

UCLA

UCLA Previously Published Works

Title

Multiset correlation and factor analysis enables exploration of multi-omics data

Permalink

<https://escholarship.org/uc/item/8cd9h24h>

Journal

Cell Genomics, 3(8)

ISSN

2666-979X

Authors

Brown, Brielin C

Wang, Collin

Kasela, Silva

et al.

Publication Date

2023-08-01

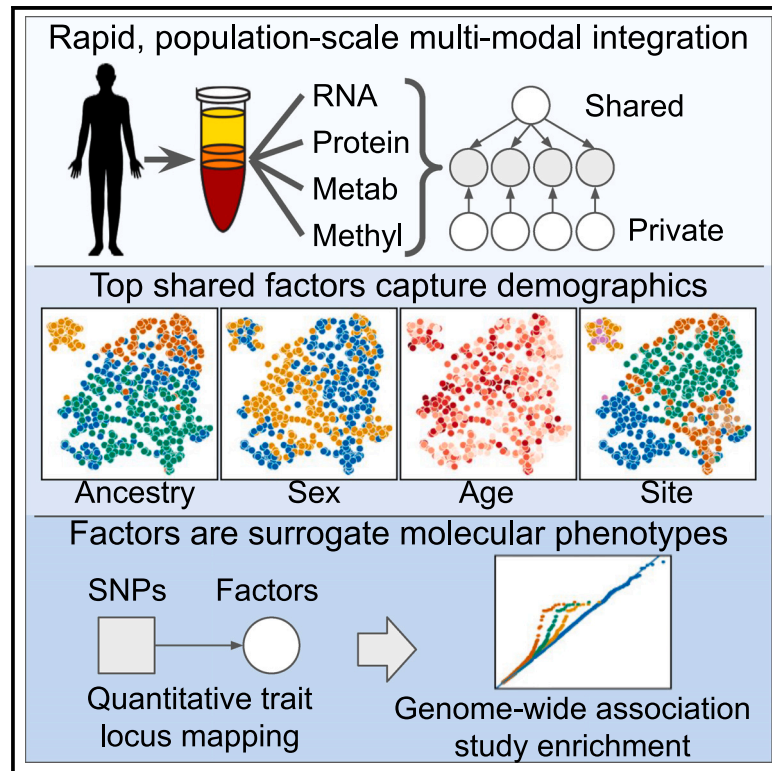
DOI

10.1016/j.xgen.2023.100359

Peer reviewed

## Multiset correlation and factor analysis enables exploration of multi-omics data

### Graphical abstract



### Authors

Brielin C. Brown, Collin Wang, Silva Kasela, ..., Kristin G. Ardlie, David A. Knowles, Tuuli Lappalainen

### Correspondence

bbrown@nygenome.org

### In brief

Brown, Wang et al. introduce MCFA, an approach to multi-modal dataset integration that generalizes canonical correlation analysis. MCFA is broadly applicable to data integration challenges but has been designed to handle issues in population-scale multi-omics data. A variety of analyses on the TOPMed/MESA multi-omics pilot demonstrate the power of this method.

### Highlights

- Rapid, unsupervised multi-modal data integration with self-inferred tuning parameters
- 614 ancestry-diverse participants from MESA/TOPMed with 5 omics types
- Top shared components capture ancestry, even without genetic information
- Further components are enriched for GWAS hits and related to metabolic disease



## Short article

# Multiset correlation and factor analysis enables exploration of multi-omics data

Brielin C. Brown,<sup>1,2,21,22,\*</sup> Collin Wang,<sup>1,3,21</sup> Silva Kasela,<sup>1,4</sup> François Aguet,<sup>5,6</sup> Daniel C. Nachun,<sup>7</sup> Kent D. Taylor,<sup>8</sup> Russell P. Tracy,<sup>9</sup> Peter Durda,<sup>9</sup> Yongmei Liu,<sup>10</sup> W. Craig Johnson,<sup>11</sup> David Van Den Berg,<sup>12</sup> Namrata Gupta,<sup>6</sup> Stacy Gabriel,<sup>6</sup> Joshua D. Smith,<sup>13</sup> Robert Gerzsten,<sup>14</sup> Clary Clish,<sup>6</sup> Quenna Wong,<sup>11</sup> George Papanicolau,<sup>15</sup> Thomas W. Blackwell,<sup>16</sup> Jerome I. Rotter,<sup>8</sup> Stephen S. Rich,<sup>17</sup> R. Graham Barr,<sup>18</sup> Kristin G. Ardlie,<sup>6</sup> David A. Knowles,<sup>1,2,3,4,20</sup> and Tuuli Lappalainen<sup>1,4,19,20</sup>

<sup>1</sup>New York Genome Center, New York, NY, USA

<sup>2</sup>Data Science Institute, Columbia University, New York, NY, USA

<sup>3</sup>Department of Computer Science, Columbia University, New York, NY, USA

<sup>4</sup>Department of Systems Biology, Columbia University, New York, NY, USA

<sup>5</sup>Illumina Incorporated, San Francisco, CA, USA

<sup>6</sup>The Broad Institute of MIT and Harvard, Boston, MA, USA

<sup>7</sup>Department of Pathology, Stanford University, Stanford, CA, USA

<sup>8</sup>Department of Pediatrics, The Institute for Translational Genomics and Population Sciences, The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance, CA, USA

<sup>9</sup>Department of Pathology and Laboratory Medicine, Larner College of Medicine, University of Vermont, Burlington, VT, USA

<sup>10</sup>Department of Medicine, Duke University Medical Center, Durham, NC, USA

<sup>11</sup>Department of Biostatistics, University of Washington, Seattle, WA, USA

<sup>12</sup>Department of Clinical Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA

<sup>13</sup>Northwest Genomics Center, University of Washington, Seattle, WA, USA

<sup>14</sup>Beth Israel Deaconess Medical Center, Division of Cardiovascular Medicine, Boston, MA, USA

<sup>15</sup>Division of Cardiovascular Sciences, National Heart, Lung, and Blood Institute, Bethesda, MD, USA

<sup>16</sup>Department of Biostatistics, University of Michigan School of Public Health, Ann Arbor, MI, USA

<sup>17</sup>Center for Public Health Genomics, University of Virginia, Charlottesville, VA, USA

<sup>18</sup>Mailman School of Public Health, Columbia University, New York, NY, USA

<sup>19</sup>Science for Life Laboratory, Department of Gene Technology, KTH Royal Institute of Technology, Stockholm, Sweden

<sup>20</sup>Senior author

<sup>21</sup>These authors contributed equally

<sup>22</sup>Lead contact

\*Correspondence: [bbrown@nygenome.org](mailto:bbrown@nygenome.org)  
<https://doi.org/10.1016/j.xgen.2023.100359>

## SUMMARY

Multi-omics datasets are becoming more common, necessitating better integration methods to realize their revolutionary potential. Here, we introduce multi-set correlation and factor analysis (MCFA), an unsupervised integration method tailored to the unique challenges of high-dimensional genomics data that enables fast inference of shared and private factors. We used MCFA to integrate methylation markers, protein expression, RNA expression, and metabolite levels in 614 diverse samples from the Trans-Omics for Precision Medicine/ Multi-Ethnic Study of Atherosclerosis multi-omics pilot. Samples cluster strongly by ancestry in the shared space, even in the absence of genetic information, while private spaces frequently capture dataset-specific technical variation. Finally, we integrated genetic data by conducting a genome-wide association study (GWAS) of our inferred factors, observing that several factors are enriched for GWAS hits and *trans*-expression quantitative trait loci. Two of these factors appear to be related to metabolic disease. Our study provides a foundation and framework for further integrative analysis of ever larger multi-modal genomic datasets.

## INTRODUCTION

Recent years have seen an explosion in multi-omics data, with studies simultaneously profiling RNA expression, protein levels, chromatin accessibility, and more.<sup>1</sup> By providing complementary views into the underlying biology, these datasets promise to illuminate molecular processes and disease states that cannot

be gleaned from any lone modality.<sup>2</sup> However, joint inference methods are lacking in either the number or type of modes that can be used or in flexibility and efficiency.<sup>1</sup> Multi-omics data bring substantial challenges: distributions differ between modes, the sample size is typically small relative to features, efficient algorithms are needed, and each mode has contributions from factors that are shared between modes and unique to itself.<sup>3,4</sup>



Canonical correlation analysis (CCA) is a statistical technique that infers shared factors between two data modes by finding correlated linear combinations of the features in each.<sup>5</sup> CCA has enjoyed substantial attention in genomics<sup>6–9</sup>; however, extending CCA to additional modes is fraught: at least 10 different formulations are equivalent in the two-mode case,<sup>10</sup> and many are challenging to fit.<sup>11</sup> Equivalently, CCA can be conceptualized as a probabilistic model (pCCA), revealing a connection to factor analysis.<sup>12</sup>

We have developed multi-set correlation and factor analysis (MCFA; [Figures 1A and S1](#)), an unsupervised integration method that generalizes pCCA and factor analysis, enabling fast inference of shared and private factors in multi-modal data. MCFA is designed to overcome challenges that are common with genomics data such as the large number of features relative to the sample size, the disparate data types, and the unknown contributions of dataset-specific technical factors. MCFA is based on two insights: (1) unlike traditional CCA, pCCA has only one natural extension to multi-modal data, which is both conceptually elegant and efficient to fit, and (2) after fitting pCCA, the residual in a mode represents private structure, which is well modeled by factor analysis. Our method combines these insights to fit factors that are shared across modalities and are private to each simultaneously. For efficiency and regularization, MCFA uses the top principal components (PCs) of each mode.<sup>6,7</sup> It allows the use of random matrix techniques<sup>13</sup> to choose the shared dimensionality and number of PCs, eliminating tuning parameters. Finally, MCFA is a natural approach to integration: as detailed in [Methods S1](#), there is a theoretical connection between our model and multi-set CCA.

We have applied MCFA to 614 ancestry-diverse individuals from the Multi-Ethnic Study of Atherosclerosis (MESA).<sup>14</sup> The Trans-Omics for Precision Medicine (TOPMed)<sup>15</sup> program instituted a multi-omics pilot study to evaluate the utility of long-term stored samples for discovery related to heart, lung, blood, and sleep disorders. MESA provided samples for five omics types: (1) whole-genome sequencing (WGS), (2) RNA sequencing of peripheral blood mononuclear cells (PBMCs), (3) DNA methylation array profiling from whole blood, (4) protein mass spectrometry of blood plasma, and (5) metabolite mass spectrometry of blood plasma. In addition, MESA has collected comprehensive phenotypic metadata. These data include demographic markers such as self-reported ancestry (SRA), sex, age, and education level; morphological features including height, weight, and hip circumference; clinical measures including those related to atherosclerosis, lipid levels, kidney function, and inflammatory biomarkers; and behavioral features regarding smoking, drinking, and exercise frequency.

## RESULTS

We integrated RNA sequencing, methylation, protein, and metabolite data using MCFA, which inferred a 14-dimensional shared space. We found that shared structure explained a large proportion of the variance in each mode ([Figure 1B](#), right). Protein levels had the highest sharing with 29.2% of the variance explained (VE) by the shared space, followed by RNA and metabolite levels (16.6% and 17.1%, respectively). Methylation

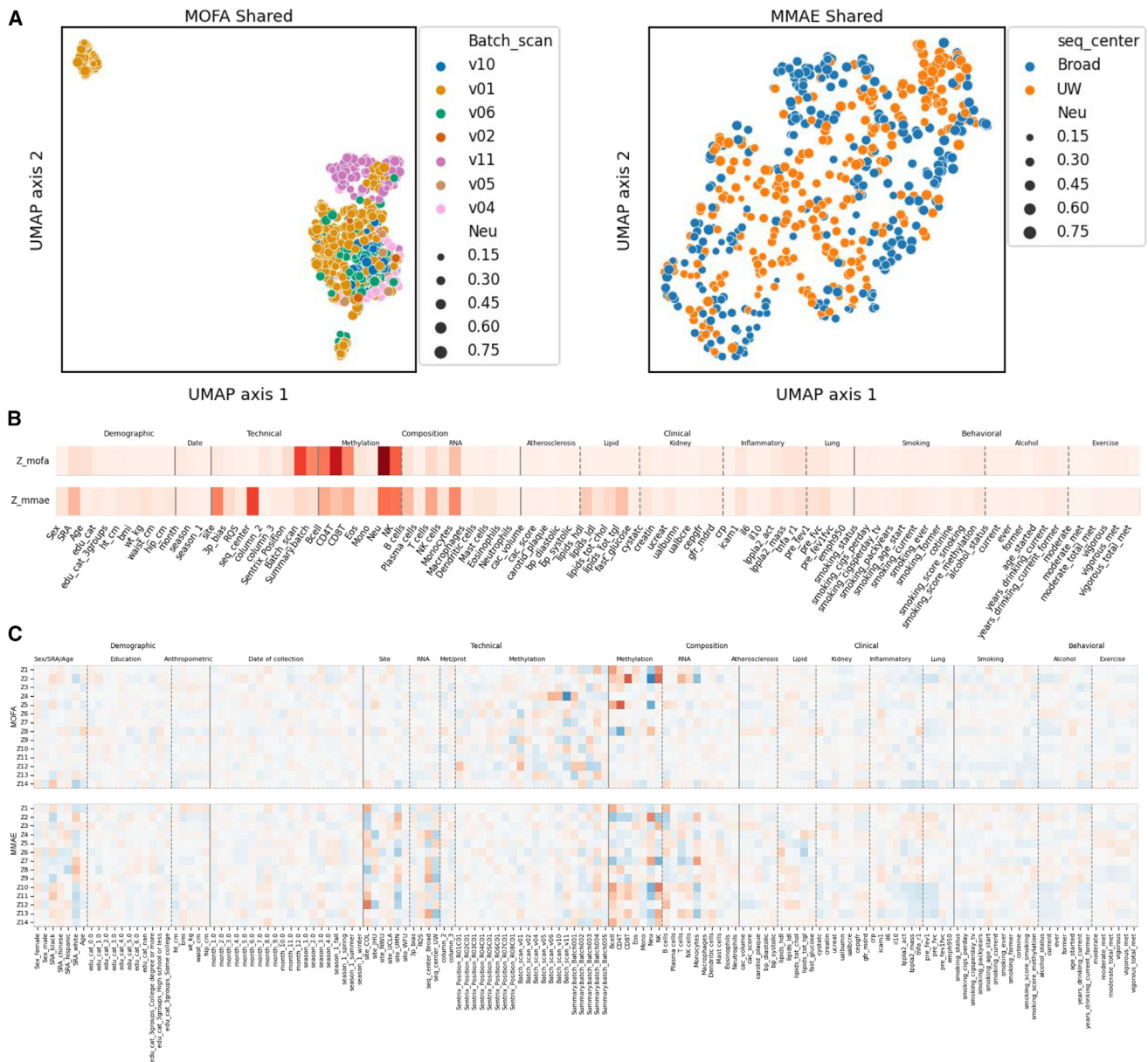
showed the least sharing, with only 8.1% VE by the shared space. Due to the high dimensionality of the data and the limited sample size, about half of the variance in each dataset is unmodeled to reduce overfitting. Using MCFA, it is possible to further infer the variance in each modality explained by the individual factors, thus determining which modalities contribute to each ([Figure 1B](#), left). Our top factor has contributions from all modalities, but their respective contributions to the other factors vary substantially.

We used uniform manifold approximation and projection (UMAP)<sup>16</sup> to construct a 2D embedding of the shared and private spaces ([Figure 1C](#)). We noticed a striking clustering of the individuals by SRA and sex in the shared space, even though the top PCs of individual modes do not cluster by these factors ([Figure S2](#)), and the shared space was inferred without genetic or sex chromosome features. Shared factor 1 separates Black and White individuals, with Hispanic individuals in between, while factor 3 separates Chinese individuals, and factor 2 differentiates by sex ([Figures S2 and S3](#)). We validated this structure via leave-one-out cross-validation, indicating our PC selection strategy mitigated over-fitting ([Figure S4](#)).

Next, we evaluated the total phenotypic VE by each of our inferred spaces ([Figures 1D and S2](#); [Tables S1, S2, and S3](#)). The shared space captured 95.3% of the variation in sex, 83.3% in site, 80.0% in SRA, and 60.2% in age. The shared space also captured anthropomorphic differences such as BMI (51.0% VE) and clinical measures including those related to kidney function (creatinine, 64.8% VE) and inflammation (tumor necrosis factor (TNF)-alpha receptor-1 69.1% VE). We used CIBERSORT<sup>17</sup> and the Houseman method<sup>18</sup> to estimate the cell-type composition of our RNA (PBMC) and methylation (whole blood) samples, respectively. Both shared and private spaces contributed to the relative proportions of PBMC-abundant cell types (e.g., T cells and natural killer (NK) cells) estimated from both data modalities, while the proportion of PBMC-depleted types (e.g., neutrophils) estimated from the methylation data was only captured by the methylation private space. Modality-private spaces frequently captured technical factors: 100% of the variance in sequencing center and 71.6% of the variance in 3' bias are captured by the RNA private space, while 76.8% of the methylation array batch is captured by its private space. Many phenotypes that are themselves measurements of metabolites were captured by the metabolite private space; however, the strongest association was with the month of sample collection (85.8% VE). We noticed no large associations between the protein private space and any of our metadata, despite several of our phenotypes being clinical protein markers; however, several of these factors are partially captured by the shared space.

We compared the results obtained on MESA using MCFA with other multi-modal analysis approaches. We focused on two alternative methods: (1) MOFA2<sup>4</sup> and (2) a multi-modal auto-encoder (MMAE, see [STAR Methods](#) and [Figure S5](#)). In the MOFA2 analysis, the methylation batch and cell-type proportions dominated the inferred shared space, likely owing to the very large number of features in that modality compared with the other modalities ([Figure 2](#)). The MMAE mitigated this over-focus on methylation somewhat and additionally captured RNA





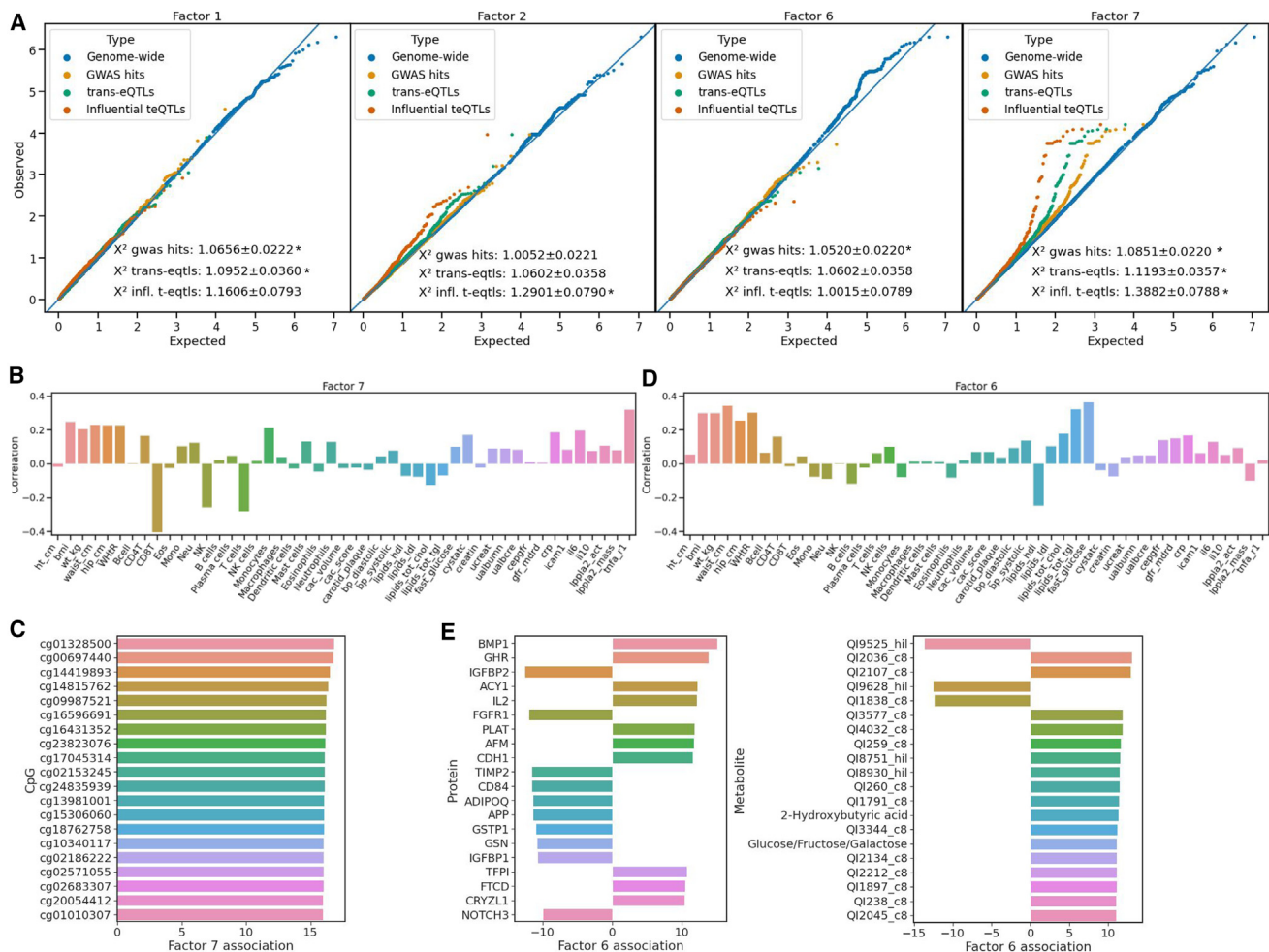
**Figure 2. Comparison of MCFA with other methods**

(A) UMAP embeddings of MOFA (left) and MMAE (right) shared space show that these methods fail to separate meaningful information from technical variation. (B) Variance in sample metadata explained by the MOFA2 (top) and MMAE (bottom) shared spaces. MOFA2 primarily learns factors related to the methylation dataset, while the MMAE additionally incorporates some factors related to RNA sequencing. (C) Correlation of each inferred factor with each metadata sample for MOFA (top) and the MMAE (bottom).

Finally, we integrated WGS data by conducting a genome-wide association study (GWAS) of the inferred factors while controlling for site, age, sex, and 11 genotype PCs. We hypothesized that genetic associations with our inferred factors, which represent major axes of molecular variation, may be enriched for known GWAS hits or *trans*-expression quantitative trait loci (eQTLs). We obtained a list of 10,174 such associations from the eQTLgen consortium,<sup>19</sup> of which 3,854 are *trans*-eQTLs, and further defined a more limited set of 1,107 “influential” *trans*-eQTLs that affect at least 10 genes. We tested the

GWAS of each factor for enrichment of these three categories and found 9 significant enrichments (mean  $\chi^2_{cat} > 1$ , false discovery rate [FDR] 5%; **Figures 3A and S6**).

Factor 7 showed the strongest enrichment for reported GWAS hits and *trans*-eQTLs. The top SNPs associated with factor 7 are from blood lipid studies and are located primarily around the FADS1 and FADS2 genes, which are known to regulate lipid metabolism.<sup>20</sup> These include rs174541 ( $p = 4.3 \times 10^{-5}$  for factor 7 association), which is also reported in GWASs of type 2 diabetes<sup>21</sup>; rs174549 ( $p = 5.6 \times 10^{-5}$ ), which is also reported in



**Figure 3. Factor interpretation and integration with GWAS data**

(A) QQ-plot of a GWAS for factors 1, 2, 6, and 7. Genetic associations with these factors are enriched for known GWAS loci (1, 6, and 7), *trans*-eQTLs (1 and 7), or highly influential *trans*-eQTLs (2 and 7).

(B and C) Correlation of factors 6 (B) and 7 (C) with morphological, immune-composition, and clinical metadata reveals that factor 6 is related to body composition and lipid profile, while factor 7 is related to body composition, inferred blood cell-type composition, and inflammatory biomarkers.

(D) Z-transformed correlation of individual protein and metabolite data with factor 6 reveals genes and metabolites related to insulin resistance and metabolic syndrome.

(E) Z-transformed correlation of individual methylation values with factor 7. Many genes collocated to these CpGs are involved in lipid metabolism.

GWASs of white blood cell count<sup>22</sup>; and rs1535 ( $p = 8.3 \times 10^{-5}$ ), which is also reported in a GWAS of inflammatory bowel disease<sup>23</sup> (Table S4). Factor 7 explains 6.7% of the modeled variation in methylation, the largest of any factor, and is anti-correlated with sample proportion of CD8<sup>+</sup> T cells and NK cells estimated from methylation data ( $\rho = -0.41$  and  $-0.25$ ), and correlated with BMI ( $\rho = 0.25$ ) and measures of inflammation including TNF-R1 ( $\rho = 0.33$ ) and interleukin-6 ( $\rho = 0.20$ ) (Figure 3B).

To assess the contribution of individual CpGs, we calculated the Z-transformed correlation of individual CpG values with factor 7 (Figure 3C). As epigenome-wide association studies remain small, generally little is known about the effects of individual CpGs and their associations with traits. Instead, we linked each gene to the CpGs falling in a window from 1.5 kb upstream

of the transcription start site to the transcription termination site. Many of the genes collocated to CpGs with high weights for factor 7 have been implicated in lipid metabolism GWASs including IQCG and TMEM178A (cg01328500 and cg02571055; phosphatidylcholine levels<sup>24</sup>), DSCAML1 (cg02571055; triglyceride levels<sup>25</sup>), PTK2 (cg02153245; ApoB and low-density lipoprotein [LDL] levels<sup>26</sup>), TULP4 (cg02571055; lipoprotein A levels<sup>27</sup>), and C7orf50 (cg20054412; LDL, high-density lipoprotein [HDL], and total cholesterol levels<sup>28</sup>). Interestingly, our second strongest hit, cg00697440, is collocated with CD86. Recent work has suggested that B7 molecules including CD86 play an important role in regulating CD8<sup>+</sup> T cell population dynamics.<sup>29</sup> While further research is needed to establish causal relationships of these genetic effects and methylation patterns in *cis* and *trans* on gene regulation and diverse traits, DNA methylation patterns

have been previously associated with lipid metabolism and metabolic disease.<sup>30,31</sup> Further research is required to determine whether the immune-cell component of this factor is related to the lipid metabolism component or whether these are simply independent biological functions captured by the same factor.

We used the same strategy to interpret factor 6. Factor 6 is correlated with fasting glucose, waist circumference, and triglycerides ( $\rho = 0.36, 0.34$ , and  $0.32$ , respectively) and anti-correlated with HDL cholesterol ( $\rho = -0.25$ ; Figure 3D). Factor 6 explains 6% of the variance in protein levels and 4.1% of the variance in metabolite levels. Many of the top-weighted metabolites are uncharacterized products from untargeted metabolomics, but the two top characterized targets are 2-hydroxybutyric acid, a known marker of insulin resistance and glucose intolerance,<sup>32,33</sup> and glucose itself (Figure 3E). Several of the top-weighted proteins in this factor have known roles in growth and development including BMP1, GHR, IGFBP2, and FGFR1. GWASs have implicated BMP1 in coronary artery disease,<sup>34,35</sup> IGFBP2 in type 2 diabetes and BMI,<sup>36</sup> and FGFR1 in triglyceride levels<sup>28</sup> and waist-hip ratio.<sup>37</sup> Other notable highly weighted proteins include TFPI, which is involved in blood coagulation and is associated with BMI-adjusted waist-hip ratio,<sup>38</sup> and ADIPOQ, which is involved in regulating glucose levels<sup>39</sup> (Figure 3E). Many of the top GWAS hits associated with this factor corroborate these observations, including rs4805885, which is associated with adiponectin (ADIPOQ) levels<sup>40</sup>; rs9787485, which is associated with insulin-carbohydrate interaction<sup>41</sup>; and rs7679, which is associated with HDL, LDL, and triglyceride levels<sup>42</sup> (Table S5).

Interestingly, the strongest genetic association with this factor comes from GWASs of schizophrenia (rs112973353;  $p = 1.6 \times 10^{-4}$  for factor 6 association), and we find 5 independent schizophrenia risk loci with factor 6 association  $p$  values below 0.01 (Table S5). Insulin resistance and schizophrenia have been consistently associated for nearly 100 years,<sup>43</sup> and while the association signal of each locus with factor 6 is relatively weak, the probability of finding 5 independent loci with these  $p$  values under the null is approximately  $4 \times 10^{-13}$ . While further research is needed, our results suggest that these particular loci may confer schizophrenia risk via insulin resistance. Another notable signal in our GWAS associations is related to erythrocyte and platelet traits. These hits include rs12451471 ( $p = 8 \times 10^{-4}$ ; mean corpuscular hemoglobin concentration<sup>44</sup>; platelet count<sup>45</sup>) and rs13224082 ( $p = 9 \times 10^{-4}$ ; platelet distribution width, platelet count, plateletcrit<sup>44</sup>), among others (Table S5). Again, further research is required to establish causality and direction of effect between genetics, metabolite and protein levels, and traits, but we note that there is an established link between insulin resistance and platelet dysfunction.<sup>46</sup>

## DISCUSSION

MCFA has several advantages compared with other multi-omics integration approaches. Compared with group factor analysis methods,<sup>4</sup> MCFA separates modality-specific from dataset-shared factors. Compared with non-negative matrix factorization-based methods that share a feature weight set across modalities,<sup>3</sup> MCFA is able to use all data types. As we have shown,

MCFA is also substantially faster and is able to handle datasets with unbalanced numbers of features across the modes.

While our top factors captured ancestry and sex, these factors are usually observed and considered confounding in clinical applications. In that context, one could fit the model conditional on known confounding factors. Since we see exploratory data analysis as a primary application of MCFA, our goal instead was to map the primary axes of biological variation contained within these population-scale multi-omics data. It is important that these factors are a primary driver of variation within such data, as it implies that sampling across race and sex is critical for equitable discovery in medical genomics. Still, because these factors are captured by the top components, and the components themselves are orthogonal, further components can still capture disease-relevant information.

Integration with GWAS is biased toward well-powered studies that will typically have more hits, some of which may be acting indirectly through another phenotype.<sup>47</sup> Interpretability of factors is also biased toward the metadata collected in the study. In MESA, the goal was evaluation of risk factors for heart disease, and thus MESA focused metadata collection on lipid phenotypes, inflammatory biomarkers, and body morphology. It is therefore unsurprising that we are most easily able to interpret factors related to metabolic syndrome, lipid metabolism, and immune function in this study. Still, the ability of MCFA to produce results that are correlated with these factors demonstrates the utility of broad-scale sample metadata when interpreting results from multi-omics studies.

Careful consideration is required when analyzing multi-omics datasets that include WGS or genotype data. There are two primary ways that one can think about integrating these data: (1) include genetic information as a mode in the fit model, interpretable as inferring a latent state that affects genotype as well as molecular factors, or (2) look for genetic associations with inferred molecular factors, interpretable as mapping QTLs for inferred molecular phenotypes. In this study, we chose the latter due to the improved causal interpretation and to demonstrate the utility of surrogate molecular phenotypes. In other cases, for example the analysis of genetic copy-number variation data in tumor samples, the former analysis approach may be preferred. Future work with larger sample sizes may allow for network inference and Mendelian randomization methods to generate directed hypotheses.<sup>47,48</sup> Genetic associations are particularly valuable in this, with the inferred axes of molecular variation providing promising future traits for GWAS and phenotype-wide association studies. TOPMed is among the most ambitious current efforts to collect multi-omics population-level data; thus, given the results of this pilot analysis, we expect future integration studies in this cohort to be fruitful.

## Limitations of the study

Due to the use of observational data and unsupervised methods, all analyses should be considered exploratory; they can find structure in the data while generating hypotheses but cannot be used to make causal claims and may reflect technical properties of the underlying data. For example, in MESA, the sample collection site is strongly correlated with SRA. We repeated our analysis of the VE by the learned space while additionally



controlling for site (Table S3) and noticed a small decrease in the proportion of VE in SRA (from 80.0% to 71.6%).

We observed that estimated cell-type composition had a strong association with both shared and private spaces. Since cell-type composition was inferred from the data, there may be circularity in composition estimation itself. In addition, complex interactions exist between cell-type composition in tissue samples and clinical, environmental factors as well as technical factors related to biospecimen collection. Thus, caution is necessary for biological interpretation in this aspect of the analysis.

## STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- **KEY RESOURCES TABLE**
- **RESOURCE AVAILABILITY**
  - Lead contact
  - Materials availability
  - Data and code availability
- **METHOD DETAILS**
  - Multiset correlation and factor analysis
  - Model initialization
  - High dimensionality and selection of hyperparameters
  - Calculating the variance explained
  - Calculating relative feature importance
  - SNP set enrichment analysis
  - The MESA multi-omics pilot
  - Cross-validation
  - Comparison to MOFA2 and MMAE
- **QUANTIFICATION AND STATISTICAL ANALYSIS**

## SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.xgen.2023.100359>.

## ACKNOWLEDGMENTS

B.C.B. would like to thank Lior Pachter and Nicholas Bray for numerous insightful conversations about CCA over the years. B.C.B. would also like to thank Andrew Stirn for reviewing the auto-encoder code. Funding for D.A.K. and B.C.B. is provided by NIA U01AG068880. Funding for B.C.B. is provided by NHGRI K99HG012373 and the Columbia Data Science Institute. Funding for T.L. and S.K. is provided by National Heart, Lung, and Blood Institute (NHLBI) R01HL142028. Funding for T.L. is provided by NIH R01AG057422 and NIMH R01MH106842. Funding for MESA lung measures is provided by NHLBI R01HL077612 and R01HL093081. WGS for the TOPMed program was supported by the NHLBI. WGS for “NHLBI TOPMed: Multi-Ethnic Study of Atherosclerosis (MESA)” (phs001416.v1.p1) was performed at the Broad Institute of MIT and Harvard (3U54HG003067-13S1). Centralized read mapping and genotype calling, along with variant quality metrics and filtering, were provided by the TOPMed Informatics Research Center (3R01HL-117626-02S1). Phenotype harmonization, data management, sample-identity quality control (QC), and general study coordination were provided by the TOPMed Data Coordinating Center (3R01HL-120393-02S1) and TOPMed MESA Multi-Omics (HHSN2682015000031/HSN26800004). The MESA projects are conducted and supported by the National Heart, Lung, and Blood Institute (NHLBI) in collaboration with MESA investigators (accession number phs000209.v13.p3). Support for the MESA projects are conducted and supported by the NHLBI in collaboration with MESA investigators. Support for

MESA is provided by contracts 75N92020D00001, HHSN2682015000031, N01-HC-95159, 75N92020D00005, N01-HC-95160, 75N92020D00002, N01-HC-95161, 75N92020D00003, N01-HC-95162, 75N92020D00006, N01-HC-95163, 75N92020D00004, N01-HC-95164, 75N92020D00007, N01-HC-95165, N01-HC-95166, N01-HC-95167, N01-HC-95168, N01-HC-95169, UL1-TR-000040, UL1-TR-001079, UL1-TR-001420, UL1TR001881, DK063491, and R01HL105756. The authors thank the other investigators, the staff, and the participants of the MESA study for their valuable contributions. A full list of participating MESA investigators and institutes can be found at <http://www.mesa-nhlbi.org>.

## AUTHOR CONTRIBUTIONS

Conceptualization, B.C.B., D.A.K., and T.L.; methodology, B.C.B. and C.W.; software, B.C.B. and C.W.; validation, B.C.B.; formal analysis, B.C.B., C.W., and D.A.K.; investigation, B.C.B. and C.W.; resources, S.K., F.A., D.C.N., K.D.T., R.P.T., P.D., Y.L., W.C.J., D.V.D.B., N.G., S.G., J.D.S., R.G., C.C., Q.W., G.P., T.W.B., J.I.R., S.S.R., R.G.B., K.G.A., D.A.K., and T.L.; data curation, B.C.B., S.K., F.A., D.C.N., K.D.T., R.P.T., P.D., Y.L., W.C.J., D.V.D.B., N.G., S.G., J.D.S., R.G., C.C., Q.W., G.P., T.W.B., J.I.R., S.S.R., R.G.B., K.G.A., D.A.K., and T.L.; writing – original draft, B.C.B.; writing – review & editing, B.C.B., C.W., S.K., J.I.R., S.S.R., D.A.K., and T.L.; supervision, J.I.R., S.S.R., D.A.K., and T.L.; funding acquisition, B.C.B., R.G.B., D.A.K., and T.L.

## DECLARATION OF INTERESTS

T.L. is a paid adviser or consultant of GSK, Pfizer, and Goldfinch Bio and has equity in Variant Bio. F.A. is an employee and shareholder of Illumina, Inc.

Received: September 28, 2022

Revised: April 26, 2023

Accepted: June 14, 2023

Published: July 10, 2023

## REFERENCES

1. Krassowski, M., Das, V., Sahu, S.K., and Misra, B.B. (2020). State of the Field in Multi-Omics Research: From Computational Needs to Data Mining and Sharing. *Front. Genet.* *11*, 610798. <https://doi.org/10.3389/FGENE.2020.610798>.
2. Hasin, Y., Seldin, M., and Lusis, A. (2017). Multi-omics approaches to disease. *Genome Biol.* *18*, 83. <https://doi.org/10.1186/S13059-017-1215-1>.
3. Welch, J.D., Kozareva, V., Ferreira, A., Vanderburg, C., Martin, C., and Macosko, E.Z. (2019). Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain Cell Identity. *Cell* *177*, 1873–1887.e17. <https://doi.org/10.1016/j.cell.2019.05.006>.
4. Argelaguet, R., Velten, B., Arnol, D., Dietrich, S., Zenz, T., Marioni, J.C., Buettner, F., Huber, W., and Stegle, O. (2018). Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets. *Mol. Syst. Biol.* *14*, e8124. <https://doi.org/10.15252/msb.20178124>.
5. Hotelling, H. (1936). Relations Between Two Sets of Variates. *Biometrika* *28*, 321–377. <https://doi.org/10.2307/2333955>.
6. Brown, B.C., Bray, N.L., and Pachter, L. (2018). Expression reflects population structure. *PLoS Genet.* *14*, e1007841. <https://doi.org/10.1371/journal.pgen.1007841>.
7. Sonesson, C., Liljebjörn, H., Fioretos, T., and Fontes, M. (2010). Integrative analysis of gene expression and copy number alterations using canonical correlation analysis. *BMC Bioinformatics* *11*, 1–20. <https://doi.org/10.1186/1471-2105-11-191>.
8. Naylor, M.G., Lin, X., Weiss, S.T., Raby, B.A., and Lange, C. (2010). Using Canonical Correlation Analysis to Discover Genetic Regulatory Variants. *PLoS One* *5*, e10395. <https://doi.org/10.1371/JOURNAL.PONE.0010395>.
9. Butler, A., Hoffman, P., Smibert, P., Papalexi, E., and Satija, R. (2018). Integrating single-cell transcriptomic data across different conditions,

- technologies, and species. *Nat. Biotechnol.* 36, 411–420. <https://doi.org/10.1038/nbt.4096>.
10. Kettnering, J.R. (1971). Canonical analysis of several sets of variables. *Biometrika* 58, 433–451.
  11. Asendorf, N.A. (2015). Informative Data Fusion: Beyond Canonical Correlation Analysis. <https://deepblue.lib.umich.edu/handle/2027.42/113419>.
  12. Bach, F.R., and Jordan, M.I. (2005). A Probabilistic Interpretation of Canonical Correlation Analysis. <https://www.di.ens.fr/~fbach/probaccpa.pdf>.
  13. Marčenko, V.A., and Pastur, L.A. (1967). Distribution of Eigenvalues for Some Sets of Random Matrices. *Math. USSR. Sb.* 1, 457–483. <https://doi.org/10.1070/sm1967v001n04abeh001994>.
  14. Bild, D.E., Bluemke, D.A., Burke, G.L., Detrano, R., Diez Roux, A.V., Folsom, A.R., Greenland, P., Jacob, D.R., Jr., Kronmal, R., Liu, K., et al. (2002). Multi-Ethnic Study of Atherosclerosis: objectives and design. *Am. J. Epidemiol.* 156, 871–881. <https://doi.org/10.1093/AJE/KWF113>.
  15. Taliun, D., Harris, D.N., Kessler, M.D., Carlson, J., Szpiech, Z.A., Torres, R., Taliun, S.A.G., Corvelo, A., Gogarten, S.M., Kang, H.M., et al. (2021). Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* 590, 290–299. <https://doi.org/10.1038/s41586-021-03205-y>.
  16. McInnes, L., Healy, J., and Melville, J. (2020). UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1802.03426>.
  17. Newman, A.M., Liu, C.L., Green, M.R., Gentles, A.J., Feng, W., Xu, Y., Hoang, C.D., Diehn, M., and Alizadeh, A.A. (2015). Robust enumeration of cell subsets from tissue expression profiles. *Nat. Methods* 12, 453–457. <https://doi.org/10.1038/nmeth.3337>.
  18. Houseman, E.A., Accomando, W.P., Koestler, D.C., Christensen, B.C., Marsit, C.J., Nelson, H.H., Wiencke, J.K., and Kelsey, K.T. (2012). DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinformatics* 13, 86. <https://doi.org/10.1186/1471-2105-13-86>.
  19. Vösa, U., Claringbould, A., Westra, H.J., Bonder, M.J., Deelen, P., Zeng, B., Kirsten, H., Saha, A., Kreuzhuber, R., Yazar, S., et al. (2021). Large-scale cis- and trans-eQTL analyses identify thousands of genetic loci and polygenic scores that regulate blood gene expression. *Nat. Genet.* 53, 1300–1310. <https://doi.org/10.1038/s41588-021-00913-z>.
  20. Schaeffer, L., Gohlke, H., Müller, M., Heid, I.M., Palmer, L.J., Kompauer, I., Demmelmair, H., Illig, T., Koletzko, B., and Heinrich, J. (2006). Common genetic variants of the FADS1 FADS2 gene cluster and their reconstructed haplotypes are associated with the fatty acid composition in phospholipids. *Hum. Mol. Genet.* 15, 1745–1756. <https://doi.org/10.1093/HMG/DDL117>.
  21. Dupuis, J., Langenberg, C., Prokopenko, I., Saxena, R., Soranzo, N., Jackson, A.U., Wheeler, E., Glazer, N.L., Bouatia-Naji, N., Gloyn, A.L., et al. (2010). New genetic loci implicated in fasting glucose homeostasis and their impact on type 2 diabetes risk. *Nat. Genet.* 42, 105–116. <https://doi.org/10.1038/NG.520>.
  22. Astle, W.J., Elding, H., Jiang, T., Allen, D., Ruklisa, D., Mann, A.L., Mead, D., Bouman, H., Riveros-Mckay, F., Kostadima, M.A., et al. (2016). The Allelic Landscape of Human Blood Cell Trait Variation and Links to Common Complex Disease. *Cell* 167, 1415–1429.e19. <https://doi.org/10.1016/J.CELL.2016.10.042>.
  23. Liu, J.Z., Van Sommeren, S., Huang, H., Ng, S.C., Alberts, R., Takahashi, A., Ripke, S., Lee, J.C., Jostins, L., Shah, T., et al. (2015). Association analyses identify 38 susceptibility loci for inflammatory bowel disease and highlight shared genetic risk across populations. *Nat. Genet.* 47, 979–986. <https://doi.org/10.1038/NG.3359>.
  24. Rhee, E.P., Ho, J.E., Chen, M.H., Shen, D., Cheng, S., Larson, M.G., Ghorbani, A., Shi, X., Helenius, I.T., O'Donnell, C.J., et al. (2013). A genome-wide association study of the human metabolome in a community-based cohort. *Cell Metab.* 18, 130–143. <https://doi.org/10.1016/J.CMET.2013.06.013>.
  25. Pollin, T.I., Damcott, C.M., Shen, H., Ott, S.H., Shelton, J., Horenstein, R.B., Post, W., McLenithan, J.C., Bielak, L.F., Peyser, P.A., et al. (2008). A null mutation in human APOC3 confers a favorable plasma lipid profile and apparent cardioprotection. *Science* 322, 1702–1705. <https://doi.org/10.1126/SCIENCE.1161524>.
  26. Richardson, T.G., Sanderson, E., Palmer, T.M., Ala-Korpela, M., Ference, B.A., Davey Smith, G., and Holmes, M.V. (2020). Evaluating the relationship between circulating lipoprotein lipids and apolipoproteins with risk of coronary heart disease: A multivariable Mendelian randomisation analysis. *PLoS Med.* 17, e1003062. <https://doi.org/10.1371/JOURNAL.PMED.1003062>.
  27. Sinnott-Armstrong, N., Tanigawa, Y., Amar, D., Mars, N., Benner, C., Aguirre, M., Venkataraman, G.R., Wainberg, M., Ollila, H.M., Kiiskinen, T., et al. (2021). Genetics of 35 blood and urine biomarkers in the UK Biobank. *Nat. Genet.* 53, 185–194. <https://doi.org/10.1038/s41588-020-00757-z>.
  28. Graham, S.E., Clarke, S.L., Wu, K.H.H., Kanoni, S., Zajac, G.J.M., Ramdas, S., Surakka, I., Ntalla, I., Vedantam, S., Winkler, T.W., et al. (2021). The power of genetic diversity in genome-wide association studies of lipids. *Nature* 600, 675–679. <https://doi.org/10.1038/S41586-021-04064-3>.
  29. Zenke, S., Palm, M.M., Braun, J., Gavrillo, A., Meiser, P., Böttcher, J.P., Beyersdorf, N., Ehl, S., Gerard, A., Lämmermann, T., et al. (2020). Quorum Regulation via Nested Antagonistic Feedback Circuits Mediated by the Receptors CD28 and CTLA-4 Confers Robustness to T Cell Population Dynamics. *Immunity* 52, 313–327.e7. <https://doi.org/10.1016/J.IMMUNI.2020.01.018>.
  30. Mittelstraß, K., and Waldenberger, M. (2018). DNA methylation in human lipid metabolism and related diseases. *Curr. Opin. Lipidol.* 29, 116–124. <https://doi.org/10.1097/MOL.0000000000000491>.
  31. Gomez-Alonso, M.D.C., Kretschmer, A., Wilson, R., Pfeiffer, L., Karhunen, V., Seppälä, I., Zhang, W., Mittelstraß, K., Wahl, S., Matias-Garcia, P.R., et al. (2021). DNA methylation and lipid metabolism: an EWAS of 226 metabolic measures. *Clin. Epigenetics* 13, 7. <https://doi.org/10.1186/s13148-020-00957-8>.
  32. Gall, W.E., Beebe, K., Lawton, K.A., Adam, K.P., Mitchell, M.W., Nakhle, P.J., Ryals, J.A., Milburn, M.V., Nannipieri, M., Camastra, S., et al. (2010).  $\alpha$ -Hydroxybutyrate Is an Early Biomarker of Insulin Resistance and Glucose Intolerance in a Nondiabetic Population. *PLoS One* 5, e10883. <https://doi.org/10.1371/JOURNAL.PONE.0010883>.
  33. Ferrannini, E., Natali, A., Camastra, S., Nannipieri, M., Mari, A., Adam, K.P., Milburn, M.V., Kastenmüller, G., Adamski, J., Tuomi, T., et al. (2013). Early Metabolic Markers of the Development of Dysglycemia and Type 2 Diabetes and Their Physiological Significance. *Diabetes* 62, 1730–1737. <https://doi.org/10.2337/DB12-0707>.
  34. Van Der Harst, P., and Verweij, N. (2018). Identification of 64 Novel Genetic Loci Provides an Expanded View on the Genetic Architecture of Coronary Artery Disease. *Circ. Res.* 122, 433–443. <https://doi.org/10.1161/CIRCRESAHA.117.312086>.
  35. Aragam, K.G., Jiang, T., Goel, A., Kanoni, S., Wolford, B.N., Atri, D.S., Weeks, E.M., Wang, M., Hindy, G., Zhou, W., et al. (2022). Discovery and systematic characterization of risk variants and genes for coronary artery disease in over a million participants. *Nat. Genet.* 54, 1803–1815. <https://doi.org/10.1038/S41588-022-01233-6>.
  36. Zhao, W., Rasheed, A., Tikkanen, E., Lee, J.J., Butterworth, A.S., Howson, J.M.M., Assimes, T.L., Chowdhury, R., Orho-Melander, M., Damrauer, S., et al. (2017). Identification of new susceptibility loci for type 2 diabetes and shared etiological pathways with coronary heart disease. *Nat. Genet.* 49, 1450–1457. <https://doi.org/10.1038/NG.3943>.
  37. Pulit, S.L., Stoneman, C., Morris, A.P., Wood, A.R., Glastonbury, C.A., Tyrrell, J., Yengo, L., Ferreira, T., Marouli, E., Ji, Y., et al. (2019). Meta-analysis of genome-wide association studies for body fat distribution in 694 649 individuals of European ancestry. *Hum. Mol. Genet.* 28, 166–174. <https://doi.org/10.1093/HMG/DDY327>.
  38. Justice, A.E., Karaderi, T., Highland, H.M., Young, K.L., Graff, M., Lu, Y., Turcot, V., Auer, P.L., Fine, R.S., Guo, X., et al. (2019). Protein-coding variants implicate novel genes related to lipid homeostasis contributing to body-fat distribution. *Nat. Genet.* 51, 452–469. <https://doi.org/10.1038/s41588-018-0334-2>.

39. Martinez-Huenchullan, S.F., Tam, C.S., Ban, L.A., Ehrenfeld-Slater, P., McLennan, S.V., and Twigg, S.M. (2020). Skeletal muscle adiponectin induction in obesity and exercise. *Metabolism* 102, 154008. <https://doi.org/10.1016/j.metabol.2019.154008>.
40. Dastani, Z., Hivert, M.F., Timpson, N., Perry, J.R.B., Yuan, X., Scott, R.A., Henneman, P., Heid, I.M., Kizer, J.R., Lytikäinen, L.P., et al. (2012). Novel loci for adiponectin levels and their influence on type 2 diabetes and metabolic traits: a multi-ethnic meta-analysis of 45,891 individuals. *PLoS Genet.* 8, e1002607. <https://doi.org/10.1371/JOURNAL.PGEN.1002607>.
41. Zheng, J.S., Arnett, D.K., Lee, Y.C., Shen, J., Parnell, L.D., Smith, C.E., Richardson, K., Li, D., Borecki, I.B., Ordovás, J.M., et al. (2013). Genome-wide contribution of genotype by environment interaction to variation of diabetes-related traits. *PLoS One* 8, e77442. <https://doi.org/10.1371/JOURNAL.PONE.0077442>.
42. Kathiresan, S., Willer, C.J., Peloso, G.M., Demissie, S., Musunuru, K., Schadt, E.E., Kaplan, L., Bennett, D., Li, Y., Tanaka, T., et al. (2009). Common variants at 30 loci contribute to polygenic dyslipidemia. *Nat. Genet.* 41, 56–65. <https://doi.org/10.1038/NG.291>.
43. Henkel, N.D., Wu, X., O'Donovan, S.M., Devine, E.A., Jiron, J.M., Rowland, L.M., Sarnyai, Z., Ramsey, A.J., Wen, Z., Hahn, M.K., et al. (2022). Schizophrenia: a disorder of broken brain bioenergetics. *Mol. Psychiatry* 27, 2393–2404. <https://doi.org/10.1038/s41380-022-01494-x>.
44. Vuckovic, D., Bao, E.L., Akbari, P., Lareau, C.A., Mousas, A., Jiang, T., Chen, M.H., Raffield, L.M., Tardaguila, M., Huffman, J.E., et al. (2020). The Polygenic and Monogenic Basis of Blood Traits and Diseases. *Cell* 182, 1214–1231.e11. <https://doi.org/10.1016/j.cell.2020.08.008>.
45. Chen, M.H., Raffield, L.M., Mousas, A., Sakaue, S., Huffman, J.E., Moscatti, A., Trivedi, B., Jiang, T., Akbari, P., Vuckovic, D., et al. (2020). Trans-ethnic and Ancestry-Specific Blood-Cell Genetics in 746,667 Individuals from 5 Global Populations. *Cell* 182, 1198–1213.e14. <https://doi.org/10.1016/j.cell.2020.06.045>.
46. Vinik, A.I., Erbas, T., Park, T.S., Nolan, R., and Pittenger, G.L. (2001). Platelet Dysfunction in Type 2 Diabetes. *Diabetes Care* 24, 1476–1485. <https://doi.org/10.2337/DIACARE.24.8.1476>.
47. Brown, B.C., and Knowles, D.A. (2020). Phenome-scale causal network discovery with bidirectional mediated Mendelian randomization. Preprint at bioRxiv. <https://doi.org/10.1101/2020.06.18.160176>.
48. Brown, B.C., and Knowles, D.A. (2021). Welch-weighted Egger regression reduces false positives due to correlated pleiotropy in Mendelian randomization. *Am. J. Hum. Genet.* 108, 2319–2335. <https://doi.org/10.1016/J.AJHG.2021.10.006>.
49. Dempster, A.P., Laird, N.M., and Rubin, D.B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *J. R. Stat. Soc. Ser. B Methodol.* 39, 1–22.
50. Parra, L.C. (2018). Multiset Canonical Correlation Analysis simply explained. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1802.03759>.
51. Witten, D.M., and Tibshirani, R.J. (2009). Extensions of Sparse Canonical Correlation Analysis with Applications to Genomic Data. *Stat. Appl. Genet. Mol. Biol.* 8, Article28. <https://doi.org/10.2202/1544-6115.1470>.
52. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine Learning in Python. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1201.0490>.
53. McFadden, D. (1973). Conditional logit analysis of qualitative choice behavior. In *Frontiers in Econometrics*, pp. 105–142.
54. Wu, D., and Smyth, G.K. (2012). Camera: A competitive gene set test accounting for inter-gene correlation. *Nucleic Acids Res.* 40, e133. <https://doi.org/10.1093/nar/gks461>.
55. Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J. (2015). Second-generation PLINK: Rising to the challenge of larger and richer datasets. *GigaScience* 4, 7–16. <https://doi.org/10.1186/s13742-015-0047-8>.
56. Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A., Vukcevic, D., Delaneau, O., O'Connell, J., et al. (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature* 562, 203–209. <https://doi.org/10.1038/s41586-018-0579-z>.
57. Kasela S., Aguet F., Kim-Hellmuth S., Brown B.C., Nachun D.C., Tracy R.P., Durda P., Liu Y., Taylor K.D., Johnson W.C., et al. Interaction molecular QTL mapping discovers cellular and environmental modifiers of genetic regulatory effects. bioRxiv 2022. doi:10.1101/2023.06.26.546528. <https://www.biorxiv.org/content/10.1101/2023.06.26.546528v1>
58. Stilp, A.M., Emery, L.S., Broome, J.G., Buth, E.J., Khan, A.T., Laurie, C.A., Wang, F.F., Wong, Q., Chen, D., D'Augustine, C.M., et al. (2021). A System for Phenotype Harmonization in the National Heart, Lung, and Blood Institute Trans-Omics for Precision Medicine (TOPMed) Program. *Am. J. Epidemiol.* 190, 1977–1992. <https://doi.org/10.1093/aje/kwab115>.

## STAR★METHODS

### KEY RESOURCES TABLE

| REAGENT or RESOURCE                           | SOURCE  | IDENTIFIER  |
|---|---------|---|
| Deposited data                                |         |   |
| MESA TOPMed multi-omics pilot data            | dbGaP   | dbGaP: phs001416.v3.p1  |
| Software and algorithms                       |         |   |
| Multiset Correlation and Factor Analysis      | Zenodo  | <a href="https://doi.org/10.5281/zenodo.7951370">https://doi.org/10.5281/zenodo.7951370</a>   |
| MOFA2   | github  | <a href="https://github.com/bioFAM/MOFA2">https://github.com/bioFAM/MOFA2</a>                 |
| eQTLgen <i>trans</i> -eQTL summary statistics | eQTLgen | <a href="https://www.eqtlgen.org/trans-eqtl.html">https://www.eqtlgen.org/trans-eqtl.html</a> |

### RESOURCE AVAILABILITY

#### Lead contact

Further information and requests for resources should be directed to and will be fulfilled by the lead contact, Dr. Brielin Brown ([bbrown@nygenome.org](mailto:bbrown@nygenome.org)).

#### Materials availability

This study did not generate new unique reagents.

#### Data and code availability

The MESA TOPMed multi-omics pilot data have been deposited on dbGap and are publicly available as of the date of publication. The accession number is listed in the [key resources table](#). All original code has been deposited on zenodo and is publicly available as of the date of publication. The DOI is listed in the [key resources table](#). The code is also available on github at <https://github.com/collinwa/MCFA>. Any additional information required to reanalyze the data reported in this paper is available from the [lead contact](#) upon request.

### METHOD DETAILS

#### Multiset correlation and factor analysis

Let  $Y = \{Y_m\}_{m=1}^M$  be a set of  $N \times p_m$  observed data matrices:  $N$  individuals measured in  $M$  data modalities consisting of  $p_m$  features each. We model each observed mode as having contributions from two low-dimensional hidden factors ([Figures 1A](#) and [S8](#))

$$z_n \sim N(0, I_d)$$

$$x_n^m \sim N(0, I_{k_m})$$

$$y_n^m \sim N(W_m z_n + L_m x_n^m, \Psi_m)$$

where  $d$  is the shared hidden dimensionality,  $k_m$  are the dataset-private hidden dimensionalities,  $W_m$  are  $p_m \times d$  shared space loading matrices,  $L_m$  are  $p_m \times k_m$  private space loading matrices and  $\Psi_m = \text{diag}(\psi_m^1, \dots, \psi_m^{p_m})$  are the diagonal residual covariance matrices. Given  $Y$ ,  $d$  and  $k_m$ , our goal is to infer the hidden factors  $Z$  and  $X_m$  and loading matrices  $W_m$  and  $L_m$ . This can be accomplished using a straightforward application of expectation maximization (EM).<sup>49</sup> For a derivation of the EM update equations, as well as a more detailed exposition including the relationship to pCCA, factor analysis and other multiset CCA (MCCA) methods, see [Methods S1](#). In practice, we center and scale all data variables. This is not strictly required, however it enables simple estimation of the number of PCs to include and simplifies explained variance calculations, see below.

#### Model initialization

An important aspect of EM optimization is choosing a good initialization. We benchmarked three approaches to initializing  $W$ : random initialization and two versions of MCCA that correspond to maximizing the sum of pairwise correlations with the average variance and average norm constraints. These MCCA formulations can be solved via simple eigendecompositions. We found that the sum of

pairwise correlations with average variance constraint produced the best initial estimates (Figure S7). This can be solved with a simple two step procedure: 1) whiten each data matrix using the singular value decomposition (SVD), 2) perform a second SVD on the concatenated whitened data matrices<sup>50</sup>:

**Input:**  $Y_1, \dots, Y_M, d$   
**Result:**  $\widehat{W} = [W_1^T : \dots : W_M^T]^T$   
 $U_{all} \leftarrow \text{concatenate}(\text{SVD}(Y_1). U, \dots, \text{SVD}(Y_M). U);$   
 $\widehat{W} \leftarrow \text{SVD}(U_{all}). V[:, 0 : d];$   
 $\widehat{\rho} \leftarrow \text{SVD}(U_{all}). \lambda[0 : d];$   
**return**  $\widehat{W}, \widehat{\rho}$

We initialize  $L$  and  $\Psi$  using probabilistic PCA on the residual data matrices after fitting MCCA. Specifically:

**Input:**  $Y_i, W_i, N, k_i$   
**Result:**  $\widehat{L}_i, \widehat{\Psi}_i$   
 $\Sigma_i^\perp \leftarrow Y_i^T Y_i / N - W_i W_i^T;$   
 $\widehat{L}_i \leftarrow \text{eigh}(\Sigma_i^\perp). V[:, 0 : k_i];$   
 $\sigma^2 \leftarrow \text{mean}(\text{eigh}(\Sigma_i^\perp). \lambda[k_i:]);$   
 $\widehat{\Psi}_i \leftarrow \sigma^2 \mathbf{1}_{k_i};$   
**return**  $\widehat{L}_i, \widehat{\Psi}_i$

### High dimensionality and selection of hyperparameters

There are two primary approaches to control for over-fitting in applications of CCA-type methods to high-dimensional ( $N \ll p$ ) problems. The first is to use penalized optimization techniques, where the objective function additionally contains an  $l_1$  constraint on the weight matrices.<sup>51</sup> The second is to project each dataset onto its informative principal components.<sup>6,7,11</sup> In this application, we choose the latter approach in order to find components with broad effects on the structure of the data, rather than specific effects on small numbers of molecular features.<sup>11</sup> We choose the number of principal components of each dataset using the Marchenko-Pasteur law,<sup>13</sup> which states that for mean 0, variance 1 data, principal components with corresponding eigenvalues above  $\lambda_m = 1 + \sqrt{\rho_m/N}$  should be considered non-noise. We are not aware of a corresponding law for the cross-covariance matrices used in CCA, however, the empirical spectral distribution of the cross-covariance of matrices of random noise can be easily estimated in practice:

**Input:**  $N, k = \{k_m\}_{m=1}^M, n_{it}$   
**Result:**  $\rho$   
**for**  $it \leftarrow 0$  **to**  $n_{it}$  **do**  
  **for**  $k_m \in k$  **do**  
     $[Y_m]_{i=1}^{N, k_m} \sim N(0, 1);$   
  **end**  
   $\rho[it] \leftarrow \max(\text{InitializeMCFA}(Y_1, \dots, Y_M). \rho)$   
**end**  
**return**  $\text{mean}(\rho)$

Then we keep all components where  $\rho_{init} > \rho$ .

### Calculating the variance explained

The linear-Gaussian nature of the model simplifies estimation of the variance explained. That is, if the features of each mode  $Y_m^{(:,j)}$  are normalized to variance 1, the model  $Y_m^{(:,j)} = \sum_d W_m^{(j,d)} Z^{(:,d)} + \sum_{k_m} L_m^{(j,k_m)} X_m^{(:,k_m)} + \epsilon$  implies that the variance in feature  $j$  of mode  $m$  explained by shared factor  $d$  is  $W_m^{(j,d)2}$ . Likewise, the variance explained by the  $k_m$ -th private factor of mode  $m$  is  $L_m^{(j,k_m)2}$ . The total variance in mode  $m$  explained by a given shared factor  $d$  (respectively, private factor  $k_m$ ) is thus given by  $\sum_j W_m^{(j,d)2}$  (respectively,  $\sum_j L_m^{(j,k_m)2}$ ), and the total variance in the mode explained by the factors are  $\sum_{j,d} W_m^{(j,d)2}$  and  $\sum_{j,k_m} L_m^{(j,k_m)2}$ , respectively. Note that when working in PC-space, the raw  $W$  and  $L$  features correspond to variance in PCs explained, rather than modality features. Thus, we calculate the variance explained after projecting back into the original feature space  $W_m \leftarrow V_m W_m, L_m \leftarrow V_m L_m$  where  $V_m$  are the right singular vectors of mode  $m$ .

To calculate the variance in a metadata feature explained by a particular space, we regressed the trait value  $T$  on the shared or private space,  $T \sim Z$  or  $T \sim X_m$ . For continuous-valued traits we used linear regression as implemented in SciKitLearn v1.0 `linear_model.LinearRegression` and report the coefficient of determination.<sup>52</sup> For discrete-valued traits, we used multinomial logistic regression as implemented in SciKitLearn v1.0 `linear_model.LogisticRegression`.<sup>52</sup> We fit two models: a null model including only intercept or intercept and site, and one including the factor variables. We report the variance explained as the McFadden pseudo- $R^2$ ,  $1 - \frac{\ln_{null}}{\ln_{alt}}$ , with  $\ln_{null}$  and  $\ln_{alt}$  being the model negative log likelihood for the null and alternative model respectively.<sup>53</sup>

### Calculating relative feature importance

Feature importance in traditional CCA is defined by the correlation of the variables in the reduced space  $\rho = \text{cor}(Y_1f_1, Y_2f_2)$ . Unfortunately this notion breaks down in higher dimensions. As we discuss further in Methods S1, the degree of sharing in MCCA is defined by functions of the cross-correlation matrix in the reduced space,

$$S = \text{cor}(Y_1f_1, \dots, Y_mf_m) \in \mathbb{R}^{m \times m}.$$

We seek to define an analogous quantity for our graphical model. In MCFA, the data in the reduced (shared) space is given by the posterior mean of  $Z$ ,  $\hat{Z} = E[Z|W, \Psi, L, Y] = Y(WW^T + LL^T\Psi)^{-1}W$ . We can also calculate the posterior mean of  $Z$  conditional on observing a single mode,  $\hat{Z}_m = E[Z|W_m, \Psi_m, L_m, Y_m] = Y_m(W_mW_m^T + L_mL_m^T\Psi_m)^{-1}W_m$ . This latter quantity is analogous to the reduced variables  $Y_mf_m$  in MCCA. Thus we can summarize the importance of each dimension of the shared space by calculating functions of the cross-correlation of columns of  $\hat{Z}_m$ ,

$$S_d = \text{cor}(\hat{Z}_1^{(:,d)}, \dots, \hat{Z}_m^{(:,d)}).$$

The relevant function in our model is the generalized variance  $|S|$ , see [Methods S1](#). The determinant of a correlation matrix is bounded between 0 and 1, with lower values indicating more correlation, and higher values less. Thus to aid interpretability, we report  $\rho_d = -\log|S_d|$  and reorder columns of  $Z$  and  $W$  with decreasing  $\rho_d$ .

### SNP set enrichment analysis

For SNP set enrichment analysis, we broadly follow the approach of CAMERA.<sup>54</sup> In brief, enrichment statistics can be inflated due to correlations in the sample - in this case, linkage disequilibrium between two GWAS SNPs. This results in an under-estimate of the standard error of the enrichment test statistic and an increase in false positives. We calculate the variance inflation factor by using plink v 1.9<sup>55</sup> to estimate linkage disequilibrium between annotation SNPs in 337,781 unrelated individuals from the UK Biobank.<sup>56</sup> The variance inflation factor is  $\nu = 1 + (\rho_A - 1)\rho_A$ , with  $\rho_A$  the average person correlation between features in set A. We test the known GWAS mean  $\chi^2$  statistic  $h_0 : \chi_A^2 = 1$  against the alternative  $h_1 : \chi_A^2 > 1$ . The standard error of the test statistic is  $\sigma_t = \sigma \sqrt{\frac{\nu}{\rho_A} - \frac{1}{\rho_m}}$  with  $\sigma$  the pooled empirical standard deviation of the test statistics.

### The MESA multi-omics pilot

The Multi-Ethnic Study of Atherosclerosis (MESA) is a prospective cohort study with the goal to identify progression of subclinical atherosclerosis.<sup>14</sup> MESA recruited 6,814 participants, ages 45–84 years and free of clinical cardiovascular disease, during 2000–2002. The participants are 53% female, 38% non-Hispanic white, 28% Black, 22% Hispanic and 12% Asian-American. The Multi-Omics pilot dataset includes 30x whole genome sequencing (WGS) through the Trans-Omics for Precision Medicine (TOPMed) Project.<sup>15</sup>

Blood samples for multi-omic analysis of participants were collected at two time points (exam 1 and exam 5). RNA expression was profiled using poly-A RNA sequencing of PBMCs, and methylation was quantified by the Illumina 750K EPIC array in whole blood. The levels of 1,305 proteins were measured from plasma samples using the standard SOMAscan DNA aptamer-based platform, and metabolite levels were determined from targeted and untargeted mass spectrometry of blood plasma. The MESA Multi-Omics pilot biospecimen collection, molecular phenotype data production and quality control (QC) are described in detail in Kasela et al.<sup>57</sup>

### Cross-validation

We used leave-one-out cross-validation (CV) to evaluate our model. The primary reason we chose leave-one-out CV over  $k$ -fold CV is that our hyperparameter selection method depends on the sample size. With  $n - 1$  individuals, the same parameters used for the full inference procedure are likely to be valid. For small  $k$ , fitting with  $\frac{k-1}{k}n$  individuals while using the same number of PCs may result in over-fitting in the training set, and using a smaller number of PCs may not capture the same variation as the full model.

To perform cross-validation we hold out a set of individuals, fit the MCFA model, then project the held out individuals into the learned space. If  $W_{tr}, L_{tr}$  and  $\Psi_{tr}$  are the model parameters learned from the training set, the projections of the test data into the learned spaces are given by

$$\hat{Z}_{te} = Y_{te}(W_{tr}W_{tr}^T + L_{tr}L_{tr}^T\Psi_{tr})^{-1}W_{tr} \quad \hat{X}_{te} = Y_{te}(W_{tr}W_{tr}^T + L_{tr}L_{tr}^T\Psi_{tr})^{-1}L_{tr}$$

The full data reconstruction is

$$\hat{Y}_{te} = \hat{Z}_{te}W_{tr}^T + \hat{X}_{te}L_{tr}^T$$

We evaluate model fit by calculating the normalized root mean squared error (NRMSE). In order to provide a fair evaluation across modes with a highly variable number of features, we calculate NRMSE on a per mode basis

$$NRMSE = \sqrt{\frac{1}{p_m} \sum_{i=1}^{p_m} \frac{(Y_m^{(:,i)} - \hat{Y}_m^{(:,i)})^2}{\text{var}Y_m^{(:,i)}}$$

and potential over-fitting can be assessed by comparing the median training set NRMSE against the median test set NRMSE over many cross-validation iterations.

### Comparison to MOFA2 and MMAE

We installed MOFA2 version 0.6.7 using `pip install mofapy2`. We used the options `scale_groups = False`, `scale_views = False`, `ard_weights = True` and `spikeslab_weights = True`. We set the convergence tolerance to `convergence_mode = 'medium'`. For comparison purposes we set the number of factors equal to the hidden dimensionality inferred by MCFA (`factors = 14`).

Our multi-modal auto-encoder architecture is visualized in [Figure S4](#). We used two hidden layers per dataset, with the first layer having dimensionality equal to 8 times that modalities MCFA-inferred number of PCs, and the second layer having dimensionality equal to that modalities MCFA-inferred number of PCs. These layers are then concatenated, and sent through an additional hidden layer with 8 times the MCFA-inferred number of shared dimensions to the final 14-dimensional encoded representation. All layers except the final encoder layer consist of a linear transform followed by ReLU activation, while the final encoder layer omits the ReLU activation. The decoder had identical architecture to the encoder only reversed. The network was implemented in `pytorch` v1.11.0 and optimized with `Adagrad` using 10 batches per epoch until the NRMSE change relative to the total loss was less than  $10^{-6}$ .

### QUANTIFICATION AND STATISTICAL ANALYSIS

We analyzed individuals from Exam 1 where all five data types were collected and passed QC. All data modalities were inverse rank normalized prior to sample filtering based on the availability of other data types. There were 614 individuals with observations of WGS, RNA-seq, methylation, metabolomics and proteomics that all pass QC. We further removed all features (CpGs, genes, proteins) located on sex-chromosomes, 0-variance features, CpGs with missing data, and CpGs where the probe was within 5 bases of an SNP, leaving us with 6,042 metabolites, 1,222 proteins, 19,034 genes, and 724,210 CpGs. We analyzed 28 PCs of RNA expression, 39 PCs of methylation, 27 PCs of protein expression and 63 PCs of metabolite, as determined using the aforementioned method. For sample metadata, we leveraged the rich phenotype data available in MESA that were harmonized by the TOPMed Data Coordinating Center.<sup>58</sup> For details on the estimation of sample cell-type proportions from methylation and RNA-seq data, see Kasela et al.<sup>57</sup> Genetic association analyses were conducted using `plink` v 1.9<sup>55</sup> while controlling for site, age, sex and 11 genotype PCs; reported *p*-values are uncorrected and tested against a null of 0 effect. SNP set enrichment significance was defined as having an FDR *q*-value below 0.05 when corrected for 3 tested sets across 14 factors tested against the null hypothesis that the mean  $\chi^2$  test statistic is 1.