

Multi-omics and Multi-organ Aging Clocks Digitize Human Aging

Junhao Wen^{1,2,3,4,*}

¹Laboratory of AI and Biomedical Science (LABS), Columbia University, New York, NY, USA

²New York Genome Center (NYGC), New York, NY, USA

³Center for Innovation in Imaging Biomarkers and Integrated Diagnostics (CIMBID), Department of Radiology, Columbia University, New York, NY, USA

⁴Data Science Institute (DSI), Columbia University, New York, NY, USA

*Corresponding author:

Junhao Wen, junhao.wen89@gmail.com

622 W 168th St, New York, NY 10032

Abstract

Multi-organ biological aging clocks derived from clinical phenotypes and neuroimaging have emerged as valuable tools for studying human aging and disease^{1,2,3,4}. Plasma proteomics provides an additional molecular dimension to enrich these clocks⁵. Here, we used 2448 plasma proteins from 43,498 participants in the UK Biobank to develop 11 multi-organ proteome-based biological age gaps (ProtBAG). We compared them to 9 multi-organ phenotype-based biological age gaps (PhenoBAG¹) regarding genetics, causal associations with 525 disease endpoints (DE) from FinnGen and PGC, and their clinical promise to predict 14 disease categories and mortality. We highlighted critical clinical and methodological considerations for generating ProtBAG, including the need for age bias correction⁶ and addressing protein organ specificity to enhance model performance and generalizability. Genetic analyses revealed overlap between ProtBAGs and PhenoBAGs, including shared loci, genetic correlations, and colocalization signals. A three-layer causal network linked ProtBAG, PhenoBAG, and DE, exemplified by the pathway of obesity→renal PhenoBAG→renal ProtBAG to holistically understand human aging and disease. Combining features across multiple organs improved predictions for disease categories and mortality. These findings provide a framework for integrating multi-omics and multi-organ biological aging clocks in biomedicine. All results are publicly disseminated at <https://labs-laboratory.com/medicine/>.

Main

Multi-organ biological aging clocks, derived from neuroimaging and clinical phenotypes, are increasingly being explored in clinical research and computational neuroscience as tools to understand human aging, disease, and mortality^{1,2,4,7}. These clocks provide a comprehensive view of biological age, reflecting the functional and structural changes across different organs. While significant advancements have been made in leveraging phenotypic data for such models, there remains a growing interest in incorporating molecular-level data, such as plasma proteomics⁵, epigenetics⁸, and metabolomics⁹, to enrich the landscape of the multi-organ biological age. Plasma proteomics from different platforms (e.g., Olink¹⁰ and SomaScan¹¹) offers the unique ability to identify and quantify proteins and post-translational modifications with high sensitivity, potentially uncovering novel insights into organ-specific aging and its relationship with health and disease¹².

Despite its promise, deriving proteome-based biological age biomarkers presents several challenges and unresolved questions. One common practice observed in neuroimaging-derived brain age is to correct the age bias in an age prediction model, which can distort associations between the biological age gap (BAG) and disease outcomes if not properly corrected^{13,14,6,15}. That is, brain age tends to be overestimated for younger individuals and underestimated for older individuals, while predictions are most accurate for those whose ages are closer to the mean of the training dataset (**Fig. 1b**). Furthermore, the lack of organ specificity of plasma proteins (analogous to pleiotropy in genetics), where a protein is over-expressed in multiple organ tissues may complicate model development, leading to overfitting and reduced interpretability. Previous studies identified similar overfitting issues and addressed them by employing data-driven feature selection methods to mitigate the problem^{5,16}. In addition, key factors that influence model performance and generalizability, such as the type of omics data, the sample size and population demographic and disease status of the training sample, and the balance between the closeness of model fit and the clinical power of BAG, have not been systematically evaluated. These challenges highlight the need for systematic and reproducible evaluations of proteome-derived BAGs (i.e., ProtBAG)¹⁷. Addressing these gaps is essential to unlocking the full potential of plasma proteomics in aging research and its clinical applications.

Phenome-wide BAGs (PhenoBAG) and ProtBAG represent two essential aspects of human aging and disease causal pathways, connecting genetics→transcriptomics→proteomics (ProtBAG)→endophenotypes (PhenoBAG)→disease outcomes (DE). Our prior studies^{1,3} have examined the genetic architecture of 9 multi-organ PhenoBAG through genome-wide association studies (GWAS) and post-GWAS validations, such as genetic correlation¹⁸, polygenic risk scores¹⁹, and causal inference²⁰. A comprehensive framework to explore the overlap and distinctions between ProtBAG and PhenoBAG is currently lacking. Addressing this gap requires connecting genetics, ProtBAG, PhenoBAG, and disease endpoints (DE). Such an integrative approach is essential for developing a holistic understanding of the causal pathways for potential therapeutic development.

Multi-organ and multi-omics approaches^{21,3,1,22,23,7,24,5,25,26,27} are gaining prominence in modeling human aging and disease, driven by the hypothesis that integrating insights across multiple spatial and temporal scales better captures underlying disease-related neurobiological processes, thus enhancing diagnostic and prognostic biomarker discovery. For instance, Zhao et al.²⁸ demonstrated improved cognitive prediction by integrating brain and heart MRI features with PRS. Similarly, our prior work on AI/ML-derived brain disease subtypes showed enhanced systemic disease prediction when combining these brain imaging-derived biomarkers with

PRS²⁹. However, the potential of multi-omics and multi-organ BAGs as complementary biomarkers for disease and mortality remains unexplored.

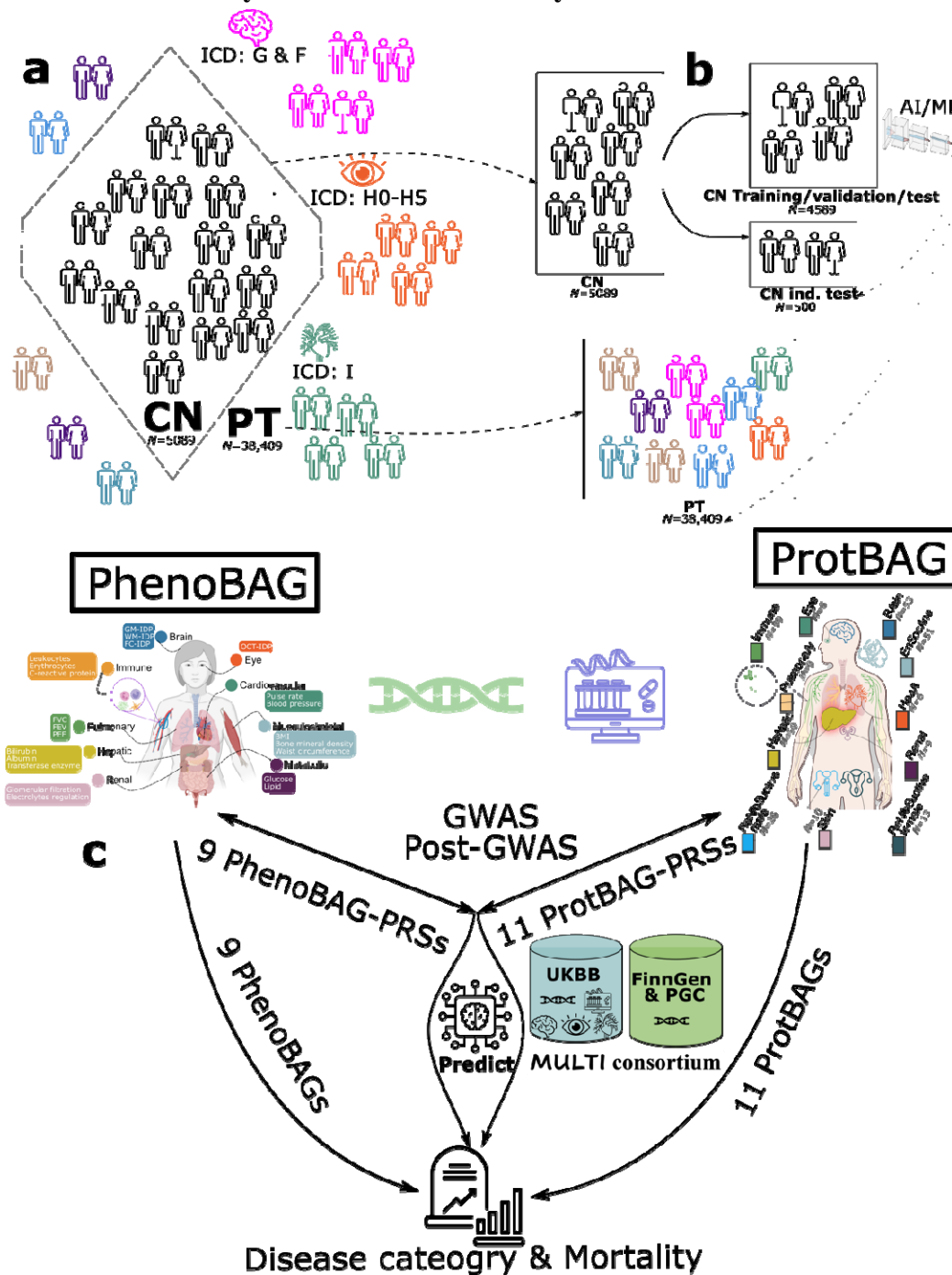
This study used 2448 Olink plasma proteins from 43,498 UK Biobank participants (**UKBB** and **Supplementary eTable 1**) to develop 11 organ-specific ProtBAGs (**Method 1**). We systematically compared the 11 ProtBAGs with 9 PhenoBAGs derived from our previous studies^{1,3} (**Method 2-3**). We evaluated the influence of key methodological components (**Method 4**) on model performance and clinical interpretation using the 11 ProtBAGs. Subsequently, we examine their genetic architecture and causal relationships with 525 DEs from FinnGen³⁰ and PGC³¹ (**Method 5**). Finally, we assessed the potential of ProtBAGs, PhenoBAGs, and their PRSs for predicting disease categories and mortality (**Method 6**). All results and pre-trained AI/ML models are publicly disseminated at the MEDICINE portal: <https://labs-laboratory.com/medicine/>.

Results

Biological age prediction performance of the 11 ProtBAGs derived from three AI/ML models

To rigorously evaluate the performance of biological age prediction models, we partitioned the 5089 healthy control (CN, without any pathologies) participants into the CN training/validation/test ($N=4589$) and independent test (ind. test; $N=500$) datasets. **Extended Data Fig. 1** details this study's population selection and overall workflow. The CN training set was used for model development and nested cross-validation when applicable, while the independent test set provided an unbiased assessment of model performance (**Supplementary eTable 1**).

Extended Data Figure 1: Schematic diagram of the definition of populations to derive ProtBAG and overall analytic workflow of the study

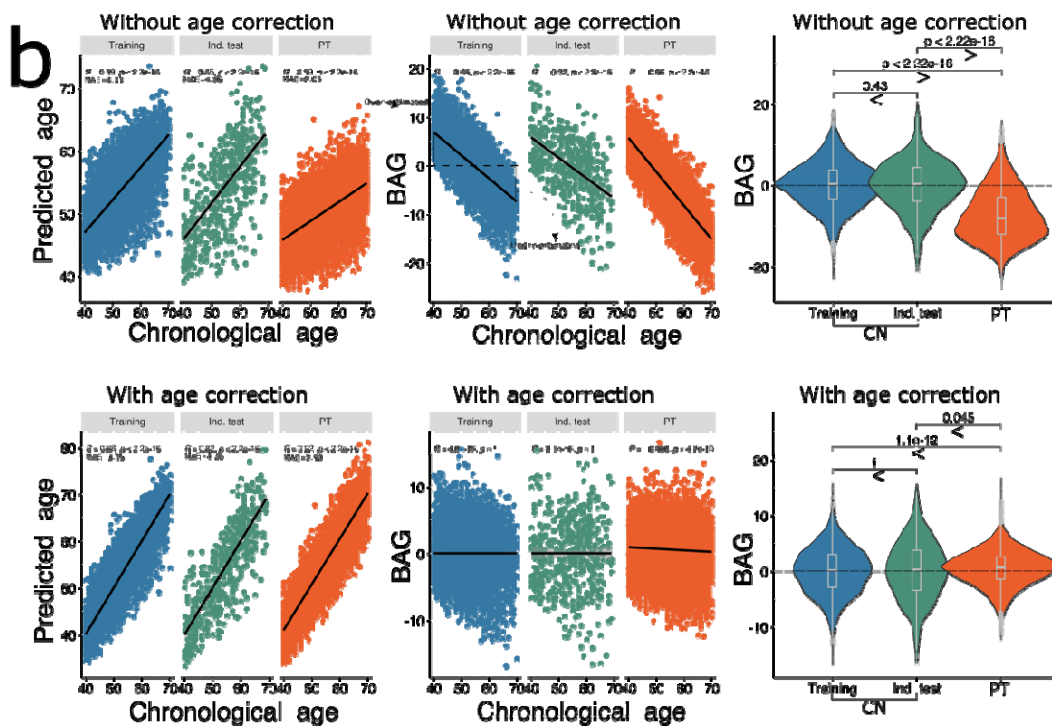
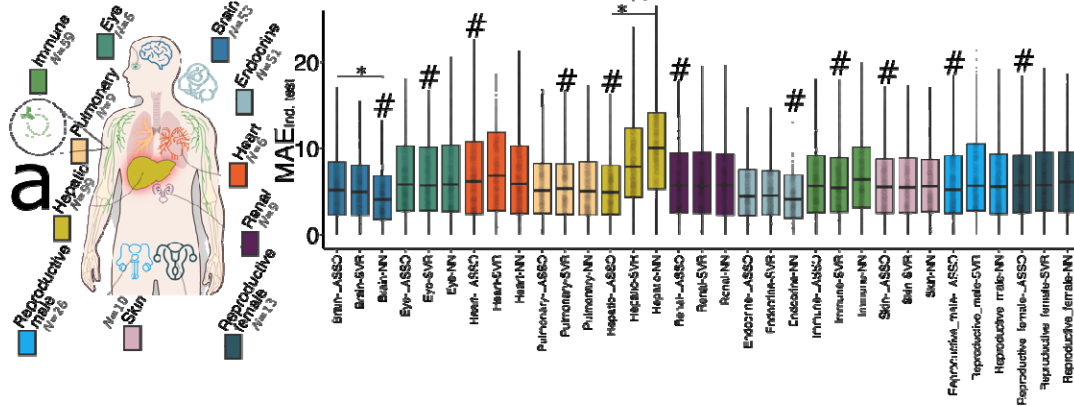


a) We first split the entire proteomics population in the UK Biobank into 5089 healthy control (CN) and 38,049 patient (PT) populations based on the ICD-10 code and other clinical history information. **b)** To derive the 11 ProtBAGs, we trained the three AI/ML models using only the CN training/validation/test population ($N=4589$) with a (nested) cross-validation procedure to select the optimal model. The CN independent test (ind. test; $N=500$) and the PT population ($N=38,409$) were used as independent test datasets. **c)** The analytical workflow of this study

involved deriving 11 ProtBAGs, integrating them with 9 PhenoBAGs, and conducting GWAS and post-GWAS analyses. The ProtBAGs, PhenoBAGs, and their PRSs were then evaluated for their predictive power across 14 systemic disease categories and mortality outcomes.

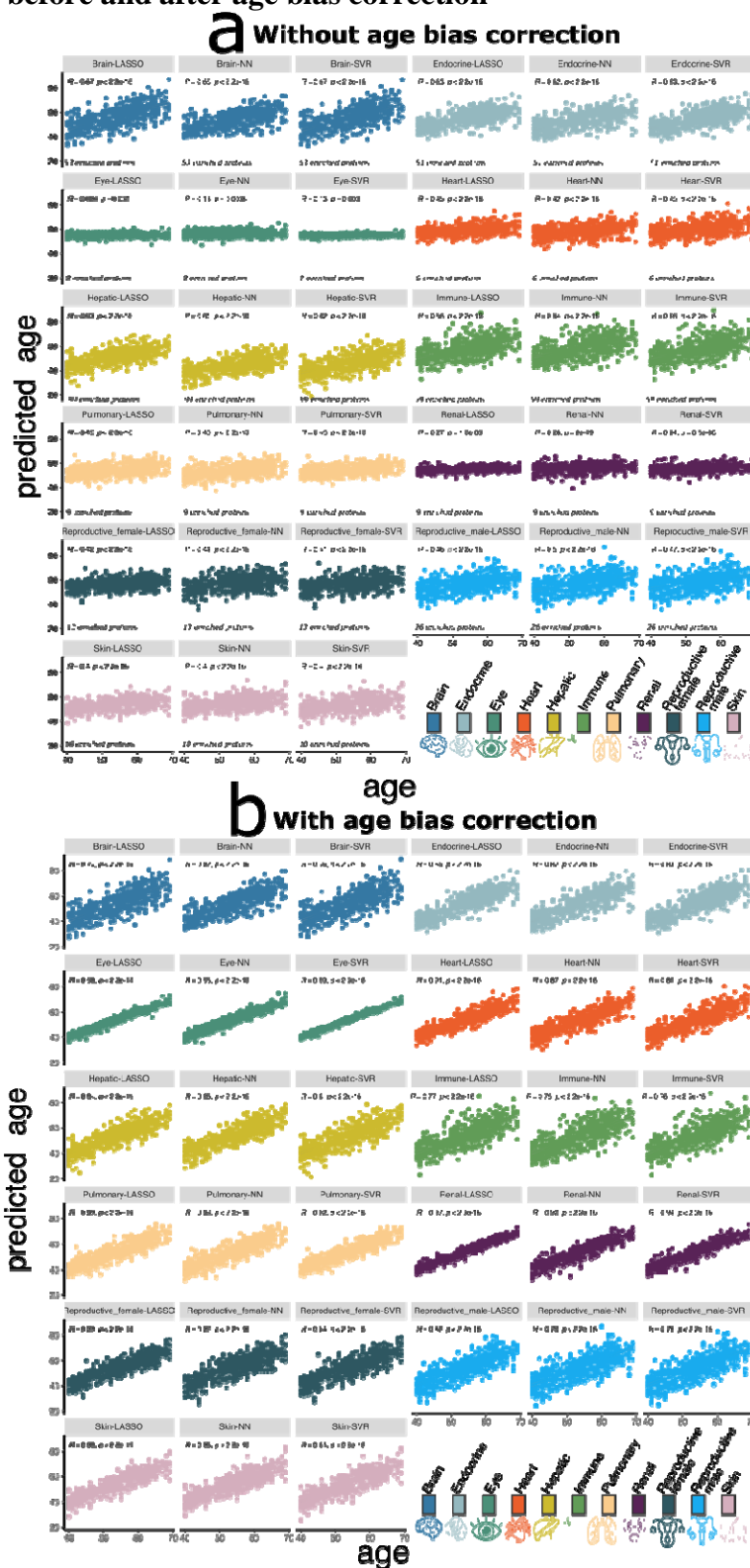
When fitting the organ-specific proteins (**Method 3**) to three AI/ML models [i.e., Lasso regression, support vector regressor (SVR), and neural network (NN)], we observed marginal variability in model performance, with no single model consistently outperforming the others (**Fig. 1a** and **Extended Data Fig. 2**). For instance, the Lasso model outperformed NN and SVR for the hepatic ProtBAG (P-value $< 2.27 \times 10^{-6}$, though the standard t-test may be permissive³² in a complex cross-validation setting). On the other hand, the brain ProtBAG derived from the NN obtained a lower MAE than the Lasso regression and SVR model (P-value $< 2.31 \times 10^{-3}$). Across the different organ systems, the best model performance, before applying the age bias correction⁶, was achieved for the brain ProBAG via the NN (ind. test MAE=4.86; Pearson's $r=0.65$); the highest MAE was achieved for the hepatic ProtBAG via the NN (MAE=10.19; $r=0.61$). Notably, we found instances where MAE and r coefficient were not aligned – a lower MAE (reflecting the magnitude of errors) did not always correspond to a higher r (indicating the strength and direction of predictions), as these metrics capture different aspects of the model performance and can serve as a potential bias-variance tradeoff and the nonlinear dynamics of proteomics aging³³. For example, the hepatic ProtBAG predicted using the NN exhibited a high ($r=0.61$) despite a substantial MAE (MAE=10.19), while the eye ProtBAG using the same model achieved a lower MAE (MAE=6.78) but a much weaker ($r=0.13$). **Supplementary eTable 2** presents detailed statistics for the age prediction tasks before and after the age bias correction⁶. **Extended Data Fig. 2** shows the Pearson's r coefficient between predicted and chronological age. **Supplementary eNote 1** presents the detailed tissue-enriched proteins in each organ to train the 11 multi-organ ProtBAGs in the primary results (**Fig. 1a** and **Method 3c**).

Figure 1: Three AI/ML models to derive the 11 multi-organ ProtBAGs



a) Age prediction performance quantified by the mean absolute error (MAE for the independent test data) across 3 AI models and 11 organ systems using Olink plasma proteomics from UKBB. The Human Protein Atlas project determined the organ-specific proteins (i.e., enriched genes for at least four-fold higher mRNA level in the tissue of interest than other tissues; <https://www.proteinatlas.org/humanproteome/tissue>). The # symbol denotes the model achieving the lowest MAE; the * symbol indicates statistical significance (P -value <0.05) using a two-sample t-test between two models. The dots present the model performance for the 50 repetitions. **b)** Age prediction performance should be reported using metrics before applying age bias correction¹⁴. Age bias correction should be explicitly applied for downstream clinical applications, or age should be at least included as a covariate. Without applying age bias correction, we demonstrated that downstream group comparisons between the healthy control (CN: training/validation and independent test) group and the patient (PT) group could lead to biased conclusions. Abbreviations: Ind. test: independent test; BAG: biological age gap.

Extended Data Figure 2: The scatter plot between AI/ML-predicted biological age and chronological age before and after age bias correction



a) The scatter plot between the AI/ML-derived biological age and chronological age without applying the age bias correction. **b)** The scatter plot between the AI/ML-derived biological age and chronological age after the age bias correction is applied.

Critical considerations for the use of ProtBAG

Ikram recently discussed the use and misuse of biological aging in biomedicine from a clinical perspective³⁴. This study provided additional critical considerations regarding methodology and clinical interpretation in deriving the 11 multi-organ ProtBAGs (**Method 4**).

Age bias correction should be applied

The age bias correction was commonly practiced in the brain imaging-derived age prediction model⁶, leading to a lower MAE and a higher r coefficient (**Fig. 1b**). One consideration in biological age research is reporting metrics before applying age bias correction⁶. Reporting uncorrected metrics ensures consistency in comparing model performance across studies, preventing potential confusion or misapplication from comparing model performance across studies. Additionally, age bias correction is necessary for downstream clinical associations to avoid false conclusions. In our analysis comparing brain ProtBAG between the healthy control (CN) and patient (PT) groups, we found that, without age bias correction¹⁴, the PT group exhibited a lower brain ProtBAG than the CN group (P-value= 2.22×10^{-16}). However, after applying age bias correction, we observed a reversed and more clinically plausible trend, with the PT group showing a higher brain ProtBAG than the CN group (P-value=0.045) (**Fig. 1b**). While including age as a covariate in downstream analyses is standard practice, applying age bias correction remains essential.

Biologically-driven feature selection based on protein organ specificity can alleviate model overfitting

Previous ProtBAG studies have demonstrated that feature selection algorithms can help mitigate model overfitting when applying AI/ML models to unseen test data. For example, Oh et al.⁵ utilized L1 regularization in aggregated Lasso models to address overfitting. Similarly, Argentieri et al.¹⁶ applied the Boruta feature selection algorithm, revealing that the most relevant 204 proteins achieved comparable performance to models trained on the complete set of 2,897 proteins.

Here, we demonstrated the generalizability of AI/ML models to independent test data diminished further when using less organ-specific proteins (e.g., tissue-elevated proteins) compared to a smaller subset of highly organ-specific proteins (e.g., tissue-enriched proteins). **Method 3c** details the definition of different levels of organ specificity. In our experiments, we found that restricting the model to brain tissue-enriched proteins ($N=53$) resulted in better model generalizability from the training/validation/test dataset to the independent test dataset (Cohen's $D=0.15$) than the other two conditions. That is, this discrepancy was significantly larger when models included 146 tissue-enhanced proteins (P-value $<2.22 \times 10^{-16}$; Cohen's $D=1.24$), 255 tissue-elevated proteins (P-value $<2.22 \times 10^{-16}$; Cohen's $D=1.46$), and all the 2448 proteins (P-value $<2.22 \times 10^{-16}$; Cohen's $D=3.52$) (**Fig. 2a**).

Model overfitting can be alleviated by increasing the sample size of the training dataset

Argentieri et al.¹⁶ reported an MAE of 2.24 years and an r of 0.94 in their holdout test data using UKBB data. Our approach differs from Argentieri et al. in several ways. For instance, we used

only 4,589 CN participants for training, whereas Argentieri et al. included a much larger training sample (i.e., 31,808 participants from the general population, including diseased participants). To investigate this, we performed additional analyses to evaluate the effect of training sample size (SS) on model generalizability, using 2,448 proteins as input features. We randomly selected varying SS values (4,589, 10,000, 20,000, 30,000, and 31,808) from the general population to train the model and assessed their generalizability to unseen data. As shown in **Fig. 2b**, increasing the SS improved the model's performance on independent test data, as measured by Cohen's D values.

A tightly-fitted model does not provide higher statistical power to predict cognition than a moderately-fitted model

We underscore that the primary objective of developing ProtBAG, or any biological age biomarker, is not to achieve a highly tightly-fitted model (e.g., a lower MAE), as this can come at the cost of overfitting and reduced power for cross-domain prediction (**Fig. 2c**). Instead, the focus should be on ensuring that the ProtBAGs demonstrate strong statistical associations with cross-domain clinical variables, such as disease diagnoses and cognitive performance. Our experiments observed that the NN model achieved a lower MAE as the number of training epochs increased. However, when assessing the association between the brain ProtBAG and the symbol digit substitution score using a linear regression model, the model at Epoch 2500 ($|\beta|=0.027$) demonstrated a smaller β coefficient compared to the model at Epoch 1000 ($|\beta|=0.035$), albeit this did not achieve statistical significance with a permutation test (P-value=0.34) (**Fig. 2d**).

The demographics of the training dataset are important for model performance and clinical interpretation

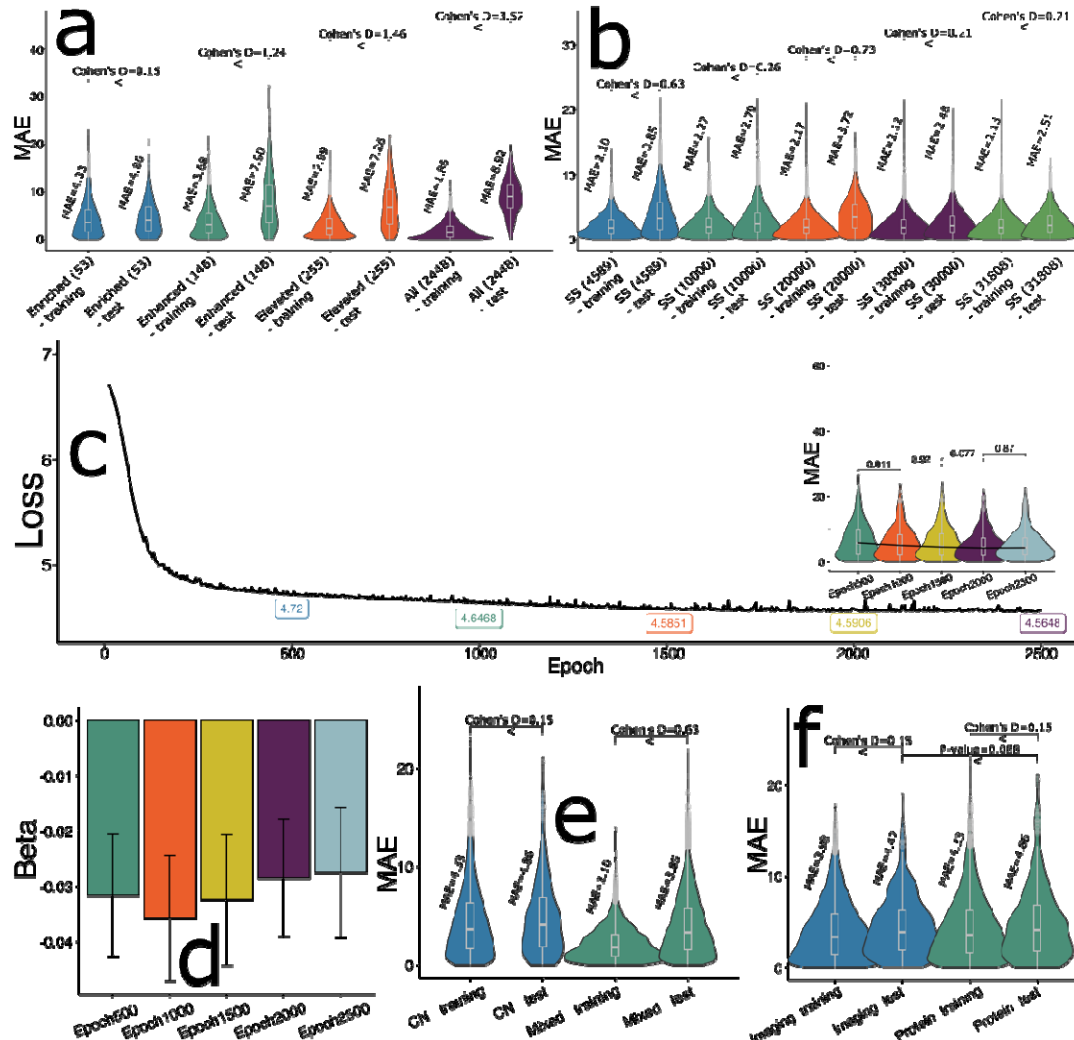
Critically, our AI/ML models (NN for the brain ProtBAG) were trained exclusively on a CN population. Our approach follows the practice in brain neuroimaging-based BAG models, where training is conducted on a healthy population to establish a normative reference for brain aging. This framework allows deviations in the brain PhenoBAG to be associated with pathological factors when the model is applied to external populations with pathologies, facilitating its clinical interpretability.

We conducted a comparative experiment with varied training populations to examine how disease diagnosis influences model performance and generalizability. Models trained on the CN group showed higher MAE and less overfitting, while mixed-population models achieved lower MAE with moderate overfitting. This may be due to increased heterogeneity/variability and extreme features tied to pathology, which risk capturing noise over generalizable signals.

Neuroimaging-derived brain PhenoBAG and brain ProtBAG achieved comparable predictive performance

Finally, we compared the brain PhenoBAG (ind. test MAE=4.47), generated from 119 MRI-derived brain imaging features³, with the brain ProtBAG (ind. test MAE=4.86), constructed using 53 brain tissue-enriched proteins, and found their performance comparable (P-value=0.088) (**Fig. 2f**).

Figure 2: The impact of key components on model performance and generalizability via the brain ProtBAG



a) Different levels of protein organ specificity serve as a means of feature selection, which refers to protein-coding genes with elevated expression levels in a specific tissue or organ, categorized as *i*) tissue-enriched genes, *ii*) tissue-enhanced genes, and *iii*) tissue-elevated genes (<https://www.proteinatlas.org/humanproteome/brain/human+brain>). The results showed that training the neural network (NN) using proteins with lower organ specificity (i.e., incorporating more proteins as features) resulted in poor generalizability ability. **b)** The issue of poor generalizability was alleviated by increasing the training sample size. In this experiment, we expanded the training population to include a mixed cohort encompassing individuals with ICD-based disease diagnoses^{16,5}, rather than restricting it to the CN population, a common practice in the neuroimaging-based brain PhenoBAG. **c)** The loss of the validation dataset for training the NN to predict the chronological age at epochs 500, 1000, 1500, 2000, and 2500. The MAE of the age prediction task at epochs 500, 1000, 1500, 2000, and 2500. **d)** A more “tightly-fitted” model did not result in higher statistical power to predict cognition (i.e., symbol digit substitution) than a “moderately-fitted” model. The β coefficient from the linear regression model associating brain ProtBAG with the cognitive score was evaluated at epochs 500, 1000, 1500, 2000, and 2500. While Epoch 1000 exhibited a trend toward a larger effect size than Epoch 2500, the permutation

test yielded a P-value of 0.34 ($N=10,000$ times). **e**) The brain ProtBAG model trained on a mixed population (comprising both CN and PT) demonstrated a lower MAE compared to the model trained exclusively on the CN population (sample size=4589). **f**) The brain PhenoBAG and ProtBAG models achieved comparable performance using brain imaging and plasma protein features, respectively.

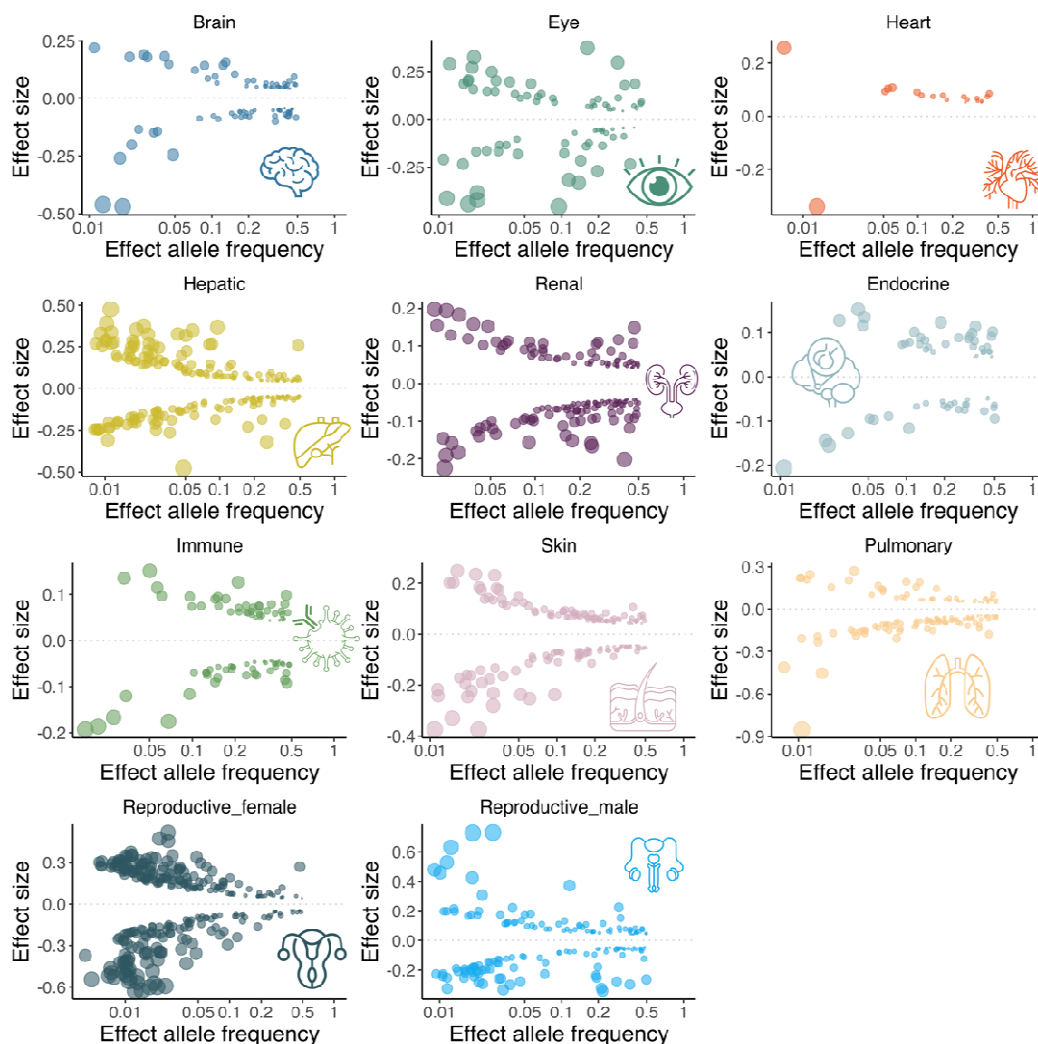
The genetic overlap between ProtBAG and PhenoBAG

We first conducted GWAS for the 11 ProtBAGs to identify shared genomic loci and regions with the 9 PhenoBAGs from our previous study (**Method 5a**).

For the 20 GWASs using European ancestry populations, we identified 129 ($P\text{-value}<5\times 10^{-8}/11$) and 308 ($P\text{-value}<5\times 10^{-8}/9$) genomic locus-BAG pairs for the 11 ProtBAGs and 9 PhenoBAGs, respectively. We denoted the genomic loci using their top lead SNPs defined by FUMA (**Supplementary eNote 2**) considering linkage disequilibrium (LD); the genomic loci are presented in **Supplementary eTable 3**. We visually present the shared genomic loci annotated by cytogenetic regions based on the GRCh37 cytoband (**Fig. 3a**). Manhattan and QQ plots, as well as the genomic inflation factor (λ) of the 11 ProtBAG and 9 PhenoBAG GWASs, are presented in our MEDICINE portal (e.g., hepatic ProtBAG: https://labs-laboratory.com/medicine/hepatic_protbag). The LDSC intercept ($LDSC_b=1.02$ [0.99, 1.03]) of the 11 ProtBAG GWASs was close to 1, indicating no severe population stratification observed. **Extended Data Fig. 3** presents the trumpet plots of the effective allele frequency vs. the β coefficients of the 11 ProtBAG GWASs.

We then computed the pairwise genetic correlation (g_c) and phenotypic correlation (p_c) between the 11 ProtBAGs and 9 PhenoBAGs (**Method 5b**). We observed strong associations between the renal PhenoBAG with multiple ProtBAG at both genetic and phenotypic levels, including the immune ProtBAG ($g_c=0.21$; $p_c=0.33$) and pulmonary ProtBAG ($g_c=0.30$; $p_c=0.28$). Additionally, within-organ associations were not consistently observed; for instance, the eye exhibited neither significant nor phenotypic correlations between the eye PhenoBAG and ProtBAG (**Fig. 3b**). **Supplementary eTable 4** presents detailed statistics on genetic and phenotypic correlations.

Extended Data Figure 3: Trumpet plots of the effect allele frequency vs. the β coefficient of the 11 ProtBAG GWASs



The trumpet plots display the inverse relationship between the alternative (effect) allele frequency and the effect size (β coefficient) for the 11 ProtBAGs. We present the independent significant SNPs defined in FUMA. The dot size corresponds to the effect size, while the transparency of the dot is proportional to its statistical significance.

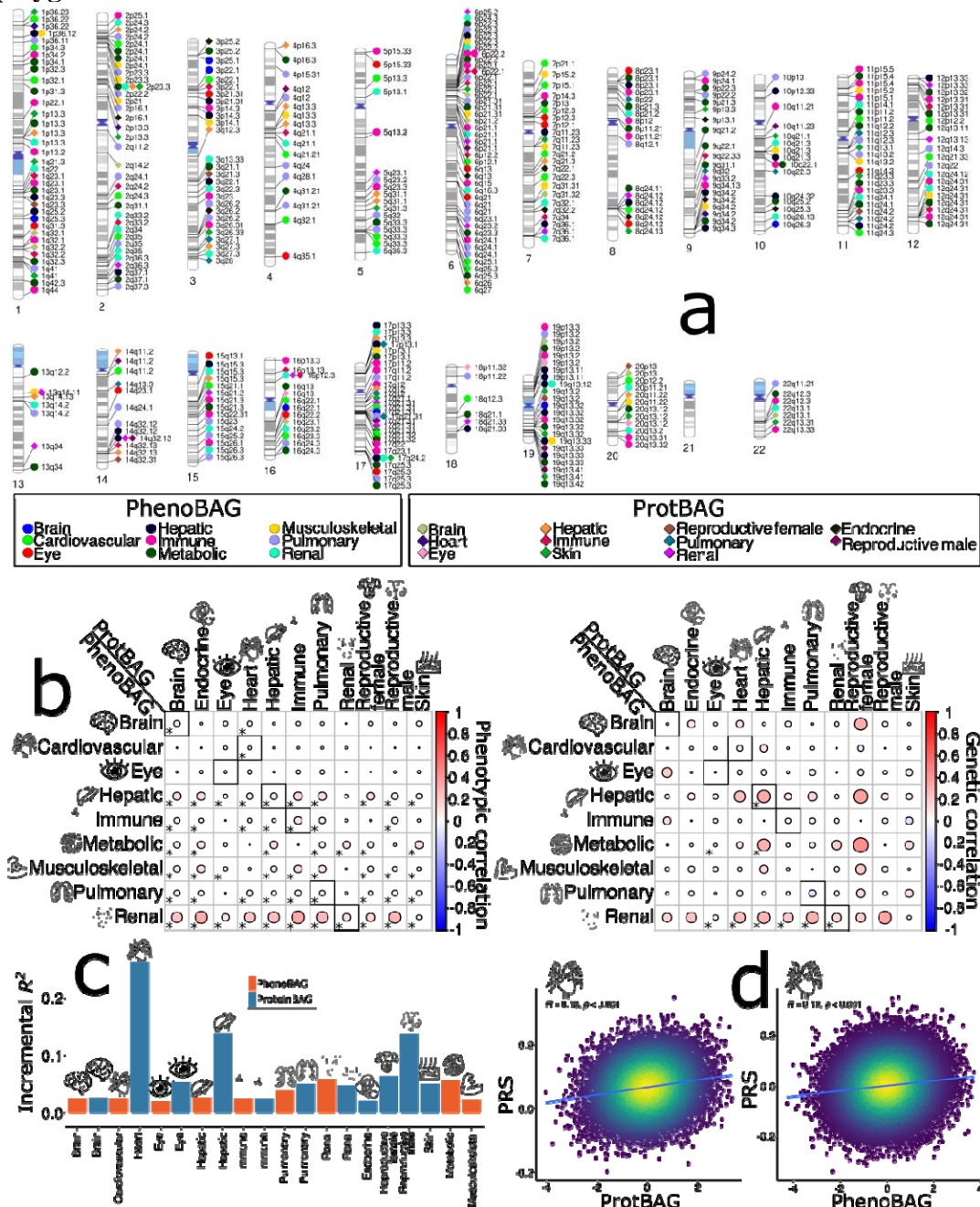
The polygenic risk score of ProtBAG is more predictive than PhenoBAG

We conducted split-sample GWAS to develop the PRS model, using split1 GWAS for training and split2 GWAS for testing, ensuring the two splits had similar age and sex distributions. We evaluated the predictive power of the PRS for the 11 ProtBAG and 9 PhenoBAG by measuring the incremental R^2 gained when predicting the BAG with the PRS as a feature on top of age and sex (**Method 5c**).

All the PRSs demonstrated significant associations with the BAGs ($P\text{-value} < 4.58 \times 10^{-81}$). The 11 ProtBAG-PRSs showed larger predictive power (incremental R^2 ranging from 2.03% to 26.3%) than the 9 PhenoBAG-PRSs (incremental R^2 ranging from 2.01% to 5.91%) when predicting the BAGs (**Fig. 3c**). For instance, the heart ProtBAG exhibited a higher Pearson's correlation coefficient with ProtBAG-PRS ($r=0.18$) compared to the heart PhenoBAG and

PhenoBAG-PRS ($r=0.12$) (**Fig. 3d**). **Supplementary eTable 5** presents detailed statistics of the PRS analyses.

Figure 3: Genetic overlap between PhenoBAG and ProtBAG and the prediction power of their polygenic risk score



a) Cytogenetic regions where the genomic region was jointly linked to PhenoBAG and ProtBAG. Bonferroni correction was applied to denote significant genomic loci associated with PhenoBAG (P -value $< 5 \times 10^{-8}/9$) and ProtBAG (P -value $< 5 \times 10^{-8}/9$). **b**) Phenotypic (pc) and genetic (gc) associations were evaluated between each pair of the 9 PhenoBAGs and 11 ProtBAGs. Statistically significant associations after Bonferroni correction ($0.05/9/11$) are marked with an asterisk (*), and within-organ associations (e.g., between the brain PhenoBAG and ProtBAG) are highlighted with black squares. **c**) The bar plot shows the incremental R^2 (i.e., the R^2 of the alternative model minus that of the null model) for the polygenic risk score (PRS) of each

PhenoBAG and ProtBAG. The PRS was calculated using the split2 target GWAS data, with split1 GWAS data serving as the training set for the PRSCs model. **d)** The scatter plot shows the relationship between the heart ProtBAG, cardiovascular PhenoBAG, and their corresponding PRS, including the P-value and Pearson's r . Notably, the relationship between PRS and PhenoBAG/ProtBAG is likely not linear (although a linear model was fitted), as PRS inherently accounts for only a small proportion of the variance in the phenotypes of interest. GWAS results are publicly disseminated at <https://labs-laboratory.com/medicine/>.

The causal relationship between the 11 ProtBAGs, 9 PhenoBAGs, and 525 DEs

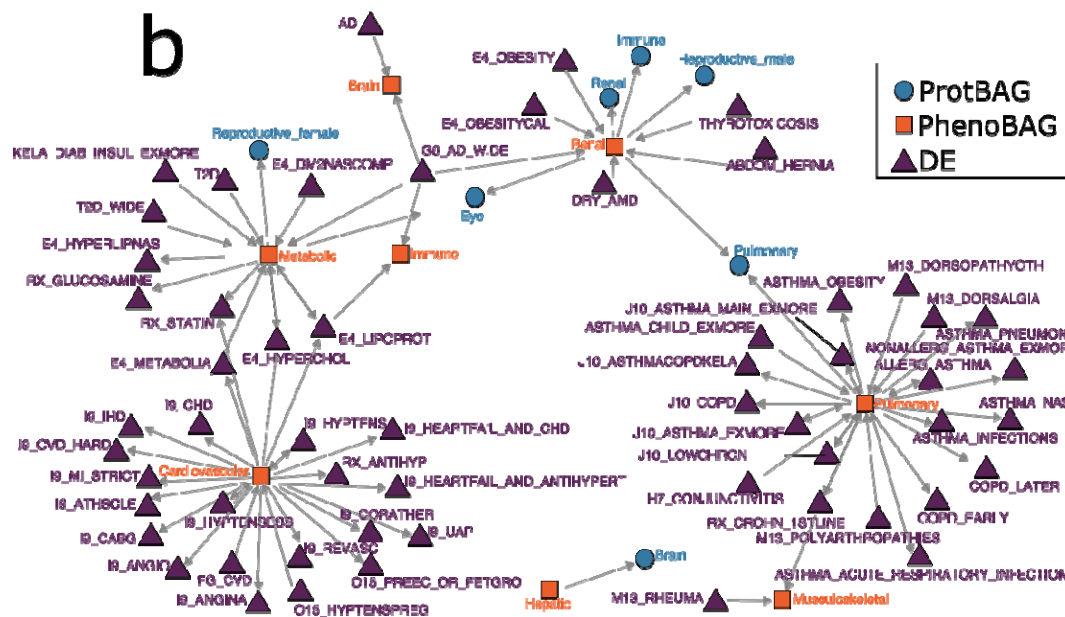
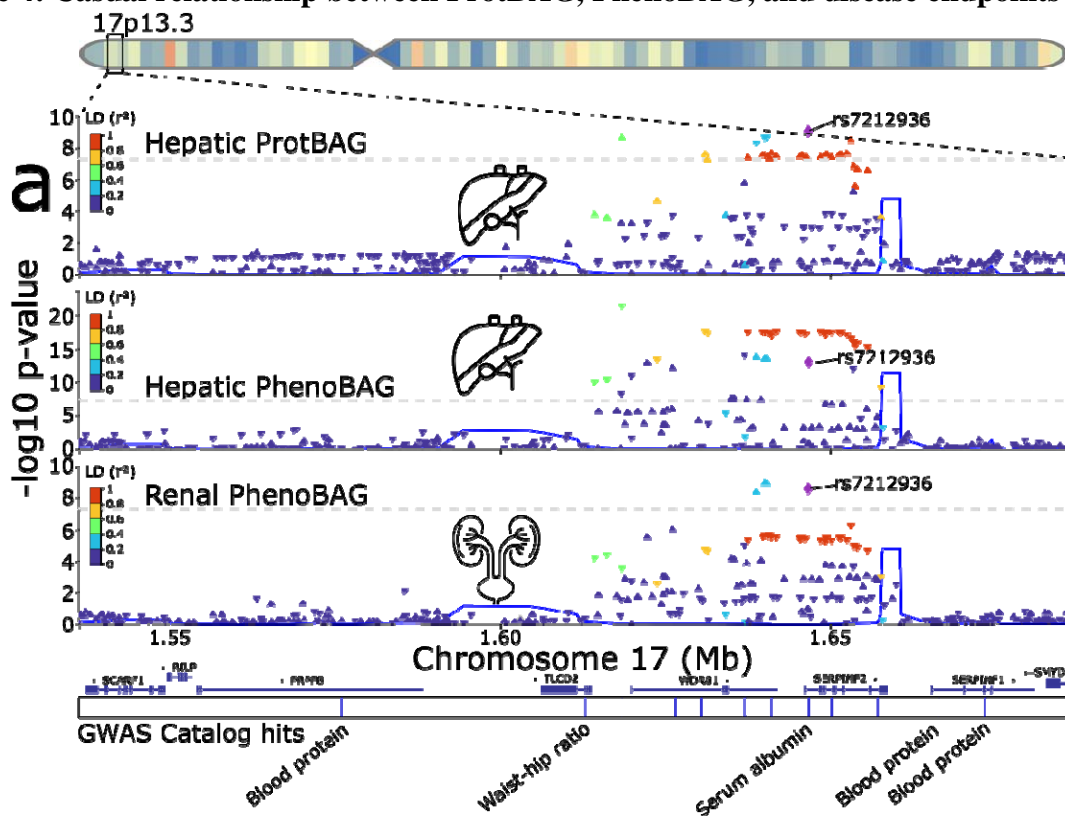
We employed two computational genomics methods to explore the causal relationships among the 11 ProtBAGs, 9 PhenoBAGs, and 525 DEs: *i)* Bayesian colocalization (**Method 5d**) and *ii)* Mendelian randomization (**Method 5e**).

Guided by the strong genetic correlation between the hepatic ProtBAG, hepatic PhenoBAG ($g_c=0.32$), and renal PhenoBAG ($g_c=0.29$), we investigated the shared causal variants between two traits via Approximate Bayes Factor colocalization³⁵ analyses. We demonstrated one genomic locus where the hepatic ProtBAG shared a potential causal variant with both the hepatic PhenoBAG and renal PhenoBAG (**Fig. 4a**). The shared causal variant (rs7212936 at 17p13.3) showed a PP.H4.ABF (Approximate Bayes Factor)=0.99, which examines the posterior probability (PP) to evaluate the hypothesis of a single shared causal variant associated with both traits within this genomic locus. This causal SVN was mapped to the *SERPINF2* gene and had prior links to traits such as serum albumin levels and urate measurements. Other variants within this locus have been connected to various traits, including blood protein levels and waist-to-hip ratio.

Using bi-directional, two-sample Mendelian randomization analyses, we subsequently established a three-layer causal network that linked ProtBAG, PhenoBAG, and DE (**Fig. 4b**). The *ProtBAG2PhenoBAG* network did not show any significant causal signals (P-value<0.05/10 exposure variables). The *PhenoBAG2ProtBAG* network found 9 causal relationships, including from the renal PhenoBAG to the renal ProtBAG [P-value= 4.11×10^{-3} <0.05/11; OR (95% CI)=1.18 (1.05, 1.31); number of IVs=46] and from the hepatic PhenoBAG to the brain ProtBAG [P-value= 3.44×10^{-3} ; OR (95% CI)=1.12 (1.04, 1.21); number of IVs=41]. The *PhenoBAG2DE* network found 41 causal relationships, including from the cardiovascular PhenoBAG to hypertension [FinnGen code: I9_HYPTENS; P-value= 3.00×10^{-7} <0.05/455; OR (95% CI)=1.73 (1.37, 2.17); number of IVs=37] and from the pulmonary PhenoBAG to chronic obstructive pulmonary disease [FinnGen code: J10_COPD; P-value= 1.48×10^{-19} ; OR (95% CI)=1.79 (1.58, 2.03); number of IVs=58]. Finally, for the *DE2PhenoBAG* network, we found 40 causal relationships, including from AD (PGC) to the brain PhenoBAG [P-value= 5.00×10^{-5} <0.05/179; OR (95% CI)=1.06 (1.03, 1.09); number of IVs=20]. This was further strengthened by the causal link from AD (FinnGen code: G6_AD_WIDE) to the brain PhenoBAG [P-value= 3.10×10^{-5} ; OR (95% CI)=1.10 (1.06, 1.14); number of IVs=8], as well as other PhenoBAG (e.g., immune and renal PhenoBAGs) (**Fig. 4b**). We highlighted a causal pathway connecting three layers: obesity→renal PhenoBAG→renal ProtBAG. Obesity (FinnGen code: E4_OBESITY) demonstrated a positive causal relationship with the renal PhenoBAG [P-value= 2.18×10^{-8} ; OR (95% CI)=1.11 (1.07, 1.15); number of IVs=19], which subsequently exerted a causal effect on the renal ProtBAG [P-value= 4.11×10^{-3} ; OR (95% CI)=1.18 (1.05, 1.31); number of IVs=46], among other ProtBAGs (i.e., eye, immune, male reproductive, and pulmonary) (**Fig. 4b**).

Mendelian randomization relies on stringent assumptions that can sometimes be violated. We conducted comprehensive sensitivity analyses for the significant signals identified to scrutinize this. **Extended Data Fig. 4** provides the results of these analyses for the abovementioned causal pathway, with a detailed discussion available in **Supplementary eNote 3**. Detailed statistics for all five estimators are presented in **Supplementary eTable 6**, and the results of the sensitivity analyses are presented in **Supplementary eFolder 1**.

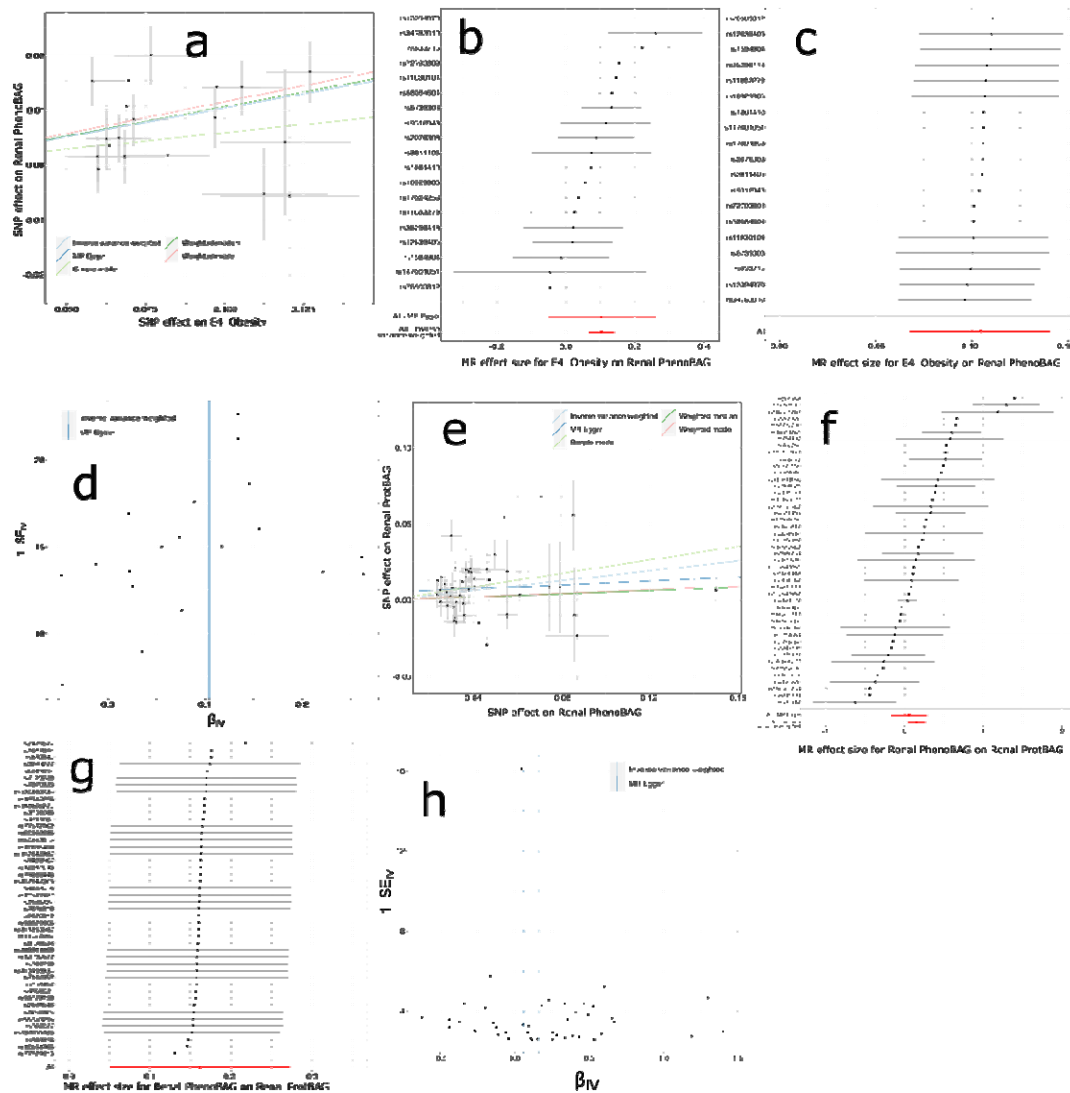
Figure 4: Casual relationship between ProtBAG, PhenoBAG, and disease endpoints



a) Genetic colocalization was evidenced at one locus (17p13.3) between the hepatic ProtBAG, hepatic PhenoBAG, and renal PhenoBAG. The signed PP.H4.ABF (>0.8) denotes the posterior probability (PP) of hypothesis H4, which suggests that both traits share the same causal SNP (rs7212936). Representative GWAS hits are annotated based on previous studies available on the NHGRI-EBI GWAS Catalog. **b)** We constructed a three-layer (ProtBAG-PhenoBAG-DE) causal network by employing bi-directional two-sample Mendelian randomization, following a rigorous

quality control procedure to select exposure and instrumental variables (number of IVs>7), corrected for multiple comparisons (based on either the number of exposure or outcome variables whichever is larger), and performed sensitivity analyses (e.g., horizontal pleiotropy and removing overlap populations) to scrutinize the robustness of our results. Four causal networks were analyzed: i) ProtBAG-to-PhenoBAG, ii) PhenoBAG-to-ProtBAG, iii) PhenoBAG-to-DE, and iv) DE-to-PhenoBAG. Notably, the ProtBAG GWASs ($N>40,000$) were underpowered compared to the PhenoGWASs ($N>11,000$ for body PhenoBAG), providing no evidence of established causality from ProtBAG to PhenoBAG; Instrumental variables were selected via clumping for these genome-wide significant SNPs considering LD. The arrows indicate the direction of the established causal relationship from the exposure variable to the outcome variable. The interactive network visualization is also available at https://labs-laboratory.com/medicine/protbag_mr. Abbreviations: DE: disease endpoint; LD: linkage disequilibrium. It is crucial to approach the interpretation of these potential causal relationships with caution despite our thorough efforts in conducting multiple sensitivity checks to assess any potential violations of underlying assumptions.

Extended Data Figure 4: Sensitivity check analyses for the causal pathway of “obesity→renal PhenoBAG→renal ProtBAG”



a) Scatter plot for the MR effect sizes of the SNP-obesity association (x -axis, log OR) and the SNP-renal PhenoBAG associations (y -axis, log OR) with standard error bars. The slopes of the five lines correspond to the causal effect sizes estimated by the five MR estimators, respectively. **b)** Forest plot for the single-SNP MR results. Each dot represents the MR effect (log OR), and the error bar displays the 95% CI for Obesity on renal PhenoBAG using only one SNP; the red line shows the MR effect using all SNPs together for IVW and MR Egger estimators. **c)** Leave-one-SNP-out analysis of obesity on renal PhenoBAG. Each dot represents the MR effect (log OR), and the error bar displays the 95% CI by excluding that SNP from the analysis. The red line depicts the IVW estimator using all SNPs. **d)** Funnel plot for the relationship between the causal effect of obesity on renal PhenoBAG. Each dot represents MR effect sizes estimated using each SNP as a separate instrument against the inverse of the standard error of the causal estimate. **e)** Scatter plot for the MR effect sizes of the SNP-renal PhenoBAG association (x -axis, log OR) and the SNP-renal ProtBAG associations (y -axis, SD units) with standard error bars. The slopes of the five lines correspond to the causal effect sizes estimated by the five MR estimators, respectively. **f)** Forest plot for the single-SNP MR results. Each dot represents the MR effect (log

OR)), and the error bar displays the 95% CI for renal PhenoBAG on renal ProtBAG using only one SNP; the red line shows the MR effect using all SNPs together for IVW and MR Egger estimators. **g**) Leave-one-SNP-out analysis of renal PhenoBAG on renal ProtBAG. Each dot represents the MR effect (log OR), and the error bar displays the 95% CI by excluding that SNP from the analysis. The red line depicts the IVW estimator using all SNPs. **h**) Funnel plot for the relationship between the causal effect of renal PhenoBAG on renal ProtBAG. Each dot represents MR effect sizes estimated using each SNP as a separate instrument against the inverse of the standard error of the causal estimate.

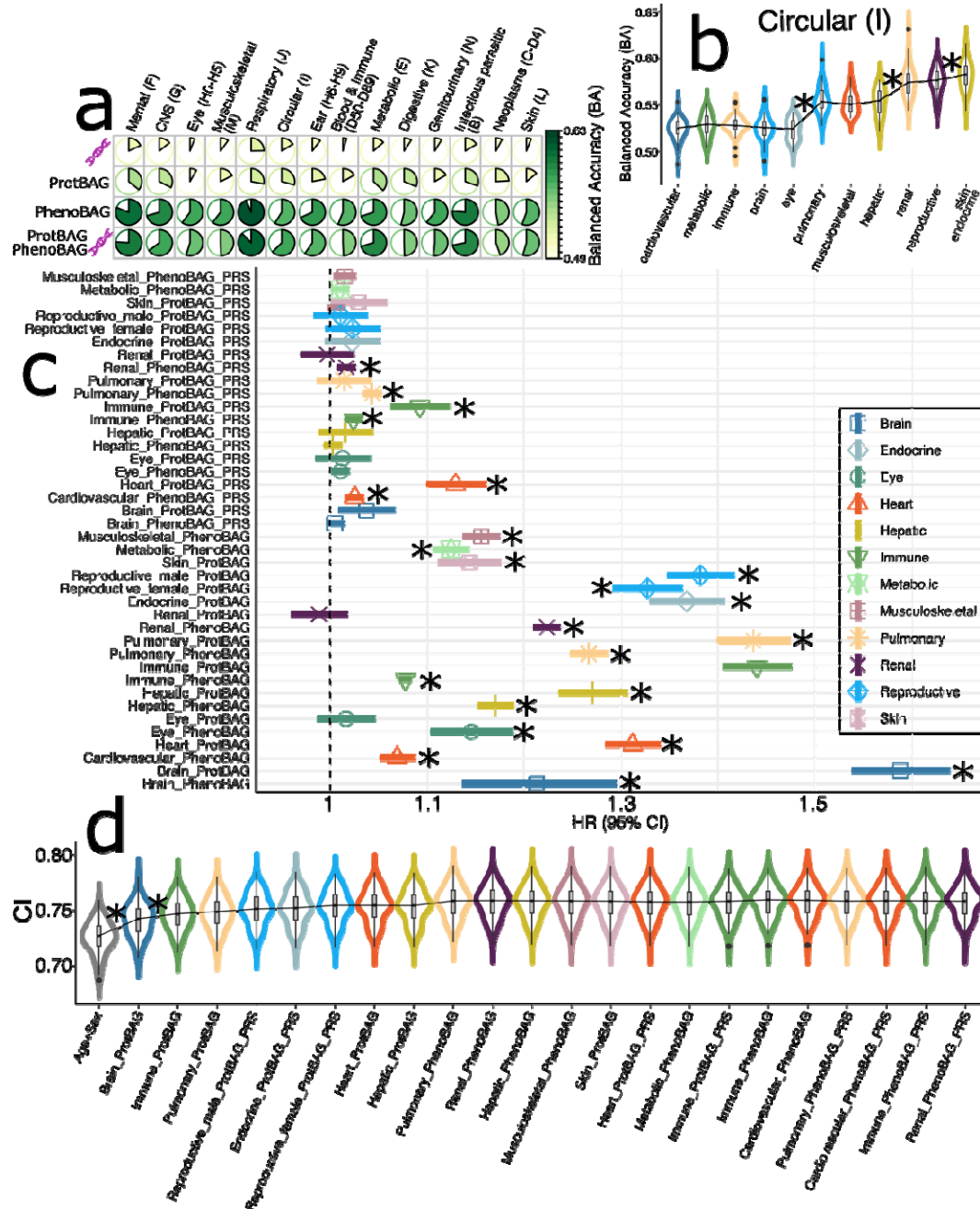
The clinical promise of the 11 ProtBAGs, 9 PhenoBAGs, and 20 PRSs

We demonstrate the clinical promise of the 11 ProtBAGs, 11 ProBAG-PRSs, 9 PhenoBAGs, and 9 PhenoBAG-PRSs in predicting various clinical outcomes through binary classification and survival analysis: *i*) the classification of 14 systemic disease categories and *ii*) the risk of mortality (**Method 6a-b**).

We assessed the prediction ability of support vector machines (SVM) at the individual level to classify the 14 disease categories (**Method 6a**). The highest performance was observed for the respiratory disease category (ICD-codes: J; balanced accuracy (BA)=0.62). The PRS and ProtBAG individually exhibited lower predictive accuracy for disease categories than PhenoBAG. Furthermore, combining all three feature sets failed to outperform the PhenoBAG alone (**Fig. 5a**). Adding age and sex enhanced the classification accuracy (**Supplementary eFigure 1**). Furthermore, we used the circulatory system disease categories (ICD code: I) as an example (**Fig. 5b**) and demonstrated that adding cross-organ features can improve classification performance. The full evaluation metrics of the cross-validated results are presented in **Supplementary eTable 7**.

We also used the 40 BAGs to predict mortality risk using UKBB data (**Method 6b**). Our analysis revealed that 24 BAGs, including ProtBAGs, PhenoBAGs, and their PRSs, showed significant associations ($P\text{-value} < 0.05/9/11$) with mortality. The brain ProtBAG showed the highest mortality risks [HR (95% CI)=1.58 (1.54, 1.63); $P\text{-value}=7.09 \times 10^{-176}$], followed by the immune ProtBAG [HR (95% CI)=1.44 (1.40, 1.48); $P\text{-value}=3.07 \times 10^{-181}$], and pulmonary ProtBAG [HR (95% CI)=1.43 (1.40, 1.47); $P\text{-value}=1.98 \times 10^{-156}$]. Among the 9 PhenoBAGs, the renal PhenoBAG [HR (95% CI)=1.22 (1.21, 1.24); $P\text{-value}=1.85 \times 10^{-252}$] and brain PhenoBAG [HR (95% CI)=1.21 (1.14, 1.30); $P\text{-value}=8.63 \times 10^{-9}$] showed the highest risks. For the 20 PRSs, the highest mortality risk was achieved with the heart ProtBAG-PRS [HR (95% CI)=1.13 (1.10, 1.16); $P\text{-value}=1.99 \times 10^{-18}$] (**Fig. 5c**). Given the population differences among ProtBAGs, PhenoBAGs, and PRSs, comparing hazard ratios (HR) directly is not advisable, as variations in baseline hazard could affect the interpretation. We conducted a cumulative prediction analysis based on the substantial associations identified in the 22 significant BAGs (excluding the brain and eye PhenoBAGs due to their limited sample sizes). This analysis demonstrated that combining these features provided additional predictive power beyond age and sex, achieving an average concordance index of 0.76 ± 0.014 (**Fig. 5d**). The brain and immune ProtBAGs contributed most significantly to this improvement. Comprehensive statistics, including HRs, P -values, and sample sizes, are available in **Supplementary eTable 8**.

Figure 5: ProtBAG, PhenoBAG, and their PRS predict systemic disease categories and mortality



a) The classification balanced accuracy (BA) for 14 ICD-based disease categories was evaluated using PRS, ProtBAG, and PhenoBAG as features within a support vector machine (SVM) framework employing a nested cross-validation (CV) approach (training/validation/test datasets). Balanced accuracy results from the CV are presented, with additional metrics provided in the Supplement. Overall, PhenoBAG demonstrated greater predictive power than other omics data, and simply combining ProtBAG, PhenoBAG, and PRS did not enhance classification performance. The brain and eye PhenoBAG were excluded because merging them with the populations of other features resulted in a very small sample size ($N < 1000$). **b)** The cumulative inclusion of organ-specific features enhanced classification performance in predicting circulatory system diseases (ICD code: I). The * symbol indicates statistical significance (< 0.05) from a two-sample t-test comparing CV test accuracy between two SVM models; however, a standard t-

test is liberal³² and should be interpreted cautiously. **c)** ProtBAG, PhenoBAG, and their PRS show significant associations with the risk of mortality. Age and sex were included as covariates in the Cox proportional hazard model. The symbol * indicates significant results that survived the Bonferroni correction ($<0.05/9/11$). It is important to note that the population sample sizes for ProtBAG and PhenoBAG differ, making their HRs not directly comparable. **d)** The significant ProtBAG, PhenoBAG, and PRS were cumulatively included as features for mortality risk prediction. The * symbol indicates statistical significance (<0.05) from a two-sample t-test comparing results between two Cox models. HR: hazard ratio; CI: concordance index.

Discussion

This study systematically benchmarks the age prediction performance across 11 multi-organ ProtBAGs, revealing insights into the factors influencing model performance and generalizability to unseen data. Inspired by common practices in brain age research⁶, we introduced critical methodological considerations to enhance rigor and clinical interpretability in multi-organ aging research. Subsequently, we comprehensively compared the genetic overlap between the 11 multi-organ ProtBAGs and the 9 PhenoBAGs. By constructing a three-layer causal network, we connected genetics, proteomics, imaging/phenotypic endophenotypes, and disease outcomes, providing an integrative framework for understanding these complex interactions. Finally, we delivered compelling evidence of the clinical potential of the ProtBAGs, PhenoBAGs, and their PRSs in predicting disease categories and mortality, positioning these biomarkers as powerful tools for translational medicine.

Reproducible and systematic evaluation of ProtBAG generation

We addressed several critical considerations for developing and applying ProtBAG. First, we emphasized the importance of age bias correction, a technique that enhances the clinical relevance of ProtBAG models. In neuroimaging-based brain age research, age bias correction has been extensively investigated^{13,14,6,15}. We provided specific scenarios using proteomics data to emphasize the importance of practicing this in ProtBAG. For instance, Oh et al.⁵ and Argentieri et al.¹⁶ did not explicitly correct this bias, although they included age as a covariate in their downstream association analyses.

Our findings also demonstrated the significance of biologically-driven feature selection in alleviating overfitting. Focusing on organ-specific proteins, such as brain tissue-enriched proteins, we achieved better generalizability to unseen data than models using broader, less specific protein sets. Methodologically-driven feature selection algorithms, such as the Boruta algorithm used by Argentieri et al.¹⁶, offer valuable tools for refining predictive models. However, several critical considerations must be addressed. First, complex feature selection should be incorporated within the (nested) cross-validation framework to prevent potential "data leakage," as highlighted in prior research on AD classification³⁶. Second, integrating feature selection within cross-validation can complicate the application of trained models to unseen data, as the features selected may vary across different folds. Moreover, increasing the training sample size reduced overfitting, emphasizing the importance of large and diverse training populations for enhancing model performance. However, diseased populations may obscure clinical interpretation, and increased data heterogeneity remains a critical area for further investigation³⁷. In addition, we noted that a tighter model fit, reflected in lower MAE, does not necessarily equate to stronger clinical associations, as shown in our analysis of cognitive prediction using the brain ProtBAG. This observation aligns with findings from a previous study that reported similar results using neuroimaging-derived brain age models³⁸.

The genetic overlap and associations between the 11 ProtBAGs, 9 PhenoBAGs, and 525 DEs

Our findings underscore the substantial genetic overlap between ProtBAGs and PhenoBAGs, offering perspectives on the shared and distinct genetic architectures underlying proteomics-driven and phenotypic aging profiles. The identification of hundreds of significant genomic loci linked to these BAGs, along with strong cross-omics and cross-organ genetic correlations, emphasizes the interconnected nature of systemic and organ-specific processes in aging^{1,2,5,4}.

Notably, the observed associations, such as those between the renal PhenoBAG and immune and pulmonary ProtBAGs, suggest the existence of genetic networks that transcend traditional organ boundaries. Our previous research^{3,1} explored the genetic overlap across organs among the 9 PhenoBAGs. Building on that foundation, the current study expands this scope by integrating 11 ProtBAGs with cross-omics data spanning multiple organs, offering a comprehensive multi-scale framework for understanding human aging and disease.

The superior predictive performance of ProtBAG-PRSs compared to PhenoBAG-PRSs underscores the potential of proteomics-based approaches to advance precision medicine in genetic aging research^{10,39,11,40}. The observed differences suggest that ProtBAG may capture distinct genetic signals with stronger biological relevance. This supports the growing recognition of proteomics as a critical component in aging studies, offering deeper insights into novel biomarkers and pathways that may remain elusive through traditional phenotypic analyses. Since proteomics is more closely linked to the underlying genetics and etiology of aging, it offers a valuable molecular layer for studying human aging.

Causal inference analyses provided further insights into the intricate relationships between BAGs and DEs. The colocalization signal of a shared causal variant in the hepatic and renal BAGs exemplifies how integrating proteomic and phenotypic dimensions can uncover biologically relevant loci with translational potential. Similarly, the causal pathway linking obesity, renal PhenoBAGs, and renal ProtBAGs highlights the systemic impact of metabolic factors on organ-specific aging processes. These findings emphasize our understanding of how systemic and organ-specific factors drive age-related phenotypes and diseases.

In summary, we demonstrated the value of integrative analyses for BAGs for uncovering the genetic and causal underpinnings of aging across multiple scales. Expanding sample sizes and incorporating diverse ancestries will be critical to enhancing the generalizability of these findings. In addition, exploring the functional consequences of shared loci and causal pathways may provide actionable insights for therapeutic interventions targeting age-related conditions⁴¹.

The prediction power of the 11 ProtBAGs, 9 PhenoBAGs, and their PRSs

The observed differences in predictive power for systemic disease categories between PhenoBAG, ProtBAG, and PRS can be attributed to the nature of the data we integrate and how they relate to disease categories versus mortality outcomes. For disease category prediction, PhenoBAG, which incorporates phenotypic traits directly linked to specific diseases, is likely more predictive because these traits often represent the clinical manifestation of disease, offering immediate and tangible insights into disease risk. Clinical features such as biomarkers, imaging data, and medical history are more directly associated with disease effects, which makes phenotypic data more informative for predicting disease outcomes. In contrast, PRS, based on genetic predisposition, and ProtBAGs, which rely on proteomic data, may not effectively capture disease-specific features. In particular, the current study focused exclusively on common genetic variants, excluding rare ones typically associated with larger effect sizes⁴². These omics layers provide broader insights into genetic risk and molecular pathways, but their relationships to specific disease categories may be more complex and indirect, making them less predictive for disease classification. Similarly, a recent study showed that multi-omics data and biomarkers can be effectively integrated to outperform PRS in disease predictions⁴³.

For mortality prediction, however, ProtBAG and PhenoBAGs show strong predictive power. This is likely because a complex interplay of molecular and clinical factors influences mortality. ProtBAG, which captures proteomic profiles, offers a more direct measure of the

molecular processes that underlie aging and disease, such as inflammation, cellular stress, and metabolic dysfunction. These processes are key contributors to mortality, especially in aging populations^{16,44}. PhenoBAG, incorporating clinical traits, also reflects the cumulative effects of health deterioration and is strongly correlated with mortality outcomes¹³. PRS, while valuable for predicting genetic susceptibility, may not fully capture the dynamic and multifactorial nature of mortality risk, which involves genetic predisposition, lifestyle factors, physiological markers, and environmental factors⁴⁵.

Interestingly, combining multi-omics BAGs did not significantly improve disease prediction, suggesting that integrating multiple omic layers does not necessarily lead to enhanced performance for disease categories. This may be because disease prediction requires biomarkers specifically relevant to each disease or the broad category, and the multi-omics approach may still lack the necessary disease-specific biomarkers²⁹. However, when predicting mortality, the multi-organ BAGs and PRS improved prediction, highlighting the importance of integrating different biological layers across multiple organs. Mortality is a more complex outcome that involves systemic processes across the entire body, making multi-organ and multi-omic approaches more effective. This suggests that combining various molecular layers across organs/omics for comprehensive risk prediction is crucial for capturing the full spectrum of biological processes that influence aging and mortality.

Outlook

This study investigates several pivotal aspects of biological age research. Future research should expand on this foundation by integrating epigenetic, transcriptomic, and metabolomic data. This will enrich the causal pathways from genetics to disease outcomes, providing a more holistic view of human aging and disease^{46,47}.

Methods

Method 1: The MULTI consortium

The MULTI consortium is an ongoing initiative to integrate and consolidate multi-organ data (e.g., brain and heart MRI and eye OCT) with multi-omics data, including imaging, genetics, and proteomics. Building on existing consortia and studies, MULTI aims to curate and harmonize the data to model human aging and disease across the lifespan. This study used individual-level and summary-level multi-omics data from UKBB, FinnGen, and PGC to derive the multi-omics and multi-organ BAGs. **Supplementary eTable1** details the sample characteristics.

UK Biobank

UKBB⁴⁸ is a population-based research initiative comprising around 500,000 individuals from the United Kingdom between 2006 and 2010. Ethical approval for the UKBB study has been secured, and information about the ethics committee can be found here: <https://www.ukbiobank.ac.uk/learn-more-about-uk-biobank/governance/ethics-advisory-committee>. This study used brain MRI, eye OCT, and clinical phenotypes (e.g., physiological and physical biomarkers) to derive the 9 PhenoBAGs; our previous studies^{1,2} detailed the generation and the phenotypes used for each organ-specific PhenoBAGs. The 11 ProtBAGs were derived from 2448 plasma proteomics data derived from the Olink platform. Imputed genotype data covering the populations of ProtBAG and PhenoBAG were used for all genetic analyses.

FinnGen

The FinnGen³⁰ study is a large-scale genomics initiative that has analyzed over 500,000 Finnish biobank samples and correlated genetic variation with health data to understand disease mechanisms and predispositions. The project is a collaboration between research organizations and biobanks within Finland and international industry partners. For the benefit of research, FinnGen generously made their GWAS findings accessible to the wider scientific community (https://www.finnngen.fi/en/access_results). This research utilized the publicly released GWAS summary statistics (version R9), which became available on May 11, 2022, after harmonization by the consortium. No individual data were used in the current study.

FinnGen published the R9 version of GWAS summary statistics via REGENIE software (v2.2.4)⁴⁹, covering 2272 DEs, including 2269 binary traits and 3 quantitative traits. The GWAS model encompassed covariates like age, sex, the initial 10 genetic principal components, and the genotyping batch. Genotype imputation was referenced on the population-specific SISu v4.0 panel. We included GWAS summary statistics for 521 FinnGen DEs in our analyses.

Psychiatric Genomics Consortium

PGC³¹ is an international collaboration of researchers studying the genetic basis of psychiatric disorders. PGC aims to identify and understand the genetic factors contributing to various psychiatric disorders such as schizophrenia, bipolar disorder, major depressive disorder, and others. The GWAS summary statistics were acquired from the PGC website (<https://pgc.unc.edu/for-researchers/download-results/>), underwent quality checks, and were harmonized to ensure seamless integration into our analysis. No individual data were used from PGC. Each study detailed its specific GWAS models and methodologies, and the consortium consolidated the release of GWAS summary statistics derived from individual studies. In the current study, we included summary data for 4 brain diseases for which allele frequencies were present.

Method 2: Phenotype analyses to derive the 9 PhenoBAGs

We derived the 9 PhenoBAGs in our previous study^{1,2}, and we also present the final included phenotypes for the 9 human organs in **Supplementary eTable 1**. In summary, we selected brain MRI, physical and physiological measures indicative of key organ systems' function, structure, or general health, including the brain (e.g., brain volume), cardiovascular (e.g., pulse rate), pulmonary (e.g., peak expiratory flow), musculoskeletal (e.g., BMI), immune (e.g., leukocytes), renal (e.g., glomerular filtration), hepatic (e.g., albumin), and metabolic systems (e.g., lipid). Data were processed to ensure reliability: averages were calculated for bilateral measures (e.g., handgrip strength), repeated tests (e.g., blood pressure), or the best performance from multiple attempts (e.g., lung function via spirometry). The eye PhenoBAG was subsequently derived in our follow-up study using eye OCT data¹.

To derive the 9 PhenoBAGs, we used a linear support vector regressor (SVR) and fit the organ-specific phenotypes as features with a 20-fold cross-validation procedure. Optimization of the SVR's hyperparameters (box constraint, kernel function, and ϵ) did not substantially improve performance. Critically, the SVR models were trained exclusively on healthy individuals, defined as those without self-reported or healthcare-documented lifetime chronic medical conditions. This approach supports the clinical interpretation of the trained models when applied to disease groups, with deviations in these PhenoBAGs presumed to reflect specific pathological factors.

Method 3: Proteomics analyses to derive the 11 ProtBAGs

(a) Additional quality checks: We downloaded the original data (Category code: 1838), which were analyzed and made available to the community by the UKB-PPP⁵⁰. The initial quality check was detailed in the original work⁵¹; we performed additional quality check steps as below. We focused our analysis on the first instance of the proteomics data ("instance"=0). Subsequently, we merged the Olink files containing coding information, batch numbers, assay details, and limit of detection (LOD) data (Category ID: 1839) to match the ID of the proteomics dataset. We eliminated Normalized Protein eXpression (NPX) values below the protein-specific LOD. Furthermore, we restricted our analysis to proteins with sample sizes exceeding 10,000. This resulted in 2448 proteins in 43,498 participants.

(b) Missing protein NPX imputation: We observed a substantial missing rate for the 2448 proteins (1229 proteins with > 10% missing values), which made it challenging to employ downstream AI/ML models for age prediction because many of these models do not directly handle missing features. We used the AutoComplete⁵² deep learning algorithm to impute the missing proteins to overcome this. In the original paper, the authors have thoroughly evaluated the impact of the missing rate on the imputation accuracy. Here, we followed the same approach proposed in the paper, assessed the impact of the probability of an individual's being masked during training, and found that this impact is minimal for the imputation accuracy (**Supplementary eFigure 2**). We observed a mean R^2 value of 0.45 between the imputed values and the ground truth for the 2448 proteins, showing improved model performance compared to the original study on cardiometabolic and psychiatric phenotypes ($0.14 < R^2 < 0.30$).

(c) Organ-specific profiles of the 2448 plasma proteins: We used the Human Protein Atlas (HPA) project (<https://www.proteinatlas.org/humanproteome/tissue>) to profile the over-

expression of a specific protein at both RNA-seq and protein levels. HPA highlights the expression profiles of genes in human tissues at both the mRNA and protein levels. Protein expression data from 44 normal human tissue types were obtained through antibody-based protein profiling using conventional and multiplex immunohistochemistry. Accompanying the resource are annotated protein expression levels and all images of immunohistochemically stained tissues. Protein data encompass 15,302 genes (76%) with available antibodies. Additionally, mRNA expression data were generated through RNA sequencing (RNA-seq) of 40 different normal tissue types. In our primary analyses to derive the 11 organ-specific ProtBAGs, we considered defining whether a protein is over-expressed in a particular organ/tissue using the following criterion: tissue-enriched genes are characterized by mRNA expression levels at least four times higher in the tissue or organ of interest compared to all other tissues. This approach aligns with the definition employed by Oh et al.⁵, which relied solely on data from the Genotype-Tissue Expression (GTEx) project. In contrast, the Human Protein Atlas (HPA) integrates resources from multiple consortia, extending beyond GTEx data.

An important yet unexplored question is the relative lack of organ specificity in plasma proteins circulating throughout the human body compared to clinical phenotypes, such as brain MRI features. Many proteins are frequently over-expressed across multiple tissues or organs, akin to the pleiotropic effects observed in genetics. This observation is biologically plausible, as proteomics is more closely linked to underlying genetic mechanisms, whereas clinical phenotypes are more directly associated with disease endpoints. Additionally, Argentieri et al.¹⁶ demonstrated through feature selection that 204 out of 2,897 proteins from the UKBB Olink platform could accurately predict chronological age. However, the impact of protein organ-specificity definitions on model overfitting remains an unresolved question. In this study, we explored this issue by systematically relaxing the organ-specific profiles of proteins under three distinct scenarios. Using the brain ProtBAG as an example, we assessed how varying the number of proteins included in the training of AI/ML models influences performance and overfitting phenomena.

- **Tissue-enriched genes/proteins:** At least four-fold higher mRNA level in the tissue of interest than in other tissues ($N=53$ proteins).
- **Tissue-enhanced genes/proteins:** At least four-fold higher mRNA level in the tissue of interest compared to the average level in all other tissues ($N=146$ proteins).
- **Tissue-elevated genes/proteins:** tissue-enriched genes (including group-enriched genes) and tissue-enhanced genes ($N=255$ proteins).

(d) Three AI/ML models: We systematically benchmarked age prediction performance using 4 AI/ML models on multi-modal brain MRI features in our previous study³. Using the same methodology, we assessed the performance of models in deriving the 11 ProtBAGs using two linear approaches (Lasso regression and SVR) and one non-linear method (neural networks). For the linear models, hyperparameter selection (e.g., the C parameter for SVR) was conducted through nested, repeated hold-out cross-validation¹⁷ with 50 repetitions (80% training/validation and 20% testing). Nested cross-validation was not applied to the neural network due to the impracticality of exhaustively exploring hyperparameter combinations.

(e) Population selections: To rigorously train the AI/ML models, we have split the CN data ($N=5089$) into the following datasets (**Supplementary eFigure 1** and **eTable 1**):

- **CN independent test dataset:** 500 participants were randomly drawn from the CN population;
 - **CN training/validation dataset:** 80% of the remaining 4589 CN were used for the inner loop 10-fold CV for hyperparameter selection;
 - **CN cross-validation test dataset:** 20% of the remaining 4589 CN were used for the outer loop 50 repetitions;
 - **PT dataset:** 38,409 participants that have at least one ICD-10-based diagnosis.
- Model evaluation metrics included mean absolute error (MAE) and Pearson's r .

Importantly, consistent with our prior studies, only healthy control participants were included in the training/validation dataset, while individuals with any disease diagnosis were reserved for the independent test dataset.

Method 4: Influence of key components in deriving the brain ProtBAG

We systematically evaluated key factors influencing model performance using the brain ProtBAG as a case study. These factors included *i*) the choice of AI/ML models (i.e., SVR, Lasso, and neural networks), *ii*) the impact of age bias correction on downstream clinical applications, such as group differences between CN and PT groups, *iii*) the effect of protein organ specificity on model overfitting, comparing enriched, enhanced, and elevated gene categories, *iv*) the influence of model fitting tightness on cross-domain prediction, particularly associations with cognitive outcomes at various epochs (i.e., 500, 1000, 1500, 2000, and 2500 epochs), and *v*) the impact of feature type on model performance, comparing brain imaging-derived features with brain over-expressed plasma proteins. These analyses provide practical guidance for using plasma proteins to develop ProtBAGs while enhancing clinical interpretability and methodological rigor.

Method 5: Genetic analyses

We used the imputed genotype data for all genetic analyses. Our quality check pipeline focused on European ancestry in UKBB (6,477,810 SNPs passing quality checks), and the quality-checked genetic data were merged with respective organ-specific populations for GWAS. We summarize our genetic quality check steps. First, we skipped the step for family relationship inference⁵³ because the linear mixed model via fastGWA⁵⁴ inherently addresses population stratification, encompassing additional cryptic population stratification factors. We then removed duplicated variants from all 22 autosomal chromosomes. Individuals whose genetically identified sex did not match their self-acknowledged sex were removed. Other excluding criteria were: *i*) individuals with more than 3% of missing genotypes; *ii*) variants with minor allele frequency (MAF; dosage mode) of less than 1%; *iii*) variants with larger than 3% missing genotyping rate; *iv*) variants that failed the Hardy-Weinberg test at 1×10^{-10} . To further adjust for population stratification,⁵⁵ we derived the first 40 genetic principle components using the FlashPCA software⁵⁶. Details of the genetic quality check protocol are described elsewhere^{57,3,1,58,59}.

(a) GWAS:

We applied a linear mixed model regression to the European ancestry populations using fastGWA⁵⁴ implemented in GCTA⁶⁰.

PhenoBAG GWAS: In our initial investigation, we conducted GWAS for the 9 PhenoBAGs using a linear model in PLINK, with fastGWA employed as a sensitivity analysis. For

consistency with the ProtBAG GWASs in this study, we used the fastGWA summary statistics for the 9 PhenoBAGs in all post-GWAS analyses. The fastGWA GWAS accounted for key confounders, including age, dataset status (training/validation/test or independent test), age-squared, sex, interactions of age with sex, and the first 40 genetic principal components. For the brain BAG GWAS specifically, additional covariates for total intracranial volume and brain position in the scanner were included. A genome-wide significance threshold (5×10^{-9}), was applied.

ProtBAG GWAS: We used fastGWA to perform the 11 ProtBAGs, adjusting age, dataset status (training/validation/test or independent test), age-squared, sex, interactions of age with sex, systolic/diastolic blood pressure, BMI, waist circumference, standing height, weight, and the first 40 genetic principal components. We applied a genome-wide significance threshold (5×10^{-11}) to annotate the significant independent genomic loci.

Annotation of genomic loci: For all GWASs, genomic loci were annotated using FUMA⁶¹. For genomic loci annotation, FUMA initially identified lead SNPs (correlation $r^2 \leq 0.1$, distance < 250 kilobases) and assigned them to non-overlapping genomic loci. The lead SNP with the lowest P-value (i.e., the top lead SNP) represented the genomic locus. Further details on the definitions of top lead SNP, lead SNP, independent significant SNP, and candidate SNP can be found in **Supplementary eMethod 1**. For visualization purposes in **Fig. 3**, we have mapped the top lead SNP of each locus to the cytogenetic regions based on the GRCh37 cytoband.

(b) Genetic correlation: We estimated the genetic correlation (g_c) between each PhenoBAG-ProtBAG pair using the LDSC software. We employed precomputed LD scores from the 1000 Genomes of European ancestry, maintaining default settings for other parameters in LDSC. It's worth noting that LDSC corrects for sample overlap, ensuring an unbiased genetic correlation estimate⁶². We also computed the pairwise Pearson's r correlation coefficient to understand whether the genetic correlation largely mirrors the phenotypic correlation (p_c). Statistical significance was determined using Bonferroni correction (0.05/9/11).

(c) PRS calculation: PRS was computed using split-sample sensitivity GWASs (split1 and split2) for the PhenoBAG and ProtBAG GWASs. The PRS weights were established using split1/discovery GWAS data as the base/training set, while the split2/replication GWAS summary statistics served as the target/testing data. Both base and target data underwent rigorous quality control procedures involving several steps: *i*) excluding duplicated and ambiguous SNPs in the base data; *ii*) excluding high heterozygosity samples in the target data; and *v*) eliminating duplicated, mismatching, and ambiguous SNPs in the target data.

After completing the QC procedures, PRS for the split2 group was calculated using the PRS-CS⁶³ method. PRC-CS applies a continuous shrinkage prior, which adjusts the SNP effect sizes based on their LD structure. SNPs with weaker evidence are "shrunk" toward zero, while those with stronger evidence retain larger effect sizes. This avoids overfitting and improves prediction performance. No clumping was performed because the method takes LD into account. The shrinkage parameter was not set, and the algorithm learned it via a fully Bayesian approach.

(d) Bayesian colocalization: We used the R package (*coloc*) to investigate the genetic colocalization signals between two traits (i.e., hepatic ProtBAG vs. hepatic PhenoBAG, and

hepatic ProtBAG vs. renal PhenoBAG) at each genomic locus. We employed the Fully Bayesian colocalization analysis using Bayes Factors (*coloc.abf*). This method examines the posterior probability (PP.H4.ABF: Approximate Bayes Factor) to evaluate hypothesis *H4*, which suggests the presence of a single shared causal variant associated with both traits within a specific genomic locus. To determine the significance of the *H4* hypothesis, we set a threshold of $PP.H4.ABF > 0.8^{35}$. All other parameters (e.g., the prior probability of p_{12}) were set as default. For each pair of traits, the genomic locus ($N > 100$ SNPs) was defined by default from FUMA for one trait, and then the *coloc* package extracted and harmonized the GWAS summary statistics within this locus for the other trait.

(e) Two-sample bidirectional Mendelian randomization: We constructed a multi-layer causal network linking ProtBAG, PhenoBAG, and DE using a bi-directional Mendelian randomization approach. In total, 4 bi-directional causal networks were established: i) *ProtBAG2PhenoBAG*, ii) *PhenoBAG2ProtBAG*, iii) *PhenoBAG2DE*, and iv) *DE2PhenoBAG*. These networks used summary statistics from our ProtBAG and PhenoBAG GWAS in the UKBB, the FinnGen³⁰, and the PGC³¹ study for the 525 DEs. For example, the *ProtBAG2PhenoBAG* causal network employed the 11 ProtBAGs as exposure variables and the 9 PhenoBAGs as outcome variables. The systematic quality-checking procedures to ensure unbiased exposure/outcome variable and instrumental variable (IVs) selection are detailed below.

We used a two-sample Mendelian randomization approach implemented in the *TwoSampleMR* package⁶⁴ to infer the causal relationships within these networks. We employed five distinct Mendelian randomization methods, presenting the results of the inverse variance weighted (IVW) method in the main text and the outcomes of the other four methods (Egger, weighted median, simple mode, and weighted mode estimators) in the supplement. The STROBE-MR Statement⁶⁵ guided our analyses to increase transparency and reproducibility, encompassing the selection of exposure and outcome variables, reporting statistics, and implementing sensitivity checks to identify potential violations of underlying assumptions. First, we performed an unbiased quality check on the GWAS summary statistics. Notably, the absence of population overlapping bias⁶⁶ was confirmed, given that FinnGen and UKBB participants largely represent populations of European ancestry without explicit overlap. PGC GWAS summary data were ensured to exclude UKBB participants. For the *ProtBAG2PhenoBAG* and *PhenoBAG2ProtBAG* networks from UKBB, we reran the ProtBAG GWAS and ensured no overlapping populations with PhenoBAG. Furthermore, all consortia's GWAS summary statistics were based on or lifted to GRCh37. Subsequently, we selected the effective exposure variables by assessing the statistical power of the exposure GWAS summary statistics in terms of instrumental variables (IVs), ensuring that the number of IVs exceeded 7 before harmonizing the data. Crucially, the function "*clump_data*" was applied to the exposure GWAS data, considering LD. The function "*harmonise_data*" was then used to harmonize the GWAS summary statistics of the exposure and outcome variables. Bonferroni correction was applied to all tested traits based on the number of effective ProtBAGs, PhenoBAGs, or DEs, whichever was larger.

Finally, we conducted multiple sensitivity analyses. First, we conducted a heterogeneity test to scrutinize potential violations in the IV's assumptions. To assess horizontal pleiotropy, which indicates the IV's exclusivity assumption⁶⁷, we utilized a funnel plot, single-SNP Mendelian randomization methods, and the Egger estimator. Furthermore, we performed a leave-one-out analysis, systematically excluding one instrument (SNP/IV) at a time, to gauge the sensitivity of the results to individual SNPs.

Method 6: Prediction analyses for 14 systemic disease categories and the risk of mortality

We investigated the clinical promise of the 11 ProtBAGs, 9 PhenoBAGs, and their PRSs in two sets of prediction analyses: i) classification tasks for predicting 14 systemic disease categories based on the ICD-10 code (**Supplementary eTable 7**) and ii) survival analysis for the risk of all-cause mortality.

(a) Support vector machines to classify patients of disease categories vs. controls: We applied SVM with ProtBAG, PhenoBAG, and their PRS, implementing a nested cross-validation procedure¹⁷ to optimize the hyperparameter C and predict individual-level outcomes. Unlike previous studies²⁹, we did not set aside an independent test dataset due to the relatively small sample size of the control population without any disease diagnoses ($N=1651$); patients for each disease category were defined by the ICD-10 code (Field-ID: 41270). Brain and eye PhenoBAGs were excluded from the analysis due to insufficient sample sizes after integrating all features. **Figure 5a-b** reports the balanced accuracy (BA) obtained from the nested test data. The nested cross-validation procedure involved an outer loop repeated 50 times, with 80% of the data randomly allocated for training/validation and 20% for testing. Within the inner loop, the training/validation data underwent a 10-fold split for model optimization. **Supplementary eTable 7** provides detailed metrics, including balanced accuracy, sensitivity, specificity, negative predictive value, positive predictive value, and sample sizes for the training/validation/test datasets.

(b) Survival analysis for mortality risk: we employed a Cox proportional hazard model while adjusting for covariates (i.e., age and sex) to test the associations of the 11 ProtBAGs, 9 PhenoBAGs, and their PRS with all-cause mortality. The covariates were included as additional right-side variables in the model. The hazard ratio (HR), $\exp(\beta_R)$, was calculated and reported as the effect size measure that indicates the influence of each biomarker on the risk of mortality. To train the model, the "time" variable was determined by calculating the difference between the date of death (Field ID: 40000) for cases (or the censoring date for non-cases) and the date attending the assessment center (Field ID: 53). Participants who passed away after enrolling in the study were classified as cases.

Data Availability

The GWAS summary statistics corresponding to this study are publicly available on the MEDICINE knowledge portal (<https://labs-laboratory.com/medicine/>). Our study used data generated by the human protein atlas (HPA: <https://www.proteinatlas.org>). GWAS summary data for the DEs were downloaded from the official websites of FinnGen (R9: https://www.finnngen.fi/en/access_results) and PGC (<https://pgc.unc.edu/for-researchers/download-results/>). Individual data from UKBB can be requested with proper registration at <https://www.ukbiobank.ac.uk/>. Certain sensitive data (e.g., allele frequency information) supporting the findings are also available from the author upon request.

Code Availability

The software and resources used in this study are all publicly available:

- MLNI: <https://github.com/anbai106/mlni>, ProtBAG generation and classification for disease categories;
- AutoComplete: <https://github.com/sriramlab/AutoComplete>, Proteomics imputation;
- FUMA: <https://fuma.ctglab.nl/>, Gene mapping, genomic locus annotation;
- GCTA: <https://yanglab.westlake.edu.cn/software/gcta/#Overview>, fastGWA;
- LDSC: <https://github.com/bulik/ldsc>, genetic correlation
- TwoSampleMR: <https://mrcieu.github.io/TwoSampleMR/index.html>, Mendelian randomization;
- PRSs: <https://github.com/getian107/PRSs>, PRS calculation;
- Lifelines: <https://lifelines.readthedocs.io/en/latest/>, Survival analysis;
- coloc: <https://github.com/chr1swallace/coloc>; Bayesian colocalization.

Competing Interests

None

References

1. Wen, J. *et al.* The genetic architecture of biological age in nine human organ systems. *Nat Aging* 1–18 (2024) doi:10.1038/s43587-024-00662-8.
2. Tian, Y. E. *et al.* Heterogeneous aging across multiple organ systems and prediction of chronic disease and mortality. *Nat Med* 1–11 (2023) doi:10.1038/s41591-023-02296-6.
3. Wen, J. *et al.* The genetic architecture of multimodal human brain age. *Nat Commun* **15**, 2604 (2024).
4. Reicher, L. *et al.* Phenome-wide associations of human aging uncover sex-specific dynamics. *Nat Aging* **4**, 1643–1655 (2024).
5. Oh, H. S.-H. *et al.* Organ aging signatures in the plasma proteome track health and disease. *Nature* **624**, 164–172 (2023).
6. de Lange, A.-M. G. & Cole, J. H. Commentary: Correction procedures in brain-age prediction. *Neuroimage Clin* **26**, 102229 (2020).
7. Nie, C. *et al.* Distinct biological ages of organs and systems identified from a multi-omics study. *Cell Reports* **38**, 110459 (2022).
8. Dabrowski, J. K. *et al.* Probabilistic inference of epigenetic age acceleration from cellular dynamics. *Nat Aging* **4**, 1493–1507 (2024).
9. Kuiper, L. M. *et al.* Epigenetic and Metabolomic Biomarkers for Biological Age: A Comparative Analysis of Mortality and Frailty Risk. *The Journals of Gerontology: Series A* **78**, 1753–1762 (2023).
10. Sun, B. B. *et al.* Plasma proteomic associations with genetics and health in the UK Biobank. *Nature* 1–10 (2023) doi:10.1038/s41586-023-06592-6.

11. Eldjarn, G. H. *et al.* Large-scale plasma proteomics comparisons through genetics and disease associations. *Nature* 1–11 (2023) doi:10.1038/s41586-023-06563-x.
12. Carrasco-Zanini, J. *et al.* Proteomic signatures improve risk prediction for common and rare diseases. *Nat Med* **30**, 2489–2498 (2024).
13. Cole, J. H. *et al.* Brain age predicts mortality. *Mol Psychiatry* **23**, 1385–1392 (2018).
14. Beheshti, I., Nugent, S., Potvin, O. & Duchesne, S. Bias-adjustment in neuroimaging-based brain age frameworks: A robust scheme. *NeuroImage: Clinical* **24**, 102063 (2019).
15. Zhang, B., Zhang, S., Feng, J. & Zhang, S. Age-level bias correction in brain age prediction. *NeuroImage: Clinical* **37**, 103319 (2023).
16. Argentieri, M. A. *et al.* Proteomic aging clock predicts mortality and risk of common age-related diseases in diverse populations. *Nat Med* **30**, 2450–2460 (2024).
17. Wen, J. *et al.* Convolutional neural networks for classification of Alzheimer’s disease: Overview and reproducible evaluation. *Medical Image Analysis* **63**, 101694 (2020).
18. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat Genet* **47**, 291–295 (2015).
19. Kachuri, L. *et al.* Principles and methods for transferring polygenic risk scores across global populations. *Nat Rev Genet* **25**, 8–25 (2024).
20. Pingault, J.-B. *et al.* Using genetic data to strengthen causal inference in observational research. *Nat Rev Genet* **19**, 566–580 (2018).
21. Zhao, B. *et al.* Heart-brain connections: Phenotypic and genetic insights from magnetic resonance images. *Science* **380**, abn6598 (2023).
22. McCracken, C. *et al.* Multi-organ imaging demonstrates the heart-brain-liver axis in UK Biobank participants. *Nat Commun* **13**, 7839 (2022).

23. Wen, J. *et al.* Neuroimaging-AI Endophenotypes of Brain Diseases in the General Population: Towards a Dimensional System of Vulnerability. 2023.08.16.23294179 Preprint at <https://doi.org/10.1101/2023.08.16.23294179> (2023).
24. Liu, Y. *et al.* Genetic architecture of 11 organ traits derived from abdominal MRI using deep learning. *eLife* **10**, e65554 (2021).
25. Jaggi, A. *et al.* A structural heart-brain axis mediates the association between cardiovascular risk and cognitive function. *Imaging Neuroscience* **2**, 1–18 (2024).
26. Wen, J. Multiorgan biological age shows that no organ system is an island. *Nat Aging* 1–2 (2024) doi:10.1038/s43587-024-00690-4.
27. Boquet-Pujadas, A. *et al.* Brain-heart-eye axis revealed by multi-organ imaging genetics and proteomics. 2025.01.04.25319995 Preprint at <https://doi.org/10.1101/2025.01.04.25319995> (2025).
28. Zhao, B. *et al.* Heart-brain connections: Phenotypic and genetic insights from magnetic resonance images. *Science* **380**, abn6598 (2023).
29. Wen, J. *et al.* Nine Neuroimaging-AI Endophenotypes Unravel Disease Heterogeneity and Partial Overlap across Four Brain Disorders: A Dimensional Neuroanatomical Representation. 2023.08.16.23294179 Preprint at <https://doi.org/10.1101/2023.08.16.23294179> (2024).
30. Kurki, M. I. *et al.* FinnGen provides genetic insights from a well-phenotyped isolated population. *Nature* **613**, 508–518 (2023).
31. O’Donovan, M. C. What have we learned from the Psychiatric Genomics Consortium. *World Psychiatry* **14**, 291–293 (2015).

32. Nadeau, C. & Bengio, Y. Inference for the Generalization Error. *Machine Learning* **52**, 239–281 (2003).
33. Shen, X. *et al.* Nonlinear dynamics of multi-omics profiles during human aging. *Nat Aging* **4**, 1619–1634 (2024).
34. Ikram, M. A. The use and misuse of ‘biological aging’ in health research. *Nat Med* **30**, 3045–3045 (2024).
35. Giambartolomei, C. *et al.* Bayesian Test for Colocalisation between Pairs of Genetic Association Studies Using Summary Statistics. *PLOS Genetics* **10**, e1004383 (2014).
36. Wen, J. *et al.* Reproducible evaluation of diffusion MRI features for automatic classification of patients with Alzheimer’s disease. 51.
37. Varoquaux, G. Cross-validation failure: Small sample sizes lead to large error bars. *Neuroimage* **180**, 68–77 (2018).
38. Bashyam, V. M. *et al.* MRI signatures of brain age and disease over the lifespan based on a deep brain network and 14□468 individuals worldwide. *Brain* **143**, 2312–2324 (2020).
39. Dhindsa, R. S. *et al.* Rare variant associations with plasma protein levels in the UK Biobank. *Nature* 1–9 (2023) doi:10.1038/s41586-023-06547-x.
40. Lehallier, B. *et al.* Undulating changes in human plasma proteome profiles across the lifespan. *Nat Med* **25**, 1843–1850 (2019).
41. Konopka, G. & Bhaduri, A. Functional genomics and systems biology in human neuroscience. *Nature* **623**, 274–282 (2023).
42. Li, W. Rare-variant genetic architecture. *Nat Genet* **55**, 327–327 (2023).
43. Garg, M. *et al.* Disease prediction with multi-omics and biomarkers empowers case–control genetic discoveries in the UK Biobank. *Nat Genet* **56**, 1821–1831 (2024).

44. Luo, H. *et al.* Association of plasma proteomics with mortality in individuals with and without type 2 diabetes: Results from two population-based KORA cohort studies. *BMC Medicine* **22**, 420 (2024).
45. Walter, S. *et al.* Genetic, Physiological, and Lifestyle Predictors of Mortality in the General Population. *Am J Public Health* **102**, e3–e10 (2012).
46. Chatsirisupachai, K., Lesluyes, T., Paraoan, L., Van Loo, P. & de Magalhães, J. P. An integrative analysis of the age-associated multi-omic landscape across cancers. *Nat Commun* **12**, 2345 (2021).
47. Mavromatis, L. A. *et al.* Multi-omic underpinnings of epigenetic aging and human longevity. *Nat Commun* **14**, 2236 (2023).
48. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
49. Mbatchou, J. *et al.* Computationally efficient whole-genome regression for quantitative and binary traits. *Nat Genet* **53**, 1097–1103 (2021).
50. Li, W. UK Biobank pharma proteomics resource. *Nat Genet* **55**, 1781–1781 (2023).
51. Sun, B. B. *et al.* Plasma proteomic associations with genetics and health in the UK Biobank. *Nature* **622**, 329–338 (2023).
52. An, U. *et al.* Deep learning-based phenotype imputation on population-scale biobank data increases genetic discoveries. *Nat Genet* **55**, 2269–2276 (2023).
53. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867–2873 (2010).
54. Jiang, L. *et al.* A resource-efficient tool for mixed model association analysis of large-scale data. *Nat Genet* **51**, 1749–1755 (2019).

55. Price, A. L., Zaitlen, N. A., Reich, D. & Patterson, N. New approaches to population stratification in genome-wide association studies. *Nat Rev Genet* **11**, 459–463 (2010).
56. Abraham, G., Qiu, Y. & Inouye, M. FlashPCA2: principal component analysis of Biobank-scale genotype datasets. *Bioinformatics* **33**, 2776–2778 (2017).
57. Wen, J. *et al.* Genomic loci influence patterns of structural covariance in the human brain. *Proceedings of the National Academy of Sciences* **120**, e2300842120 (2023).
58. Wen, J. *et al.* Characterizing Heterogeneity in Neuroimaging, Cognition, Clinical Symptoms, and Genetics Among Patients With Late-Life Depression. *JAMA Psychiatry* (2022) doi:10.1001/jamapsychiatry.2022.0020.
59. Wen, J. *et al.* Genetic and clinical correlates of two neuroanatomical AI dimensions in the Alzheimer’s disease continuum. *Transl Psychiatry* **14**, 1–14 (2024).
60. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: A Tool for Genome-wide Complex Trait Analysis. *Am J Hum Genet* **88**, 76–82 (2011).
61. Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping and annotation of genetic associations with FUMA. *Nat Commun* **8**, 1826 (2017).
62. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits. *Nat Genet* **47**, 1236–1241 (2015).
63. Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C. A. & Smoller, J. W. Polygenic prediction via Bayesian regression and continuous shrinkage priors. *Nat Commun* **10**, 1776 (2019).
64. Hemani, G. *et al.* The MR-Base platform supports systematic causal inference across the human phenome. *eLife* **7**, e34408 (2018).

65. Skrivankova, V. W. *et al.* Strengthening the Reporting of Observational Studies in Epidemiology Using Mendelian Randomization: The STROBE-MR Statement. *JAMA* **326**, 1614–1621 (2021).
66. Sanderson, E. *et al.* Mendelian randomization. *Nat Rev Methods Primers* **2**, 1–21 (2022).
67. Bowden, J. *et al.* A framework for the investigation of pleiotropy in two-sample summary data Mendelian randomization. *Stat Med* **36**, 1783–1802 (2017).

Acknowledgments

The MULTI consortium aims to integrate multi-organ imaging with multi-omics data to advance our understanding of human aging and disease mechanisms. This research has been conducted using data from UK Biobank, a major biomedical database, under Application Number 647044, 35148, and 60698. We want to express our sincere gratitude to the UK Biobank team for their invaluable contribution to advancing clinical research in our field (<https://www.ukbiobank.ac.uk/>). We also acknowledge the UKB-PPP consortium (<https://registry.opendata.aws/ukbPPP/>; Category code: 1838) to share the returned data with the community. We want to acknowledge the participants and investigators of the FinnGen study and the PGC consortium, and we thank FinnGen (<https://www.finnngen.fi/en>) and PGC (<https://pgc.unc.edu/>) for their generosity in sharing the GWAS summary statistics with the scientific community. We acknowledge the leadership of the Brain Imaging Genetics (BIG) workgroup, led by Dr. Tavia Evans, Dr. Natalia Vilor-Tejedor, and Dr. Junhao Wen, within the International Society to Advance Alzheimer's Research and Treatment (ISTAART) community, for advocating brain imaging genetics in Alzheimer's and aging research.