shown in the Principal Component Analysis (PCA) plots in Figure 1A,
242 H/L samples are not well-represented in TCGA project. For the six individuals with two primary
243 tumors (in the contralateral breasts), the tumors were considered independent tumors
244 (Supplemental Table 1) which was borne out by different somatic-mutation profiles. The majority
245 of the women were diagnosed with Stage I (44%) or II (43%) tumors (Table 1). There were 22
246 recurrences and 10 deaths during the time of follow-up. Of the 146 tumors, 83% were ER-
247 positive, 72% were PR-positive, and 17% were HER-2 positive and these proportions were
248 similar to White women in TCGA¹. Germline pathogenic variants in breast cancer
249 predisposition genes were identified in six participants including one *BRCA1* exon 9-12 deletion,
250 four *CHEK2* L236P, and one *NF1* Y408X variants of which the *BRCA1* and *CHEK2* variants are
251 of Indigenous-American ancestry³⁵.

252

253 Somatic mutations. We observed a total of 4510 true somatic mutations in 3391 genes in the
254 146 primary breast tumors (Supplemental Table 2). The number of mutations per individual
255 varied from 2 to 225. Using MutSigCV, we found that mutations in *PIK3CA*, *TP53*, *GATA3*,
256 *MAP3K1*, *CDH1*, *CBFB*, *PTEN*, and *RUNX1* were significant (FDR < 0.05) cancer driver
257 mutations. To identify additional, potentially significantly mutated genes in H/L, we merged the
258 mutation data from our cohort with a previously published study of Mexican breast-cancer
259 patients (N = 135)¹⁹. Within the aggregated mutation data of this combined cohort (N = 281), we
260 re-ran MutSigCV and identified one more significantly mutated gene, *AKT1*, which only occurred
261 twice in our 146 primary breast tumors. Using the statistically significantly mutated genes

262 obtained from the aggregated cohort, we visualized the mutational profiles within our cohort
263 (Figure 2a) and the variant locations for *PIK3CA* and *GATA3* (Figure 2b). For *MAP3K1* and
264 *RUNX1*, at least one tumor harbored multiple mutations in the same gene. Furthermore, in
265 *GATA3*, eight tumors had the identical splice mutation
266 (NM_002051.2:exon5:r.spl;NM_001002295.1:exon5:r.spl) that affected expression (data not
267 shown). Other genes of interest that did not meet the significance threshold (FDR < 0.05) but
268 which have been identified as significant in prior studies and were mutated in our dataset
269 included *MLL3* (aka *KMTC2*) (6%), *PTPRD* (3%), *MAP2K4* (2%), *PIK3R1* (2%), *NF1* (1%), *RB1*
270 (1%), *TBX3* (1%), *FOXA1* (1%), *PADI4* (1%), *CDKN1B* (1%), *CTCF* (1%), and *NCOR1* (1%). In
271 addition, we found mutations in *MET* (4.1%) which is not generally considered a breast-cancer
272 gene but is a known driver in other cancer types³⁶.

273
274 The frequency of mutations in genes known to be significantly mutated in breast cancer,
275 including *PIK3CA*, *MAP3K1*, *GATA3*, *CBFB*, and *MLL3/KMT2C*, were not significantly different
276 in tumors from H/L compared to tumors from White women in TCGA (FDR $q > 0.05$,
277 Supplemental Table 3). Similar to tumors from Whites, *PIK3CA* and *TP53* were the most
278 commonly mutated genes. We identified *AKT1* mutations in 2 of 146 tumors (1.4%), including
279 the E17K hotspot mutation which was found to be mutated in 8% of patients among Mexican
280 women¹⁹. After correction for multiple hypothesis testing, we found no somatic mutations
281 significantly associated with genetic ancestry.

282
283 Mutational signature analysis. To investigate the mutational processes in H/L breast-cancer
284 tumors and the association between PAM50 subtypes and mutational patterns, we adopted the
285 non-negative matrix factorization approach as proposed by Alexandrov et al.³⁷ for mutational
286 signature analysis of tumors. Signature calling revealed five major contributing signatures in the

287 146 tumors corresponding to the COSMIC signatures 1,2,3,13 and 16 (Figure 3a; Supplemental
288 Table 4). Signature 1 was detected in all 146 tumors. The contribution of COSMIC signature 1
289 was greater in luminal A and B subtypes than HER2 and basal subtypes ($p < 0.05$, two-sided
290 Wilcoxon rank sum test) (Figure 3b). Signatures 2 and 13, attributed to activity of the
291 AID/APOBEC family of cytidine deaminases, were found in tandem in 16% ($n = 23$) of the
292 tumors and were more common in tumors with HER2 subtype compared to luminal A and B
293 subtypes (Figure 3b). We found that 13 tumors were homozygous and 29 tumors were
294 heterozygous for a common 29.5kbp germline deletion spanning most of *APOBEC3B*. Tumors
295 with the deletion had a higher proportion of COSMIC signatures 2 ($p = 0.0005$, Wilcoxon rank
296 sum test) and 13 ($p = 0.0008$, Wilcoxon rank sum test). Signature 3, attributed to defects of
297 homologous recombination double-stranded DNA break-repair, was found significantly more
298 often in basal subtypes than the other PAM50 subtypes ($p < 0.05$, two-sided Wilcoxon rank sum
299 test) including the tumor with the germline *BRCA1* exon 9-12 deletion. We observed a group of
300 tumors ($N=40$, 27.4%) with more than 5% COSMIC signature 16 contributions. Since this was
301 not previously reported in other breast tumor studies, we re-examined other datasets, using the
302 same analytic pipeline used herein. We found that signature 16 was present in 20 (19.6%)
303 tumors in a previous study of Mexican breast-cancer patients¹⁹ which was not significantly
304 different than the proportion in our dataset ($p = 0.18$, Fisher's exact test). The proportion in
305 tumors from TCGA White women ($N=75$; 8.9%) was significantly lower than in our dataset ($p <$
306 0.001 , Fisher's exact test) (Figure 3c) and in the Romero-Cordoba et al. dataset ($p < 0.0001$,
307 Fisher's exact test). The percentage of this signature was significantly higher in luminal A and B
308 subtypes compared to HER2 and basal tumors ($p < 0.05$, two-sided Wilcoxon rank sum test)
309 (Figure 3b).

310

311 Somatic CNAs (SCNAs). Using GISTIC2, we identified chromosome arm-level SCNAs that were
312 significantly ($q < 0.05$) amplified at 1q, 8q, 6p, 1p, 6q, 16p, 20q, 8p, 12q and deleted at 22q,

313 16p, 17p, 8p (Supplemental Table 5). In addition to these broad SCNAs, we identified
314 significantly ($q < 0.05$) amplified or deleted focal regions including 29 peak regions of
315 amplification and 48 regions of deletion (Figure 4A). Seven recurrently amplified regions contain
316 common oncogenes (*FGFR1*, *MYC*, *CCND1*, *MDM2*, *IGF1R*, *ERBB2*, and *ZFP217*); one
317 recurrently deleted region contains *TP53* (Figure 4A). By integrative analysis of RNA-seq gene-
318 expression data and copy-number data, we observed that greater than 90% of expression
319 outliers (defined by robust Z-score greater than 3.0) in *ERBB2*, *FGFR1*, *IGF1R*, and *MDM2*
320 were associated with copy-number gain (Figure 4B). Therefore, we sought to identify
321 expression outliers from 1,121 genes contained in the 29 copy-number amplification peak
322 regions for the 146 H/L breast tumor samples and 452 White TCGA breast tumor samples. Of
323 1,121 genes in the 29 regions, over 90% of expression outliers were associated with copy-
324 number gain in 214 genes, including 88 genes from the 146 H/L samples, 62 genes from the
325 452 TCGA White samples, and 64 genes from both sample groups (Supplemental Table 6).
326 Eighteen of 214 genes had significant ($FDR < 0.05$) difference in frequency of expression
327 outliers copy number between the 146 H/L and 452 TCGA White tumor samples (Table 2 and
328 the top 18 rows in Supplemental Table 6). Expression outliers from those genes were more
329 prevalent in the 146 H/L than in the 452 White tumors because we focused on the 29 copy-
330 number regions (Figure 4A) found in H/L (Table 2; Supplemental Table 6).

331
332 Using this combined copy-number and gene-expression analysis approach, we identified
333 *KIAA0100*, also known as Breast Cancer Overexpressed Gene 1 (*BCOX1*), as the top gene that
334 was systematically different between Whites (TCGA) and our H/L cohort. (Figure 5A). Since this
335 gene is within ~11 megabases of *ERBB2* on chromosome 17q, we investigated whether it was
336 part of the *ERBB2* GISTIC amplification peak. We observed that the peaks for copy-number
337 amplifications (Figure 5B) were distinct for *KIAA0100* and *ERBB2* are at 17q11.2 and 17q12.

338

339 **Discussion**

340

341 We analyzed tumor-germline sequencing data combined with RNA-seq data from 146 tumors
342 from 140 self-identified H/L recruited from a single center in the Los Angeles region. As
343 expected, the majority were of mixed European and IA ancestries. Since TCGA has extremely
344 limited samples of breast cancer from H/L and particularly of H/L of mixed IA ancestry, our
345 report fills a critical gap in the landscape of somatic mutations and copy-number alterations in
346 this increasing US population. Together, our analyses and the recent paper focused on Mexican
347 women living in Mexico ¹⁹ substantially enhance the data in the public domain for women of H/L
348 heritage.

349

350 The most commonly mutated gene in our population was *PIK3CA* which is the most commonly
351 mutated gene in TCGA White samples. For women with advanced ER+/HER2- breast cancers,
352 alpelisib is a currently approved therapy, and our results suggest that this therapy should be
353 useful in a large fraction of H/L women. The Romero-Cordoba et. al. study identified a high
354 frequency (8%) of the E17K activating *AKT1* mutation indicating such women may benefit from
355 AKT inhibitors. We only identified two tumors with mutations in *AKT1* and only one with the
356 E17K mutation. The difference between our results and those of Romero-Cordoba may be due
357 to chance, differences in selection criteria between the two cohorts, and/or differences in
358 environmental exposures between the two cohorts. Since the ancestry of our population is
359 similar, it is unlikely that the differences we observed are due to germline-genetic differences
360 between the two cohorts.

361

362 We performed analyses of the somatic-mutational signatures and compared them to the TCGA
363 dataset. Our analysis identified COSMIC signature 16 (contribution > 5%) in a significant
364 fraction of tumors (27.4%) in our dataset with similar rates in the data from Romero-Cordoba et

365 al. who analyzed breast tumors from Mexican patients. Because Romero-Cordoba et al. used a
366 contribution cutoff in their mutation-signature-analysis pipeline, they did not report this signature.
367 However, in our analysis, we implemented the non-negative matrix factorization algorithm and
368 no contribution cutoff was applied such that signature 16 was observed. There were significantly
369 lower rates of this signature in TCGA White women ($p < 0.001$). We do not believe our finding is
370 a technical artifact from FFPE because this signature was found in frozen tissue in the Romero-
371 Cordoba et al data. No known genetic or environmental exposures that predispose to this
372 signature have been reported and prior studies have not found this mutational signature in
373 breast cancer, although it has been reported to be common in liver cancers³⁷.

374

375 Other COSMIC signatures were the same as those previously reported in TCGA. We found
376 signatures 2 and 13 associated with APOBEC loss as a relatively common finding, associated
377 with HER2-amplified tumors and specifically with the germline APOBEC copy-number variant
378 similar to previous reports³⁸. The APOBEC3B common 29.5-kbp germline deletion results in the
379 fusion of APOBEC3A and the 3'UTR of APOBEC3B³⁹. This fusion generates a more stable
380 APOBEC3A mRNA, resulting in increased expression of APOBEC3A, higher overall mutation
381 burden, and a higher odds ratio of developing breast cancer^{40, 41}. We also found Signature 3,
382 associated with defects in homologous recombination repair as a common signature, which is
383 over-represented in basal-like tumors as previously reported^{37, 42}.

384

385 Our copy-number analyses identified copy-number gains, i.e., 1q, 8q, 17q which are common in
386 breast cancer in other populations^{1, 2}. We also identified several known CNAs which were
387 recurrently gained in our dataset. In combined analysis of copy-number alterations and gene
388 expression, we identified KIAA0100 (BCOX1) as a recurrently amplified region with high gene
389 expression which was more common in tumors from H/L than tumors from White women in
390 TCGA. KIAA0100 was originally identified in a screen for genes that were more frequently found

391 in breast tumor than in normal breast tissue⁴³ and increased expression was associated with
392 poor prognosis^{43, 44}. Knock-down of *KIAA0100* by siRNA in the breast-cancer cell line MDA-MB-
393 231 reduced cell aggregation, reattachment, cell metastasis and invasion⁴⁵. Thus, KIAA0100
394 may be of interest for further study in understanding the biology of tumors in H/L and stratifying
395 women for risk of recurrence.

396

397 Our study has several limitations. We included only women who did not have neoadjuvant
398 therapy prior to surgical resection. We chose this subset of women to avoid effects possibly
399 induced by neoadjuvant chemotherapy such as new mutations and/or selection for resistant
400 subclones. However, because neoadjuvant therapy is more likely to be given to patients with
401 large tumors and/or tumors with poor prognosis⁴⁶, tumors included in our study may have some
402 differences in comparison with prior studies due to these selection criteria. For example,
403 because most triple-negative breast tumors are first treated with neoadjuvant therapy, the
404 proportion of triple-negative tumors in our study was lower than previously reported⁴⁷. Our
405 analysis of tumor copy-number alterations was based on WES data. Although WES and other
406 forms of targeted sequencing are used for CNA analysis, it makes it difficult to conduct one-to-
407 one comparisons to array-based or whole genome sequencing-based analyses. Therefore, we
408 limited our analyses to copy-number events that also demonstrated gene-expression
409 differences across populations. Finally, although our study substantially increases the number
410 of tumors analyzed by WES in H/L, the overall numbers are still substantially lower than in
411 White women. In particular, we are likely underpowered to discover low frequency, ethnic
412 and/or ancestry-specific drivers that may be unique to this population. There also were too few
413 recurrences and deaths for statistical analyses.

414

415 In summary, we conducted a comprehensive characterization of somatic mutations, CNAs, and
416 gene expression in 146 breast tumors from 140 H/L from Los Angeles County, California. We

417 found that COSMIC signature 16 was more common in our dataset and a recently published
418 dataset of Mexican women living in Mexico, suggesting that this signature may be important in
419 self-reported H/L/Hispanic women and potentially useful to understand differences at diagnosis
420 and for outcome. Finally, our combined CNA and gene-expression analysis suggested that
421 KIAA0100 may be a possible driver of breast-cancer aggressiveness in a subset of our sample.
422 These results should be useful to understanding the biology and guiding therapy for breast
423 cancer among H/L.

424

425 Acknowledgments:

426 This work was funded by the National Cancer Institute (R01CA184585, K24CA169004), the
427 National Institute on Minority Health and Health Disparities Division of Intramural Research, and
428 the California Initiative to Advance Precision Medicine (OPR18111). Research reported in this
429 publication included work performed in the City of Hope Integrative Genomics Core and the
430 Pathology Core supported by the National Cancer Institute of the National Institutes of Health
431 under grant number P30CA033572. The content and views are solely the responsibility of the
432 authors and should not be construed to represent the views of the National Institutes of Health.
433 SLN and this research were partially funded by the Morris and Horowitz Families Professorship.
434 CDA is supported by the National Heart, Lung, and Blood Institute (NHLBI T32HL007118)
435 through the training Program in Molecular and Integrative Physiological Sciences at the Harvard
436 T.H. Chan School of Public Health. LF is supported by R01CA204797. JNW was supported by
437 NIH RC4 CA153828; Breast Cancer Research Foundation (#20-172), and American Society of
438 Clinical Oncology Conquer Cancer® Research Professorship in Breast Cancer Disparities.

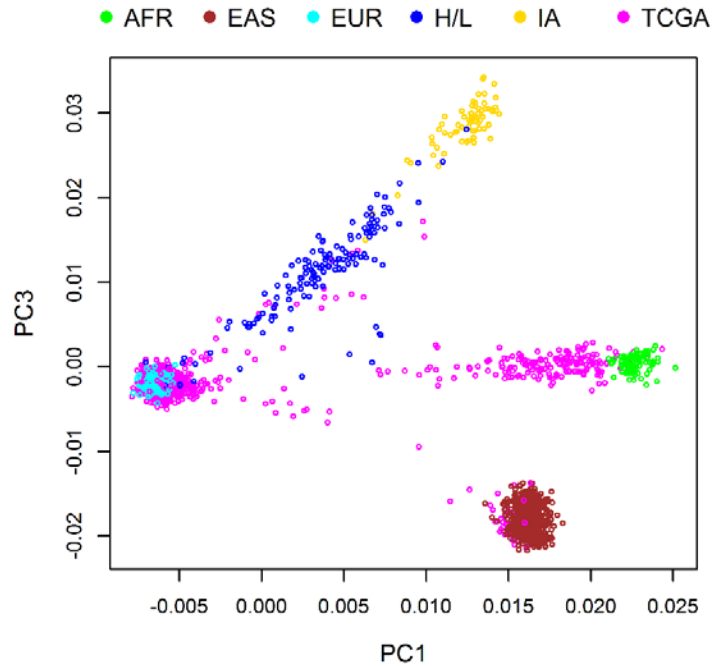
439

440 Conflicts of interest. JNW is a speaker for the Bureau for AstraZeneca, and an employee at
441 Natera. No other conflicts of interest from authors

442 Data are being deposited in dbGAP and will be made available at the time of acceptance.

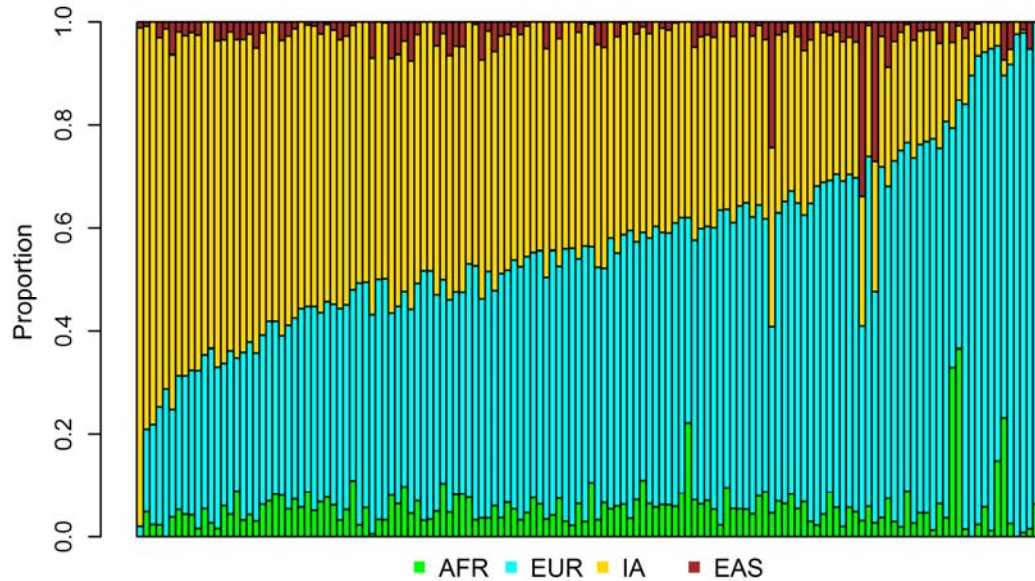
443 **Figure 1**

444 A.



445

446 B.



447

448 **Figure 1.** Ancestry of the cohort. Results of principal components analysis comparing the

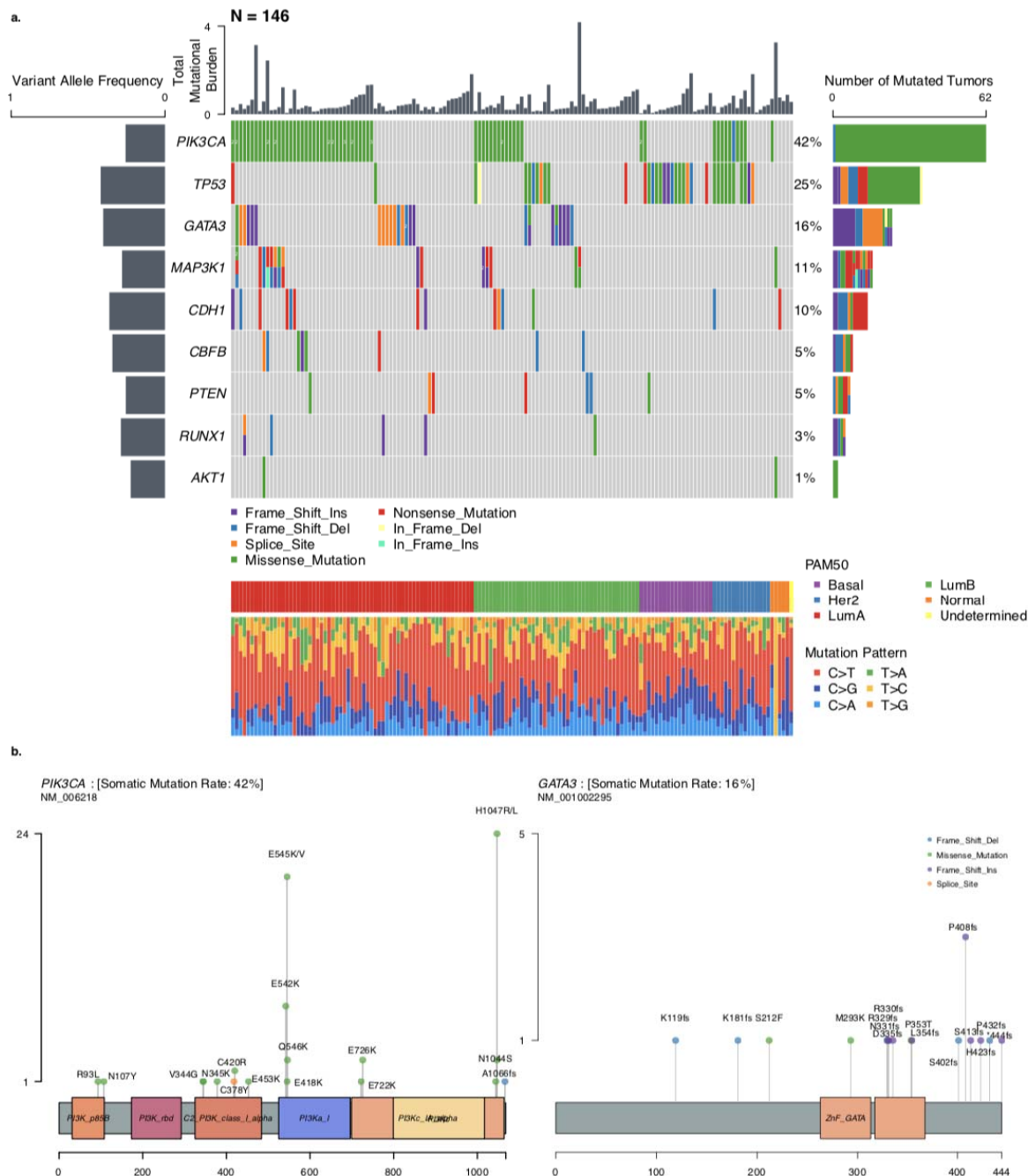
449 values for samples on principal component (PC) 1 (x-axis) and PC3 (y-axis) (A). Each dot

450 represents the results from one individual: Hispanic/Latina (H/L) (dark blue); TCGA (pink); and

451 reference populations including African (AFR), Yoruban individuals from Nigeria from HapMap
452 (light green); East Asians (EAS), Han Chinese from HapMap (brown); European American
453 (EUR), CEBP from HapMap (light blue); and Indigenous American (IA) (yellow) from Mexico.
454 PC2 (not shown) captures individuals of Asian and Indigenous-American ancestry. (B). Results
455 from ADMIXTURE analysis. Each vertical bar represents estimate of ancestry from one
456 individual. Ancestry is assigned for each individual as a fraction of either African (green), Asian
457 (brown), European (light blue) or Indigenous American (yellow) ancestry.
458

459 **Figure 2**

Figure 2



460
 461 **Figure 2.** Tumor mutational burdens and somatic-mutational profiles. a. Mutation plot of nine
 462 significantly mutated genes in the 146 tumors. Different mutation classifications are color-coded.
 463 Numbers are shown where multiple mutations of the same classification were detected. Total

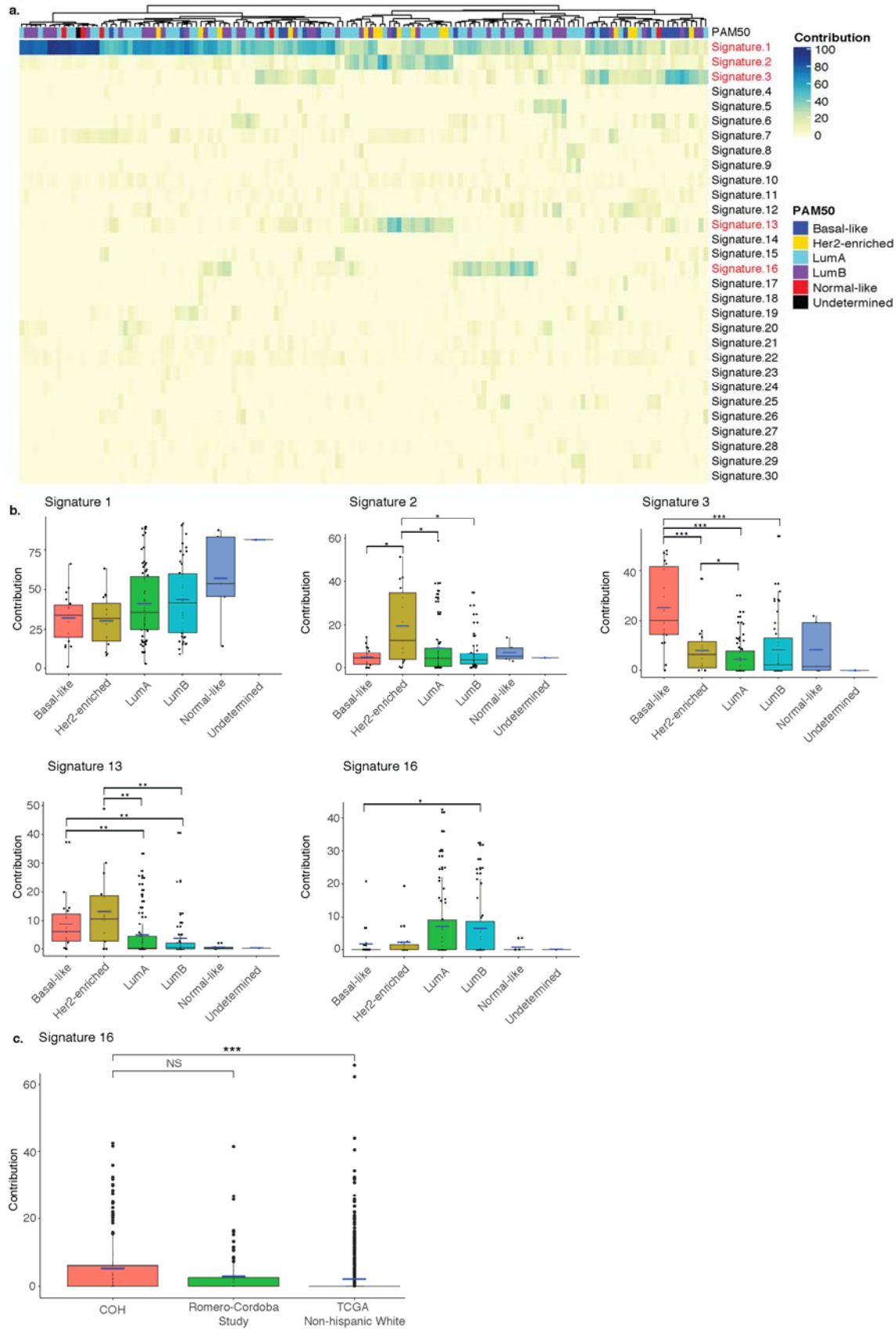
464 mutational burden for each tumor is shown as a bar chart on top. The mean variant allelic
465 frequency is shown for each gene on the left. PAM50 subtype and mutation pattern for each
466 tumor are shown at the bottom. b. Lollipop plot of *PIK3CA* and *GATA3* mutations within the 146
467 tumors. Mutation classifications are color coded and amino-acid changes are specified for each
468 mutation.

469

470

471

Figure 3



473 **Figure 3.** Mutational Signatures. a. Unsupervised clustered heatmap of contributions from each
474 mutational signatures for the 146 tumors. Significant signatures are highlighted in red. PAM50
475 subtype for each tumor is shown on top of the heatmap. b. Box-plots comparison of the
476 contributions of the five significant mutational signatures (Signature 1, 2, 3, 13, 16) across the
477 PAM50 subtypes. Statistical significance levels are indicated within the box plots (*: $p < 0.05$; **: $p < 0.01$; ***: $p < 0.001$, Wilcoxon Rank-sum test). c. Box plot of Signature 16 contributions in
478 the 146 tumors from the Hispanic-Mexican cohort (COH), Romero-Cordoba Study and the Non-
479 Hispanic White tumors in the TCGA dataset. Statistical significance levels are indicated within
480 the box plot (NS: not significant, $p > 0.05$; *: $p < 0.05$; **: $p < 0.01$; ***: $p < 0.001$, Wilcoxon
481 Rank-sum test).
482

483

484

485

486

487

488

489

490

491

492

493

494

495

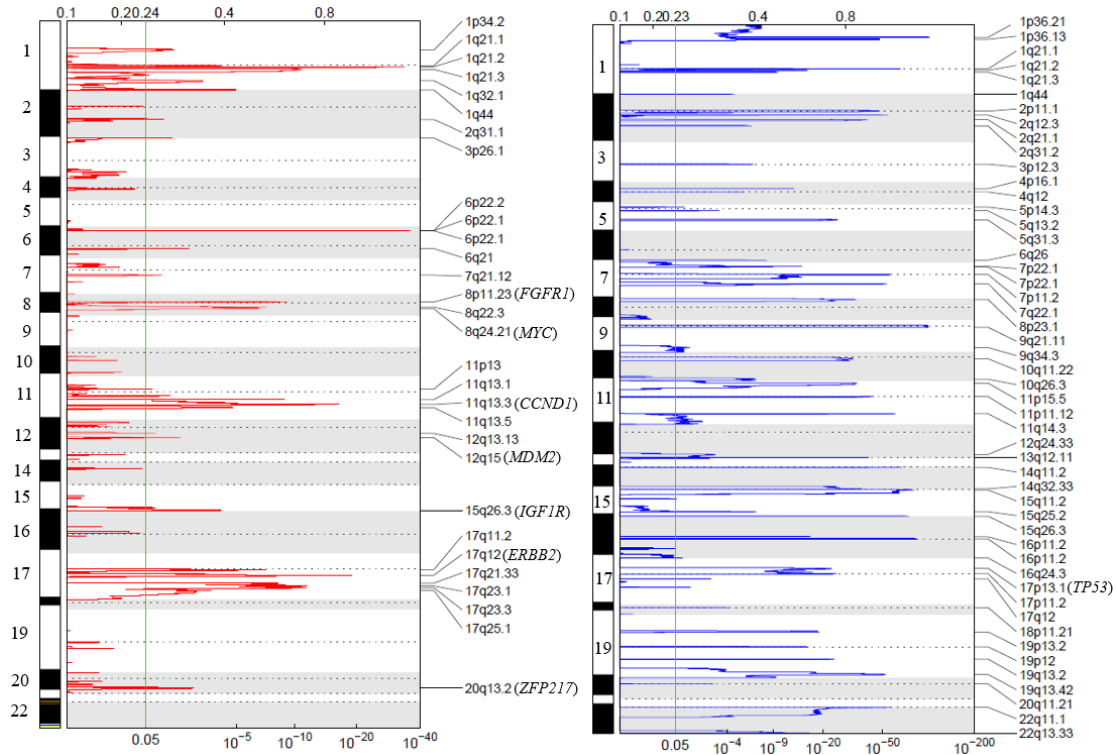
496

497

498

499 Figure 4.

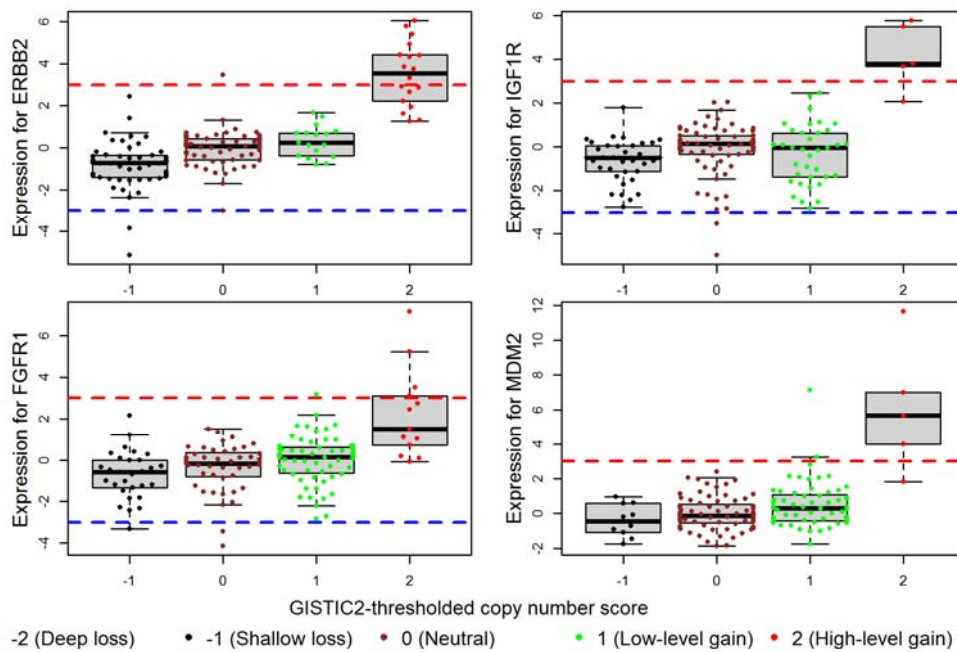
500 **A.**



501

502

503 **B.**



504

505

506 **Figure 4:** Copy-number alterations.

507 a. Genomic regions of significant copy-number gain (left) and loss (right) identified by GISTIC2.

508 Common oncogenes and tumor suppressor genes are in parentheses next to the corresponding

509 cytobands. The green vertical line marks the GISTIC2 q value of 0.05 (bottom x-axis). b.

510 Outlying gene expression and copy-number gain in four genes in 146 H/L breast tumor

511 samples. Gene-expression values on the y-axis are Z-scores estimated by robust

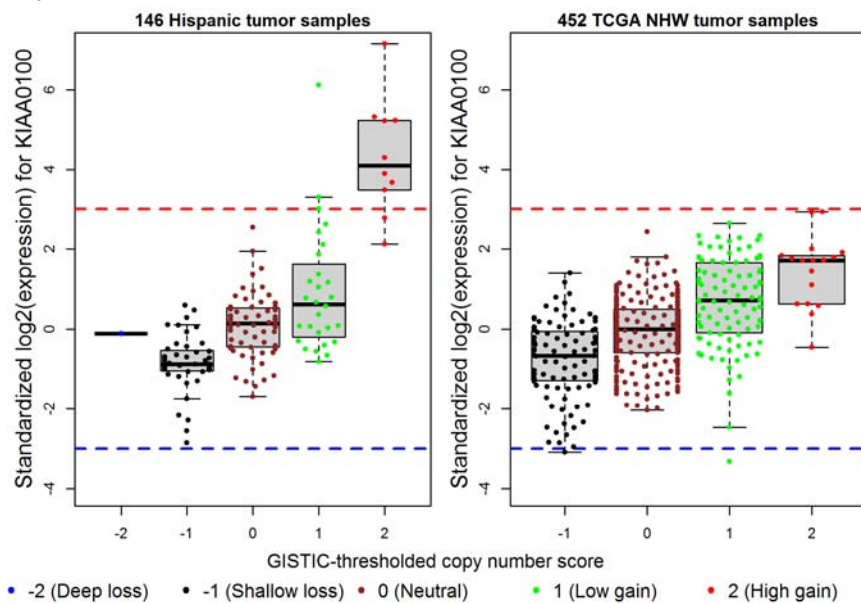
512 standardization; the red dash line of Z-score = 3 and blue dash line of Z-score = -3 are cutoff

513 values for outliers of over-expression and under-expression, respectively.

514

515 **Figure 5.**

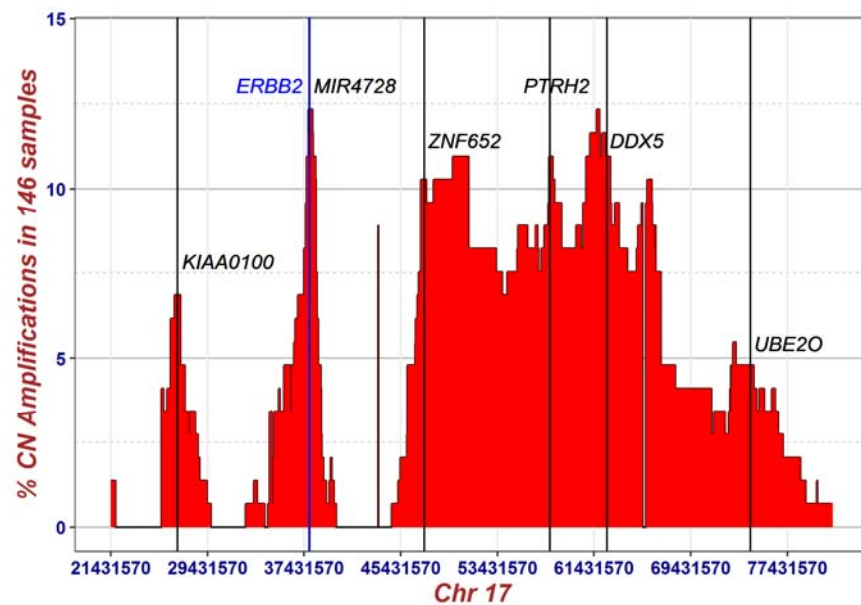
516 **A.**



517

518

519 **B.**



520

521 **Figure 5.** Expression outliers and copy-number gain in KIAA0100.

522 a. Distribution of gene expression and GISTIC2-thresholded copy-number scores in *KIAA0100*

523 for 146 breast tumor samples from Hispanic whites and 452 breast tumor samples from TCGA

524 Non-Hispanic whites (A). The y-axis is standardized gene-expression values (Z-scores)

525 estimated robustly based on the corrected median absolute deviation (MAD). Red and blue
526 dashed lines represent Z-score of 3 and -3, respectively. b. Distribution of proportion of high-
527 level copy-number gain for 950 genes spanning the 6 amplified regions of 17q11.2, 17q12,
528 17q21.33, 17q23.1, 17q23.3, and 17q25.1. Y-axis is the percentage of the 146 Hispanic
529 samples with GISTIC2-thresholded copy-number score of 2; x-axis is genomic boundaries
530 (Chr17: 21431570 – 81188573, hg19) for the six significantly amplified regions determined by
531 GISTIC2. The vertical lines mark the genomic locations of *KIAA0100* (*BCOX1*, 17q11.2) at
532 Chr17:26941457 – 26972177, *ERBB2* (17q12) at Chr17: 37844336 – 37873910, *MIR4728*
533 (microRNA 4728, 17q12) at Chr17: 37882747 – 37882814, *ZNF652* (17q21.33) at Chr17:
534 47366567 – 47439476, *PTRH2* (17q23.1) at Chr17: 57774666 – 57784959, *DDX5* (17q23.3) at
535 Chr17: 62494371 – 62503156, and *UBE2O* (17q25.1) at Chr17: 74385612 – 74449288.

536

537

538

539 **Table 1.** Patient and tumor characteristics of 140 H/L breast-cancer cases and their 146 breast
 540 tumors

Patient characteristics	Mean	Range	Median	
Age at diagnosis (years)	48.7	31-75	48	
Breastfeeding (months)	7.2	0-84	2	
Parity (number children)	2.3	0-8	2	
Age at menarche (years)	12.6	9-18	12	
Tumor characteristics	Positive	Negative	Unknown	Equivocal
Estrogen receptor	120 (82%)	25 (17%)	1 (0.7%)	
Progesterone receptor	104 (72%)	41 (28%)	1 (0.7%)	
HER2	25 (17%)	116 (80%)	1 (0.7%)	4 (3%)
Stage at diagnosis	I	II	III	IV
	63 (44%)	63 (43%)	17 (12%)	3(2%)

541

542

543 Table 2: Frequency difference in expression outliers driven by copy number gain
 544 between 146 tumors from H/L and 452 tumors from TCGA White

Gene	GISTIC2 gain region	Specific to H/L*	GISTIC2 q value	# of Outliers in 146 H/L	Freq of Outliers in 146 H/L	# of Outliers in 452 White	Freq of Outliers in 452 White	Fisher Exact p value**	BH adjusted p value
KIAA0100	17q11.2	yes	7.85E-08	11	0.08	0	0	1.37E-07	2.93E-05
DSCC1	8q24.21	<u>yes</u>	1.18E-06	7	0.05	0	0	4.63E-05	4.95E-03
C4BPA	1q32.1	yes	6.24E-04	10	0.07	4	0.01	2.31E-04	9.88E-03
C4BPB	1q32.1	yes	6.24E-04	6	0.04	0	0	1.96E-04	9.88E-03
RNF169	11q13.5	yes	1.90E-05	12	0.08	6	0.01	1.41E-04	9.88E-03
POLDIP2	17q11.2	yes	7.85E-08	10	0.07	5	0.01	5.48E-04	1.95E-02
FOXJ3	1p34.2	yes	8.16E-03	7	0.05	2	0	1.05E-03	2.94E-02
MIR4728	17q12	no	1.02E-19	12	0.08	9	0.02	1.10E-03	2.94E-02
MYBPH	1q32.1	yes	6.24E-04	8	0.05	4	0.01	2.13E-03	3.95E-02
SAP30BP	17q25.1	yes	2.60E-04	10	0.07	7	0.02	2.22E-03	3.95E-02
SDF2	17q11.2	yes	7.85E-08	8	0.05	4	0.01	2.13E-03	3.95E-02
UBE2O	17q25.1	yes	2.60E-04	12	0.08	10	0.02	1.91E-03	3.95E-02
AHCTF1	1q44	no	1.16E-05	5	0.03	1	0	3.96E-03	4.71E-02
GSDMC	8q24.21	<u>yes</u>	1.18E-06	10	0.07	8	0.02	3.95E-03	4.71E-02
MTF1	1p34.2	yes	8.16E-03	4	0.03	0	0	3.44E-03	4.71E-02
PIGS	17q11.2	yes	7.85E-08	6	0.04	2	0	3.49E-03	4.71E-02
QSER1	11p13	no	3.38E-02	9	0.06	6	0.01	3.13E-03	4.71E-02
UNC13D	17q25.1	yes	2.60E-04	5	0.03	1	0	3.96E-03	4.71E-02

545 H/L:Hispanic/Latino; White: Non-hispanic White; Freq: frequency; GISTIC2: GISTIC2 algorithm for copy number
 546 analysis.

547 *GISTIC2 gain regions are identified in the 146 HW samples, but not in the 663 TCGA Caucasian samples based on
 548 GISTIC2 results published by Romero-Cordoba SL, et al[9]; the 8q24.21 region was identified in both groups,
 549 however, the wide-peak boundary for the 663 TCGA Caucasian samples (chr8:128657453-128779930) was narrower
 550 than that for the 146 HW samples (chr8:114449162-130760646), therefore, DSCC1 and GSDMC are included in
 551 8q24.21 from the 146 H/L samples, but not in the 8q24.21 from the 663 TCGA Caucasian samples.

552 **frequency difference in the number of expression outliers between H/L and White group was tested using the
 553 Fisher's exact method.

554

555 References:

- 556 1. Cancer Genome Atlas N. Comprehensive molecular portraits of human breast tumours.
557 Nature. 2012;490(7418):61-70. Epub 20120923. doi: 10.1038/nature11412. PubMed PMID:
558 23000897; PMCID: PMC3465532.
- 559 2. Curtis C, Shah SP, Chin SF, Turashvili G, Rueda OM, Dunning MJ, Speed D, Lynch AG,
560 Samarajiwa S, Yuan Y, Graf S, Ha G, Haffari G, Bashashati A, Russell R, McKinney S, Group
561 M, Langerod A, Green A, Provenzano E, Wishart G, Pinder S, Watson P, Markowitz F, Murphy
562 L, Ellis I, Purushotham A, Borresen-Dale AL, Brenton JD, Tavare S, Caldas C, Aparicio S. The
563 genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups.
564 Nature. 2012;486(7403):346-52. Epub 20120418. doi: 10.1038/nature10983. PubMed PMID:
565 22522925; PMCID: PMC3440846.
- 566 3. Pereira B, Chin SF, Rueda OM, Vollan HK, Provenzano E, Bardwell HA, Pugh M, Jones
567 L, Russell R, Sammut SJ, Tsui DW, Liu B, Dawson SJ, Abraham J, Northen H, Peden JF,
568 Mukherjee A, Turashvili G, Green AR, McKinney S, Oloumi A, Shah S, Rosenfeld N, Murphy L,
569 Bentley DR, Ellis IO, Purushotham A, Pinder SE, Borresen-Dale AL, Earl HM, Pharoah PD,
570 Ross MT, Aparicio S, Caldas C. The somatic mutation profiles of 2,433 breast cancers refines
571 their genomic and transcriptomic landscapes. Nat Commun. 2016;7:11479. Epub 20160510.
572 doi: 10.1038/ncomms11479. PubMed PMID: 27161491; PMCID: PMC4866047.
- 573 4. Andre F, Ciruelos E, Rubovszky G, Campone M, Loibl S, Rugo HS, Iwata H, Conte P,
574 Mayer IA, Kaufman B, Yamashita T, Lu YS, Inoue K, Takahashi M, Papai Z, Longin AS, Mills D,
575 Wilke C, Hirawat S, Juric D, Group S-S. Alpelisib for PIK3CA-Mutated, Hormone Receptor-
576 Positive Advanced Breast Cancer. N Engl J Med. 2019;380(20):1929-40. doi:
577 10.1056/NEJMoa1813904. PubMed PMID: 31091374.
- 578 5. Huang SF, Liu HP, Li LH, Ku YC, Fu YN, Tsai HY, Chen YT, Lin YF, Chang WC, Kuo
579 HP, Wu YC, Chen YR, Tsai SF. High frequency of epidermal growth factor receptor mutations
580 with complex patterns in non-small cell lung cancers related to gefitinib responsiveness in
581 Taiwan. Clin Cancer Res. 2004;10(24):8195-203. doi: 10.1158/1078-0432.CCR-04-1245.
582 PubMed PMID: 15623594.
- 583 6. Arrieta O, Cardona AF, Martin C, Mas-Lopez L, Corrales-Rodriguez L, Bramuglia G,
584 Castillo-Fernandez O, Meyerson M, Amieva-Rivera E, Campos-Parra AD, Carranza H, Gomez
585 de la Torre JC, Powazniak Y, Aldaco-Sarvide F, Vargas C, Trigo M, Magallanes-Maciél M,
586 Otero J, Sanchez-Reyes R, Cuello M. Updated Frequency of EGFR and KRAS Mutations in
587 NonSmall-Cell Lung Cancer in Latin America: The Latin-American Consortium for the
588 Investigation of Lung Cancer (CLICaP). J Thorac Oncol. 2015;10(5):838-43. doi:
589 10.1097/JTO.0000000000000481. PubMed PMID: 25634006.
- 590 7. Carrot-Zhang J, Soca-Chafre G, Patterson N, Thorner AR, Nag A, Watson J, Genovese
591 G, Rodriguez J, Gelbard MK, Corrales-Rodriguez L, Mitsuishi Y, Ha G, Campbell JD, Oxnard
592 GR, Arrieta O, Cardona AF, Gusev A, Meyerson M. Genetic Ancestry Contributes to Somatic
593 Mutations in Lung Cancers from Admixed Latin American Populations. Cancer Discov.
594 2021;11(3):591-8. Epub 20201202. doi: 10.1158/2159-8290.CD-20-1165. PubMed PMID:
595 33268447; PMCID: PMC7933062.
- 596 8. Li J, Xu C, Lee HJ, Ren S, Zi X, Zhang Z, Wang H, Yu Y, Yang C, Gao X, Hou J, Wang
597 L, Yang B, Yang Q, Ye H, Zhou T, Lu X, Wang Y, Qu M, Yang Q, Zhang W, Shah NM,
598 Pehrsson EC, Wang S, Wang Z, Jiang J, Zhu Y, Chen R, Chen H, Zhu F, Lian B, Li X, Zhang Y,
599 Wang C, Wang Y, Xiao G, Jiang J, Yang Y, Liang C, Hou J, Han C, Chen M, Jiang N, Zhang D,
600 Wu S, Yang J, Wang T, Chen Y, Cai J, Yang W, Xu J, Wang S, Gao X, Wang T, Sun Y. A
601 genomic and epigenomic atlas of prostate cancer in Asian populations. Nature.
602 2020;580(7801):93-9. Epub 20200325. doi: 10.1038/s41586-020-2135-x. PubMed PMID:
603 32238934.

- 604 9. Yuan J, Hu Z, Mahal BA, Zhao SD, Kensler KH, Pi J, Hu X, Zhang Y, Wang Y, Jiang J,
605 Li C, Zhong X, Montone KT, Guan G, Tanyi JL, Fan Y, Xu X, Morgan MA, Long M, Zhang Y,
606 Zhang R, Sood AK, Rebbeck TR, Dang CV, Zhang L. Integrated Analysis of Genetic Ancestry
607 and Genomic Alterations across Cancers. *Cancer Cell*. 2018;34(4):549-60 e9. doi:
608 10.1016/j.ccell.2018.08.019. PubMed PMID: 30300578; PMCID: PMC6348897.
- 609 10. Carrot-Zhang J, Chambwe N, Damrauer JS, Knijnenburg TA, Robertson AG, Yau C,
610 Zhou W, Berger AC, Huang KL, Newberg JY, Mashl RJ, Romanel A, Sayaman RW, Demichelis
611 F, Felau I, Frampton GM, Han S, Hoadley KA, Kemal A, Laird PW, Lazar AJ, Le X, Oak N, Shen
612 H, Wong CK, Zenklusen JC, Ziv E, Cancer Genome Atlas Analysis N, Cherniack AD, Beroukhim
613 R. Comprehensive Analysis of Genetic Ancestry and Its Molecular Correlates in Cancer. *Cancer*
614 *Cell*. 2020;37(5):639-54 e6. doi: 10.1016/j.ccell.2020.04.012. PubMed PMID: 32396860;
615 PMCID: PMC7328015.
- 616 11. DeSantis CE, Miller KD, Goding Sauer A, Jemal A, Siegel RL. Cancer statistics for
617 African Americans, 2019. *CA Cancer J Clin*. 2019;69(3):211-33. Epub 20190214. doi:
618 10.3322/caac.21555. PubMed PMID: 30762872.
- 619 12. Huo D, Hu H, Rhie SK, Gamazon ER, Cherniack AD, Liu J, Yoshimatsu TF, Pitt JJ,
620 Hoadley KA, Troester M, Ru Y, Lichtenberg T, Sturtz LA, Shelley CS, Benz CC, Mills GB, Laird
621 PW, Shriver CD, Perou CM, Olopade OI. Comparison of Breast Cancer Molecular Features and
622 Survival by African and European Ancestry in The Cancer Genome Atlas. *JAMA Oncol*.
623 2017;3(12):1654-62. doi: 10.1001/jamaoncol.2017.0595. PubMed PMID: 28472234; PMCID:
624 PMC5671371.
- 625 13. John EM, Phipps AI, Davis A, Koo J. Migration history, acculturation, and breast cancer
626 risk in Hispanic women. *Cancer Epidemiol Biomarkers Prev*. 2005;14(12):2905-13. doi:
627 10.1158/1055-9965.EPI-05-0483. PubMed PMID: 16365008.
- 628 14. Fejerman L, Ahmadiyah N, Hu D, Huntsman S, Beckman KB, Caswell JL, Tsung K, John
629 EM, Torres-Mejia G, Carvajal-Carmona L, Echeverry MM, Tuazon AM, Ramirez C, Consortium
630 C, Gignoux CR, Eng C, Gonzalez-Burchard E, Henderson B, Le Marchand L, Kooperberg C,
631 Hou L, Agalliu I, Kraft P, Lindstrom S, Perez-Stable EJ, Haiman CA, Ziv E. Genome-wide
632 association study of breast cancer in Latinas identifies novel protective variants on 6q25. *Nat*
633 *Commun*. 2014;5:5260. Epub 20141020. doi: 10.1038/ncomms6260. PubMed PMID: 25327703;
634 PMCID: PMC4204111.
- 635 15. Hendrick RE, Monticciolo DL, Biggs KW, Malak SF. Age distributions of breast cancer
636 diagnosis and mortality by race and ethnicity in US women. *Cancer*. 2021;127(23):4384-92.
637 Epub 20210824. doi: 10.1002/cncr.33846. PubMed PMID: 34427920.
- 638 16. Primm KM, Zhao H, Hernandez DC, Chang S. A Contemporary Analysis of Racial and
639 Ethnic Disparities in Diagnosis of Early-Stage Breast Cancer and Stage-Specific Survival by
640 Molecular Subtype. *Cancer Epidemiol Biomarkers Prev*. 2022;31(6):1185-94. doi:
641 10.1158/1055-9965.EPI-22-0020. PubMed PMID: 35314859.
- 642 17. Fejerman L, Hu D, Huntsman S, John EM, Stern MC, Haiman CA, Perez-Stable EJ, Ziv
643 E. Genetic ancestry and risk of mortality among U.S. Latinas with breast cancer. *Cancer Res*.
644 2013;73(24):7243-53. Epub 20131031. doi: 10.1158/0008-5472.CAN-13-2014. PubMed PMID:
645 24177181; PMCID: PMC3881587.
- 646 18. Marker KM, Zavala VA, Vidaurre T, Lott PC, Vasquez JN, Casavilca-Zambrano S,
647 Calderon M, Abugattas JE, Gomez HL, Fuentes HA, Picoaga RL, Cotrina JM, Neciosup SP,
648 Castaneda CA, Morante Z, Valencia F, Torres J, Echeverry M, Bohorquez ME, Polanco-
649 Echeverry G, Estrada-Florez AP, Serrano-Gomez SJ, Carmona-Valencia JA, Alvarado-Cabrero
650 I, Sanabria-Salas MC, Velez A, Donado J, Song S, Cherry D, Tamayo LI, Huntsman S, Hu D,
651 Ruiz-Cordero R, Balassanian R, Ziv E, Zabaleta J, Carvajal-Carmona L, Fejerman L,
652 Consortium C. Human Epidermal Growth Factor Receptor 2-Positive Breast Cancer Is
653 Associated with Indigenous American Ancestry in Latin American Women. *Cancer Res*.

- 654 2020;80(9):1893-901. Epub 20200403. doi: 10.1158/0008-5472.CAN-19-3659. PubMed PMID:
655 32245796; PMCID: PMC7202960.
- 656 19. Romero-Cordoba SL, Salido-Guadarrama I, Rebollar-Vega R, Bautista-Pina V,
657 Dominguez-Reyes C, Tenorio-Torres A, Villegas-Carlos F, Fernandez-Lopez JC, Uribe-
658 Figueroa L, Alfaro-Ruiz L, Hidalgo-Miranda A. Comprehensive omic characterization of breast
659 cancer in Mexican-Hispanic women. *Nat Commun.* 2021;12(1):2245. Epub 20210414. doi:
660 10.1038/s41467-021-22478-5. PubMed PMID: 33854067; PMCID: PMC8046804.
- 661 20. Chen Y, Lun AT, Smyth GK. From reads to genes to pathways: differential expression
662 analysis of RNA-Seq experiments using Rsubread and the edgeR quasi-likelihood pipeline.
663 *F1000Res.* 2016;5:1438. Epub 20160620. doi: 10.12688/f1000research.8987.2. PubMed PMID:
664 27508061; PMCID: PMC4934518.
- 665 21. Zhao X, Rodland EA, Tibshirani R, Plevritis S. Molecular subtyping for clinically defined
666 breast cancer subgroups. *Breast Cancer Res.* 2015;17:29. Epub 20150226. doi:
667 10.1186/s13058-015-0520-4. PubMed PMID: 25849221; PMCID: PMC4365540.
- 668 22. Li R, Qu H, Wang S, Wei J, Zhang L, Ma R, Lu J, Zhu J, Zhong WD, Jia Z.
669 GDCRNATools: an R/Bioconductor package for integrative analysis of lncRNA, miRNA and
670 mRNA data in GDC. *Bioinformatics.* 2018;34(14):2515-7. doi: 10.1093/bioinformatics/bty124.
671 PubMed PMID: 29509844.
- 672 23. Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, Grody WW, Hegde M, Lyon
673 E, Spector E, Voelkerding K, Rehm HL, Committee ALQA. Standards and guidelines for the
674 interpretation of sequence variants: a joint consensus recommendation of the American College
675 of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet Med.*
676 2015;17(5):405-24. Epub 20150305. doi: 10.1038/gim.2015.30. PubMed PMID: 25741868;
677 PMCID: PMC4544753.
- 678 24. Landrum MJ, Lee JM, Benson M, Brown GR, Chao C, Chitipiralla S, Gu B, Hart J,
679 Hoffman D, Jang W, Karapetyan K, Katz K, Liu C, Maddipatla Z, Malheiro A, McDaniel K,
680 Ovetsky M, Riley G, Zhou G, Holmes JB, Kattman BL, Maglott DR. ClinVar: improving access to
681 variant interpretations and supporting evidence. *Nucleic Acids Res.* 2018;46(D1):D1062-D7. doi:
682 10.1093/nar/gkx1153. PubMed PMID: 29165669; PMCID: PMC5753237.
- 683 25. Spear ML, Hu D, Pino-Yanes M, Huntsman S, Eng C, Levin AM, Ortega VE, White MJ,
684 McGarry ME, Thakur N, Galanter J, Mak ACY, Oh SS, Ampleford E, Peters SP, Davis A, Kumar
685 R, Farber HJ, Meade K, Avila PC, Serebrisky D, Lenoir MA, Brigino-Buenaventura E, Cintron
686 WR, Thyne SM, Rodriguez-Santana JR, Ford JG, Chapela R, Estrada AM, Sandoval K, Seibold
687 MA, Winkler CA, Bleecker ER, Myers DA, Williams LK, Hernandez RD, Torgerson DG,
688 Burchard EG. A genome-wide association and admixture mapping study of bronchodilator drug
689 response in African Americans with asthma. *Pharmacogenomics J.* 2019;19(3):249-59. Epub
690 20180912. doi: 10.1038/s41397-018-0042-4. PubMed PMID: 30206298; PMCID: PMC6414286.
- 691 26. Alexander DH, Lange K. Enhancements to the ADMIXTURE algorithm for individual
692 ancestry estimation. *BMC Bioinformatics.* 2011;12:246. Epub 20110618. doi: 10.1186/1471-
693 2105-12-246. PubMed PMID: 21682921; PMCID: PMC3146885.
- 694 27. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation
695 PLINK: rising to the challenge of larger and richer datasets. *Gigascience.* 2015;4:7. Epub
696 20150225. doi: 10.1186/s13742-015-0047-8. PubMed PMID: 25722852; PMCID: PMC4342193.
- 697 28. Van der Auwera GA, O'Connor BD. *enomics in the Cloud: Using Docker, GATK, and*
698 *WDL in Terra (1st Edition).* O'Reilly Media; 2020.
- 699 29. Tate JG, Bamford S, Jubb HC, Sondka Z, Beare DM, Bindal N, Boutselakis H, Cole CG,
700 Creatore C, Dawson E, Fish P, Harsha B, Hathaway C, Jupe SC, Kok CY, Noble K, Ponting L,
701 Ramshaw CC, Rye CE, Speedy HE, Stefancsik R, Thompson SL, Wang S, Ward S, Campbell
702 PJ, Forbes SA. COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res.*
703 2019;47(D1):D941-D7. doi: 10.1093/nar/gky1015. PubMed PMID: 30371878; PMCID:
704 PMC6323903.

- 705 30. Lawrence MS, Stojanov P, Polak P, Kryukov GV, Cibulskis K, Sivachenko A, Carter SL,
706 Stewart C, Mermel CH, Roberts SA, Kiezun A, Hammerman PS, McKenna A, Drier Y, Zou L,
707 Ramos AH, Pugh TJ, Stransky N, Helman E, Kim J, Sougnez C, Ambrogio L, Nickerson E,
708 Shefler E, Cortes ML, Auclair D, Saksena G, Voet D, Noble M, DiCara D, Lin P, Lichtenstein L,
709 Heiman DI, Fennell T, Imielinski M, Hernandez B, Hodis E, Baca S, Dulak AM, Lohr J, Landau
710 DA, Wu CJ, Melendez-Zajgla J, Hidalgo-Miranda A, Koren A, McCarroll SA, Mora J, Crompton
711 B, Onofrio R, Parkin M, Winckler W, Ardlie K, Gabriel SB, Roberts CWM, Biegel JA, Stegmaier
712 K, Bass AJ, Garraway LA, Meyerson M, Golub TR, Gordenin DA, Sunyaev S, Lander ES, Getz
713 G. Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature*.
714 2013;499(7457):214-8. Epub 20130616. doi: 10.1038/nature12213. PubMed PMID: 23770567;
715 PMCID: PMC3919509.
- 716 31. Shen R, Seshan VE. FACETS: allele-specific copy number and clonal heterogeneity
717 analysis tool for high-throughput DNA sequencing. *Nucleic Acids Res*. 2016;44(16):e131. Epub
718 20160607. doi: 10.1093/nar/gkw520. PubMed PMID: 27270079; PMCID: PMC5027494.
- 719 32. Mermel CH, Schumacher SE, Hill B, Meyerson ML, Beroukhi R, Getz G. GISTIC2.0
720 facilitates sensitive and confident localization of the targets of focal somatic copy-number
721 alteration in human cancers. *Genome Biol*. 2011;12(4):R41. Epub 20110428. doi: 10.1186/gb-
722 2011-12-4-r41. PubMed PMID: 21527027; PMCID: PMC3218867.
- 723 33. Ciriello G, Gatza ML, Beck AH, Wilkerson MD, Rhie SK, Pastore A, Zhang H, McLellan
724 M, Yau C, Kandoth C, Bowlby R, Shen H, Hayat S, Fieldhouse R, Lester SC, Tse GM, Factor
725 RE, Collins LC, Allison KH, Chen YY, Jensen K, Johnson NB, Oesterreich S, Mills GB,
726 Cherniack AD, Robertson G, Benz C, Sander C, Laird PW, Hoadley KA, King TA, Network TR,
727 Perou CM. Comprehensive Molecular Portraits of Invasive Lobular Breast Cancer. *Cell*.
728 2015;163(2):506-19. doi: 10.1016/j.cell.2015.09.033. PubMed PMID: 26451490; PMCID:
729 PMC4603750.
- 730 34. Blokzijl F, Janssen R, van Boxtel R, Cuppen E. MutationalPatterns: comprehensive
731 genome-wide analysis of mutational processes. *Genome Med*. 2018;10(1):33. Epub 20180425.
732 doi: 10.1186/s13073-018-0539-0. PubMed PMID: 29695279; PMCID: PMC5922316.
- 733 35. Weitzel JN, Neuhausen SL, Adamson A, Tao S, Ricker C, Maoz A, Rosenblatt M,
734 Nehoray B, Sand S, Steele L, Unzeitig G, Feldman N, Blanco AM, Hu D, Huntsman S, Castillo
735 D, Haiman C, Slavin T, Ziv E. Pathogenic and likely pathogenic variants in PALB2, CHEK2, and
736 other known breast cancer susceptibility genes among 1054 BRCA-negative Hispanics with
737 breast cancer. *Cancer*. 2019;125(16):2829-36. Epub 20190617. doi: 10.1002/cncr.32083.
738 PubMed PMID: 31206626; PMCID: PMC7376605.
- 739 36. Lorenzato A, Olivero M, Patane S, Rosso E, Oliaro A, Comoglio PM, Di Renzo MF.
740 Novel somatic mutations of the MET oncogene in human carcinoma metastases activating cell
741 motility and invasion. *Cancer Res*. 2002;62(23):7025-30. PubMed PMID: 12460923.
- 742 37. Alexandrov LB, Kim J, Haradhvala NJ, Huang MN, Tian Ng AW, Wu Y, Boot A,
743 Covington KR, Gordenin DA, Bergstrom EN, Islam SMA, Lopez-Bigas N, Klimczak LJ,
744 McPherson JR, Morganella S, Sabarinathan R, Wheeler DA, Mustonen V, Group PMSW, Getz
745 G, Rozen SG, Stratton MR, Consortium P. The repertoire of mutational signatures in human
746 cancer. *Nature*. 2020;578(7793):94-101. Epub 20200205. doi: 10.1038/s41586-020-1943-3.
747 PubMed PMID: 32025018; PMCID: PMC7054213.
- 748 38. Kanu N, Cerone MA, Goh G, Zalmas LP, Bartkova J, Dietzen M, McGranahan N, Rogers
749 R, Law EK, Gromova I, Kschischo M, Walton MI, Rossanese OW, Bartek J, Harris RS,
750 Venkatesan S, Swanton C. DNA replication stress mediates APOBEC3 family mutagenesis in
751 breast cancer. *Genome Biol*. 2016;17(1):185. Epub 20160915. doi: 10.1186/s13059-016-1042-
752 9. PubMed PMID: 27634334; PMCID: PMC5025597.
- 753 39. Komatsu A, Nagasaki K, Fujimori M, Amano J, Miki Y. Identification of novel deletion
754 polymorphisms in breast cancer. *Int J Oncol*. 2008;33(2):261-70. PubMed PMID: 18636146.

- 755 40. Caval V, Suspene R, Shapira M, Vartanian JP, Wain-Hobson S. A prevalent cancer
756 susceptibility APOBEC3A hybrid allele bearing APOBEC3B 3'UTR enhances chromosomal
757 DNA damage. *Nat Commun.* 2014;5:5129. Epub 20141009. doi: 10.1038/ncomms6129.
758 PubMed PMID: 25298230.
- 759 41. Xuan D, Li G, Cai Q, Deming-Halverson S, Shrubsole MJ, Shu XO, Kelley MC, Zheng
760 W, Long J. APOBEC3 deletion polymorphism is associated with breast cancer risk among
761 women of European ancestry. *Carcinogenesis.* 2013;34(10):2240-3. Epub 20130528. doi:
762 10.1093/carcin/bgt185. PubMed PMID: 23715497; PMCID: PMC3786378.
- 763 42. Nik-Zainal S, Davies H, Staaf J, Ramakrishna M, Glodzik D, Zou X, Martincorena I,
764 Alexandrov LB, Martin S, Wedge DC, Van Loo P, Ju YS, Smid M, Brinkman AB, Morganella S,
765 Aure MR, Lingjaerde OC, Langerod A, Ringner M, Ahn SM, Boyault S, Brock JE, Broeks A,
766 Butler A, Desmedt C, Dirix L, Dronov S, Fatima A, Foekens JA, Gerstung M, Hooijer GK, Jang
767 SJ, Jones DR, Kim HY, King TA, Krishnamurthy S, Lee HJ, Lee JY, Li Y, McLaren S, Menzies
768 A, Mustonen V, O'Meara S, Pauporte I, Pivot X, Purdie CA, Raine K, Ramakrishnan K,
769 Rodriguez-Gonzalez FG, Romieu G, Sieuwerts AM, Simpson PT, Shepherd R, Stebbings L,
770 Stefansson OA, Teague J, Tommasi S, Treilleux I, Van den Eynden GG, Vermeulen P, Vincent-
771 Salomon A, Yates L, Caldas C, van't Veer L, Tutt A, Knappskog S, Tan BK, Jonkers J, Borg A,
772 Ueno NT, Sotiriou C, Viari A, Futreal PA, Campbell PJ, Span PN, Van Laere S, Lakhani SR,
773 Eyfjord JE, Thompson AM, Birney E, Stunnenberg HG, van de Vijver MJ, Martens JW,
774 Borresen-Dale AL, Richardson AL, Kong G, Thomas G, Stratton MR. Landscape of somatic
775 mutations in 560 breast cancer whole-genome sequences. *Nature.* 2016;534(7605):47-54. Epub
776 20160502. doi: 10.1038/nature17676. PubMed PMID: 27135926; PMCID: PMC4910866.
- 777 43. Song J, Yang W, Shih le M, Zhang Z, Bai J. Identification of BCOX1, a novel gene
778 overexpressed in breast cancer. *Biochim Biophys Acta.* 2006;1760(1):62-9. Epub 20051025.
779 doi: 10.1016/j.bbagen.2005.09.017. PubMed PMID: 16289875.
- 780 44. Liu T, Zhang XY, He XH, Geng JS, Liu Y, Kong DJ, Shi QY, Liu F, Wei W, Pang D. High
781 levels of BCOX1 expression are associated with poor prognosis in patients with invasive ductal
782 carcinomas of the breast. *PLoS One.* 2014;9(1):e86952. Epub 20140128. doi:
783 10.1371/journal.pone.0086952. PubMed PMID: 24489812; PMCID: PMC3904964.
- 784 45. Zhong Z, Pannu V, Rosenow M, Stark A, Spetzler D. KIAA0100 Modulates Cancer Cell
785 Aggression Behavior of MDA-MB-231 through Microtubule and Heat Shock Proteins. *Cancers*
786 (Basel). 2018;10(6). Epub 20180604. doi: 10.3390/cancers10060180. PubMed PMID:
787 29867023; PMCID: PMC6025110.
- 788 46. Thompson AM, Moulder-Thompson SL. Neoadjuvant treatment of breast cancer. *Ann*
789 *Oncol.* 2012;23 Suppl 10:x231-6. doi: 10.1093/annonc/mds324. PubMed PMID: 22987968;
790 PMCID: PMC6278992.
- 791 47. Zavala VA, Bracci PM, Carethers JM, Carvajal-Carmona L, Coggins NB, Cruz-Correa
792 MR, Davis M, de Smith AJ, Dutil J, Figueiredo JC, Fox R, Graves KD, Gomez SL, Llera A,
793 Neuhausen SL, Newman L, Nguyen T, Palmer JR, Palmer NR, Perez-Stable EJ, Piawah S,
794 Rodriguez EJ, Sanabria-Salas MC, Schmit SL, Serrano-Gomez SJ, Stern MC, Weitzel J, Yang
795 JJ, Zabaleta J, Ziv E, Fejerman L. Cancer health disparities in racial/ethnic minorities in the
796 United States. *Br J Cancer.* 2021;124(2):315-32. Epub 20200909. doi: 10.1038/s41416-020-
797 01038-6. PubMed PMID: 32901135; PMCID: PMC7852513.
- 798