

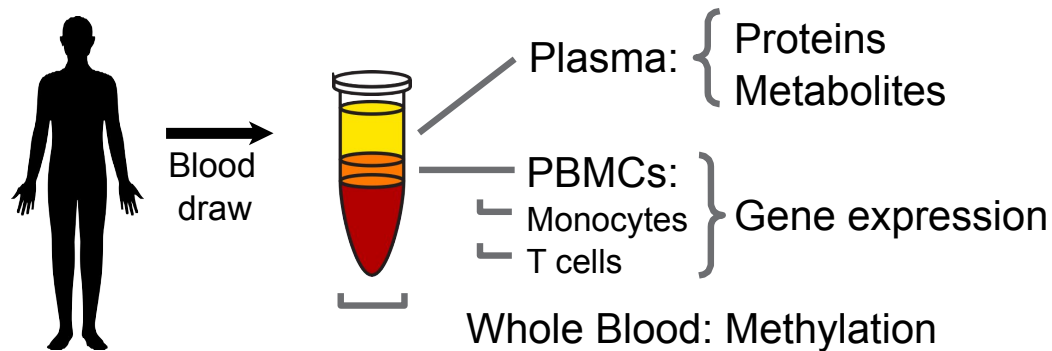
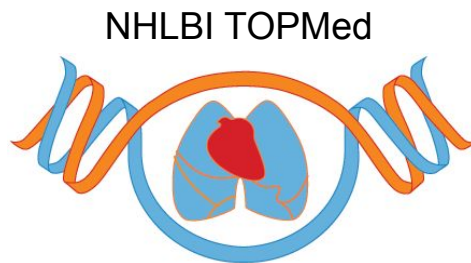
Multi-set correlation and factor analysis of multi-omic data

Brielin C Brown PhD

Postdoctoral researcher
Knowles and Lappalainen labs
New York Genome Center

DSI Fellow
Data Science Institute
Columbia University

Multi-omics data promise to revolutionize biomedical research



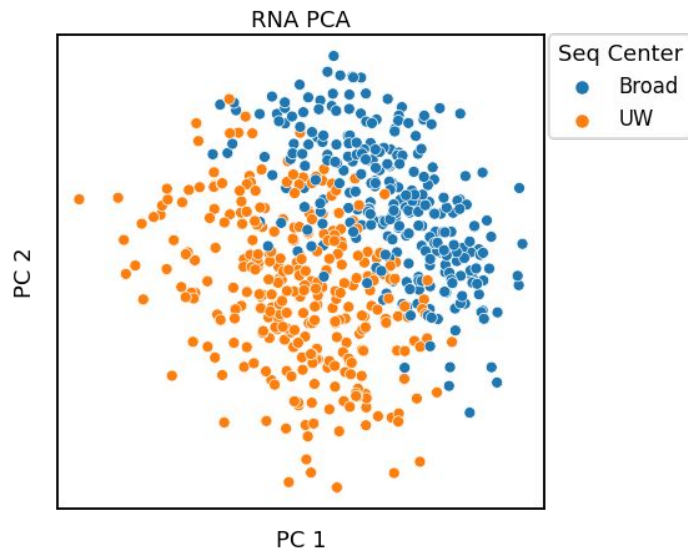
Multi-*omic* but also

- Multi-*phenotype*
- Multi-*tissue*
- Multi-*ethnic*
- Multi-*center*

NHLBI TOPMed

- Generating multi-omic data for tens of thousands of patients
- MESA multi-omic pilot

How can we think about exploring genomic data?



How can we think about exploring genomic data?

RNA levels

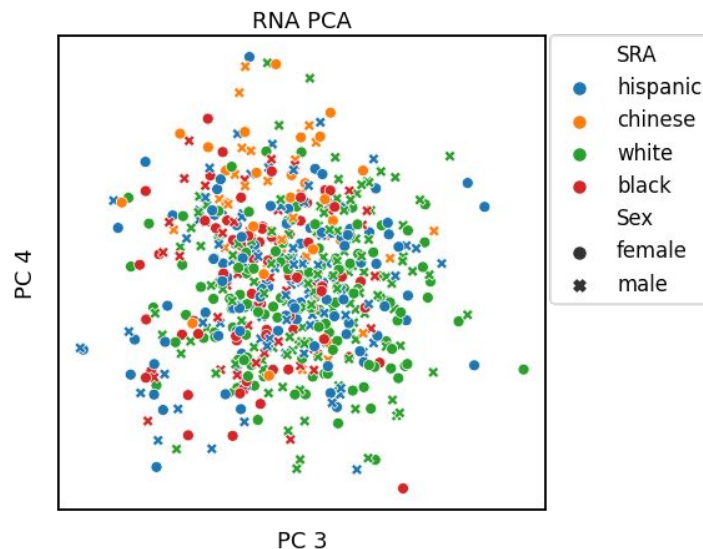
RNA loadings

$$Y = XL^T + \epsilon$$

Low-dim features Residual

Includes

- Principal Components Analysis (PCA)
- Factor Analysis (FA)



How can we simultaneously explore two datasets?

RNA levels $Y_1 = ZW_1^\top + \epsilon_1$

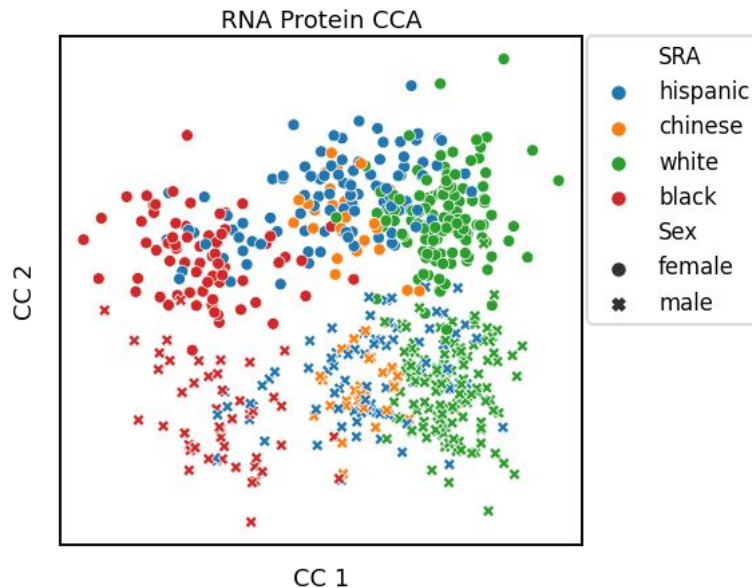
Protein levels $Y_2 = ZW_2^\top + \epsilon_2$

Feature loadings Residual

Low-dim shared factors

Includes

- Group Factor Analysis (GFA)
- Canonical Correlation Analysis (CCA)
- *Probabilistic CCA*

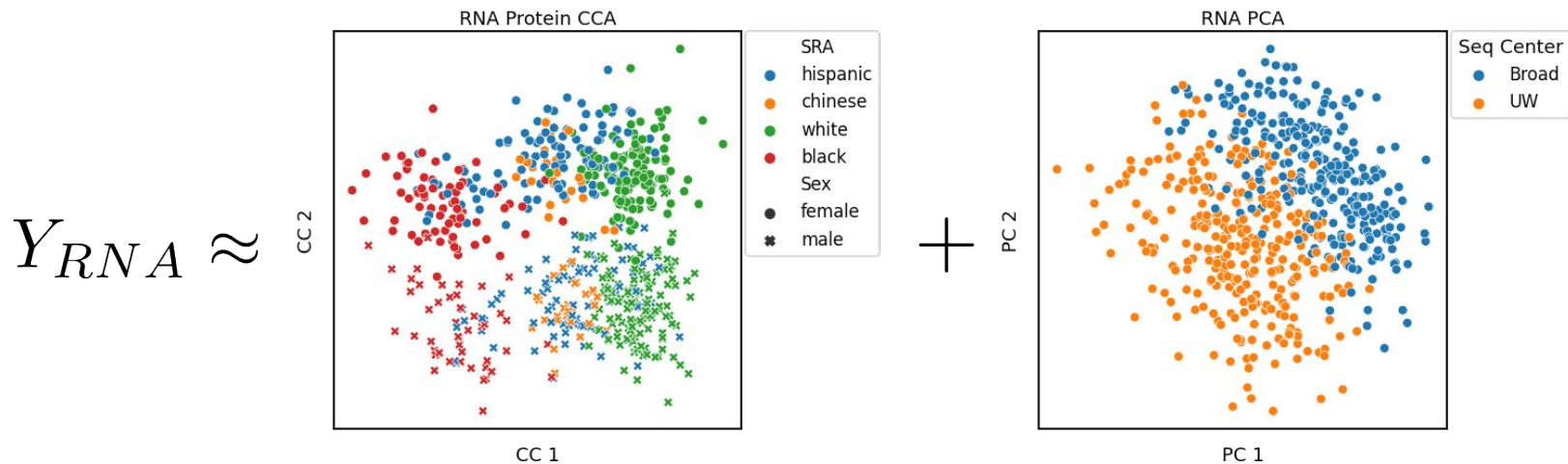


Multiset Correlation and Factor Analysis combines these

$$\begin{array}{c}
 Y_1 \approx ZW_1^\top + X_1L_1^\top + \epsilon_1 \\
 \vdots \\
 Y_m \approx ZW_m^\top + X_mL_m^\top + \epsilon_m
 \end{array}$$

Multi-omic data

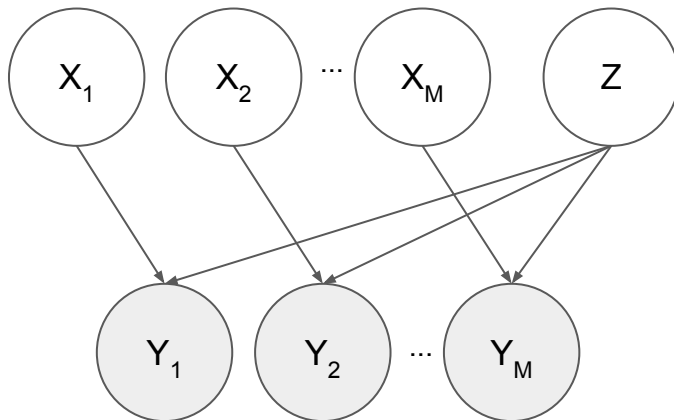
Shared factors and loadings Private factors and loadings Residual



Multiset Correlation and Factor Analysis combines these

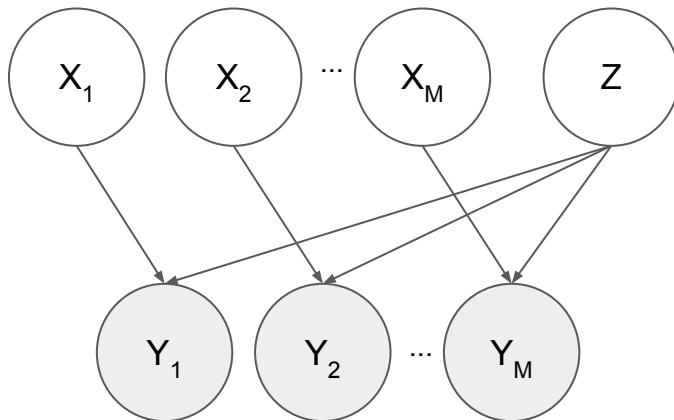
Multi-omic data

$$\begin{matrix} Y_1 \\ \vdots \\ Y_m \end{matrix} \approx \underbrace{ZW_1^\top \vdots W_m^\top}_{\text{Shared factors and loadings}} + \underbrace{X_1 L_1^\top \vdots X_m L_m^\top}_{\text{Private factors and loadings}} + \underbrace{\epsilon_1 \vdots \epsilon_m}_{\text{Residual}}$$

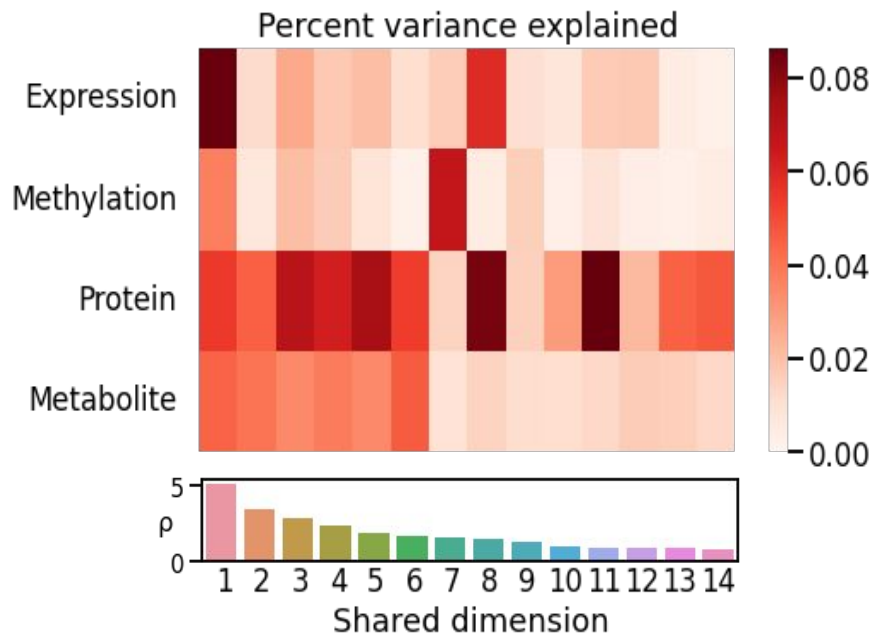
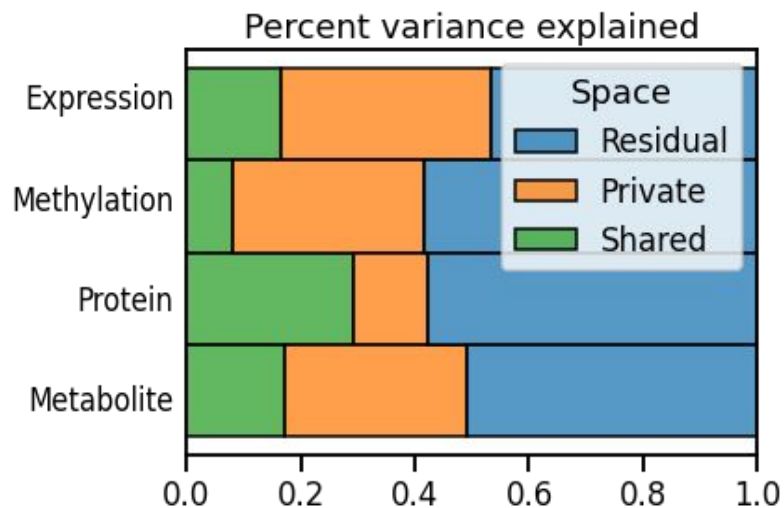


MCFA has numerous advantages

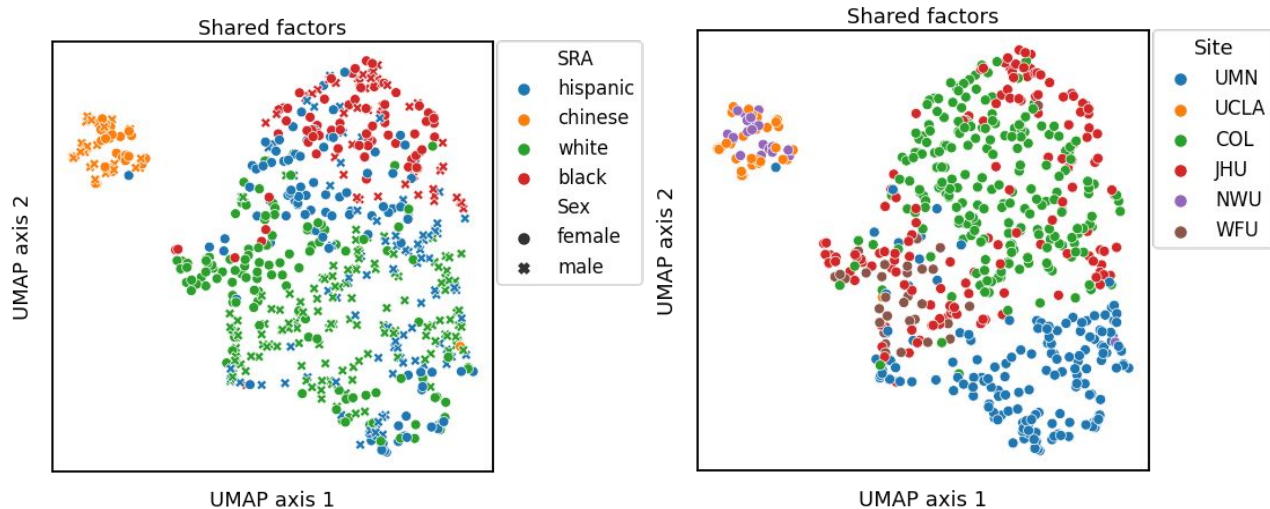
- Elegant approach to multiset CCA
 - Strong theoretical connection
- Calculations done in PC-space
 - Extremely efficient, controls over-fitting
- No hyperparameters to tune
 - Self-selected with random matrix techniques
- Works with any data type
 - No “gene scores”
- Factor loadings are interpretable
 - Similar to regression coefficients
- Fully unsupervised!
 - Supports exploratory analysis



MCFA finds structure in 614 multi-omic samples

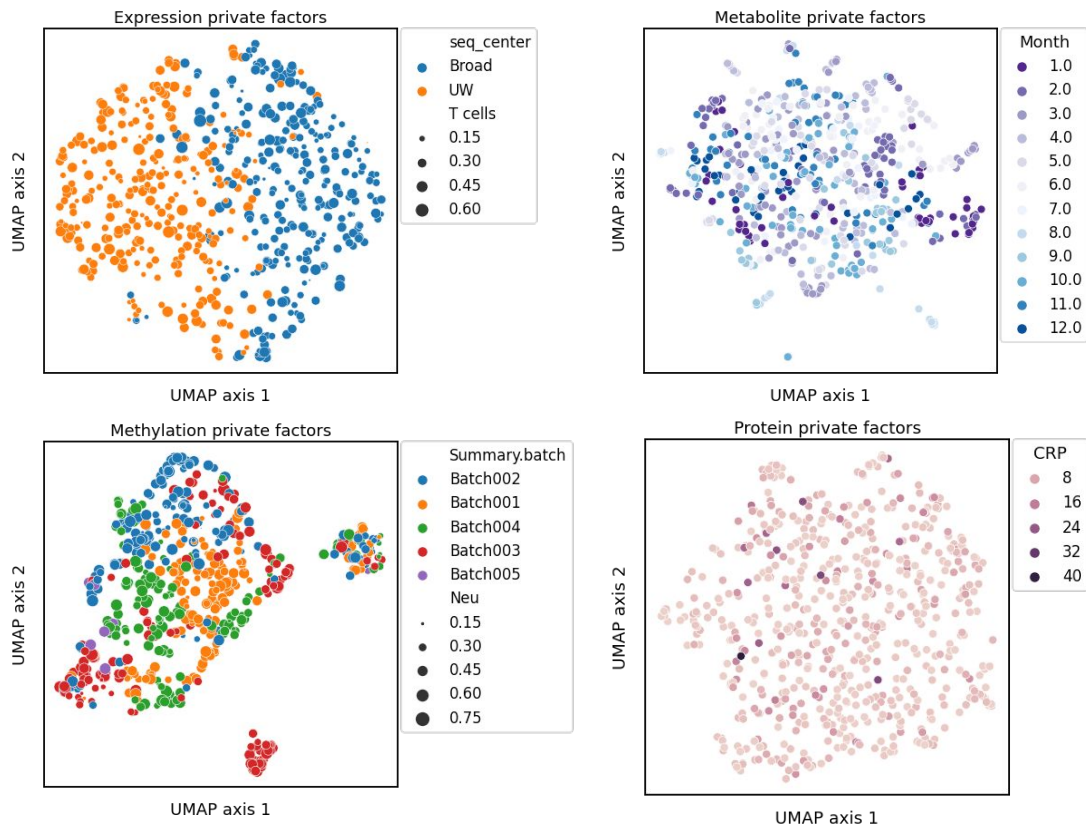


Shared structure reflects demographics and site

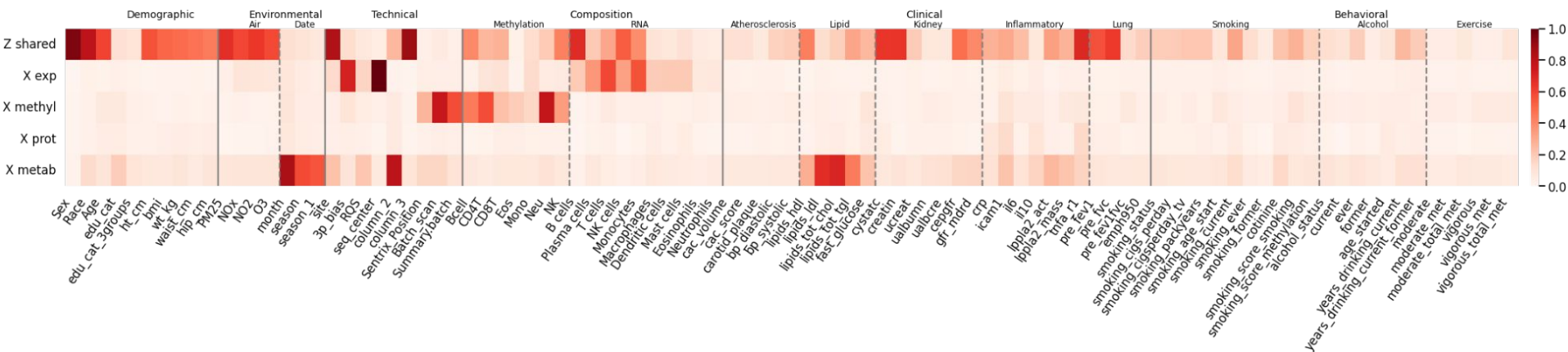


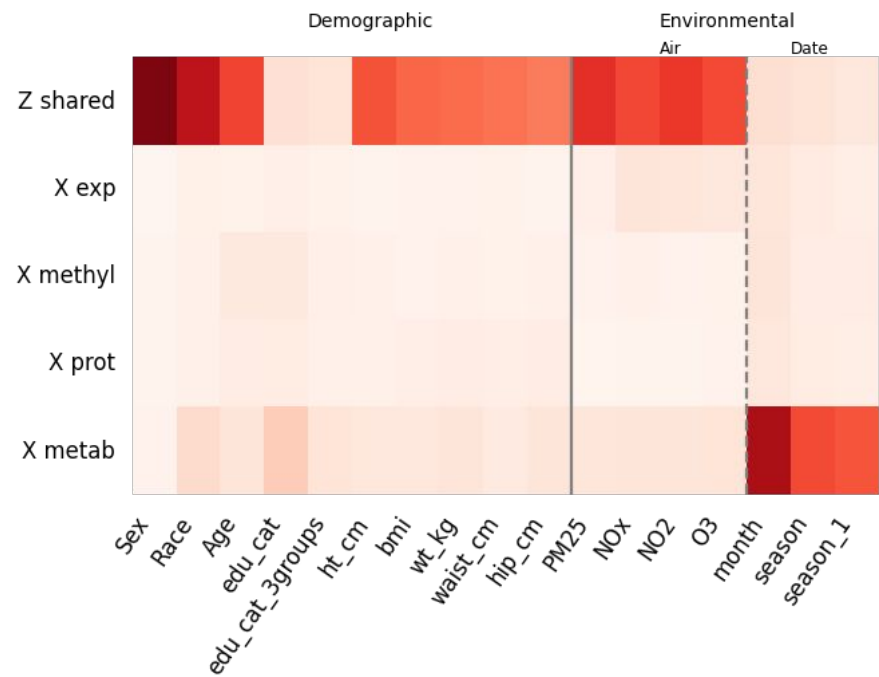
This is without genetic information

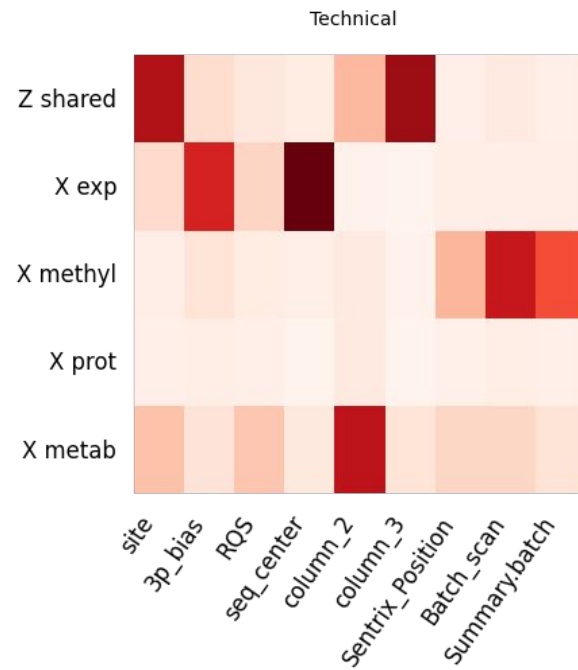
Private structure captures omic-specific effects

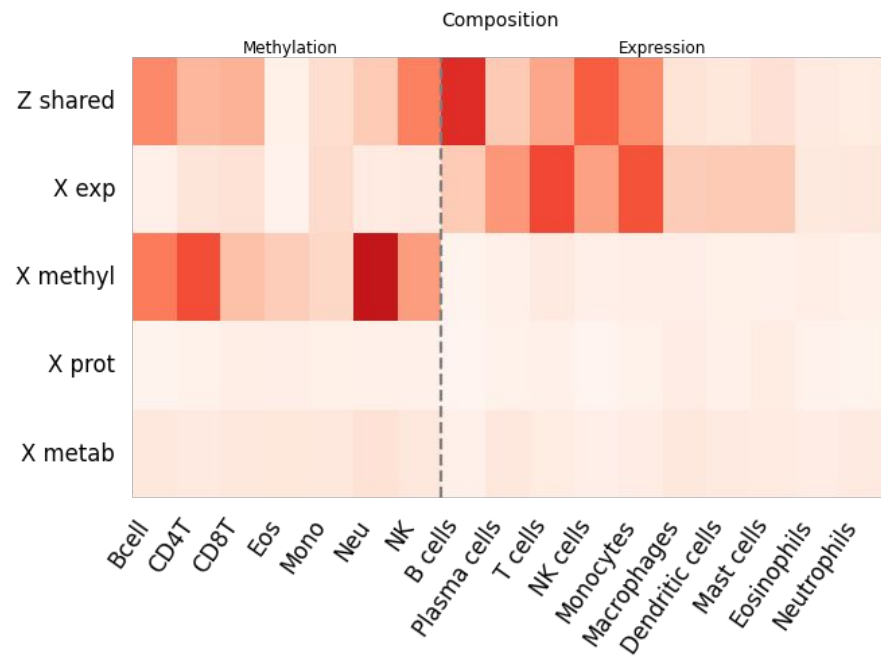


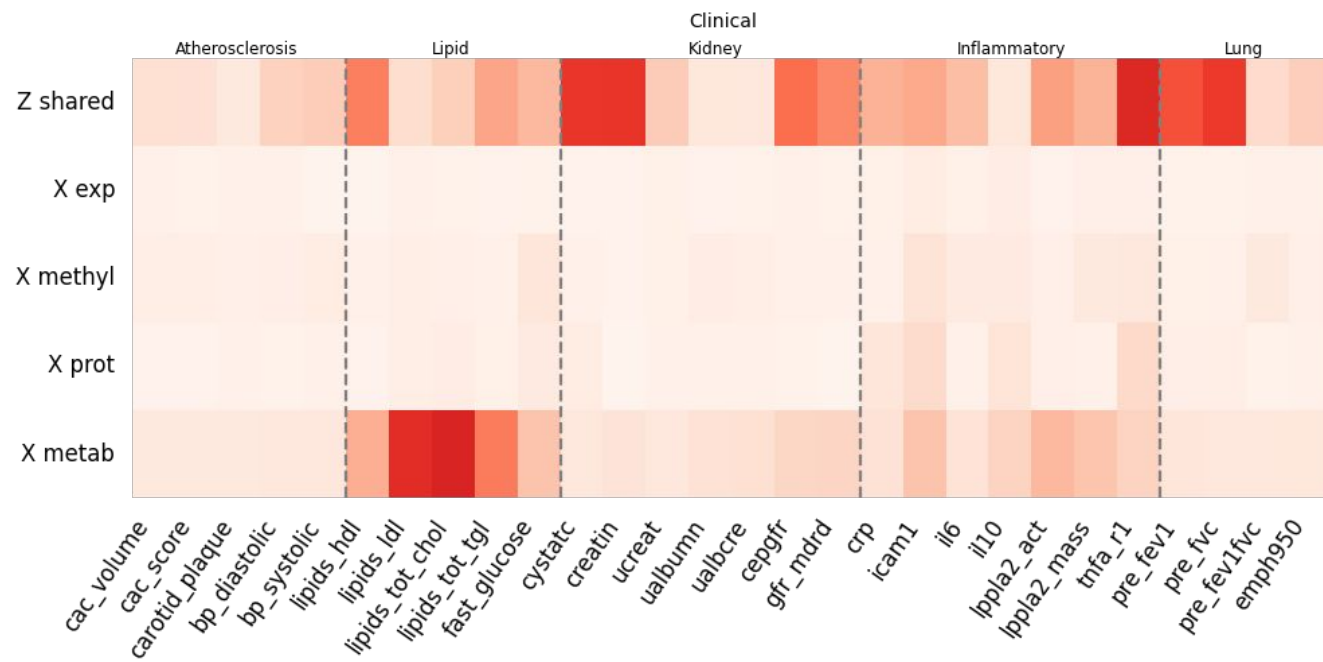
MCFA further captures environmental differences and clinical biomarkers







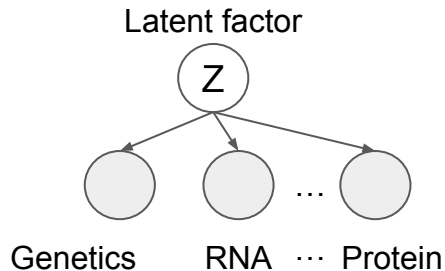




There are multiple ways to integrate genetic data

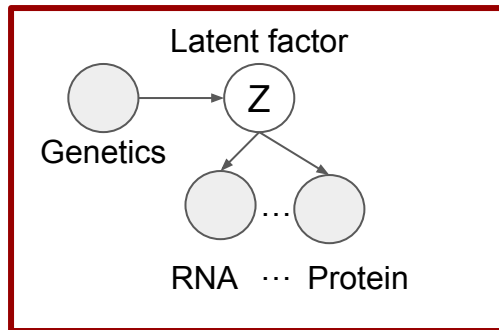
1. Include genetics in the MCFA model

- Somatic variation
- Copy number variation
- Known functional SNPs (eQTLs, etc)



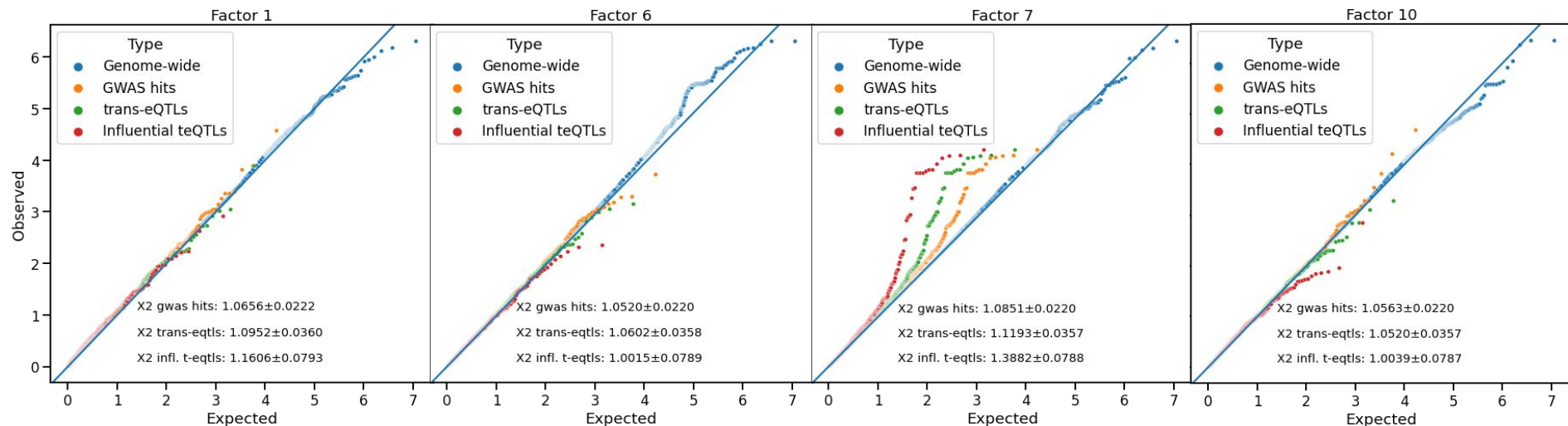
2. Think of factors as “latent phenotypes”

- Factors learned on molecular data may represent coordinated biological activity
- These may be affected by genetic variation
- Look for genetic associations with latent phenotypes (GWAS)



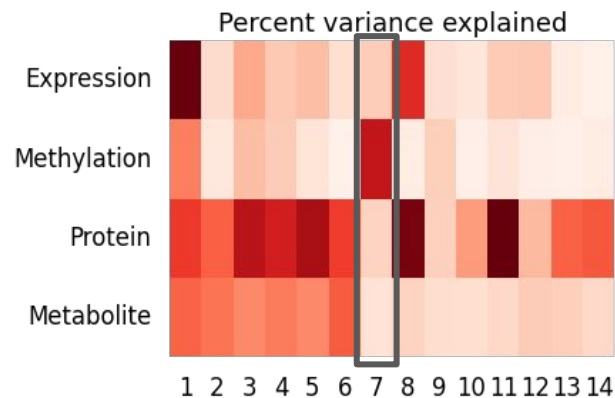
GWAS can be used to integrate genetic data

Are genetic associations with shared factors enriched for known GWAS hits or *trans*-eQTLs?



Top SNPs implicate lipid metabolism near FADS1, FADS2 genes

	Trait	rs ID	Chr	Position	Gene	p-val_F7
	Red blood cell fatty acid levels	rs174541	11	61565908	FADS2	0.000043
	Metabolite levels	rs174549	11	61571382	FADS1, FADS2	0.000056
	Blood metabolite levels	rs174556	11	61580635	FADS1, FADS2	0.000057
	Height	rs7184046	15	75866150	PTPN9	0.000061
	Plasma omega-6 polyunsaturated fatty...	rs174555	11	61579760	FADS1, FADS2	0.000064
	Phospholipid levels (plasma)	rs1535	11	61597972	FADS2	0.000083
	Plasma omega-6 polyunsaturated fatty...	rs174537	11	61552680	MYRF, TMEM258	0.000109
	Phospholipid levels (plasma)	rs174536	11	61551927	MYRF, TMEM258	0.000109
	Cholesterol, total	rs174554	11	61579463	FADS1, FADS2	0.000117
	Plasma omega-6 polyunsaturated fatty...	rs174547	11	61570783	FADS1, FADS2	0.000127
	Plasma omega-6 polyunsaturated fatty...	rs174550	11	61571478	FADS1, FADS2	0.000127
	Plasma omega-6 polyunsaturated fatty...	rs174546	11	61569830	FADS1, FADS2	0.000127
	Red blood cell fatty acid levels	rs174545	11	61569306	FADS1, FADS2	0.000127
	Blood metabolite ratios	rs174548	11	61571348	FADS1, FADS2	0.000253
	Triglycerides	rs174529	11	61543961	MYRF, TMEM258	0.000264
	Systemic lupus erythematosus	rs4852324	2	74202578	DGUOK-AS1	0.000343
	Trans fatty acid levels	rs174583	11	61609750	FADS2	0.000405
	Mature red cell;HGB	rs1256061	14	64703593	ESR2	0.000464
	IgG glycosylation	rs2186369	22	24170996	SMARCB1	0.000505
	Plasma omega-6 polyunsaturated fatty...	rs2727270	11	61603237	FADS2	0.000774



GWAS SNPs associated with Factor 7 are also associated with methylation level

	Trait	rsid	Chr	Position	Gene	p-val_trait	p-val_F7	CpG	p-val_mqtl	n_snps
	Red blood cell fatty acid levels	rs174541	11	61565908	FADS2	3.000000e-19	0.000043	cg19610905	3.964333e-115	29
	Genes & Nutrition				PTPN9	2.000000e-10	0.000061	cg01268058	1.333136e-48	1

Losol et al. *Genes & Nutrition* (2019) 14:20
<https://doi.org/10.1186/s12263-019-0644-8>


OPEN ACCESS Freely available online

PLOS ONE

RESEARCH

Open

Effect of gestational oily fish intake on the risk of allergy in children may be influenced by *FADS1/2*, *ELOVL5* expression and DNA methylation

Purevsuren Losol^{1,2}, Faisal I. Rezwan¹, Veeresh K. Patil^{3,4}, Carina Venter³, Susan Ewart⁵, Hongmei Zhang⁶, S. Hasan Arshad^{3,4}, Wilfried Karmaus⁶ and John W. Holloway^{1,4*} 

Gomez-Alonso et al. *Clin Epigenet* (2021) 13:7
<https://doi.org/10.1186/s13148-020-00957-8>

Supplementation with N-3 Long-Chain Polyunsaturated Fatty Acids or Olive Oil in Men and Women with Renal Disease Induces Differential Changes in the DNA Methylation of FADS2 and ELOVL5 in Peripheral Blood Mononuclear Cells

Samuel P. Hoile¹, Rebecca Clarke-Harris¹, Rae-Chi Huang², Philip C. Calder^{1,3}, Trevor A. Mori⁴, Lawrence J. Beilin⁴, Karen A. Lillycrop⁵, Graham C. Burdge^{1*}

Clinical Epigenetics

RESEARCH

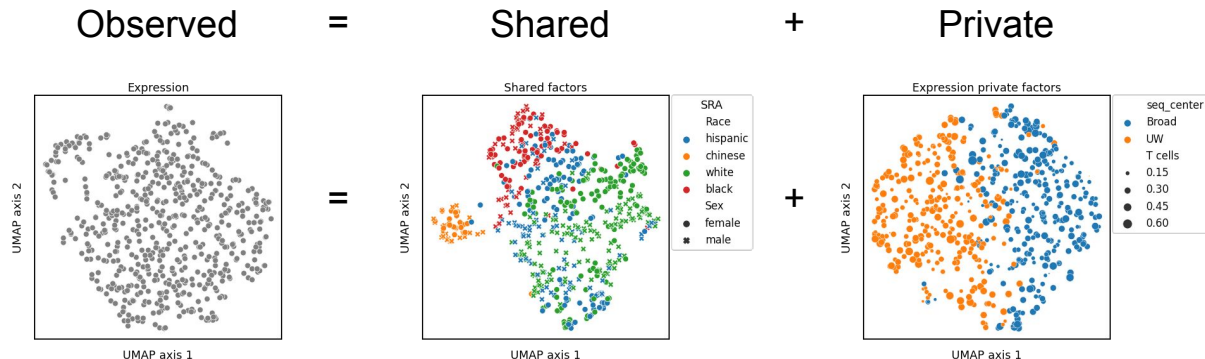
Open Access

DNA methylation and lipid metabolism: an EWAS of 226 metabolic measures



Monica del C. Gomez-Alonso^{1,2†}, Anja Kretschmer^{1,2†}, Rory Wilson^{1,2}, Liliane Pfeiffer^{1,2}, Ville Karhunen^{3,10}, Ilkka Seppälä⁴, Weihua Zhang^{3,5}, Kirstin Mittelstraß^{1,2}, Simone Wahl^{1,2}, Pamela R. Matias-Garcia^{1,2}, Holger Prokisch^{6,7}, Sacha Horn^{1,2}, Thomas Meitinger^{6,7,8}, Luis R. Serrano-Garcia^{1,2,9}, Sylvain Sebert¹⁰, Olli Raitakari^{11,12,13}, Marie Loh^{3,14}, Wolfgang Rathmann^{15,16}, Martina Müller-Nurasyid^{17,18,30}, Christian Herder^{16,19,20}, Michael Roden^{16,19,20}, Mikko Hurme²¹, Marjo-Riitta Jarvelin^{3,10,22,23}, Mika Ala-Korpela^{10,24,25,26}, Jaspal S. Kooner^{5,27,28,29}, Annette Peters², Terho Lehtimäki⁴, John C. Chambers^{3,5,14,28,29}, Christian Gieger^{1,2}, Johannes Kettunen^{10,25,26} and Melanie Waldenberger^{1,2,8*} 

MCFA is a powerful tool for exploratory analysis



- Unsupervised and exploratory
 - Interpret with caution
- MCFA remains in active development
 - Numerous additional applications!
- NHLBI TOPMed is among the most ambitious multi-omic data collection efforts
 - We look forward to future analyses with larger sample sizes

Thank you!

Development and analysis

- Collin Wang *
- Silva Kasela
- Francois Aguet
- Daniel Nachun
- Stephen Montgomery
- David Knowles *
- Tuuli Lappalainen *



MESA/TOPMED Collaborators

- Kent Taylor
- Russ Tracy
- Peter Durda
- Yongmei Liu
- W. Craig Johnson
- David Van Den Berg
- Namrata Gupta
- Stacy Gabriel
- Josh Smith
- Robert Gerzsten
- Clary Clish
- Quenna Wong
- George Papanicolaou
- Jerome I. Rotter
- Stephen S. Rich