

1 **Gene level germline contributions to clinical risk of recurrence scores in Black and White breast cancer**  
2 **patients**

3 Achal Patel<sup>1</sup>, Montserrat García-Closas<sup>2,3</sup>, Andrew F. Olshan<sup>1,4</sup>, Charles M. Perou<sup>4,5,6</sup>, Melissa A. Troester<sup>1,6</sup>,  
4 Michael I. Love<sup>5,7</sup>, Arjun Bhattacharya<sup>8,9\*</sup>

5  
6 1. Department of Epidemiology, Gillings School of Global Public Health, University of North Carolina-Chapel Hill,  
7 Chapel Hill, NC, USA

8 2. Division of Cancer Epidemiology and Genetics, National Cancer Institute, Bethesda, USA

9 3. Division of Genetics and Epidemiology, Institute of Cancer Research, London, UK

10 4. Lineberger Comprehensive Cancer Center, University of North Carolina-Chapel Hill, Chapel Hill, USA

11 5. Department of Genetics, University of North Carolina-Chapel Hill, Chapel Hill, NC, USA

12 6. Department of Pathology and Laboratory Medicine, University of North Carolina-Chapel Hill, Chapel Hill, NC,  
13 USA

14 7. Department of Biostatistics, Gillings School of Global Public Health, University of North Carolina-Chapel Hill,  
15 Chapel Hill, NC, USA

16 8. Department of Pathology and Laboratory Medicine, David Geffen School of Medicine, University of California-  
17 Los Angeles, Los Angeles, CA, USA

18 9. Institute for Quantitative and Computational Biosciences, David Geffen School of Medicine, University of  
19 California-Los Angeles, Los Angeles, CA, USA

20

21 Running Title: multi-ancestry GReX study of continuous risk-of-recurrence

22

23 Correspondence can be directed to Arjun Bhattacharya

24 Mailing Address: 3746 Mentone Avenue, Apartment 202, Los Angeles, CA 90034

25 Phone: +1 (919) 742-0101

26 Email: [abtbhatt@ucla.edu](mailto:abtbhatt@ucla.edu)

27

28 **CONFLICT OF INTEREST STATEMENT**

- 1 CMP is an equity stock holder, consultant, and board of directors member of BioClassifier LLC and GeneCentric
- 2 Diagnostics. CMP is also listed as an inventor on patent applications on the Breast PAM50 assay. The other
- 3 authors declare no potential conflicts of interest.

1 **ABSTRACT**

2 Continuous risk of recurrence scores (CRS) based on tumor gene expression are vital prognostic tools for breast  
3 cancer (BC). Studies have shown that Black women (BW) have higher CRS than White women (WW). Although  
4 systemic injustices contribute substantially to BC disparities, evidence of biological and germline contributions is  
5 emerging. In this study, we investigated germline genetic associations with CRS and CRS disparity using  
6 approaches modeled after transcriptome-wide association studies (TWAS). In the Carolina Breast Cancer Study,  
7 using race-specific predictive models of tumor expression from germline genetics, we performed race-stratified  
8 (N=1,043 WW, 1,083 BW) linear regressions of three CRS (ROR-S: PAM50 subtype score; Proliferation Score;  
9 ROR-P: ROR-S plus Proliferation Score) on imputed Genetically-Regulated tumor eXpression (GReX). Bayesian  
10 multivariate regression and adaptive shrinkage tested GReX-prioritized genes for associations with tumor PAM50  
11 expression and subtype to elucidate patterns of germline regulation underlying GReX-CRS associations. At FDR-  
12 adjusted  $P < 0.10$ , 7 and 1 GReX-prioritized genes among WW and BW, respectively. Among WW, CRS were  
13 positively associated with *MCM10*, *FAM64A*, *CCNB2*, and *MMP1* GReX and negatively associated with *VAV3*,  
14 *PCSK6*, and *GNG11* GReX. Among BW, higher *MMP1* GReX predicted lower Proliferation score and ROR-P.  
15 GReX-prioritized gene and PAM50 tumor expression associations highlighted potential mechanisms for GReX-  
16 prioritized gene to CRS associations. Among BC patients, differential germline associations with CRS were found  
17 by race, underscoring the need for larger, diverse datasets in molecular studies of BC. These findings also  
18 suggest possible germline trans-regulation of PAM50 tumor expression, with potential implications for CRS  
19 interpretation in clinical settings.

20

21 **SIGNIFICANCE**

22 This study identifies race-specific genetic associations with breast cancer risk of recurrence scores and suggests  
23 mediation of these associations by PAM50 subtype and expression, with implications for clinical interpretation of  
24 these scores.

25

26 *Keywords:* breast cancer recurrence, risk of recurrence, transcriptome-wide association study, molecular subtype,  
27 trans-eQTL mapping

1 **ABBREVIATIONS**

2 BW Black Women

3 CBCS Carolina Breast Cancer Study

4 CRS Continuous Risk of recurrence Score

5 eQTL expression Quantitative Trait Locus

6 ER Estrogen Receptor

7 FDR False Discovery Rate

8 GReX Genetically-Regulated tumor eXpression

9 GWAS Genome-Wide Association Study

10 HR Hormone Receptor

11 LFSR Local False Sign Rate

12 LumA Luminal A

13 LumB Luminal B

14 NC North Carolina

15 ROR Risk of Recurrence

16 SCC Subtype-Centroid Correlations

17 SNP Single Nucleotide Polymorphism

18 TCGA The Cancer Genome Atlas

19 TWAS Transcriptome-Wide Association Study

20 WW White Women

## 1 INTRODUCTION

2 Tumor expression-based molecular profiling has improved clinical classification of breast cancer (1-3). One tool is  
3 the PAM50 assay, which integrates tumor expression of 50 genes (derived from a set of 1,900 subtype-specific  
4 genes identified in microarray studies) to determine PAM50 intrinsic molecular subtypes: Luminal A (LumA),  
5 Luminal B (LumB), Human epidermal growth factor 2-enriched (HER2-enriched), Basal-like, and Normal-like (1,4).  
6 Intrinsic molecular subtypes are strong prognostic factors for breast cancer outcomes, including recurrence and  
7 mortality. For instance, Basal-like breast cancer has substantially higher recurrence and mortality risk compared  
8 to LumA breast cancer (5-8). In recent years, continuous risk of recurrence scores (CRS) have gained traction as  
9 a potential clinical tool that encapsulates prognostic differences of breast cancer intrinsic molecular subtypes into  
10 a singular measure that can be used to guide treatment decisions. CRS include ROR-S, Proliferation score, ROR-  
11 P, and ROR-PT (1,9). ROR-P, for instance, is determined by combining ROR-S (PAM50 tumor expression-based  
12 subtype score) and Proliferation score (tumor expression of 11 PAM50 genes). ROR-PT further integrates ROR-P  
13 with information on tumor size. Studies show that CRS offer significant prognostic information beyond clinical  
14 variables (e.g., nodal status, tumor grade, age, hormonal therapy), improve adjuvant treatment decisions, and  
15 offer robust risk stratification for distant (5-10 years post diagnosis) recurrence (10-12).

16  
17 In the Carolina Breast Cancer Study (CBCS), Black women (BW) with breast cancer have disproportionately  
18 higher CRS than White Women (9), and similar disparities have been noted in Oncotype Dx recurrence score  
19 (9,13). Systemic injustices, like disparities in healthcare access, explain a substantial proportion of breast cancer  
20 outcome disparities (14-17). Recent studies additionally suggest that germline genetic variation is associated with  
21 breast cancer outcomes, and these associations vary across ancestry groups (18-21). In The Cancer Genome  
22 Atlas (TCGA), BW had substantially higher polygenic risk scores for the more aggressive ER-negative subtype  
23 than WW, suggesting differential genetic contributions for susceptibility for breast cancer, especially ER-negative  
24 breast cancer (21). In a transcriptome-wide association study (TWAS) of breast cancer mortality, germline-  
25 regulated gene expression (GRex) of four genes was associated with mortality among BW and gene expression  
26 for no genes was associated among WW (18). However, the role of germline genetic variation in recurrence,  
27 CRS, and CRS disparity remains a critical knowledge gap. Studying genetic associations with breast cancer

1 outcomes in BW is necessary to ensure advancements in breast cancer genetics are not limited to or  
2 generalizable in only White populations, thus aiding in decreasing health disparities.

3  
4 As racially-diverse genetic datasets typically have small samples of BW, gene-level association tests can increase  
5 study power. These approaches include TWAS, which integrates relationships between single nucleotide  
6 polymorphisms (SNP) and gene expression with genome-wide association studies (GWAS) to prioritize gene-trait  
7 associations (22,23). TWAS aids in interpreting genetic associations by mapping significant GWAS associations  
8 to tissue-specific expression of individual genes. In cancer applications, TWAS has identified susceptibility genes  
9 at loci previously undetected through GWAS, highlighting its improved power and interpretability (24-26). Previous  
10 studies show that stratification of the entire TWAS (model training, imputation, and association testing) is  
11 preferable in diverse populations, as models may perform poorly across ancestry groups and methods for TWAS  
12 in admixed populations are unavailable (18,27).

13  
14 Here, using data from the CBCS, which includes a large sample of Black breast cancer patients with tumor gene  
15 expression data, we study race-specific germline genetic associations for CRS using a gene-based association  
16 testing approach that borrows from TWAS methodology. CRS included in this study are ROR-S, Proliferation  
17 score, and ROR-P. Using race-specific predictive models for tumor expression from germline genetics, we identify  
18 sets of GReX-prioritized genes (i.e. genes whose GReX is associated with CRS) across BW and WW. We  
19 additionally investigate ROR-P specific GReX-prioritized genes for associations with PAM50 subtype and  
20 subtype-specific tumor gene expression to elucidate germline contributions to PAM50 subtype, and how these  
21 mediate GReX-prioritized gene and CRS associations. Unlike previous studies that correlated tumor gene  
22 expression (as opposed to germline-regulated tumor gene expression) with subtype or subtype-specific tumor  
23 gene expression, TWAS enables directional interpretation of observed associations (22,23).

24

## 25 **MATERIALS AND METHODS**

### 26 ***Data collection***

### 27 *Study population*

1 The CBCS is a population-based study of North Carolina (NC) breast cancer patients, enrolled in three phases;  
2 study details have been previously described (28,29). Patients aged 20 to 74 were identified using rapid case  
3 ascertainment with the NC Central Cancer Registry with randomized recruitment to oversample self-identified  
4 Black and young women (ages 20-49) (9,29). Demographic and clinical data (age, menopausal status, body mass  
5 index, hormone receptor status, tumor stage, study phase, recurrence) were obtained through questionnaires and  
6 medical records. The study was approved by the Office of Human Research Ethics at the University of North  
7 Carolina at Chapel Hill, and written informed consent was obtained from each participant.

8

9 *CBCS genotype data*

10 Genotypes were assayed on the OncoArray Consortium's custom SNP array (Illumina Infinium OncoArray) (30)  
11 and imputed using the 1000 Genomes Project (v3) as a reference panel for two-step phasing and imputation  
12 using SHAPEIT2 and IMPUTEv2 (31-34). The DCEG Cancer Genomics Research Laboratory conducted  
13 genotype calling, quality control, and imputation (30). We excluded variants with less than 1% minor allele  
14 frequency and deviations from Hardy-Weinberg equilibrium at  $P < 10^{-8}$  (35,36). We intersected genotyping  
15 panels for BW and WW samples, resulting in 5,989,134 autosomal variants and 334,391 variants on the X  
16 chromosome (37). We only consider the autosomal variants in this study.

17

18 *CBCS gene expression data*

19 Paraffin-embedded tumor blocks were assayed for gene expression of 406 breast cancer-related and 11  
20 housekeeping genes using NanoString nCounter at the Translational Genomics Laboratory at UNC-Chapel Hill  
21 (4,9). These 406 breast cancer-related genes include genes part of the PAM50, P53, E2, IGF, and EGFR  
22 signatures, among others (**Supplementary Table S1**). As described previously, we eliminated samples with  
23 insufficient data quality using NanoStringQCPro (18,38), scaled distributional difference between lanes with  
24 upper-quartile normalization (39), and removed two dimensions of unwanted technical and biological variation,  
25 estimated from housekeeping genes using RUVSeq (39,40). The current analysis included 1,199 samples with  
26 both genotype and gene expression data (628 BW, 571 WW).

27

28 **Statistical analysis**

## 1 *Overview of GReX and TWAS*

2 We adopted TWAS methodology to construct GReX (exposure of interest in this study). GReX for a given gene  
3 represents the portion of tumor expression explained by *cis*-genetic regulation; GReX was constructed for the  
4 aforementioned set of BC-related genes (**Supplementary Table S1**). Briefly, TWAS integrates expression data  
5 with GWAS to prioritize gene-level germline-trait associations through a two-step analysis (**Figure 1A-B**). First,  
6 using germline and transcriptomic data, we trained predictive models of tumor gene expression using all SNPs  
7 within 0.5 Megabase of the gene (18,23). Second, we used these models to impute the GReX of a gene by  
8 multiplying the SNP-gene weights from the predictive model with the dosages of each SNP. Associations between  
9 GReX (for a given gene) and trait (CRS, for instance) in regression analyses identify gene-trait relationships that  
10 are a consequence of germline variation. If sufficiently heritable genes are assayed in the correct tissue, TWAS-  
11 based GReX analyses increase power to detect germline-trait associations and aids interpretability of results, as  
12 associations are mapped from germline genetics to individual genes (23,41).

13

## 14 *GReX analysis of CRS in CBCS*

15 We adopted techniques from FUSION to train predictive models of tumor expression from *cis*-germline genotypes  
16 (18,23). Motivated by strong associations between germline genetics and tumor expression in CBCS (18), for  
17 genes with non-zero *cis*-heritability at nominal  $P < 0.10$ , we trained predictive models for covariate-residualized  
18 tumor expression with all *cis*-SNPs within 0.5 Megabase using linear mixed modeling or elastic net regression  
19 (**Supplementary Methods, Supplementary Materials**) (42,43). Here, we used the 628 BW samples and 571  
20 WW samples with both genotype and expression data to train these race-specific expression models. We  
21 selected models with five-fold cross-validation adjusted  $R^2 > 0.01$  between predicted and observed expression  
22 values, resulting in 59 and 45 models for WW and BW, respectively. Further details on these models, including  
23 heritability and cross-validation performance are available at **Supplementary Table S2**. These models also  
24 showed sufficiently strong predictive performance in external validation using TCGA data (18).

25

26 Using only germline genetics as input, we imputed GReX in 1,043 WW and 1,083 BW, respectively, in CBCS. For  
27 samples not present in the training dataset, we multiplied the SNP weights from the predictive models with the  
28 SNP dosages to construct GReX. For samples in both the training and imputation datasets, GReX was imputed



1 via cross-validation to minimize data leakage. We tested GReX for associations with ROR-S, Proliferation Score,  
2 and ROR-P using multiple linear regression adjusted for age, estrogen receptor (ER) status, tumor stage, and  
3 study phase (1). We corrected for test-statistic bias and inflation using a Bayesian bias and inflation adjustment  
4 method *bacon*, as TWAS are prone to bias and inflation of test statistics (44). We then adjusted for multiple  
5 testing using the Benjamini-Hochberg procedure (44,45). As a comparison for the germline effect of GReX-  
6 prioritized genes, we additionally assessed the effect of total (germline-regulated and post-transcriptional) tumor  
7 expression of those GReX-prioritized genes on CRS using similar linear models. We were underpowered to study  
8 time-to-recurrence, as recurrence data was collected only in CBCS Phase 3 (635 WW, 742 BW with GReX and  
9 recurrence data; 183 WW, 283 BW with tumor expression and recurrence data). For significant GReX-prioritized  
10 genes for CRS (FDR-adjusted  $P < 0.10$ ), we conducted follow-up permutation tests: we shuffle the SNP-gene  
11 weights in the predictive model 5,000 times to generate a null distribution and compare the original GReX-CRS  
12 associations to this null distribution. This permutation test assessed whether the GReX association provides more  
13 tissue-specific expression context, beyond any strong SNP-CRS associations at the genetic locus (23).

14

#### 15 *PAM50 assay and ROR-S, Proliferation score, and ROR-P calculation*

16 As described previously (1), using partition-around-medoid clustering, we calculated the correlation with each  
17 subtype's centroid for study individuals based on PAM50 expressions (10 PAM50 genes per subtype). The largest  
18 subtype-centroid correlation defined the individual's molecular subtype. ROR-S was determined via a linear  
19 combination of the PAM50 subtype-centroid correlations (SCCs); the coefficients to the PAM50 SCCs in the linear  
20 combination are positive for Luminal B, HER2-enriched, and Basal-like and negative for Luminal A (1).  
21 Proliferation score was computed using log-scale expression of 11 PAM50 genes, while ROR-P was computed by  
22 combining ROR-S and Proliferation score.

23

24 Assignment of PAM50 gene to subtype was based on PAM50 gene centroid values for each subtype; a PAM50  
25 gene is assigned to the subtype with the largest positive centroid value. Subtype assignment through this "greedy  
26 algorithm" are specific to this study and represent a simplified reality (e.g., *ESR1* classified as part of Luminal A  
27 subtype only even though *ESR1* expression correlates with both Luminal A and to a slightly lesser degree  
28 Luminal B subtype). Moreover, subtype assignment for this portion of analyses was conducted only for visual

1 comparison of patterns of associations between GReX-prioritized genes and PAM50 tumor gene expressions  
2 (i.e., subtype assignment in this portion of analyses had no bearing on continuous ROR score calculations or  
3 subtype-centroid correlations).

4

#### 5 *Bayesian multivariate regressions and multivariate adaptive shrinkage*

6 As previously noted (1), CRS are functions of PAM50 SCCs and gene expression profiles. Thus, we followed up  
7 on CRS-associated GReX-prioritized genes by studying their associations with PAM50 SCCs and gene  
8 expression. We assessed GReX-prioritized genes (for ROR-P) in relation to SCCs and PAM50 tumor gene  
9 expression (**Figure 1C**). Importantly, consistent with the original formulation of ROR-S, we did not consider  
10 normal-like subtype and normal-like subtype specific genes; subtype-specific genes were determined using a  
11 greedy assignment algorithm, described in the previous section. This classification scheme offers analytic  
12 simplicity but is an oversimplification for some PAM50 genes. We found that none of our GReX-prioritized genes  
13 were within 1 Megabase of PAM50 genes and that most GReX-prioritized genes were not on the same  
14 chromosome as PAM50 genes (**Supplementary Table S3**).

15

16 Existing gene-based mapping techniques for *trans*-expression quantitative trait loci (eQTL) (SNP and gene are  
17 separated by more than 1 Megabase) mapping include *trans*-PrediXcan and GBAT (46,47). We employed  
18 Bayesian multivariate linear regression (BtQTL) to account for correlation in multivariate outcomes (SCCs and  
19 PAM50 gene expression) in association testing. BtQTL improves power to detect significant *trans*-associations,  
20 especially when considering multiple genes with highly correlated (>0.5) expression (**Supplementary Figures**  
21 **S1-S2**). Lastly, we conducted adaptive shrinkage on BtQTL estimates using mashr, an empirical Bayes method to  
22 estimate patterns of similarity and improve accuracy in associations tests across multiple outcomes (48). mashr  
23 outputs revised posterior means, standard deviations, and corresponding measures of significance (local false  
24 sign rates, or LFSR).

25

#### 26 *Associations of genetic ancestry and race with tumor expression and GReX of GReX-prioritized genes*

27 Prior studies using CBCS have reported concordance between self-reported race and genetic ancestry (first  
28 principal component of combined genotype matrix) (49). In an effort to further contextualize CRS associations

1 across race and to disentangle race from genetic ancestry in our study population (specifically, whether race,  
2 which captures both genetic ancestry and socioeconomic context, is a proxy for genetic ancestry in our study  
3 population), we investigated: 1) association between genetic ancestry and tumor expression of GReX-prioritized  
4 genes; 2) association between genetic ancestry and GReX of GReX-prioritized genes; 3) association between  
5 race and tumor expression of GReX-prioritized genes; 4) association between race and GReX of GReX-prioritized  
6 genes. Genetic ancestry was computed by aggregating across local ancestry, as determined through the RFMix  
7 pipeline (50).

8

### 9 **Availability of data and materials**

10 Expression data from CBCS is available on NCBI GEO with accession number GSE148426. CBCS genotype  
11 datasets analyzed in this study are not publicly available as many CBCS patients are still being followed and  
12 accordingly CBCS data is considered sensitive; the data is available from M.A.T upon reasonable request.

13 Supplementary Data includes summary statistics for eQTL results, tumor expression models, and relevant R code  
14 for training expression models in CBCS and are freely available at [https://github.com/bhattacharya-a-  
15 bt/CBCS\\_TWAS\\_Paper/](https://github.com/bhattacharya-a-bt/CBCS_TWAS_Paper/). Scripts utilized in this analysis are provided at [https://github.com/APUNC/CBCS---Risk-  
16 of-Recurrence-Paper](https://github.com/APUNC/CBCS---Risk-of-Recurrence-Paper).

17

## 18 **RESULTS**

### 19 **Race-specific associations between GReX and CRS**

20 We performed race-specific GReX analysis for CRS to investigate the role of germline genetic variation in CRS  
21 and CRS racial disparity. We identified 8 genes (*MCM10*, *FAM64A*, *CCNB2*, *MMP1*, *VAV3*, *PCSK6*, *NDC80*,  
22 *MLPH*), 8 genes (*MCM10*, *FAM64A*, *CCNB2*, *MMP1*, *VAV3*, *NDC80*, *MLPH*, *EXO1*), and 10 genes (*MCM10*,  
23 *FAM64A*, *CCNB2*, *MMP1*, *VAV3*, *PCSK6*, *GNG11*, *NDC80*, *MLPH*, *EXO1*) whose GReX was associated with  
24 ROR-S, proliferation, and ROR-P, respectively, in WW, and 1 gene (*MMP1*) whose GReX was associated with  
25 proliferation and ROR-P in BW at FDR-adjusted  $P < 0.10$  (**Figure 2A, 2B**). No associations were detected  
26 between GReX and ROR-S among BW. We refer to genes with statistically significant GReX analysis  
27 associations (FDR-adjusted  $P < 0.10$ ) as GReX-prioritized genes. Among these identified genes, only genes that  
28 are not part of the PAM50 panel (i.e., excluding *NDC80*, *MLPH*, *EXO1*) were considered in downstream

1 permutation and GReX-prioritized gene follow up analyses (**Figure 1C**), as we wished to focus investigation on  
2 relationship between non-PAM50 GReX-prioritized genes and PAM50 (tumor) genes. **Supplementary Figure S3**  
3 shows results from a sensitivity analysis comparing the effect sizes for the GReX-CRS associations within  
4 samples used in training, not used in training, and the overall associations using all training and non-training  
5 samples. In general, we see concordance in the direction of association across these three splits of data, though  
6 some of the associations detected within only training or non-training samples intersect the null.

7

8 Among WW, increased GReX of *MCM10*, *FAM64A*, *CCNB2*, and *MMP1* were associated with higher CRS while  
9 increased GReX of *VAV3*, *PCSK6*, and *GNG11* were associated with lower CRS (**Figure 2A**).

10 Among BW, increased GReX of *MMP1* was associated with lower CRS (Proliferation, ROR-P, but not ROR-S)  
11 (**Figure 2A**). **Supplementary Figure S4** shows the nominal differences in eQTL architecture across BW and WW  
12 for these genes. In particular, for *MMP1*, we found differences in the standardized effects across WW and BW: a  
13 sizable proportion of shared eQTLs had discordant effects across WW and BW (**Supplementary Figure S5**). The  
14 LD structure for eQTLs differed across WW and BW, with eQTL effect size peaks ( $-\log_{10}$  p-values: 4.73 (WW);  
15 3.17 (BW)) at differing genomic locations (**Supplementary Figure S5**).

16

17 Briefly, to contextualize the functions of these GReX-prioritized genes, *MCM10* is involved in DNA replication,  
18 *FAM64A* and *CCNB2* are implicated in progression and regulation of the cell cycle, and *MMP1*, like the broader  
19 *MMP* family, is involved in the breakdown of the extracellular matrix (51-55). *GNG11* and *VAV3* are involved in  
20 signal transduction: *GNG11* as a component of a transmembrane G-protein and *VAV3* as a guanine nucleotide  
21 exchange factor for GTPases (56,57).

22

23 Associations between tumor expression of GReX-prioritized genes and CRS were concordant, in terms of  
24 direction of association to germline-only effects among WW; findings were discordant among BW where higher  
25 tumor expression of *MMP1* was associated with higher CRS (**Table 1, Supplementary Table S4**). We found  
26 differences in the pattern of associations between genetic ancestry and race with tumor expression and GReX of  
27 GReX-prioritized genes (**Supplementary Figure S6**). For instance, while higher African ancestry was associated

1 with higher tumor expression of *MCM10*, higher African ancestry was instead associated with lower GReX of  
2 *MCM10*.

3

#### 4 **Permutation testing provides context to GReX-prioritized gene and CRS associations**

5 To assess the statistical significance for the observed variance in CRS explained by significant GReX-prioritized  
6 genes, we conducted two permutation analyses. First, we assessed the per-gene significance of the GReX-CRS  
7 associations, conditional on the SNP-trait effects at the locus, by generating a null distribution for the GReX-CRS  
8 association via shuffling the SNP-gene weights from the predictive models 5,000 times. We generated a  
9 permutation P-value for the GReX-CRS association by comparing to this null distribution. Here, we found that all  
10 GReX-CRS associations showed significance in permutation testing at FDR-adjusted  $P < 0.05$  (**Table 1**). These  
11 per-GReX-prioritized gene permutation tests show that GReX (of GReX-prioritized genes) adds more context  
12 beyond the genetic architecture at the locus and provide evidence that germline genetics to GReX-prioritized  
13 gene expression relationship mediates, to some level, the complex genetic effects on CRS.

14

15 Next, we quantified the percent variance explained of CRS by the GReX-prioritized genes, in aggregate, by  
16 calculating the model adjusted- $R^2$  for a regression of covariate-residualized CRS on GReX all GReX-prioritized  
17 genes. To context these model adjusted- $R^2$ , we conducted two permutation tests. First, we permuted the sample  
18 labels for covariate-residualized CRS 10,000 times and computed the model adjusted  $R^2$  at each iteration to  
19 generate a null distribution for adjusted  $R^2$  between GReX-prioritized genes and CRS. Across WW and BW, the  
20 observed  $R^2$  of GReX-prioritized genes against CRS (7-10% among WW and 1% among BW) were statistically  
21 significant against the respective null distributions (P-values and distributions in **Figure 2B**). To further  
22 contextualize the proportion of variance in CRS explained by GReX-prioritized genes, we computed race-specific  
23 heritability estimates using GCTA (58). Given the limited sample size for which CRS data were available, we  
24 computed the heritability based on typed SNPs. Moreover, heritability estimates for CRS were stratified by race.  
25 Among WW, heritability ranged from 0.13 (SE: 0.23) for ROR-S to 0.21 (SE: 0.23) for Proliferation score. Among  
26 BW, heritability was much lower and ranged from 0.01 (SE: 0.12) for Proliferation score to 0.02 (SE: 0.14) for  
27 ROR-P. However, we note that heritability estimates from GCTA were imprecise due to limited sample size.

1 Permutation tests for analyses of tumor expression of GReX-prioritized genes and CRS are available in  
2 **Supplementary Figure S7.**

3  
4 Second, we wanted to assess if the GReX of these sets of GReX-prioritized genes (7 in WW and 1 in BW)  
5 explained more of the variance in CRS than the GReX of a randomly selected set of genes of the same size.  
6 Previous studies have shown that the tumor expression of a set randomly selected genes is likely to be predictive  
7 of breast cancer outcomes; we wished to investigate this phenomenon on the GReX level (59,60). Over 10,000  
8 repetitions, we randomly selected 7 and 1 genes in WW and BW subjects, respectively, ran a multivariable  
9 regression, and calculated the model adjusted- $R^2$  to generate another null distribution. Here again, we found that  
10 the true model  $R^2$  outperformed the null distribution, all showing permutation  $P < 0.05$  in these settings (**Figure**  
11 **2B**). These permutation tests show that our GReX-prioritized genes, taken together, appreciably explain  
12 differences in CRS.

#### 14 **Associations between GReX-prioritized genes and PAM50 subtype correlations and gene expression**

15 As CRS are constructed from PAM50 subtype-specific correlations and gene expression profiles, we further  
16 studied associations between GReX of GReX-prioritized genes and PAM50 SCCs and gene expression to  
17 understand how PAM50 subtype and gene expression mediate GReX-prioritized gene and CRS associations.  
18 Among WW, a one standard deviation increase in *FAM64A* and *CCNB2* GReX resulted in significantly increased  
19 Basal-like SCC while an identical increase in *VAV3*, *PCSK6*, and *GNG11* GReX resulted in significantly increased  
20 Luminal A SCC. The magnitude of increase in correlation for respective subtypes per GReX-prioritized gene was  
21 approximately 0.05, and most estimates had credible intervals that did not intersect the null. Among WW,  
22 associations between HER2-like SCC and GReX-prioritized genes followed similar patterns to associations for the  
23 Basal-like subtype, although associations for HER2 were more precise (**Figure 3A**). We found predominantly null  
24 associations between GReX-prioritized genes and Luminal B SCC among WW (**Figure 3A**). Unlike in WW, for  
25 BW, an increase in *MMP1* GReX was not associated with Luminal A, HER2 or Basal-like SCCs. Instead, among  
26 BW, *MMP1* GReX was significantly negatively associated with Luminal B SCC. Estimates from univariate  
27 regressions are provided in **Supplementary Tables S5-S8.**

28

1 For both WW and BW, the pattern of associations between GReX-prioritized genes and PAM50 tumor expression  
2 were predominantly congruent with observed associations between GReX-prioritized genes and PAM50 SCCs as  
3 well as GReX-prioritized genes and CRS (**Figure 3, Supplementary Tables S9-S12**). In WW, a one standard  
4 deviation increase in *CCNB2* GReX was associated with significantly increased *ORC6L*, *PTTG1*, and *KIF2C*  
5 (Basal-like genes) expression and *UBE2T* and *MYBL2* (LumB genes) expression. By contrast, a one standard  
6 deviation increase in *PCSK6* GReX significantly increased *BAG1*, *FOXA1*, *MAPT*, and *NAT1* (LumA genes)  
7 expression (**Figure 3B**). While increased *MMP1* GReX was associated with significantly increased expression of  
8 *ORC6L* (basal-like gene), *MYBL2*, and *BIRC5* (LumB genes) among WW, this was not the case among BW.  
9 Instead, increased *MMP1* GReX among BW was significantly associated with increased expression of *SLC39A6*  
10 (LumA gene) and decreased expression of *ACTR3B*, *PTTG1*, and *EXO1* (Basal-like genes) (**Figure 3B**).  
11 Associations between GReX-prioritized genes and PAM50 genes provide a granular, gene interaction level view  
12 into the mediation of the GReX-prioritized gene and CRS association, suggesting that *trans*-regulation of subtype-  
13 specific PAM50 genes by GReX-prioritized genes in breast tumors could be a possible contributor to subtypes  
14 and, subsequently, CRS and recurrence.

15

## 16 **DISCUSSION**

17 Through a GReX analysis, we identified 7 and 1 genes among WW and BW, respectively, for which genetically-  
18 regulated breast tumor expression was associated with CRS and underlying PAM50 gene expression and  
19 subtype. Among WW, these 7 GReX-prioritized genes explained between 7-10% of the variation in CRS, a large  
20 and statistically significant proportion of variance. Among BW, the singular GReX prioritized gene explained a  
21 statistically significant ~1% of the variation in Proliferation score and ROR-P. The magnitudes of these estimates  
22 were concordant with race-specific heritability estimates for CRS (13-21% for WW; 1-2% or BW) in this study  
23 population and suggest higher germline genetic contribution to CRS among WW compared to BW and as  
24 substantial contribution of GReX-prioritized genes to race-specific CRS heritability. There are two key novel  
25 aspects to this study. First, existing literature on associations between tumor gene expression and recurrence (for  
26 which CRS are a proxy) cannot distinguish between genetic and non-genetic components of effect (61), whereas,  
27 here, we estimate the contribution of the genetic component. Second, GReX analysis allows directional  
28 interpretation of observed associations that are not possible when correlating tumor gene expression and

1 recurrence. For instance, prior studies report *CCNB2* is upregulated in triple-negative breast cancers (TNBC) but  
2 were unable to determine whether increased *CCNB2* expression contributes to development or maintenance of  
3 TNBC or is part of the molecular response to cancer progression (62,63). By contrast, GReX is a function of only  
4 genetic variation. As such, we can confidently rule out that differences in *CCNB2* GReX are not direct  
5 consequences of subtype (and by extension recurrence); however, our observed associations of *CCNB2* GReX  
6 and subtype suggest a potential directional relationship for further study. Thus, GReX analysis allows a  
7 directional, potentially causal interpretation, subject to effective control for population stratification, minimal  
8 horizontal pleiotropy, and assumptions of independent assortment of alleles (22,23).

9  
10 Our GReX-prioritized gene and subtype associations among WW are consistent with literature on the association  
11 between tumor (i.e., genetic and non-genetic) expression of our GReX-prioritized genes and subtype. Prior  
12 investigations in cohorts of primarily European ancestry have reported that *MCM10*, *FAM64A*, and *CCNB2*  
13 expression is higher in ER-negative compared to ER-positive tumors (62-64). In studies that compared triple-  
14 negative and non-triple negative subtypes, higher *MCM10*, *FAM64A*, and *CCNB2* expression was detected in  
15 triple-negative breast cancer (62,63). Histologically, HER2-enriched and Basal-like subtypes are typically ER-  
16 negative, and triple-negatives are similar to Basal-like subtypes (9,65). Moreover, our findings among WW that  
17 GReX of *PCSK6* and *VAV3* associated with LumA subtype and LumA-specific gene expression are also  
18 consistent with previous results of *PCSK6* and *VAV3* upregulation in ER-positive subtypes (66,67). Importantly,  
19 our associations suggest directional mechanisms: from germline variation, to GReX of GReX-prioritized gene, and  
20 ultimately, to subtype.

21  
22 Presently, little is known about germline genetic regulation of PAM50 tumor expression. In CBCS, we found that  
23 tumor expression of most PAM50 genes is not *cis*-heritable (18). Instead, observed GReX-prioritized gene and  
24 PAM50 gene expression associations may implicate *trans*-gene regulation of the PAM50 signature. For instance,  
25 we found that *VAV3* GReX is significantly positively associated with tumor expression of *BAG1*, *FOXA1*, *MAPT*,  
26 and *NAT1* and nominally with increased tumor *ESR1* expression, all of which correspond well with LumA  
27 signature. Such *trans*-genic regulation signals, especially in the case of *ESR1*, pose significant clinical and  
28 therapeutic implication if confirmed under experimental conditions. For example, *VAV3* has been shown to



1 activate *RAC1*, which upregulates *ESR1* (68,69), but such mechanistic evidence is sparse for other putative  
2 GReX-prioritized gene to PAM50 associations. More generally, two of the GReX-prioritized genes among WW  
3 have been found to activate transcription factors; *FAM64A* enhances oncogenic nuclear factor-kappa B (NF-κB)  
4 signaling while both *FAM64A* and *PCSK6* activate oncogenic *STAT3* signaling (70-72).

5  
6 Interestingly, we found *MMP1* GReX has divergent associations with CRS across race. There are a few potential  
7 explanations. While heritability and proportion of variance in *MMP1* expression were similar across WW and BW  
8 predictive models, we found that the range of *MMP1* GReX was manifold among WW than BW. Potential  
9 differences in influence of germline genetics on tumor expression and CRS by race could be an artifact of  
10 divergent somatic or epigenetic factors that CBCS has not assayed (73-76). Second, while studies generally  
11 report that *MMP1* tumor expression is higher in triple-negative and Basal-like breast cancer, one study reported  
12 that *MMP1* expression in tumor cells does not significantly differ by subtype (77-79). Instead, Bostrom *et al.*  
13 reported that *MMP1* expression differs in stromal cells of patients with different subtypes (79). There is evidence  
14 to suggest that tumor composition, including stromal and immune components, may influence breast cancer  
15 progression in a subtype-specific manner. Future studies should consider expression predictive models that  
16 integrate greater detail on tumor cell-type composition to disentangle potential race-specific tumor composition  
17 effects on race-specific GReX associations (80,81).

18  
19 In this study, race (derived from self-report) captures genetic ancestry and additionally, socioeconomic context.  
20 Prior investigations using CBCS data have reported concordance between self-reported race and the first  
21 principal component of the combined (i.e. WW and BW) genotype matrix. In our analysis of local-ancestry derived  
22 global ancestry estimates and self-reported race, we found a similar, high level of concordance. In the absence of  
23 available methods that allow stratification or adjustments based on genetic ancestry across the GReX analytic  
24 framework, the use of race as a stratifying variable is intended to serve as a proxy for stratification by genetic  
25 ancestry. We acknowledge the limitation that race may not be a viable proxy across other populations outside  
26 CBCS, and that it is challenging to parse effects seen across race into effects of genetic ancestry and effects of  
27 socioeconomic context.

28

1 We found marked differences in the pattern of associations between genetic ancestry and race with tumor  
2 expression and GReX of GReX-prioritized genes, highlighting potential differences in contributions of germline  
3 and non-germline components to tumor expression across European and African ancestry groups. One particular  
4 example is *MCM10*. In the literature, higher *MCM10* tumor expression is correlated with Basal-like subtype, which  
5 is more prevalent among BW. The spectrum of our observations suggest that higher *MCM10* tumor expression is  
6 associated with Basal-like subtype across both BW and WW, but that the germline-regulated component of this  
7 expression may be stronger among WW. Similar patterns were seen for *FAM64A* and *CCNB2*. Analyses by race  
8 instead of genetic ancestry yielded associations similar in magnitude and direction. Racial differences in non-  
9 germline components of tumor expression, including tumor methylation and somatic alternations, may partly  
10 explain race-specific differences in GReX-prioritized genes (18,73-76,82,83). Other factors that warrant further  
11 investigation include potential greater contribution of *trans*-regulation in tumor gene expression in BW (methods  
12 for capturing *trans*-regulation in gene expression predictive models are not as well-developed as those for *cis*-  
13 regulation) (18). These factors should be investigated further as transcriptomic and epigenomic datasets for  
14 racially-diverse cohorts of breast cancer patients become available.

15  
16 There are a few limitations to this study. First, as CBCS used a Nanostring nCounter probeset for mRNA  
17 expression quantification of genes relevant for breast cancer, we could not analyze the whole human  
18 transcriptome. While this probeset may exclude several *cis*-heritable genes, CBCS contains one of the largest  
19 breast tumor transcriptomic datasets for Black women, allowing us to build well-powered race-specific predictive  
20 models, a pivotal step in ancestry-specific GReX analysis. Second, CBCS lacked data on somatic amplifications  
21 and deletions, inclusion of which could enhance the performance of predictive models of tumor expression (84).  
22 Third, as recurrence data was collected in a small subset with few recurrence events, we were unable to make a  
23 direct comparison between CRS and recurrence results, which may affect clinical generalizability. However, to our  
24 knowledge, CBCS is the largest resource of PAM50-based CRS data.

25  
26 Our analysis provides evidence of race-specific putative germline associations to CRS, mediated through  
27 associations between genetically-regulated tumor expression of GReX-prioritized genes and PAM50 expressions  
28 and subtype. This work underscores the need for larger and more diverse cohorts for genetic epidemiology

1 studies of breast cancer. Future studies should consider subtype-specific genetics (i.e., stratification by subtype in  
2 predictive model training and association analyses) to elucidate heritable gene expression effects on breast  
3 cancer outcomes both across and within subtype, which may yield further hypotheses for more fine-tuned clinical  
4 intervention.

5

## 6 **ACKNOWLEDGEMENTS**

7 We thank the Carolina Breast Cancer Study participants and volunteers. We also thank Colin Begg, Jianwen Cai,  
8 Katherine Hoadley, Yun Li, and Bogdan Pasaniuc for valuable discussion during the research process. We thank  
9 Erin Kirk and Jessica Tse for their invaluable support during the research process. We thank the DCEG Cancer  
10 Genomics Research Laboratory and acknowledge the support from Stephen Chanock, Rose Yang, Meredith  
11 Yeager, Belynda Hicks, and Bin Zhu.

12

## 13 **FUNDING**

14 This work was supported by Susan G. Komen® for the Cure for CBCS study infrastructure. Funding was provided  
15 by the National Institutes of Health, National Cancer Institute P01-CA151135, P50-CA05822, and U01-CA179715  
16 to AFO, CMP, and MAT. AP is supported by T32ES007018. MIL is supported by R01-HG009937, R01-  
17 MH118349, P01-CA142538, and P30-ES010126. The Translational Genomics Laboratory is supported in part by  
18 grants from the National Cancer Institute (3P30CA016086) and the University of North Carolina at Chapel Hill  
19 University Cancer Research Fund. Genotyping was done at the DCEG Cancer Genomics Research Laboratory  
20 using funds from the NCI Intramural Research Program. This content is solely the responsibility of the authors  
21 and does not necessarily represent the official views of the National Institutes of Health. The funder had no role in  
22 study design, data collection, analysis or interpretation, or writing of the manuscript.

23

24 Funding for BCAC came from: Cancer Research UK [grant numbers C1287/A16563,  
25 C1287/A10118, C1287/A10710, C12292/A11174, C1281/A12014, C5047/A8384, C5047/A15007, C5047/A10692,  
26 C8197/A16565], the European Union's Horizon 2020 Research and Innovation

1 Programme (grant numbers 634935 and 633784 for BRIDGES and B-CAST respectively), the European  
2 Community's Seventh Framework Programme under grant agreement n° 223175 [HEALTHF2-2009-223175]  
3 (COGS), the National Institutes of Health [CA128978] and Post-Cancer GWAS initiative [1U19  
4 CA148537, 1U19 CA148065-01 (DRIVE) and 1U19 CA148112 - the GAME-ON initiative], the Department of  
5 Defence [W81XWH-10-1-0341], and the Canadian Institutes of Health Research CIHR) for the CIHR Team in  
6 Familial Risks of Breast Cancer [grant PSR-SIIRI-701]. All studies and funders as listed in Michailidou K *et al*  
7 (2013 and 2015) and in Guo Q *et al* (2015) are acknowledged for their contributions.

8

### 9 **AUTHOR CONTRIBUTIONS**

10 Conceptualization: AP, MAT, MIL, AB. Data curation: MG, AFO, CMP, MAT. Formal analysis: AP, MAT, MIL, AB.  
11 Funding acquisition: AP, MG, AFO, CMP, MAT, MIL. Methodology: AP, MIL, AB. Project administration: MAT,  
12 MIL, AB. Resources: MG, AFO, CMP, MAT, MIL. Supervision: MAT, MIL, AB. Visualization: AP, AB. Writing –  
13 original draft: AP, AB. Writing – reviewing and editing: AP, MG, AFO, CMP, MAT, MIL, AB.

14

### 15 **REFERENCES**

- 16 1. Parker JS, Mullins M, Cheang MC, Leung S, Voduc D, Vickery T, *et al*. Supervised risk predictor of breast  
17 cancer based on intrinsic subtypes. *J Clin Oncol* **2009**;27:1160-7
- 18 2. Wallden B, Storhoff J, Nielsen T, Dowidar N, Schaper C, Ferree S, *et al*. Development and verification of  
19 the PAM50-based Prosigna breast cancer gene signature assay. *BMC Med Genomics* **2015**;8:54
- 20 3. Paik S, Shak S, Tang G, Kim C, Baker J, Cronin M, *et al*. A multigene assay to predict recurrence of  
21 tamoxifen-treated, node-negative breast cancer. *N Engl J Med* **2004**;351:2817-26
- 22 4. Geiss GK, Bumgarner RE, Birditt B, Dahl T, Dowidar N, Dunaway DL, *et al*. Direct multiplexed  
23 measurement of gene expression with color-coded probe pairs. *Nat Biotechnol* **2008**;26:317-25
- 24 5. Carey LA, Perou CM, Livasy CA, Dressler LG, Cowan D, Conway K, *et al*. Race, breast cancer subtypes,  
25 and survival in the Carolina Breast Cancer Study. *Jama* **2006**;295:2492-502
- 26 6. O'Brien KM, Cole SR, Tse CK, Perou CM, Carey LA, Foulkes WD, *et al*. Intrinsic breast tumor subtypes,  
27 race, and long-term survival in the Carolina Breast Cancer Study. *Clin Cancer Res* **2010**;16:6100-10

- 1 7. Shim HJ, Kim SH, Kang BJ, Choi BG, Kim HS, Cha ES, *et al.* Breast cancer recurrence according to  
2 molecular subtype. *Asian Pac J Cancer Prev* **2014**;15:5539-44
- 3 8. van Maaren MC, de Munck L, Strobbe LJA, Sonke GS, Westenend PJ, Smidt ML, *et al.* Ten-year  
4 recurrence rates for breast cancer subtypes in the Netherlands: A large population-based study. *Int J*  
5 *Cancer* **2019**;144:263-72
- 6 9. Troester MA, Sun X, Allott EH, Geradts J, Cohen SM, Tse CK, *et al.* Racial Differences in PAM50  
7 Subtypes in the Carolina Breast Cancer Study. *J Natl Cancer Inst* **2018**;110:176-82
- 8 10. Dowsett M, Sestak I, Lopez-Knowles E, Sidhu K, Dunbier AK, Cowens JW, *et al.* Comparison of PAM50  
9 Risk of Recurrence Score With Oncotype DX and IHC4 for Predicting Risk of Distant Recurrence After  
10 Endocrine Therapy. *Journal of Clinical Oncology* **2013**;31:2783-90
- 11 11. Sestak I, Buus R, Cuzick J, Dubsy P, Kronenwett R, Denkert C, *et al.* Comparison of the Performance of  
12 6 Prognostic Signatures for Estrogen Receptor-Positive Breast Cancer: A Secondary Analysis of a  
13 Randomized Clinical Trial. *JAMA Oncol* **2018**;4:545-53
- 14 12. Ohnstad HO, Borgen E, Falk RS, Lien TG, Aaserud M, Sveli MAT, *et al.* Prognostic value of PAM50 and  
15 risk of recurrence score in patients with early-stage breast cancer with long-term follow-up. *Breast Cancer*  
16 *Res* **2017**;19:120
- 17 13. Albain KS, Gray RJ, Makower DF, Faghih A, Hayes DF, Geyer CE, *et al.* Race, ethnicity and clinical  
18 outcomes in hormone receptor-positive, HER2-negative, node-negative breast cancer in the randomized  
19 TAILORx trial. *J Natl Cancer Inst* **2020**
- 20 14. Reeder-Hayes KE, Anderson BO. Breast Cancer Disparities at Home and Abroad: A Review of the  
21 Challenges and Opportunities for System-Level Change. *Clin Cancer Res* **2017**;23:2655-64
- 22 15. Durham DD, Robinson WR, Lee SS, Wheeler SB, Reeder-Hayes KE, Bowling JM, *et al.* Insurance-Based  
23 Differences in Time to Diagnostic Follow-up after Positive Screening Mammography. *Cancer Epidemiol*  
24 *Biomarkers Prev* **2016**;25:1474-82
- 25 16. Wheeler SB, Reeder-Hayes KE, Carey LA. Disparities in breast cancer treatment and outcomes:  
26 biological, social, and health system determinants and opportunities for research. *Oncologist*  
27 **2013**;18:986-93

- 1 17. Ko NY, Hong S, Winn RA, Calip GS. Association of Insurance Status and Racial Disparities With the  
2 Detection of Early-Stage Breast Cancer. *JAMA Oncology* **2020**;6:385-92
- 3 18. Bhattacharya A, García-Closas M, Olshan AF, Perou CM, Troester MA, Love MI. A framework for  
4 transcriptome-wide association studies in breast cancer in diverse study populations. *Genome Biol*  
5 **2020**;21:42
- 6 19. Escala-Garcia M, Guo Q, Dörk T, Canisius S, Keeman R, Dennis J, *et al.* Genome-wide association study  
7 of germline variants and breast cancer-specific mortality. *Br J Cancer* **2019**;120:647-57
- 8 20. Muranen TA, Khan S, Fagerholm R, Aittomäki K, Cunningham JM, Dennis J, *et al.* Association of  
9 germline variation with the survival of women with BRCA1/2 pathogenic variants and breast cancer. *NPJ*  
10 *Breast Cancer* **2020**;6:44
- 11 21. Huo D, Hu H, Rhie SK, Gamazon ER, Cherniack AD, Liu J, *et al.* Comparison of Breast Cancer Molecular  
12 Features and Survival by African and European Ancestry in The Cancer Genome Atlas. *JAMA Oncol*  
13 **2017**;3:1654-62
- 14 22. Gamazon ER, Wheeler HE, Shah KP, Mozaffari SV, Aquino-Michaels K, Carroll RJ, *et al.* A gene-based  
15 association method for mapping traits using reference transcriptome data. *Nat Genet* **2015**;47:1091-8
- 16 23. Gusev A, Ko A, Shi H, Bhatia G, Chung W, Penninx BW, *et al.* Integrative approaches for large-scale  
17 transcriptome-wide association studies. *Nat Genet* **2016**;48:245-52
- 18 24. Zhong J, Jermusyk A, Wu L, Hoskins JW, Collins I, Mocchi E, *et al.* A Transcriptome-Wide Association  
19 Study Identifies Novel Candidate Susceptibility Genes for Pancreatic Cancer. *J Natl Cancer Inst*  
20 **2020**;112:1003-12
- 21 25. Wu L, Shi W, Long J, Guo X, Michailidou K, Beesley J, *et al.* A transcriptome-wide association study of  
22 229,000 women identifies new candidate susceptibility genes for breast cancer. *Nat Genet* **2018**;50:968-  
23 78
- 24 26. Mancuso N, Gayther S, Gusev A, Zheng W, Penney KL, Kote-Jarai Z, *et al.* Large-scale transcriptome-  
25 wide association study identifies new prostate cancer risk regions. *Nat Commun* **2018**;9:4079
- 26 27. Keys KL, Mak ACY, White MJ, Eckalbar WL, Dahl AW, Mefford J, *et al.* On the cross-population  
27 generalizability of gene expression prediction models. *PLoS Genet* **2020**;16:e1008927

- 1 28. Hair BY, Hayes S, Tse CK, Bell MB, Olshan AF. Racial differences in physical activity among breast  
2 cancer survivors: implications for breast cancer care. *Cancer* **2014**;120:2174-82
- 3 29. Newman B, Moorman PG, Millikan R, Qaqish BF, Geradts J, Aldrich TE, *et al.* The Carolina Breast  
4 Cancer Study: integrating population-based epidemiology and molecular biology. *Breast Cancer Res*  
5 *Treat* **1995**;35:51-60
- 6 30. Amos CI, Dennis J, Wang Z, Byun J, Schumacher FR, Gayther SA, *et al.* The OncoArray Consortium: A  
7 Network for Understanding the Genetic Architecture of Common Cancers. *Cancer Epidemiol Biomarkers*  
8 *Prev* **2017**;26:126-35
- 9 31. Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, *et al.* A global reference for human  
10 genetic variation. *Nature* **2015**;526:68-74
- 11 32. O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, *et al.* A general approach for  
12 haplotype phasing across the full spectrum of relatedness. *PLoS Genet* **2014**;10:e1004234
- 13 33. Delaneau O, Marchini J, Zagury JF. A linear complexity phasing method for thousands of genomes. *Nat*  
14 *Methods* **2011**;9:179-81
- 15 34. Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method for the next  
16 generation of genome-wide association studies. *PLoS Genet* **2009**;5:e1000529
- 17 35. Wigginton JE, Cutler DJ, Abecasis GR. A note on exact tests of Hardy-Weinberg equilibrium. *Am J Hum*  
18 *Genet* **2005**;76:887-93
- 19 36. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, *et al.* PLINK: a tool set for whole-  
20 genome association and population-based linkage analyses. *Am J Hum Genet* **2007**;81:559-75
- 21 37. Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, *et al.* dbSNP: the NCBI database of  
22 genetic variation. *Nucleic Acids Res* **2001**;29:308-11
- 23 38. Bhattacharya A, Hamilton AM, Furberg H, Pietzak E, Purdue MP, Troester MA, *et al.* An approach for  
24 normalization and quality control for NanoString RNA expression data. *Brief Bioinform* **2020**
- 25 39. Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biol*  
26 **2010**;11:R106
- 27 40. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with  
28 DESeq2. *Genome Biol* **2014**;15:550

- 1 41. Ding B, Cao C, Li Q, Wu J, Long Q. Power analysis of transcriptome-wide association study. *bioRxiv*  
2 **2020**:2020.07.19.211151
- 3 42. Endelman JB. Ridge Regression and Other Kernels for Genomic Selection with R Package rrBLUP. *The*  
4 *Plant Genome* **2011**;4
- 5 43. Friedman J, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear Models via Coordinate  
6 Descent. *J Stat Softw* **2010**;33:1-22
- 7 44. van Iterson M, van Zwet EW, Heijmans BT. Controlling bias and inflation in epigenome- and  
8 transcriptome-wide association studies using the empirical null distribution. *Genome Biol* **2017**;18:19
- 9 45. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to  
10 Multiple Testing. *Journal of the Royal Statistical Society Series B (Methodological)* **1995**;57:289-300
- 11 46. Wheeler HE, Ploch S, Barbeira AN, Bonazzola R, Andaleon A, Fotuhi Siahpirani A, *et al.* Imputed gene  
12 associations identify replicable trans-acting genes enriched in transcription pathways and complex traits.  
13 *Genetic Epidemiology* **2019**;43:596-608
- 14 47. Liu X, Mefford JA, Dahl A, He Y, Subramaniam M, Battle A, *et al.* GBAT: a gene-based association test  
15 for robust detection of trans-gene regulation. *Genome Biology* **2020**;21:211
- 16 48. Urbut SM, Wang G, Carbonetto P, Stephens M. Flexible statistical methods for estimating and testing  
17 effects in genomic studies with multiple conditions. *Nat Genet* **2019**;51:187-95
- 18 49. Bhattacharya A, García-Closas M, Olshan AF, Perou CM, Troester MA, Love MI. A framework for  
19 transcriptome-wide association studies in breast cancer in diverse study populations. *Genome Biology*  
20 **2020**;21:42
- 21 50. Maples BK, Gravel S, Kenny EE, Bustamante CD. RFMix: a discriminative modeling approach for rapid  
22 and robust local-ancestry inference. *American journal of human genetics* **2013**;93:278-88
- 23 51. Watase G, Takisawa H, Kanemaki MT. Mcm10 plays a role in functioning of the eukaryotic replicative  
24 DNA helicase, Cdc45-Mcm-GINS. *Curr Biol* **2012**;22:343-9
- 25 52. Zhao W-m, Coppinger JA, Seki A, Cheng X-l, Yates JR, Fang G. RCS1, a substrate of APC/C, controls  
26 the metaphase to anaphase transition. *Proceedings of the National Academy of Sciences*  
27 **2008**;105:13415-20



- 1 53. Daldello EM, Luong XG, Yang C-R, Kuhn J, Conti M. Cyclin B2 is required for progression through  
2 meiosis in mouse oocytes. *Development* **2019**;146:dev172734
- 3 54. Draetta G, Luca F, Westendorf J, Brizuela L, Ruderman J, Beach D. Cdc2 protein kinase is complexed  
4 with both cyclin A and B: evidence for proteolytic inactivation of MPF. *Cell* **1989**;56:829-38
- 5 55. Page-McCaw A, Ewald AJ, Werb Z. Matrix metalloproteinases and the regulation of tissue remodelling.  
6 *Nat Rev Mol Cell Biol* **2007**;8:221-33
- 7 56. Rao S, Lyons LS, Fahrenholtz CD, Wu F, Farooq A, Balkan W, *et al.* A novel nuclear role for the Vav3  
8 nucleotide exchange factor in androgen receptor coactivation in prostate cancer. *Oncogene* **2012**;31:716-  
9 27
- 10 57. Hossain MN, Sakemura R, Fujii M, Ayusawa D. G-protein gamma subunit GNG11 strongly regulates  
11 cellular senescence. *Biochem Biophys Res Commun* **2006**;351:645-50
- 12 58. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait analysis. *Am J*  
13 *Hum Genet* **2011**;88:76-82
- 14 59. Venet D, Dumont JE, Detours V. Most random gene expression signatures are significantly associated  
15 with breast cancer outcome. *PLoS Comput Biol* **2011**;7:e1002240
- 16 60. Shimoni Y. Association between expression of random gene sets and survival is evident in multiple  
17 cancer types and may be explained by sub-classification. *PLoS Comput Biol* **2018**;14:e1006026
- 18 61. Parada H, Jr., Sun X, Fleming JM, Williams-DeVane CR, Kirk EL, Olsson LT, *et al.* Race-associated  
19 biological differences among luminal A and basal-like breast cancers in the Carolina Breast Cancer  
20 Study. *Breast Cancer Res* **2017**;19:131
- 21 62. Prat A, Adamo B, Cheang MC, Anders CK, Carey LA, Perou CM. Molecular characterization of basal-like  
22 and non-basal-like triple-negative breast cancer. *Oncologist* **2013**;18:123-33
- 23 63. Zhang C, Han Y, Huang H, Min L, Qu L, Shou C. Integrated analysis of expression profiling data identifies  
24 three genes in correlation with poor prognosis of triple-negative breast cancer. *Int J Oncol* **2014**;44:2025-  
25 33
- 26 64. Mahadevappa R, Neves H, Yuen SM, Jameel M, Bai Y, Yuen HF, *et al.* DNA Replication Licensing  
27 Protein MCM10 Promotes Tumor Progression and Is a Novel Prognostic Biomarker and Potential  
28 Therapeutic Target in Breast Cancer. *Cancers (Basel)* **2018**;10

- 1 65. Hagemann IS. Molecular Testing in Breast Cancer: A Guide to Current Practices. Arch Pathol Lab Med  
2 **2016**;140:815-24
- 3 66. Thakkar AD, Raj H, Chakrabarti D, Ravishankar, Saravanan N, Muthuvelan B, *et al.* Identification of gene  
4 expression signature in estrogen receptor positive breast carcinoma. Biomark Cancer **2010**;2:1-15
- 5 67. Aguilar H, Urruticoechea A, Halonen P, Kiyotani K, Mushiroda T, Barril X, *et al.* VAV3 mediates resistance  
6 to breast cancer endocrine therapy. Breast Cancer Res **2014**;16:R53
- 7 68. Zeng L, Sachdev P, Yan L, Chan JL, Trenkle T, McClelland M, *et al.* Vav3 mediates receptor protein  
8 tyrosine kinase signaling, regulates GTPase activity, modulates cell morphology, and induces cell  
9 transformation. Mol Cell Biol **2000**;20:9212-24
- 10 69. Rosenblatt AE, Garcia MI, Lyons L, Xie Y, Maiorino C, Désiré L, *et al.* Inhibition of the Rho GTPase,  
11 Rac1, decreases estrogen receptor levels and is a novel therapeutic strategy in breast cancer. Endocr  
12 Relat Cancer **2011**;18:207-19
- 13 70. Xu Z-S, Zhang H-X, Li W-W, Ran Y, Liu T-T, Xiong M-G, *et al.* FAM64A positively regulates STAT3  
14 activity to promote Th17 differentiation and colitis-associated carcinogenesis. Proceedings of the National  
15 Academy of Sciences **2019**;116:10447-52
- 16 71. Jiang H, Wang L, Wang F, Pan J. Proprotein convertase subtilisin/kexin type 6 promotes in vitro  
17 proliferation, migration and inflammatory cytokine secretion of synovial fibroblast-like cells from  
18 rheumatoid arthritis via nuclear- $\kappa$ B, signal transducer and activator of transcription 3 and extracellular  
19 signal regulated 1/2 pathways. Mol Med Rep **2017**;16:8477-84
- 20 72. Jiang L, Ren L, Zhang X, Chen H, Chen X, Lin C, *et al.* Overexpression of PIMREG promotes breast  
21 cancer aggressiveness via constitutive activation of NF- $\kappa$ B signaling. EBioMedicine **2019**;43:188-200
- 22 73. Shang L, Smith JA, Zhao W, Kho M, Turner ST, Mosley TH, *et al.* Genetic Architecture of Gene  
23 Expression in European and African Americans: An eQTL Mapping Study in GENOA. Am J Hum Genet  
24 **2020**;106:496-512
- 25 74. Wang S, Dorsey TH, Terunuma A, Kittles RA, Ambs S, Kwabi-Addo B. Relationship between tumor DNA  
26 methylation status and patient characteristics in African-American and European-American women with  
27 breast cancer. PLoS One **2012**;7:e37928

- 1 75. Conway K, Edmiston SN, Tse CK, Bryant C, Kuan PF, Hair BY, *et al.* Racial variation in breast tumor  
2 promoter methylation in the Carolina Breast Cancer Study. *Cancer Epidemiol Biomarkers Prev*  
3 **2015**;24:921-30
- 4 76. Chen Y, Sadasivan SM, She R, Datta I, Taneja K, Chitale D, *et al.* Breast and prostate cancers harbor  
5 common somatic copy number alterations that consistently differ by race and are associated with survival.  
6 *BMC Med Genomics* **2020**;13:116
- 7 77. Wang QM, Lv L, Tang Y, Zhang L, Wang LF. MMP-1 is overexpressed in triple-negative breast cancer  
8 tissues and the knockdown of MMP-1 expression inhibits tumor cell malignant behaviors in vitro. *Oncol*  
9 *Lett* **2019**;17:1732-40
- 10 78. McGowan PM, Duffy MJ. Matrix metalloproteinase expression and outcome in patients with breast  
11 cancer: analysis of a published database. *Ann Oncol* **2008**;19:1566-72
- 12 79. Boström P, Söderström M, Vahlberg T, Söderström KO, Roberts PJ, Carpén O, *et al.* MMP-1 expression  
13 has an independent prognostic value in breast cancer. *BMC Cancer* **2011**;11:348
- 14 80. Acerbi I, Cassereau L, Dean I, Shi Q, Au A, Park C, *et al.* Human breast cancer invasion and aggression  
15 correlates with ECM stiffening and immune cell infiltration. *Integr Biol (Camb)* **2015**;7:1120-34
- 16 81. González LO, Corte MD, Junquera S, González-Fernández R, del Casar JM, García C, *et al.* Expression  
17 and prognostic significance of metalloproteases and their inhibitors in luminal A and basal-like  
18 phenotypes of breast carcinoma. *Hum Pathol* **2009**;40:1224-33
- 19 82. Gravel S. Population genetics models of local ancestry. *Genetics* **2012**;191:607-19
- 20 83. Nelson D, Kelleher J, Ragsdale AP, Moreau C, McVean G, Gravel S. Accounting for long-range  
21 correlations in genome-wide simulations of large cohorts. *PLoS Genet* **2020**;16:e1008619
- 22 84. Xia Y, Fan C, Hoadley KA, Parker JS, Perou CM. Genetic determinants of the molecular portraits of  
23 epithelial cancers. *Nat Commun* **2019**;10:5666

24  
25

1 **TABLES**

2 **Table 1:** Race-specific associations between germline-regulated tumor gene expression (GReX) of GReX-  
 3 prioritized genes and CRS. Effect estimates correspond to change in CRS per 1 standard deviation increase in  
 4 GReX, adjusted for age, estrogen receptor status, stage, and CBCS study phase. 95% confidence intervals of  
 5 effect sizes are provided. All GReX-prioritized gene and CRS pairs shown here showed overall association FDR-  
 6 adjusted  $P < 0.10$ , and FDR-adjusted permutation  $P < 0.05$  (across 5,000 permutations of the SNP-gene  
 7 weights). We also provide signatures that include these genes as reference (**Supplementary Table S1**).

Gene	Signature	WW (N = 1,043)			BW (N = 1,083)		
		ROR-S	Proliferation	ROR-P	ROR-S	Proliferation	ROR-P
<b>MCM10</b>	IGF	3.03 (1.73, 4.33)	0.06 (0.03, 0.08)	3.11 (1.72, 4.50)	-	-	-
<b>FAM64A</b>	IGF	2.57 (1.28, 3.86)	0.05 (0.02, 0.07)	2.64 (1.26, 4.02)	-	-	-
<b>CCNB2</b>	Estradiol	2.69 (1.40, 3.98)	0.05 (0.02, 0.08)	2.71 (1.33, 4.09)	-	-	-
<b>MMP1</b>	Estradiol	2.73 (1.45, 4.01)	0.05 (0.02, 0.07)	2.58 (1.21, 3.96)	-1.84 (-3.12, -0.56)	-0.04 (-0.07, -0.02)	-2.21 (-3.56, -0.87)
<b>VAV3</b>	Other	-2.22 (-3.51, -0.93)	-0.04 (-0.07, -0.02)	-2.40 (-3.79, -1.03)	-	-	-
<b>PCSK6</b>	IGF	-2.16 (-3.45, -0.88)	-0.03 (-0.06, 0.00)	-1.88 (-3.25, -0.50)	-	-	-
<b>GNG11</b>	Claudin-low	-1.27 (-2.56, 0.02)	-0.02 (-0.05, 0.00)	-1.42 (-2.80, -0.05)	-	-	-

8

9 **FIGURE LEGENDS**

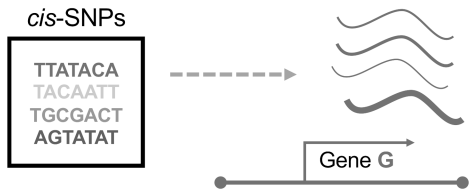
10 **Figure 1.** Schematic of study analytic approach. A) In CBCS, constructed race-stratified predictive models of  
 11 tumor gene expression from *cis*-SNPs. B) In CBCS, imputed GReX at individual-level using genotypes and tested  
 12 for associations between GReX and CRS in race-stratified linear models; only GReX of genes with significant *cis*-  
 13  $h^2$  and high cross validation performance ( $R^2 > 0.01$  between observed and predicted expression) considered for  
 14 race-stratified association analyses. C) Follow-up analyses on GReX-prioritized genes (i.e., genes whose GReX  
 15 were significantly associated with CRS at FDR  $< 0.10$ ). In race-stratified models, PAM50 SCCs and PAM50 tumor  
 16 expressions were regressed against GReX-prioritized genes under a Bayesian multivariate regression and  
 17 multivariate adaptive shrinkage approach.

18

1 **Figure 2.** *Permutation tests and associations between GReX-prioritized genes and CRS for WW and BW.* A) Effect estimates correspond to change in ROR-S, Proliferation score, and ROR-P per one standard deviation increase in GReX-prioritized gene expression (i.e., one standard deviation increase in GReX of gene). Triangle denotes WW and circle denotes BW. B) Boxplots correspond to null distributions (shuffled GReX-sample labels on left, random set of genes on right) of covariates residualized-R2 for regressions of CRS on GReX-prioritized genes. Null distributions are provided for 10,000 permutations of the GReX-sample labels and 10,000 random sets of genes. Dashed horizontal lines correspond to observed covariates residualized-R2.

8  
9 **Figure 3.** *Associations between GReX-prioritized genes and PAM50 SCCs and gene expression.* A) Among BW (top) and WW (bottom), associations between GReX-prioritized genes and PAM50 SCCs using Bayesian multivariate regression and multivariate adaptive shrinkage. Effect estimates show change in SCCs (range -1 to 1) for one standard deviation increase in GReX-prioritized gene GReX. Circle, triangle, and square denote corresponding LFSR intervals for effect sizes. B) Heatmap of change in  $\log_2$  normalized PAM50 tumor expression for one standard deviation increase in GReX-Prioritized gene GReX. \*, \*\*, \*\*\* denote FDR intervals for effect sizes.

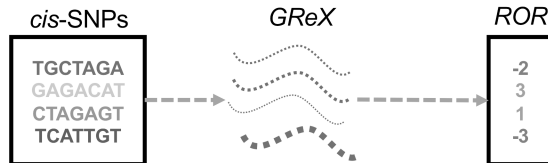
Figure 1

**A. Train expression models with CBCS**

Train race-stratified models of breast tumor expression using FUSION (Gusev *et al* 2016, Nature Genetics)

**B. Perform race-stratified GReX analysis of ROR scores**

Impute GReX (genetically-regulated expression) using individual genotypes from CBCS  
 Test for associations between GReX and ROR-scores

**C. Follow-up on GReX-prioritized genes**

Bayesian multivariate regression analysis to find associations between GReX and PAM50 gene expression and subtype correlations



50 PAM50 genes

1.4  
 5.2  
 0.8  
 2.2

1.4  
 5.2  
 0.8  
 2.2

PAM50 gene Expression  
(assayed in tumor)

5 PAM50 subtypes

0.7  
 -0.4  
 0.3  
 0.5

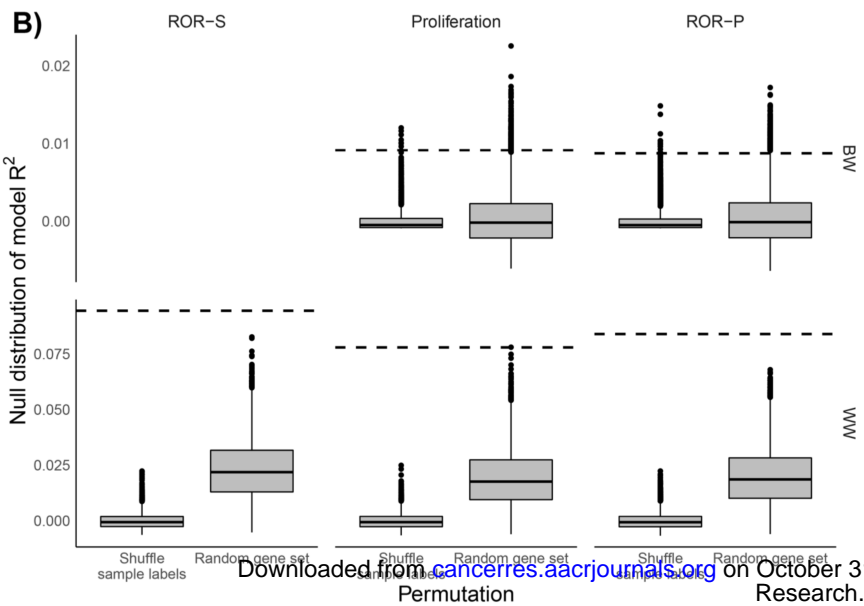
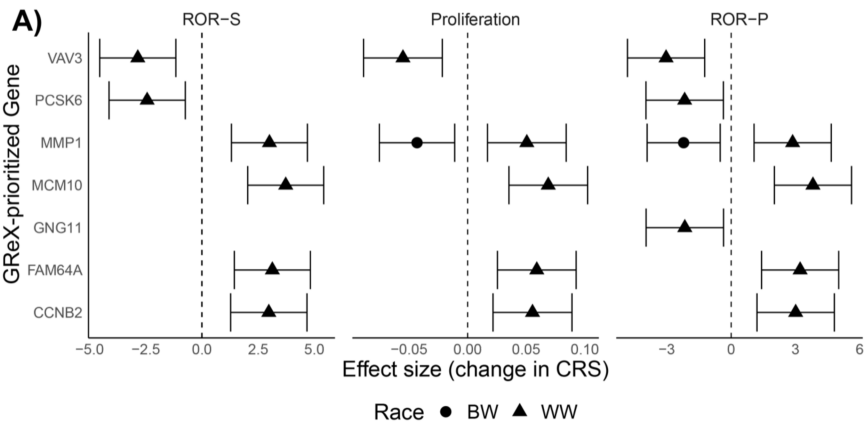
0.6  
 -0.1  
 -0.6  
 0.8

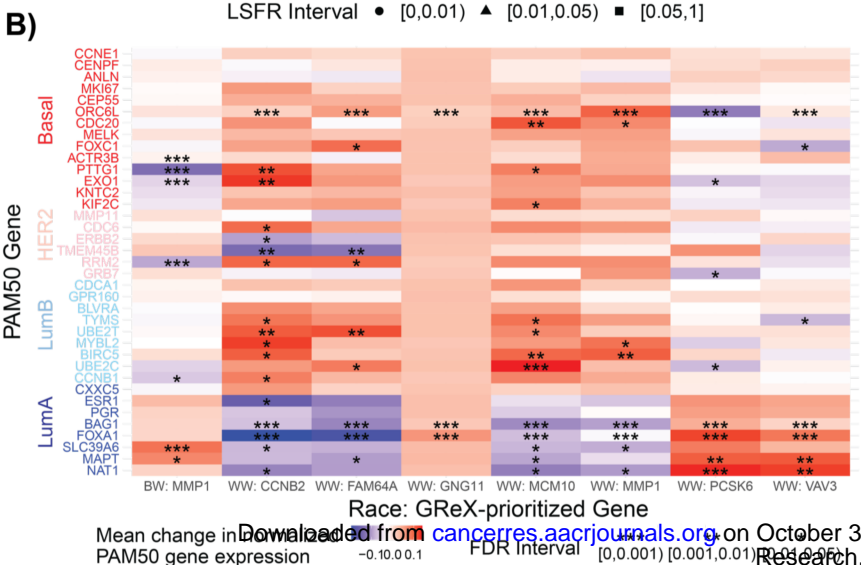
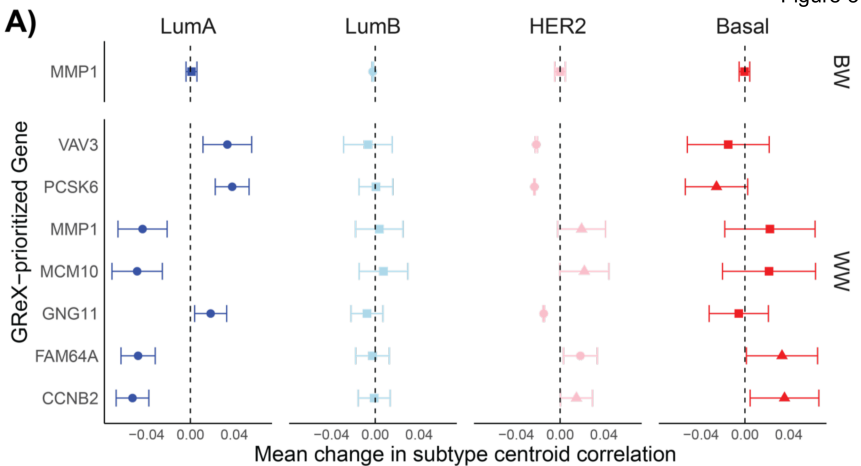
PAM50 subtype correlations  
(from tumor expression)

-2  
 3  
 1  
 -3

ROR  
(from subtype correlations)

Figure 2







# Cancer Research

The Journal of Cancer Research (1916–1930) | The American Journal of Cancer (1931–1940)

## Gene level germline contributions to clinical risk of recurrence scores in Black and White breast cancer patients

Achal Patel, Montserrat Garcia-Closas, Andrew F Olshan, et al.

*Cancer Res* Published OnlineFirst October 28, 2021.

<b>Updated version</b>	Access the most recent version of this article at: doi: <a href="https://doi.org/10.1158/0008-5472.CAN-21-1207">10.1158/0008-5472.CAN-21-1207</a>
<b>Supplementary Material</b>	Access the most recent supplemental material at: <a href="http://cancerres.aacrjournals.org/content/suppl/2021/10/28/0008-5472.CAN-21-1207.DC1">http://cancerres.aacrjournals.org/content/suppl/2021/10/28/0008-5472.CAN-21-1207.DC1</a>
<b>Author Manuscript</b>	Author manuscripts have been peer reviewed and accepted for publication but have not yet been edited.

**E-mail alerts** [Sign up to receive free email-alerts](#) related to this article or journal.

**Reprints and Subscriptions** To order reprints of this article or to subscribe to the journal, contact the AACR Publications Department at [pubs@aacr.org](mailto:pubs@aacr.org).

**Permissions** To request permission to re-use all or part of this article, use this link <http://cancerres.aacrjournals.org/content/early/2021/10/27/0008-5472.CAN-21-1207>. Click on "Request Permissions" which will take you to the Copyright Clearance Center's (CCC) Rightslink site.