

EPITHYR

EPITHYR is a collaborative consortium¹ comprising seven population-based case-control studies focused on thyroid cancer research. These studies include three conducted in Metropolitan France (CATHY, YOUNG-Thyr, and E3N), two in the South Pacific Islands (French Polynesia and New Caledonia), one in Cuba, and one in the Gomel region of Belarus, which was contaminated following the Chernobyl disaster. The CATHY study was carried out across three French administrative regions, involving patients aged 25 years and older diagnosed with DTC between 2002 and 2007, with DNA samples collected from 516 cases and 569 controls. The New Caledonia study encompassed the entire country, including DTC patients diagnosed between 1993 and 1999, with available DNA samples from 224 cases and 261 controls. The YOUNG-Thyr study focused on individuals under 35 years old, born after January 1, 1971, and diagnosed with DTC between 2002 and 2006, with DNA samples collected from 715 cases and 692 controls. The Cuban study included patients aged 17-60 years, residing in Havana, who were diagnosed and treated for DTC between 2000 and 2011, consisting of 199 cases and 193 controls. The French Polynesia study involved patients diagnosed between 1983 and 2003, younger than 56 years, and born and living in French Polynesia, including 144 cases and 231 controls. The Chernobyl study examined individuals from the most contaminated areas of Belarus, including 83 cases of papillary thyroid cancer (PTC) and 324 matched controls who were under 15 years old at the time of the Chernobyl accident. The E3N study was a prospective cohort study comprising 98,995 women in France born between 1925 and 1950, with sufficient DNA material available for 568 participants. Genotyping in EPITHYR was conducted using the Infinium OncoArray-500K BeadChip (Illumina), with quality control performed on 4,553 genotyped individuals, excluding those with a call rate below 95%, sex discordance, duplicates, and related individuals (π -hat threshold of 0.2²). We filtered by ethnicity to include only individuals of European ancestry, and ultimately, the study included 1,552 incident DTC cases and 1,954 controls. SNP quality control reduced the dataset to 401,525 SNPs by removing variants with poor call rates ($n=8,327$), duplicates ($n=814$), monomorphic SNPs ($n=13,832$), and those not meeting Hardy-Weinberg equilibrium ($P < 10^{-7}$ in controls, $P < 10^{-12}$ in cases) or minor allele frequency (<1%) thresholds. Imputation used the 1000 Genomes Project dataset, employing SHAPEIT2 and IMPUTEv2 for European populations, resulting in 9,665,802 SNPs. Logistic regression analysis, adjusted for sex, region, and the top 10 principal components, was used to assess DTC associations.

EPIC

The EPIC cohort³ is a large multicenter prospective study initiated in 1992 across 7 European countries (Denmark, France, Germany, Italy, Spain, Netherlands, and United Kingdom) to investigate links between nutrition, lifestyle, metabolism, genetic factors, and cancer risk. It includes approximately 370,000 women and 150,000 men. Thyroid cancer cases were identified among initially cancer-free participants who donated blood and were later diagnosed with TC during follow-up. DTC cases were classified using ICD-10 codes, with papillary cases corresponding to morphologic codes 8050, 8130, 8260, 8340-8344, and 8350, while follicular cases were assigned codes 8290, 8330-8335. Controls were selected from living, cancer-free cohort members at the time of diagnosis and matched based on recruitment center, sex, age, date, time, fasting status at blood collection, and follow-up duration. Genotyping and quality control in the EPIC study followed the same protocols as EPITHYR, leading to the analysis of 9,652,732 SNPs. Logistic regression adjusted for sex, country, and the top 10 principal components was employed to measure associations, with the study including 345 incident DTC cases and 783 controls.

deCODE genetics

deCODE genetics study⁴ examined genetic associations with thyroid cancer across European ancestry populations, including individuals from Iceland, the Netherlands, Spain, and the United States. All participants were of European descent. In Iceland, thyroid cancer cases were identified through the Icelandic Cancer Registry (ICD-10: C73), with genotypic data available for 1,003 individuals, predominantly women. The Dutch cohort comprised 85 non-medullary thyroid cancer cases and 4,956 cancer-free individuals, primarily women, with a mean diagnosis age of 40 years. The Columbus, USA cohort included 1,580 cases with a mean diagnosis age of 43 years and 1,628 controls with a mean age of 45 years. The Houston, USA cohort consisted of 250 PTC cases and 363 cancer-free controls, with mean ages at diagnosis of 45 and 53 years, respectively. The Spanish cohort included 83 non-medullary thyroid cancer cases and 1,612 controls, with a mean diagnosis age of 49 years. The Icelandic thyroid cancer GWAS dataset was generated using Illumina technology with an average sequencing depth of 34X, whereas genotyping for the Dutch, Spanish, and US samples was conducted using Omni-1 Quad-bead chips. Variants were excluded based on a yield below 94%, minor allele frequency under 1%, failure in Hardy-Weinberg equilibrium testing ($P < 1.0 \times 10^{-6}$), or significant genotype batch differences ($P < 1.0 \times 10^{-6}$), yielding a final dataset of 9,863,645 SNPs. Logistic regression, adjusted for study center, sex, age, and the top 10 principal components,

assessed genetic associations. Results across study groups were combined using a Mantel-Haenszel model⁵. The study included 3,001 incident DTC cases and 287,550 controls.

UK Biobank

The UK Biobank⁶ is a large population-based health research resource comprising approximately 500,000 participants aged 40-69 years, recruited between 2006 and 2010 across the UK. DTC cases were identified using ICD-10 (Data field: 40006, code C730-739) and ICD-9 (Data field: 40013, code 1930-1939) codes, covering histological subtypes such as papillary, follicular, and oxyphilic cell thyroid cancer (Data field: 40011). Only first primary cancer occurrences were considered. Controls were cancer-free (except for non-melanoma skin cancer) throughout follow-up. Genotyping was conducted using the UK Biobank Axiom Array, measuring approximately 850,000 variants, with over 90 million imputed using the Haplotype Reference Consortium and UK10K + 1000 Genomes reference panels. Only European participants were analyzed using PCA. Logistic regression, adjusted for sex, age at recruitment, and the top 10 principal components, assessed genetic associations. The study included 518 DTC cases (262 incident cases) and 358,640 controls.

Italian Study

The Italian study⁷ enrolled 701 consecutive DTC cases from the Department of Endocrinology at the University Hospital of Pisa. Controls (n=499) were healthy individuals recruited from blood donors and hospital workers during routine surveillance. Both cases and controls were required to be at least 18 years old, of Caucasian origin, and unrelated. Controls were excluded if they had any history of malignancy, chronic inflammatory disease, or benign thyroid disorders. Genotyping was performed using Illumina BeadChips, with genotype calling via Illumina GenomeStudio 2010. Samples with SNP call rates below 95% were excluded, and markers were filtered out if they had a genotype call rate under 95%, minor allele frequency below 5%, or violated Hardy-Weinberg equilibrium ($P < 5.0 \times 10^{-5}$). Additional study details were previously published alongside the GWAS⁷. A total of 632 cases and 430 controls were included.

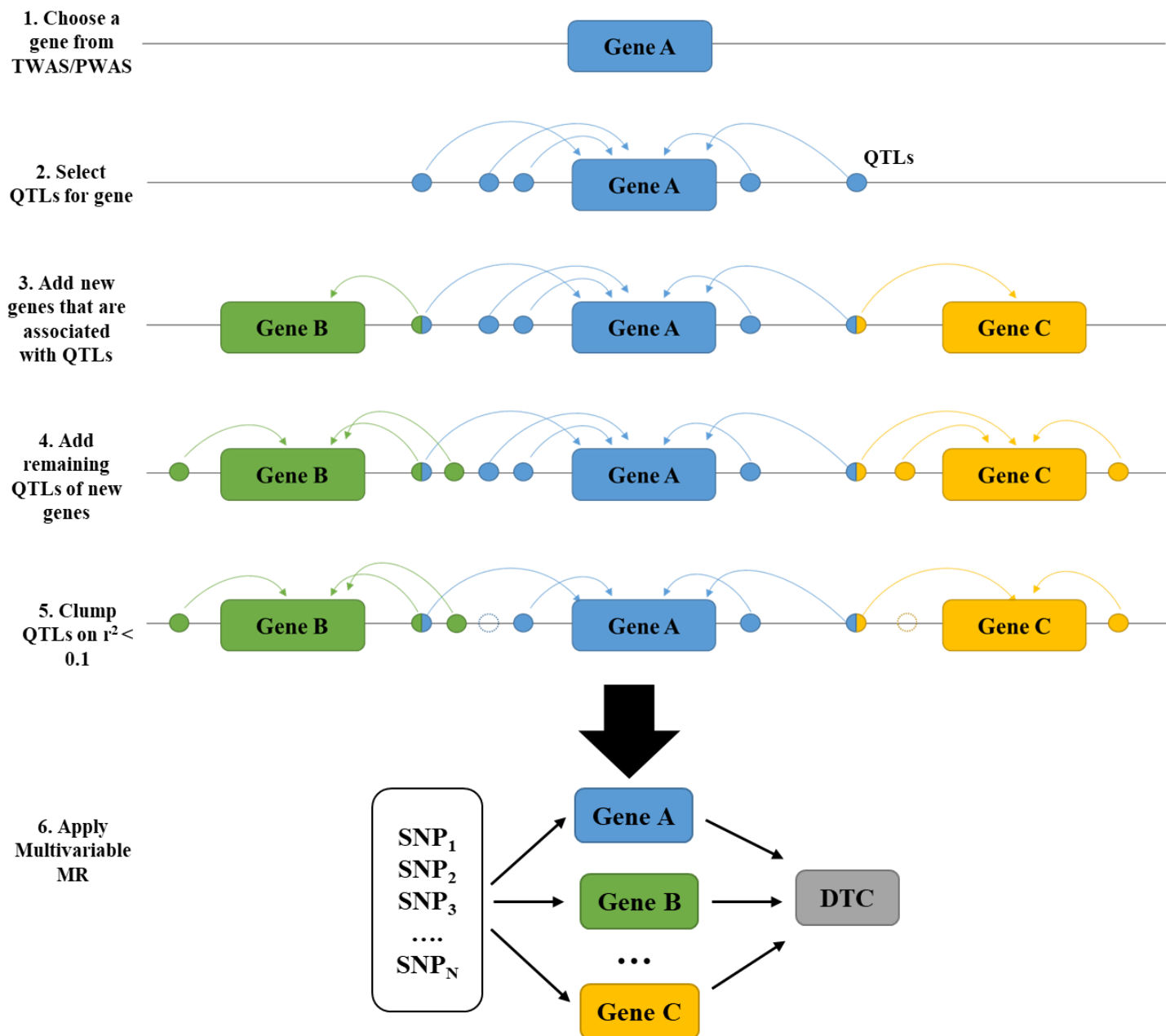
FinnGen

The FinnGen⁸ R10 dataset comprises 412,181 individuals of Finnish descent, sourced from Finnish biobanks and digital health registries. Study details are available at (https://www.finnngen.fi/en/access_results). Exclusion criteria included ambiguous gender, genotype missingness above 5%, extreme heterozygosity (± 4 SD), and non-Finnish ancestry. SNP genotyping was conducted using Illumina and Affymetrix arrays. The dataset included 1,471 papillary and 162 follicular

thyroid carcinoma cases, with cancer-free controls identified using 'EXALLC' variables that excluded other cancers. GWAS analysis applied GAIGE, a mixed model logistic regression, adjusting for sex, age, genotyping batch, and the top 10 principal components.

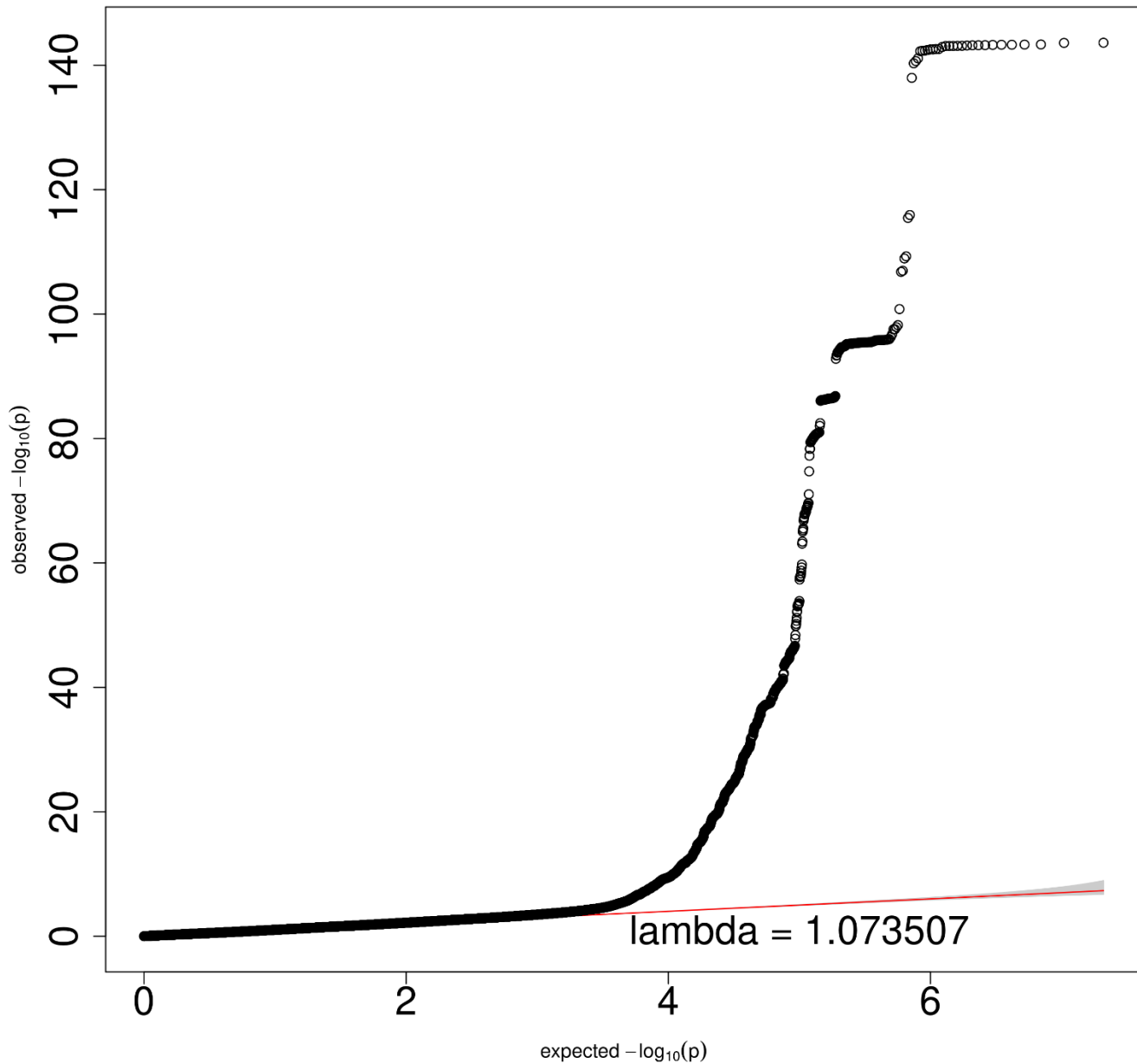
Supplementary Methods 2 Development of a polygenic risk score framework incorporating genome-, transcriptome-, and proteome-wide association studies for multi-omics risk prediction.

To construct the polygenic risk score, we utilized summary statistics from genome-wide association studies (GWAS), transcriptome-wide association studies (TWAS), and proteome-wide association studies (PWAS). The polygenic risk score (PRS) was constructed by first using significant single nucleotide polymorphisms (SNPs) identified in GWAS. Each SNP was weighted based on its effect size, and the weighted sum of risk alleles was computed for each individual. The formula for PRS can be expressed as $PRS_i = \sum_{j=1}^m \gamma_j x_{ij}$, where γ_j is the effect size of SNP j , x_{ij} is the risk allele count for SNP j in individual i , and m represents the total number of SNPs included in the PRS model. For the polygenic TWAS-based risk score (PTRS), we applied the PrediXcan framework to estimate genetically regulated gene expression (GReX) for each individual. The PTRS was calculated by summing the predicted gene expression levels, weighted by the effect sizes of each gene, derived from the S-PrediXcan⁹ framework. Predicted gene expression levels were calculated as $GeneExp_j = \sum_{j=1}^m X_j \beta_j$, where $GeneExp_j$ is the predicted expression level of gene j , X_j is the count of effect alleles for variant j , and β_j is the variant's effect on gene expression derived from a regression model across m variants. And then PTRS were calculated as $PTRS = \sum_{j=1}^g GeneExp_j \times Z_j$, where Z_j represents the TWAS effect size of gene j , $GeneExp_j$ is the predicted expression for gene j , and g is the total number of genes included¹⁰. To integrate proteomic information, we incorporated PWAS results by considering genetically predicted protein abundances for genes identified in brain and plasma blood tissues. The polygenic PWAS-based risk score (PPRS) was constructed similarly to the PTRS approach, by weighting the predicted protein levels by the effect sizes from proteomic data. We applied the FUSION¹¹ framework to estimate genetically regulated protein abundance expression for each individual. The formula for PPRS is $PPRS = \sum_{j=1}^g ProteinExp_j \times Z_j$, where $ProteinExp_j$ represents the predicted protein abundance for gene j , and Z_j is its corresponding PWAS effect size across g genes. For all models, gene selection was refined using a p-value threshold corrected for multiple testing via Bonferroni correction. Genes overlapping with GWAS loci were excluded to avoid redundancy, ensuring that the final risk scores reflected distinct contributions from each omics layer. Finally, we evaluated the predictive performance of the PRS by using receiver operating characteristic (ROC) curves and area under the curve (AUC).

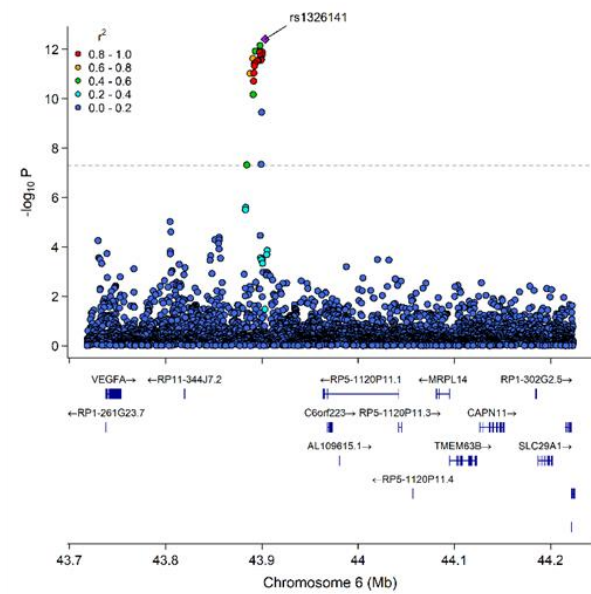
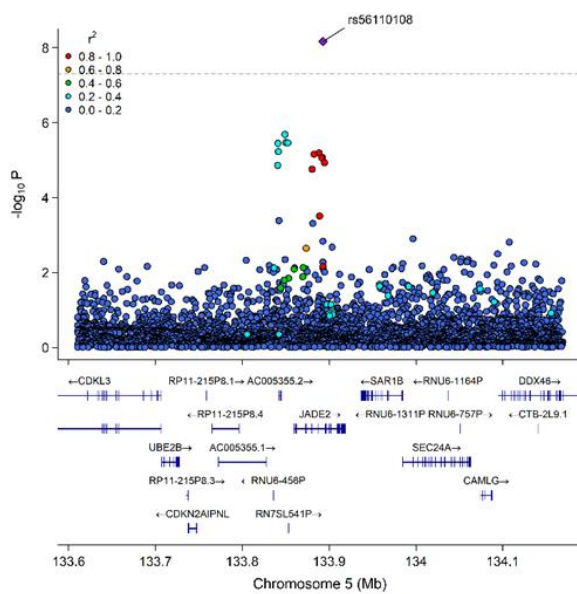
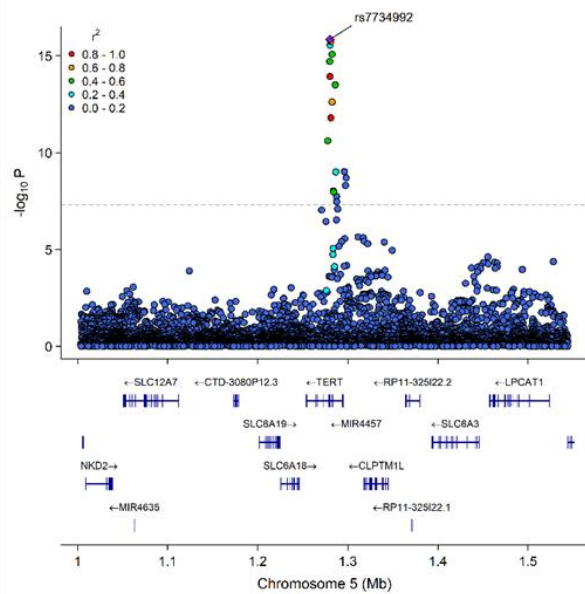
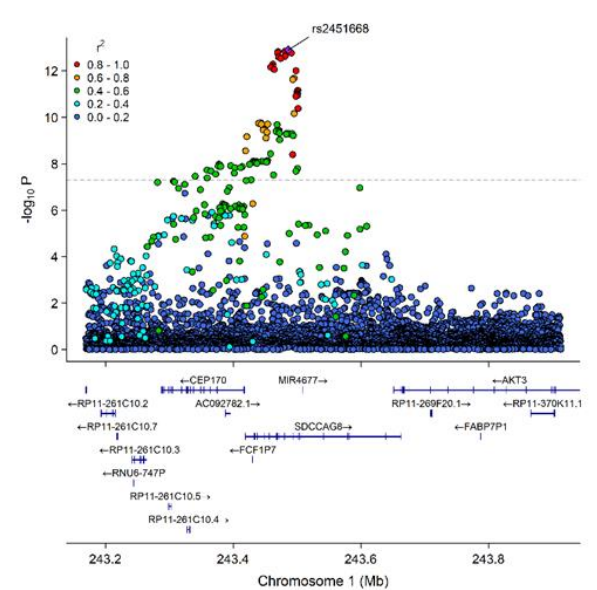
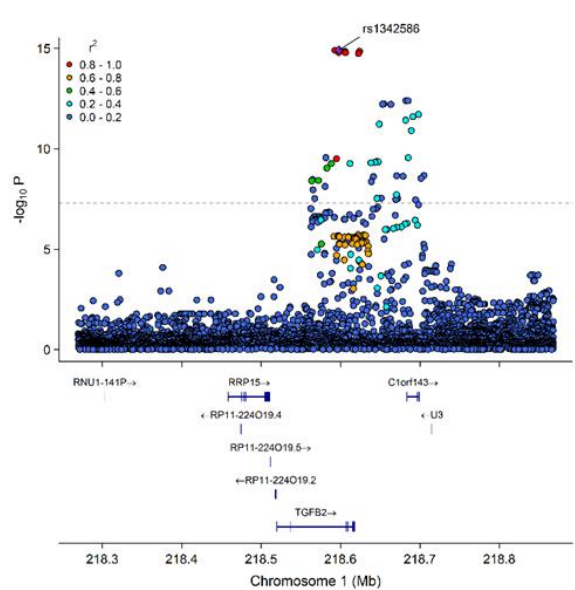
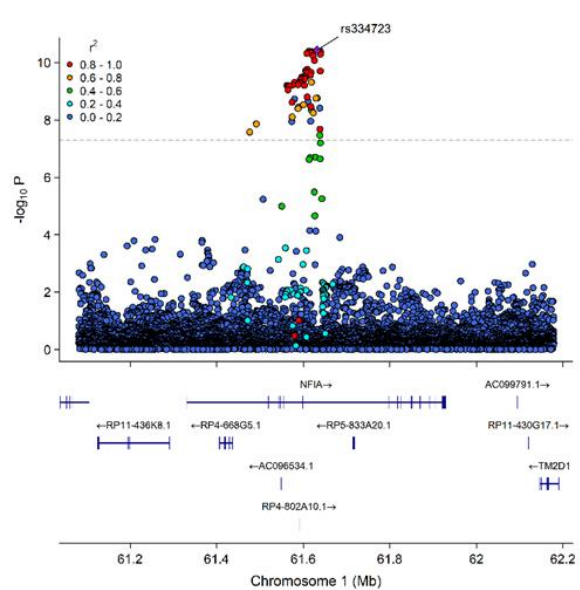


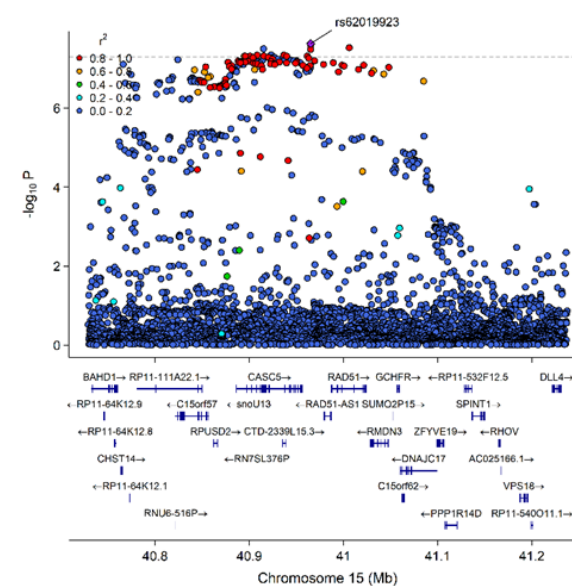
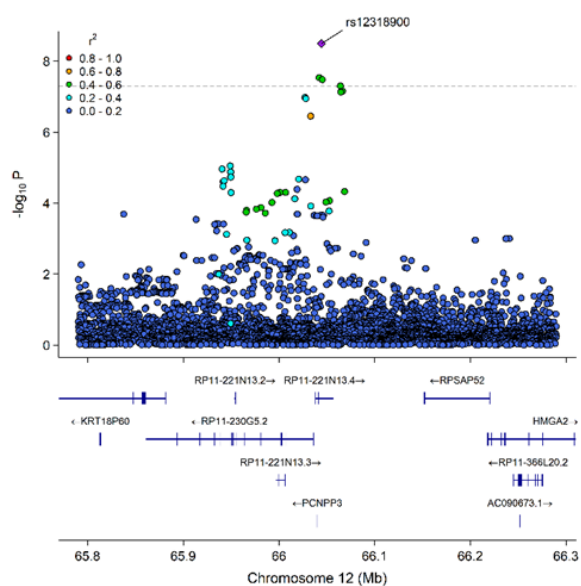
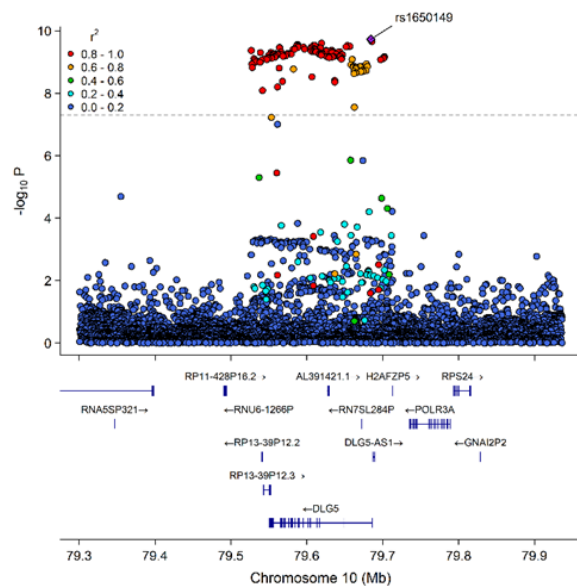
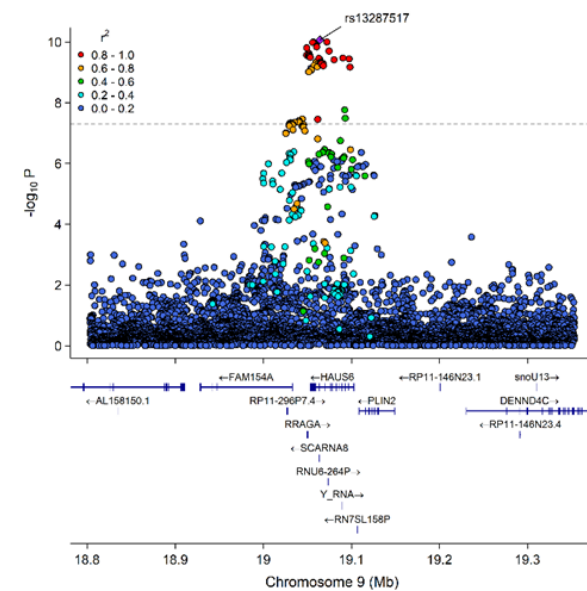
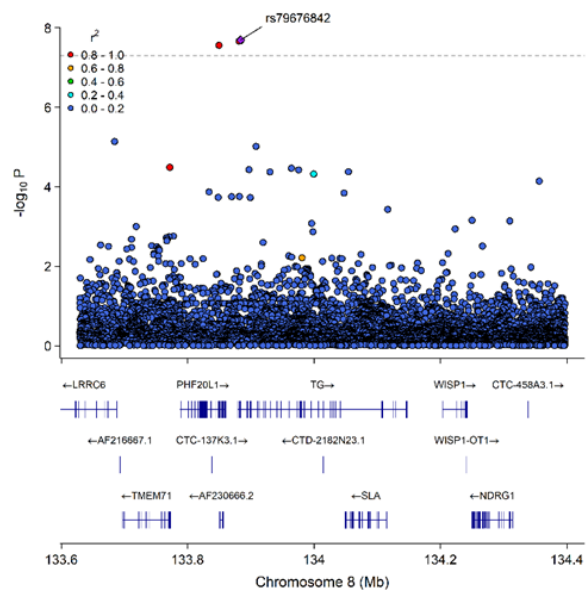
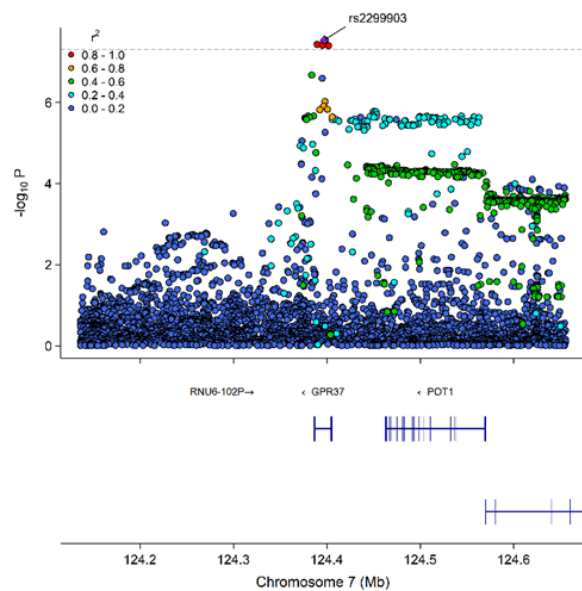
Supplementary Figure 1 A Mendelian Randomization framework utilizing genes as exposures.

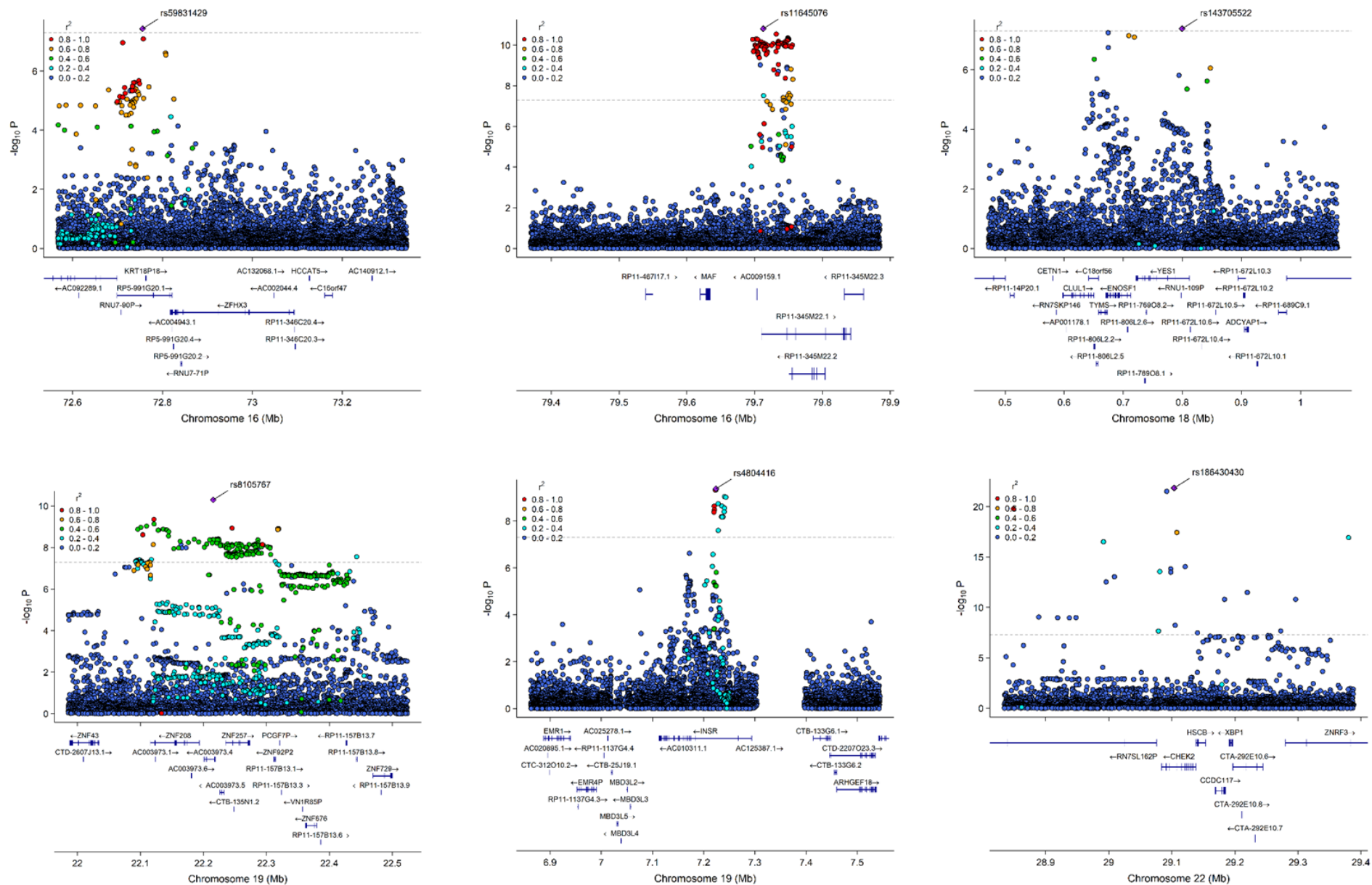
QQ plot of p-values



Supplementary Figure 2 Quantile-Quantile plots of DTC GWAS including 7,681 cases and 963,550 controls of European descent. The red diagonal line represents expected distribution under the assumption that there is no inflation of the chi-square statistics.

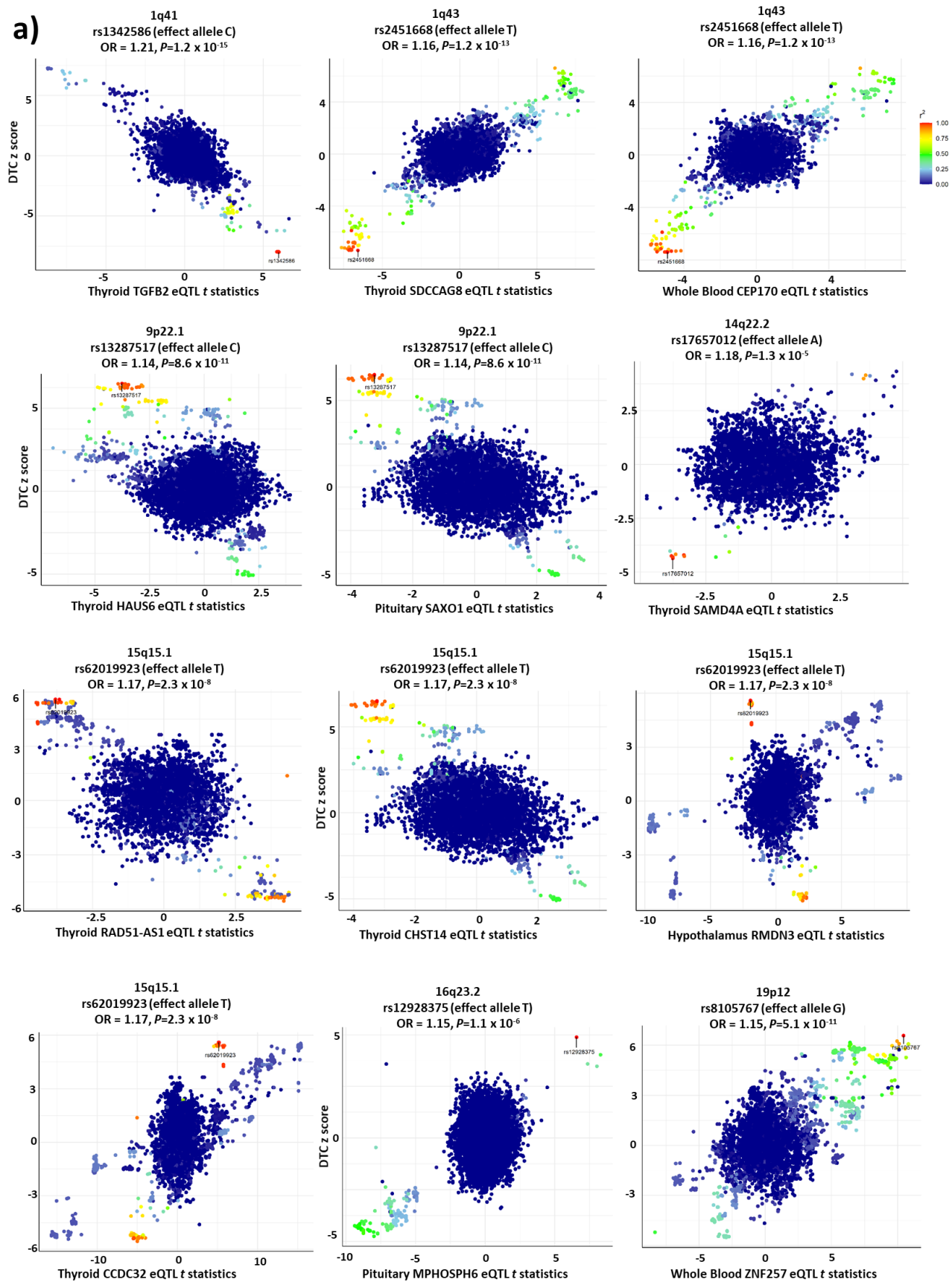




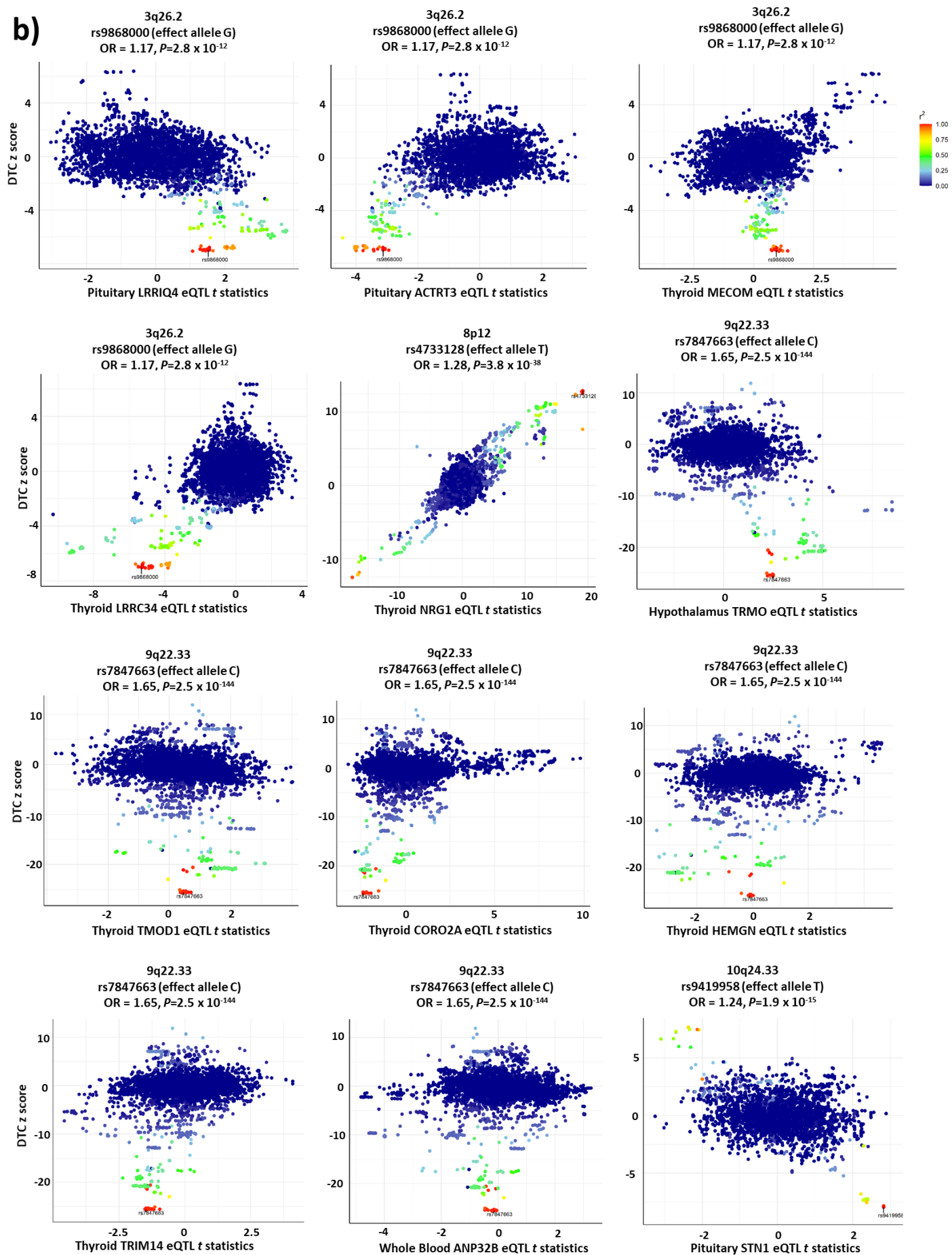


Supplementary Figure 3 Regional plots of the DTC GWAS are displayed for all significant novel loci at 1p31.3, 1q41, 1q43, 5p15.33, 5q31.1, 6p21.1, 7q31.33, 8q24.22, 9p22.1, 10q22.3, 12q14.3, 15q15.q, 16q22.2, 16q23.2, 18p11.32, 19p13.2, 19p12 and 22q12.1. The x-axis represents the position of SNPs and gene on the chromosome, while the y-axis represents the $-\log_{10}(P)$ -value. The color of each SNP dots indicates its correlation with LD.

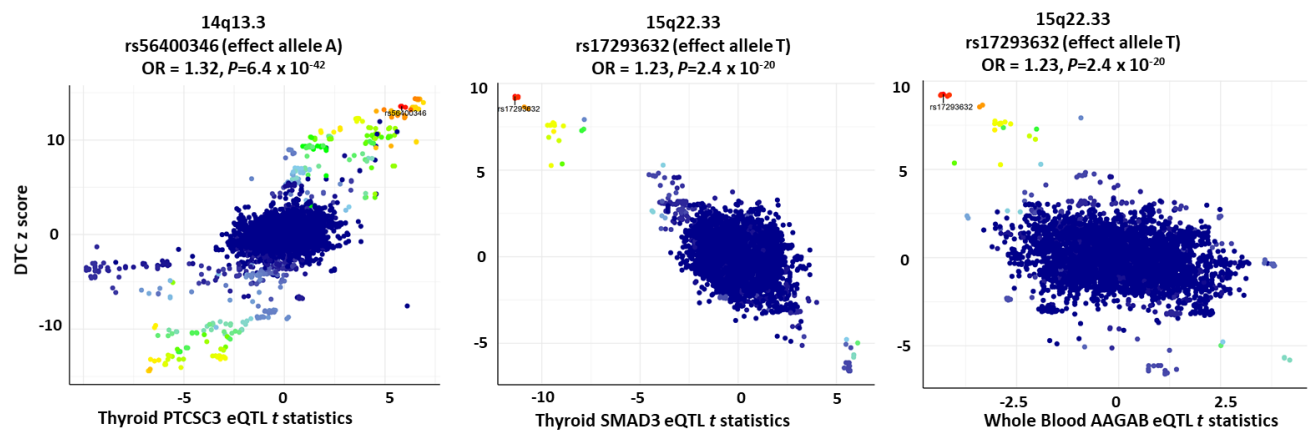
a)



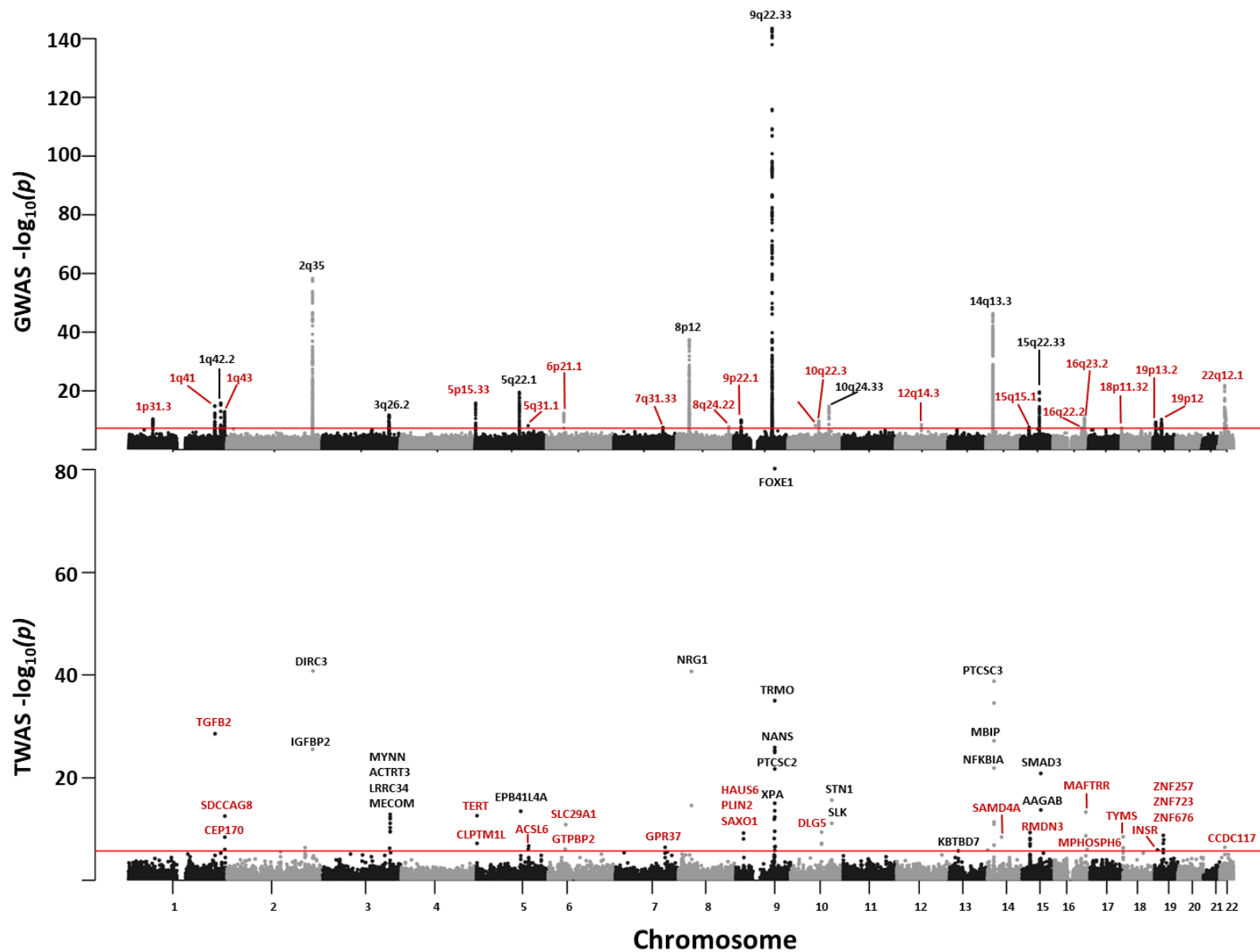
b)



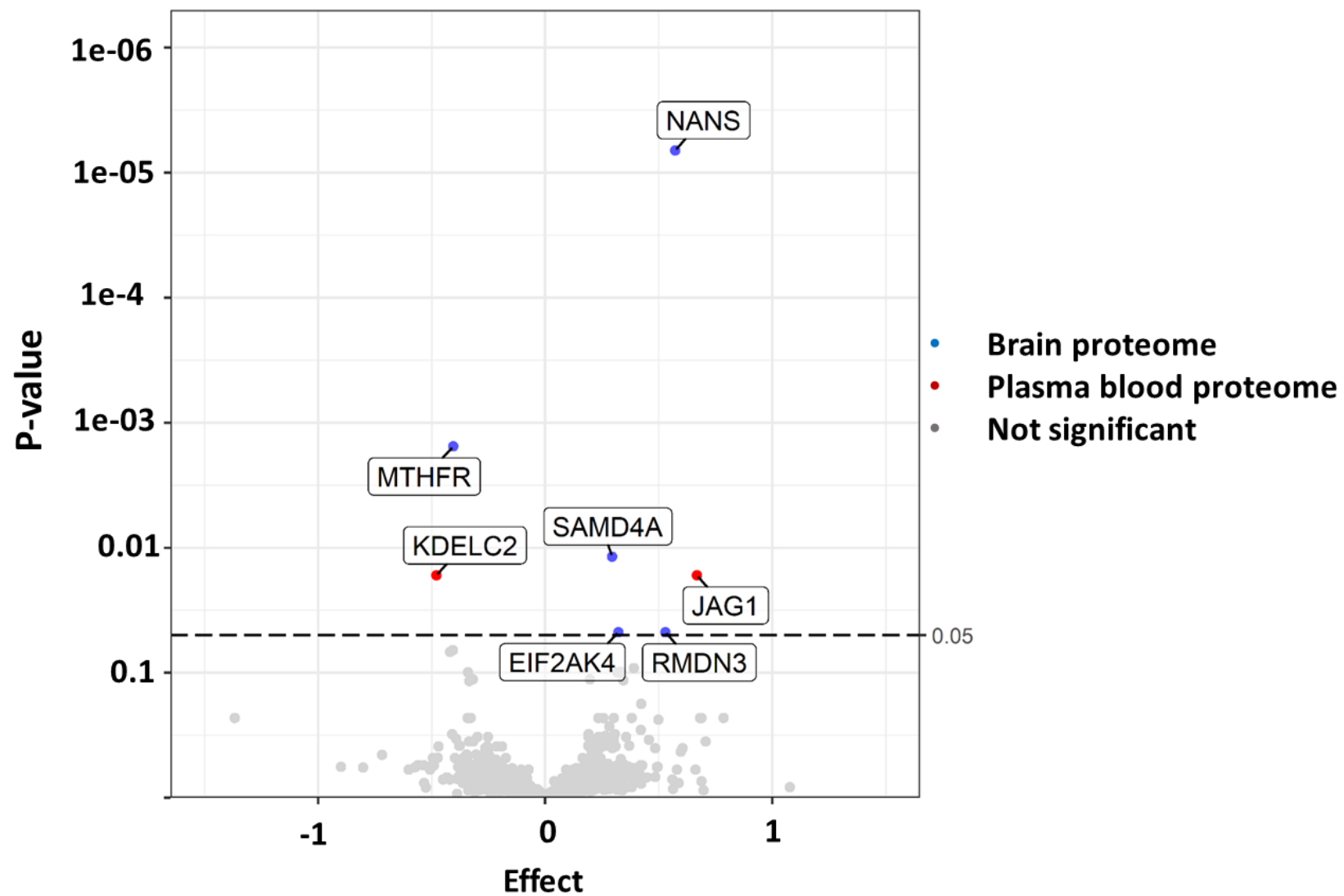
Cont.



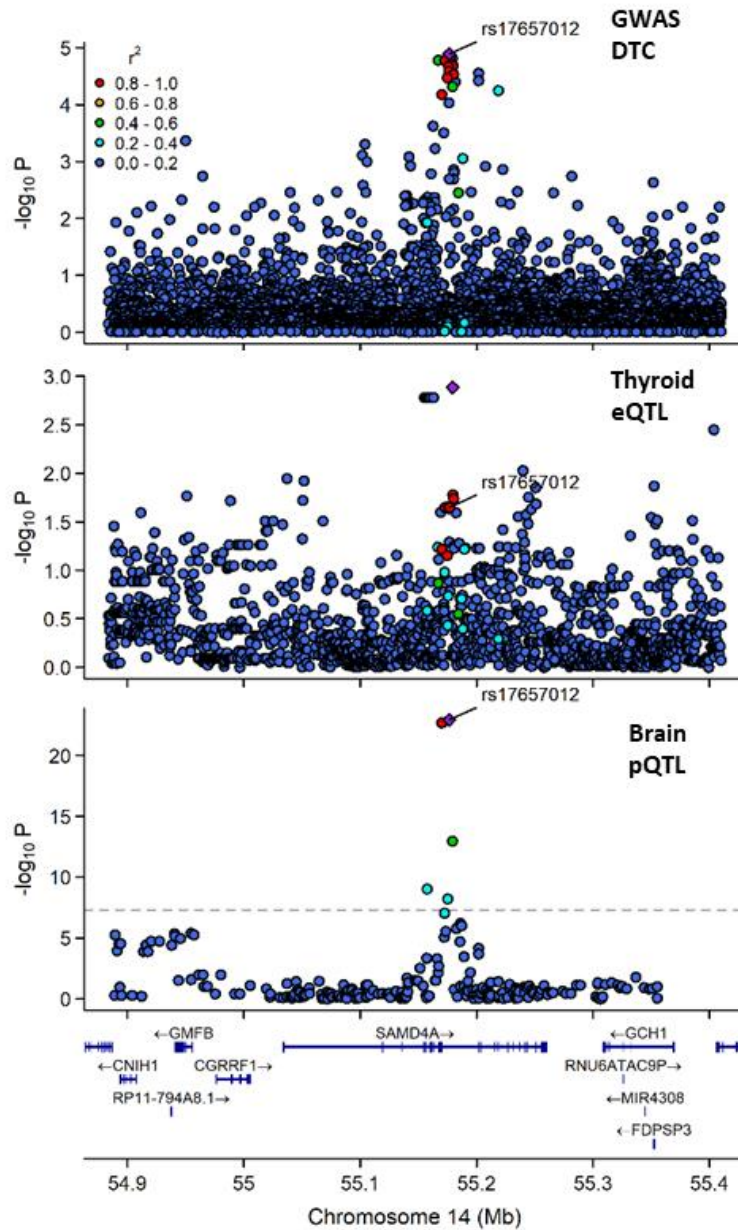
Supplementary Figure 4 Scatterplots of the a) **novel** and b) **known** susceptibility loci and identified by TWAS using thyroid-relevant tissues, showing their associations with DTC and corresponding cis-eQTLs in the tissue. The y-axis represents the association of variants with DTC, while the x-axis shows their association with gene expression levels. Each variant is color-coded according to its LD correlation.



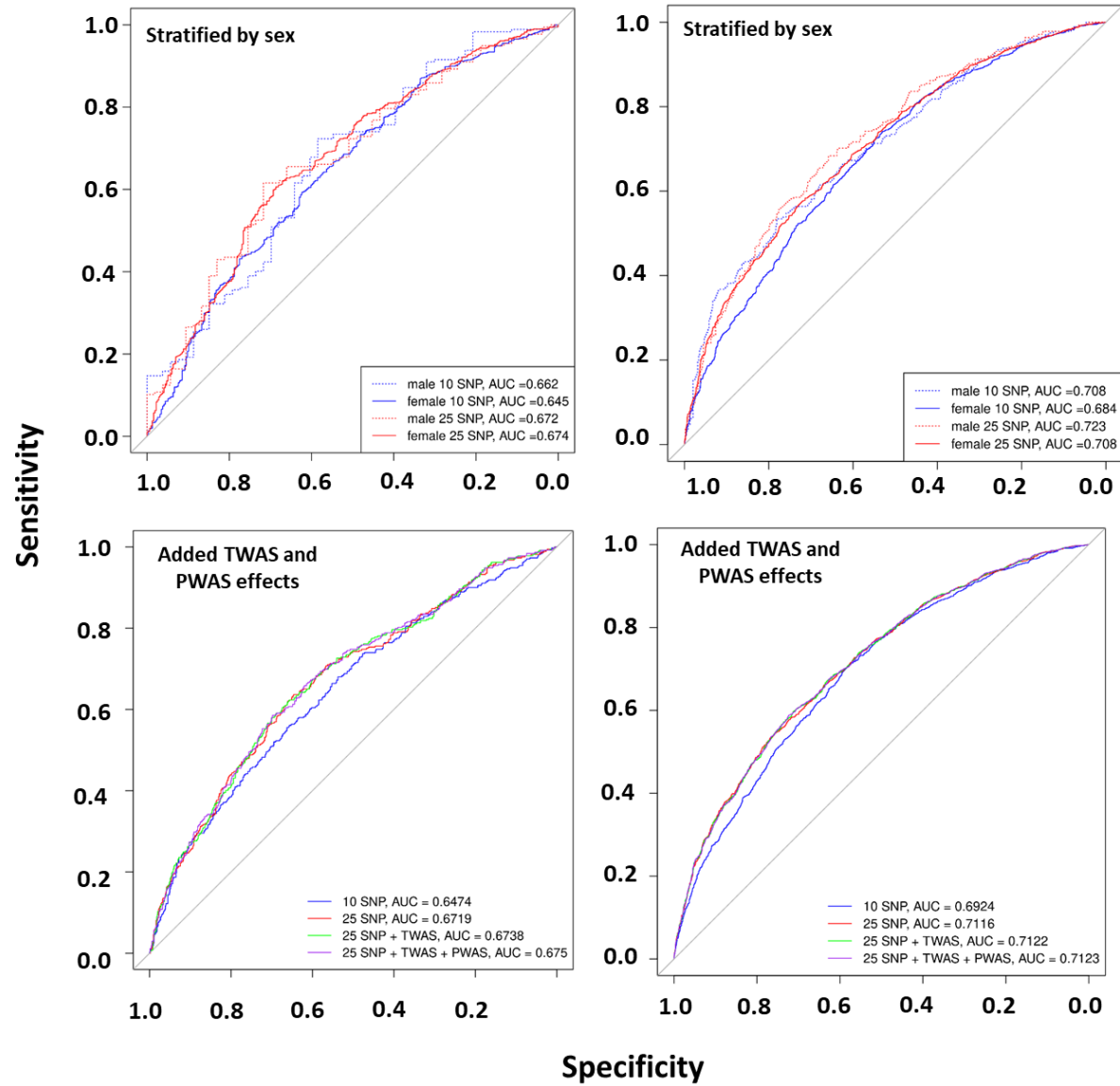
Supplementary Figure 5 The Manhattan plot displays the results of the DTC GWAS and TWAS analyses. A fixed-effect meta-analysis was performed on 20.8 million variants, including 7,681 cases and 963,550 controls of European descent. Previously identified risk loci are marked in black, while novel risk loci are highlighted in red, with the associated genes indicated in the TWAS Manhattan plot. Each dot in TWAS Manhattan plot corresponds to an association test between genetically predicted gene expression in all 49 tissues available from GTEx and differentiated thyroid cancer. The red line indicates the significant boundary for GWAS and TWAS after Bonferroni correction (5.0×10^{-8} and 2.45×10^{-6} respectively).



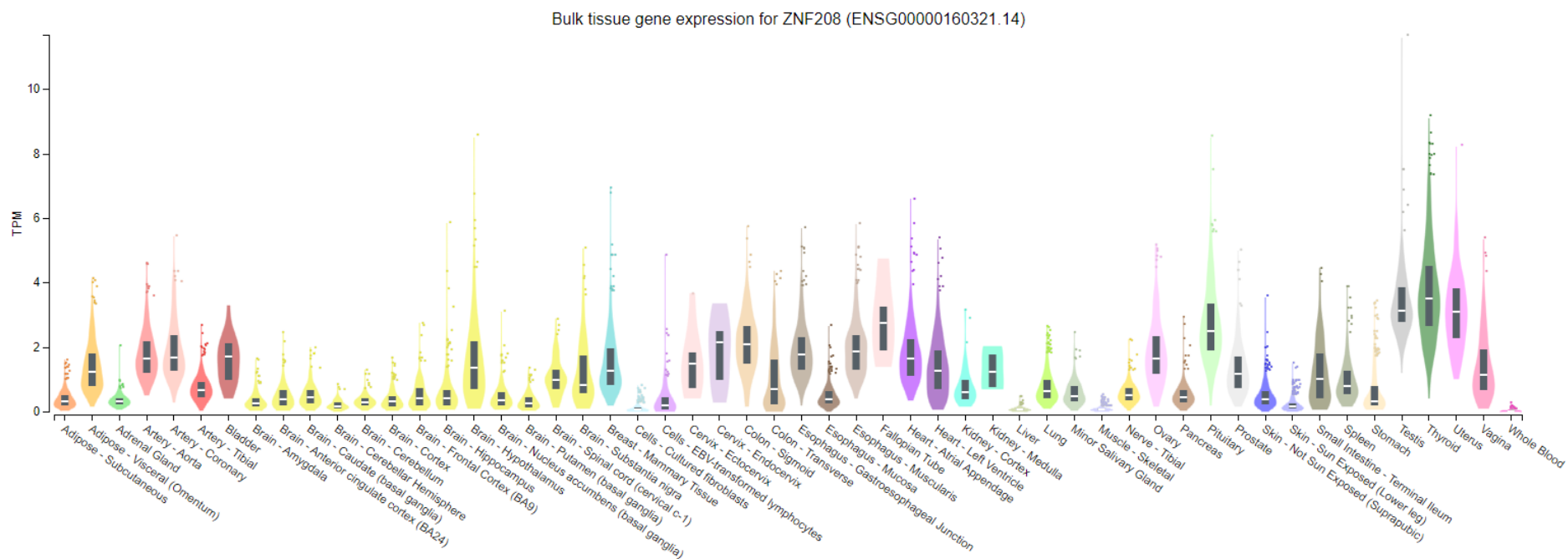
Supplementary Figure 6 Volcano plot showing results from proteome-wide association study (PWAS)



Supplementary Figure 7 A regional locus plots are shown from GWAS of differentiated thyroid cancer, gene expression (eQTL) and brain proteome levels (pQTL). *SAMD4A* genes were both identified in TWAS using thyroid-relevant tissue (thyroid tissue and brain hypothalamus respectively) as well as in brain proteome PWAS. *SAMD4A* genes displayed significant colocalization with high posterior probability (92% for eQTL and 99.1% for pQTL). SNPs are colored based on their LD with the GWAS lead variant (rs17657012).



Supplementary Figure 9 A receiver operating characteristic (ROC) curve evaluating the discriminative ability of PRS models based on 25 GWAS-identified significant SNPs compared to a previously reported model using 10 SNPs in the EPIC and EPITHYR study. Both 25 and 10 SNPs models were stratified by sex. The PRS models were further enhanced by incorporating TWAS and PWAS results to assess their effect on predictive performance.



Supplementary Figure 10 The expression levels of *ZNF208* at 19q12 across various tissues from the GTEx database. Expression values are presented in TPM (Transcripts Per Million), calculated from a gene model with isoforms collapsed into a single gene. Recently GTEx has adopted TPM instead of RPKM as the unit for comparing RNA-seq samples. TPM can be converted from RPKM using the formula: $TPM = RPKM / \text{sum}(RPKM) * 1.0 \times 10^{-6}$ for each sample/column.

Reference

1. Truong, T. *et al.* Multiethnic genome-wide association study of differentiated thyroid cancer in the EPITHYR consortium. *International Journal of Cancer* **148**, 2935–2946 (2021).
2. Anderson, C. A. *et al.* Data quality control in genetic case-control association studies. *Nat Protoc* **5**, 1564–1573 (2010).
3. Riboli, E. *et al.* European Prospective Investigation into Cancer and Nutrition (EPIC): study populations and data collection. *Public Health Nutr* **5**, 1113–1124 (2002).
4. Gudmundsson, J. *et al.* A genome-wide association study yields five novel thyroid cancer risk loci. *Nat Commun* **8**, 14517 (2017).
5. Mantel, N. & Haenszel, W. Statistical aspects of the analysis of data from retrospective studies of disease. *J Natl Cancer Inst* **22**, 719–748 (1959).
6. Sudlow, C. *et al.* UK Biobank: An Open Access Resource for Identifying the Causes of a Wide Range of Complex Diseases of Middle and Old Age. *PLoS Med* **12**, e1001779 (2015).
7. Köhler, A. *et al.* Genome-wide association study on differentiated thyroid cancer. *J Clin Endocrinol Metab* **98**, E1674–1681 (2013).
8. Kurki, M. I. *et al.* FinnGen provides genetic insights from a well-phenotyped isolated population. *Nature* **613**, 508–518 (2023).
9. Barbeira, A. N. *et al.* Exploring the phenotypic consequences of tissue specific gene expression variation inferred from GWAS summary statistics. *Nat Commun* **9**, 1825 (2018).
10. Pain, O. *et al.* Imputed gene expression risk scores: a functionally informed component of polygenic risk. *Hum Mol Genet* **30**, 727–738 (2021).
11. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies. *Nat Genet* **48**, 245–252 (2016).