

Supplementary Information

Table of Contents

Supplementary Figures	1
Supplementary Table Legends	31

Supplementary Figures

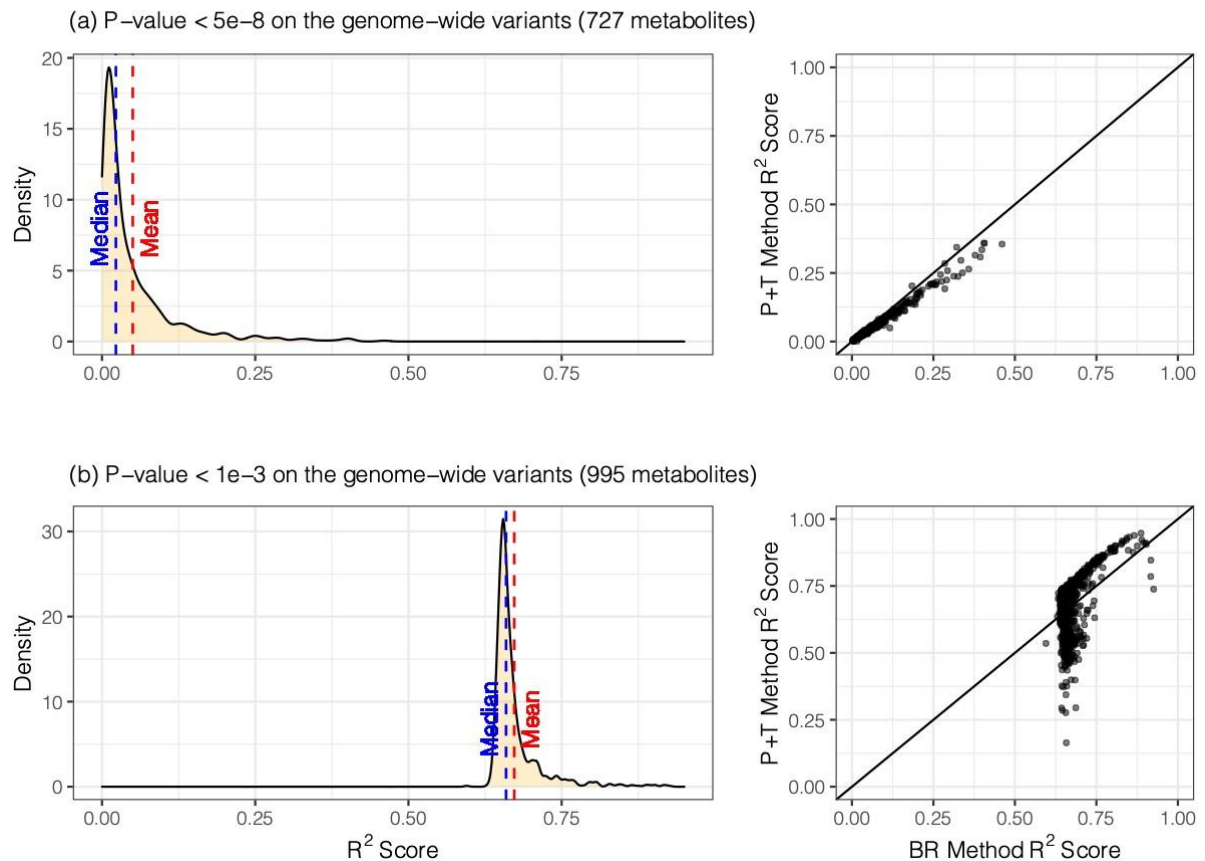


Figure S1: Performance (R^2) distribution of genetic scores and performance comparison (R^2) between Bayesian Ridge (BR) and pruning and thresholding (P+T) methods for Metabolon traits in internal validation. The density plots show the distributions of R^2 performance for genetic scores developed using BR method on different variant sets. P+T constructs genetic scores using weighted sum of a selected genetic variant set, where GWAS effect sizes of these variants are used as their weights. The scatter plots compare the performance of genetic scores developed using BR and P+T on different variant sets. It is noted that the variant set with p-value < 1e-3 resulted in an overfitting problem (see **Figure S4** for details). P-values in the GWAS for omic traits were derived by t-test in linear regression and all tests were two-sided.

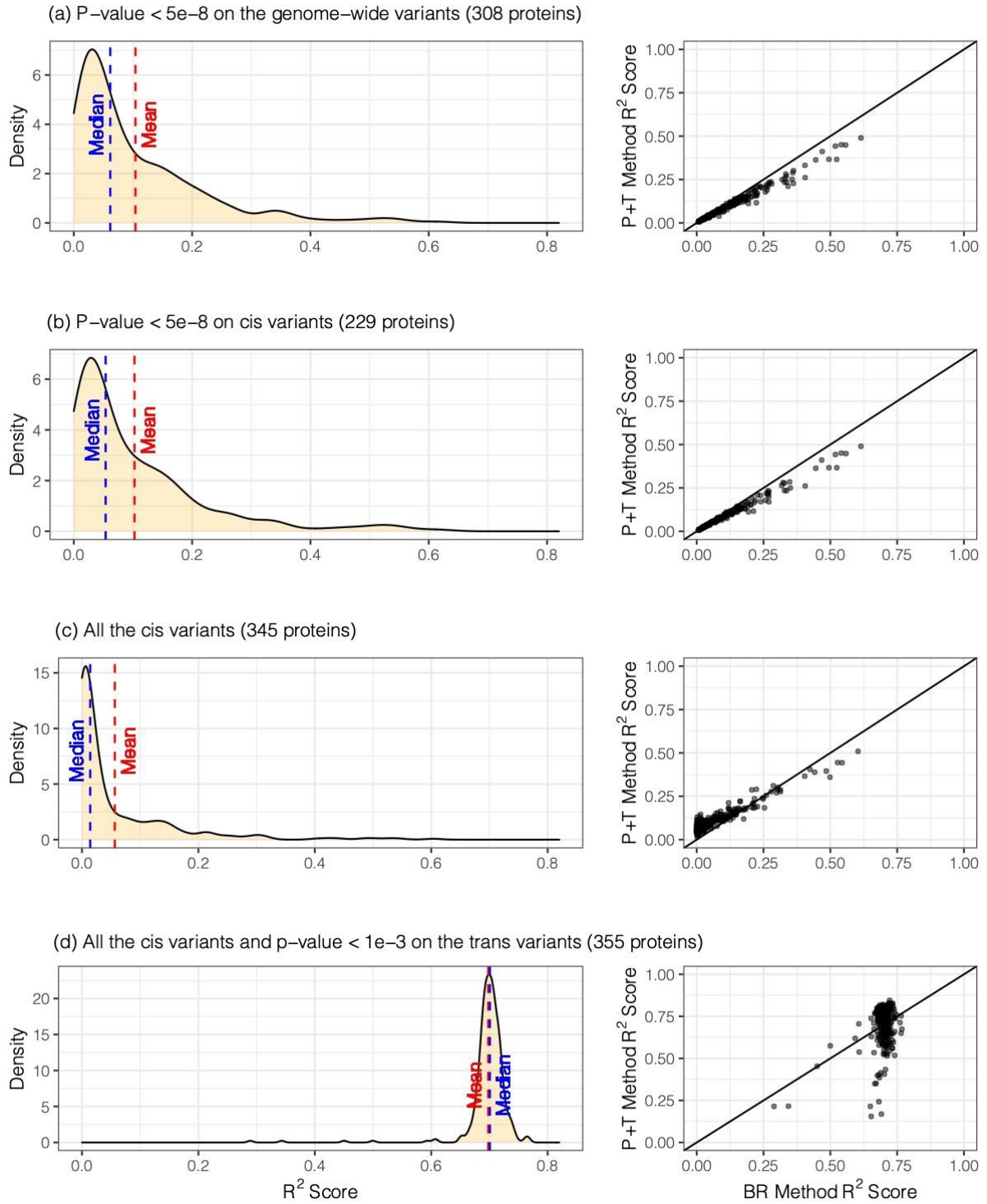


Figure S2: Performance (R^2) distribution of genetic scores and performance comparison (R^2) between BR and P+T methods for Olink traits in internal validation. The density plots show the distributions of R^2 performance for genetic scores developed using BR method on different variant sets. The scatter plots compare the performance of genetic scores developed using BR and P+T on different variant sets. P-values in the GWAS for omic traits were derived by t-test in linear regression and all tests were two-sided.

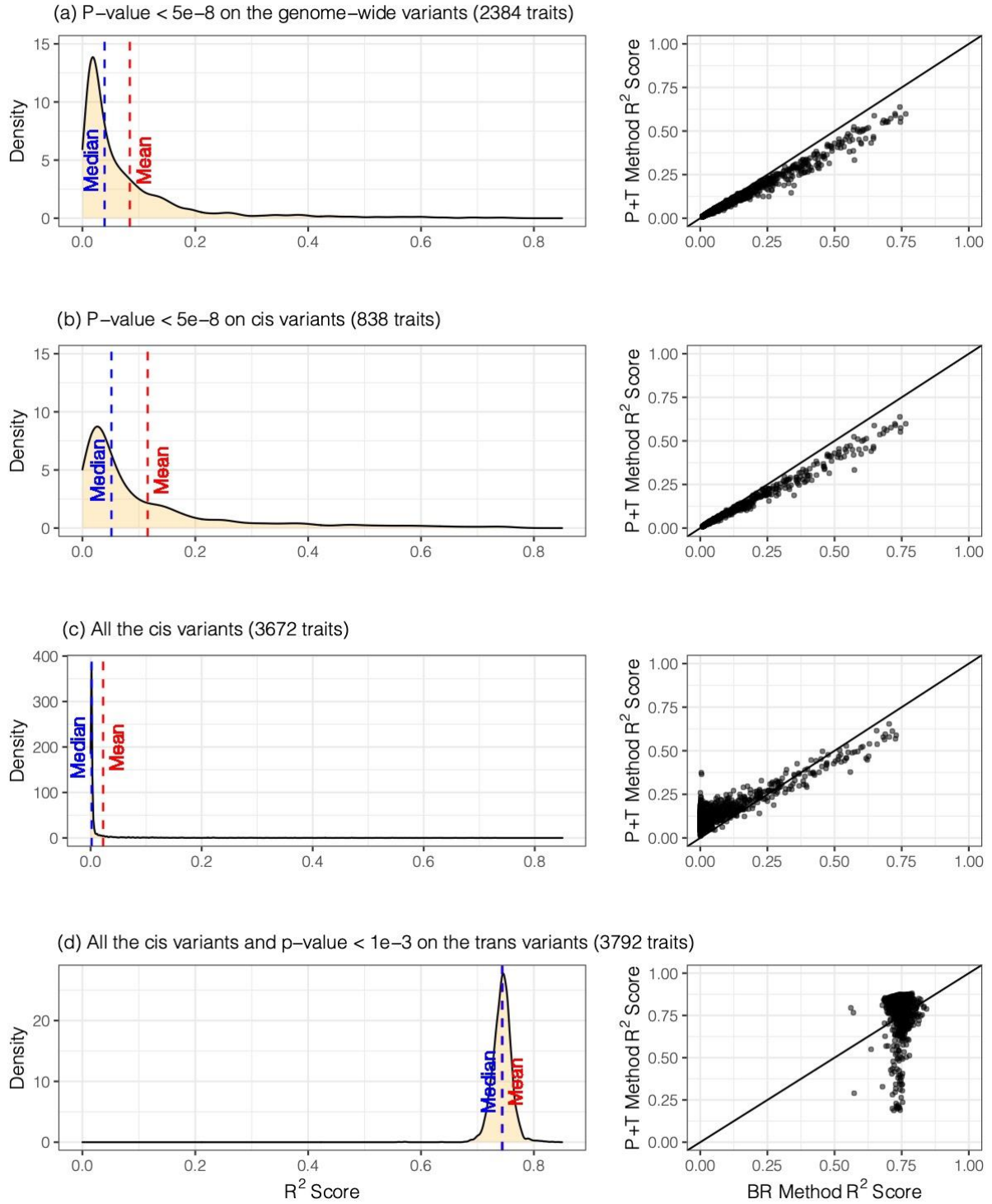


Figure S3: Performance (R^2) distribution of genetic scores and performance comparison (R^2) between BR and P+T methods for SomaScan traits in internal validation. The density plots show the distributions of R^2 performance for genetic scores developed using BR method on different variant sets. The scatter plots compare the performance of genetic scores developed using BR and P+T on different variant sets. P-values in the GWAS for omic traits were derived by t-test in linear regression and all tests were two-sided.

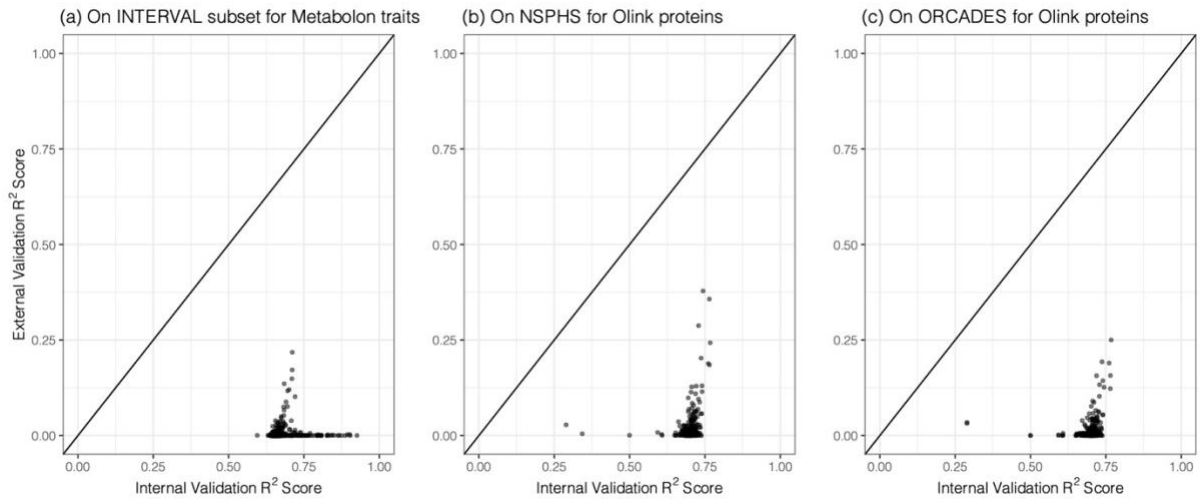


Figure S4: Performance (R^2) comparison between internal and external validation for genetic scores of Metabolon and Olink traits. The genetic scores were constructed using BR method on the set of genome wide variants with p -value $< 1 \times 10^{-3}$ for Metabolon traits and the set of all *cis* variants + p -value $< 1 \times 10^{-3}$ on the *trans* variants for Olink traits. P-values in the GWAS for omic traits were derived by t-test in linear regression and all tests were two-sided.

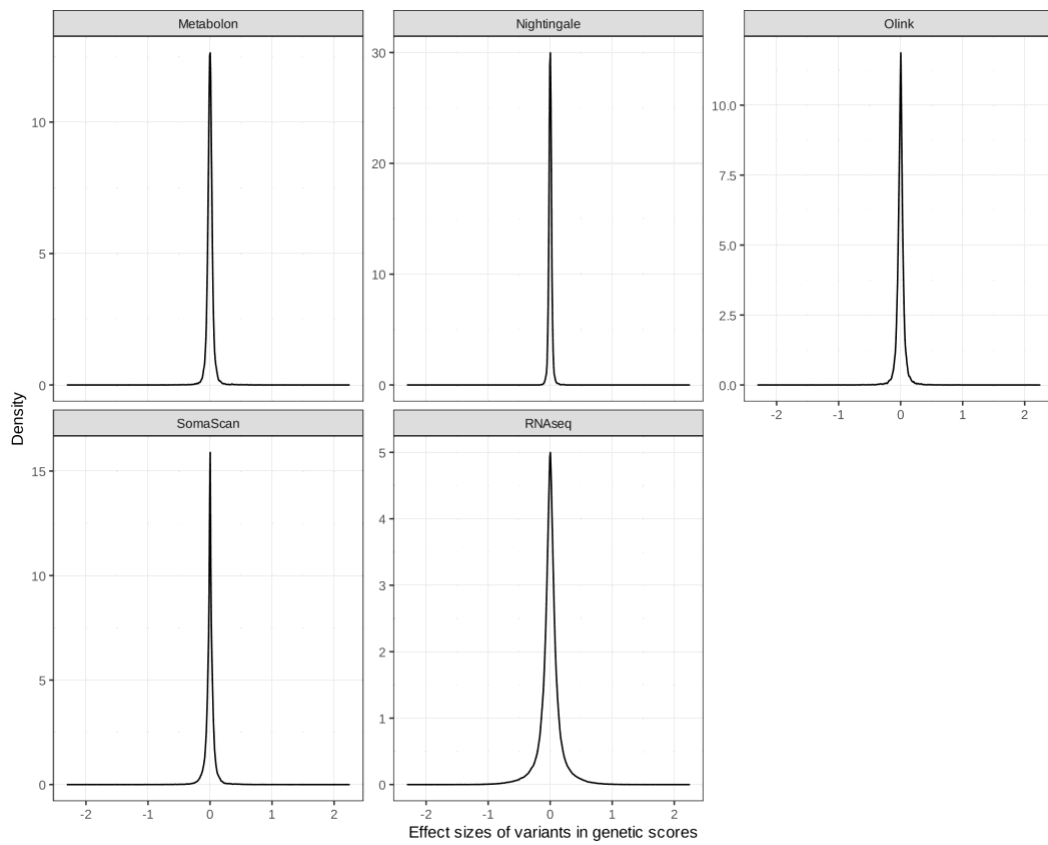


Figure S5: Distribution of effect sizes of genetic scores at each platform. This analysis took into consideration all the composing variants of developed genetic scores at each platform, and the figure shows the distribution of their effect sizes by platform.

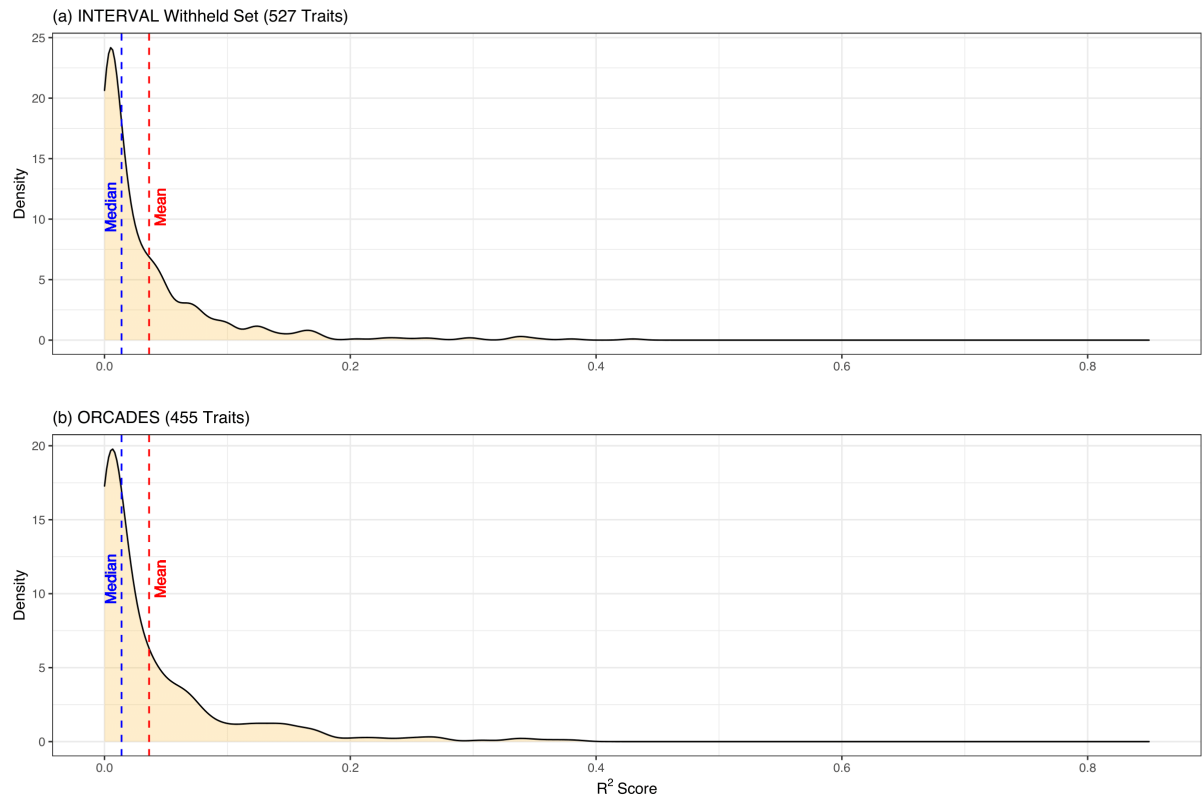


Figure S6: Distribution of R^2 performance in external validation for genetic scores of Metabolon traits. This analysis included all the traits validated in the external cohort or the INTERVAL withheld set.

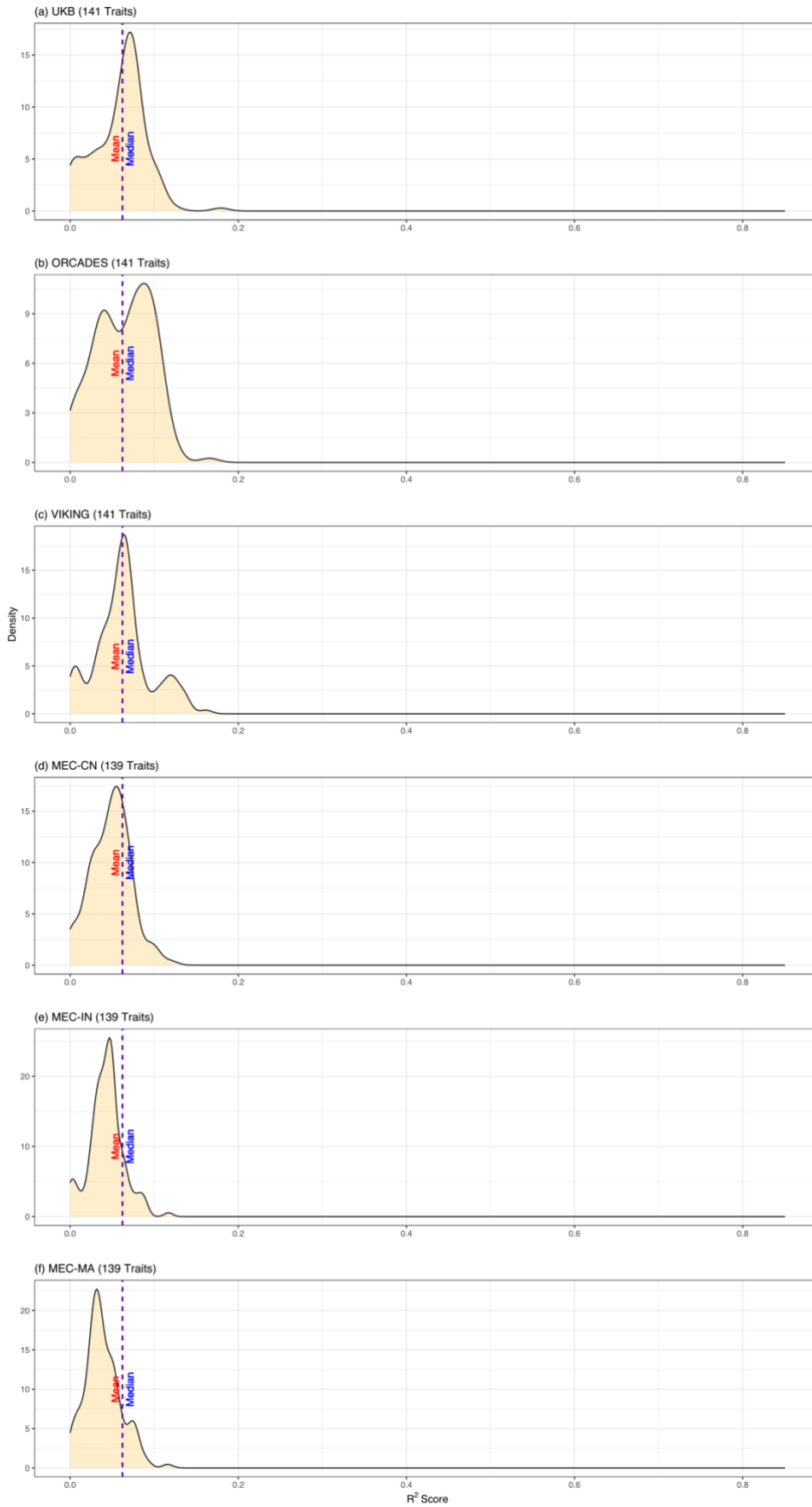


Figure S7: Distribution of R^2 performance in external validation for genetic scores of Nightingale traits. This analysis included all the traits validated in each external cohort.

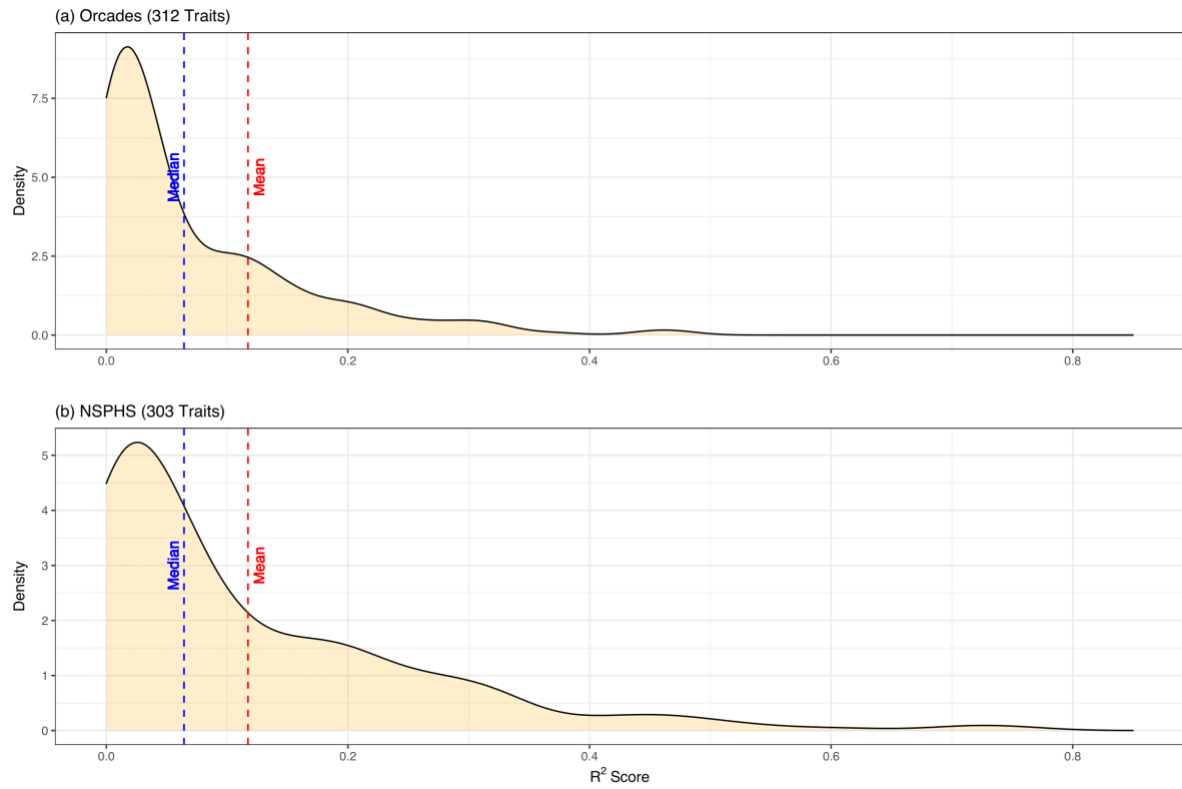


Figure S8: Distribution of R^2 performance in external validation for genetic scores of Olink traits. This analysis included all the traits validated in each external cohort.

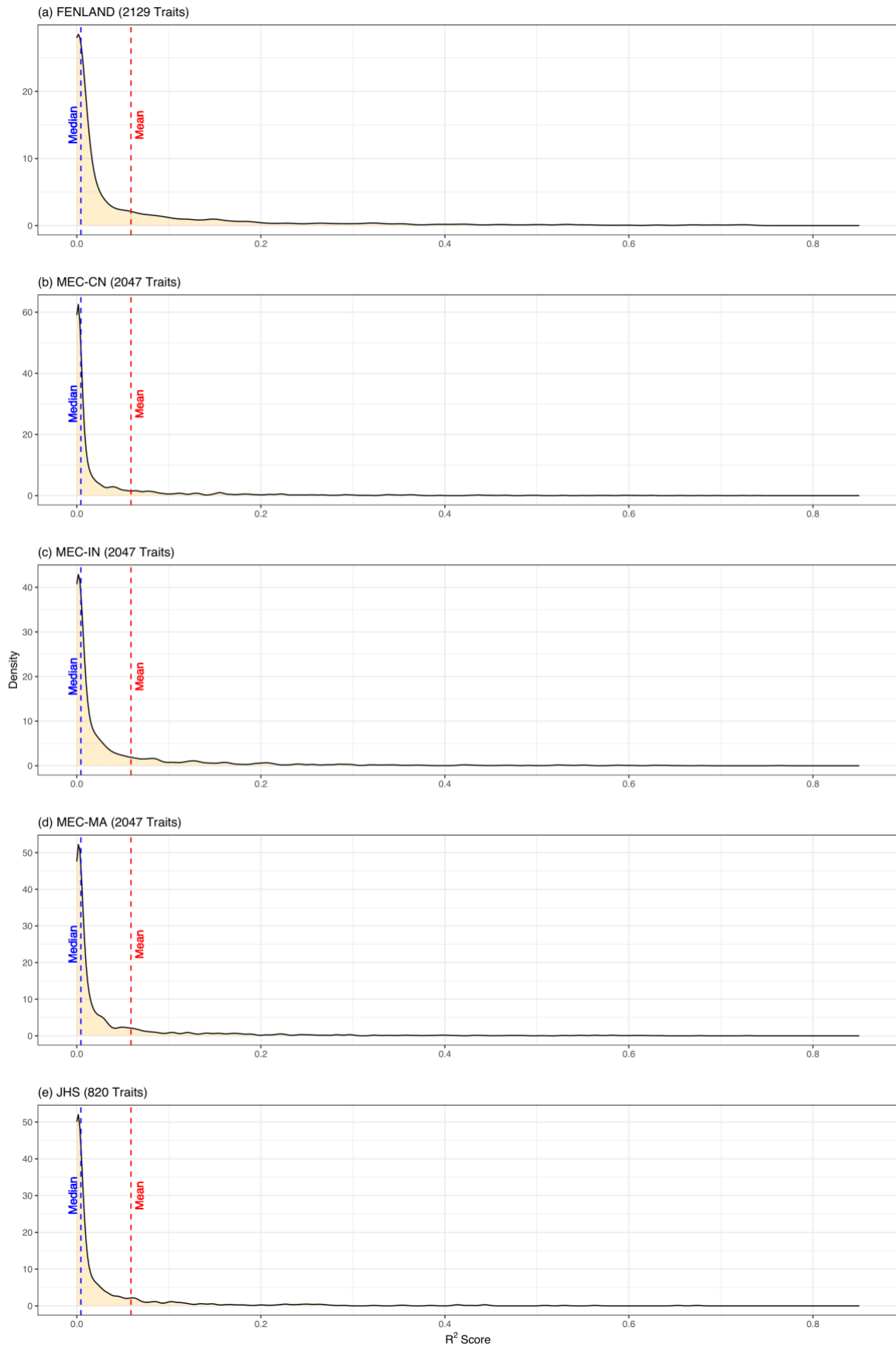


Figure S9: Distribution of R^2 performance in external validation for genetic scores of SomaScan traits. This analysis included all the traits validated in each external cohort.

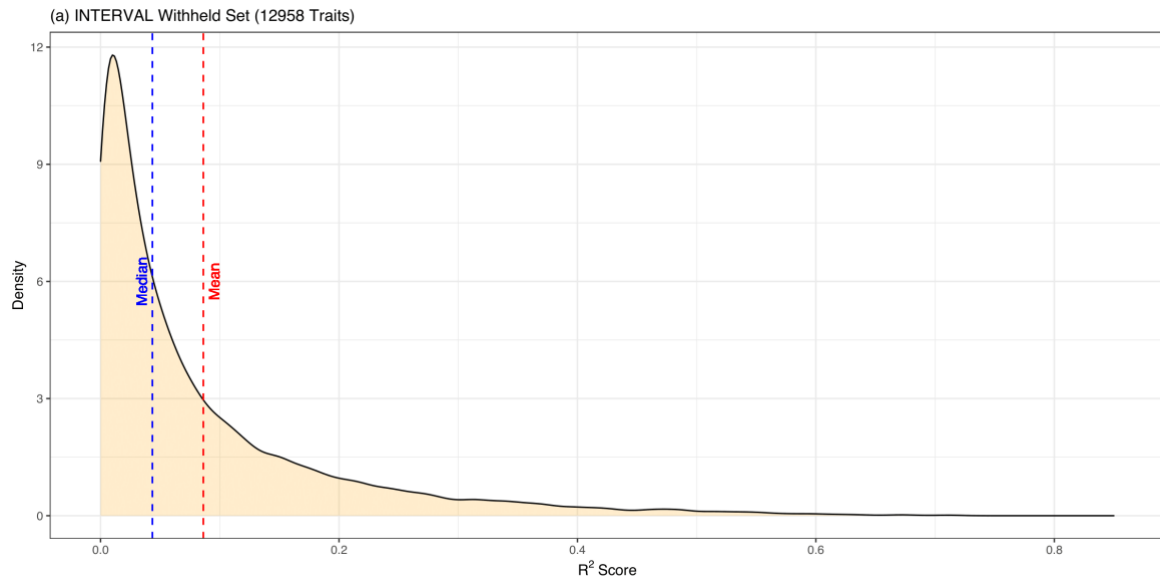


Figure S10: Distribution of R^2 performance in external validation for genetic scores of gene expression traits. This analysis included all the traits validated in the INTERVAL withheld set.

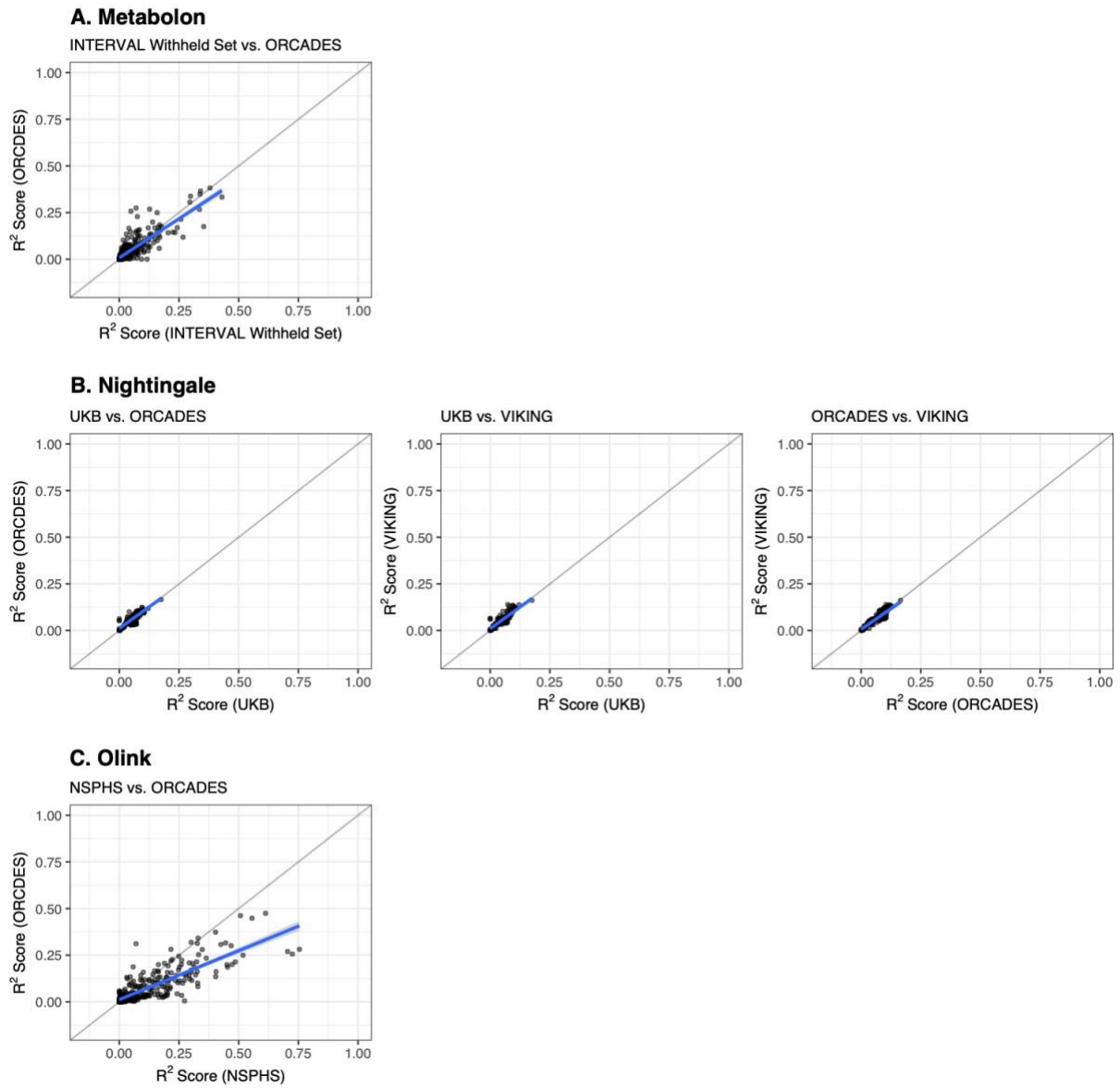


Figure S11: R^2 performance comparison of genetic scores between external European cohorts. The analyses included all the overlapped traits between two external validations at each platform. The blue line shows the linear models fitting all the performance comparison points.

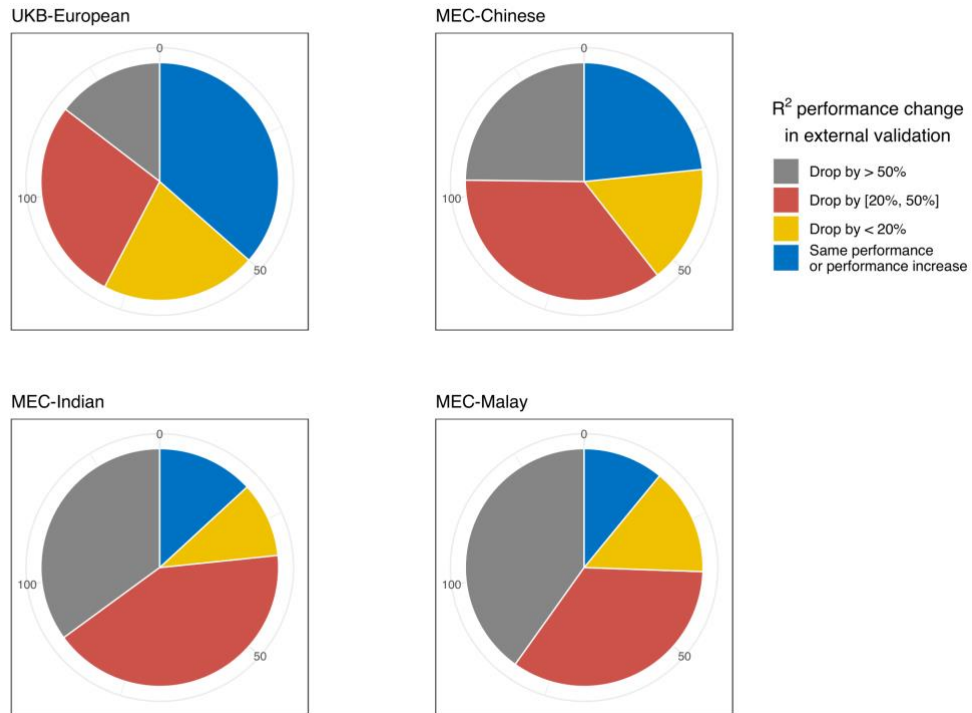


Figure S12: The number of genetic scores by R² performance change between internal and external validation for Nightingale traits. This figure shows the number of genetic scores for Nightingale traits by different levels of R² performance change between internal and external validation at each external population. This analysis only included genetic scores passing Bonferroni-adjusted significance threshold in internal validation.

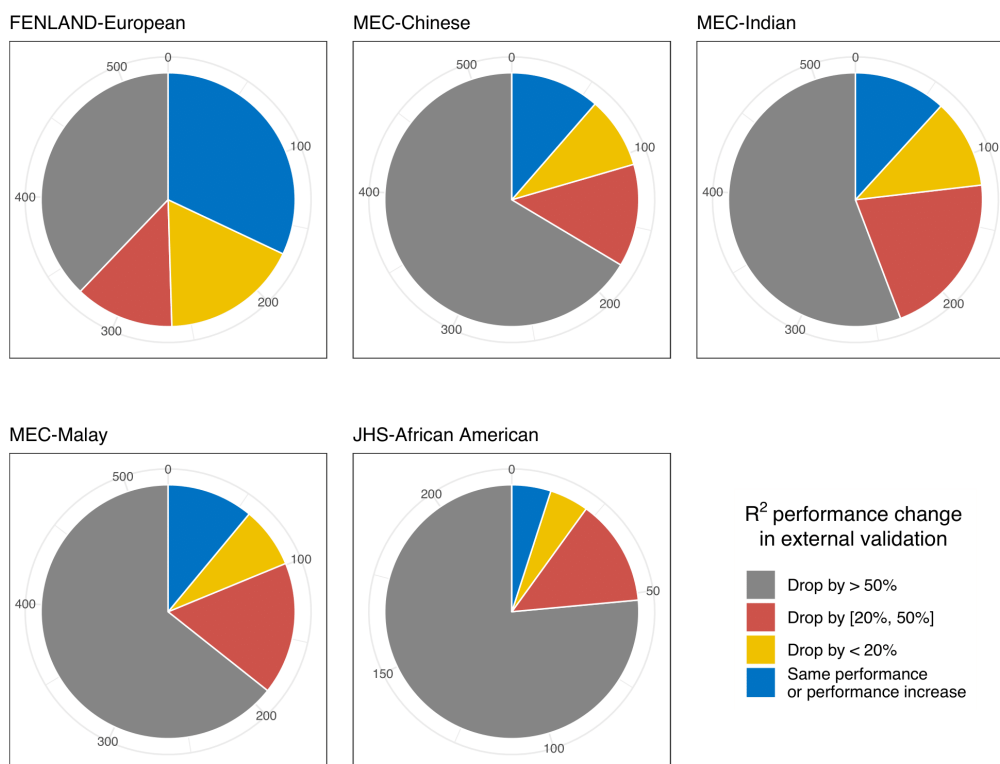


Figure S13: The number of genetic scores by R^2 performance change between internal and external validation for SomaScan traits. This figure shows the number of genetic scores for SomaScan traits by different levels of R^2 performance change between internal and external validation at each external population. This analysis only included genetic scores passing Bonferroni-adjusted significance threshold in internal validation.

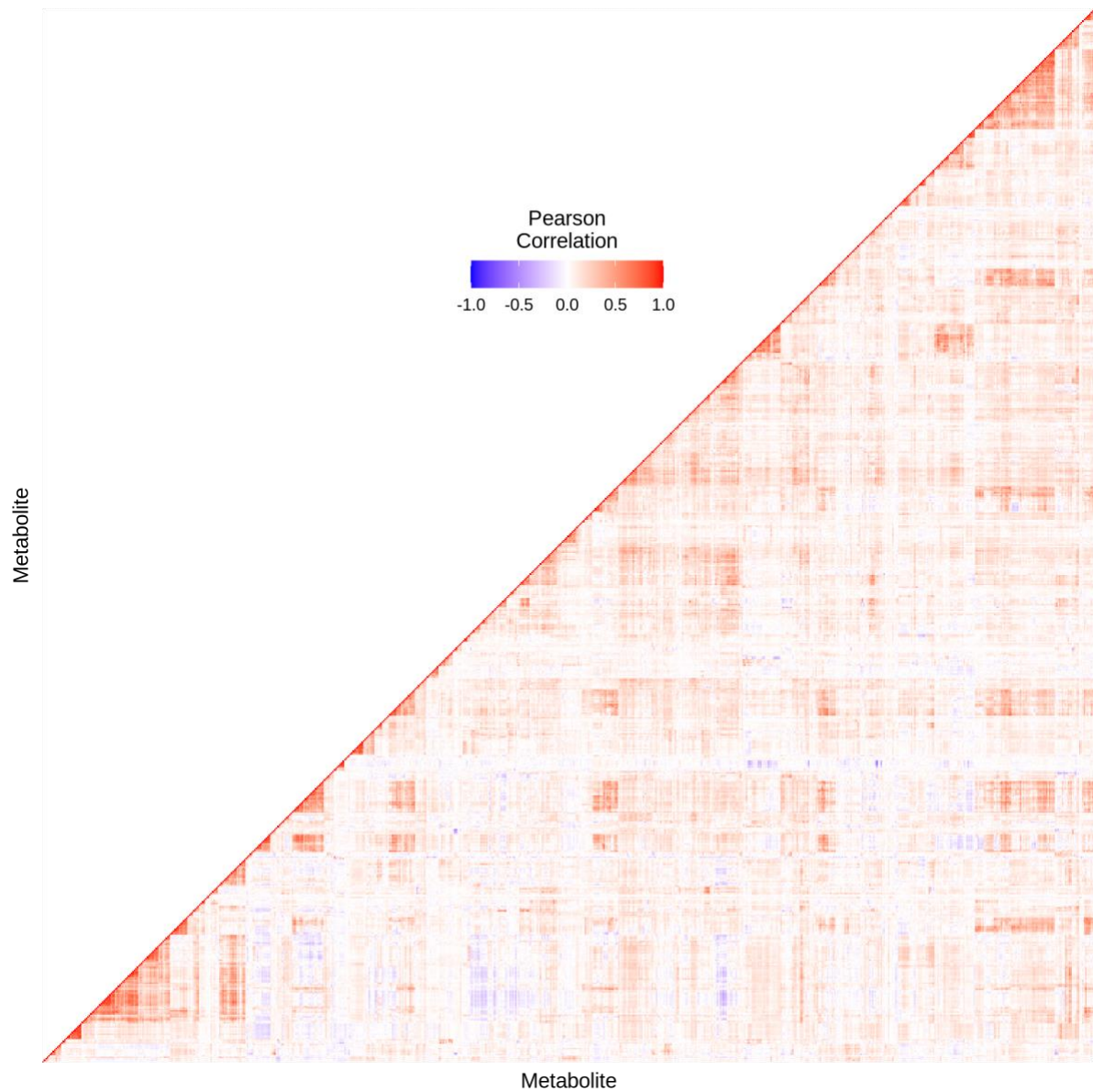


Figure S14: Correlation between Metabolite trait levels in INTERVAL. This analysis included all Metabolite traits qualified for genetic score development in this study and calculated Pearson's correlations between these traits using INTERVAL training samples.

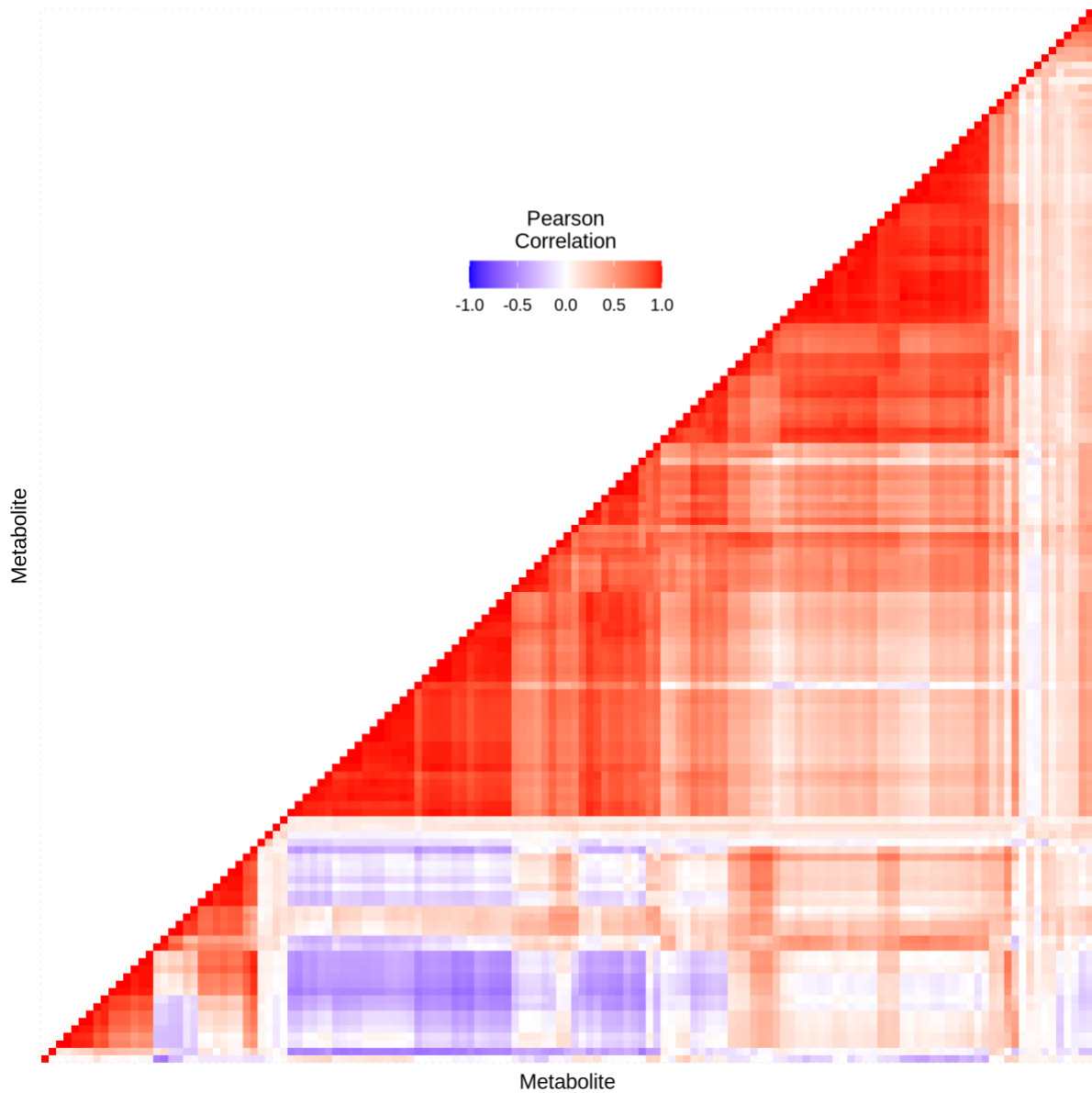


Figure S15: Correlation between Nightingale trait levels in INTERVAL. This analysis included all Nightingale traits qualified for genetic score development in this study and calculated Pearson's correlations between these traits using INTERVAL training samples.

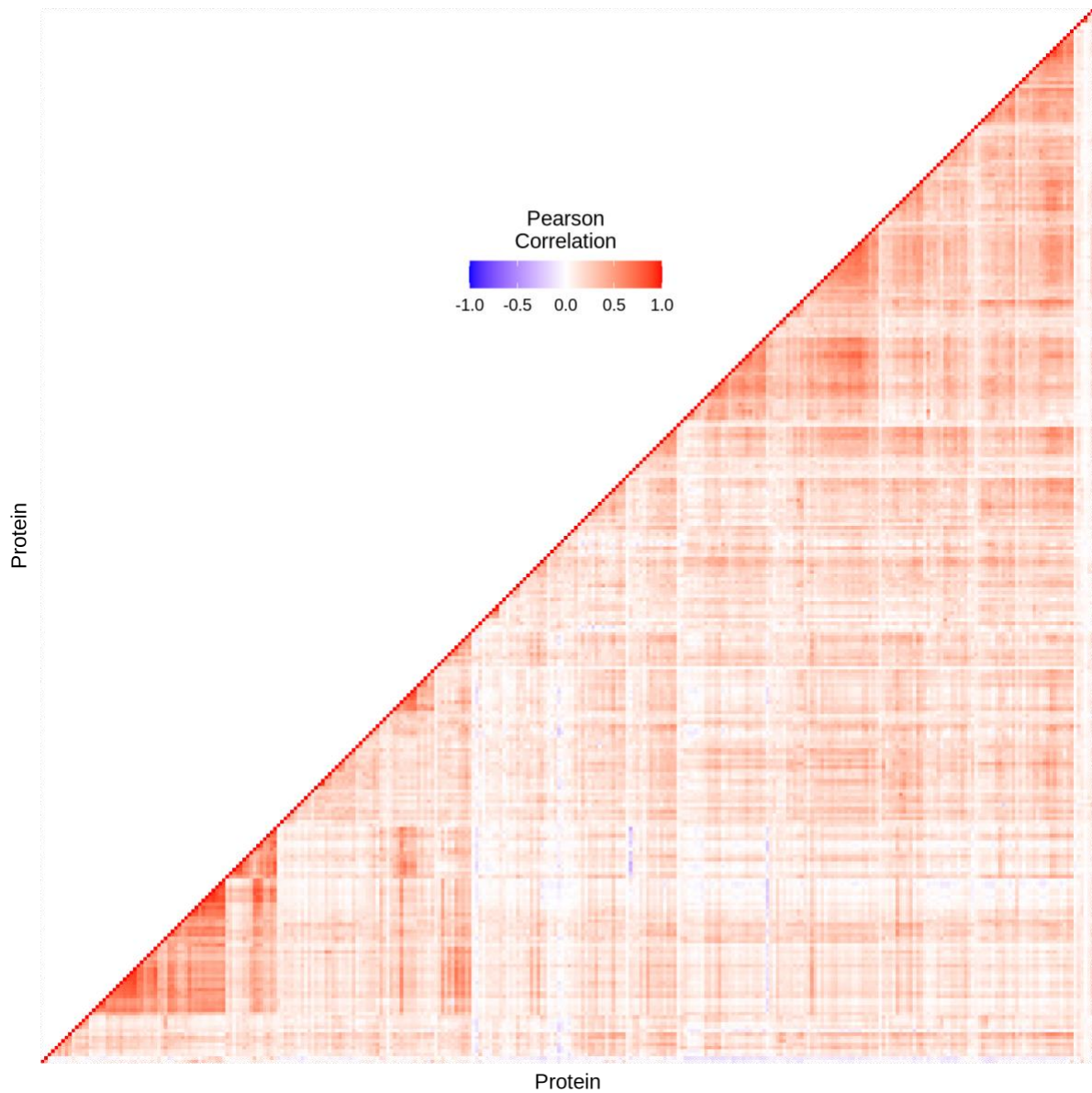


Figure S16: Correlation between Olink trait levels in INTERVAL. This analysis included all Olink proteins qualified for genetic score development in this study and calculated Pearson's correlations between these traits using INTERVAL training samples.



Figure S17: Correlation between SomaScan trait levels in INTERVAL. This analysis included all SomaScan proteins qualified for genetic score development in this study and calculated Pearson's correlations between these traits using INTERVAL training samples.



Figure S18: Correlation between RNAseq trait levels in INTERVAL. This analysis included all gene expression traits qualified for genetic score development in this study and calculated Pearson's correlations between these traits using INTERVAL training samples. Gene expression traits were grouped by chromosome in this analysis due to the large size.

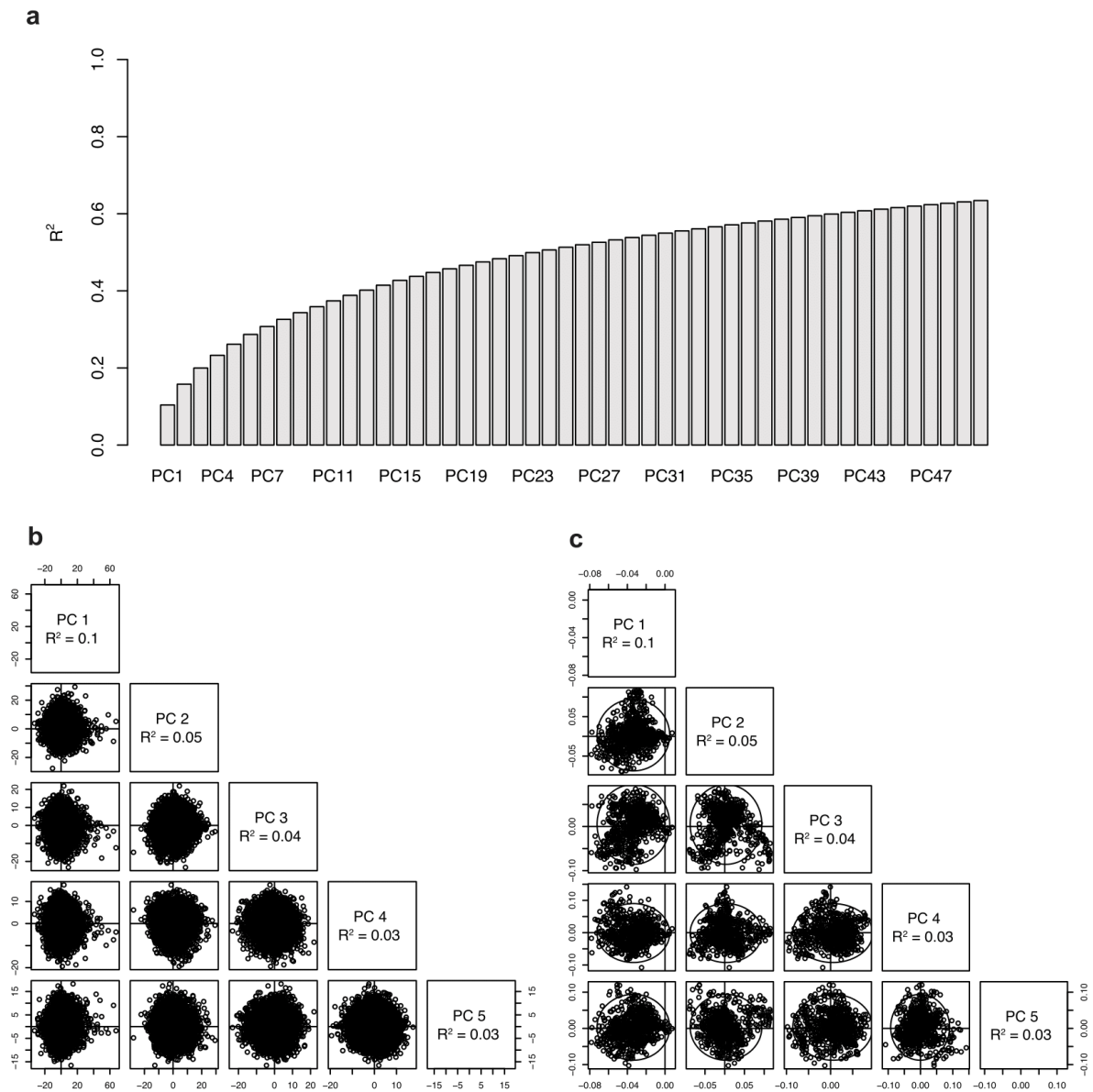


Figure S19: Principal component analysis (PCA) of Metabolon data in INTERVAL. a, Cumulative explained variance (R^2) by the top 50 PCs. **b,** Top 5 PC score comparison of the samples. **c,** Top 5 PC loading comparison of the traits. In b and c, the ellipse shows the Hotelling's 95% confidence. This analysis included all Metabolon traits qualified for genetic score development in this study.

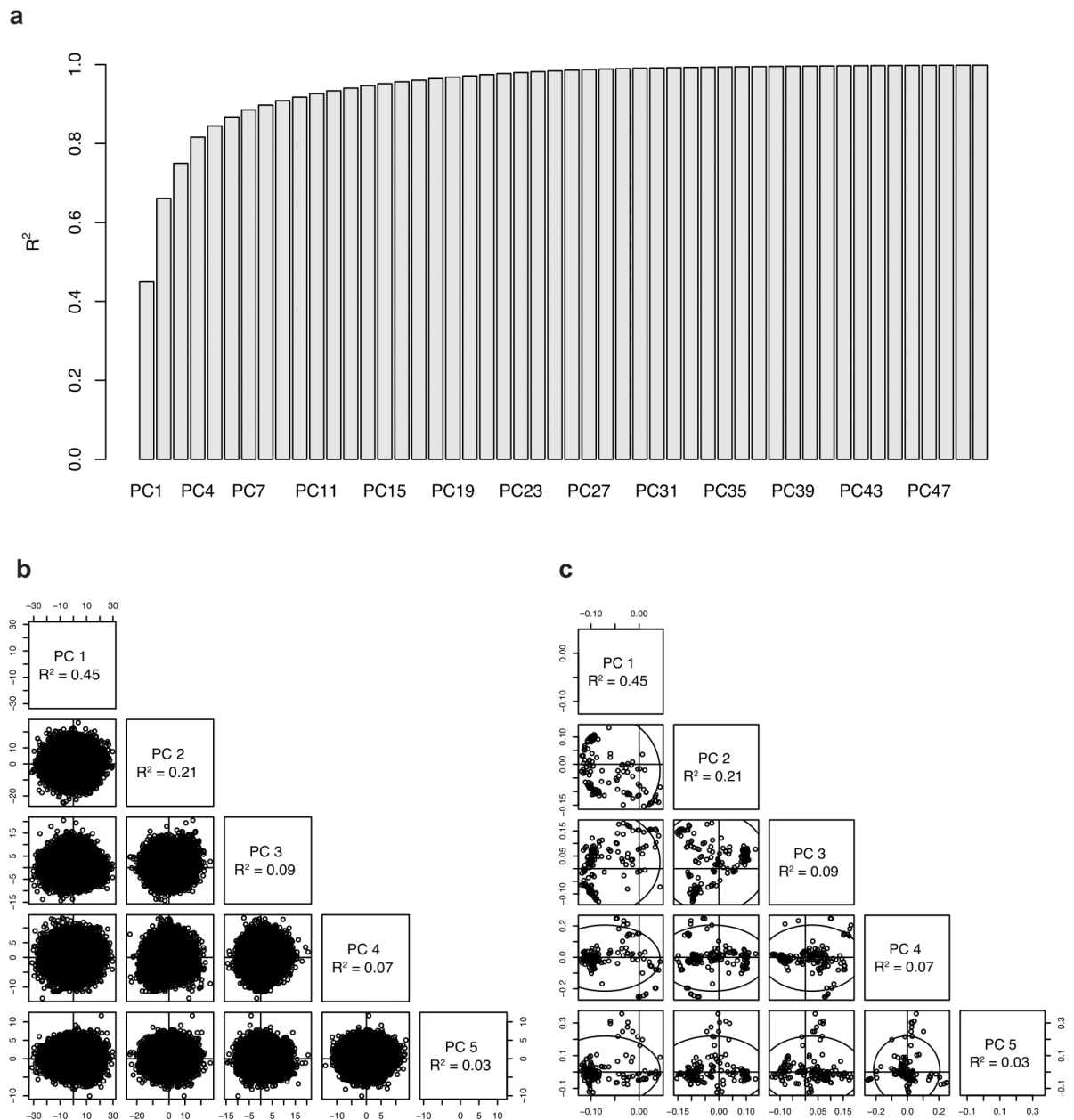


Figure S20: Principal component analysis of Nightingale data in INTERVAL. a, Cumulative explained variance (R^2) by the top 50 PCs. b, Top 5 PC score comparison of the samples. c, Top 5 PC loading comparison of the traits.

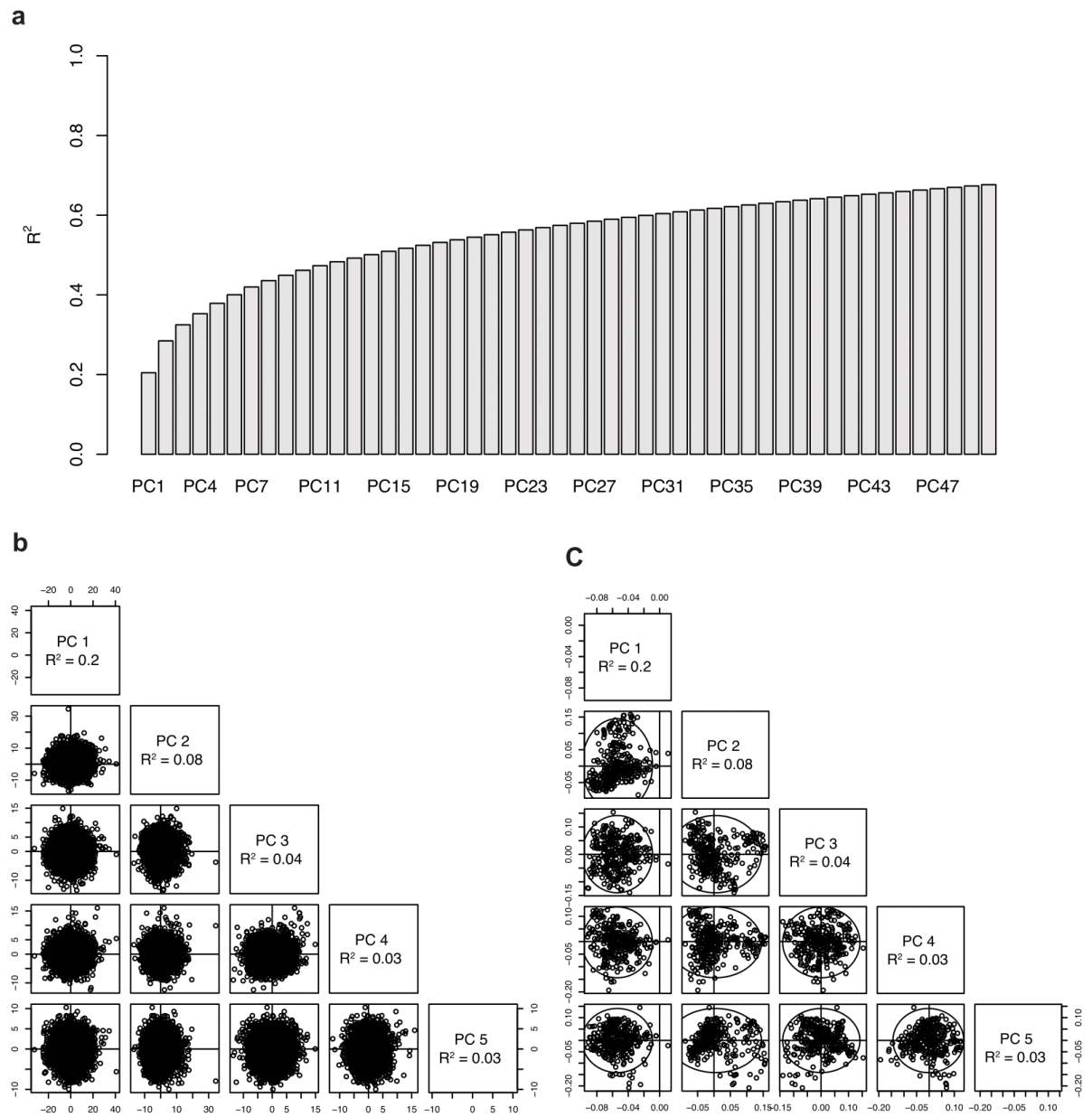


Figure S21: Principal component analysis of Olink data in INTERVAL. a, Cumulative explained variance (R^2) by the top 50 PCs. **b,** Top 5 PC score comparison of the samples. **c,** Top 5 PC loading comparison of the traits.

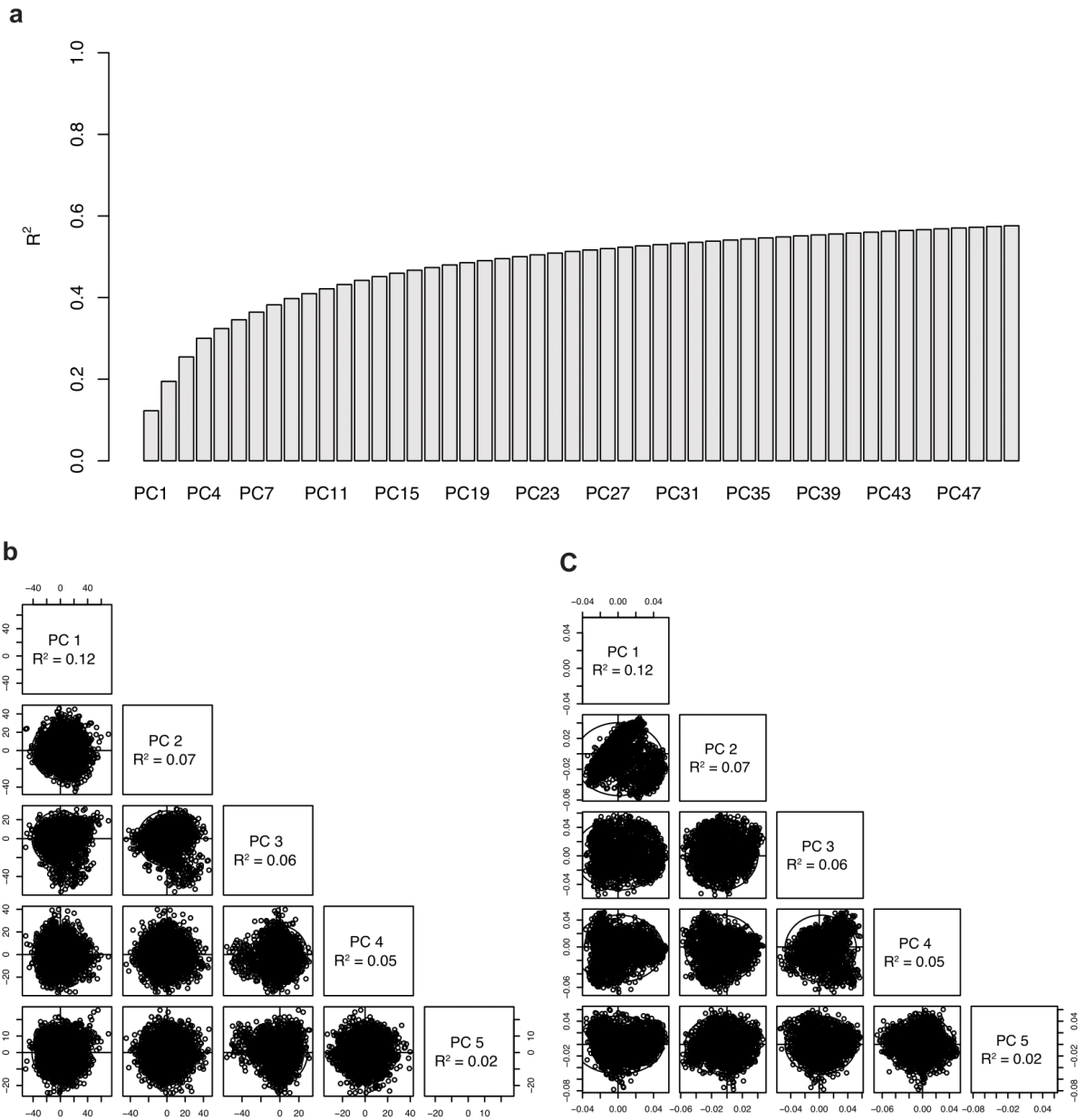


Figure S22: Principal component analysis of SomaScan data in INTERVAL. a, Cumulative explained variance (R^2) by the top 50 PCs. **b**, Top 5 PC score comparison of the samples. **c**, Top 5 PC loading comparison of the traits.

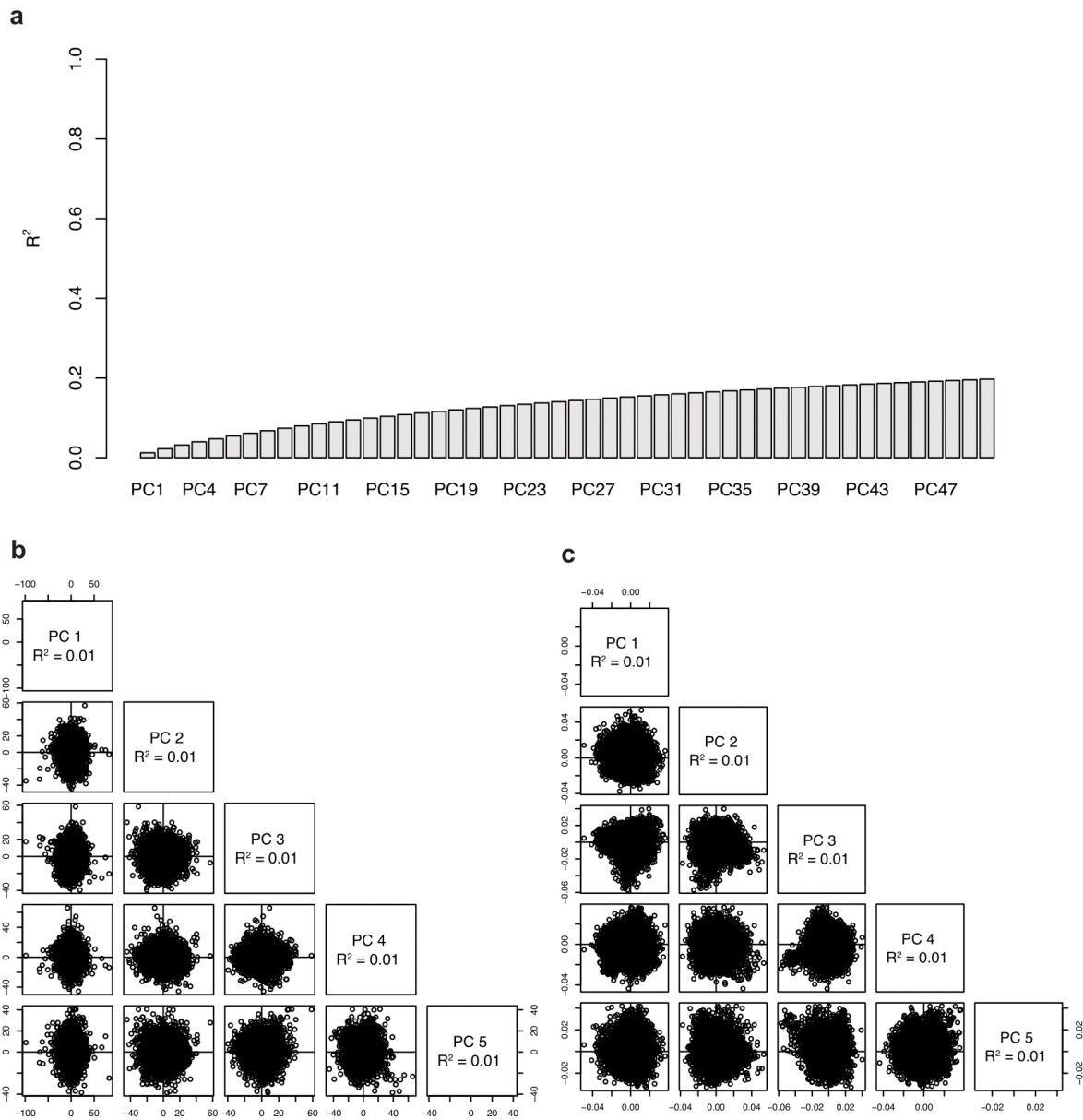


Figure S23: Principal component analysis of RNAseq data in INTERVAL. a, Cumulative explained variance (R^2) by the top 50 PCs. **b,** Top 5 PC score comparison of the samples. **c,** Top 5 PC loading comparison of gene expressions.

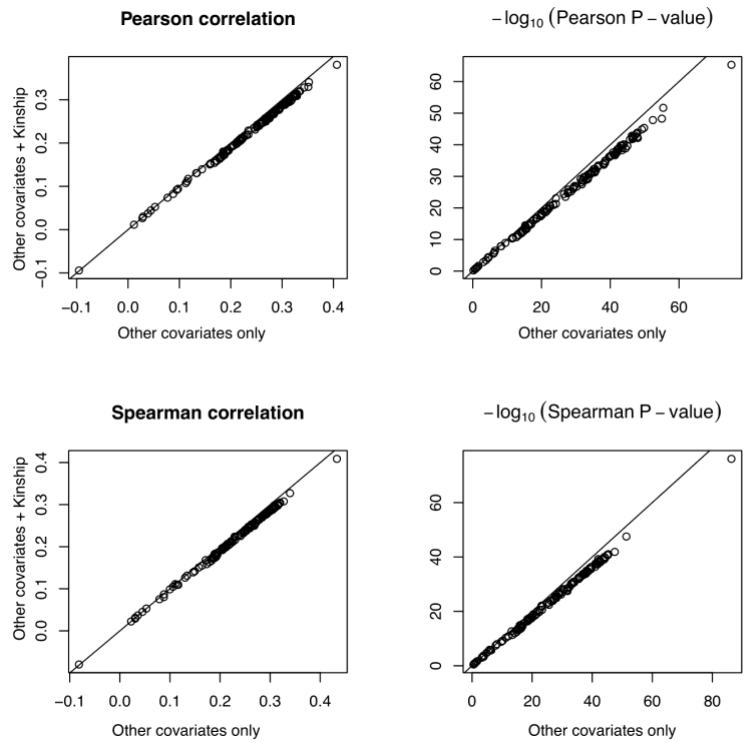


Figure S24: The impact of adjustment for family structure on genetic score validation for Nightingale traits in ORCADES. The figures compare the validation results for Nightingale traits in ORCADES using traits levels with or without adjustment for kinship. The x-axis shows the validation results using traits levels adjusted for sex, age, BMI, season of venepuncture, year of venepuncture, genotyping array and top 20 genetic principal components only and the y-axis shows the validation results with traits levels adjusted for the same set of covariates + kinship. P-values in this analysis were derived by t-test (Pearson's correlation coefficient) or Mann-Whitney U test (Spearman's correlation coefficient), and all tests were two-sided.



Figure S25: R^2 performance of Bayesian ridge in internal validation with different hyperparameter settings for SomaScan subset. This figure presents the R^2 performance of Bayesian ridge in internal validation with different hyperparameter settings, i.e. α_1 , α_2 , λ_1 and λ_2 for 20 randomly selected SomaScan traits. The x-axis shows all possible combinations of α_1 , α_2 , λ_1 and λ_2 taken from $\{0, 10^{-10}, 10^{-5}, 10^{-3}, 10^{-1}, 10, 10^3, 10^5, 10^{10}\}$, and the red points are all these combinations of α_1 , α_2 , λ_1 and λ_2 taken from $\{0, 10^{-10}, 10^{-5}, 10^{-3}\}$. This analysis used the variant set of p -value $< 5 \times 10^{-8}$ on genome-wide variants for BR (two-sided t-test in linear regression; **Methods**).



Figure S26: R^2 performance of Bayesian ridge in internal validation with different hyperparameter settings for Olink subset. This analysis used the variant set of p -value $< 5 \times 10^{-8}$ on genome-wide variants for BR (two-sided t-test in linear regression; **Methods**).

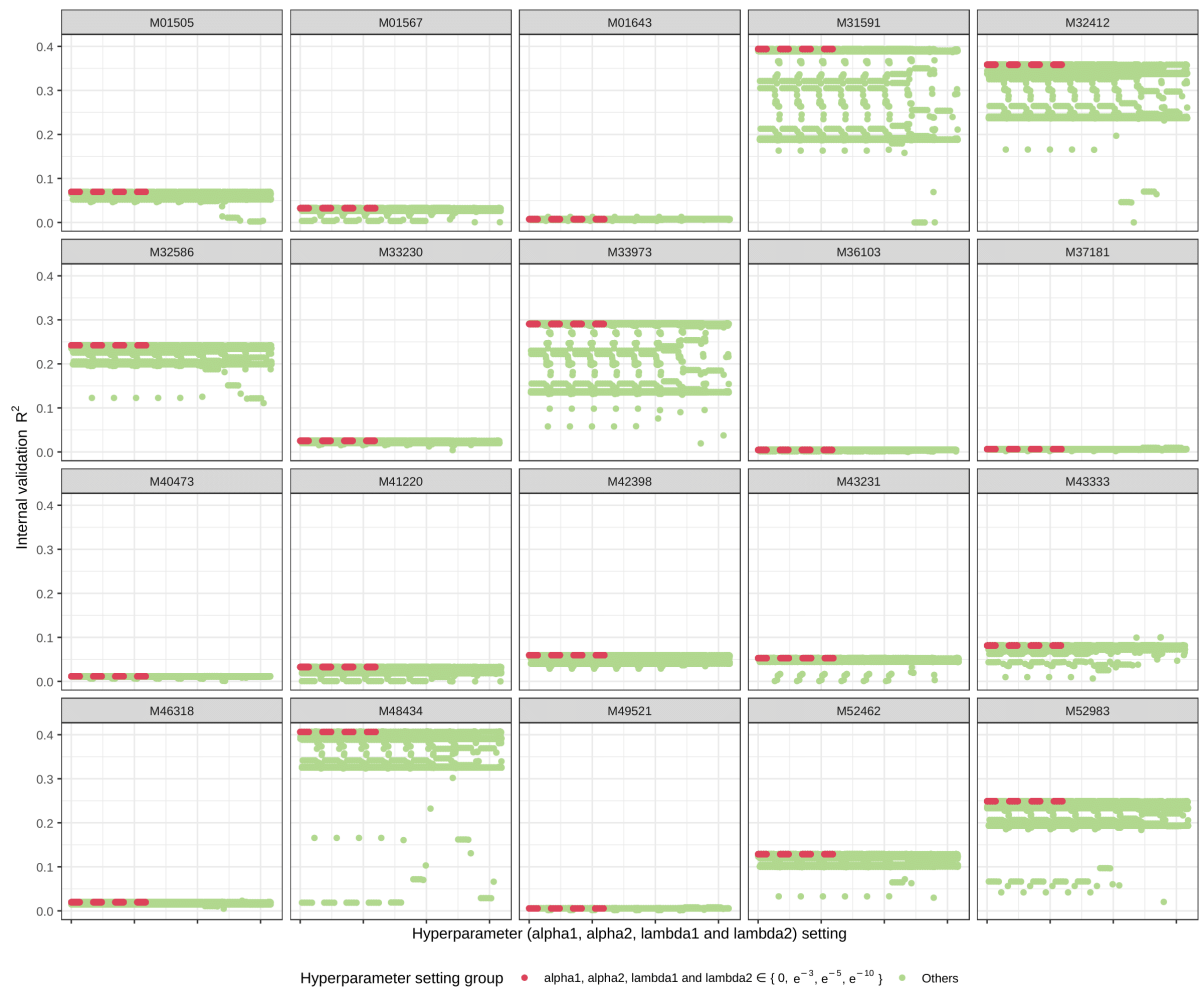


Figure S27: R^2 performance of Bayesian ridge in internal validation with different hyperparameter settings for Metabolon subset. This analysis used the variant set of $p\text{-value} < 5 \times 10^{-8}$ on genome-wide variants for BR (two-sided t-test in linear regression; **Methods**).

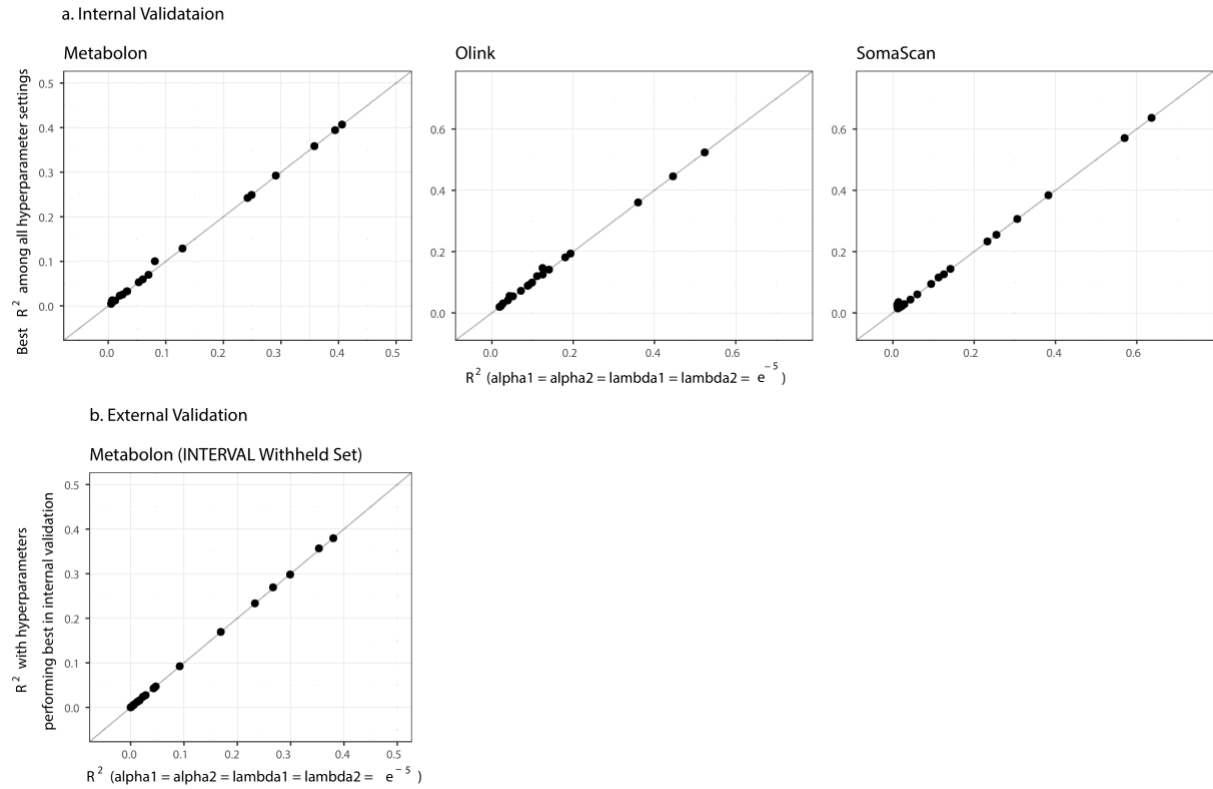


Figure S28: R^2 performance comparison of Bayesian ridge with the best performing hyperparameters through extensive search and the hyperparameters $\alpha_1 = \alpha_2 = \lambda_1 = \lambda_2 = 10^{-5}$ in internal (a) and external validation (b). Figure (a) compares the R^2 performance of Bayesian ridge with the best performing α_1 , α_2 , λ_1 and λ_2 (extensive search from $\{0, 10^{-10}, 10^{-5}, 10^{-3}, 10^{-1}, 10, 10^3, 10^5, 10^{10}\}$; y-axis) and $\alpha_1 = \alpha_2 = \lambda_1 = \lambda_2 = 10^{-5}$ (x-axis) in internal validation for 20 randomly selected traits in each platform (Metabolon, Olink and SomaScan). Figure (b) further compares the R^2 performance of Bayesian ridge with the best performing α_1 , α_2 , λ_1 and λ_2 (in internal validation) and $\alpha_1 = \alpha_2 = \lambda_1 = \lambda_2 = 10^{-5}$ for 20 randomly selected Metabolon traits in external validation (INTERVAL withheld set). This analysis used the variant set of p -value $< 5 \times 10^{-8}$ on genome-wide variants for BR (two-sided t-test in linear regression; **Methods**).

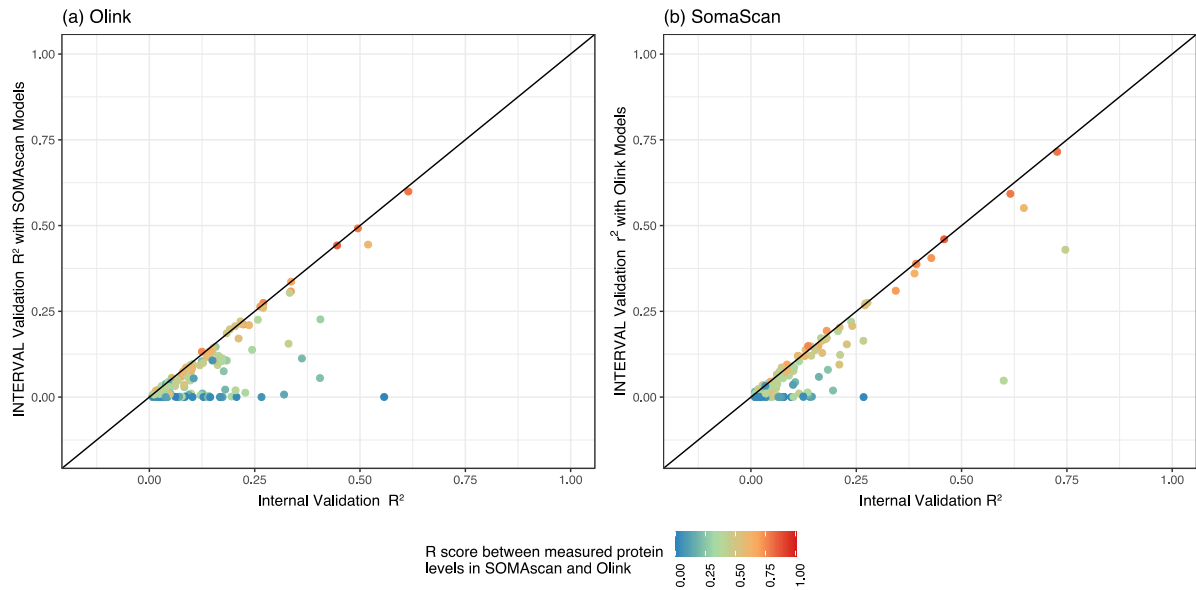


Figure S29: R^2 Performance comparison of genetic scores for shared proteins between SomaScan and Olink in INTERVAL. We compared the internal validation R^2 performance of 169 shared proteins on SomaScan (or Olink) with R^2 performance of their corresponding genetic scores trained on Olink (or SomaScan) in predicting protein levels on SomaScan (or Olink) using all the INTERVAL training samples. The points are coloured by the Pearson's r score between the actual proteins levels of a protein measured by SomaScan and Olink for those samples who were assayed with both platforms.

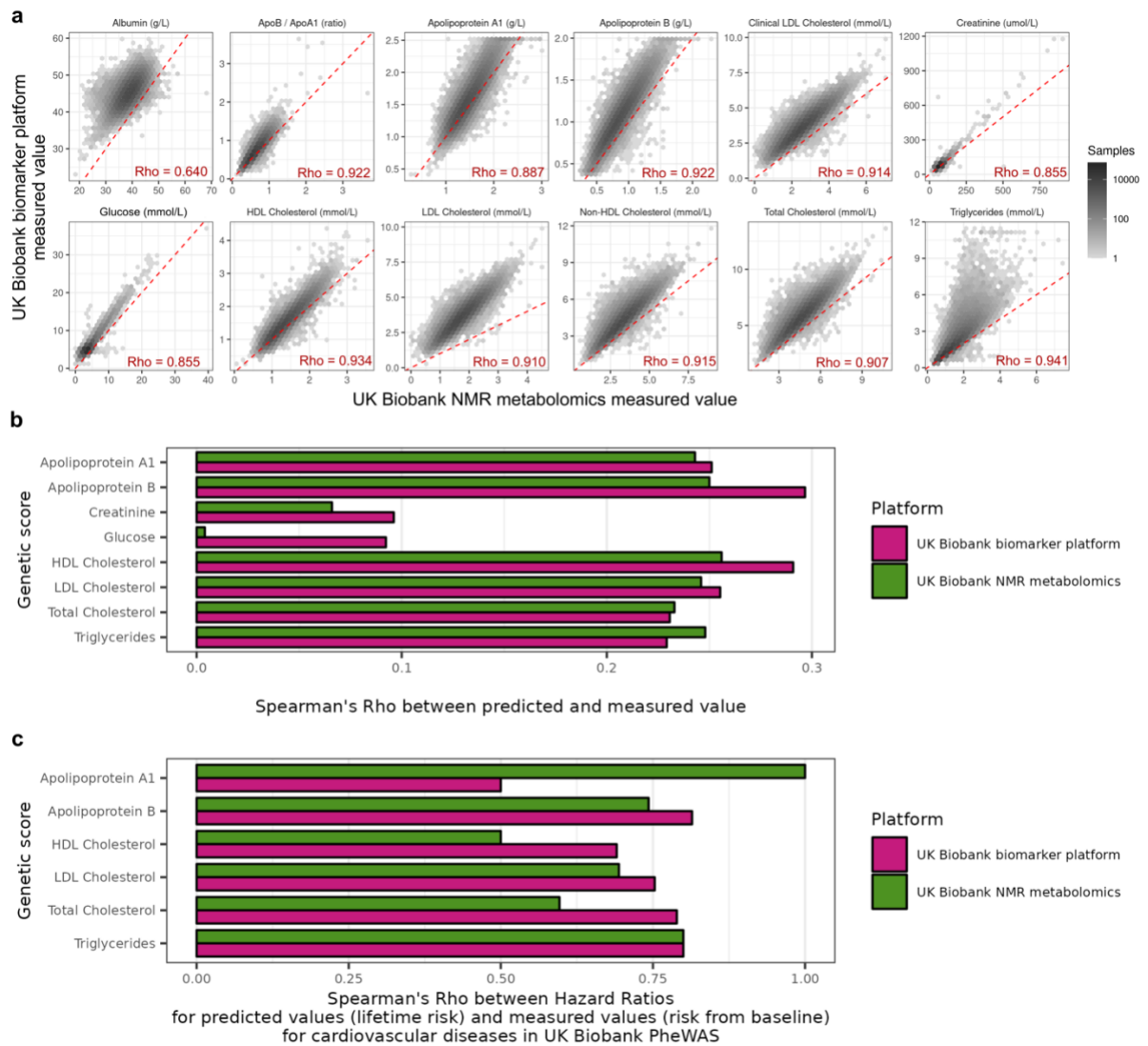


Figure S30: NMR metabolomics genetic scores are robust to measurement technology. (a) Comparison of NMR metabolomics measurements to clinical chemistry measurements in 116,472 UK Biobank participants for the 12 metabolites and ratios quantified by both platforms. Hexagonal bins show the number of participants (on a log₁₀ scale) whose biomarker concentrations quantified by the NMR metabolomics platform (x-axis) and clinical chemistry biomarker platform (y-axis) are comparable at each paired x-y interval. The dashed red line shows x=y where biomarker concentrations quantified by both platforms are identical. (b) Comparison of genetic scores performance for predicting biomarker concentrations quantified by the NMR metabolomics platform (green) and clinical biochemistry platform (pink) in 97,088 UK Biobank participants in the UK Biobank defined White British genetic ancestry cluster. Biomarkers present in panel A, but not panel B, are those without genetic scores. (c) Comparison of lifetime risk of cardiovascular diseases (phecodes 400-500) predicted by genetic scores to incident cardiovascular disease risk (from baseline assessment) predicted by biomarker concentrations quantified by NMR metabolomics and clinical biochemistry platforms.

Supplementary Table Legends

Table S1: Summary information of Metabolon traits and validation results of their genetic scores. This table lists the summary information of Metabolon traits and genetic scores, as well as shows the performances (explained variance R^2 and Spearman correlation score Rho) of Metabolon genetic scores in internal and external validations (INTERVAL withheld set and ORCADES). P-values were estimated using two-sided t-test for R^2 and two-sided Mann-Whitney U test for Rho. The column “OMICSPRED ID” gives the unique identifier of a genetic score in the OmicsPred online portal; the column “#SNP” shows the number of variants comprising the genetic score.

Table S2: Summary information of Nightingale traits and validation results of their genetic scores. This table lists the summary information of Nightingale traits and genetic scores, as well as shows the performances (explained variance R^2 and Spearman correlation score Rho) of Nightingale genetic scores in internal and external validations (UKB, ORCADES, VIKING, MEC-Chinese, MEC-Malay and MEC-Indian). P-values were estimated using two-sided t-test for R^2 and two-sided Mann-Whitney U test for Rho. The column “OMICSPRED ID” gives the unique identifier of a genetic score in the OmicsPred online portal; the column “#SNP” shows the number of variants comprising the genetic score.

Table S3: Summary information of Olink traits and validation results of their genetic scores. This table lists the summary information of Olink traits and genetic scores, as well as shows the performances (explained variance R^2 and Spearman correlation score Rho) of Olink genetic scores in internal and external validations (ORCADES and NSPHS). P-values were estimated using two-sided t-test for R^2 and two-sided Mann-Whitney U test for Rho. Note that the array may target (1) multiple proteins e.g. a protein complex; (2) proteins encoded by multiple genes or (3) a combination of both, in which multiple UniProt IDs and gene names are listed with separator ";". The column “OMICSPRED ID” gives the unique identifier of a genetic score in the OmicsPred online portal; the column “#SNP” shows the number of variants comprising the genetic score.

Table S4: Summary information of SomaScan traits and validation results of their genetic scores. This table lists the summary information of SomaScan traits and genetic scores, as well as shows the performances (explained variance R^2 and Spearman correlation score Rho) of SomaScan genetic scores in internal and external validations (Fenland, MEC-Chinese, MEC-Malay, MEC-Indian and JHS). P-values were estimated using two-sided t-test for R^2 and two-sided Mann-Whitney U test for Rho. Note that some aptamers may target (1) multiple proteins e.g. a protein complex; (2) proteins encoded by multiple genes or (3) a combination of both, in which multiple UniProt IDs and gene names are listed with separator "|". More than one aptamer can target the same protein, in which the same UniProt ID and Gene name are used. The column “OMICSPRED ID” gives the unique identifier of a genetic score in the OmicsPred online portal; the column “#SNP” shows the number of variants comprising the genetic score.

Table S5: Summary information of RNAseq traits and validation results of their genetic scores. This table lists the summary information of RNAseq traits and genetic scores, as well as shows the performances (explained variance R^2 and Spearman correlation score Rho) of RNAseq genetic scores in internal and independent validations (withheld INTERVAL subset). P-values were estimated using two-sided t-test for R^2 and two-sided Mann-Whitney U test for Rho. The column “OMICSPRED ID” gives the unique identifier of a genetic score in the OmicsPred online portal; the column “#SNP” shows the number of variants comprising the genetic score.

Table S6: Significant associations detected in PheWAS using UK Biobank data. This table lists all the significant associations (two-sided Wald test and FDR-corrected p-value < 0.05 for 11,576 tested traits) identified in the PheWAS with UKB Biobank data. The column "Internal Validation R²" gives the R² performance of trait genetic scores in internal validation with INTERVAL training data. Hazard ratio and its 95% confidence interval of the genetic score are given in the table for the associations.

Table S7: Summary statistics of phenotypes tested in PheWAS with UK biobank. This table lists the summary information of phenotypes tested in the PheWAS with UK biobank data and the number of significant associations (two-sided Wald test and FDR-corrected p-value < 0.05 for 11,576 tested traits) between genetic scores in each platform and each phenotype.

Table S8: Settings of GWAS and genetic score training for omic traits across different platforms in INTERVAL. This table summarises the key steps of data pre-processing and analyses in GWAS and genetic score training for omic traits in each platform in INTERVAL.

Table S9: Demographic statistics of samples by Olink panel in each cohort. This table lists the demographic statistics of samples for Olink traits by panel in INTERVAL, ORCADES and NSPHS. *The mean of validation results for a protein overlapped between panels was taken as the validation performance of the protein.

Table S10: Groups of traits that are highly correlated in multi-omic platforms. We consider traits in each platform as vertices of an undirected graph and vertices are connected via edges if traits are correlated with Pearson $r > 0.9$ (based on the trait levels in INTERVAL training data). Then, subgraphs in this graph are used to identify groups of highly correlated traits in each platform. This analysis identified 2,225, 299, 700, 29, 13,663 (in total 16,916 groups out of 17,227 traits) highly correlated groups of traits in SomaScan, Olink, Metabolon, Nightingale and RNAseq respectively.

Table S11: Overview of the genetic and omic data used for genetic score validation in external cohorts and withheld INTERVAL subsets. This table summarises the key information on the genetic data and omic data used for validation in external cohorts (or withheld INTERVAL samples).