

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

3 Achal Patel, MPH¹, Montserrat García-Closas, MD, DrPH^{2,3}, Andrew F. Olshan, PhD^{1,4}, Charles M. Perou,
4 PhD^{4,5,6}, Melissa A. Troester, PhD^{1,6}, Michael I. Love, PhD^{5,7}, Arjun Bhattacharya, PhD^{8*}

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

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

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

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

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

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

19

20 *Correspondence can be directed to AB (abtbhatt@ucla.edu)

21

22 **CONFLICT OF INTEREST STATEMENT**

23 CMP is an equity stock holder, consultant, and board of directors member of BioClassifier LLC and

24 GeneCentric Diagnostics. CMP is also listed as an inventor on patent applications on the Breast PAM50

25 assay. The other 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
3 breast cancer (BC). Studies have shown that Black women (BW) have higher CRS than White women
4 (WW). Although systemic injustices contribute substantially to BC disparities, evidence for biological and
5 germline contributions is emerging. We investigated germline genetic associations with CRS and CRS
6 disparity through a Transcriptome-Wide Association Study (TWAS). In the Carolina Breast Cancer Study,
7 using race-specific predictive models of tumor expression from germline genetics, we performed race-
8 stratified (N=1,043 WW, 1083 BW) linear regressions of three CRS (ROR-S: PAM50 subtype score;
9 Proliferation Score; ROR-P: ROR-S plus Proliferation Score) on imputed Genetically-Regulated tumor
10 eXpression (GReX). Using Bayesian multivariate regression and adaptive shrinkage, we tested TWAS-
11 significant genes for associations with PAM50 tumor expression and subtype to elucidate patterns of
12 germline regulation underlying TWAS-gene and CRS associations. At FDR-adjusted $P < 0.10$, we
13 detected 7 TWAS-genes among WW and 1 TWAS-gene among BW. Among WW, CRS showed positive
14 associations with *MCM10*, *FAM64A*, *CCNB2*, and *MMP1* GReX and negative associations with *VAV3*,
15 *PCSK6*, and *GNG11* GReX. Among BW, higher *MMP1* GReX predicted lower Proliferation score and
16 ROR-P. TWAS-gene and PAM50 tumor expression associations highlighted potential mechanisms for
17 TWAS-gene to CRS associations. Among BC patients, we find differential germline associations with
18 three CRS by race, underscoring the need for larger, more diverse datasets in molecular studies of BC.
19 Our findings also suggest possible germline *trans*-regulation of PAM50 tumor expression, with potential
20 implications for interpreting CRS in clinical settings.

21
22 **SIGNIFICANCE**

23 We find race-specific genetic associations with breast cancer risk-of-recurrence scores (CRS). Follow-up
24 analyses suggest mediation of these associations by PAM50 molecular subtype and gene expression,
25 with implications for clinical interpretation of CRS.

26
27 *Keywords:* breast cancer recurrence, risk of recurrence, transcriptome-wide association study, molecular
28 subtype, 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 (Manuscript Word Count = 3978/5000)**

2 Tumor expression-based molecular profiling has improved clinical classification of breast cancer (1-3).
3 One tool is the PAM50 assay, which integrates tumor expression of 50 genes (derived from a set of 1,900
4 subtype-specific genes identified in microarray studies) to determine PAM50 intrinsic molecular subtypes:
5 Luminal A (LumA), Luminal B (LumB), Human epidermal growth factor 2-enriched (HER2-enriched),
6 Basal-like, and Normal-like (1,4). Intrinsic molecular subtypes are strong prognostic factors for breast
7 cancer outcomes, including recurrence and mortality. For instance, Basal-like breast cancer has
8 substantially higher recurrence and mortality risk compared to LumA breast cancer (5-8). In recent years,
9 continuous risk of recurrence scores (CRS) have gained traction as a potential clinical tool that
10 encapsulates prognostic differences of breast cancer intrinsic molecular subtypes into a singular measure
11 that can be used to guide treatment decisions. CRS include ROR-S, Proliferation score, ROR-P, and
12 ROR-PT (1,9). ROR-P, for instance, is determined by combining ROR-S (PAM50 tumor expression-based
13 subtype score) and Proliferation score (tumor expression of 11 PAM50 genes). ROR-PT further integrates
14 ROR-P with information on tumor size. Studies show that CRS offer significant prognostic information
15 beyond clinical variables (e.g., nodal status, tumor grade, age, hormonal therapy), improve adjuvant
16 treatment decisions, and offer robust risk stratification for distant (5-10 years post diagnosis) recurrence
17 (10-12).

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

1 was associated among WW (18). However, the role of germline genetic variation in recurrence, CRS, and
2 CRS disparity remains a critical knowledge gap. Studying genetic associations with breast cancer
3 outcomes in BW is necessary to ensure advancements in breast cancer genetics are not limited to or
4 generalizable in only White populations, thus aiding in decreasing health disparities.

5

6 As racially-diverse genetic datasets typically have small samples of BW, gene-level association tests can
7 be used to increase study power. These approaches include TWAS, which integrates relationships
8 between single nucleotide polymorphisms (SNP) and gene expression with genome-wide association
9 studies (GWAS) to prioritize gene-trait associations (24,25). TWAS aids in interpreting genetic
10 associations by mapping significant GWAS associations to tissue-specific expression of individual genes.
11 Furthermore, TWAS can also prioritize novel tissue-specific gene-trait associations by using germline
12 genetics as a causal anchor, as, under assumptions of independent assortment of alleles, or barring that,
13 effective control for population stratification, germline genetics cannot be affected by gene expression or
14 confounding exposures. In cancer applications, TWAS has identified susceptibility genes at loci
15 previously undetected through GWAS, highlighting its improved power and interpretability (26-28).

16 Previous studies show that stratification of the entire TWAS (model training, imputation, and association
17 testing) is preferable in diverse populations, as models may perform poorly across ancestry groups and
18 methods for TWAS in admixed populations are unavailable (18,29).

19

20 Here, using data from the CBCS, which includes a large sample of Black breast cancer patients with
21 tumor gene expression data, we study race-specific germline genetic associations for CRS using TWAS.
22 CRS included in this study are ROR-S, Proliferation score, and ROR-P. Using race-specific predictive
23 models for tumor expression from germline genetics, we identify sets of TWAS-prioritized genes (TWAS-
24 genes) associated with these CRS across BW and WW. We additionally investigate TWAS-genes for
25 ROR-P for associations with PAM50 subtype and subtype-specific tumor gene expression to elucidate
26 germline contributions to PAM50 subtype, and how these mediate TWAS-gene and CRS associations.
27 Unlike previous studies that correlated tumor gene expression (as opposed to germline-regulated tumor

1 gene expression) with subtype or subtype-specific tumor gene expression, TWAS enables directional
2 interpretation of observed associations by ruling out reverse causality (24,25).

3

4 **MATERIALS AND METHODS**

5 ***Data collection***

6 *Study population*

7 The CBCS is a population-based study of North Carolina (NC) breast cancer patients, enrolled in three
8 phases; study details have been previously described (30,31). Patients aged 20 to 74 were identified
9 using rapid case ascertainment with the NC Central Cancer Registry with randomized recruitment to
10 oversample self-identified Black and young women (ages 20-49) (9,31). Demographic and clinical data
11 (age, menopausal status, body mass index, hormone receptor status, tumor stage, study phase,
12 recurrence) were obtained through questionnaires and medical records. Time-to-event recurrence data
13 were only available for CBCS Phase 3. The study was approved by the Office of Human Research Ethics
14 at the University of North Carolina at Chapel Hill, and informed consent was obtained from each
15 participant.

16

17 *CBCS genotype data*

18 Genotypes were assayed on the OncoArray Consortium's custom SNP array (Illumina Infinium
19 OncoArray) (32) and imputed using the 1000 Genomes Project (v3) as a reference panel for two-step
20 phasing and imputation using SHAPEIT2 and IMPUTEv2 (33-36). The DCEG Cancer Genomics
21 Research Laboratory conducted genotype calling, quality control, and imputation (32). We excluded
22 variants with less than 1% minor allele frequency and deviations from Hardy-Weinberg equilibrium at $P <$
23 10^{-8} (37,38). We intersected genotyping panels for BW and WW samples, resulting in 5,989,134
24 autosomal variants and 334,391 variants on the X chromosome (39). We only consider the autosomal
25 variants in this study.

26

27 *CBCS gene expression data*

1 Paraffin-embedded tumor blocks were assayed for gene expression of 406 breast cancer-related and 11
2 housekeeping genes using NanoString nCounter at the Translational Genomics Laboratory at UNC-
3 Chapel Hill (4,9). These 406 breast cancer-related genes include genes part of the PAM50, P53, E2, IGF,
4 and EGFR signatures, among others (**Supplementary Table S1**). As described previously, we eliminated
5 samples with insufficient data quality using NanoStringQCPro (18,40), scaled distributional difference
6 between lanes with upper-quartile normalization (41), and removed two dimensions of unwanted technical
7 and biological variation, estimated from housekeeping genes using RUVSeq (41,42). The current analysis
8 included 1,199 samples with both genotype and gene expression data (628 BW, 571 WW).

9

10 **Statistical analysis**

11 *Overview of TWAS*

12 TWAS integrates expression data with GWAS to prioritize gene-trait associations through a two-step
13 analysis (**Figure 1A-B**). First, using genetic and transcriptomic data, we trained predictive models of
14 tumor gene expression using all SNPs within 0.5 Megabase of the gene (18,25). Second, we used these
15 models to impute the Genetically-Regulated tumor eXpression (GReX) of a gene into an external GWAS
16 panel by multiplying the SNP-gene weights from the predictive model with the dosages of each SNP in
17 the GWAS panel. The GReX represents the portion of tumor expression explained by *cis*-genetic
18 regulation. We test for gene-trait associations with an outcome by running a linear regression using GReX
19 as the primary predictor of interest. By focusing on genetically regulated expression, TWAS generally
20 avoids instances of expression-trait association that are not consequences of genetic variation but are
21 driven by the effect of traits on expression. If sufficiently heritable genes are assayed in the correct tissue,
22 TWAS increases power to detect gene-trait associations and aids interpretability of results, as
23 associations are mapped from germline genetics to individual genes (25,43).

24

25 *CRS TWAS in CBCS*

26 We adopted techniques from FUSION to train predictive models of tumor expression from *cis*-germline
27 genotypes, as discussed previously (18,25). Motivated by strong associations between germline genetics
28 and tumor expression in CBCS (18), for genes with non-zero *cis*-heritability at nominal $P < 0.10$, we

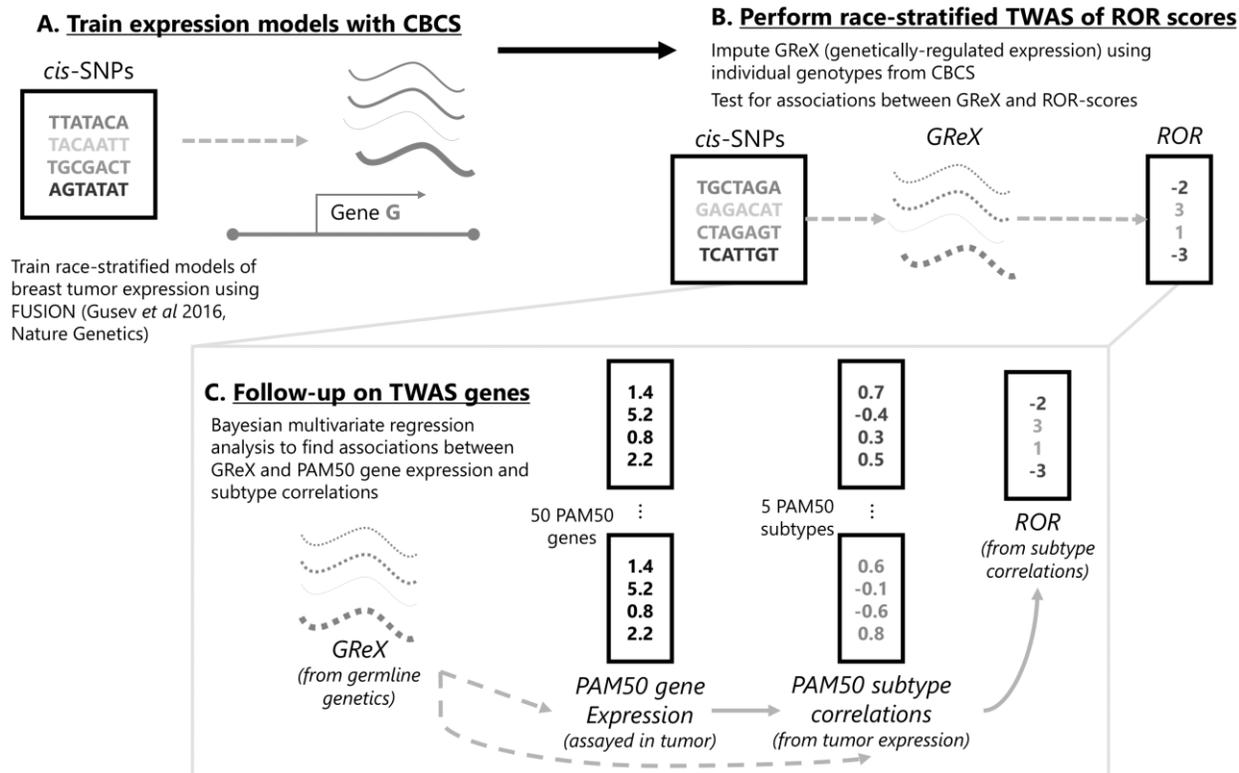


Figure 1. Schematic of study analytic approach. A) In CBCS, constructed race-stratified predictive models of tumor gene expression from *cis*-SNPs. B) In CBCS, imputed GReX at individual-level using genotypes and tested for associations between GReX and CRS in race-stratified linear models; only GReX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

- 1 trained predictive models for covariate-residualized tumor expression with all *cis*-SNPs within 0.5
- 2 Megabase using linear mixed modeling or elastic net regression (**Supplementary Methods,**
- 3 **Supplementary Materials**) (44,45). We selected models with five-fold cross-validation adjusted $R^2 >$
- 4 0.01 between predicted and observed expression values, resulting in 59 and 45 models for WW and BW,
- 5 respectively (**Supplementary Data**). These models also showed sufficiently strong predictive
- 6 performance in external validation using TCGA data (18).
- 7
- 8 Using only germline genetics as an input, we imputed GReX in 1,043 WW and 1,083 BW, respectively, in
- 9 CBCS. For samples that were not present in the training dataset, we multiplied the SNP-weights from the
- 10 predictive models with the SNP dosages to construct GReX. For samples in both the training and

1 imputation datasets, GReX was imputed via cross-validation to minimize data leakage. We tested GReX
2 for associations with ROR-S, Proliferation Score, and ROR-P using multiple linear regression adjusted for
3 age, estrogen receptor (ER) status, tumor stage, and study phase (1). We corrected for test-statistic bias
4 and inflation using a Bayesian bias and inflation adjustment method *bacon*, as TWAS are prone to bias
5 and inflation of test statistics (46). We then adjusted for multiple testing using the Benjamini-Hochberg
6 procedure (46,47). To compare germline effects with total (germline and post-transcriptional) effects on
7 ROR, we assessed relationships between total tumor expression of TWAS genes and CRS using similar
8 linear models. We were underpowered to study time-to-recurrence due to small sample size, as
9 recurrence data was collected only in CBCS Phase 3 (635 WW, 742 BW with GReX and recurrence data;
10 183 WW, 283 BW with tumor expression and recurrence data). For significant TWAS-identified genes for
11 CRS (FDR-adjusted $P < 0.10$), we conducted a follow-up permutation test: we shuffle the SNP-gene
12 weights in the predictive model 5,000 times to generate a null distribution and compare the original
13 GReX-CRS associations to this null distribution. This permutation test assessed whether the TWAS
14 association provides more tissue-specific expression context, beyond any strong SNP-CRS associations
15 at the genetic locus (25).

16

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

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

25

26 Assignment of PAM50 gene to subtype was based on PAM50 gene centroid values for each subtype; the
27 subtype assigned to a PAM50 gene corresponded to the largest positive centroid value across subtypes
28 for that gene. Subtype assignment through this "greedy algorithm" are specific to this study and represent

1 a simplified reality (e.g., ESR1 classified as part of Luminal A subtype only even though ESR1 expression
2 correlates with both Luminal A and to a slightly lesser degree Luminal B subtype). Moreover, subtype
3 assignment for this portion of analyses was conducted only for visual comparison of patterns of
4 associations between TWAS-genes and PAM50 tumor gene expressions (i.e., subtype assignment in this
5 portion of analyses had no bearing on continuous ROR score calculations or subtype-centroid
6 correlations).

7

8 *Bayesian multivariate regressions and multivariate adaptive shrinkage*

9 As previously noted (1), CRS are functions of PAM50 SCCs and gene expression profiles. Thus, we
10 followed up on CRS-associated TWAS-genes by studying their associations with PAM50 SCCs and gene
11 expression. We assessed TWAS-genes (for ROR-P) in relation to SCCs and PAM50 tumor gene
12 expression (**Figure 1C**). Importantly, consistent with the original formulation of ROR-S, we did not
13 consider normal-like subtype and normal-like subtype specific genes; subtype-specific genes were
14 determined using a greedy assignment algorithm where the highest centroid value across subtypes for a
15 given PAM50 gene determined that PAM50 gene's subtype. This classification scheme offers analytic
16 simplicity but is an oversimplification for some PAM50 genes such as ESR1, which shows high centroid
17 values for both LumA and LumB, albeit highest for LumA (ESR1 was assigned as a LumA gene in this
18 study). We found that none of our TWAS-genes were within 1 Megabase of PAM50 genes and that most
19 TWAS-genes were not on the same chromosome as PAM50 genes (**Supplementary Table S2**). Existing
20 gene-based mapping techniques for *trans*-expression quantitative trait loci (eQTL) (SNP and gene are
21 separated by more than 1 Megabase) mapping include *trans*-PrediXcan and GBAT (48,49). We employed
22 Bayesian multivariate linear regression (BtQTL) to account for correlation in multivariate outcomes (SCCs
23 and PAM50 gene expression) in association testing. BtQTL improves power to detect significant *trans*-
24 associations, especially when considering multiple genes with highly correlated (>0.5) expression
25 (**Supplementary Methods, Supplementary Figures S1-S2, Supplementary Materials**). Lastly, we
26 conducted adaptive shrinkage on BtQTL estimates using mashr, an empirical Bayes method to estimate
27 patterns of similarity and improve accuracy in associations tests across multiple outcomes (50). mashr

1 outputs revised posterior means, standard deviations, and corresponding measures of significance (local
2 false sign rates, or LFSR).

3

4 **RESULTS**

5 **Race-specific associations between GReX and CRS**

6 We performed race-specific TWAS for CRS to investigate the role of germline genetic variation in CRS
7 and CRS racial disparity. We identified 8 genes (*MCM10*, *FAM64A*, *CCNB2*, *MMP1*, *VAV3*, *PCSK6*,
8 *NDC80*, *MLPH*), 8 genes (*MCM10*, *FAM64A*, *CCNB2*, *MMP1*, *VAV3*, *NDC80*, *MLPH*, *EXO1*), and 10
9 genes (*MCM10*, *FAM64A*, *CCNB2*, *MMP1*, *VAV3*, *PCSK6*, *GNG11*, *NDC80*, *MLPH*, *EXO1*) whose GReX
10 was associated with ROR-S, proliferation, and ROR-P, respectively, in WW, and 1 gene (*MMP1*) whose
11 GReX was associated with proliferation and ROR-P in BW at FDR-adjusted $P < 0.10$ (**Figure 2A, 2B**). No
12 associations were detected between GReX and ROR-S among BW. We refer to genes with statistically
13 significant TWAS associations (FDR-adjusted $P < 0.10$) as TWAS-genes. Among these identified genes,
14 only genes that are not part of the PAM50 panel (i.e., excluding *NDC80*, *MLPH*, *EXO1*) were considered
15 in downstream permutation and TWAS-gene follow up analyses (**Figure 1C**), as we wished to focus
16 investigation on relationship between non-PAM50 TWAS-genes and PAM50 (tumor) genes.

17

18 Among WW, increased GReX of *MCM10*, *FAM64A*, *CCNB2*, and *MMP1* were associated with higher
19 CRS while increased GReX of *VAV3*, *PCSK6*, and *GNG11* were associated with lower CRS (**Figure 2A**).
20 Among BW, increased GReX of *MMP1* was associated with lower CRS (Proliferation, ROR-P, but not
21 ROR-S) (**Figure 2A**). Briefly, to contextualize the functions of these TWAS-identified genes, *MCM10* is
22 involved in DNA replication, *FAM64A* and *CCNB2* are implicated in progression and regulation of the cell
23 cycle, and *MMP1*, like the broader *MMP* family, is involved in the breakdown of the extracellular matrix
24 (51-55). *GNG11* and *VAV3* are involved in signal transduction: *GNG11* as a component of a
25 transmembrane G-protein and *VAV3* as a guanine nucleotide exchange factor for GTPases (56,57).

26

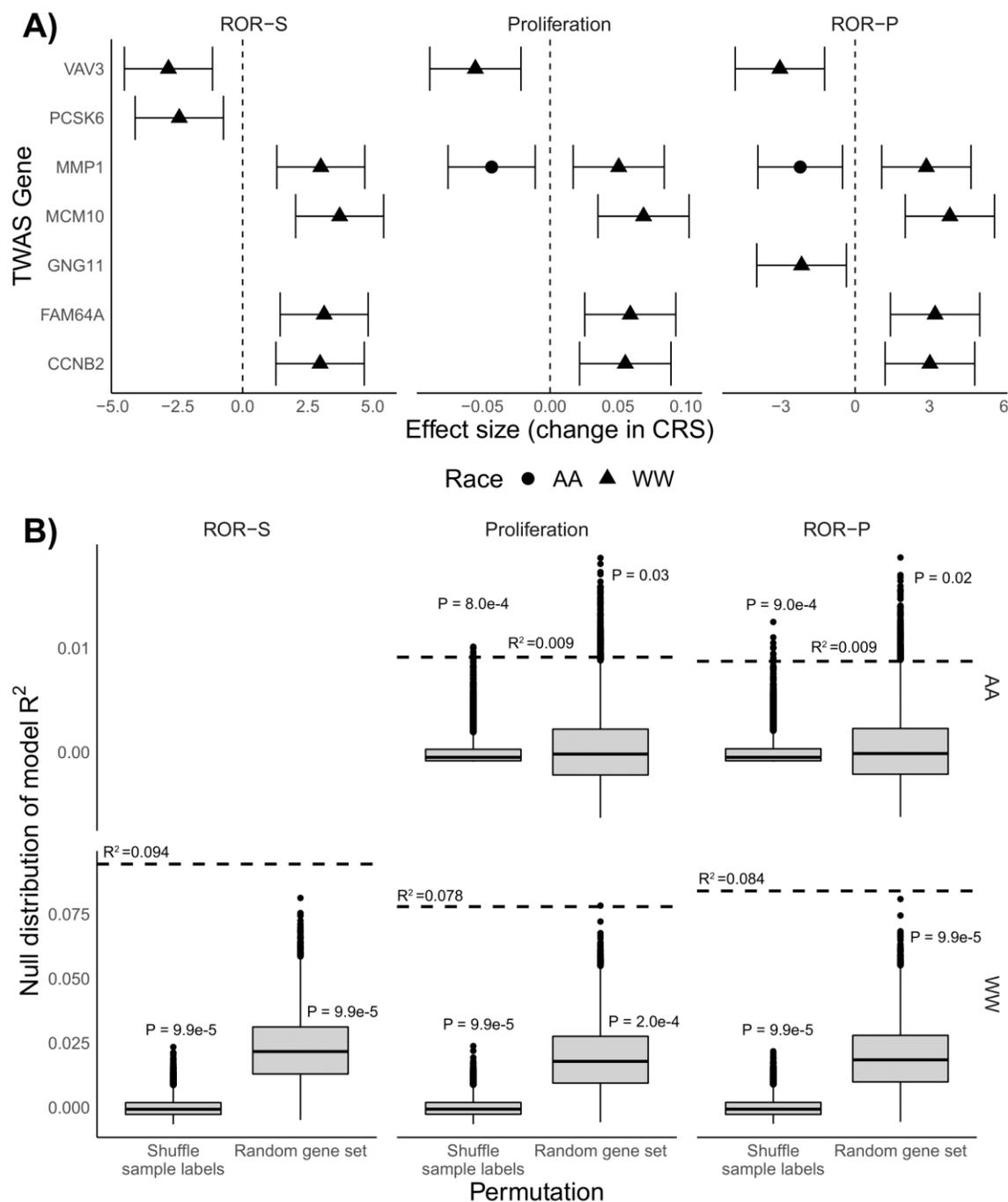


Figure 2. Permutation tests and associations between TWAS-genes and CRS for WW and BW. A) Effect estimates correspond to change in ROR-S, Proliferation score, and ROR-P per one standard deviation increase in TWAS-gene expression (i.e., one standard deviation increase in GReX of gene). Triangle denotes 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- R^2 for regressions of CRS on TWAS-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- R^2 .

- 1 Associations between tumor expression of TWAS-genes and CRS were concordant, in terms of direction
- 2 of association to germline-only effects among WW; findings were discordant among BW where higher

1 tumor expression of MMP1 was associated with higher CRS (**Table 1, Supplementary Table S3**).

2

3 **Permutation testing provides context to TWAS-gene and CRS associations**

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

14

15 Next, we quantified the percent variance explained of CRS by the TWAS genes, in aggregate, by
16 calculating the model adjusted- R^2 for a regression of covariate-residualized CRS on GReX all TWAS-
17 genes. To context these model adjusted- R^2 , we conducted two permutation tests. First, we permuted the
18 sample labels for covariate-residualized CRS 10,000 times and computed the model adjusted R^2 at each
19 iteration to generate a null distribution for adjusted R^2 between TWAS-genes and CRS. Across WW and
20 BW, the observed R^2 of TWAS-genes against CRS (7-10% among WW and 1% among BW) were
21 statistically significant against the respective null distributions (P-values and distributions in **Figure 2B**).

22 Permutation tests for analyses of tumor expression of TWAS-genes and CRS are available in
23 **Supplementary Figure S3**. Second, we wanted to assess if the GReX of these sets of TWAS-genes (7 in
24 WW and 1 in AA) explained more of the variance in CRS than the GReX of a randomly selected set of
25 genes of the same size. Previous studies have shown that the tumor expression of a set randomly
26 selected genes is likely to be predictive of breast cancer outcomes; we wished to investigate this
27 phenomenon on the GReX level (58,59). Over 10,000 repetitions, we randomly selected 7 and 1 genes in
28 WW and AA subjects, respectively, ran a multivariable regression, and calculated the model adjusted- R^2

1 to generate another null distribution. Here again, we found that the true model R^2 outperformed the null
2 distribution, all showing permutation $P < 0.05$ in these settings (**Figure 2B**). These permutation tests
3 show that our TWAS-prioritized genes, taken together, appreciably explain differences in CRS.

4

5 **Associations between TWAS-genes and PAM50 subtype correlations and gene expression**

6 As CRS are constructed from PAM50 subtype-specific correlations and gene expression profiles, we
7 further studied associations between GReX of TWAS-genes and PAM50 SCCs and gene expression to
8 understand how PAM50 subtype and gene expression mediate TWAS-gene and CRS associations.

9 Among WW, a one standard deviation increase in *FAM64A* and *CCNB2* GReX resulted in significantly
10 increased Basal-like SCC while an identical increase in *VAV3*, *PCSK6*, and *GNG11* GReX resulted in
11 significantly increased Luminal A SCC. The magnitude of increase in correlation for respective subtypes
12 per TWAS-gene was approximately 0.05, and most estimates had credible intervals that did not intersect
13 the null. Among WW, associations between HER2-like SCC and TWAS-genes followed similar patterns to
14 associations for the Basal-like subtype, although associations for HER2 were more precise (**Figure 3A**).

15 We found predominantly null associations between TWAS-genes and Luminal B SCC among WW
16 (**Figure 3A**). Unlike in WW, for BW, an increase in *MMP1* GReX was not associated with Luminal A,
17 HER2 or Basal-like SCCs. Instead, among BW, *MMP1* GReX was significantly negatively associated with
18 Luminal B SCC. Estimates from univariate regressions are provided in **Supplementary Tables S4-S7**.

19

20 For both WW and BW, the pattern of associations between TWAS-genes and PAM50 tumor expression
21 were predominantly congruent with observed associations between TWAS-genes and PAM50 SCCs as
22 well as TWAS-genes and CRS (**Figure 4, Supplementary Tables S8-S11**). In WW, a one standard
23 deviation increase in *CCNB2* GReX was associated with significantly increased *ORC6L*, *PTTG1*, and
24 *KIF2C* (Basal-like genes) expression and *UBE2T* and *MYBL2* (LumB genes) expression. By contrast, a
25 one standard deviation increase in *PCSK6* GReX significantly increased *BAG1*, *FOXA1*, *MAPT*, and
26 *NAT1* (LumA genes) expression (**Figure 3B**). While increased *MMP1* GReX was associated with
27 significantly increased expression of *ORC6L* (basal-like gene), *MYBL2*, and *BIRC5* (LumB genes) among
28 WW, this was not the case among BW. Instead, increased *MMP1* GReX among BW was significantly

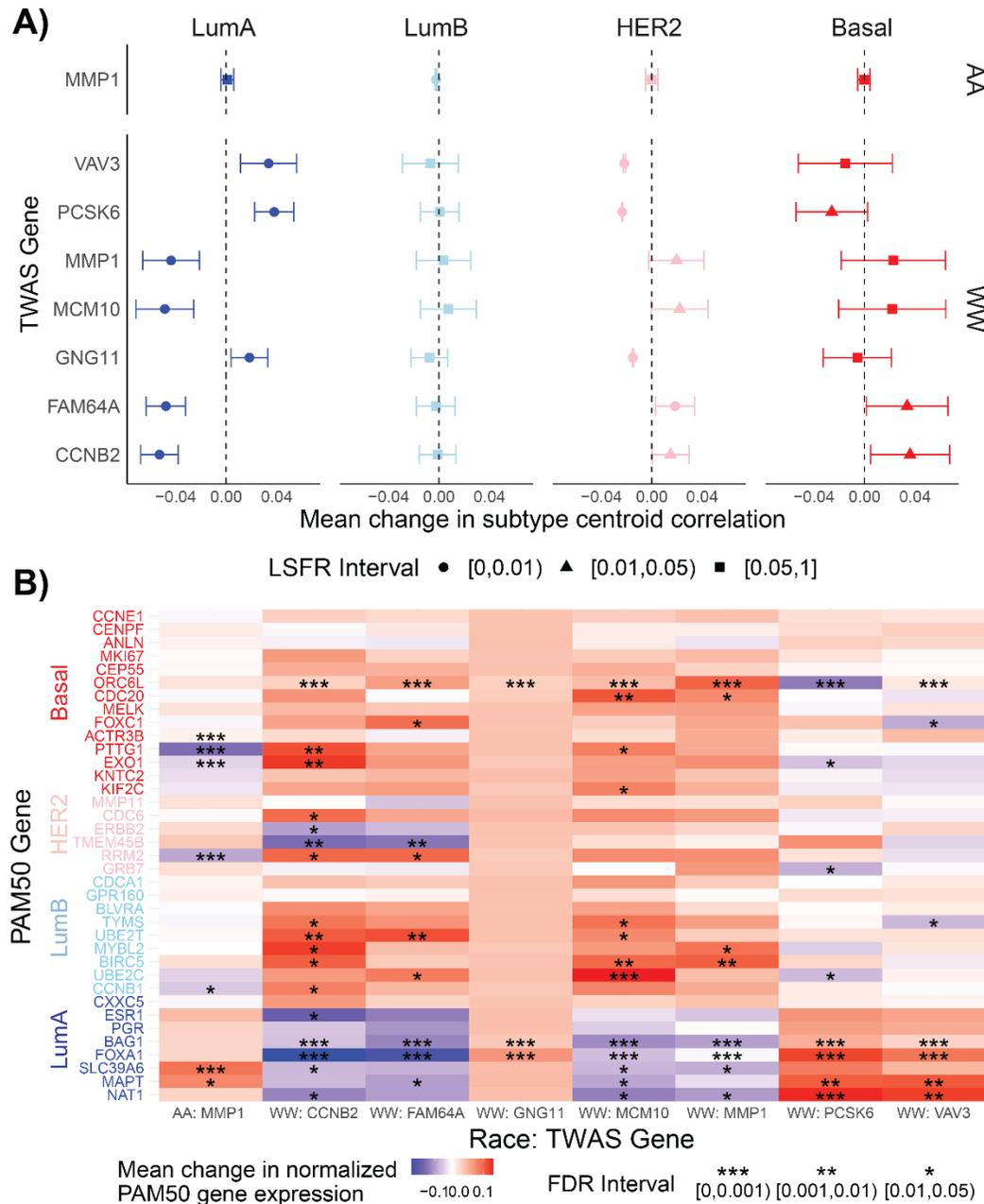


Figure 3. Associations between TWAS-genes and PAM50 SCCs and gene expression. A) Among AA (top) and WW (bottom), associations between TWAS-genes and PAM50 SCCs using Bayesian multivariate regression and multivariate adaptive shrinkage. Effect estimates show change in SCCs (range -1 to 1) for one standard deviation increase in TWAS-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 TWAS-GReX. *, **, *** denote FDR intervals for effect sizes.

- 1 associated with increased expression of *SLC39A6* (LumA gene) and decreased expression of *ACTR3B*,
- 2 *PTTG1*, and *EXO1* (Basal-like genes) (**Figure 3B**). Associations between TWAS-genes and PAM50
- 3 genes provide a granular, gene interaction level view into the mediation of the TWAS-gene and CRS

1 association, suggesting that *trans*-regulation of subtype-specific PAM50 genes by TWAS-genes in breast
2 tumors could be a possible contributor to subtype formation/maintenance and, subsequently, CRS and
3 recurrence.

4

5 **DISCUSSION**

6 Through TWAS, we identified 7 and 1 genes among WW and BW, respectively, for which genetically-
7 regulated breast tumor expression was associated with CRS and underlying PAM50 gene expression and
8 subtype. Among WW, these 7 TWAS-genes explained between 7-10% of the variation in CRS, a large
9 and statistically significant proportion of variance. Among BW, the singular TWAS-gene explained a
10 statistically significant ~1% of the variation in Proliferation score and ROR-P. There are two key novel
11 aspects to this study. First, existing literature on associations between tumor gene expression and
12 recurrence (for which CRS are a proxy) cannot distinguish between genetic and non-genetic components
13 of effect (60), whereas, here, we estimate the contribution of the genetic component. Second, TWAS
14 allows directional interpretation of observed associations that are not possible when correlating tumor
15 gene expression and recurrence. For instance, prior studies report *CCNB2* is upregulated in triple-
16 negative breast cancers (TNBC) but were unable to determine whether increased *CCNB2* expression
17 contributes to development or maintenance of TNBC or is part of the molecular response to cancer
18 progression (61,62). By contrast, GReX is a function of only genetic variation. As such, we can
19 confidently rule out that differences in *CCNB2* GReX are not direct consequences of subtype (and by
20 extension recurrence); however, our observed associations of *CCNB2* GReX and subtype suggest a
21 potential directional relationship for further study. Thus, TWAS allows a directional, causal interpretation,
22 subject to effective control for population stratification, minimal horizontal pleiotropy, and assumptions of
23 independent assortment of alleles (24,25).

24

25 Our TWAS-gene and subtype associations among WW are consistent with literature on the association
26 between tumor (i.e., genetic and non-genetic) expression of our TWAS-genes and subtype. Prior
27 investigations in cohorts of primarily European ancestry have reported that *MCM10*, *FAM64A*, and
28 *CCNB2* expression is higher in ER-negative compared to ER-positive tumors (61-63). In studies that

1 compared triple-negative and non-triple negative subtypes, higher *MCM10*, *FAM64A*, and *CCNB2*
2 expression was detected in triple-negative breast cancer (61,62). Histologically, HER2-enriched and
3 Basal-like subtypes are typically ER-negative, and triple-negatives are similar to Basal-like subtypes
4 (9,64). Moreover, our findings among WW that GReX of *PCSK6* and *VAV3* associated with LumA
5 subtype and LumA-specific gene expression are also consistent with previous results of *PCSK6* and
6 *VAV3* upregulation in ER-positive subtypes (65,66). Importantly, our associations suggest directional
7 mechanisms: from germline variation, to GReX of TWAS-gene, and ultimately, to subtype.

8
9 Presently, little is known about germline genetic regulation of PAM50 tumor expression. In CBCS, we
10 found that tumor expression of most PAM50 genes is not *cis*-heritable (18). Instead, observed TWAS-
11 gene and PAM50 gene expression associations may implicate *trans*-gene regulation of the PAM50
12 signature. For instance, we found that *VAV3* GReX is significantly positively associated with tumor
13 expression of *BAG1*, *FOXA1*, *MAPT*, and *NAT1* and nominally with increased tumor *ESR1* expression, all
14 of which are Luminal A-specific genes. Such *trans*-genetic regulation signals, especially in the case of
15 *ESR1*, pose significant clinical and therapeutic implication if confirmed under experimental conditions. For
16 example, *VAV3* activates *RAC1* which upregulates *ESR1* but such mechanistic evidence is sparse for
17 other putative TWAS-gene to PAM50 gene associations (67,68). More generally, two of the TWAS-genes
18 among WW have been found to activate transcription factors; *FAM64A* enhances oncogenic nuclear
19 factor-kappa B (NF- κ B) signaling while both *FAM64A* and *PCSK6* activate oncogenic *STAT3* signaling
20 (69-71).

21
22 Differences in the number and effect of identified TWAS-genes by race point to factors that warrant
23 further investigation: (1) potential greater contribution of *trans*-regulation in tumor gene expression in BW
24 (methods for capturing *trans*-regulation in gene expression predictive models are not as well-developed
25 as those for *cis*-regulation) (18) and (2) potential racial differences in tumor methylation and somatic
26 alternations, which were not assayed in CBCS (18,72-77). These factors should be investigated further as
27 transcriptomic and epigenomic datasets for racially-diverse cohorts of breast cancer patients become
28 available. Interestingly, we found *MMP1* GReX has divergent associations with ROR across race. There

1 are a few potential explanations. First, the range of *MMP1* GReX was manifold among WW than BW,
2 suggesting sparser *cis*-eQTL architecture of *MMP1* in BW and more influence from *trans*-acting signals.
3 Potential differences in influence of germline genetics on tumor expression and ROR by race could be an
4 artifact of divergent somatic or epigenetic factors that CBCS has not assayed (74-77). Second, while
5 studies generally report that *MMP1* tumor expression is higher in triple-negative and Basal-like breast
6 cancer, one study reported that *MMP1* expression in tumor cells does not significantly differ by subtype
7 (78-80). Instead, Bostrom *et al.* reported that *MMP1* expression differs in stromal cells of patients with
8 different subtypes (80). There is evidence to suggest that tumor composition, including stromal and
9 immune components, may influence breast cancer progression in a subtype-specific manner. Future
10 studies should consider expression predictive models that integrate greater detail on tumor cell-type
11 composition to disentangle potential race-specific tumor composition effects on race-specific TWAS
12 associations (81,82).

13
14 There are a few limitations to this study. First, as CBCS used a Nanostring nCounter probeset for mRNA
15 expression quantification of genes relevant for breast cancer, we could not analyze the whole human
16 transcriptome. While this probeset may exclude several *cis*-heritable genes, CBCS contains one of the
17 largest breast tumor transcriptomic datasets for Black women, allowing us to build well-powered race-
18 specific predictive models, a pivotal step in ancestry-specific TWAS. Second, CBCS lacked data on
19 somatic amplifications and deletions, inclusion of which could enhance the performance of predictive
20 models of tumor expression (83). Third, as recurrence data was collected in a small subset with few
21 recurrence events, we were unable to make a direct comparison between CRS and recurrence results,
22 which may affect clinical generalizability. However, to our knowledge, CBCS is the largest resource of
23 PAM50-based CRS data.

24
25 Our analysis provides evidence of race-specific putative germline associations to CRS, mediated through
26 associations between genetically-regulated tumor expression of TWAS-genes and PAM50 expressions
27 and subtype. This work underscores the need for larger and more diverse cohorts for genetic
28 epidemiology studies of breast cancer. Future studies should consider subtype-specific genetics (i.e.,

1 stratification by subtype in predictive model training and association analyses) to elucidate heritable gene
2 expression effects on breast cancer outcomes both across and within subtype, which may yield further
3 hypotheses for more fine-tuned clinical intervention.

4

5 **ACKNOWLEDGEMENTS**

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

12

13 **FUNDING**

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

24

25 Funding for BCAC and iCOGS came from: Cancer Research UK [grant numbers C1287/A16563,
26 C1287/A10118, C1287/A10710, C12292/A11174, C1281/A12014, C5047/A8384, C5047/A15007,
27 C5047/A10692, 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-
3 223175] (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
5 of Defence [W81XWH-10-1-0341], and the Canadian Institutes of Health Research CIHR) for the CIHR
6 Team in Familial Risks of Breast Cancer [grant PSR-SIIRI-701]. All studies and funders as listed in
7 Michailidou K *et al* (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,
11 MIL, AB. Funding acquisition: AP, MG, AFO, CMP, MAT, MIL. Methodology: AP, MIL, AB. Project
12 administration: MAT, MIL, AB. Resources: MG, AFO, CMP, MAT, MIL. Supervision: MAT, MIL, AB.
13 Visualization: AP, AB. Writing – original draft: AP, AB. Writing – reviewing and editing: AP, MG, AFO,
14 CMP, MAT, MIL, AB.

15

16 **AVAILABILITY OF DATA AND MATERIALS**

17 Expression data from CBCS is available on NCBI GEO with accession number GSE148426. CBCS
18 genotype datasets analyzed in this study are not publicly available as many CBCS patients are still being
19 followed and accordingly CBCS data is considered sensitive; the data is available from M.A.T upon
20 reasonable request. Supplementary Data includes summary statistics for eQTL results, tumor expression
21 models, and relevant R code for training expression models in CBCS and are freely available
22 at https://github.com/bhattacharya-a-bt/CBCS_TWAS_Paper/. R code for analyses provided in this paper
23 are available at <https://github.com/APUNC/CBCS---Risk-of-Recurrence-Paper>. iCOGs summary statistics
24 are available online at <http://bcac.ccge.medschl.cam.ac.uk/bcacdata/icogs-complete-summary-results>.

25

26 **REFERENCES**

27 1. Parker JS, Mullins M, Cheang MC, Leung S, Voduc D, Vickery T, *et al*. Supervised risk predictor of
28 breast cancer based on intrinsic subtypes. *J Clin Oncol* **2009**;27:1160-7

- 1 2. Wallden B, Storhoff J, Nielsen T, Dowidar N, Schaper C, Ferree S, *et al.* Development and
2 verification of the PAM50-based Prosigna breast cancer gene signature assay. *BMC Med Genomics*
3 **2015**;8:54
- 4 3. Paik S, Shak S, Tang G, Kim C, Baker J, Cronin M, *et al.* A multigene assay to predict recurrence of
5 tamoxifen-treated, node-negative breast cancer. *N Engl J Med* **2004**;351:2817-26
- 6 4. Geiss GK, Bumgarner RE, Birditt B, Dahl T, Dowidar N, Dunaway DL, *et al.* Direct multiplexed
7 measurement of gene expression with color-coded probe pairs. *Nat Biotechnol* **2008**;26:317-25
- 8 5. Carey LA, Perou CM, Livasy CA, Dressler LG, Cowan D, Conway K, *et al.* Race, breast cancer
9 subtypes, and survival in the Carolina Breast Cancer Study. *Jama* **2006**;295:2492-502
- 10 6. O'Brien KM, Cole SR, Tse CK, Perou CM, Carey LA, Foulkes WD, *et al.* Intrinsic breast tumor
11 subtypes, race, and long-term survival in the Carolina Breast Cancer Study. *Clin Cancer Res*
12 **2010**;16:6100-10
- 13 7. Shim HJ, Kim SH, Kang BJ, Choi BG, Kim HS, Cha ES, *et al.* Breast cancer recurrence according to
14 molecular subtype. *Asian Pac J Cancer Prev* **2014**;15:5539-44
- 15 8. van Maaren MC, de Munck L, Strobbe LJA, Sonke GS, Westenend PJ, Smidt ML, *et al.* Ten-year
16 recurrence rates for breast cancer subtypes in the Netherlands: A large population-based study. *Int J*
17 *Cancer* **2019**;144:263-72
- 18 9. Troester MA, Sun X, Allott EH, Geradts J, Cohen SM, Tse CK, *et al.* Racial Differences in PAM50
19 Subtypes in the Carolina Breast Cancer Study. *J Natl Cancer Inst* **2018**;110:176-82
- 20 10. Dowsett M, Sestak I, Lopez-Knowles E, Sidhu K, Dunbier AK, Cowens JW, *et al.* Comparison of
21 PAM50 Risk of Recurrence Score With Oncotype DX and IHC4 for Predicting Risk of Distant
22 Recurrence After Endocrine Therapy. *Journal of Clinical Oncology* **2013**;31:2783-90
- 23 11. Sestak I, Buus R, Cuzick J, Dubsy P, Kronenwett R, Denkert C, *et al.* Comparison of the
24 Performance of 6 Prognostic Signatures for Estrogen Receptor-Positive Breast Cancer: A Secondary
25 Analysis of a Randomized Clinical Trial. *JAMA Oncol* **2018**;4:545-53
- 26 12. Ohnstad HO, Borgen E, Falk RS, Lien TG, Aaserud M, Sveli MAT, *et al.* Prognostic value of PAM50
27 and risk of recurrence score in patients with early-stage breast cancer with long-term follow-up.
28 *Breast Cancer Res* **2017**;19:120

- 1 13. Albain KS, Gray RJ, Makower DF, Faghiih A, Hayes DF, Geyer CE, *et al.* Race, ethnicity and clinical
2 outcomes in hormone receptor-positive, HER2-negative, node-negative breast cancer in the
3 randomized TAILORx trial. *J Natl Cancer Inst* **2020**
- 4 14. Reeder-Hayes KE, Anderson BO. Breast Cancer Disparities at Home and Abroad: A Review of the
5 Challenges and Opportunities for System-Level Change. *Clin Cancer Res* **2017**;23:2655-64
- 6 15. Durham DD, Robinson WR, Lee SS, Wheeler SB, Reeder-Hayes KE, Bowling JM, *et al.* Insurance-
7 Based Differences in Time to Diagnostic Follow-up after Positive Screening Mammography. *Cancer*
8 *Epidemiol Biomarkers Prev* **2016**;25:1474-82
- 9 16. Wheeler SB, Reeder-Hayes KE, Carey LA. Disparities in breast cancer treatment and outcomes:
10 biological, social, and health system determinants and opportunities for research. *Oncologist*
11 **2013**;18:986-93
- 12 17. Ko NY, Hong S, Winn RA, Calip GS. Association of Insurance Status and Racial Disparities With the
13 Detection of Early-Stage Breast Cancer. *JAMA Oncology* **2020**;6:385-92
- 14 18. Bhattacharya A, García-Closas M, Olshan AF, Perou CM, Troester MA, Love MI. A framework for
15 transcriptome-wide association studies in breast cancer in diverse study populations. *Genome Biol*
16 **2020**;21:42
- 17 19. Escala-Garcia M, Guo Q, Dörk T, Canisius S, Keeman R, Dennis J, *et al.* Genome-wide association
18 study of germline variants and breast cancer-specific mortality. *Br J Cancer* **2019**;120:647-57
- 19 20. Muranen TA, Khan S, Fagerholm R, Aittomäki K, Cunningham JM, Dennis J, *et al.* Association of
20 germline variation with the survival of women with BRCA1/2 pathogenic variants and breast cancer.
21 *NPJ Breast Cancer* **2020**;6:44
- 22 21. Huo D, Hu H, Rhie SK, Gamazon ER, Cherniack AD, Liu J, *et al.* Comparison of Breast Cancer
23 Molecular Features and Survival by African and European Ancestry in The Cancer Genome Atlas.
24 *JAMA Oncol* **2017**;3:1654-62
- 25 22. Pitt JJ, Riester M, Zheng Y, Yoshimatsu TF, Sanni A, Oluwasola O, *et al.* Characterization of Nigerian
26 breast cancer reveals prevalent homologous recombination deficiency and aggressive molecular
27 features. *Nature Communications* **2018**;9:4181

- 1 23. Figueroa JD, Davis Lynn BC, Edusei L, Titiloye N, Adjei E, Clegg-Lampitey J-N, *et al.* Reproductive
2 factors and risk of breast cancer by tumor subtypes among Ghanaian women: A population-based
3 case–control study. *International Journal of Cancer* **2020**;147:1535-47
- 4 24. Gamazon ER, Wheeler HE, Shah KP, Mozaffari SV, Aquino-Michaels K, Carroll RJ, *et al.* A gene-
5 based association method for mapping traits using reference transcriptome data. *Nat Genet*
6 **2015**;47:1091-8
- 7 25. Gusev A, Ko A, Shi H, Bhatia G, Chung W, Penninx BW, *et al.* Integrative approaches for large-scale
8 transcriptome-wide association studies. *Nat Genet* **2016**;48:245-52
- 9 26. Zhong J, Jermusyk A, Wu L, Hoskins JW, Collins I, Mocci E, *et al.* A Transcriptome-Wide Association
10 Study Identifies Novel Candidate Susceptibility Genes for Pancreatic Cancer. *J Natl Cancer Inst*
11 **2020**;112:1003-12
- 12 27. Wu L, Shi W, Long J, Guo X, Michailidou K, Beesley J, *et al.* A transcriptome-wide association study
13 of 229,000 women identifies new candidate susceptibility genes for breast cancer. *Nat Genet*
14 **2018**;50:968-78
- 15 28. Mancuso N, Gayther S, Gusev A, Zheng W, Penney KL, Kote-Jarai Z, *et al.* Large-scale
16 transcriptome-wide association study identifies new prostate cancer risk regions. *Nat Commun*
17 **2018**;9:4079
- 18 29. Keys KL, Mak ACY, White MJ, Eckalbar WL, Dahl AW, Mefford J, *et al.* On the cross-population
19 generalizability of gene expression prediction models. *PLoS Genet* **2020**;16:e1008927
- 20 30. Hair BY, Hayes S, Tse CK, Bell MB, Olshan AF. Racial differences in physical activity among breast
21 cancer survivors: implications for breast cancer care. *Cancer* **2014**;120:2174-82
- 22 31. Newman B, Moorman PG, Millikan R, Qaqish BF, Geradts J, Aldrich TE, *et al.* The Carolina Breast
23 Cancer Study: integrating population-based epidemiology and molecular biology. *Breast Cancer Res*
24 *Treat* **1995**;35:51-60
- 25 32. Amos CI, Dennis J, Wang Z, Byun J, Schumacher FR, Gayther SA, *et al.* The OncoArray Consortium:
26 A Network for Understanding the Genetic Architecture of Common Cancers. *Cancer Epidemiol*
27 *Biomarkers Prev* **2017**;26:126-35

- 1 33. Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, *et al.* A global reference for
2 human genetic variation. *Nature* **2015**;526:68-74
- 3 34. O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, *et al.* A general approach for
4 haplotype phasing across the full spectrum of relatedness. *PLoS Genet* **2014**;10:e1004234
- 5 35. Delaneau O, Marchini J, Zagury JF. A linear complexity phasing method for thousands of genomes.
6 *Nat Methods* **2011**;9:179-81
- 7 36. Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method for the next
8 generation of genome-wide association studies. *PLoS Genet* **2009**;5:e1000529
- 9 37. Wigginton JE, Cutler DJ, Abecasis GR. A note on exact tests of Hardy-Weinberg equilibrium. *Am J*
10 *Hum Genet* **2005**;76:887-93
- 11 38. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, *et al.* PLINK: a tool set for
12 whole-genome association and population-based linkage analyses. *Am J Hum Genet* **2007**;81:559-75
- 13 39. Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, *et al.* dbSNP: the NCBI database
14 of genetic variation. *Nucleic Acids Res* **2001**;29:308-11
- 15 40. Bhattacharya A, Hamilton AM, Furberg H, Pietzak E, Purdue MP, Troester MA, *et al.* An approach for
16 normalization and quality control for NanoString RNA expression data. *Brief Bioinform* **2020**
- 17 41. Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biol*
18 **2010**;11:R106
- 19 42. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data
20 with DESeq2. *Genome Biol* **2014**;15:550
- 21 43. Ding B, Cao C, Li Q, Wu J, Long Q. Power analysis of transcriptome-wide association study. *bioRxiv*
22 **2020**:2020.07.19.211151
- 23 44. Endelman JB. Ridge Regression and Other Kernels for Genomic Selection with R Package rrBLUP.
24 *The Plant Genome* **2011**;4
- 25 45. Friedman J, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear Models via
26 Coordinate Descent. *J Stat Softw* **2010**;33:1-22
- 27 46. van Iterson M, van Zwet EW, Heijmans BT. Controlling bias and inflation in epigenome- and
28 transcriptome-wide association studies using the empirical null distribution. *Genome Biol* **2017**;18:19

- 1 47. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach
2 to Multiple Testing. *Journal of the Royal Statistical Society Series B (Methodological)* **1995**;57:289-
3 300
- 4 48. Wheeler HE, Ploch S, Barbeira AN, Bonazzola R, Andaleon A, Fotuhi Siahipirani A, *et al.* Imputed
5 gene associations identify replicable trans-acting genes enriched in transcription pathways and
6 complex traits. *Genetic Epidemiology* **2019**;43:596-608
- 7 49. Liu X, Mefford JA, Dahl A, He Y, Subramaniam M, Battle A, *et al.* GBAT: a gene-based association
8 test for robust detection of trans-gene regulation. *Genome Biology* **2020**;21:211
- 9 50. Urbut SM, Wang G, Carbonetto P, Stephens M. Flexible statistical methods for estimating and testing
10 effects in genomic studies with multiple conditions. *Nat Genet* **2019**;51:187-95
- 11 51. Watase G, Takisawa H, Kanemaki MT. Mcm10 plays a role in functioning of the eukaryotic replicative
12 DNA helicase, Cdc45-Mcm-GINS. *Curr Biol* **2012**;22:343-9
- 13 52. Zhao W-m, Coppinger JA, Seki A, Cheng X-l, Yates JR, Fang G. RCS1, a substrate of APC/C,
14 controls the metaphase to anaphase transition. *Proceedings of the National Academy of Sciences*
15 **2008**;105:13415-20
- 16 53. Daldello EM, Luong XG, Yang C-R, Kuhn J, Conti M. Cyclin B2 is required for progression through
17 meiosis in mouse oocytes. *Development* **2019**;146:dev172734
- 18 54. Draetta G, Luca F, Westendorf J, Brizuela L, Ruderman J, Beach D. Cdc2 protein kinase is
19 complexed with both cyclin A and B: evidence for proteolytic inactivation of MPF. *Cell* **1989**;56:829-38
- 20 55. Page-McCaw A, Ewald AJ, Werb Z. Matrix metalloproteinases and the regulation of tissue
21 remodelling. *Nat Rev Mol Cell Biol* **2007**;8:221-33
- 22 56. Rao S, Lyons LS, Fahrenholtz CD, Wu F, Farooq A, Balkan W, *et al.* A novel nuclear role for the
23 Vav3 nucleotide exchange factor in androgen receptor coactivation in prostate cancer. *Oncogene*
24 **2012**;31:716-27
- 25 57. Hossain MN, Sakemura R, Fujii M, Ayusawa D. G-protein gamma subunit GNG11 strongly regulates
26 cellular senescence. *Biochem Biophys Res Commun* **2006**;351:645-50
- 27 58. Venet D, Dumont JE, Detours V. Most random gene expression signatures are significantly
28 associated with breast cancer outcome. *PLoS Comput Biol* **2011**;7:e1002240

- 1 59. Shimoni Y. Association between expression of random gene sets and survival is evident in multiple
2 cancer types and may be explained by sub-classification. *PLoS Comput Biol* **2018**;14:e1006026
- 3 60. Parada H, Jr., Sun X, Fleming JM, Williams-DeVane CR, Kirk EL, Olsson LT, *et al.* Race-associated
4 biological differences among luminal A and basal-like breast cancers in the Carolina Breast Cancer
5 Study. *Breast Cancer Res* **2017**;19:131
- 6 61. Prat A, Adamo B, Cheang MC, Anders CK, Carey LA, Perou CM. Molecular characterization of basal-
7 like and non-basal-like triple-negative breast cancer. *Oncologist* **2013**;18:123-33
- 8 62. Zhang C, Han Y, Huang H, Min L, Qu L, Shou C. Integrated analysis of expression profiling data
9 identifies three genes in correlation with poor prognosis of triple-negative breast cancer. *Int J Oncol*
10 **2014**;44:2025-33
- 11 63. Mahadevappa R, Neves H, Yuen SM, Jameel M, Bai Y, Yuen HF, *et al.* DNA Replication Licensing
12 Protein MCM10 Promotes Tumor Progression and Is a Novel Prognostic Biomarker and Potential
13 Therapeutic Target in Breast Cancer. *Cancers (Basel)* **2018**;10
- 14 64. Hagemann IS. Molecular Testing in Breast Cancer: A Guide to Current Practices. *Arch Pathol Lab*
15 *Med* **2016**;140:815-24
- 16 65. Thakkar AD, Raj H, Chakrabarti D, Ravishankar, Saravanan N, Muthuvelan B, *et al.* Identification of
17 gene expression signature in estrogen receptor positive breast carcinoma. *Biomark Cancer* **2010**;2:1-
18 15
- 19 66. Aguilar H, Urruticoechea A, Halonen P, Kiyotani K, Mushiroda T, Barril X, *et al.* VAV3 mediates
20 resistance to breast cancer endocrine therapy. *Breast Cancer Res* **2014**;16:R53
- 21 67. Zeng L, Sachdev P, Yan L, Chan JL, Trenkle T, McClelland M, *et al.* Vav3 mediates receptor protein
22 tyrosine kinase signaling, regulates GTPase activity, modulates cell morphology, and induces cell
23 transformation. *Mol Cell Biol* **2000**;20:9212-24
- 24 68. Rosenblatt AE, Garcia MI, Lyons L, Xie Y, Maiorino C, Désiré L, *et al.* Inhibition of the Rho GTPase,
25 Rac1, decreases estrogen receptor levels and is a novel therapeutic strategy in breast cancer.
26 *Endocr Relat Cancer* **2011**;18:207-19

- 1 69. Xu Z-S, Zhang H-X, Li W-W, Ran Y, Liu T-T, Xiong M-G, *et al.* FAM64A positively regulates STAT3
2 activity to promote Th17 differentiation and colitis-associated carcinogenesis. Proceedings of the
3 National Academy of Sciences **2019**;116:10447-52
- 4 70. Jiang H, Wang L, Wang F, Pan J. Proprotein convertase subtilisin/kexin type 6 promotes in vitro
5 proliferation, migration and inflammatory cytokine secretion of synovial fibroblast-like cells from
6 rheumatoid arthritis via nuclear- κ B, signal transducer and activator of transcription 3 and extracellular
7 signal regulated 1/2 pathways. Mol Med Rep **2017**;16:8477-84
- 8 71. Jiang L, Ren L, Zhang X, Chen H, Chen X, Lin C, *et al.* Overexpression of PIMREG promotes breast
9 cancer aggressiveness via constitutive activation of NF- κ B signaling. EBioMedicine **2019**;43:188-200
- 10 72. Gravel S. Population genetics models of local ancestry. Genetics **2012**;191:607-19
- 11 73. Nelson D, Kelleher J, Ragsdale AP, Moreau C, McVean G, Gravel S. Accounting for long-range
12 correlations in genome-wide simulations of large cohorts. PLoS Genet **2020**;16:e1008619
- 13 74. Shang L, Smith JA, Zhao W, Kho M, Turner ST, Mosley TH, *et al.* Genetic Architecture of Gene
14 Expression in European and African Americans: An eQTL Mapping Study in GENOA. Am J Hum
15 Genet **2020**;106:496-512
- 16 75. Wang S, Dorsey TH, Terunuma A, Kittles RA, Ambts S, Kwabi-Addo B. Relationship between tumor
17 DNA methylation status and patient characteristics in African-American and European-American
18 women with breast cancer. PLoS One **2012**;7:e37928
- 19 76. Conway K, Edmiston SN, Tse CK, Bryant C, Kuan PF, Hair BY, *et al.* Racial variation in breast tumor
20 promoter methylation in the Carolina Breast Cancer Study. Cancer Epidemiol Biomarkers Prev
21 **2015**;24:921-30
- 22 77. Chen Y, Sadasivan SM, She R, Datta I, Taneja K, Chitale D, *et al.* Breast and prostate cancers
23 harbor common somatic copy number alterations that consistently differ by race and are associated
24 with survival. BMC Med Genomics **2020**;13:116
- 25 78. Wang QM, Lv L, Tang Y, Zhang L, Wang LF. MMP-1 is overexpressed in triple-negative breast
26 cancer tissues and the knockdown of MMP-1 expression inhibits tumor cell malignant behaviors in
27 vitro. Oncol Lett **2019**;17:1732-40

- 1 79. McGowan PM, Duffy MJ. Matrix metalloproteinase expression and outcome in patients with breast
2 cancer: analysis of a published database. *Ann Oncol* **2008**;19:1566-72
- 3 80. Boström P, Söderström M, Vahlberg T, Söderström KO, Roberts PJ, Carpén O, *et al.* MMP-1
4 expression has an independent prognostic value in breast cancer. *BMC Cancer* **2011**;11:348
- 5 81. Acerbi I, Cassereau L, Dean I, Shi Q, Au A, Park C, *et al.* Human breast cancer invasion and
6 aggression correlates with ECM stiffening and immune cell infiltration. *Integr Biol (Camb)*
7 **2015**;7:1120-34
- 8 82. González LO, Corte MD, Junquera S, González-Fernández R, del Casar JM, García C, *et al.*
9 Expression and prognostic significance of metalloproteases and their inhibitors in luminal A and
10 basal-like phenotypes of breast carcinoma. *Hum Pathol* **2009**;40:1224-33
- 11 83. Xia Y, Fan C, Hoadley KA, Parker JS, Perou CM. Genetic determinants of the molecular portraits of
12 epithelial cancers. *Nat Commun* **2019**;10:5666

13

14 **FIGURE LEGENDS**

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

23

24 **Figure 2.** *Permutation tests and associations between TWAS-genes and CRS for WW and BW.* A) Effect
25 estimates correspond to change in ROR-S, Proliferation score, and ROR-P per one standard deviation
26 increase in TWAS-gene expression (i.e., one standard deviation increase in GReX of gene). Triangle
27 denotes WW and circle denotes BW. B) Boxplots correspond to null distributions (shuffled GReX-sample
28 labels on left, random set of genes on right) of covariates residualized-R2 for regressions of CRS on

1 TWAS-genes. Null distributions are provided for 10,000 permutations of the GReX-sample labels and
2 10,000 random sets of genes. Dashed horizontal lines correspond to observed covariates residualized-
3 R2.

4
5 **Figure 3.** *Associations between TWAS-genes and PAM50 SCCs and gene expression.* A) Among AA
6 (top) and WW (bottom), associations between TWAS-genes and PAM50 SCCs using Bayesian
7 multivariate regression and multivariate adaptive shrinkage. Effect estimates show change in SCCs
8 (range -1 to 1) for one standard deviation increase in TWAS-gene GReX. Circle, triangle, and square
9 denote corresponding LFSR intervals for effect sizes. B) Heatmap of change in log₂ normalized PAM50
10 tumor expression for one standard deviation increase in TWAS-GReX. *, **, *** denote FDR intervals for
11 effect sizes.

1 **TABLES**

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

Gene	Signature	WW (N = 1,043)			AA (N = 1,083)		
		ROR-S	Proliferation	ROR-P	ROR-S	Proliferation	ROR-P
MCM10	IGF	3.03 (1.73, 4.33)	0.06 (0.03, 0.08)	3.11 (1.72, 4.50)	-	-	-
FAM64A	IGF	2.57 (1.28, 3.86)	0.05 (0.02, 0.07)	2.64 (1.26, 4.02)	-	-	-
CCNB2	Estradiol	2.69 (1.40, 3.98)	0.05 (0.02, 0.08)	2.71 (1.33, 4.09)	-	-	-
MMP1	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)
VAV3	Other	-2.22 (-3.51, -0.93)	-0.04 (-0.07, -0.02)	-2.40 (-3.79, -1.03)	-	-	-
PCSK6	IGF	-2.16 (-3.45, -0.88)	-0.03 (-0.06, 0.00)	-1.88 (-3.25, -0.50)	-	-	-
GNG11	Claudin-low	-1.27 (-2.56, 0.02)	-0.02 (-0.05, 0.00)	-1.42 (-2.80, -0.05)	-	-	-