

RESEARCH

Open Access

Ancestry-specific predisposing germline variants in cancer



Ninad Oak^{1,2}, Andrew D. Cherniack^{3,4}, R. Jay Mashl^{5,6}, TCGA Analysis Network, Fred R. Hirsch⁷, Li Ding^{5,6,8,9}, Rameen Beroukhi^{3,4,10}, Zeynep H. Gümüş^{11,12}, Sharon E. Plon^{2,13,14} and Kuan-lin Huang^{11,12,15*} 

Abstract

Background: Distinct prevalence of inherited genetic predisposition may partially explain the difference of cancer risks across ancestries. Ancestry-specific analyses of germline genomes are required to inform cancer genetic risk and prognosis of diverse populations.

Methods: We conducted analyses using germline and somatic sequencing data generated by The Cancer Genome Atlas. Collapsing pathogenic and likely pathogenic variants to cancer predisposition genes (CPG), we analyzed the association between CPGs and cancer types within ancestral groups. We also identified the predisposition-associated two-hit events and gene expression effects in tumors.

Results: Genetic ancestry analysis classified the cohort of 9899 cancer cases into individuals of primarily European ($N = 8184$, 82.7%), African ($N = 966$, 9.8%), East Asian ($N = 649$, 6.6%), South Asian ($N = 48$, 0.5%), Native/Latin American ($N = 41$, 0.4%), and admixed ($N = 11$, 0.1%) ancestries. In the African ancestry, we discovered a potentially novel association of *BRCA2* in lung squamous cell carcinoma (OR = 41.4 [95% CI, 6.1–275.6]; FDR = 0.002) previously identified in Europeans, along with a known association of *BRCA2* in ovarian serous cystadenocarcinoma (OR = 8.5 [95% CI, 1.5–47.4]; FDR = 0.045). In the East Asian ancestry, we discovered one previously known association of *BRIP1* in stomach adenocarcinoma (OR = 12.8 [95% CI, 1.8–90.8]; FDR = 0.038). Rare variant burden analysis further identified 7 suggestive associations in African ancestry individuals previously described in European ancestry, including *SDHB* in pheochromocytoma and paraganglioma, *ATM* in prostate adenocarcinoma, *VHL* in kidney renal clear cell carcinoma, *FH* in kidney renal papillary cell carcinoma, and *PTEN* in uterine corpus endometrial carcinoma. Most predisposing variants were found exclusively in one ancestry in the TCGA and gnomAD datasets. Loss of heterozygosity was identified for 7 out of the 15 African ancestry carriers of predisposing variants. Further, tumors from the *SDHB* or *BRCA2* carriers showed simultaneous allelic-specific expression and low gene expression of their respective affected genes, and *FH* splice-site variant carriers showed mis-splicing of *FH*.

Conclusions: While several CPGs are shared across patients, many pathogenic variants are found to be ancestry-specific and trigger somatic effects. Studies using larger cohorts of diverse ancestries are required to pinpoint ancestry-specific genetic predisposition and inform genetic screening strategies.

* Correspondence: kuan-lin.huang@mssm.edu

¹¹Department of Genetics and Genomics, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA

¹²Icahn Institute for Data Science and Genomic Technology, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA

Full list of author information is available at the end of the article



© The Author(s). 2020 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

Background

Cancer risk differs across ancestries. According to the National Cancer Institute's Surveillance, Epidemiology, and End Results (SEER) program, the cancer incidence per 100,000 ranges from 449 in race/ethnicity population self-identified as Whites, 453 in Blacks, 298 in Asian/Pacific Islanders, 315 in American Indian/Alaskan Natives, and 336 in Hispanics in the USA between 2011 and 2015 [1, 2]. While some of these differences may be attributed to non-genetic factors such as access to health care or diet, much can likely be explained by differences in the genomic architecture of these ancestries and differing frequencies of inherited genetic predisposition. Previous studies revealed different carrier rates of pathogenic variants across ancestries, albeit often in a limited panel of genes or selected cancer types [3–5].

While multiple large-scale genome-wide association studies have investigated the common risk variants contributing to cancer [6–10], fewer studies have interrogated rare pathogenic variants in non-European ancestries [5, 11–15]. A 2019 systematic review of cancer sequencing studies found a total of only 764 reported non-European (minority) cases in 27 published studies with reported race/ethnicity [9]. Consequently, germline genetic testing in non-White patients often results in higher rates of variants of unknown significance (VUSs) [16]. Ongoing efforts are bridging the knowledge gap of cancer genetic predisposition in under-studied populations [17–19]. Meanwhile, systematic cross-ancestry investigations of predisposing variants across cancer types are urgently needed to inform genetic testing for each ancestral group.

Herein, we analyzed germline variant data of 9899 cancer cases across 33 cancer types from the Cancer Genome Atlas Project (TCGA) [20] to identify ancestry-specific cancer-gene associations where the genes show an excess of pathogenic/likely pathogenic germline variants the TCGA samples. In samples of African ancestry, we identified two associations, *BRCA2* in lung squamous cell carcinoma (LUSC) and ovarian serous cystadenocarcinoma (OV). In analyses of individuals with East Asian ancestry, we identified an association for *BRIP1* in stomach adenocarcinoma (STAD). Using a rare-variant association analysis, we identified seven additional suggestive cancer gene associations. Evidence of a somatic second hit event (i.e., loss of heterozygosity [LOH] or a biallelic mutation) was found in two thirds of the tumors with germline predisposing variants. Many carriers of ancestry-specific predisposition variants showed altered expression of the affected genes, including allelic-specific expression (ASE), mis-splicing, and reduced tumor suppressor gene expression, further supporting these genetic variants' contribution to cancer predisposition.

Methods

Study cohort and genetic ancestry assignment

We used the clinical data provided by TCGA PanCanAtlas and restricted analyses to those with pass-QC blood/normal sequencing data. In addition to excluding cases with PanCanAtlas blacklisted germline BAM-files, cases with less than 60% genotype concordance between sequencing variant calls and SNP-genotype data were eliminated, where 10,389 cases were left [20]. We further overlapped with the cases included in the PanCanAtlas Ancestry Informative Markers (AIM) genetic ancestry assignment, resulting in the final set of 9899 samples. The detailed descriptions of ancestry assignment procedures are available in the marker publication [21].

Briefly, consensus genetic ancestry for each TCGA case was determined as the majority of ancestry assignments that were independently determined by five methods across four institutions. These methods include those based on SNP-array genotypes used by Broad Institute, University of California San Francisco (UCSF), and Washington University (WashU), as well as those based on whole-exome sequencing data used by University of Trento and ExAC/Broad Institute. The five methods conducted variations of principal component analyses (PCA) on TCGA normal samples to infer genetic ancestry. We further provide the PCA plots showing the alignment of the major PCs in the UCSF and WashU analyses with the AIM-group consensus genetic ancestry in Additional file 1: Fig. S1.

For each sample, the percentage of global ancestry of African, European, East Asian, Native/Latin American, and South Asian ($k = 5$) was further estimated using ADMIXTURE [22] version 1.23 based on the common SNP markers (1000 genomes allele frequency (AF) > 1%) in the Broad Institute analysis. Samples with the proportion of the secondary ancestry greater than 20% were considered as admixed samples (Additional file 2: Table S1). Sensitivity analyses revealed increased power by including admix samples in this cohort. Thus, cases with admixed ancestry assignments were grouped to their nearest neighbors (e.g., afr_admix to afr) for downstream analyses.

Pathogenic and likely pathogenic germline variant calls

We downloaded the overall and predisposing germline variant calls previously reported by the PanCanAtlas Germline Analyses Working Group (<https://gdc.cancer.gov/about-data/publications/PanCanAtlas-Germline-AWG>) [20]. The detailed description of variant calling and classification procedures are available in the TCGA PanCanAtlas germline publication [20].

Briefly, germline SNVs were identified using the union of variant calls between Varscan [23] and GATK [24]. Germline indels were identified using Varscan, GATK,

and Pindel [25], and we only retained variants called by at least two out of the three callers or high-confidence Pindel-unique calls (at least 30× coverage and 20% variant allele fraction [VAF]). We used the GRCh37-lite reference. We further required the variants to have an allelic depth (AD) ≥ 5 for the alternative allele. We then used bam-readcount to quantify the number of reference and alternative alleles in both normal and tumor samples. We required the variants to have at least 5 counts of the alternative allele and an alternative allele frequency of at least 20%. Of these, we included those rare variants with $\leq 0.05\%$ allele frequency in 1000 Genomes and ExAC (release r0.3.1). We subsequently retained only cancer-relevant pathogenic variants, based on whether they were found in the curated cancer variant databases or a 152 curated cancer predisposing gene list. Finally, we manually reviewed all variants using integrative genomics viewer (IGV) and filtered out variants with poor support sequence reads.

The variants defined by the above pipeline were then classified using an automatic pipeline termed CharGer [26] (<https://github.com/ding-lab/CharGer>) that adopts the American College of Medical Genetics and Genomics/Association of Molecular Pathology (ACMG/AMP) variant classification guidelines which are designed for assessment of germline variants in Mendelian disorders [27]. For the CharGer classification pipeline, we defined 12 pathogenic evidence levels and 4 benign evidence levels using a number of datasets, including ExAC and ClinVar. The pathogenic evidence adds points, whereas benign evidence subtracts points that amount to pathogenicity (pathogenic requires the variant to be described as pathogenic by the reviewed clinical significance in ClinVar (not including variants showing “conflicting interpretations of pathogenicity”) or other cancer predisposition gene databases, likely pathogenic requires CharGer score > 8). To acquire enough CharGer points to be classified as likely pathogenic, the variants typically need to be predicted to result in truncation in cancer predisposition genes where the loss of function (LOF) is a known disease mechanism and harbor variants with a dominant (evidence level PVS1, + 8 points) or a recessive (evidence level PSC1, + 4 points) mode of inheritance. Additionally, evidence level PS1, + 7 points are scored if the variant results in the same peptide sequence change as an established pathogenic variant. All other modules will each add ≤ 2 points.

Principal component analysis (PCA)

Birdseed genotype files were downloaded from Genomic Data Commons (GDC) in the legacy (hg19) archive onto Institute for System Biology-Cancer Genome Cloud (ISB-CGC), converted to individual VCF files, and then merged into a combined VCFs containing 11,459

samples and 522,606 variants. We conducted PCA as implemented by PLINK (v1.9) [28]. Specifically, we retained 298,004 variants with AF > 0.15 for population structure analysis. The resulting eigenvalues and eigenvectors were then recorded. PC1 and PC2 accounted for 51.6% and 29.2% of the variations across the first 20 PCs, and none of the trailing PCs accounted for more than 3.2%. Thus, we subsequently controlled for PC1 and PC2 in ancestry-specific cancer predisposing gene analysis (Additional file 1: Fig. S1).

Multivariate regression to identify the enrichment of pathogenic variants

For each cancer type within each ancestry, we conducted multivariate logistic regression analyses considering the case status of the cancer type as the dependent variable (using all other cancer cohorts as controls) and the carrier status of each predisposing gene as an independent variable. The model corrected for age at the initial pathologic diagnosis, gender, and the first two principal components (accounted for 80.8% variations across the first 20 PCs). All ancestry cohorts are called using the same variant calling pipeline, thus avoiding the potential danger of comparing this population against other cohorts such as ExAC. We collapsed predisposing (pathogenic and likely pathogenic) germline variants to the gene level. Only ancestry-cancer combinations with at least 20 cases and predisposing genes with at least two individuals with predisposing variants within the cohort are tested. In total, we tested 33 cancers in European Ancestry, 15 cancers in African Ancestry, and 8 cancers in East Asian ancestry that met this criterion. No cohorts of the Native/Latin American and South Asian ancestry have sufficient sample sizes in TCGA for testing. Among these tested cancers, we tested a total of 114 cancer-gene combinations for multivariate regression analysis, of which 101 were within European ancestry, 9 were in African ancestry, and 4 were in East Asian ancestry. *P* values were calculated using the Wald test and adjusted to FDR using the standard Benjamini-Hochberg procedure.

Burden testing of pathogenic variants

We conducted burden testing of the cohort within each ancestry as defined by the TCGA AIM working group. Specifically, we adopted the Total Frequency Test (TFT) [29] by collapsing predisposing (pathogenic and likely pathogenic) germline variants to the gene level. For each cancer type with at least 20 cases of the tested ancestry with at least one predisposing variant carrier, we tested the burden of predisposing variants for each gene against all other cancer cohorts as controls. Among the cancers that met the sample size criteria described above, we tested a total of 120 cancer-gene combinations using

rare variant burden testing, of which 104 were within European ancestry, 11 were in African ancestry, and 5 were in East Asian ancestry. The resulting *P* values were adjusted to FDR using the standard Benjamini-Hochberg procedure.

gnomAD analysis

We analyzed the gene-level and variant-level frequency of the identified genetic predisposition using the non-cancer subset of the genome aggregation database (gnomAD-non-cancer) cohort (118,479 WES and 15,708 WGS samples) [30, 31] (<http://gnomad.broadinstitute.org>). For the gene-level analysis, we retained rare variants with ancestry-specific minor allele frequency < 0.5%. We further retained pathogenic and likely pathogenic variants per ACMG/AMP criteria as ascertained by InterVar [32] and annotated using ANNOVAR [33]. Allele frequencies were summarized at gene-level within each sub-population in gnomAD using total allele counts and maximum allele numbers within each group.

The lollipop diagrams in Fig. 2 were constructed and modified using protein paint (<https://pecan.stjude.cloud/proteinpaint>).

Expression analysis

TCGA level-3 normalized RNA expression data were downloaded from Firehose (2016/1/28 analysis archive). The tumor expression percentile of individual genes in each cancer cohort was calculated using the empirical cumulative distribution function (ecdf), as implemented in R. We annotated germline carriers of predisposition variants with extreme mRNA tumor expression (> 80th or < 20th percentile) of the affected gene. For samples within the same ancestry and same cancer cohort, we then used the two-sample Kolmogorov-Smirnov test to compare the expression percentile distribution between variants of oncogenes and tumor suppressors. The resulting *P* values were adjusted to false discovery rate (FDR) using the standard Benjamini-Hochberg procedure.

For the ancestry-specific variants, we recorded the RNA VAF of the mutant allele in the RNA-Seq bam files. For splice site variants, we assessed the mis-splicing of the transcript and variants using IGV.

Power and downsampling analysis

Post hoc power analyses were performed using R-package SKAT [34] and the `power_logistic` function to calculate the number of samples for rare variant association with causal percentage = 80%, minor allele frequency < 0.1%, and using odds ratio (OR) > 1 through OR < 10. Each calculation was performed using 100 simulations over a target 5 kb region.

Additionally, we performed a downsampling analysis for each tumor type by random sampling of subsets of samples with incremental sizes from zero to the total number of samples in that tumor type. We identified the number of significantly mutated genes as described above within each subset and plotted a smoothed function (loess method) against the subset size. Each calculation was performed at ten iterations (Additional file 1: Fig. S2).

Results

Ancestry demographics of TCGA cohort

We classified the 9899 TCGA cases with pass-QC germline data across 33 cancer types by genotype-defined ancestries defined by the PanCanAtlas Ancestry Informative Markers (AIM) working group (Additional file 1: Fig. S1, the “Methods” section, Table 1). The European ancestry contained 82.68% (*n* = 8184) of individuals in this cohort. The remainder of the cohort consisted of 9.76% (*n* = 966) African ancestry, 6.56% (*n* = 649) East Asian ancestry, 0.48% (*n* = 48) South Asian ancestry, 0.41% (*n* = 41) Native/Latin American ancestry, and 0.11% (*n* = 11) mixed ancestry. The largest ancestry-specific tumor cohorts are breast invasive carcinoma (BRCA) for the European ancestry (*n* = 811) and African ancestry (*n* = 180), liver hepatocellular carcinoma (LIHC) for the East Asian ancestry (*n* = 162), and thyroid carcinoma (THCA) for the Native/Latin American ancestry (*n* = 11) and the South Asian ancestry (*n* = 11).

Ancestry-specific cancer predisposing genes

Acknowledging the limited power to assess ancestry-specific associations as shown by the post hoc power analyses (Additional file 1: Fig. S2), we sought to identify cancer predisposing genes within each ancestry. We considered cancer predisposing genes as those statistically enriched for pooled pathogenic and likely pathogenic variants (referred to here as predisposing variants) as previously classified [20]). For each ancestry-cancer type pair, we conducted multivariate regression analyses correcting for onset age, gender, and the first two principal components.

Along with 36 cancer-gene associations (FDR < 0.05, Wald test) found in the European ancestry, we identified two specific cancer-gene associations in the African ancestry: *BRCA2* in ovarian cancer (OV) (OR = 8.5 [95% CI, 1.5–47.4]; FDR = 0.045) and *LUSC* (OR = 41.4 [95% CI, 6.1–275.6]; FDR = 0.002). We also identified one association in the East Asian ancestry, *BRIP1* in STAD (OR = 12.8 [95% CI, 1.8–90.84]; FDR = 0.038) (Fig. 1, Additional file 2: Table S2a). While the association of *BRCA2* and *LUSC* is first described in African-American ancestry here, *BRCA2* was previously found to be associated with non-small cell lung cancer (including LUAD

Table 1 The demographic distribution of TCGA PanCanAtlas cohort

	Cancer	European	African	Native/Latin American	East Asian	South Asian	Admix	Total
Total Number (%)		8184 (82.68%)	966 (9.76%)	41 (0.41%)	649 (6.56%)	48 (0.48%)	11 (0.11%)	9899 (100%)
ACC	Adrenocortical Carcinoma	82	2	0	2	0	0	86
BLCA	Bladder Urothelial Carcinoma	329	21	1	43	2	0	396
BRCA	Breast Invasive Carcinoma	811	180	5	53	8	1	1058
CESC	Cervical Squamous Cell Carcinoma and Endocervical Adenocarcinoma	208	31	1	22	0	0	262
CHOL	Cholangiocarcinoma	30	2	0	2	0	0	34
COAD	Colon Adenocarcinoma	343	58	0	12	0	0	413
DLBC	Lymphoid Neoplasm Diffuse Large B-cell Lymphoma	23	0	0	15	1	0	39
ESCA	Esophageal Carcinoma	126	4	0	44	0	0	174
GBM	Glioblastoma Multiforme	292	37	0	2	1	0	332
HNSC	Head and Neck Squamous Cell Carcinoma	439	50	6	7	5	1	508
KICH	Kidney Chromophobe	56	4	0	1	1	0	62
KIRC	Kidney Renal Clear Cell Carcinoma	306	54	2	7	1	0	370
KIRP	Kidney Renal Papillary Cell Carcinoma	207	63	0	6	1	0	277
LAML	Acute Myeloid Leukemia	125	14	0	2	0	0	141
LGG	Brain Lower Grade Glioma	455	23	4	10	3	0	495
LIHC	Liver Hepatocellular Carcinoma	179	18	0	162	1	1	361
LUAD	Lung Adenocarcinoma	416	52	1	9	0	0	478
LUSC	Lung Squamous Cell Carcinoma	455	29	0	11	0	0	495
MESO	Mesothelioma	78	0	0	0	1	0	79
OV	Ovarian Serous Cystadenocarcinoma	348	30	1	7	5	0	391
PAAD	Pancreatic Adenocarcinoma	158	9	0	11	0	1	179
PCPG	Pheochromocytoma and Paraganglioma	146	20	0	3	4	0	173
PRAD	Prostate Adenocarcinoma	411	59	0	9	2	1	482
READ	Rectum Adenocarcinoma	135	6	0	1	0	0	142
SARC	Sarcoma	217	18	0	6	0	2	243
SKCM	Skin Cutaneous Melanoma	448	1	2	12	0	0	463
STAD	Stomach Adenocarcinoma	294	15	0	90	0	0	399
TGCT	Testicular Germ Cell Tumors	107	4	0	3	0	0	114
THCA	Thyroid Carcinoma	359	33	11	52	11	3	469
THYM	Thymoma	96	8	0	12	0	0	116
UCEC	Uterine Corpus Endometrial Carcinoma	382	112	7	30	1	1	533
UCS	Uterine Carcinosarcoma	43	9	0	3	0	0	55
UVM	Uveal Melanoma	80	0	0	0	0	0	80

and LUSC) and ovarian cancer (OV) in the European ancestry [35–37]. The association of *BRIP1* predisposition to STAD in the East Asian ancestry was also

previously reported for the European ancestry [38]. These findings (including novel associations) in a large heterogeneous cancer population build on older studies

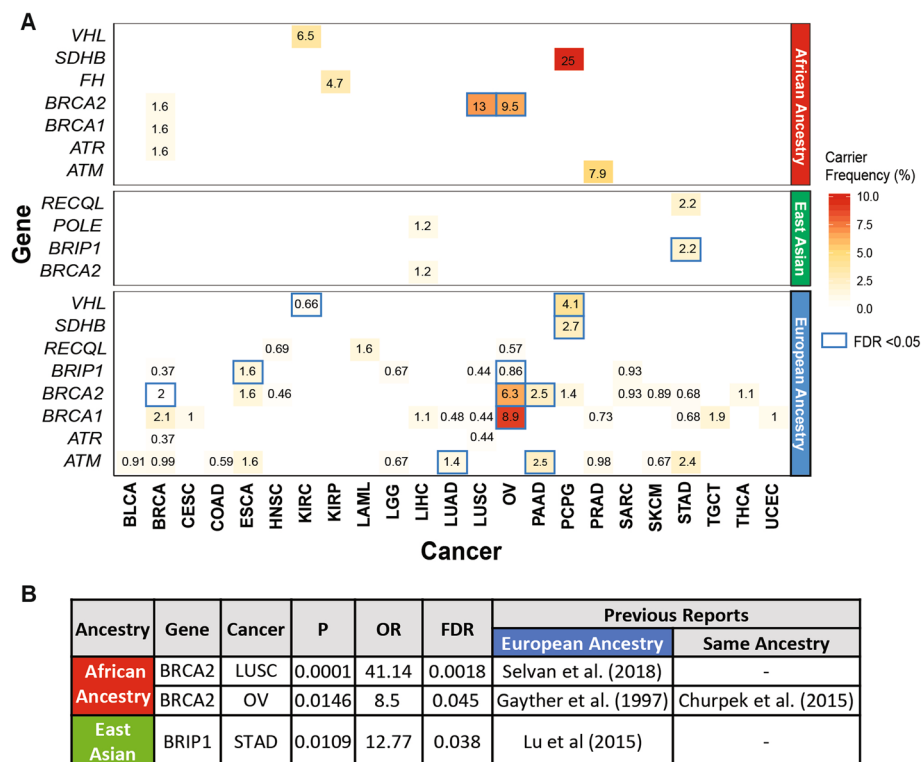


Fig. 1 Cancer predisposing genes identified in each ancestry across 9899 TCGA cases across cancer types in the African ancestry, East Asian, and European ancestries. **a** Ancestry-specific cancer-gene pairs from TCGA dataset containing cancer predisposing variants as identified by multivariate logistic regression analyses. Each number represents carrier frequencies of predisposing genes within that cancer cohort. Genes with significant associations (Wald test $FDR < 0.05$) are highlighted with blue boxes. **b** Significant cancer-predisposing gene associations ($FDR < 0.05$) identified in the African and East Asian ancestries

that evaluated individual cancer predisposition genes and cancer risk across ancestries.

The top associated predisposing genes and their carrier frequency vary widely across ancestries (Fig. 1a). For genes with a significant association in the African ancestry, we observed a higher carrier frequency compared to other ancestries. For example, in LUSC, *BRCA2* predisposing variants were found in 2 of the 29 African ancestry samples (6.9%), whereas we only found 1 *BRCA2* carrier out of the 455 European-ancestry samples (0.44%).

We next investigated whether the cross-ancestry differences in predisposing gene frequencies were also observed in other cohorts. Specifically, we examined the gene-level rates of individuals carrying pathogenic and likely pathogenic variants in the gnomAD non-cancer cohort [30, 31] (118,479 WES and 15,708 WGS samples, the “Methods” section, Additional file 2: Table S3). *BRCA2* showed the highest frequency in the African ancestry (0.072%) than all other defined ancestries, including non-Finnish European (0.048%) and East Asian (0.047%). *BRIP1* also showed higher frequency in the East Asian ancestry (0.068%) than all ancestries ($\leq 0.045\%$) except for the non-Finnish European ancestry (0.099%).

To generate hypotheses for future targeted studies, we investigated additional ancestry-implicated genes using total frequency testing (TFT) of predisposing variants, fully acknowledging potential confounders using this method (Additional file 2: Table S2b). We identified 7 suggestive ($FDR < 0.05$ in the TFT analysis) ancestry-specific cancer-gene associations in the African ancestry, 6 of which have been previously described including *SDHB* in PCPG [39], *ATM* in PRAD [40, 41], *FH* in KIRP [42], *VHL* in KIRC [43], *PTEN* in UCEC [44], and *BRCA2* in OV [12]. We also re-discovered the *BRCA2* in LUSC described above. In the East Asian ancestry, we identified 3 borderline-suggestive associations ($FDR = 0.32$): *RECQL* in STAD, *BRIP1* in STAD, and *POLE* in LIHC. In STAD, *RECQL* and *BRIP1* each affected 2 of the 90 East Asian ancestry cases, but none of the 294 European-ancestry cases. In LIHC, two protein-truncating variants were seen in *POLE* among 162 East Asian ancestry cases compared to none in 179 European-ancestry cases. These suggestive associations remain to be established and are only used to identify potential predisposing variants with supporting somatic evidence.

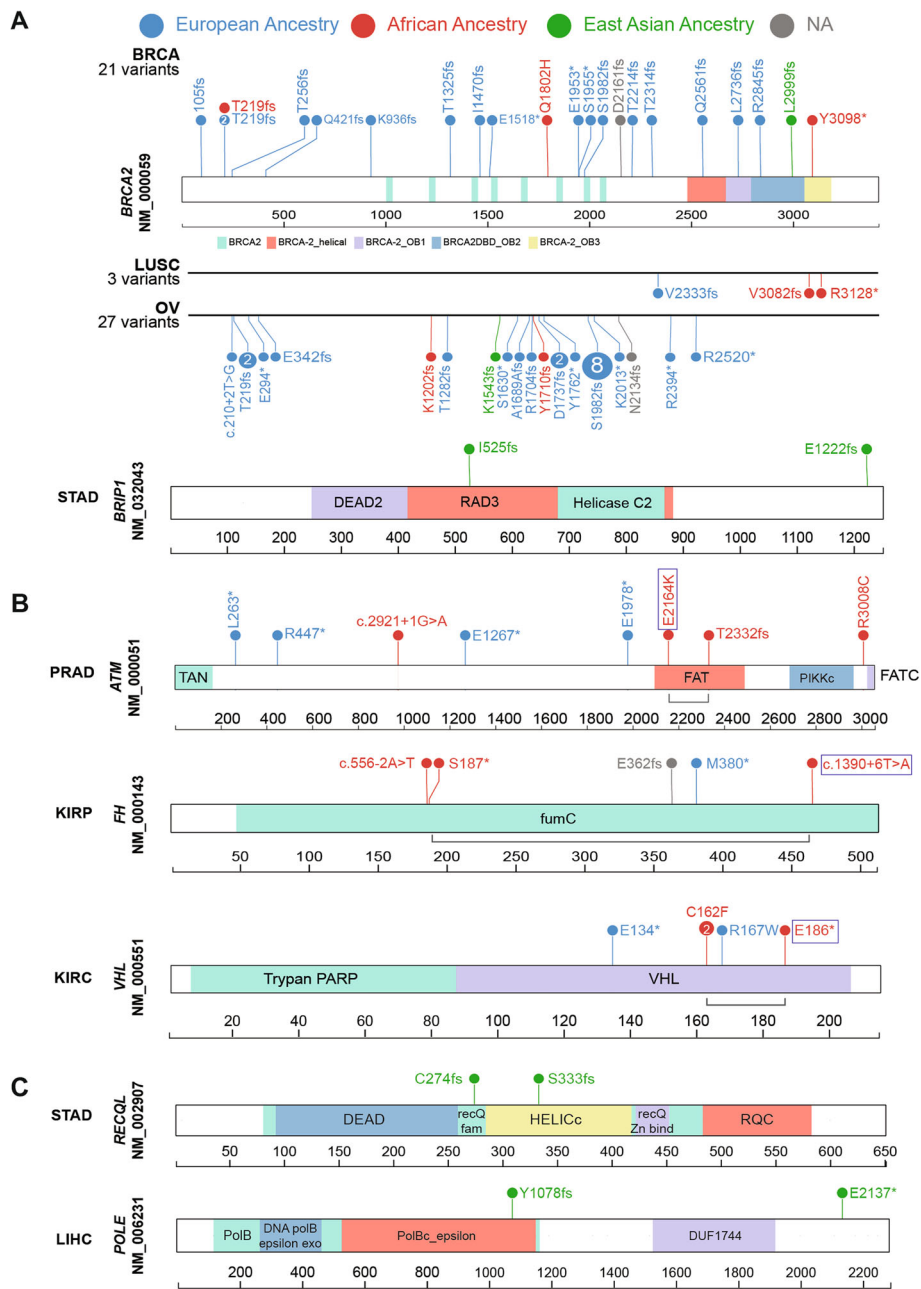


Fig. 2 Ancestry-specific predisposing germline variants. Predisposing variants in the significant (regression analysis; **a**) and suggestive (rare variant burden testing; **b**) cancer-gene associations are shown. The variants are labeled with carrier counts and colored by their respective carriers' ancestry (European Ancestry: blue, African ancestry: red, East Asian Ancestry: green). **a** Significant predisposing variants identified in the African and East Asian ancestries are shown across respective cancer types. For *BRCA2*, predisposing variants across all cancers are shown (top) in comparison with the two cancer types with significant associations in the African ancestry (LUSC and OV, bottom). Similarly, predisposing variants contributing to the significant association of *BRIP1* in STAD in the East Asian ancestry are shown. **b** Suggestive predisposing variants identified in the African ancestry are shown for *ATM*, *FH*, and *VHL* genes within their associated cancer types. Bi-allelic events in each carrier are linked by a grey line bracket where the somatic second-hit mutations are marked with a box. **c** Borderline-suggestive predisposing variants identified in the East Asian ancestry are shown for *RECQL* in STAD and *POLE* in LIHC

Ancestry-specific predisposing variants

We next examined ancestry-specific predisposition at the variant level (Fig. 2, Additional file 2: Table S4) for the 3 significant associations from the multivariate

logistic regression analyses and the 7 suggestive associations from the TFT analysis. The cancer-gene pairs included 15 predisposing variants within the African ancestry and another 6 within the East Asian ancestry.

None of the above variants discovered in the African ancestry were observed in any other ancestry within that cancer type (Fig. 2). Across the pan-cancer TCGA cohort, all of the *BRCA2* frameshift variants found in LUSC and OV were unique to the African ancestry. For other associated genes in the African ancestry, including *ATM* (PRAD), *FH* (KIRP), and *VHL* (KIRC), the predisposing variants differ between the African and European ancestries (Fig. 2b). The African ancestry-specific predisposing variants include splice site variants *ATM* c.2921+1G>A and *FH* c.556-2A>T, protein-truncating variants *ATM* p.T2333fs and *FH* p.S187*, and missense variants *ATM* p.R3008C. *VHL* p.C162F is the only recurrent variant found in two KIRC cases.

In the East Asian ancestry, we assessed predisposing variants in *BRIP1* (STAD), *POLE* (LIHC), and *RECQL* (STAD) (Fig. 2a and c). These include two *BRIP1* variants p.I525fs and p.E1222fs and two protein-truncating variants in *POLE* and *RECQL*, respectively. All six predisposing variants were not shared with any other ancestry in the TCGA cohort (Fig. 2c).

We further analyzed the frequency of these variants of the gnomAD non-cancer dataset [30, 31]. Among the African ancestry-specific predisposing variants, splice-site variant *ATM* c.2921+1G>A (African ancestry allelic count [AC]/total allele number [AN] = 1/14,878; allelic frequency [AF] = 0.0067%) and *BRCA2* p.R3128* (African ancestry AC/AN = 4/23,610; AF = 0.016%) were the only variants present in the African and non-Finnish European ancestries in gnomAD-non-cancer dataset. All other variants were absent within African ancestry and most other ancestries in gnomAD except *SDHB* p.R46* (Finnish European ancestry AC/AN = 2/25,066; AF = 0.007%) and *ATM* p.R3008C (East Asian ancestry AC/AN = 1/17,688; AF = 0.005%). Similarly, only two of the six East Asian ancestry-specific predisposing variants, *BRIP1* p.E1222Gfs (East Asian ancestry AC/AN = 11/19,232; AF = 0.05%) and *POLE* p.Tyr1078fs (East Asian ancestry AC/AN = 1/17,692; AF = 0.005%), were present exclusively in the East Asian ancestry of gnomAD-non-cancer dataset. Of note, 7 of the 15 predisposing variants, including *BRCA2* variants in OV (p.Y1710fs, p.K1202fs) and in LUSC (p.V3082fs), were not found in ClinVar [45]. While *VHL* p.C162F lacks a ClinVar record, the co-localizing p.C162W showed three reports of pathogenicity and one report of uncertain significance.

We also investigated the presence of the six predisposing variants in the East Asian ancestry from the gnomAD non-cancer dataset. Only the *POLE* p.Y1078fs (AC/AN = 1/17,692, AF = 0.0056%) and *BRIP1* p.E1222fs (AC/AN = 11/19,232, AF = 0.057%) were present exclusively in the East Asian ancestry of gnomAD-non-cancer dataset. All other East Asian-ancestry variants were not

detected in this dataset. Of note, none of the six variants were previously reported in ClinVar [45].

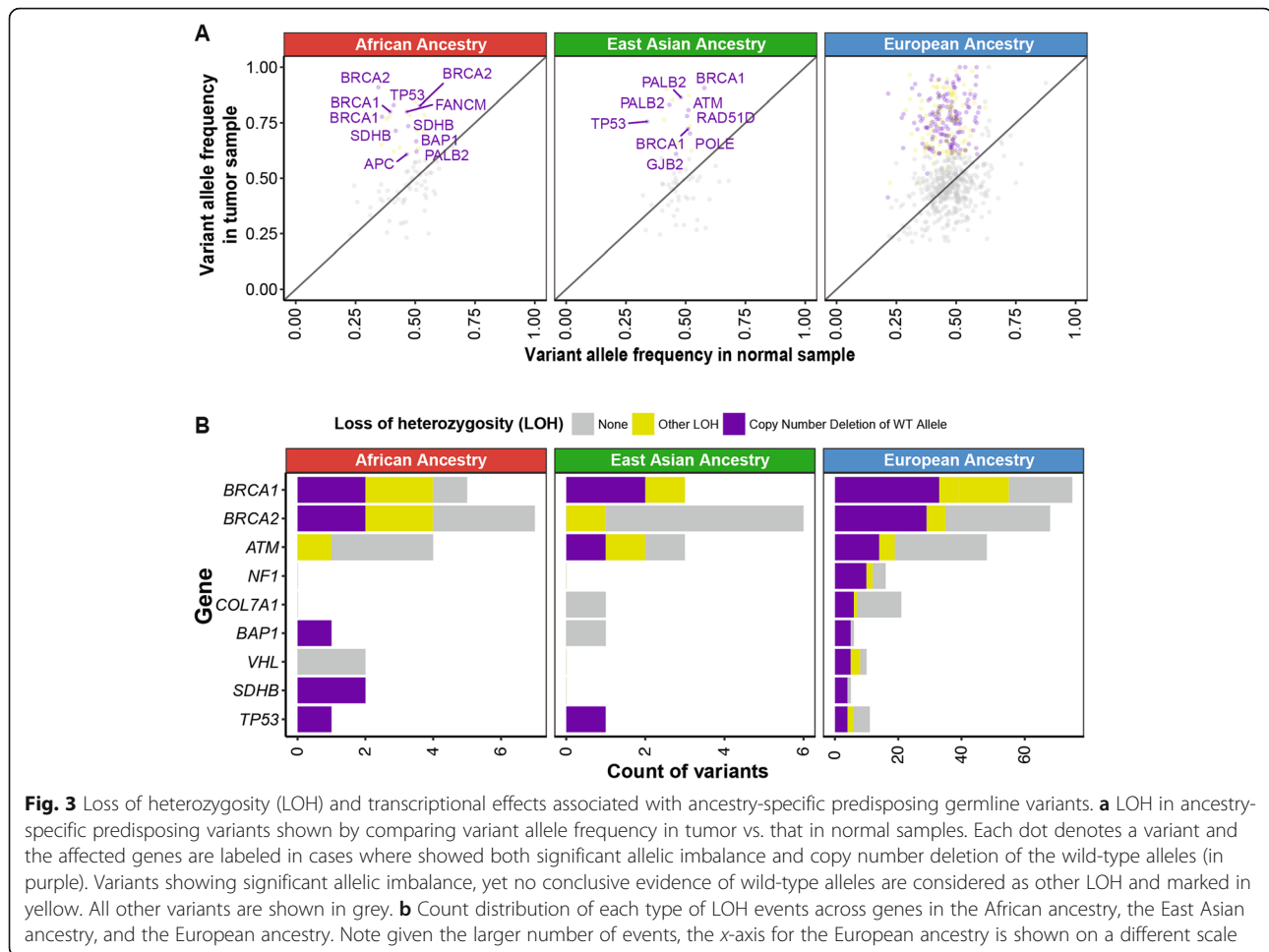
Germline-somatic two-hit events

We next examined the two-hit hypothesis, whereby a somatic second hit of the same gene is found in carriers of the germline predisposing variants [46, 47]. First, we investigated the extent of loss of heterozygosity (LOH) of the predisposing variants using our previously developed statistical test [38] (the “Methods” section) that compares the variant allele fractions in tumor vs. normal samples. Among the variants observed in the African ancestry, we observed significant LOH (FDR < 0.05) for both truncating variants in *SDHB* p.R116fs and p.R46* in PCPG (Fig. 3a). Three additional variants exhibited significant LOH, including *BRCA2* p.R3128* (LUSC), *BRCA2* p.K1202fs (OV), and *FH* p.S187* (KIRP). We also observed suggestive LOH (FDR < 0.15 or tumor VAF > 0.6) for *ATM* c.2921+1G>A (PRAD) and *BRCA2* p.Y1710fs (OV) (Fig. 3b). Among the six predisposing variants in the East Asian ancestry, only *POLE* p.E2137* (LIHC) showed significant LOH (Fig. 3a).

As an alternative mechanism of a somatic second hit, we identified three biallelic mutations where the rare germline predisposing variant was coupled with a second somatic mutation of the same gene, all found in African ancestry carriers (labeled in Fig. 2b, Additional file 2: Table S4b). In a PRAD carrier of *ATM*, the germline p.L2332fs variant was coupled with a somatic p.E2164K mutation; in the KIRC carrier of *VHL*, the germline p.C162F variant was coupled with somatic p.E186* mutation. In a KIRP carrier of *FH*, whose *FH* gene expression is low (Fig. 4a), germline p.S187* variant was coupled with a somatic splice-site mutation c.1390+6T>A. Analysis of RNA from the KIRP tumor revealed that the somatic *FH*: c.1390+6T>A causes missplicing of 27.6% of the transcripts in tumor RNA, as indicated by the number of reads spanning consensus splice site ($n = 68$) and the new cryptic splice site ($n = 26$) (case 2 in Fig. 4b). None of the six carriers of the predisposing variants in East Asian ancestry harbored a biallelic somatic mutation. Overall, the assessment of LOH and biallelic mutation supports the variants' contribution to oncogenesis through the two-hit model.

Expression changes in predisposing genes

To examine the transcriptional effects of the predisposing variants, we investigated the gene expression in tumor samples of the predisposing variant carriers (Fig. 4a). We observed 154 overall and 27 non-European ancestry-specific predisposing variants co-occurring with an extreme expression (> 80% or < 20% in the same cancer cohort) of the respective gene, although the current sample sizes preclude us from discovering significantly



associated genes compared to non-carriers within each ancestry-cancer cohort (Additional file 2: Table S5a).

All of the expression-associated variants were germline heterozygous variants at the DNA level. The degree of their variant allele fraction in the tumor RNAseq data (RNA VAF) thus indicates the degree of allelic-specific expression (ASE). The African carriers of *SDHB* truncating variants p.R116fs (the corresponding gene's expression ranks at the bottom 0.5 percentile among all PCPG cases [0.5%], RNA VAF = 0.25 and p.R46* (9% in PCGP, RNA VAF = 0.80) showed low *SDHB* expression. The African carriers of *BRCA2* p.Y1710fs (6% in OV, RNA VAF = 0) and p.3082fs (15% in LUSC, RNA VAF = 0) also exhibited low *BRCA2* (Fig. 4c). In the OV case, the germline *BRCA2* p.Y1710fs is coupled with a somatic LOH event, resulting in nearly complete loss of *BRCA2* expression.

Both of the African ancestry carriers of *FH* predisposing variants, *FH* p.S187* (2% in KIRP, RNA VAF = 0.13) and *FH:c.556-2A>T* (2% in KIRP, RNA VAF = 0.50), showed low *FH* expression. In addition to the biallelic somatic *FH:c.1390+6T>A* mutation in the carrier of

germline *FH* p.S187* described earlier, we also observed a mis-splicing event in a different case carrying germline *FH:c.556-2A>T* at the RNA level (case 1 in Fig. 4b).

For other ancestries, the tumor from one predisposing variant carrier of the Native/Latin American ancestry, *NF1* p.Y489C, showed low *NF1* mRNA expression (2% in BRCA, RNA VAF = 0). Overall, RNA VAF of the majority of protein-truncating variants not accompanied by LOH varied between 0 and 0.25 (Additional file 2: Table S5a), suggesting degradation of the mutant allele.

Many predisposing truncating variants of tumor suppressors are assumed to lead to loss of gene expression through mechanisms such as nonsense-mediated decay (NMD). Using the NMD Classifier [48], we revealed all frameshift variants found in the African and East Asian ancestries were located in the NMD-competent region (Additional file 1: Fig. S3). These results support that a fraction of predisposing variants likely result in reduced gene products of tumor suppressors in ancestral groups.

Conversely, for the rare tumors with germline variants in oncogenes, the two predisposing *RET* variants are coupled with elevated *RET* expression in their African

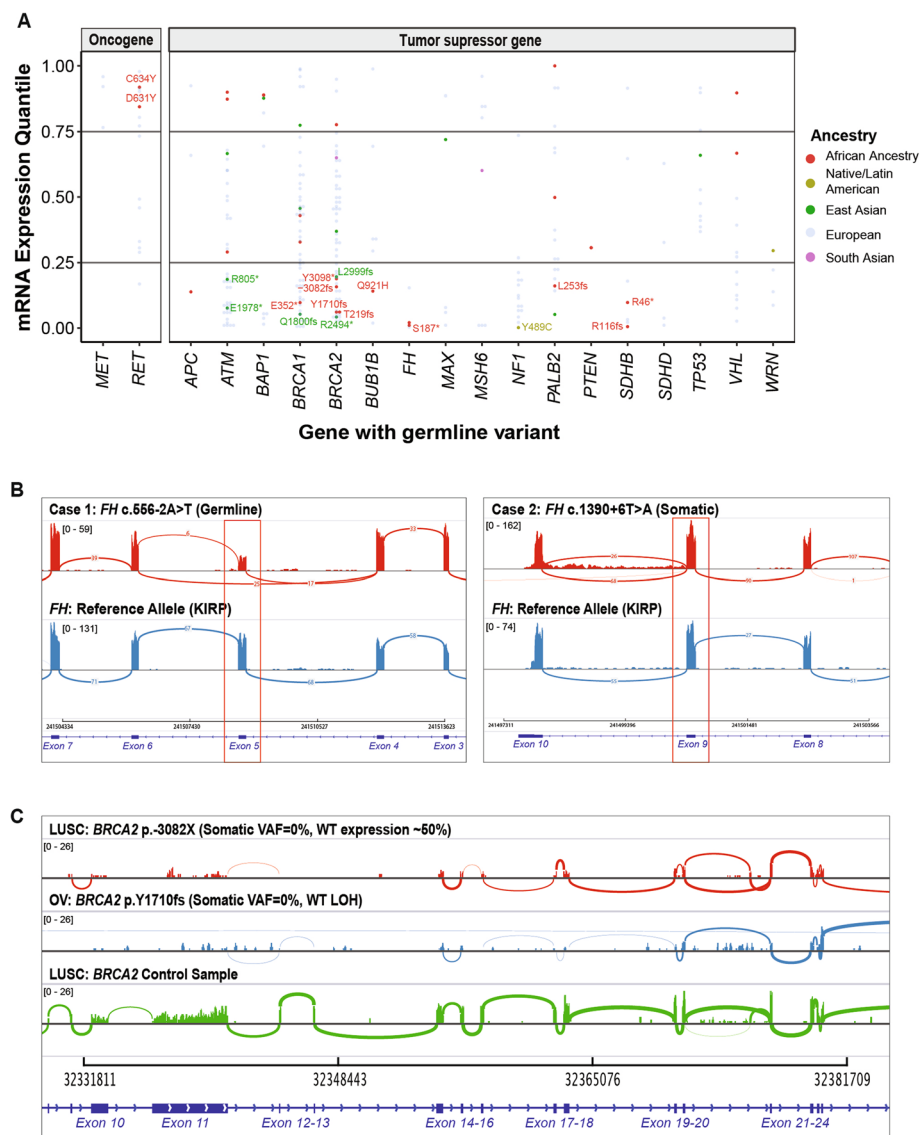


Fig. 4 Expression changes associated with the predisposing variants. **a** mRNA gene expression of the affected genes in the carriers of ancestry-specific variants as quantiles in their respective cancer cohort. Each dot denotes the gene expression level of a predisposing variant carrier colored by ancestry. Non-European variants corresponding to the bottom 25% expression in affected tumor suppressor genes and top 25% expression in affected oncogenes are further labeled. **b** Tumor RNA expression highlighting (red box) mis-spliced exon 5 with germline or somatic splice site variants in two cases with *FH* splice site variants as visualized using the integrated genome viewer (IGV). **c** Tumor RNA expression for the *BRCA2* gene. The first two rows correspond to samples with a germline predisposing variant coupled with or without somatic LOH event, respectively. The third row corresponds to an unrelated sample without any *BRCA2* alteration. All three coverage plots are group-scaled to show lower expression in the two samples harboring *BRCA2* alterations

ancestry carriers, including p.C631Y (84% in KIRC) and p.D634Y (91% in PCGP).

Power consideration for predisposing gene discovery

Given the currently limited sample sizes in most of the minority cohorts, we sought to identify the required numbers of samples to discover novel cancer predisposing genes. We performed post hoc power analyses to detect a rare-variant association in an aggregation test

using SKAT [34]. We assumed that a high proportion (80%) of variants are casual when focusing on prioritized predisposing variants in accordance with ACMG/AMP guidelines (Additional file 2: Table S6a, see the “Methods” section) [26, 27, 32]. The detection of rare variants (AF < 0.01) with moderate effect sizes (odds ratio [OR] > 5) with at least 80% power requires sample sizes exceeding 1000 samples ($n = 1014$) per cancer type (Additional file 1: Fig. S2A).

The sample size requirement suggests limited power for ancestry-specific analyses using TCGA, one of the largest cancer sequencing cohorts to date. For the largest ancestry subgroup in the study, European-ancestry BRCA cases ($n = 811$), there is 67% power to detect genes with smaller effect sizes ($OR < 3$). For all other ancestries, their respective largest cohorts afford inadequate power to detect genes with large effect sizes ($OR = 9$), including the African ancestry BRCA cohort ($n = 180$, power = 36%), the East Asian-ancestry LIHC cohort ($n = 162$, power = 24.5%), and the Native/Latin American-ancestry THCA cohort ($n = 11$, power = < 1%). As a reference, most known cancer predisposing genes, including *ATM*, *PTEN*, *STK11*, *CHEK2*, *BRIP1*, and *PALB2*, have an estimated $OR < 10$. *BRCA1/BRCA2* are exceptions with an $OR > 10$ for BRCA, but also show more moderate OR for other cancer types [49]. Despite limited power, this TCGA study includes threefold more non-European cases ($n = 1715$) compared to the combined number of samples across 27 published non-TCGA sequencing studies that report race/ethnicity information from cancer cohorts ($n = 764$ non-Europeans, 10 cancer types) [9]. Moreover, the majority of these studies focused on somatic alterations, and only a handful reported ancestry-specific germline predisposition (Additional file 2: Table S7).

Standard power analyses have the caveat of assuming various unknown parameters that may be inaccurate. We thus performed a downsampling analysis using two cancer types with at least five significantly associated germline genes in the European-ancestry: pheochromocytoma and paraganglioma (PCPG) and sarcoma (SARC) [4] (Additional file 1: Fig. S2B, Additional file 2: Table S6b). We found that the sample size requirements differ for each gene and cancer cohort, likely due to varying penetrance. For example, six predisposing genes are discovered in both PCPG ($n = 146$) and SARC ($n = 217$) samples of the European ancestry, respectively, at their full cohort size. Upon downsampling the cohort size in half, we found *VHL*, *SDHB*, *RET*, and *NF1* to be still associated in 73 PCPG cases, whereas only *TP53* remained significantly associated in 108 SARC cases. Even while assuming similar penetrance of the predisposing genes across ancestries, this analysis implicates that the discovery power is still far from saturation for most ancestry-specific cohorts ($N < 100$). The different predisposition landscapes across cancer types should also be accounted for in future study designs.

Discussion

We report one of the most extensive multi-ancestry investigations of rare cancer predisposing genes to date, encompassing 9899 cancer cases across 33 cancer types. In the African ancestry, our results validated six known

predisposing genes and nominated *BRCA2* as a potential predisposing gene for LUSC (Fig. 1) previously shown only for Europeans. In the East Asian ancestry, we found predisposing variants affecting *BRIP1* in STAD that warrants further investigation. Although the number of germline predisposing variants is small, they were associated with LOH (Fig. 3), biallelic mutations (Fig. 2), and gene expression effects in the tumor samples (Fig. 4), supporting their potential contribution to cancer predisposition in carriers.

In this TCGA cohort, we found multiple significant predisposing genes for the European ancestry and seven for the African ancestry, yet lack cancer cohorts with sufficient testing samples for many other ancestries, including Native/Latin American and South Asian that each constitute a considerable fraction of the US population. Even when tested, this study likely contains false negatives in multiple smaller cancer cohorts, especially those of non-Europeans. To achieve 80% power, the post hoc power calculation showed that the detection of rare variants ($AF < 0.01$) with moderate effect sizes ($OR > 5$) requires at least 1014 samples (Additional file 1: Fig. S2), a cohort size larger than any of the TCGA non-European cohorts.

It is necessary to use caution when interpreting the ancestry-specific predisposing gene associations identified herein or previous studies of smaller sample sizes, where a handful of carriers may give rise to the association in a limited cancer cohort. Further, the suggestive associations nominated by the TFT analyses will need to be established by analyses of larger cohorts adjusted for potential confounders. Two of the associations we identified in the African ancestry were also complemented by familial studies [39, 42], providing further validation. To design future cancer genomics studies, one must note that the power considerations differ for discovering somatic driver genes and germline predisposing genes. Current detection powers have potentially reached saturation in detecting somatically mutated genes for sample sizes in multiple cancer types of TCGA [4], although racial disparities of the sequencing data could potentially limit the generalizability of findings [50–52]. We further highlighted the imbalanced dataset limits power for germline gene discovery in populations under-represented in research studies.

We observed selected predisposing genes shared across ancestries (ex. *BRCA2* in BRCA/OV and *SDHB* in PCPG for both the African and European ancestries). Predisposing variants, on the other hand, are highly ancestry-specific (Fig. 2). Many of the predisposing variants found in the African or East Asian ancestry were not identified in the much larger European-ancestry population of TCGA ($n = 8184$) or even the gnomAD non-cancer cohort ($n = 134,187$) or submitted to ClinVar

by clinical laboratories assessing patients for cancer predisposition. Rare variant classification and interpretation remain a challenge given the low frequency of observation precluding statistical associations. The identification of ancestry-specific predisposing variants further highlights this challenge in minority groups, where current germline sequencing often results in higher rates of variants of unknown significance (VUSs) [16].

Personalized medicine provides tailored disease diagnosis and treatment plans based on an individual's unique genetic profile. The knowledge of different cancer predisposing genes and prevalence across ancestries suggests that we need to provide ancestry-specific interpretations of genetic data. In particular, many of the current guidelines for when genetic testing is recommended rely on the underlying likelihood of identifying a germline variant. Thus, accurate estimates of germline prevalence may alter recommendations for different patient populations. At the current sample sizes for minority cohorts, our study is still limited in power to discover and establish ancestry-specificity of predisposing genes (Additional file 1: Fig. S2). However, we were able to discover many ancestry-specific variants not currently submitted to ClinVar. Further, much of the diverse populations within the USA, not to mention worldwide, still lack representation in existing sequencing cohorts. Ongoing sequencing projects will begin to address this disparity within US populations (e.g., CSER [17], eMERGEIII [18], Million Veteran Program [19], and the All of Us Research Program) and multiple countries in East Asia and Europe [53]. Yet, many populations, such as the diverse African ancestry [54], remain underserved although projects like H3Africa are designed to address this problem. Additional efforts will be required to deliver the promise of genome-based precision medicine for all.

TCGA provides a powerful multi-omic sequencing dataset comprising more than ten thousand adult cancer cases [55, 56]. The dataset is used not only for characterizing somatic mutations and molecular subtypes but also enables studies of rare genetic predisposition and germline-somatic interactions [20, 38, 57–59]. However, in such applications, one needs to note that TCGA is not a prospective cohort nor designed as a case-control study. Using the matched-ancestry cases of other cancer types as “controls” (the “Methods” section) is not ideal, yet they are the only available samples in the same study. The associations herein, therefore, may show biased effect sizes that require validation in carefully designed epidemiological studies. To enhance the confidence of the reported variants, we focused on identifying their somatic impacts, including LOH, ASE, and extreme gene expression levels that can be uniquely revealed in the multi-omic dataset.

To aid interpretation of low-frequency ancestry-specific variants, evidence of a somatic second hit event (i.e., loss of heterozygosity [LOH] or a biallelic mutation) in tumor samples can support functionality. Our analysis of the two-hit model identified the second somatic events in two thirds (10/15) of the African ancestry-specific predisposing variants and in one out of six of the East Asian ancestry-specific predisposing variants (Additional file 2: Table S4b). Additionally, some carriers of ancestry-specific predisposing variants showed simultaneous extreme expression of the affected genes (Fig. 3). Such evidence derived from analysis of the somatic genome or transcriptome can be further utilized to characterizing rare germline variants [60], especially when DNA-level analysis still suffers from limited sample sizes.

Our observation of somatic second hit (Figs. 2 and 3) and transcriptional effects (Fig. 4) coupled with germline variants also adds on to the current literature on germline-somatic interactions in cancer [61]. While the majority of cancer genomic studies focus exclusively on the germline or somatic genome, pathogenic germline variants are associated with different somatic mutational signatures, allele-specific imbalance, or somatic drivers [20, 38, 58, 62, 63]. The availability of germline DNA analysis and tumor genomic and transcriptomic analyses from the same individual provides critical data to the analyses described here that is not possible in many studies that only analyze germline DNA samples alone. Collectively, these findings are providing the roadmaps of how germline variants may trigger and collaborate with specific somatic mutations, eventually leading to cancer development. In this process, genomes across different ancestral populations provide different contexts for developing somatic mutations and genomic instability, even when the individual carries the same germline predisposition variant. We showcased examples of predisposition-associated LOH and gene expression changes in diverse individuals. As sample sizes of sequencing cohorts expand, analyzing germline-somatic interactions across ancestry will be pivotal to reveal potential ancestry-specific effects.

Conclusions

In summary, we identify ancestry-specific predisposing genes and variants contributing to multiple cancer types. The results provide insights into rare genetic predisposition and their somatic impacts in cases of African and East Asian ancestries. While the identified cancer predisposition genes are known, most predisposing variants are found to be exclusive within ancestries, supporting the “clan-genomics” hypothesis [64]. Continuous studies using larger ancestry cohorts will be required to enable adequately powered discovery of predisposing genes and improve genetic screening for diverse populations [65].

Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s13073-020-00744-3>.

Additional file 1: Figure S1. Principal component analyses (PCA) of germline TCGA samples to infer genetic ancestry as performed by PanCanAtlas Ancestry Informative Markers (AIM) working group. **Figure S2.** Power analysis for ancestry-specific sample sizes to discover predisposing genes. **Figure S3.** Nonsense-mediated decay prediction for predisposing frameshift variants in African and East Asian ancestries.

Additional file 2: Table S1. The demographic information of TCGA PanCanAtlas cohort with separate admixture populations. **Table S2a.** Ancestry-specific cancer-gene associations discovered from multivariate regression analyses. **Table S2b.** Ancestry-specific cancer-gene associations discovered from rare variant burden testing (Total Frequency Test-TFT). **Table S3.** Frequency of predisposing variants in TCGA PanCanAtlas and gnomAD-non-cancer subset across all ancestries. **Table S4a.** Ancestry-Specific Predisposing Variants as identified from Supp. Table.2. **Table S4b.** Summary of somatic second hit mutations in carriers of germline predisposing variants. **Table S5a.** Statistical analysis of gene expression in tumor samples of the variant carriers vs. non-carriers within each ancestry-cancer combination. **Table S5b.** Tumor RNAseq variant allele fractions and the somatic second hit events in germline predisposing variants with extreme expression within that cancer type. **Table S6a.** Post hoc power analyses to detect rare-variant associations in an aggregation test using SKAT. **Table S6a.** Down-sampling analysis for PCGP and SARC (cancers with at least 5 significantly associated germline genes in the European ancestry). **Table S7.** Prior studies that report ancestry-specific germline predisposition.

Abbreviations

ACC: Adrenocortical carcinoma; BLCA: Bladder urothelial carcinoma; BRCA: Breast invasive carcinoma; CESC: Cervical squamous cell carcinoma and endocervical adenocarcinoma; CHOL: Cholangiocarcinoma; COAD: Colon adenocarcinoma; DLBC: Lymphoid neoplasm diffuse large B cell lymphoma; ESCA: Esophageal carcinoma; GBM: Glioblastoma multiforme; HNSC: Head and neck squamous cell carcinoma; KICH: Kidney chromophobe; KIRC: Kidney renal clear cell carcinoma; KIRP: Kidney renal papillary cell carcinoma; LAML: Acute myeloid leukemia; LGG: Brain lower grade glioma; LIHC: Liver hepatocellular carcinoma; LUAD: Lung adenocarcinoma; LUSC: Lung squamous cell carcinoma; MESO: Mesothelioma; OV: Ovarian serous cystadenocarcinoma; PAAD: Pancreatic adenocarcinoma; PCPG: Pheochromocytoma and paraganglioma; PRAD: Prostate adenocarcinoma; READ: Rectum adenocarcinoma; SARC: Sarcoma; SKCM: Skin cutaneous melanoma; STAD: Stomach adenocarcinoma; TGCT: Testicular germ cell tumors; THCA: Thyroid carcinoma; THYM: Thymoma; UCEC: Uterine corpus endometrial carcinoma; UCS: Uterine carcinosarcoma; UVM: Uveal melanoma; SEER: Surveillance, Epidemiology, and End Results; TCGA: The Cancer Genome Atlas; VUS: Variants of unknown significance; LOH: Loss of heterozygosity; AIM: Ancestry informative markers; PCA: Principal component analyses; OR: Odds ratio; ExAC: Exome aggregation consortium; gnomAD: Genome aggregation dataset; IGV: Integrative Genomics Viewer; GDC: Genomic Data Commons; TFT: Total Frequency Test; ACMG-AMP: American College of Medical Genetics- Association for Molecular Pathology; FDR: False discovery rate; VAF: Variant allele fraction; AC: Allele count; AN: Allele number; AF: Allele frequency; CI: Confidence interval; WES: Whole exome sequencing; WGS: Whole genome sequencing

Acknowledgements

The authors wish to acknowledge The Cancer Genome Atlas and its participating patients and family that generously contributed the data. The authors would also like to acknowledge members of the TCGA PanCanAtlas Research network, particularly active members of the Germline Analysis Working Group and the Ancestry Informative Markers Analysis Working Group, for helpful discussions. ZHG acknowledges funds from the LUNgevity Foundation.

©TCGA Analysis Network

Jian Carrot-Zhang^{1,2}, Nyasha Chambwe³, Jeffrey S. Damrauer⁴, Theo A. Knijnenburg³, A. Gordon Robertson⁵, Christina Yau^{6,7}, Wanding Zhou⁸, Ashton C. Berger^{1,2}, Kuan-lin Huang⁹, R. Jay Mash¹⁰, Justin Newberg¹¹, Alessandro

Romanel¹², Rosalyn W. Sayaman^{13,14}, Francesca Demichelis¹², Ina Felau¹⁵, Garrett Frampton¹¹, Seunghun Han^{1,2}, Katherine A. Hoadley⁴, Anab Kemal¹⁵, Peter W. Laird⁹, Alexander J. Lazar¹⁶, Xiuning Le¹⁷, Ninad Oak^{18, 19}, Hui Shen⁸, Christopher K. Wong²⁰, Jean C. Zenklusen¹⁵, Elad Ziv^{13,14}, Francois Aguet¹, Li Ding⁹, John A. Demchok¹⁵, Michael K.A. Mensah¹⁵, Roy Tarnuzzer¹⁵, Zhining Wang¹⁵, Liming Yang¹⁵, Jessica Alfoldi¹, Konrad J. Karczewski¹, Daniel G. MacArthur¹, Garret M. Frampton¹¹, Christopher Benz⁶, Joshua M. Stuart²⁰, Andrew D. Cherniack^{1,2}, Rameen Beroukhi^{1,2,21}.

1. The Eli and Edythe L. Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, MA 02142, USA.
2. Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA 02215, USA.
3. Institute for Systems Biology, Seattle, WA 98109, USA.
4. Lineberger Comprehensive Cancer Center, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, USA.
5. British Columbia Cancer Agency, Genome Sciences Centre, Vancouver, Canada V5Z 4S6.
6. Buck Institute for Research on Aging, Novato, CA 94945, USA.
7. Department of Surgery, University of California, San Francisco, San Francisco, CA 94115, USA.
8. Van Andel Research Institute, Grand Rapids, MI 49503, USA.
9. Department of Genetics and Genomics, Icahn School of Medicine at Mount Sinai, New York, NY 2129, USA.
10. Department of Medicine, Washington University in St. Louis, St. Louis, MO 63110, USA.
11. Cancer Genomics Research, Foundation Medicine, Inc., Cambridge, MA 02141, USA.
12. Department of Cellular, Computational and Integrative Biology (CIBIO), University of Trento, Via Sommarive 9 Povo (TN) 38123 Italy.
13. Department of Laboratory Medicine, Helen Diller Family Comprehensive Cancer Center, University of California San Francisco, San Francisco, CA 94143, USA.
14. Department of Population Sciences, Beckman Research Institute, City of Hope, Duarte, CA 9210.
15. National Cancer Institute, Bethesda, MD 20892, USA.
16. Departments of Pathology, Genomic Medicine, and Translational Molecular Pathology, The University of Texas M.D. Anderson Cancer Center, Houston, TX 77030, USA.
17. Department of Thoracic and Head and Neck Medical Oncology, The University of Texas M.D. Anderson Cancer Center, Houston, TX 77030, USA.
18. Department of Oncology, St. Jude Children's Research Hospital, Memphis, TN 38105, USA.
19. Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX 77030, USA.
20. Department of Biomolecular Engineering, Center for Biomolecular Sciences and Engineering, University of California, Santa Cruz, Santa Cruz, CA 95064, USA.
21. Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA.

Authors' contributions

K.H. and N.O. conceived and designed the research and analyses. N.O., K.H., R.B., S.E.P., and L.D. acquired the data. The PanCanAtlas AIM working group, A.C., J.M., and K.H. conducted the genetic ancestry assignment. N.O. and K.H. conducted the analyses. N.O., K.H., and S.E.P. interpreted the results and drafted the manuscript. K.H. supervised the study. All authors read, edited, and approved the manuscript.

Funding

K.H. was supported by the Icahn School of Medicine at Mount Sinai seed fund to the Huang lab. This work was also funded by the ClinGen grant [1U41HG009649] supported by the National Human Genome Research Institute to SEP.

Availability of data and materials

Ancestry ascertainment for all samples in this study were obtained from the TCGA-AIM analysis working group's publication and its supplementary information files [21]. Germline variant calls used in this study are available through controlled access data release of the TCGA PanCanAtlas germline working group (<https://gdc.cancer.gov/about-data/publications/PanCanAtlas-Germline-AWG>) [20]. Somatic mutation and expression data are publicly

available at the National Cancer Institute's GDC data portal (<https://portal.gdc.cancer.gov/>). All downstream analysis results generated in this study are available in Additional files 1 and 2.

Ethics approval and consent to participate

The need for Institutional Review Board Approval at our institution (Icahn School of Medicine at Mount Sinai, New York) was waived for this study as all data used from this project had previously been generated as part of The Cancer Genome Atlas Project and none of the results reported in this manuscript can be used to identify individual patients.

Consent for publication

Not applicable.

Competing interests

S.E.P. is a member of the Baylor Genetics laboratory scientific advisory panel. A.D.C. receives research funding from Bayer. R.B. owns equity in Ampressa Therapeutics and receives research funding from Novartis. The remaining authors declare that they have no competing interests.

Author details

¹Department of Oncology, St. Jude Children's Research Hospital, Memphis, TN 38105, USA. ²Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX 77030, USA. ³The Eli and Edythe L. Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, MA 02142, USA. ⁴Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA 02215, USA. ⁵Department of Medicine, Washington University in St. Louis, St. Louis, MO 63108, USA. ⁶McDonnell Genome Institute, Washington University in St. Louis, St. Louis, MO 63108, USA. ⁷Department of Oncological Sciences, Center for Thoracic Oncology, Tisch Cancer Institute, New York, NY, USA. ⁸Department of Genetics, Washington University in St. Louis, St. Louis, MO 63108, USA. ⁹Steman Cancer Center, Washington University in St. Louis, St. Louis, MO 63108, USA. ¹⁰Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA. ¹¹Department of Genetics and Genomics, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA. ¹²Icahn Institute for Data Science and Genomic Technology, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA. ¹³Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX 77030, USA. ¹⁴Department of Pediatrics, Baylor College of Medicine, Houston, TX 77030, USA. ¹⁵Center for Transformative Disease Modeling, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA.

Received: 2 December 2019 Accepted: 7 May 2020

Published online: 29 May 2020

References

- Siegel RL, Miller KD, Jemal A. Cancer statistics, 2018. *CA Cancer J Clin.* 2018; 68:7–30.
- Noone AM, Howlander N, Krapcho M, Miller D, Brest A, Yu M, et al. SEER cancer statistics review. Bethesda: National Cancer Institute; 2018. Available from: https://seer.cancer.gov/csr/1975_2015/.
- Kurian AW. BRCA1 and BRCA2 mutations across race and ethnicity: distribution and clinical implications. *Curr Opin Obstet Gynecol.* 2010;22:72–8.
- Lawrence MS, Stojanov P, Mermel CH, Robinson JT, Garraway LA, Golub TR, et al. Discovery and saturation analysis of cancer genes across 21 tumour types. *Nature.* 2014;505:495–501.
- Huo D, Hu H, Rhie SK, Gamazon ER, Cherniack AD, Liu J, et al. Comparison of breast cancer molecular features and survival by African and European ancestry in The Cancer Genome Atlas. *JAMA Oncol.* 2017;3:1654–62.
- Easton DF, Pooley KA, Dunning AM, Pharoah PDP, Thompson D, Ballinger DG, et al. Genome-wide association study identifies novel breast cancer susceptibility loci. *Nature.* 2007;447:1087–93.
- Zheng W, Zhang B, Cai Q, Sung H, Michailidou K, Shi J, et al. Common genetic determinants of breast-cancer risk in East Asian women: a collaborative study of 23 637 breast cancer cases and 25 579 controls. *Hum Mol Genet.* 2013;22:2539–50.
- Huo D, Feng Y, Haddad S, Zheng Y, Yao S, Han Y-J, et al. Genome-wide association studies in women of African ancestry identified 3q26.21 as a novel susceptibility locus for oestrogen receptor negative breast cancer. *Hum Mol Genet.* 2016;25:4835–46.
- Nugent A, Conatser KR, Turner LL et al. Reporting of race in genome and exome sequencing studies of cancer: a scoping review of the literature. *Genet Med.* 2019;21:2676–80. <https://doi.org/10.1038/s41436-019-0558-2>.
- Lilyquist J, Ruddy KJ, Vachon CM, Couch FJ. Common genetic variation and breast cancer risk—past, present, and future. *Cancer Epidemiol Prev Biomark.* 2018;27:380–94.
- Couch FJ, Shimelis H, Hu C, Hart SN, Polley EC, Na J, et al. Associations between cancer predisposition testing panel genes and breast cancer. *JAMA Oncol.* 2017;3:1190–6.
- Churpek JE, Walsh T, Zheng Y, Moton Z, Thornton AM, Lee MK, et al. Inherited predisposition to breast cancer among African American women. *Breast Cancer Res Treat.* 2015;149:31–9.
- Fackenthal JD, Zhang J, Zhang B, Zheng Y, Hagos F, Burrill DR, et al. High prevalence of BRCA1 and BRCA2 mutations in unselected Nigerian breast cancer patients. *Int J Cancer.* 2012;131:1114–23.
- Newman B, Mu H, Butler LM, Millikan RC, Moorman PG, King MC. Frequency of breast cancer attributable to BRCA1 in a population-based series of African American women. *JAMA.* 1998;279:915–21.
- Schubert EL, Lee MK, Mefford HC, Argonza RH, Morrow JE, Hull J, et al. BRCA2 in American families with four or more cases of breast or ovarian cancer: recurrent and novel mutations, variable expression, penetrance, and the possibility of families whose cancer is not attributable to BRCA1 or BRCA2. *Am J Hum Genet.* 1997;60:1031–40.
- Caswell-Jin JL, Gupta T, Hall E, Petrovich IM, Mills MA, Kingham KE, et al. Racial/ethnic differences in multiple-gene sequencing results for hereditary cancer risk. *Genet Med.* 2018;20:234–9.
- Amendola LM, Berg JS, Horowitz CR, Angelo F, Bensen JT, Biesecker BB, et al. The clinical sequencing evidence-generating research consortium: integrating genomic sequencing in diverse and medically underserved populations. *Am J Hum Genet.* 2018;103:319–27.
- Harmonizing Clinical Sequencing and Interpretation for the eMERGE III Network. *Am J Hum Genet.* 2019;105:588–605.
- Gaziano JM, Concato J, Brophy M, Fiore L, Pyarajan S, Breeling J, et al. Million Veteran Program: a mega-biobank to study genetic influences on health and disease. *J Clin Epidemiol.* 2016;70:214–23.
- Huang KL, Mashl RJ, Wu Y, et al. Pathogenic Germline Variants in 10,389 Adult Cancers. *Cell.* 2018;173(2):355–70.e14. <https://doi.org/10.1016/j.cell.2018.03.039>.
- Carrot-Zhang J, Chambwe N, Damrauer JS, et al. Comprehensive Analysis of Genetic Ancestry and Its Molecular Correlates in Cancer. *Cancer Cell.* 2020; 37(5):639–54.e6. <https://doi.org/10.1016/j.ccell.2020.04.012>.
- Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* 2009;19:1655–64.
- Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, et al. VarScan 2: somatic mutation and copy number alteration discovery in cancer by exome sequencing. *Genome Res.* 2012;22:568–76.
- McKenna AH, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 2010;20:1297–303.
- Ye K, Schulz MH, Long Q, Apweiler R, Ning X, Pindel: a pattern growth approach to detect break points of large deletions and medium sized insertions from paired-end short reads. *Bioinformatics.* 2009;25:2865–71.
- Scott AD, Huang KL, Weerasinghe A, et al. CharGer: clinical Characterization of Germline variants. *Bioinformatics.* 2019;35(5):865–67. <https://doi.org/10.1093/bioinformatics/bty649>.
- Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, et al. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet Med.* 2015;17:405–23.
- Purcell S, Neale B, Todd-Brown K, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet.* 2007; 81(3):559–75. <https://doi.org/10.1086/519795>.
- Basu S, Pan W. Comparison of statistical tests for disease association with rare variants. *Genet Epidemiol.* 2011;35:606–19.
- Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. Variation across 141,456 human exomes and genomes reveals the spectrum of loss-of-function intolerance across human protein-coding genes. *bioRxiv.* 2019;531210. <https://doi.org/10.1101/531210>.
- Lek M, Karczewski KJ, Samocha KE, Banks E, Fennell T, O AH, et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature.* 2015;536:285–291.

32. Li Q, Wang K. InterVar: Clinical Interpretation of Genetic Variants by the 2015 ACMG-AMP Guidelines. *Am J Hum Genet.* 2017;100(2):267–80. <https://doi.org/10.1016/j.ajhg.2017.01.004>.
33. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* 2010;38:1–7.
34. Wu MC, Lee S, Cai T, Li Y, Boehnke M, Lin X. Rare-variant association testing for sequencing data with the sequence kernel association test. *Am J Hum Genet.* 2011;89:82–93.
35. Esai Selvan M, Klein RJ, Gümüş ZH. Rare, Pathogenic Germline Variants in Fanconi Anemia Genes Increase Risk for Squamous Lung Cancer. *Clin Cancer Res.* 2019;25(5):1517–25. <https://doi.org/10.1158/1078-0432.CCR-18-2660>.
36. Wang Y, McKay JD, Rafnar T, Wang Z, Timofeeva MN, Broderick P, et al. Rare variants of large effect in BRCA2 and CHEK2 affect risk of lung cancer. *Nat Genet.* 2014;46:736–41.
37. Gayther SA, Mangion J, Russell P, Seal S, Barfoot R, Ponder BAJ, et al. Variation of risks of breast and ovarian cancer associated with different germline mutations of the BRCA2 gene. *Nat Genet.* 1997;15:103–5.
38. Lu C, Xie M, Wendl MC, Wang J, McLellan MD, Leiserson MDM, et al. Patterns and functional implications of rare germline variants across 12 cancer types. *Nat Commun.* 2015;6:10086.
39. Majumdar S, Friedrich CA, Koch CA, Megason GC, Fratkin JD, Moll GW. Compound heterozygous mutation with a novel splice donor region DNA sequence variant in the succinate dehydrogenase subunit B gene in malignant paraganglioma. *Pediatr Blood Cancer.* 2010;54:473–5.
40. Beebe-Dimmer JL, Zuhlke KA, Johnson AM, Liesman D, Cooney KA. Rare germline mutations in African American men diagnosed with early-onset prostate cancer. *Prostate.* 2018;78:321–6.
41. Na R, Zheng SL, Han M, Yu H, Jiang D, Shah S, et al. Germline mutations in ATM and BRCA1/2 distinguish risk for lethal and indolent prostate cancer and are associated with early age at death. *Eur Urol.* 2017;71:740–7.
42. Wei M-H, Toure O, Glenn GM, Pithukpakorn M, Neckers L, Stolle C, et al. Novel mutations in FH and expansion of the spectrum of phenotypes expressed in families with hereditary leiomyomatosis and renal cell cancer. *J Med Genet.* 2006;43:18–27.
43. Krishnan B, Rose TL, Kardos J, Milowsky MI, Kim WY. Intrinsic genomic differences between African American and white patients with clear cell renal cell carcinoma. *JAMA Oncol.* 2016;2:664–7.
44. Sutton J, Orloff MS, Michener C, Chiesa-Vottero A, Prayson R, Nowacki AS, et al. Association of specific PTEN/10q haplotypes with endometrial cancer phenotypes in African-American and European American women. *Gynecol Oncol.* 2015;138:434–40.
45. Landrum MJ, Lee JM, Benson M, et al. ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res.* 2018;46(D1):D1062–67. <https://doi.org/10.1093/nar/gkx1153>.
46. Knudson AG Jr. Mutation and cancer: statistical study of retinoblastoma. *Proc Natl Acad Sci U S A.* 1971;68(4):820–23. <https://doi.org/10.1073/pnas.68.4.820>.
47. Knudson AG. Two genetic hits (more or less) to cancer. *Nat Rev Cancer.* 2001;1(2):157–62. <https://doi.org/10.1038/35101031>.
48. Hsu MK, Lin HY, Chen FC. NMD Classifier: A reliable and systematic classification tool for nonsense-mediated decay events. *PLoS One.* 2017;12(4):e0174798. <https://doi.org/10.1371/journal.pone.0174798>.
49. Turnbull C, Rahman N. Genetic predisposition to breast cancer: past, present, and future. *Annu Rev Genomics Hum Genet.* 2008;9:321–45. <https://doi.org/10.1146/annurev.genom.9.081307.164339>.
50. Spratt DE, Chan T, Waldron L, et al. Racial/Ethnic Disparities in Genomic Sequencing. *JAMA Oncol.* 2016;2(8):1070–4. <https://doi.org/10.1001/jamaoncol.2016.1854>.
51. Berger AC, Korkut A, Kanchi RS, et al. A Comprehensive Pan-Cancer Molecular Study of Gynecologic and Breast Cancers. *Cancer Cell.* 2018;33(4):690–705.e9. <https://doi.org/10.1016/j.ccell.2018.03.014>.
52. Yuan J, Hu Z, Mahal BA, et al. Integrated Analysis of Genetic Ancestry and Genomic Alterations across Cancers. *Cancer Cell.* 2018;34(4):549–60.e9. <https://doi.org/10.1016/j.ccell.2018.08.019>.
53. Abul-Husn NS, Kenny EE. Personalized medicine and the power of electronic health records. *Cell.* 2019;177:58–69.
54. Tishkoff SA, Reed FA, Friedlaender FR, Ehret C, Ranciaro A, Froment A, et al. The genetic structure and history of Africans and African Americans. *Science.* 2009;324:1035–44.
55. Cancer Genome Atlas Research Network, Weinstein JN, Collisson EA, Mills GB, KRM S, Ozenberger BA, et al. The Cancer Genome Atlas Pan-Cancer analysis project. *Nat Genet.* 2013;45:1113–20.
56. Ding L, Bailey MH, Porta-Pardo E, Thorsson V, Colaprico A, Bertrand D, et al. Perspective on oncogenic processes at the end of the beginning of cancer genomics. *Cell.* 2018;173:305–320.e10.
57. Park S, Supek F, Lehner B. Systematic discovery of germline cancer predisposition genes through the identification of somatic second hits. *Nat Commun.* 2018;9:2601.
58. Carter H, Marty R, Hofree M, Gross AM, Jensen J, Fisch KM, et al. Interaction landscape of inherited polymorphisms with somatic events in cancer. *Cancer Discov.* 2017;7:410–23.
59. Shahamatdar S, He MX, Reyna MA, Gusev A, AlDubayan SH, Van Allen EM, et al. Germline features associated with immune infiltration in solid tumors. *Cell Rep.* 2020;30:2900–2908.e4.
60. Walsh MF, Ritter DI, Kesserwan C, et al. Integrating somatic variant data and biomarkers for germline variant classification in cancer predisposition genes. *Hum Mutat.* 2018;39(11):1542–52. <https://doi.org/10.1002/humu.23640>.
61. Ramroop JR, Gerber MM, Toland AE. Germline variants impact somatic events during tumorigenesis. *Trends Genet.* 2019;35:515–26.
62. Campbell PJ, Getz G, Korbel JO et al. Pan-cancer analysis of whole genomes. *Nature.* 2020;578:82–93. <https://doi.org/10.1038/s41586-020-1969-6>.
63. Waszak SM, Tiao G, Zhu B, Rausch T, Muyas F, Rodríguez-Martín B, et al. Germline determinants of the somatic mutation landscape in 2,642 cancer genomes. *bioRxiv.* 2017;208330.
64. Lupski JR, Belmont JW, Boerwinkle E, Gibbs RA. Clan genomics and the complex architecture of human disease. *Cell.* 2011;147:32–43.
65. Hindorf LA, Bonham VL, Brody LC, Ginoza MEC, Hutter CM, Manolio TA, et al. Prioritizing diversity in human genomics research. *Nat Rev Genet.* 2018;19:175–85.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

