

TITLE

Isoform-level analyses of 6 cancers uncover extensive genetic risk mechanisms undetected at the gene-level

AUTHORS

Yung-Han Chang¹, S. Taylor Head², Tabitha Harrison³, Yao Yu², Chad D. Huff², Bogdan Pasaniuc⁴, Sara Lindström^{3,5}, Arjun Bhattacharya^{2,6}

AFFILIATIONS

1. Quantitative Sciences Program, The University of Texas MD Anderson Cancer Center UTHealth Houston Graduate School of Biomedical Sciences, Houston, TX, USA
2. Department of Epidemiology, University of Texas MD Anderson Cancer Center, Houston, TX, USA
3. Department of Epidemiology, School of Public Health, University of Washington, Seattle, WA, USA
4. Department of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA
5. Public Health Sciences Division, Fred Hutchinson Cancer Center, Seattle, WA, USA
6. Institute for Data Science in Oncology, University of Texas MD Anderson Cancer Center, Houston, TX, USA

Corresponding author: Arjun Bhattacharya (abhattacharya3@mdanderson.org)

ABSTRACT

Integrating genome-wide association study (GWAS) and transcriptomic datasets can help identify potential mediators for germline genetic risk of cancer. However, traditional methods have been largely unsuccessful because of an overreliance on total gene expression. These approaches overlook alternative splicing, which can produce multiple isoforms from the same gene, each with potentially different effects on cancer risk. Here, we integrate genetic and multi-tissue isoform-level gene expression data from the Genotype Tissue-Expression Project (GTEx, N = 108-574) with publicly available European-ancestry GWAS summary statistics (all N > 20,000 cases) to identify both isoform- and gene-level risk associations with six cancers (breast, endometrial, colorectal, lung, ovarian, prostate) and six related cancer subtype classifications (N = 12 total). Compared to traditional methods leveraging total gene expression, directly modeling isoform expression through transcriptome-wide association studies (isoTWAS) substantially increases discovery of transcriptomic mechanisms underlying genetic associations. Using the same RNA-seq datasets, isoTWAS identified 164% more significant unique gene associations compared to TWAS (6,163 and 2,336, respectively), with isoTWAS-prioritized genes enriched 4-fold for evolutionarily-constrained genes ($P = 6.1 \times 10^{-13}$). isoTWAS tags transcriptomic associations at 52% more independent GWAS loci compared to TWAS across the six cancers. Additionally, isoform expression mediates an estimated 63% greater proportion of cancer risk SNP heritability compared to gene expression when evaluating cis-genetic influence on isoform expression. We highlight several notable isoTWAS associations that demonstrate GWAS colocalization at the isoform level but not at the gene level, including, *CLPTM1L* (lung cancer), *LAMC1* (colorectal), and *BABAM1* (breast). These results underscore the critical importance of modeling isoform-level expression to maximize discovery of genetic risk mechanisms for cancers.

1 INTRODUCTION

2 Over the past twenty years, genome-wide association studies (GWASs) have successfully linked hundreds of
3 genetic variants associated with increased risks for common cancers, including breast, prostate, lung, and
4 colorectal cancers¹. For example, Zhang and Ahearn et al. identified 32 new loci associated with breast cancer
5 risk in European ancestry individuals², totaling over 200 loci identified for overall and subtype-specific breast
6 cancer risk. Similarly, Wang and Shen et al. identified 187 novel risk associations in a multi-ancestry prostate
7 cancer GWAS³. Despite these findings underscoring the potential of GWASs to uncover genetic risk factors for
8 cancer, a significant challenge remains as most genome-wide significant variants are in non-coding regions of
9 the human genome, making it difficult to understand the biological mechanisms underlying these genetic
10 associations.

11
12 To address this gap, transcriptome-wide association studies (TWASs) have emerged as a powerful
13 complementary approach, providing valuable insights into the functional impact of genetic variants⁴⁻⁸. By
14 focusing on the regulatory effects of genetic variants on gene expression, TWASs can identify potential target
15 genes and pathways involved in disease processes. This integrated approach enhances our ability to translate
16 genetic findings from GWASs into a deeper molecular understanding of disease etiology, which is necessary for
17 improved diagnostics, prevention strategies, and therapeutic interventions in oncology.

18
19 Despite the promise of TWASs and related approaches, recent work suggests that gene expression may be a
20 poor mediator of genetic risk of complex traits, potentially owing to differences in evolutionary pressures on
21 genetic variants with large effects on complex traits and those with large effects on molecular phenotypes^{9,10}. It
22 is possible that part of this poor mediation is due to an overreliance on total gene expression as a fundamental
23 unit of measure of the transcriptome. Total gene expression does not account for alternative splicing that can
24 produce multiple isoforms from the same gene, which are often under subtle genetic or environmental control.
25 The mRNA and protein isoforms generated through alternative processing of primary RNA transcripts can vary
26 in function and potentially affect cancer risk in different ways^{11,12}. To address this limitation, we developed
27 isoform-level TWAS (isoTWAS), a multivariate, stepwise framework that integrates genetic and isoform-level
28 transcriptomic variation with GWAS summary statistics. This approach identifies both isoform- and gene-trait
29 associations while controlling the false discovery rate, identifying transcriptomic associations at far more GWAS
30 loci compared to traditional gene-level TWAS^{13,14}.

31
32 In this study, we utilized isoTWAS to assess the association between isoform expression and risk of six cancer
33 types (breast^{15,16}, endometrial¹⁷, colorectal¹⁸, lung¹⁹, ovarian²⁰, and prostate²¹), and their subtype classifications
34 where applicable. We used publicly available European-ancestry cancer GWAS summary statistics and multi-
35 tissue isoform- and gene-level expression data from the Genotype Tissue-Expression Project (GTEx, N > 100)²².
36 isoTWAS reveals more associated loci than traditional TWAS, highlighting its potential to uncover novel genetic
37 insights and enhance our understanding of cancer biology.

RESULTS

Overview of isoTWAS

We conducted isoform- (isoTWAS) and gene-level transcriptome-wide association studies (TWAS) for a total of 12 cancer outcomes: (1-3) breast (BRCA; overall, estrogen receptor (ER)+, ER-)^{15,16}, (4) colorectal (CRC)¹⁸, (5) lung (LUNG; overall)¹⁹, (6) lung adenocarcinoma (LUAD)¹⁹, (7) lung squamous cell carcinoma (LUSC)¹⁹ (8-9) ovarian (OVCA; overall, serous)²⁰, (10-11) prostate (PRCA; overall, advanced)²¹, and (12) endometrial (UCEC)¹⁷. We integrated GWAS summary statistics (N = 63,053-228,951) with multi-tissue expression QTL data from individuals without cancer collected post-mortem from the Genotype Tissue-Expression Project (GTEx, N = 108-574). GWAS and GTEx sample sizes and assignments of relevant tissues to cancer outcomes are provided in **Supplemental Tables 1-2**. See **Methods** for further details.

We used univariate and multivariate penalized regressions (**Figure 1A**) to train gene and isoform expression models and map gene- and isoform-level risk associations for the 12 cancer outcomes using a weighted burden test (**Figure 1B**)^{5,23}. In isoTWAS, we inferred gene-level associations by combining isoform associations via aggregated Cauchy association²⁴. False discovery rate is controlled for via a Benjamini-Hochberg procedure, and family-wise error rates are controlled via Shaffer's modified sequentially rejective Bonferroni (MSRB) for correlated isoforms. Post-hoc analyses include permutation tests and fine-mapping to identify isoform sets driving significant associations and credible sets of genes or isoforms (**Figure 1C**). See **Methods** for further details.

Isoform-level analysis reveals cancer risk associations not detected at the gene-level

In total, isoTWAS identifies 11,078 significant isoform associations across 6,163 unique genes, representing a ~164% increase in gene identification as compared to traditional TWAS which identified 2,336 significant gene associations, (**Figure 2A; Supplemental Table S3-5, Supplemental Data S1-2**). Of the genes identified, 10% from isoTWAS and 7% from traditional TWAS are included in the OncoKB Cancer Gene List, representing a moderate enrichment of hallmark cancer genes among isoTWAS gene (P = 0.03)^{25,26}. For both TWAS and isoTWAS, most identified gene associations are observed for BRCA, ER+ BRCA, and PRCA, which also are the three largest GWAS. Specifically, isoTWAS identifies 3.12 times as many associations for BRCA (3,889 isoforms), 2.85 times more for PRCA (1,849 isoforms), and 2.39 times as many for ER+ BRCA (2,304 isoforms) as compared to TWAS, while showing relatively similar performance for LUNG; Manhattan plots of TWAS and isoTWAS results are shown in **Supplemental Figures S1-S12**. Given that the mean of the χ^2 distribution is linearly related to power and sample size, the percent increase in the test statistic serves as a measure of power or effective sample size. For $\chi^2 > 1$, we calculated the percent increase for isoTWAS-based associations compared to TWAS-based associations (**Figure 2B**). Across the 12 cancer outcomes, there is an average 25.3-37.4% increase in effective sample size for isoTWAS compared to TWAS, suggesting that isoTWAS may enhance discovery in real data as compared to traditional TWAS. Several key biological pathways are enriched

1 for cancer-specific sets of isoTWAS-prioritized genes, including cell cycle and mitosis regulation, DNA and RNA
2 binding, immune pathways, as well as downstream targets of cancer-relevant transcription factors like *ESR1*²⁷,
3 *RUNX2*²⁸, and *YY1*²⁹ (**Supplemental Table S6, Supplemental Figures S13-S18**).

4 We further explored if cancer risk genes identified by TWAS and isoTWAS capture true disease signals by
5 identifying genes under selective constraint. If a gene is constrained, selection will act to remove variants that
6 diminish gene function from the population, such as loss-of-function (LOF) variants. Here, we used Bayesian
7 estimates of the s_{het} measure of constraint³⁰, which, unlike traditional measures of constraint, is not biased
8 towards longer genes. In total, 19.9% of isoTWAS gene associations (1,226 of 6,163) show $s_{\text{het}} > 0.1$, compared
9 to 12.5% of TWAS gene association (293 of 2,336) with $s_{\text{het}} > 0.1$. Not only does this represent a significant
10 enrichment of high s_{het} among isoTWAS-prioritized gene associations compared to TWAS-prioritized gene
11 associations (χ^2 test $P = 3.7 \times 10^{-15}$), isoTWAS-prioritized genes are significantly enriched compared to the
12 genome-wide proportion of high s_{het} genes (14.7%, $\chi^2 P < 2.2 \times 10^{-16}$). In addition, we find significant enrichments
13 of high s_{het} genes among transcriptome-wide significant genes for ER- BRCA, LUNG, LUSC, PRCA, Adv PRCA,
14 and UCEC (χ^2 FDR-adjusted P-value < 0.05 ; **Figure 2C, Supplemental Table S4**).

15 Lastly, isoTWAS identifies 52 genes (34 undetected by TWAS) that are associated with five or more cancer
16 outcomes (**Figure 2D**), including multiple known oncogenes or tumor suppressor genes, such as *MYC*³¹,
17 *MUTYH*³², *GNAI2*^{33,34}, *ACO2*³⁵, and *BMI1*³⁶. The 34 genes undetected by TWAS are enriched for multiple salient
18 pathways: regulation of CD8-positive and T cells, regulation of membrane protein complexes, and phagocytic
19 and autolysosome function. Additionally, these 34 genes are highly enriched for downstream targets of crucial
20 oncogenic transcription factors, like *ESR1*, *GATA4*, *YY1*, and *MYC*, all of which were determined using ChIP-
21 Seq experiments in human or mouse cancer cells or tumors (**Figure 2E**)³⁷. Altogether, not only can isoTWAS
22 identify constrained susceptibility genes for multiple cancers, but it can also reveal pan-cancer risk signals that
23 TWAS is unable to detect.

24 ***Isoform expression explains more GWAS loci and overall SNP heritability of cancer risk***

25 In addition to increasing the discovery of susceptibility genes across the genome, isoTWAS can identify far more
26 transcriptomic mechanisms within independent, high-confidence genome-wide significant loci. Across all cancer
27 outcomes, we identify 622 risk-associated GWAS SNPs ($P < 5 \times 10^{-8}$) in independent linkage disequilibrium (LD)
28 blocks (see **Methods**). Among these 622 loci, 288 (46.3%) are tagged by TWAS, and 439 (70.6%) are tagged
29 by isoTWAS, with 249 (40.0%) loci identified by both TWAS and isoTWAS (**Figure 3A, Supplemental Table**
30 **S7**). In total, this represents an increase in significant associations in GWAS loci of 52.4% when using isoTWAS,
31 rather than TWAS. Additionally, isoTWAS identified 2,911 unique significant genes outside of known GWAS loci.
32

33 We also considered an orthogonal analysis to compare GWAS follow-up using both gene and isoform
34 expression. We mapped *cis*-expression quantitative trait loci (eQTLs) for gene and isoform expression and
35 conducted Bayesian colocalization analysis with GWAS signals using eCAVIAR³⁸ to estimate the CoLocalization

1 Posterior Probability (CLPP; see **Methods; Supplemental Table S8, Supplemental Data S3**). Generally,
2 isoTWAS exhibits higher CLPPs compared to TWAS, indicating a stronger likelihood of colocalization for isoform-
3 eQTLs with GWAS loci (**Figure 3B**), with the proportion of loci having isoform expression Quantitative Trait Loci
4 (isoQTL) colocalizations with CLPP > 0.01 consistently higher across the 12 cancer outcomes (**Supplemental**
5 **Figure S19**). Notably, in OVCA and UCEC, isoTWAS shows significantly higher CLPPs (median = 0.026, P =
6 0.013) than TWAS (median = 0.001, P = 0.001), suggesting that isoTWAS more effectively captures colocalized
7 signals in for these cancer outcomes. Additionally, we estimated the proportion of total SNP heritability (h^2)
8 mediated by gene- and isoform-level expression (h^2_{med}), (**Figure 3C** and **Supplemental Table S9**). Overall,
9 isoform-level expression explains 62.7% more of cancer risk SNP heritability ($19.2 \pm 10.3\%$) compared to gene-
10 level expression ($11.8 \pm 5.7\%$). Wald-type tests reveal that isoform expression mediates a significant proportion
11 of cancer h^2 , with gene expression only explaining a significant portion for ER- BRCA (isoTWAS: 0.291, FDR-
12 adjusted P = 0.005; TWAS: 0.133, P = 0.068), ER+ BRCA (isoTWAS: 0.173, P = 0.045; TWAS: 0.110, P = 0.102),
13 LUNG (isoTWAS: 0.287, P = 0.044; TWAS: 0.181, P = 0.056), PRCA (isoTWAS: 0.314, P = 0.005; TWAS: 0.108,
14 P = 0.164), and Adv PRCA (isoTWAS: 0.358, P = 0.014; TWAS: 0.243, P = 0.046). These results indicate that
15 not only does isoTWAS recapitulate an overwhelming majority of TWAS signals at GWAS loci, isoTWAS
16 substantially increases discovery of candidate GWAS mechanisms and transcriptomic features that potentially
17 mediate genetic effects on cancer risk.

18 ***isoTWAS and isoform-eQTL colocalization prioritize undetected mechanisms at GWAS loci***

19 A main goal of isoform-specific analyses is to nominate a more granular hypothesis of transcriptomic regulation.
20 Post-hoc transcript-level fine-mapping reveals nine loci where isoTWAS prioritization coincides with isoform-
21 eQTL colocalization (CLPP > 0.01) but no gene-eQTL colocalization (**Methods; Supplemental Data S3**), across
22 BRCA (**Figure 4, Supplemental Figures S20**), ER- BRCA (**Supplemental Figure S21-S22**), CRC (**Figure 5,**
23 **Supplemental Figure S23**), LUNG (**Figure 6**), PRCA (**Supplemental Figure S24**), and UCEC (**Supplemental**
24 **Figure S25**). We highlight *CLPTM1L*, *LAMC1*, and *BABAM1* here due to previously-reported pleiotropic
25 associations across multiple cancers. *CLPTM1L* is located near the *TERT* locus, known for its pleiotropic links
26 to various cancers³⁹, the *LAMC1* locus has shown associations with PRCA⁴⁰, and *BABAM1* is a recognized
27 GWAS hit with broad implications for cancer risk. These three genes are prioritized using isoTWAS models
28 trained in subcutaneous adipose tissue, with additional examples presented in the **Supplemental Results**.

29
30
31 First, four isoforms of *CLPTM1L* (Chromosome 5p15.33, $s_{het} = 0.02$, 16 total isoforms in GENCODE) are
32 significantly associated with LUNG. We also find that isoforms of *CLPTM1L* are associated with BRCA, PRCA,
33 LUAD, and LUSC. Fine-mapping prioritizes ENST00000511268.6 in the 90% credible set with Posterior Inclusion
34 Probability (PIP) = 0.90. Although there are many genome-wide significant SNPs within the gene body of
35 *CLPTM1L*, no gene-eQTL signal is observed ($P < 10^{-6}$), whereas a strong ENST00000511268.6 isoform-eQTL
36 signal colocalizes with the GWAS signal (CLPP = 0.07). Additionally, there is no strong isoform-eQTL signal for
37 the other three isoforms identified via isoTWAS in this region. The lead isoQTL rs414965 is in high LD with

1 multiple genome-wide significant SNPs (**Figure 4A**). The exons comprising ENST00000511268.6 overlap with
2 those of ENST00000503534.5 and ENST00000506641.5 with significant isoform-eQTLs flanking
3 ENST00000511268.6 (**Figure 4B**). However, rs414965 does not have significant effects on
4 ENST00000503534.5, ENST00000506641.5, or ENST00000503151.5 expression (**Figure 4C**), despite its
5 strong negative associations with LUNG and ENST00000511268.6 expression.

6
7 Next, three isoforms of *LAMC1* (Chromosome 1q25.3, $s_{\text{het}} = 0.10$, six total isoforms in GENCODE) are associated
8 with CRC. Fine-mapping prioritizes only ENST00000466964.1 in the 90% credible set with PIP = 1. Again,
9 genome-wide significant CRC-associated SNPs within the gene body show no gene-eQTL signal ($P < 10^{-6}$) but
10 a strong ENST00000466964.1 isoform-eQTL signal that colocalizes with the GWAS signal (CLPP = 0.04) is
11 observed. For the other two isoforms, ENST00000478064.1 and ENST00000495918.1, the lead isoQTL rs20558
12 is in high LD with multiple genome-wide significant SNPs (**Figure 5A**). The exons of these three isoforms are
13 generally distinct sets with significant isoform-eQTLs generally falling on the 3' end of multiple exons (**Figure**
14 **5B**). rs20558 shows a significant association with increased CRC risk, no effect on *BABAM1* gene expression,
15 and a significant decreasing effect on ENST00000466964.1 expression (**Figure 5C**). Interestingly, this same
16 SNP has a significant increasing effect on the expression of ENST00000478064.1 and ENST00000495918.1,
17 though only ENST00000478064.1 colocalizes with CRC risk at CLPP > 0.01.

18
19 Lastly, nine isoforms of *BABAM1* (Chromosome 19p13.11, $s_{\text{het}} = 0.05$, 18 total isoforms in GENCODE) are
20 associated with BRCA. In our study, isoforms of *BABAM1* are also associated with OVCA, OVCA ser, LUSC,
21 and ER- BRCA in isoTWAS. Fine-mapping prioritizes ENST00000599474.5 in the 90% credible set and a PIP
22 of 1. There are multiple genome-wide significant BRCA-associated SNPs ($P < 5 \times 10^{-8}$) within the *BABAM1* gene
23 body showing no gene-eQTL signal ($P < 10^{-6}$). Only ENST00000599474.5 (CLPP = 0.44) and
24 ENST00000359435.8 (CLPP=0.59) show strong isoQTL effects that colocalize with the GWAS signal. In
25 addition, the lead isoQTL rs34084277 for ENST00000599474.5 is in high LD ($LD > 0.8$) with multiple genome-
26 wide significant SNPs (**Figure 6A**). The exon structure of these nine isoforms reveals a group of exons at the 3'
27 end of the gene body, which is flanked by the lead isoQTLs of ENST00000599474.5 (**Figure 6B**). Specifically,
28 due to its exon structure, we followed up on this locus using a rare-variant analysis in UK Biobank whole exome
29 sequencing data⁴¹ using the Variant Annotation, Analysis, and Search Tool (VAAS2)^{42,43}. We find three BRCA-
30 associated isoforms in *BABAM1* enriched for risk-associated rare variants (MAF < 0.5%), with associations
31 mainly concentrated in exons at the 5' end and first exon at the 3' end of the transcript (**Supplemental Methods**;
32 **Supplemental Figure 26**; **Supplemental Tables S10-S11**). See **Supplemental Results** for further details.
33 Lastly, rs34084277 has a strong protective effect on *BRCA* risk, no effect on *BABAM1* gene expression, and a
34 significantly increasing effect on expression of ENST00000599474.5 (**Figure 6C**).

35 36 DISCUSSION

1 We show that integrating cancer risk GWAS summary statistics with isoform-level transcriptomic variation can
2 greatly increase discovery of susceptibility genes for six cancers and their subtype classifications. By using
3 isoTWAS rather than traditional gene-level TWAS, we identify nearly 2.5-fold more associations. More saliently,
4 isoTWAS-identified genes are significantly enriched for evolutionarily constrained genes, which are more likely
5 to contain clinically-relevant *de novo* or rare variants, predict drug toxicity, and characterize transcriptional
6 regulation^{10,44}. Isoform expression, rather than total gene expression, captures more risk-associated genetic
7 variation and exhibits stronger colocalization with GWAS signals. Additionally, though isoform expression alone
8 does not reconcile most of the SNP heritability of cancer risk, isoform expression is estimated to mediate nearly
9 twice as much heritability as gene expression, and nearly three times of the SNP heritability signal in the case
10 of PRCA. Most importantly, isoTWAS can find isoform associations at 249 of 288 GWAS loci tagged by TWAS-
11 identified genes and uniquely tag 190 additional GWAS loci that cannot be contextualized via TWAS. In
12 aggregate, these findings underscore the utility of considering the transcriptome on the transcript-isoform, rather
13 than gene-level. By modeling a different quantification of the same steady-state RNA-seq datasets with sample
14 sizes of only up to ~800, isoTWAS increases discovery specifically at GWAS loci by ~52% without additional
15 sequencing costs.

16
17 Our isoform-level fine-mapping coupled with eQTL colocalization identifies nine gene candidates with GWAS risk
18 SNPs within the gene body that colocalize with isoQTLs despite exhibiting no gene-eQTL signal at $P < 10^{-6}$. We
19 discuss three of these gene candidates. First, isoTWAS identifies an association between ENST00000511268.6,
20 an isoform of *CLPTM1L* (16 total isoforms), and LUNG. Studies spanning back 15 years have identified multiple
21 polymorphisms within the *CLPTM1L* gene body associated with LUNG, as well as pleiotropic associations with
22 other malignancies^{45–49}. In particular, an *in vitro* study provided evidence that *CLPTM1L* is oncogenic, specifically
23 for Ras-driven lung cancers, in line with its effect on protecting tumor cells from genotoxic apoptosis⁵⁰, and is a
24 promising therapeutic target for therapy-resistant tumors⁵¹. *CLPTM1L* (Chromosome 5p15.33) is also local to
25 *TERT*, which harbors pleiotropic associations with multiple cancers³⁹. Our study identifies associations between
26 *CLPTM1L* and PRCA, consistent with previous findings of SNPs in the 5p15.33 region associated with PRCA^{52,53}.
27 However, while *TERT* is known for its pleiotropic effects, we observed only gene-level associations for *TERT* in
28 our analyses, with no isoTWAS associations detected for its specific isoforms. Further work is needed to fully
29 interrogate this pleiotropy and assess tissue-specific isoform expression for both *TERT* and *CLPTM1L*,
30 especially since *TERT* shows low expression across multiple tissues in GTEx²². We find a similar pattern for
31 colorectal cancer and ENST00000466964.1, an isoform of *LAMC1* (six total isoforms), a gene that regulates cell
32 adhesion, differentiation, migration, and signaling, and lies at a pleiotropic locus for CRC, PRCA⁴⁰, and obesity
33 risk⁵⁴. A previous study incorporating splicing measures has prioritized *LAMC1* as a novel transcriptome-
34 mediated CRC locus⁵⁵. Additionally, rs34295433 in *LAMC1* was identified as a susceptibility SNP for PRCA in a
35 Taiwanese population⁴⁰. However, in our study, isoTWAS did not detect any associations between *LAMC1* and
36 PRCA. *LAMC1* belongs to the laminin family of extracellular matrix proteins that are significantly involved in
37 survival and proliferation of cancer cells, angiogenesis, migration and basement membrane breach by cancer

1 cells, and metastatic events⁵⁶. Computational analyses of laminin proteins have shown their prognostic ability in
2 colorectal cancer progression, placing higher weights for *LAMC1* compared to other constituents of the family⁵⁷.

3
4 Lastly, isoTWAS detects an association between BRCA and ENST00000599474.5, an isoform of *BABAM1* (18
5 total isoforms), a gene involved in checkpoint signaling, regulation of DNA repair, and mitosis. *BABAM1* has
6 previously been implicated as a low-penetrance risk locus that interacts with *BRCA1* in both triple-negative breast
7 cancer and ovarian cancer risk⁵⁸⁻⁶⁰. Consistent with previous studies, isoTWAS also identifies associations within
8 this region for OVCA, OVCA ser, and ER- BRCA. Additionally, our VAAST analysis for *BABAM1* in UKBB
9 indicates that a group of exons at the 3' end of the gene harbor potentially disease-causing rare variants (MAF
10 < 0.01, only in coding regions). The lead isoQTL for our isoTWAS-prioritized isoform of *BABAM1* (MAF > 0.01,
11 both coding and non-coding region) is upstream of these exons at the 3' end, potentially influencing the splicing
12 patterns specifically at the 3' end of the gene. Further research is required to investigate how these isoQTLs
13 influence splicing regulation and their specific role in tumorigenesis and cancer development, but a
14 methodological opportunity may lie in integrating splicing- and isoQTLs with rare variants to identify
15 transcriptomic mechanisms for cancer risk. These results underscore that the added resolution provided with
16 isoform expression can lead to more specific mechanistic hypotheses for cancer risk and inform follow-up with
17 functional studies, both *in silico* and experimental.

18
19 We conclude with three limitations of our work. First, the complexity of the *BABAM1* transcript structure and the
20 sheer number of *BABAM1* isoforms with strong isoQTL signals underscores a methodological opportunity. We
21 applied the FOCUS framework which leverages a non-informative prior that may be insufficient⁶¹. Fine-mapping
22 of isoforms is challenging due to horizontal pleiotropy of SNP-isoform effects shared across exons and strong
23 LD patterns. This horizontal pleiotropy can reduce power and increase false-positive rates. Incorporating classes
24 of isoforms with shared exon sets may lead to improved fine-mapping, thereby increasing coverage and
25 resolution of credible sets of isoforms within a GWAS locus.

26
27 Second, we derived isoform-level quantifications using Salmon, a method constrained by the limitations of short-
28 read RNA-seq where maximum-likelihood estimates are based on transcriptomic annotations, introducing some
29 uncertainty. The limited diversity of tissue-specific annotated isoforms may affect the accuracy of these
30 estimates. In contrast, long-read RNA-seq can capture splicing and structural variation missed by short-read
31 RNA-seq, revealing more complex and rarer transcript-isoforms. As long-read RNA-seq becomes increasingly
32 scalable, its comprehensive view of transcriptomes will allow for more precise quantification of isoforms
33 potentially improving isoTWAS performance. Third, as current eQTL datasets are mostly generated in
34 populations of European ancestry, we focused the current analyses on cancer GWAS summary statistics based
35 on individuals of European ancestry. Previous research show that expression models trained on predominantly
36 European-ancestry cohorts often fail to predict expression accurately in individuals of different ancestries⁶²⁻⁶⁵. A
37 critical emphasis of future data generation should be to increase diversity in molecular datasets and for

1 methodological research to develop computational frameworks that can model multi-ancestry samples across
2 eQTL and GWAS datasets.

3 4 **METHODS**

5 ***Data collection***

6 *GTEx genotype and transcriptomic data*

7 For 48 tissues from GTEx²², we quantified RNA-seq data (all N > 100) using Salmon v1.5.2 in mapping-based
8 mode⁶⁶. We built a decoy-aware transcriptomic index in Salmon with GENCODE v38 transcript sequences and
9 the full GRCh38 reference genome as decoy sequences⁶⁷. Salmon was then run on FASTQ files with mapping
10 validation and corrections for sequencing and GC bias. We then imported Salmon isoform-level quantifications
11 and aggregated to the gene-level using tximeta v1.16.1⁶⁸. Using edgeR, gene and isoform-level quantifications
12 underwent TMM-normalization, followed by transformation into a log-space using the variance-stabilizing
13 transformation using DESeq2 v1.38.3^{69,70}. We then residualized isoform-level and gene-level expression (as log-
14 transformed CPM) by all tissue-specific covariates (clinical, demographic, genotype principal components (PCs),
15 and expression PEER factors) used in the original QTL analyses in GTEx.

16
17 SNP genotype calls were derived from Whole Genome Sequencing data from individuals of European ancestry,
18 filtering out SNPs with minor allele frequency (MAF) less than 5% or that deviated from HWE at $P < 10^{-5}$. We
19 further filtered out SNPs with MAF less than 1% frequency among the European ancestry samples in 1000
20 Genomes Project⁷¹. Due to limited sample sizes in other ancestry groups in GTEx, we focused solely on
21 European ancestry for this study.

22 23 ***GWAS summary statistics***

24 We obtained GWAS summary statistics for risk of 12 cancer outcomes, all from samples of European-ancestry
25 individuals⁷²: overall breast cancer (BRCA; 122,977 cases/105,974 controls)¹⁵, estrogen-receptor positive breast
26 cancer (ER+ BRCA; 69,501 cases/95,042 controls)¹⁵, estrogen-receptor negative breast cancer (ER- BRCA;
27 30,882 cases/110,058 controls)¹⁶, Colorectal cancer (CRC; 55,168 cases/65,160 controls)¹⁸, overall lung cancer
28 (LUNG; 29,266 cases/56,450 controls)¹⁹, lung adenocarcinoma (LUAD; 11,273 cases/55,483 controls)¹⁹, lung
29 squamous cell Carcinoma (LUSC; 7,426 cases/55,627 controls)¹⁹, overall ovarian cancer (OVCA; 22,406
30 cases/40,951 controls)²⁰, serous ovarian cancer (OVCA ser; 19,890 cases/68,502 controls)²⁰, overall prostate
31 cancer (PRCA; 79,166 cases/61,106 controls)²¹, advanced prostate cancer (Adv PRCA; 15,167 cases/58,308
32 controls)²¹, endometrial cancer (UCEC; 12,906 cases/108,979 controls)¹⁷.

33 34 ***Statistical analysis***

35 *Code and data availability*

36 Sample scripts for analyses are available from github.com/bhattacharya-a-bt/MultiCancerIsoTWAS. Links to
37 GWAS summary statistics are provided in **Supplemental Table S1**. isoTWAS and TWAS models are available

1 from <https://zenodo.org/records/11048201>⁷³. GTEx genotyping and expression data were accessible from
2 dbGaP accession number phs000424.v9. Supplemental Data Tables S1-S3 can be downloaded from
3 <https://zenodo.org/records/14010391>.
4

5 ***Isoform- and gene-level transcriptome-wide association studies***

6 We trained predictive models of gene and isoform expression using all *cis*-SNPs within 1 Mb of the gene body
7 (**Figure 1A**). For gene expression, we implemented three methods: (1) elastic net regression⁷⁴, (2) linear mixed
8 models using a ridge regression approximation⁷⁵, and (3) Sum of Single Effects (SuSiE) regression⁷⁶ (commonly
9 used in FUSION⁵) and then selected the model with the best prediction based on 10-fold cross-validation⁵. For
10 isoform expression, we considered the matrix of expression for all isoforms of a gene and implemented four
11 multivariate methods: (1) multivariate elastic net regression, (2) multivariate LASSO with covariance estimation,
12 (3) multivariate elastic net with a stacked generalization, and (4) multivariate sparse partial least squares. If a
13 gene only had one isoform, we trained univariate penalized regressions, as done for gene expression. Again,
14 we selected the best isoform model through a 10-fold cross-validation. We only retained genes and isoforms that
15 could be predicted at cross-validation $R^2 > 0.01$.
16

17 Using these predictive models, we employed a weighted burden test to estimate the association between the
18 genetically-predicted component of a gene or isoform with cancer risk (**Figure 1B**). We denoted \hat{w} as the
19 prediction weights for a gene or isoform model. We denoted Z as the vector of standardized effect sizes from the
20 GWAS and Σ as the LD matrix between the SNPs represented in Z . We estimated the standardized effect size
21 of the genetically-predicted component of a gene or isoform's expression on cancer risk as $T = \frac{\hat{w}Z}{\hat{w}^T \Sigma \hat{w}}$. As one
22 gene is likely to carry multiple isoform associations, isoform-level mapping has an increased testing burden
23 (**Figure 1C**). Accordingly, we used the aggregated Cauchy association to combine the t test statistics T_1, \dots, T_t
24 for isoforms of a single gene to estimate the isoTwas gene-level association. For both TWAS and isoTwas, we
25 controlled for false discovery across all genes via the Benjamini-Hochberg procedure. For isoTwas, we first
26 identified isoforms for genes with an FDR-adjusted ACAT $P < 0.05$, and then employed Shaffer's modified
27 sequentially rejective Bonferroni procedure to control the within-gene family-wide error rate (FWER). At the end
28 of these two steps, isoTwas identified a set of genes and their isoforms that were associated with the trait.
29

30 For significant genes (from TWAS) and isoforms (from isoTwas), we tested whether the SNP-gene/isoform
31 effects from the predictive models add additional information beyond the SNP-risk effects from the GWAS
32 through a conservative permutation test. We permuted the SNP-gene/isoform effects in the predictive models
33 10,000 times and generated a null distribution for the gene/isoform test statistic. We used this null distribution to
34 generate a permutation-based P-value for the original test statistic for each gene/isoform. We only conducted
35 this permutation test for genes from TWAS with FDR-adjusted $P < 0.05$ and isoforms from isoTwas with FDR-
36 adjusted $P < 0.05$ and FWER-adjusted $P < 0.05$.
37

1 After permutation testing, we conducted isoform- and gene-level fine-mapping for isoTWAS and TWAS
2 associations that overlap in a 1 Mb window using methods from the FOCUS framework⁶¹. We accounted for the
3 correlation between genetically-predicted isoform or gene expression induced by LD and shared prediction
4 weights and control for certain pleiotropic effects. Through Bayesian methods, we estimated the Posterior
5 Inclusion Probability (PIP) for each isoform or gene and defined a credible set of isoforms or genes to explain
6 the signal with 90% confidence.

7 **Gene-set enrichment analysis**

9 We conducted gene-set enrichment analysis using enrichR³⁷ to investigate the gene ontologies enriched among
10 genes identified by isoTWAS. The analysis queried multiple databases, including GO Biological Process 2023,
11 GO Cellular Component 2023, GO Molecular Function 2023, KEGG 2021 Human, Reactome 2022, and ChEA
12 2022. ChEA is a curated database of transcription factor targets across multiple ChIP-chip, ChIP-seq, ChIP-PET
13 and DamID assays across a host of tissues and cell-types⁷⁷. We extracted the top 10 most significant terms from
14 each database (FDR-adjusted $P < 0.1$) and calculated the odds ratios to identify over-represented gene
15 functions.

17 **Quantitative trait locus mapping and Bayesian colocalization**

18 For the same GTEx tissues that were used in TWAS/isoTWAS, we mapped *cis*-eQTLs for all genes and isoforms
19 using all SNPs within a 1 Mb interval around the transcription start site. We used ordinary least squares to
20 estimate the allelic effect of the SNP (coded as 0, 1, or 2 copies of the risk allele) on gene or isoform expression,
21 adjusted for the same variables used in the original GTEx analyses (five principal components of the genotype
22 matrix, 30-60 PEER factors depending on sample size for the given tissue, age at death, sex). We estimated the
23 eQTL effect sizes, Wald-type standard errors, and P-values.

25 Next, using GWAS summary statistics for each of the 12 cancer outcomes, gene- and isoform-eQTL summary
26 statistics, and reference LD for European ancestry individuals from the 1000 Genomes Project⁷¹, we conducted
27 Bayesian colocalization using eCAVIAR³⁸ to estimate the CoLocalization Posterior Probability (CLPP) that the
28 same SNP is causal for both cancer risk and the gene or isoform expression. We considered a gene or isoform
29 to colocalize with a GWAS-significant locus if three conditions were met: a SNP has a (1) strong effect on cancer
30 risk (GWAS $P < 5 \times 10^{-8}$), (2) strong effect on gene or isoform expression (eQTL $P < 10^{-6}$), and (3) CLPP > 0.01 ,
31 as is proposed in the original eCAVIAR paper.

33 **Estimation of expression-mediated SNP heritability**

34 We employed mediated expression score regression (MESCR)⁹ to estimate the proportion of total SNP heritability
35 (h^2) mediated by the *cis*-genetic component of gene or isoform expression levels (h^2_{med}). First, using the
36 genotypes of all SNPs, expression of all genes or isoforms across all 48 tissues in GTEx, and the same
37 covariates used in the eQTL analysis above, we estimated eQTL effect sizes in each tissue using LASSO

1 regression (`./run_mesc.py --compute-expscore-indiv ...`). Next, we meta-analyzed expression scores
2 across the 49 tissues and estimated h^2_{med} for each of the 12 cancer outcomes by using GWAS summary statistics
3 munged to the LD score regression `.sumstats` format. This analysis provided estimates of h^2 , h^2_{med} , and
4 h^2_{med}/h^2 , each with standard errors and P-values for the hypothesis test comparing to the null value of no
5 heritability.

7 **AUTHORS' DISCLOSURES**

8 No disclosures are reported by the authors.

10 **AUTHORS' CONTRIBUTIONS**

11 YHC: Formal analysis, validation, investigation, visualization, methodology, writing—original draft, writing—review
12 and editing; STH: Formal analysis, writing—review and editing; TH: Formal analysis, writing—review and editing;
13 YY: Formal analysis, validation, visualization, methodology, writing—review and editing; CDH: Methodology,
14 writing—review and editing, funding acquisition; BP: Writing—review and editing, resources, funding acquisition;
15 SL: Writing—review and editing, funding acquisition; AB: Conceptualization, resources, supervision, funding
16 acquisition, investigation, methodology, writing—review and editing

18 **ACKNOWLEDGEMENTS**

19 Funding was provided by NIH CA293419 and CA194393. The funder had no role in study design, data collection,
20 analysis, or interpretation, or writing of the article.

1 REFERENCES

- 2 1. Sud, A., Kinnersley, B. & Houlston, R. S. Genome-wide association studies of cancer: current insights and
3 future perspectives. *Nat Rev Cancer* **17**, 692–704 (2017).
- 4 2. Zhang, H. *et al.* Genome-wide association study identifies 32 novel breast cancer susceptibility loci from
5 overall and subtype-specific analyses. *Nat Genet* **52**, 572–581 (2020).
- 6 3. Wang, A. *et al.* Characterizing prostate cancer risk through multi-ancestry genome-wide discovery of 187
7 novel risk variants. *Nat Genet* **55**, 2065–2074 (2023).
- 8 4. Gamazon, E. R. *et al.* A gene-based association method for mapping traits using reference transcriptome
9 data. *Nat Genet* **47**, 1091–1098 (2015).
- 10 5. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies. *Nat Genet*
11 **48**, 245–252 (2016).
- 12 6. Wainberg, M. *et al.* Opportunities and challenges for transcriptome-wide association studies. *Nat Genet* **51**,
13 592–599 (2019).
- 14 7. Mancuso, N. *et al.* Large-scale transcriptome-wide association study identifies new prostate cancer risk
15 regions. *Nat Commun* **9**, 4079 (2018).
- 16 8. Gao, G. *et al.* A joint transcriptome-wide association study across multiple tissues identifies candidate
17 breast cancer susceptibility genes. *The American Journal of Human Genetics* **110**, 950–962 (2023).
- 18 9. Yao, D. W., O'Connor, L. J., Price, A. L. & Gusev, A. Quantifying genetic effects on disease mediated by
19 assayed gene expression levels. *Nat Genet* **52**, 626–633 (2020).
- 20 10. Mostafavi, H., Spence, J. P., Naqvi, S. & Pritchard, J. K. Systematic differences in discovery of genetic
21 effects on gene expression and complex traits. *Nat Genet* **55**, 1866–1875 (2023).
- 22 11. Zhang, Y., Qian, J., Gu, C. & Yang, Y. Alternative splicing and cancer: a systematic review. *Sig Transduct*
23 *Target Ther* **6**, 1–14 (2021).
- 24 12. Wang, E. & Aifantis, I. RNA Splicing and Cancer. *Trends in Cancer* **6**, 631–644 (2020).
- 25 13. Bhattacharya, A. *et al.* Isoform-level transcriptome-wide association uncovers genetic risk mechanisms for
26 neuropsychiatric disorders in the human brain. *Nat Genet* **55**, 2117–2128 (2023).
- 27 14. O'Connell, K. S. *et al.* Genomics yields biological and phenotypic insights into bipolar disorder.
28 2023.10.07.23296687 Preprint at <https://doi.org/10.1101/2023.10.07.23296687> (2024).

- 1 15. K, M. *et al.* Association analysis identifies 65 new breast cancer risk loci. *Nature* **551**, (2017).
- 2 16. RI, M. *et al.* Identification of ten variants associated with risk of estrogen-receptor-negative breast cancer.
- 3 *Nature genetics* **49**, (2017).
- 4 17. O'Mara, T. A. *et al.* Identification of nine new susceptibility loci for endometrial cancer. *Nat Commun* **9**,
- 5 3166 (2018).
- 6 18. Huyghe, J. R. *et al.* Discovery of common and rare genetic risk variants for colorectal cancer. *Nat Genet*
- 7 **51**, 76–87 (2019).
- 8 19. McKay, J. D. *et al.* Large-scale association analysis identifies new lung cancer susceptibility loci and
- 9 heterogeneity in genetic susceptibility across histological subtypes. *Nat Genet* **49**, 1126–1132 (2017).
- 10 20. Phelan, C. M. *et al.* Identification of 12 new susceptibility loci for different histotypes of epithelial ovarian
- 11 cancer. *Nat Genet* **49**, 680–691 (2017).
- 12 21. Schumacher, F. R. *et al.* Association analyses of more than 140,000 men identify 63 new prostate cancer
- 13 susceptibility loci. *Nat Genet* **50**, 928–936 (2018).
- 14 22. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues.
- 15 *Science* **369**, 1318–1330 (2020).
- 16 23. Pasaniuc, B. *et al.* Fast and accurate imputation of summary statistics enhances evidence of functional
- 17 enrichment. *Bioinformatics* **30**, 2906–2914 (2014).
- 18 24. Liu, Y. *et al.* ACAT: A Fast and Powerful p Value Combination Method for Rare-Variant Analysis in
- 19 Sequencing Studies. *Am J Hum Genet* **104**, 410–421 (2019).
- 20 25. Chakravarty, D. *et al.* OncoKB: A Precision Oncology Knowledge Base. *JCO Precis Oncol* 1–16 (2017)
- 21 doi:10.1200/PO.17.00011.
- 22 26. Suehnholz, S. P. *et al.* Quantifying the Expanding Landscape of Clinical Actionability for Patients with
- 23 Cancer. *Cancer Discovery* **14**, 49–65 (2024).
- 24 27. Carroll, J. S. Mechanisms of oestrogen receptor (ER) gene regulation in breast cancer. *Eur J Endocrinol*
- 25 **175**, R41-49 (2016).
- 26 28. Pulica, R., Cohen-Solal, K. & Lasfar, A. Evaluating the Role of RUNX2 in Cancer and Its Potential as a
- 27 Therapeutic Target. in *Handbook of Cancer and Immunology* (ed. Rezaei, N.) 1–22 (Springer International
- 28 Publishing, Cham, 2022). doi:10.1007/978-3-030-80962-1_254-1.

- 1 29. Agarwal, N. & Theodorescu, D. The Role of Transcription Factor YY1 in the Biology of Cancer. *Crit Rev*
2 *Oncog* **22**, 13–21 (2017).
- 3 30. Zeng, T., Spence, J. P., Mostafavi, H. & Pritchard, J. K. Bayesian estimation of gene constraint from an
4 evolutionary model with gene features. *Nat Genet* **56**, 1632–1643 (2024).
- 5 31. Chen, H., Liu, H. & Qing, G. Targeting oncogenic Myc as a strategy for cancer treatment. *Sig Transduct*
6 *Target Ther* **3**, 1–7 (2018).
- 7 32. Paller, C. J. *et al.* Pan-Cancer Interrogation of MUTYH Variants Reveals Biallelic Inactivation and Defective
8 Base Excision Repair Across a Spectrum of Solid Tumors. *JCO Precis Oncol* e2300251 (2024)
9 doi:10.1200/PO.23.00251.
- 10 33. Gupta, S. K., Gallego, C., Johnson, G. L. & Heasley, L. E. MAP kinase is constitutively activated in gip2
11 and src transformed rat 1a fibroblasts. *Journal of Biological Chemistry* **267**, 7987–7990 (1992).
- 12 34. Ikezu, T. *et al.* Bidirectional regulation of c-fos promoter by an oncogenic gip2 mutant of G alpha i2. A novel
13 implication of retinoblastoma gene product. *Journal of Biological Chemistry* **269**, 31955–31961 (1994).
- 14 35. Wang, Z. *et al.* Pan-Cancer analysis shows that ACO2 is a potential prognostic and immunotherapeutic
15 biomarker for multiple cancer types including hepatocellular carcinoma. *Front Oncol* **12**, 1055376 (2022).
- 16 36. Xu, J., Li, L., Shi, P., Cui, H. & Yang, L. The Crucial Roles of Bmi-1 in Cancer: Implications in
17 Pathogenesis, Metastasis, Drug Resistance, and Targeted Therapies. *Int J Mol Sci* **23**, 8231 (2022).
- 18 37. Chen, E. Y. *et al.* Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC*
19 *Bioinformatics* **14**, 128 (2013).
- 20 38. Hormozdiari, F. *et al.* Colocalization of GWAS and eQTL Signals Detects Target Genes. *Am J Hum Genet*
21 **99**, 1245–1260 (2016).
- 22 39. Chen, H. *et al.* Large-scale cross-cancer fine-mapping of the 5p15.33 region reveals multiple independent
23 signals. *HGG Adv* **2**, 100041 (2021).
- 24 40. Bau, D.-T. *et al.* Genetic susceptibility to prostate cancer in Taiwan: A genome-wide association study. *Mol*
25 *Carcinog* **63**, 617–628 (2024).
- 26 41. Backman, J. D. *et al.* Exome sequencing and analysis of 454,787 UK Biobank participants. *Nature* **599**,
27 628–634 (2021).

- 1 42. Yandell, M. *et al.* A probabilistic disease-gene finder for personal genomes. *Genome Res* **21**, 1529–1542
2 (2011).
- 3 43. Hu, H. *et al.* VAAST 2.0: Improved Variant Classification and Disease-Gene Identification Using a
4 Conservation-Controlled Amino Acid Substitution Matrix. *Genet Epidemiol* **37**, 622–634 (2013).
- 5 44. Wang, X. & Goldstein, D. B. Enhancer Domains Predict Gene Pathogenicity and Inform Gene Discovery in
6 Complex Disease. *Am J Hum Genet* **106**, 215–233 (2020).
- 7 45. McKay, J. D. *et al.* Lung cancer susceptibility locus at 5p15.33. *Nat Genet* **40**, 1404–1406 (2008).
- 8 46. Rafnar, T. *et al.* Sequence variants at the TERT-CLPTM1L locus associate with many cancer types. *Nat*
9 *Genet* **41**, 221–227 (2009).
- 10 47. Shete, S. *et al.* A Genome-Wide Association Study Identifies Two Novel Susceptible Regions for
11 Squamous Cell Carcinoma of the Head and Neck. *Cancer Res* **80**, 2451–2460 (2020).
- 12 48. Lesueur, C. *et al.* Genome-wide association analyses identify new susceptibility loci for oral cavity and
13 pharyngeal cancer. *Nat Genet* **48**, 1544–1550 (2016).
- 14 49. Petersen, G. M. *et al.* A genome-wide association study identifies pancreatic cancer susceptibility loci on
15 chromosomes 13q22.1, 1q32.1 and 5p15.33. *Nat Genet* **42**, 224–228 (2010).
- 16 50. James, M. A., Vikis, H. G., Tate, E., Rymaszewski, A. L. & You, M. CRR9/CLPTM1L regulates cell survival
17 signaling and is required for Ras transformation and lung tumorigenesis. *Cancer Res* **74**, 1116–1127
18 (2014).
- 19 51. Parashar, D. *et al.* Targeted biologic inhibition of both tumor cell-intrinsic and intercellular CLPTM1L/CRR9-
20 mediated chemotherapeutic drug resistance. *npj Precis. Onc.* **5**, 1–13 (2021).
- 21 52. Kim, N. W. *et al.* Specific association of human telomerase activity with immortal cells and cancer. *Science*
22 **266**, 2011–2015 (1994).
- 23 53. Ge, M. *et al.* Functional evaluation of TERT-CLPTM1L genetic variants associated with susceptibility of
24 papillary thyroid carcinoma. *Sci Rep* **6**, 26037 (2016).
- 25 54. Gholami, M. *et al.* Association of miRNA targetome variants in LAMC1 and GNB3 genes with colorectal
26 cancer and obesity. *Cancer Med* **11**, 3923–3938 (2022).

- 1 55. Hazelwood, E. *et al.* Integrating multi-tissue expression and splicing data to prioritise anatomical subsite-
2 and sex-specific colorectal cancer susceptibility genes with therapeutic potential. 2024.09.10.24313450
3 Preprint at <https://doi.org/10.1101/2024.09.10.24313450> (2024).
- 4 56. Qin, Y., Rodin, S., Simonson, O. E. & Hollande, F. Laminins and cancer stem cells: Partners in crime?
5 *Semin Cancer Biol* **45**, 3–12 (2017).
- 6 57. Galatenko, V. V., Maltseva, D. V., Galatenko, A. V., Rodin, S. & Tonevitsky, A. G. Cumulative prognostic
7 power of laminin genes in colorectal cancer. *BMC Med Genomics* **11**, 9 (2018).
- 8 58. Bogdanova, N., Helbig, S. & Dörk, T. Hereditary breast cancer: ever more pieces to the polygenic puzzle.
9 *Hereditary Cancer in Clinical Practice* **11**, 12 (2013).
- 10 59. Antoniou, A. C. *et al.* A locus on 19p13 modifies risk of breast cancer in BRCA1 mutation carriers and is
11 associated with hormone receptor-negative breast cancer in the general population. *Nat Genet* **42**, 885–
12 892 (2010).
- 13 60. Bolton, K. L. *et al.* Common variants at 19p13 are associated with susceptibility to ovarian cancer. *Nat*
14 *Genet* **42**, 880–884 (2010).
- 15 61. Mancuso, N. *et al.* Probabilistic fine-mapping of transcriptome-wide association studies. *Nat Genet* **51**,
16 675–682 (2019).
- 17 62. Bhattacharya, A. *et al.* Best practices for multi-ancestry, meta-analytic transcriptome-wide association
18 studies: Lessons from the Global Biobank Meta-analysis Initiative. *Cell Genomics* **2**, 100180 (2022).
- 19 63. Kachuri, L. *et al.* Gene expression in African Americans, Puerto Ricans and Mexican Americans reveals
20 ancestry-specific patterns of genetic architecture. *Nat Genet* **55**, 952–963 (2023).
- 21 64. Bhattacharya, A. *et al.* A framework for transcriptome-wide association studies in breast cancer in diverse
22 study populations. *Genome Biol* **21**, 42 (2020).
- 23 65. Keys, K. L. *et al.* On the cross-population generalizability of gene expression prediction models. *PLoS*
24 *Genet* **16**, e1008927 (2020).
- 25 66. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-aware
26 quantification of transcript expression. *Nat Methods* **14**, 417–419 (2017).
- 27 67. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res*
28 **47**, D766–D773 (2019).

- 1 68. Love, M. I. *et al.* Tximeta: Reference sequence checksums for provenance identification in RNA-seq.
2 *PLOS Computational Biology* **16**, e1007664 (2020).
- 3 69. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and microarray
4 studies. *Nucleic Acids Res* **43**, e47 (2015).
- 5 70. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data
6 with DESeq2. *Genome Biol* **15**, 550 (2014).
- 7 71. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- 8 72. Lindström, S. *et al.* Genome-wide analyses characterize shared heritability among cancers and identify
9 novel cancer susceptibility regions. *JNCI: Journal of the National Cancer Institute* **115**, 712–732 (2023).
- 10 73. Bhattacharya, A. isoTwas models for 48 GTEx models (04/2024). Zenodo
11 <https://doi.org/10.5281/zenodo.11048201> (2024).
- 12 74. Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear Models via Coordinate
13 Descent. *J Stat Softw* **33**, 1–22 (2010).
- 14 75. Endelman, J. B. Ridge Regression and Other Kernels for Genomic Selection with R Package rrBLUP. *The*
15 *Plant Genome* **4**, (2011).
- 16 76. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A Simple New Approach to Variable Selection in
17 Regression, with Application to Genetic Fine Mapping. *Journal of the Royal Statistical Society Series B:*
18 *Statistical Methodology* **82**, 1273–1300 (2020).
- 19 77. Lachmann, A. *et al.* ChEA: transcription factor regulation inferred from integrating genome-wide ChIP-X
20 experiments. *Bioinformatics* **26**, 2438 (2010).

FIGURES

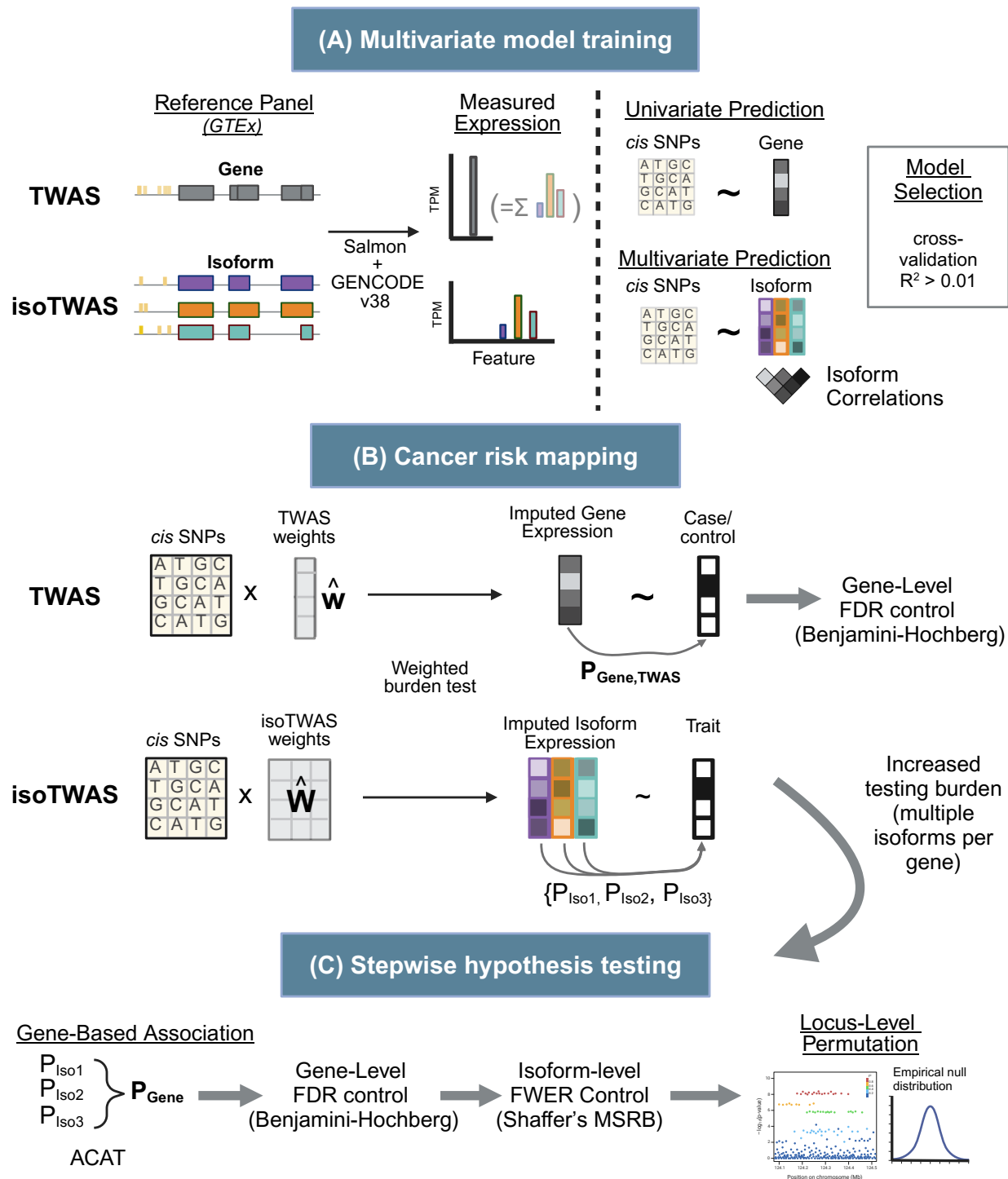


Figure 1: Overview of isoform- and gene-level TWAS. (A) Using multi-tissue GTEx data, isoform- and gene-level expression was quantified using Salmon and GENCODE v38 annotations. Univariate and multivariate predictive models for gene and isoform expression were trained. (B) Using the weighted burden test, predictive models of gene and isoform expression were used to map gene- and isoform-level cancer risk associations. For gene-level TWAS associations, FDR is controlled to 5% via Benjamini-Hochberg. (C) As isoTWAS presents an increased testing burden, stepwise hypothesis testing is conducted. Isoform P-values are aggregated to the gene-level via aggregated Cauchy aggregation. FDR is controlled at 5% via Benjamini-Hochberg procedure on gene-level P-values. For significantly-associated genes, FWER is controlled to 5% for correlated isoform associations using Shaffer's MSRB procedure. Locus-level permutation is conducted to control of local LD patterns. Modified from Bhattacharya et al 2023.

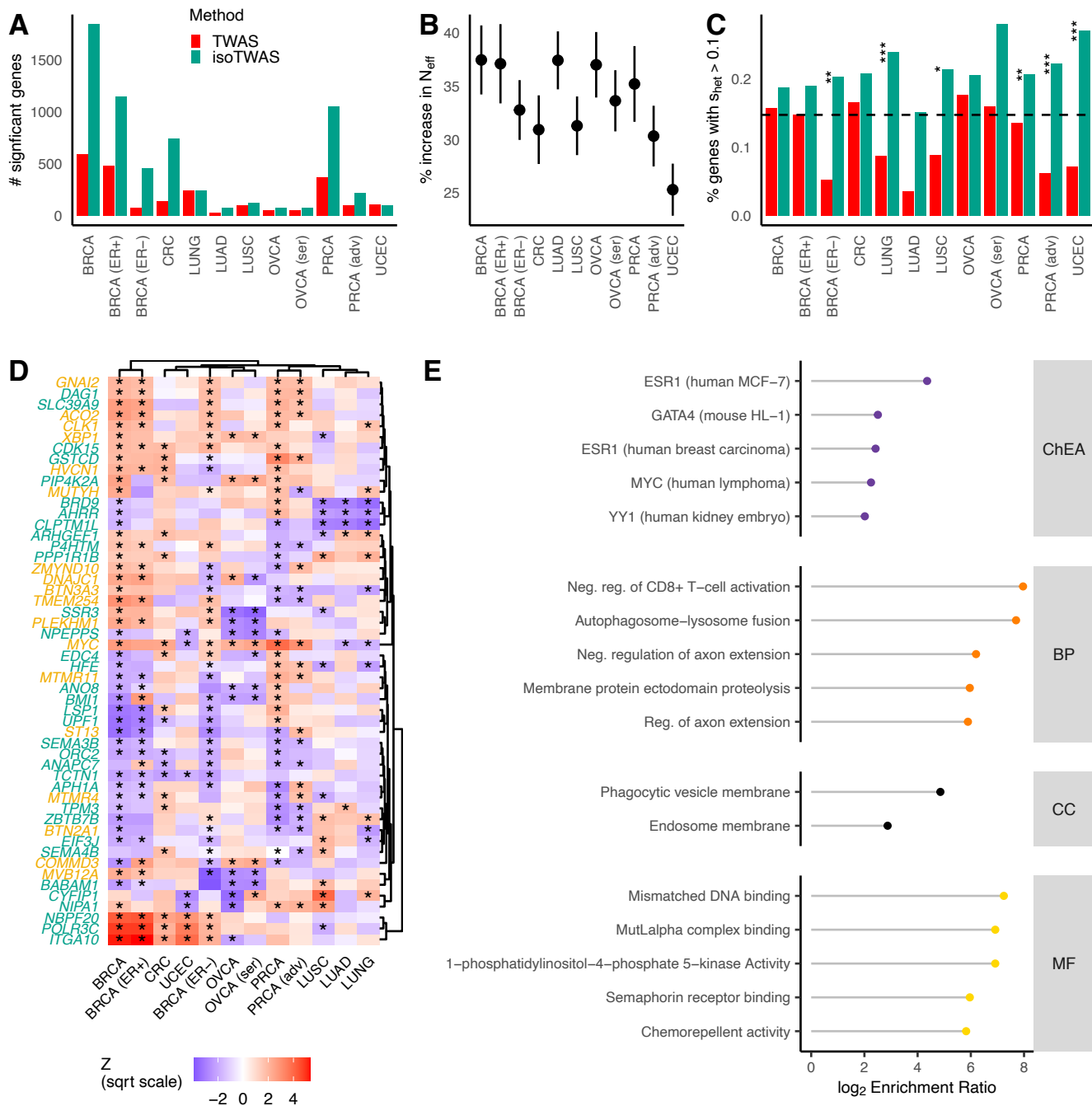


Figure 2: Isoform-level analysis identifies substantially more gene associations with cancer risk across 12 outcomes. (A) The number of unique transcriptome-wide significant genes identified with TWAS (red) and isoTWAS (green). **(B)** Percent increase in effective sample size with Wald-type 95% confidence interval using jackknife standard errors. **(C)** Proportion of transcriptome-wide significant genes identified with $s_{\text{het}} > 0.1$. Asterisks indicate FDR-adjusted P-value of χ^2 test for enrichment ratio of high pLI genes among all transcriptome-wide significant genes across method. (*) indicates FDR-adjusted $P < 0.05$, (**) $P < 0.01$, (***) $P < 0.005$. Black line shows the genome-wide proportion of genes with $s_{\text{het}} > 0.1$. **(D)** Scaled gene-level isoTWAS Z-score for 52 genes with isoform-level risk associations with at least 5 cancer outcomes across 12 outcomes. Genes are marked in green if no gene-level risk associations with any cancer outcome, and asterisk is shown if the gene association is significant for the given cancer outcome. **(E)** Log-enrichment ratio (X-axis) of over-represented gene ontologies (Y-axis), across ChIP-seq identified transcription factor targets (ChEA), biological process (BP), cellular component (CC), or molecular function (MF) pathways for 34 genes with isoform-level risk associations with at least 5 cancer outcomes and no TWAS associations.

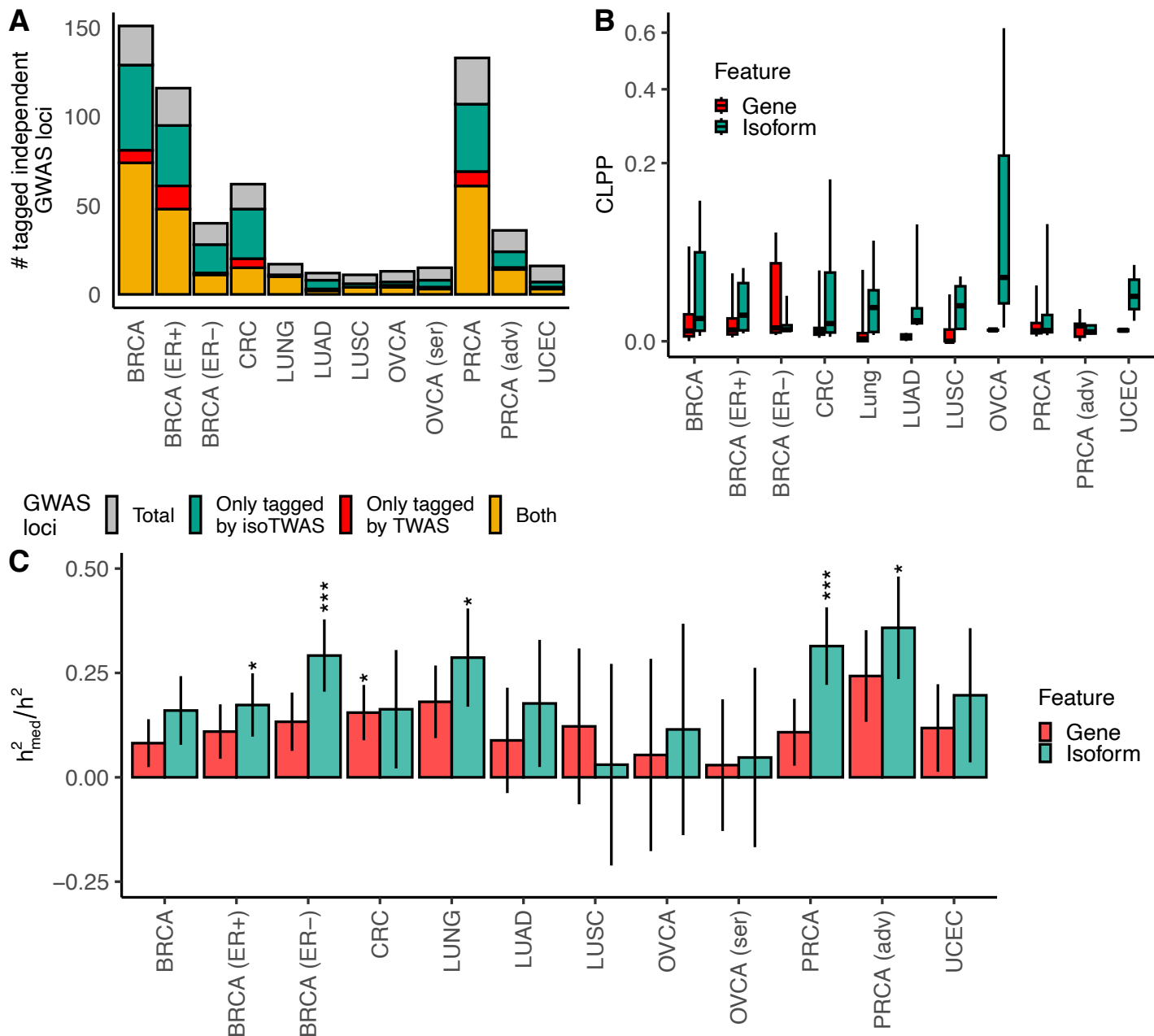


Figure 3: Isoform expression potentially mediates far more GWAS signal than gene expression. (A) The number of independent GWAS loci tagged by TWAS (red) or isoTWAS (green) or both in common (gold). **(B)** Boxplots of CoLocalization Posterior Probability (CLPP, Y-axis, square-root axis) of GWAS and gene expression QTL (red) and isoform expression QTL (green) for genomic loci with a GWAS $P < 5 \times 10^{-8}$ and QTL $P < 10^{-6}$. **(C)** Ratio of gene- (red) and isoform-level (green) expression mediated heritability (h^2_{med}) and total SNP heritability (h^2) with standard errors. Asterisks indicate FDR-adjusted P-value of Wald-type Z tests of $h^2_{med}/h^2 = 0$. (*) indicates FDR-adjusted $P < 0.05$, (**) $P < 0.01$, (***) $P < 0.005$.

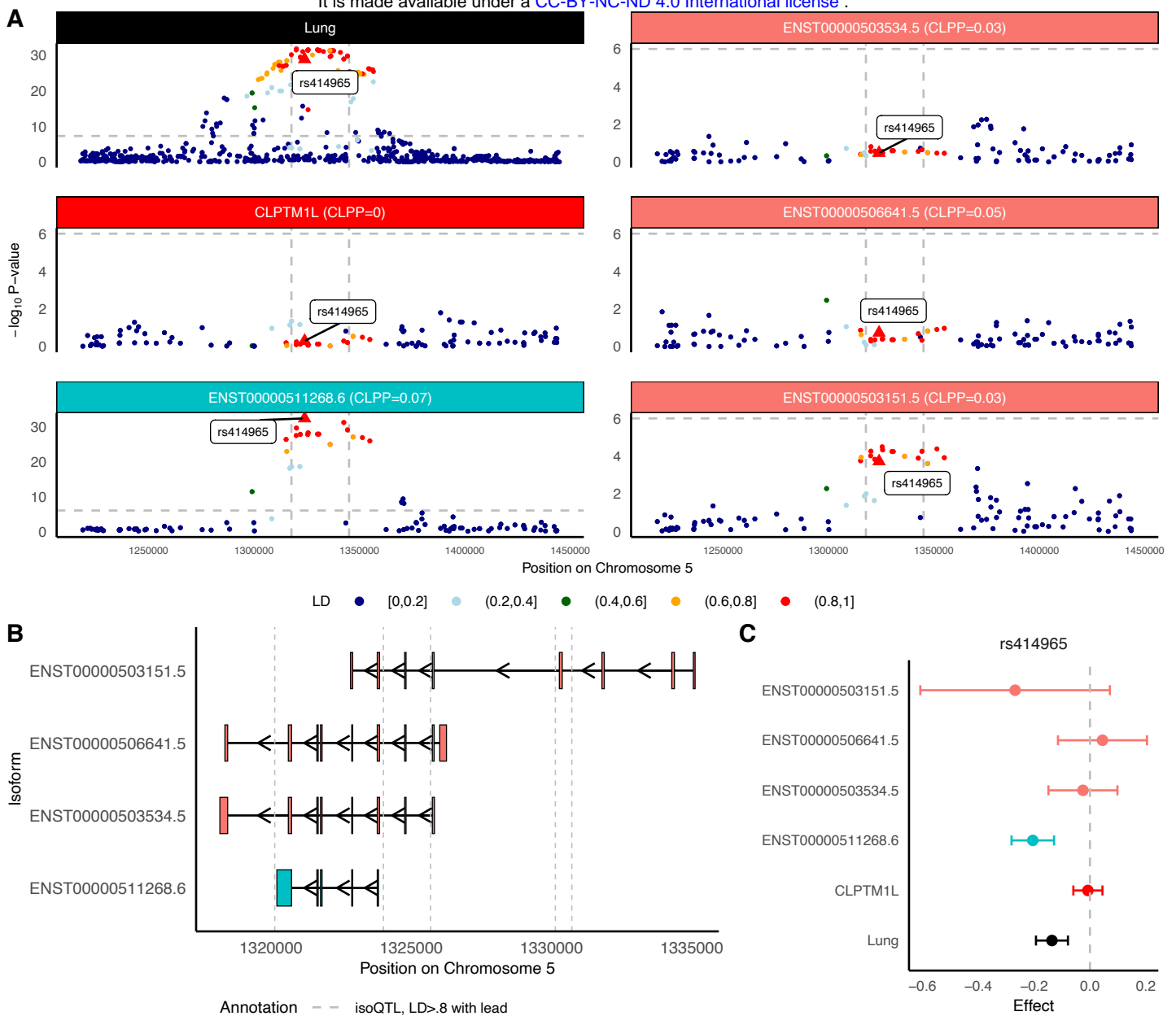


Figure 4: *CLPTM1L* isoforms may mediate lung cancer risk GWAS locus at Chromosome 5p15.33. (A) Manhattan plot of GWAS effects, *CLPTM1L* gene-eQTLs, and isoform-eQTLs for all significantly associated isoforms of *CLPTM1L*, colored by LD to rs414965, the lead isoQTL for ENST00000511268.6. **(B)** Transcript structure of significantly associated isoforms of *CLPTM1L*. Vertical lines indicated significant isoQTLs of LD > 0.8 to rs414965, strongest isoQTL of *CLPTM1L* isoforms. **(C)** SNP effect sizes on lung cancer risk (black), *CLPTM1L* gene expression (red), ENST00000511268.6 isoform expression (blue), and other expression of other isoforms (peach) for rs414965.

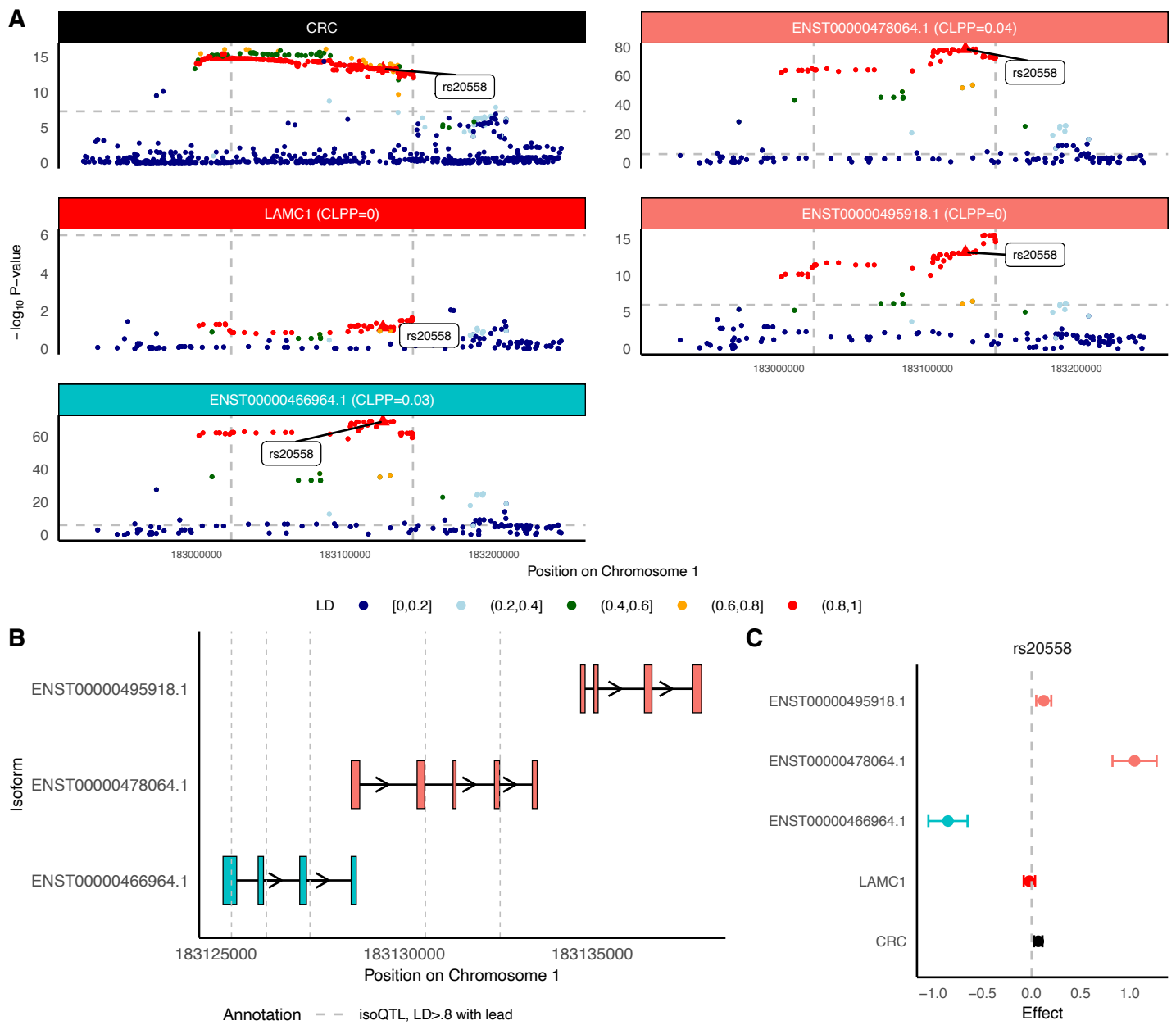


Figure 5: *LAMC1* isoforms may mediate colorectal cancer GWAS locus at Chromosome 1q25.3. (A) Manhattan plot of GWAS effects, *LAMC1* gene-eQTLs, and isoform-eQTLs for all significantly associated isoforms of *LAMC1*, colored by LD to rs20558, the lead isoQTL for ENST00000466964.1. **(B)** Transcript structure of significantly associated isoforms of *LAMC1*. Vertical lines indicated significant isoQTLs of LD > 0.8 to rs20558, strongest isoQTL of *LAMC1* isoforms. **(C)** SNP effect sizes on colorectal cancer risk (black), *LAMC1* gene expression (red), ENST00000466964.1 isoform expression (blue), and other expression of other isoforms (peach) for rs20558.

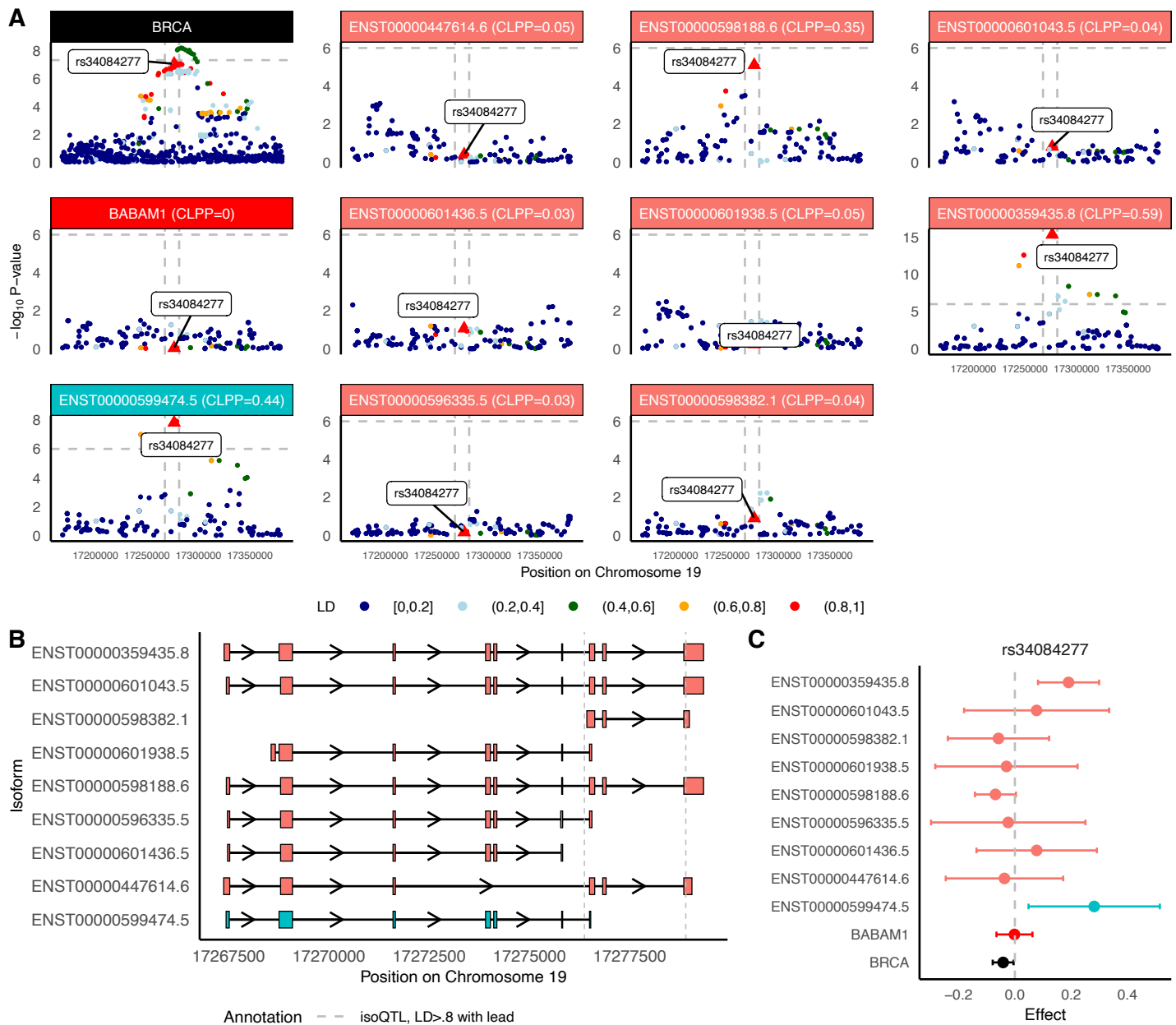


Figure 6: *BABAM1* isoforms may mediate breast cancer risk GWAS locus at Chromosome 19p13.11. **(A)** Manhattan plot of GWAS effects, *BABAM1* gene-eQTLs, and isoform-eQTLs for isoforms of *BABAM1*, either prioritized through isoTWAS or with an isoQTL with $P < 1e-6$. **(B)** Transcript structure of *BABAM1*. Vertical lines indicated significant isoQTLs of LD > 0.8 to rs34084277, strongest isoQTL of *BABAM1* isoforms. **(C)** SNP effect sizes on breast cancer risk (black), *BABAM1* gene expression (red), expression of isoTWAS-prioritized isoforms (blue), and expression of other isoforms (peach) for rs34084277 (lead isoQTL).