

# **<sup>1</sup>H-NMR metabolomics-based surrogates to impute common clinical risk factors and endpoints**

D. Bizzarri<sup>1,2</sup>, M.J.T. Reinders<sup>2,3</sup>, M. Beekman<sup>1</sup>, P.E. Slagboom<sup>1,4</sup>, BBMRI-NL<sup>5</sup> and E.B. van den Akker<sup>1,2,3, #</sup>

<sup>1</sup>Molecular Epidemiology, LUMC, Leiden, The Netherlands

<sup>2</sup>Leiden Computational Biology Center, LUMC, Leiden, The Netherlands

<sup>3</sup>Delft Bioinformatics Lab, TU Delft, Delft, The Netherlands

<sup>4</sup>Max Planck Institute for the Biology of Ageing, Cologne, Germany

<sup>5</sup>BBMRI-NL: <https://www.bbMRI.nl>; see Consortium Banner Supplement S1

## **Material & correspondence (#)**

Erik B. van den Akker, PhD; Leiden Computational Biology Center, Leiden University Medical Center, Leiden, The Netherlands; Einthovenweg 20, 2333 ZC, Leiden, The Netherlands; Tel: +31 (0)71 526 85 57; Fax: +31 (0)71 526 82 80; E-mail: [e.b.van\\_den\\_akker@lumc.nl](mailto:e.b.van_den_akker@lumc.nl); website: <http://www.lcbc.nl>

## **ABSTRACT**

Missing or incomplete phenotypic information can severely deteriorate the statistical power in epidemiological studies. High-throughput quantification of small-molecules in bio-samples, i.e. ‘metabolomics’, is steadily gaining popularity, as it is highly informative for various phenotypical characteristics. Here we aim to leverage metabolomics to impute missing data in clinical variables routinely assessed in large epidemiological and clinical studies. To this end, we have employed ~25,000 <sup>1</sup>H-NMR metabolomics samples from 28 Dutch cohorts collected within the BBMRI-NL consortium, to create 19 metabolomics-based predictors for clinical variables, including diabetes status ( $AUC_{5-Fold CV} = 0.94$ ) and lipid medication usage ( $AUC_{5-Fold CV} = 0.90$ ). Subsequent application in independent cohorts confirmed that our metabolomics-based predictors can indeed be used to impute a wide array of missing clinical variables from a single metabolomics data resource. In addition, application highlighted the potential use of our predictors to explore the effects of totally unobserved confounders in omics association studies. Finally, we show that our predictors can be used to explore risk factor profiles contributing to mortality in older participants. To conclude, we provide <sup>1</sup>H-NMR metabolomics-based models to impute clinical variables routinely assessed in epidemiological studies and illustrate their merit in scenarios when phenotypic variables are partially

**NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

## 37 INTRODUCTION

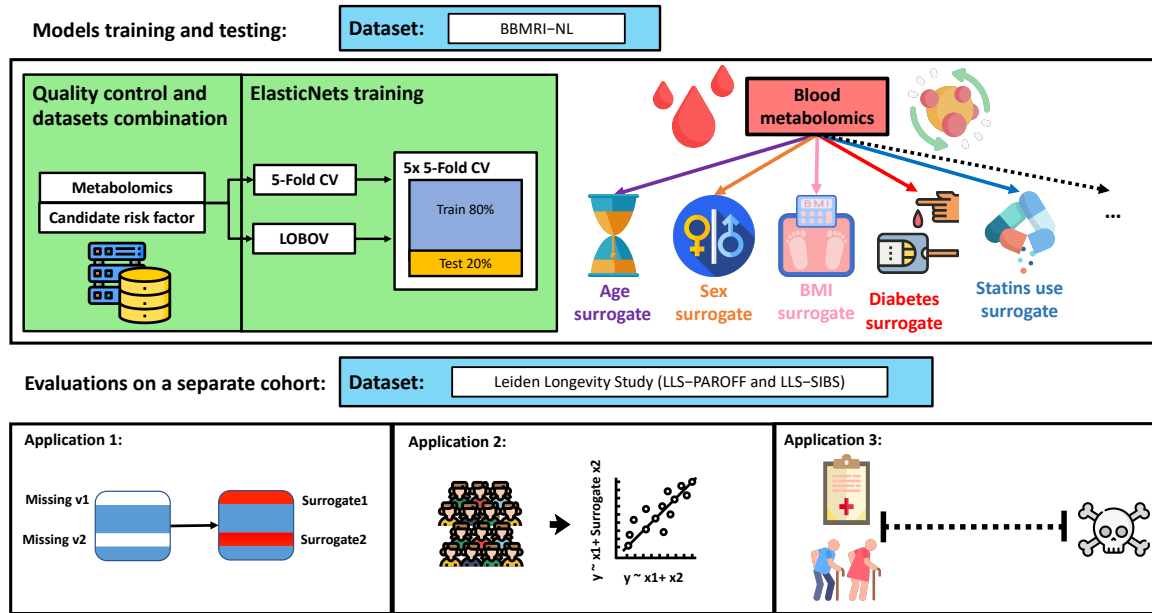
38 A major goal in biomedical research is to find faithful biomarkers of health, defined as  
39 accurate and reproducible assays that provide objective indications on the health of an  
40 individual and his/her risk of developing a disease over predefined time trajectory [1]. Over  
41 the years, many types of putative biomarkers have been proposed, ranging from environmental  
42 factors to biochemical assays, that may aid the diagnosis and prognostication of disease,  
43 including cardiovascular disease, cancer and immunological disorders. Many of these clinical  
44 variables, however, are costly or cumbersome to obtain, especially for more critical and frail  
45 participants, such as older individuals [2], [3]. Consequently, missing data frequently occurs in  
46 large epidemiological or clinical studies, potentially leading to a significant loss of statistical  
47 power, thus impeding biomarker research in studies of older individuals [4].

48 Missing phenotypic data can be handled in various ways. Often, analyses are either  
49 restricted to individuals or variables with complete data, which both may introduce potential  
50 biases [5]. Alternatively, missing data can be imputed using complete phenotypic variables [4],  
51 [6]–[8], yet these approaches work only satisfactory if the complete phenotypic variables are  
52 informative for the ones with missing observations. A third solution basically extends the  
53 second approach by leveraging informative omics data to impute missing phenotypic data.  
54 Particularly useful in this context are metabolite quantifications in minimally invasive  
55 biomaterials, such as urine, saliva or blood plasma, obtained with proton Nuclear Magnetic  
56 Resonance ( $^1\text{H-NMR}$ ) assays [9]. Although this technique only captures a modest number of  
57 analytes,  $^1\text{H-NMR}$  metabolomics data is frequently acquired in large-scale epidemiological  
58 studies, as it is a cost-efficient and reproducible data resource. The underlying motivation is  
59 that metabolite concentrations in blood seem to be direct readouts of various biological  
60 processes, incorporating cues of the environment as well as the host's genetic background, and  
61 hence may be regarded as intermediate phenotypes. Indeed, metabolomics has been shown to

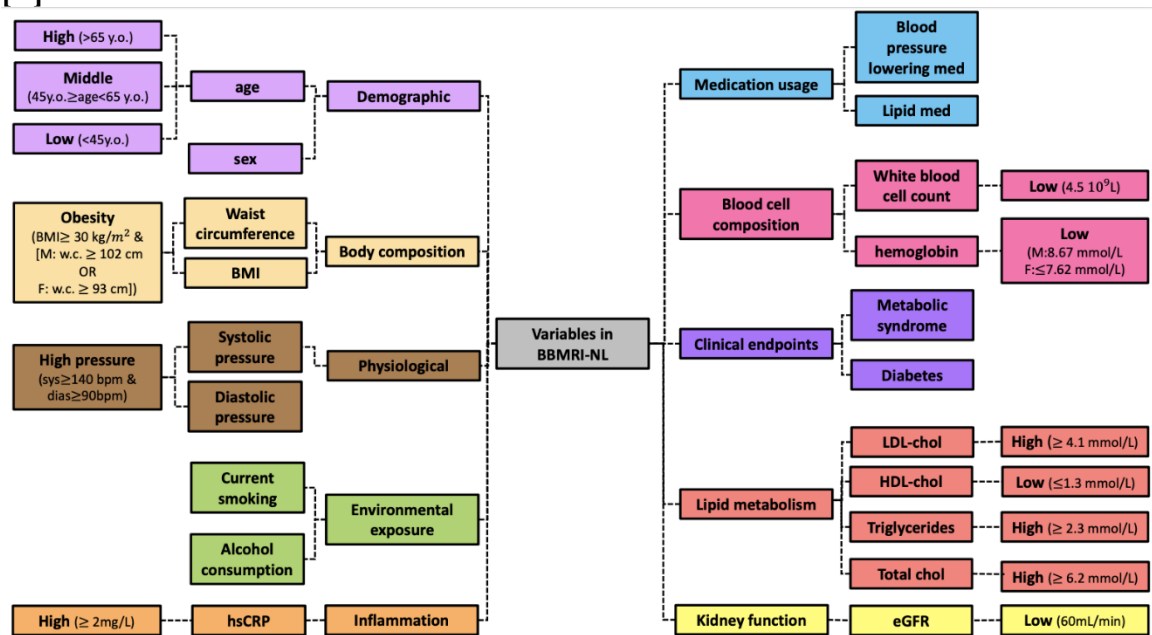
62 capture information on the effect of drug treatments [10], disease status [11]–[14], functional  
63 and cognitive decline [15], and aging [16], [17]. In addition, several studies used the blood  
64 metabolome to predict single anthropometric measures, i.e. BMI [18], or other physiological  
65 characteristics, i.e. sex [19] or age [16]. However, it remains unclear whether the blood  
66 metabolome captured by <sup>1</sup>H-NMR could represent phenotypic information over a wider set of  
67 conventional clinical variables.

68 We hypothesize that a single set of blood metabolic markers combined in multiple  
69 algorithms may represent a range of conventional clinical variables. As a proof of concept, we  
70 generated metabolic surrogates for 20 variables of general clinical and epidemiological interest  
71 available in at least 6 of the cohorts collaborating in BBMRI-NL. Here we will designate these  
72 as conventional clinical variables and they comprehend physiological measures (sex, age,  
73 blood pressure, etc.), environmental exposures (current smoking, etc.), body composition  
74 measures (BMI, etc.), inflammatory factors (hsCRP), medication usage (lipids medication,  
75 etc.), blood composition (white cell counts, etc.) lipids metabolism (LDL-cholesterol, etc.) and  
76 cardiometabolic clinical endpoints (diabetes and metabolic syndrome). Acquiring data for all  
77 these variables is costly and requires sufficient biomaterial, meaning that not every study has  
78 collected the same set of data. We further explored these methods to establish metabolic  
79 surrogate values in the Leiden Longevity Study, which we used to showcase possible  
80 applications in epidemiological research. We showed the validity of the surrogates in an  
81 external cohort comparing them to the original values, we examined their association to further  
82 clinically valuable cardiometabolic health markers, and explored whether the metabolic  
83 surrogates associate, separately or combined, to all-cause mortality.

[A]



[B]



**Figure 1: Study design.** [A] Upper panel: Training of  $^1\text{H-NMR}$  metabolomics-based predictors for routinely assessed phenotypic variables available in BBMRI.nl. This data set was created as a collaboration of 28 community and hospital-based cohorts that collected nuclear magnetic resonance ( $^1\text{H-NMR}$ ) metabolomics data (Nightingale) for  $\sim 31,000$  individuals. Upper panel left: Metabolomics-based predictors were trained using an *inner* loop of 5-fold Cross Validation (CV) (with 5 repetitions) for hyperparameter optimization and were evaluated in unseen data employing an *outer* loop of 5-fold CV or Leave-One-Biobank-Out-Validation (LOBOV). Upper panel right: using our models 19 different surrogate values can be derived from a single metabolomics data measurement to impute or complement a broad set of conventional clinical variables routinely assessed in epidemiological and clinical studies. Lower panel: Trained metabolomics-predictors were evaluated in two application scenarios using a held-out study, the Leiden Longevity Study [20]. This study is a two-generation family-based cohort consisting of highly aged parents (LLS-SIBS,  $N = 817$ , median age = 92 years) and their middle-aged offspring and the partners thereof (LLS-PAROFF,  $N = 2,280$ , median age = 59 years), for which we had access to additional detailed phenotypic information. Trained predictors were evaluated for their ability to reconstruct missing datapoints in an independent dataset (Application 1, lower left), to be used as confounder in Metabolome Wide Association Studies (Application 2, lower central), and to investigate and to explore determinants of health in older individuals (Application 3, lower right). [B] Groupings of phenotypic variables routinely assessed in epidemiological and clinical studies for which data was available in BBMRI-NL. Continuous variables are dichotomized at levels generally accepted to confer an increased risk for cardio-metabolic endpoints. As various cutoffs on chronological age are in use, in part reflecting the highly non-linear relation between chronological age and disease risk, we choose to split chronological age in three categories (I ‘young’:  $< 45$  years [TRUE/FALSE]; II ‘middle-aged’:  $\geq 45$  years [TRUE/FALSE] and III ‘old’:  $< 65$  years [TRUE/FALSE];  $\geq 65$  years). We integrated Body Mass Index, waist circumference and sex into one sex-specific measure of ‘obesity’. Similarly, we integrated diastolic blood pressure (DBP) and systolic blood pressure to arrive at one variable ‘high pressure’. Overall, we obtain data for 20 dichotomous phenotypic variables. Colors indicate groupings.

84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101  
102  
103

104 Through the paper, we indicated some of the clinical variables with the following abbreviation: BMI=Body Mass Index, med=medication,  
105 e.g.: lipid or blood pressure lowering medication, hsCRP=high-sensitivity C-Reactive Protein, eGFR=estimated Glomerular Filtration Rate,  
106 chol=cholesterol, hgb=haemoglobin, wbc=white blood cells.

## 107 **RESULTS**

### 108 **<sup>1</sup>H-NMR metabolomics can be used to successfully predict 19 out of 20 clinical variables** 109 **routinely measured in epidemiological and clinical studies**

110 Missing or incomplete phenotypic information can severely deteriorate the statistical power  
111 in epidemiological studies. Here we evaluate the ability of Nuclear Magnetic Resonance (<sup>1</sup>H-  
112 NMR) metabolomics (Nightingale Health<sup>®</sup>, Helsinki, Finland) to reconstruct conventional  
113 clinical variables. For this purpose, we trained and evaluated prediction models (**Figure 1A**)  
114 for 20 conventional clinical variables (**Figure 1B**) using data of ~31,000 individuals collected  
115 within the Dutch Biobanking and BioMolecular resources and Research Infrastructure  
116 (BBMRI-NL: <https://www.bbmri.nl/>). Out of 220 metabolomic variables measured on the  
117 platform, we employed 56 metabolic markers, selected to be the most uncorrelated [21], [22]  
118 and most successfully measured in the BBMRI studies (**Methods and Supplementary**  
119 **Materials**). Conventional clinical variables were transformed or constructed with the emphasis  
120 to be able to capture clinically relevant aspects of disease risk. For instance, we dichotomized  
121 continuous variables according to generally accepted clinical thresholds, thus obtaining for  
122 each of these clinical variables an ‘at risk’ [TRUE/FALSE] variable. For the same purpose,  
123 some variables were either merged or split. For instance, a sex-specific ‘*obesity*’  
124 [TRUE/FALSE] variable was defined using body mass index, waist circumference and sex,  
125 whereas chronological age was split into three categories (**Figure 1B**). Overall, we were able  
126 to construct and evaluate 20 variables mainly representing risk factors of cardio-metabolic  
127 health that are routinely assessed in epidemiological and clinical studies.

128 Logistic Elastic-NET regression models were trained for each of the 20 dichotomous  
129 variables, measured in at least 6 of the BBMRI studies, in both healthy and diseased  
130 individuals. Model development was performed in two loops to prevent overtraining. An *inner*  
131 loop of 5-Fold Cross Validation with 5 repetitions was used to tune the hyperparameters of the

132 model. Model performances were then evaluated in an *outer* loop of held out data, using again  
133 a 5-Fold CV or a Leave-One-Biobank-Out-Validation (LOBOV) (**Figure 1A, Methods**). We  
134 assessed model performances using the mean Area Under the Curve (AUC) of the receiver-  
135 operator curve obtained in the *outer* 5-Fold CV (**Table 1**) and considered a model's  
136 performance to be sufficiently accurate at  $AUC > 0.7$ . Overall, 19 out of 20 models passed this  
137 criterium, with only a single phenotypic variable, '*high-pressure*', that could not be accurately  
138 captured by <sup>1</sup>H-NMR ( $AUC_{5\text{-Fold CV}} = 0.68$ ). Strikingly, 9 out of 20 models achieved an  $AUC_{5\text{-Fold CV}} > 0.9$ . While some of these high performances are expected as they directly relate to  
139 metabolic markers assessed on the platform ('*Low eGFR*', '*high triglycerides*', '*high LDL*  
140 '*cholesterol*', '*high total cholesterol*', and '*low LDL cholesterol*'), this is not the case for four  
141 other high performing models: '*diabetes*' ( $AUC_{5\text{-Fold CV}} = 0.94$ ), '*metabolic syndrome*' ( $AUC_{5\text{-Fold CV}} = 0.93$ ), '*sex*' ( $AUC_{5\text{-Fold CV}} = 0.92$ ), '*lipid medication*' ( $AUC_{5\text{-Fold CV}} = 0.90$ ). Also, other  
142 important cardio-metabolic health statuses, including '*obesity*', '*high CRP*' and '*blood*  
143 '*pressure lowering medication*' were predicted at a more than satisfactory accuracy ( $AUC_{5\text{-Fold CV}} > 0.8$ ), indicating that overall, the <sup>1</sup>H-NMR metabolome can be used to impute a broad  
144 spectrum of common clinical variables.

148 As the performances of our models may vary per biobank due to study-specific  
149 characteristics, e.g. varying study inclusion criteria or protocols for sample storage, we also  
150 evaluated the variation of our model performances across biobanks. First, using a Leave-One-  
151 Biobank-Out-Validation (LOBOV), we evaluate how our models would perform when applied  
152 to data of a new unseen biobank. As expected, mean model accuracies of the LOBOV,  
153 weighted based on the size of the testing biobank, show more variation across folds (**Figure**  
154 **S2A-B**) and are generally slightly lower than the overall results of the 5-Fold CV (**Table 1**). In  
155 particular, some of the smaller studies containing diseased patients showed relatively poor  
156 accuracies (**Figure S2B**). Indeed, surrogate values do show cohort specific effects, but

157 interestingly, this does not necessarily affect its predictive performance within cohorts (**Figure**  
158 **S2C**). Overall, 14 out of the 20 models performed on average satisfactorily ( $AUC_{LOBOV} > 0.7$ )  
159 across all studies in the LOBOV setting.

Table 1   Performances of the 20 metabolic predictors					
Binary Outcomes (threshold)	# samples	# cohorts	# True positives	Results (AUC)	
				5-Fold CV (mean)	LOBOV (weighted mean)
<b>Low eGFR</b> (eGFR $\leq$ 60 ml/min [23])	21,439	23	1,196 (5.6%)	0.99 [0.98-0.99]	0.97 [0.91-0.99]
<b>High triglycerides</b> (trig $\geq$ 2.3 mmol/L [24])	13,401	11	1,645 (12.3%)	0.97 [0.97-0.98]	0.95 [0.84-0.99]
<b>High LDL cholesterol</b> (LDL $\geq$ 4.1 mmol/L [24])	13,261	11	2,051 (15.5%)	0.96 [0.96-0.97]	0.97 [0.86-0.98]
<b>High total cholesterol</b> (totchol $\geq$ 6.2 mmol/L [24])	16,586	11	3,206 (19.3%)	0.96 [0.96-0.96]	0.96 [0.83-0.99]
<b>Low HDL cholesterol</b> (HDL $\leq$ 1.3 mmol/L [24])	16,506	11	7,414 (44.9%)	0.95 [0.95-0.96]	0.95 [0.85-0.96]
<b>Diabetes</b> (TRUE/FALSE)	18,841	16	4,034 (21.4%)	0.94 [0.93-0.9]	0.86 [0.72-0.98]
<b>Metabolic syndrome</b> (TRUE/FALSE)	7,811	6	3,452 (44.2%)	0.93 [0.92-0.94]	0.86 [0.71-0.93]
<b>Sex (male)</b> (TRUE/FALSE)	21,610	23	10,281 (47.6%)	0.92 [0.92-0.93]	0.91 [0.73-0.99]
<b>Lipid medication</b> (TRUE/FALSE)	17,707	14	5,783 (32.7%)	0.91 [0.90-0.91]	0.85 [0.77-0.94]
<b>Low age</b> (age $<$ 45 y.o.)	21,519	23	3,353 (15.6%)	0.89 [0.88-0.90]	0.80 [0.55-0.85]
<b>High hsCRP</b> (hsCRP $>$ 3mg/L [25])	5,180	8	1,548 (29.9%)	0.86 [0.84-0.86]	0.81 [0.7-0.86]
<b>Blood pressure lowering medication</b> (TRUE/FALSE)	15,832	13	7,234 (45.7%)	0.82 [0.81-0.83]	0.71 [0.51-0.84]
<b>High age</b> (age $\geq$ 65 y.o.)	21,519	23	8,273 (38.4%)	0.82 [0.80-0.83]	0.73 [0.64-0.86]
<b>Obesity status</b> (BMI $\geq$ 30 kg/m <sup>2</sup> and w.c. $\geq$ 102 cm [M] BMI $\geq$ 30 kg/m <sup>2</sup> and w.c. $\geq$ 93 cm [F] [26])	19,322	18	3,135 (16.2%)	0.78 [0.75-0.80]	0.76 [0.69-0.81]
<b>Low hemoglobin</b> (hgb $\leq$ 6.67 mmol/L [M]; hgb $\leq$ 7.62 mmol/L [F] [27])	10,508	6	1,299 (12.4%)	0.76 [0.73-0.78]	0.72 [0.63-0.75]
<b>Low white blood cells</b> (wbc $\leq$ 4.5x10 <sup>9</sup> L [27])	9,496	6	818 (8.6%)	0.73 [0.69-0.76]	0.61 [0.5-0.71]
<b>Current smoking</b> (TRUE/FALSE)	21,662	23	8,276 (38.2%)	0.71 [0.70-0.72]	0.63 [0.48-0.78]
<b>Alcohol consumption</b> (TRUE/FALSE)	16,430	13	11,763 (71.6%)	0.71 [0.70-0.73]	0.60 [0.48-0.70]
<b>Middle age</b> (45 y.o. $\geq$ Age $<$ 65 y.o.)	21,519	23	9,893 (46.0%)	0.71 [0.70-0.72]	0.58 [0.50-0.69]
<b>High pressure</b> (systolic $\geq$ 140 mmHg and diastolic $\geq$ 90 mmHg [24])	17,509	12	7,765 (44.3%)	0.68 [0.66-0.69]	0.60 [0.52-0.76]



# samples: number of participants; # cohorts: number of cohorts that we could use for training the models and for the evaluation using 5-Fold Cross Validation (5-Fold CV); # TRUE POSITIVES: the number of samples with original variable equal to TRUE; Results (AUC): 5-Fold CV= the mean AUCs of the 5-Fold CV) and LOBOV= the mean AUCs of the Leave One Biobank Out Validation weighted based on the size of the testing biobank. [M]: male, or [F]: female specific criteria, eGFR=estimated glomerular filtration rate, w.c=waist circumference, hgb=haemoglobin, wbc=white blood cells.

160

## 161 **Metabolic surrogates show dependencies mimicking the conventional clinical variables**

162 Given that all models are trained on a relatively limited set of metabolic markers, we

163 investigated to what extent the produced models and predictions show mutual dependencies.

164 For this purpose, we first visualized the coefficients (betas) of the logistic Elastic-NETs

165 (**Figure 2**) to show the relative importance of the metabolites within each of the prediction

166 models. While the selection of variables for the models shows a distinct pattern, we also note

167 some similarities, as quantified by the correlations between the model coefficients (**Figure S3**).

168 Overall, we note a clear preference for the models to include metabolites of the classes

169 “Lipoproteins” and “Lipids and related measures” over “Amino Acids”. In addition, we note

170 that the models of related phenotypes also display some resemblances in the employed features,

171 for instance ‘*lipid medication*’ and ‘*blood\_pressure\_lowering\_medication*’ share some model

172 characteristics.

173 We next evaluated correlations between the outputs of our models, from here on referred

174 to as the ‘metabolic surrogates’ (**Figure 3**) and compared these to correlations between the

175 original clinical variables (**Figure S1B**) in the BBMRI.nl data set. Overall, we observe that the

176 model outputs show correlation patterns and groupings that largely mimic that of the original

177 variables. For instance, model outputs trained on variables related to weight problems, i.e.

178 ‘*obesity*’, ‘*diabetes*’, ‘*metabolic syndrome*’, show high mutual correlations, and moreover are

179 grouped with model outputs trained on medication usage, i.e. ‘*lipid medication*’ and ‘*blood*

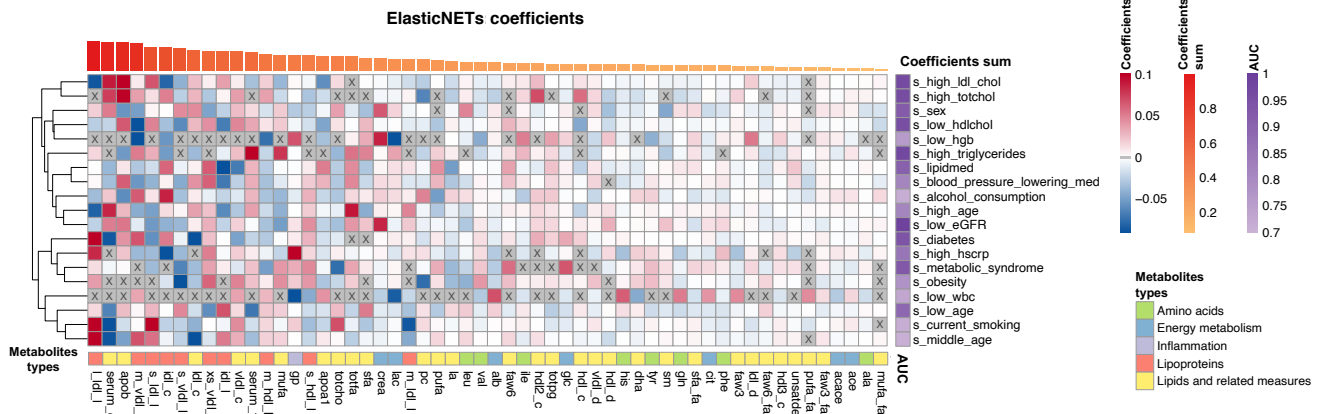
180 *pressure lowering medication*’. Although we observe some correlations between the outputs

181 of our different age predictors i.e. ‘*low age*’, ‘*middle age*’, ‘*high age*’, we observe that ‘*high*

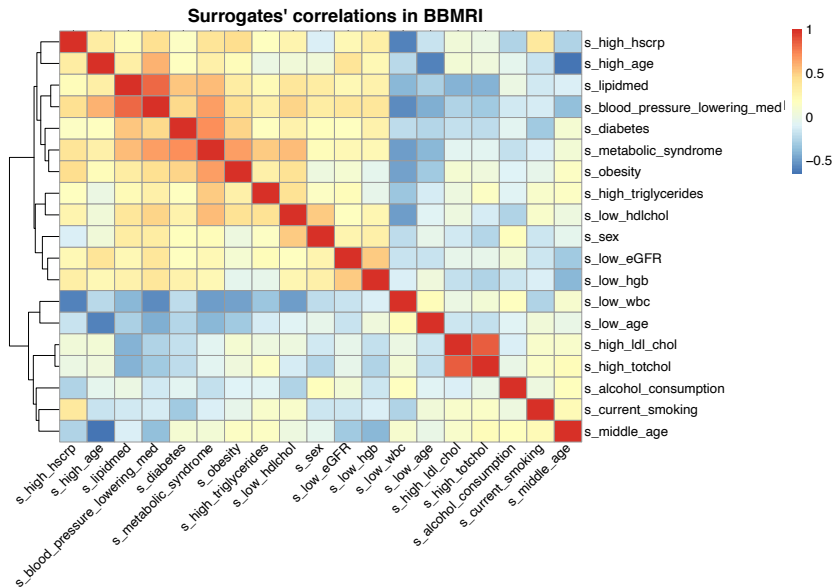
182 *age*’ is grouped with the models for ‘*high hscrp*’, ‘*lipid medication*’ and ‘*blood pressure*

183 *lowering medication*’, while ‘*middle age*’ is grouped with ‘*current smoking*’ and ‘*alcohol use*’

184 and ‘low age’ with ‘low white blood cell count’. This suggests that at different ages, different  
 185 conventional clinical variables play a role in physiology; an aspect well-known from literature  
 186 [28]–[30]. Overall, this indicates that our models show mutual dependencies similar as we  
 187 observe for the original clinical variables.



188  
 189 **Figure 2: ElasticNETs metabolites relative importance.** The heatmap reports the relative importance of the metabolites  
 190 (columns) in each of the trained models (rows). Prior to visualization, metabolite coefficients were scaled per model by  
 191 dividing them by the coefficients’ sum in each model to create the relative importance per model. Top: Metabolites were  
 192 then ordered based on the sum of their importance across all models. In addition, the models are clustered on the similarity  
 193 between relative importance. Bottom: Categorized metabolic measures: “Amino acids”, “energy metabolism”,  
 194 “inflammation”, “lipoproteins” and “lipids and related measures”. Right: Mean AUCs of the 5-FoldCV in a scale of  
 195 purple.



196  
 197 **Figure 3: Heatmap of pairwise correlations of the metabolic based surrogate markers calculated in BBMRI.**  
 198 The heatmap of correlations of the metabolic surrogate values of the 19 successful models, clustered based on the  
 199 correlation levels, between the imputed metabolic surrogate levels within BBMRI-NL.

201 **Projection in an independent study demonstrates model accuracy**

202 We performed a more extensive evaluation of the surrogate values by employing  
 203 Nightingale <sup>1</sup>H-NMR metabolomics and phenotypic data of the Leiden Longevity Study, a

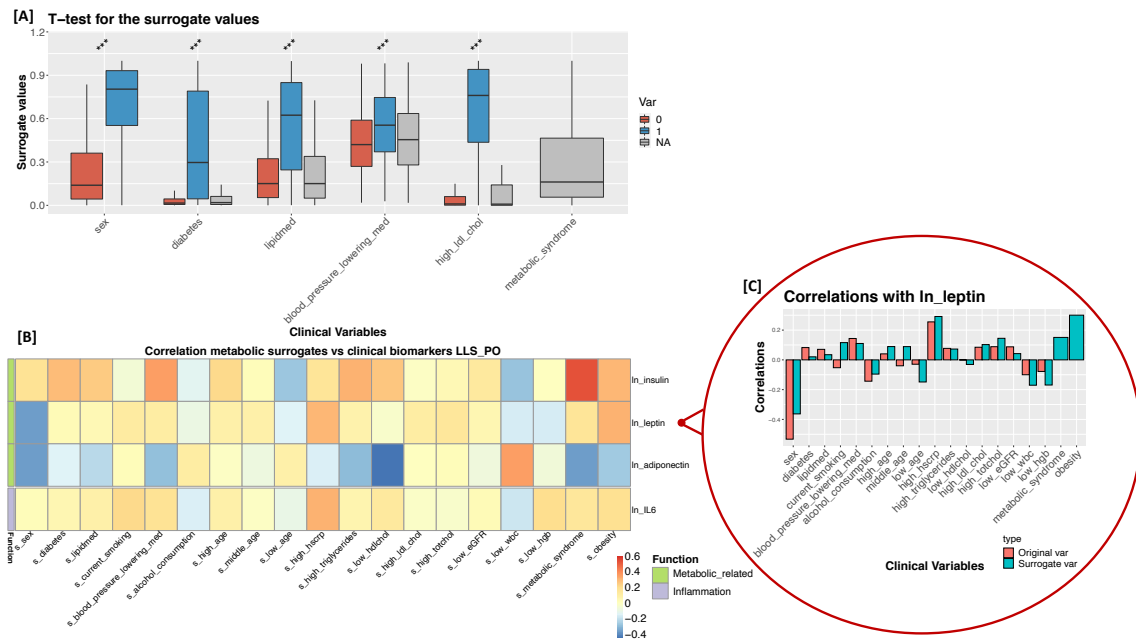
204 cohort excluded from the training and testing sets [20]. The Leiden Longevity Study is a two-  
205 generation family-based cohort consisting of highly aged parents (LLS-SIBS, 851 individuals,  
206 age median = 92 years old) their middle-aged offspring and their partners (LLS-PAROFF,  
207 2,307 individuals, age median = 59 years old). Using our models to project metabolic  
208 surrogates in the LLS-PAROFF gave an independent confirmation that conventional clinical  
209 variables can be readily captured by <sup>1</sup>H-NMR metabolomics. Splitting the surrogate values by  
210 the actual labels of the corresponding binary phenotypes generally showed a good separation  
211 for important cardio-metabolic variables like ‘sex’, ‘diabetes status’, ‘lipid medication’,  
212 ‘blood\_pressure\_lowering\_medication’ and ‘high LDL cholesterol’ (**Figure 4A** and **S4**),  
213 emphasizing the suitability of our models for quality control purposes or to impute missing  
214 data. For instance, model results for ‘sex’ could be applied to verify absence of sample mix-  
215 ups ( $t.stat = 44.58, p = 1.4 \times 10^{-313}$ ). In addition, surrogate values seem informative on the  
216 nature of the missingness of phenotypic data. For instance, participants with a missing diabetes  
217 status typically had metabolic surrogate values similar to those of participants without diabetes  
218 (diabetes:  $\mu_F = 0.05, \mu_T = 0.41, \mu_{NA} = 0.08$ ), suggesting that a missing diabetes status  
219 generally implies ‘non-diabetics’ in this cohort. Similar observations were made for medication  
220 status (lipidmed:  $\mu_F = 0.22, \mu_T = 0.45, \mu_{NA} = 0.21$  and blood\_pressure\_lowering\_med:  
221  $\mu_F = 0.44, \mu_T = 0.55, \mu_{NA} = 0.46$ ): participants with missing statuses were more similar to  
222 non-medication users than medication users. In contrast, participants with missing values in  
223 LDL cholesterol had surrogate values indicating “at risk” levels of LDL cholesterol  
224 (high\_ldl\_chol:  $\mu_F = 0.07, \mu_T = 0.66, \mu_{NA} = 0.14$ ). Lastly, our surrogates also allow for  
225 explorative analyses of totally unrecorded variables. For instance, the ‘metabolic syndrome’  
226 surrogate indicates participants who are more likely to have metabolic syndrome, a status  
227 which was not assessed in the LLS-PAROFF cohort (**Figure 4A**).

228

229 **Projection in the Leiden Longevity Study shows associations with additional cardio-**  
230 **metabolic phenotypes**

231 Within the LLS-PAROFF we had access to several additional variables pertaining to one's  
232 cardio-metabolic risk namely hormone levels of insulin, leptin, and adiponectin, as well as the  
233 levels of the inflammatory marker interleukin 6 (*IL6*) (**Figure 4B**). As expected, insulin levels  
234 correlated positively with most surrogate cardio-metabolic risk factors and endpoints,  
235 including 'diabetes' [31] ( $r = 0.28, p = 7.6 \times 10^{-42}$ ) and even more so with 'metabolic  
236 syndrome' ( $r = 0.52, p = 1.4 \times 10^{-152}$ ) [32]. Conversely, both 'low wbc' and 'low age' were  
237 inversely correlated with insulin, both reflecting the associations with decreased insulin  
238 sensitivity in those with high white blood cell counts [33] or in old age [34]. A similar analysis  
239 for the satiety hormone leptin showed the strongest positive correlations with 'obesity' [35]  
240 ( $r = 0.3, p = 5.6 \times 10^{-49}$ ), but also with 'high hscrp' [36] ( $r = 0.29, p = 4.2 \times 10^{-46}$ ). A  
241 significant correlation with leptin was also found for 'metabolic syndrome' [37]  
242 ( $r = 0.15, p = 3.9 \times 10^{-13}$ ), yet not 'diabetes' [38] ( $r = 0.02, p = 0.33$ ). In line with  
243 previous studies, higher levels of the adiponectin hormone generally correlated with lower  
244 values of the surrogates, most prominently with 'low hdlchol' ( $r = -0.47, p = 8.6 \times 10^{-123}$ )  
245 and 'high triglycerides' ( $r = 0.3, p = 1.1 \times 10^{-48}$ ) [39]. Higher adiponectin levels were  
246 positively correlated with 'low wbc' ( $r = 0.35, p = 6.1 \times 10^{-68}$ ), reproducing the previously  
247 reported association by Matsubara *et al* [40]. Levels of the inflammatory marker IL6 were most  
248 prominently positively correlated with the surrogates 'high hscrp' [41] ( $r = 0.31, p =$   
249  $8.7 \times 10^{-51}$ ), 'current smoking' [42] ( $r = 0.2, p = 3.9 \times 10^{-13}$ ), and inversely correlated with  
250 'low wbc' [43] ( $r = -0.19, p = 3.9 \times 10^{-13}$ ). When comparing these correlation patterns  
251 obtained with the surrogate values (**Figure 4B**) with those you would get when using the  
252 original values (**Figure S7A**), we generally notice highly similar trends ( $r \sim 0.83$ , **Figure**  
253 **S7F**), as exemplified for leptin (**Figure 4C**). Overall, these findings indicate that our metabolic

254 surrogates faithfully reproduce the original clinical variables in their association with insulin,  
 255 leptin, adiponectin and IL6.



256 **Figure 4: Metabolic surrogates applied to LLS-PAROFF:** [A] Paired boxplots show surrogate values split between  
 257 the TRUE/FALSE (0/1) in the original values of the clinical variables (\*\*\*)  $p$  value  $\leq 0.001$ ). For metabolic syndrome  
 258 the original variables are entirely missing, so no  $p$ -value is reported. [B] Heatmap of the correlations between the metabolomic  
 259 surrogates (columns) and four additional cardio-metabolic biomarkers (rows) available in LLS-PAROFF. Values of insulin,  
 260 leptin, adiponectin and IL6 were transformed with a natural logarithm. [C] Paired bar plots comparing the correlations  
 261 computed between leptin and the metabolic surrogates (blue) and those computed between leptin and the values of the original  
 262 clinical variables (red) the surrogates are trained to predict.  
 263

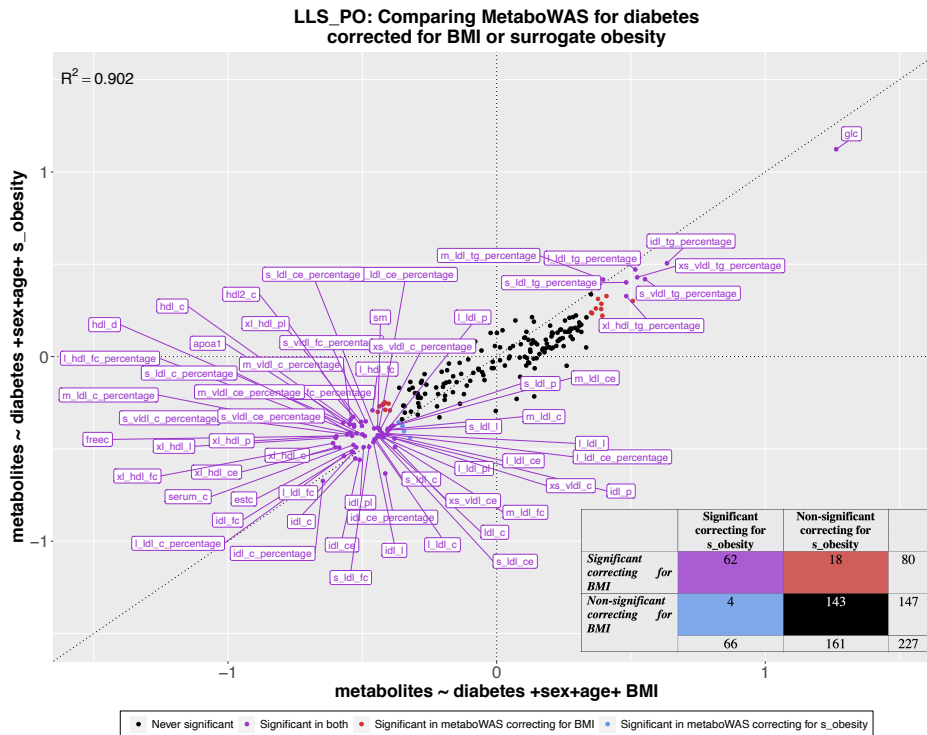
## 264 Metabolic surrogates to explore confounders in Metabolome Wide Association Studies

265 We next explored the use of  $^1\text{H-NMR}$  metabolic surrogates to complement missing  
 266 phenotypic data in metabolome-wide association studies (MetaboWAS). As an example, we  
 267 evaluated the association of metabolic markers with Type 2 Diabetes status (T2D) in absence  
 268 of information on a known potential confounder: BMI. We designed a controlled experiment  
 269 to evaluate to what extent surrogate ‘obesity’ can replace BMI, using data of 1,697 individuals  
 270 of LLS-PAROFF with complete metabolomic, BMI and diabetes status, of which 79 are  
 271 diagnosed with type 2 diabetes.

272 First, we ascertained that BMI was indeed a confounder, also within the LLS-PAROFF, by  
 273 showing that BMI associated with the outcome (Type 2 Diabetes status,  $t$ -test= -7.83,  $p$  =  
 274  $8.25 \times 10^{-15}$  **Figure S8A**), as well as many of the determinants (147 significant metabolites

275 after correction, see **Methods**) of the MetaboWAS. Concomitantly, further adjustment of the  
276 MetaboWAS on T2D for BMI drastically reduced the number of significant metabolites. To  
277 compare, when adjusting for age and sex we identified 136 metabolites significantly associated  
278 with diabetes status, whereas further adjustment for BMI identified 80 significant metabolites  
279 (**Figure S9B Comparison 1**). Next, we performed the same association analyses using the  
280 ‘*obesity*’ surrogate as confounder. Similar to BMI, also the ‘*obesity*’ surrogate is significantly  
281 higher in diabetics as compared to non-diabetics ( $t\text{-test} = -11.2, p = 2.48 \times 10^{-28}$  **Figure S8B**)  
282 and was associated with many of the metabolites (176 significant associations). Further  
283 adjusting the MetaboWAS on T2D for ‘*obesity*’ reduced the number of significant metabolites  
284 to 66 (**Figure S9B Comparison 2**).

285 We then investigated to what extent adjusting for BMI or adjusting for the ‘*obesity*’  
286 surrogate yields similar metabolite markers to T2D associations, by comparing the obtained  
287 estimates from both models (**Figure 5**). Overall, highly similar ( $r^2 = 0.902$ ) associations  
288 between metabolic markers and T2D are found for both models, with glucose being the most  
289 significantly associated marker in both ( $p = 9.43 \times 10^{-28}$  when correcting for BMI and  $p =$   
290  $1.6 \times 10^{-26}$  correcting for ‘*obesity*’). While most metabolites reported to be significantly  
291 associated with T2D overlap between the two models (62 out of 227; in purple), some  
292 discrepancies were observed, particularly at the significance threshold. When correcting for  
293 BMI, 18 significant metabolic markers were identified, that were not identified when correcting  
294 for ‘*obesity*’ (red dots, false negative rate  $\sim 0.11$ ). Conversely, 7 metabolites were deemed  
295 significantly associated with diabetes status when adjusting for ‘*obesity*’, but not when  
296 adjusting for BMI (blue dots, false positive rate  $\sim 0.027$ ). Nevertheless, overall, the differences  
297 in estimated effects remain small, indicating that metabolic surrogates may prove useful to  
298 account for missing data in epidemiological studies.



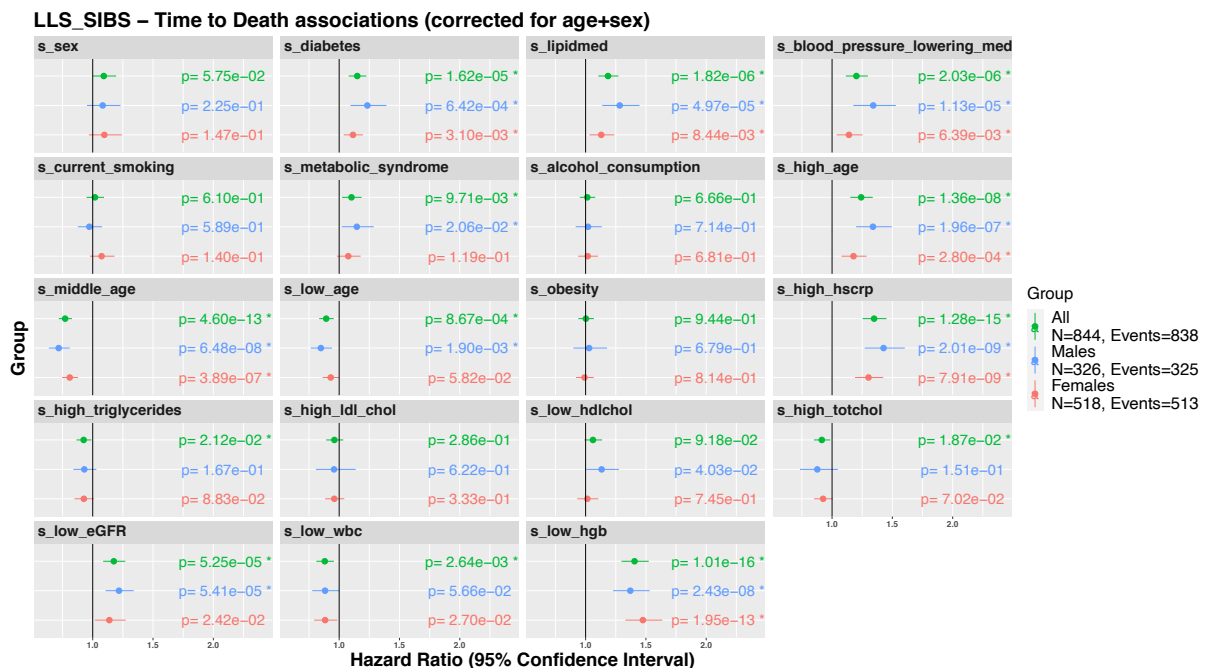
299  
300 **Figure 5: Comparison between the estimated coefficients of metaboWAS on T2D adjusted for BMI or for surrogate**  
301 **'obesity'.** On the x-axis the metaboWAS for diabetes adjusting for BMI and on the y-axis the metaboWAS for diabetes adjusted  
302 for surrogate obesity. The data set composed of 1,697 individuals, 79 of which are diabetics. Estimated coefficients for each  
303 metabolite (points) are colored based on their significance in the two models: purple: significant in both; red: significant when  
304 adjusted for BMI only; blue: significant when adjusted for surrogate obesity only; black never significant. Lower right corner:  
305 a contingency table with the number of significant and non-significant metabolites identified using the two models.

### 306 **Metabolic surrogates associate with incident all-cause mortality in older individuals**

307 Next, we evaluated whether metabolic surrogates are indicative of health at old age, by  
308 associating these with all-cause mortality in a nonagenarian subsample of the Leiden Longevity  
309 Study (LLS\_SIBS; 844 individuals, median age at baseline: 92 years old) (**Figure 6A**). Using  
310 a Cox proportional hazards model adjusted for sex and age at inclusion for each of the 19  
311 metabolic surrogates (**Materials and Methods**), we observed that 13 out of the 19 surrogates  
312 associated significantly with all-cause mortality (**Figure 6**, 'all'). In line with previous reports,  
313 we observed the largest effect sizes with the surrogate levels of 'high age', 'medications  
314 usage', 'diabetes status', 'high hscrp' and 'hemoglobin'. As previous studies have reported  
315 sex-specific associations for these clinical variables with all-cause mortality, we conducted a  
316 stratified analysis [44]–[50]. Although, the direction of association with all-cause mortality  
317 remains generally the same between men and women, the strengths and their significance are  
318 in some cases different. For instance, the surrogate 'diabetes' is associated with a higher risk

319 on mortality in men (HR = 1.23,  $p = 6.42 \times 10^{-4}$ ), than for women (HR = 1.11,  $p = 3.1 \times 10^{-3}$ ),  
 320 the same goes for ‘blood pressure lowering medication’ (men: HR = 1.3,  $p = 1.13 \times 10^{-5}$ ,  
 321 women: HR = 1.1,  $p = 6.39 \times 10^{-3}$ ). In contrast, ‘low hemoglobin’ is associated with a higher  
 322 risk in women (HR = 1.5,  $p = 1.95 \times 10^{-13}$ ), than men (HR = 1.37, FDR =  $2.42 \times 10^{-8}$ ).

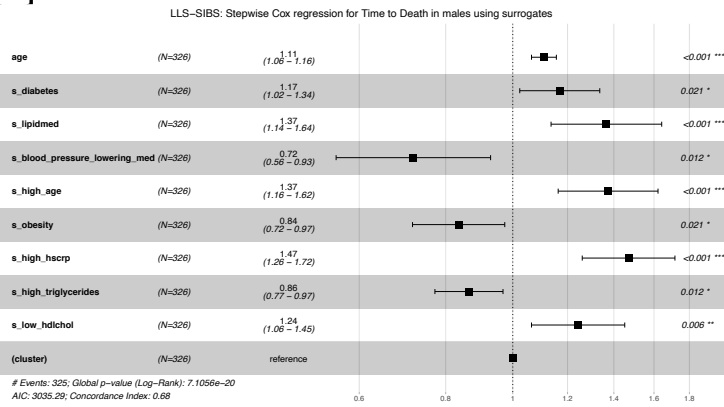
323 To identify the minimal set of metabolic surrogates independently associating with all-  
 324 cause mortality, we performed a stepwise (forward/backward) cox regression, adjusted for age  
 325 at sampling, in the LLS-SIBS dataset (Figure 7A-B), stratified for sex. The surrogates ‘high  
 326 hsCRP’ and ‘high triglycerides’ emerged as independent predictive features in both male and  
 327 female models, associated with an increased and decreased risk respectively. While eight  
 328 surrogates contributed to the male model, including ‘lipid medication’, ‘high age’, ‘high  
 329 hsCRP’ and ‘low hdlchol’, only three surrogates contributed to the mortality prediction in  
 330 females: ‘high hsCRP’, ‘high triglycerides’ and ‘low hemoglobin’. These findings are in line  
 331 with previous reports that different risk factors seem to predict survival up to the highest ages  
 332 for the different sexes [51]–[53].



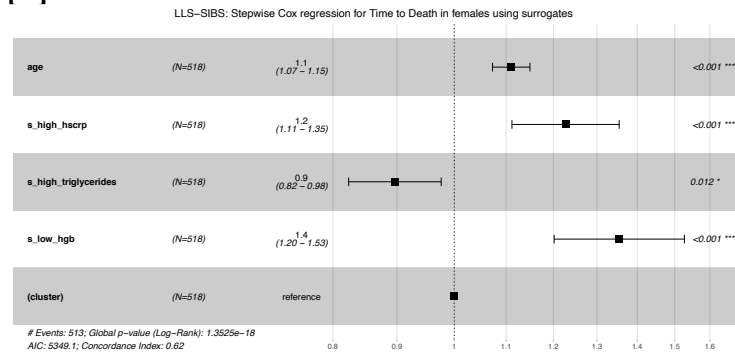
333  
 334 **Figure 6: Associations of the surrogate metabolic measures with incident all-cause mortality.** Associations of the  
 335 metabolic surrogates with time to all-cause mortality in LLS-SIBS, in groups comprising the entire set (“All”, N = 844 with  
 336 838 reported deaths), males (N = 326 with 325 reported deaths) or females (N = 518 with 513 reported deaths).  
 337



[A]



[B]



338 **Figure 7: Composite metabolomics predictors of incident all-cause mortality:** Predictors of time to death for males [A]  
 339 and females [B], in LLS-SIBS, composed using the surrogate metabolic measures, sex and age. “(cluster)” refers to the variable  
 340 controlling for family relationships (methods). Cox regression models were made using a step forward/backward selection.

## 341 **Discussion**

342 Missing phenotypic data is common in large epidemiological studies and in particular  
343 impedes biomarker research in older individuals. We employed <sup>1</sup>H-NMR metabolomics data  
344 as a single source of information to successfully impute 19 out of 20 conventional clinical  
345 variables that mainly relate to cardio-metabolic health. We highlighted the potential of our  
346 imputation models for conventional clinical variables with three application scenarios. First,  
347 we applied our models to an independent study, the Leiden Longevity Study, demonstrating  
348 that we can reconstruct conventional clinical variables at high accuracy. Secondly, we showed  
349 the value of metabolic surrogates in omics studies when data on potential confounders is  
350 missing. Finally, we exemplified how metabolic surrogates can be used to explore risk factors  
351 of health in older individuals by showing that multiple metabolic surrogates are independently  
352 predictive of all-cause mortality.

353 Using logistic ElasticNET regression models we were able to reconstruct a broad range of  
354 conventional clinical variables assessed in BBMRI-NL pertaining to physiological measures,  
355 body composition measures, environmental exposures, inflammatory factors, medication usage  
356 blood cell composition, lipids metabolism, and also clinical endpoints. For this purpose, we  
357 constructed composite variables that may better capture particular aspects of health, for  
358 instance, our '*obesity*' variable integrates body mass index, waist circumference and sex to  
359 create a sex-specific measure for overweight. In addition, we chose to construct binary  
360 representations of the continuous clinical variables for several reasons. First, we binarized  
361 continuous variables for a practical reason – to be able to judge all models on the same criteria.  
362 Secondly, predicting continuous variables using linear ElasticNET regression models  
363 emphasizes the prediction of the extremes of a phenotypic distribution, i.e. the model will fit  
364 the most atypical participants, whereas the current approach emphasizes to predict the  
365 commonly populated phenotypic range in which participants become at risk. Thirdly, our

366 models output a posterior probability that indicates the likelihood (a continuous score) of a  
367 sample belonging to one of two labels, e.g. *obese/non-obese*. In effect, these posteriors  
368 reconstitute part of the information lost when dichotomizing continuous variables, as  
369 exemplified by the observed correlation patterns between surrogates that mimic the correlation  
370 patterns between the original variables.

371 Our pre-trained models for conventional clinical variables allow for the imputation of  
372 missing datapoints in partially incomplete phenotypic variables, and moreover they offer the  
373 opportunity to explore associations with completely unobserved phenotypic variables. The  
374 latter is very much in line with the current use of PolyGenic Scores (PGSs) [54]–[56]. A PGS  
375 captures the genetic propensity of the realization of a particular polygenic phenotype. Nearly a  
376 thousand PGSs have been collected [57], which can be used to systematically explore  
377 correlations between a measured variable of interest and a wide array of phenotypes-by-proxy  
378 in genetic studies. We propose a similar use for metabolic surrogates in large metabolomics  
379 studies, yet with two noteworthy distinctions. Whereas PGSs can arguably be used to tease out  
380 causality in so-called Mendelian Randomization studies [58], metabolomic surrogates cannot.  
381 In contrast, while PGSs often only explain a very modest part of their respective phenotypes,  
382 metabolic surrogates explain a much larger part, thus enabling different types of applications.  
383 We illustrated this in our second application scenario where we showcased the use of surrogates  
384 to explore potential confounding by non-assessed phenotypic variables in omics studies. While  
385 use of actual phenotypic variables will always be preferred over metabolic surrogates, the  
386 availability of these metabolic surrogates can thus be used to direct replication efforts or to  
387 inform the design of new or follow-up studies.

388 Besides anthropometric measures and other physiological characteristics, the blood <sup>1</sup>H-  
389 NMR metabolome was previously also shown to capture aspects directly pertaining to health  
390 outcomes. In particular, we and others have previously reported <sup>1</sup>H-NMR metabolomics-based

391 risk estimators of cardiometabolic disease [59], [60], pneumonia and COVID infection [61],  
392 and all-cause mortality [17]. While this clearly illustrates the vast potential of the blood <sup>1</sup>H-  
393 NMR metabolome as a universal readout for health outcomes, it also raises the question what  
394 factors give rise to metabolomic profiles associated with adverse outcomes. Given that the find  
395 variation in the <sup>1</sup>H-NMR metabolome is the result of a complex interplay of both environmental  
396 and genetic factors, we evaluated whether our surrogates might give us a first indication. To  
397 do so, we tested which of our surrogates might be indicative of all-cause mortality in an elderly  
398 subset of the LLS-study. Intriguingly, by employing our pre-computed models as well as when  
399 we built multi-variate cox-regression models for time-to-death, we find metabolic surrogates  
400 that relate to conventional clinical risk factors known to associate with mortality risk at old  
401 age. Moreover, sex-stratified analyses recapitulate some of the known differences in mortality  
402 associations observed at old age, with for instance many more risk factors independently  
403 associated for mortality in males, as compared to females. These results illustrate that  
404 metabolic surrogates can aid in the interpretation of metabolomics-based risk estimators.

405 This study has several limitations. LOBOV analyses revealed that the trained surrogates  
406 may show study-specific effects that may relate to employed procedures of data collection or  
407 sample storage of the cohorts under investigation. While these artifacts may be addressed using  
408 batch-correction algorithms [62], or employing deep learning models for the prediction tasks,  
409 we note that differences between studies may also be due to valid biological reasons, such as  
410 differences in inclusion criteria. Secondly, the number of biomarkers captured by the targeted  
411 NMR platform is small compared to the whole human metabolome (over 19,000 according to  
412 the Human Metabolome Database [63]). Therefore, more elaborate, though typically more  
413 costly, high-throughput platforms might reach even higher accuracy levels. However,  
414 employing more biomarkers also has the danger of overfitting to the aforementioned study-  
415 related artifacts.

416 In conclusion, we have shown that the blood metabolome assayed by <sup>1</sup>H-NMR  
417 metabolomics can successfully capture a broad set of conventional clinical variables opening  
418 various possibilities to exploit surrogates of these clinical variables in in large epidemiological  
419 and clinical studies.

## 420 **MATERIALS and METHODS**

421

### 422 **1. Study populations**

423 The samples used for the current study are part of the BBMRI-NL Consortium (Dutch  
424 Biobanking and BioMolecular resources and Research Infrastructure, <https://www.bbmri.nl/>),  
425 which includes the following 28 Dutch biobanks: ALPHAOMEGA, BIOMARCS, CHARM,  
426 CHECK, CODAM, CSF, DMS, DZS\_WF, ERF, FUNCTGENOMICS, GARP, HELIUS,  
427 HOF, LIFELINES, LLS\_PARTOFFS, LLS\_SIBS, MRS, NESDA, PROSPER, RAAK, RS,  
428 STABILITEIT, STEMI\_GIPS-III, TACTICS, TOMAAT, UCORBIO, VUMC\_ADC,  
429 VUNTR. A description of the cohorts included is provided in the Supplementary Materials.  
430 Ethics committees approved the protocols for these studies in all the involved institutes, and  
431 all participants provided informed consent. The whole data set contains samples of ~31,000  
432 individuals.

433

### 434 **2. Metabolomic measurements**

435 The present study included metabolite concentrations measured in EDTA plasma  
436 samples using the high-throughput proton Nuclear Magnetic Resonance (<sup>1</sup>H-NMR)  
437 metabolomics (Brainshake Ltd./Nightingale Health<sup>®</sup>, Helsinki, Finland). This device provides  
438 the quantification of routine lipids, lipoprotein subclasses, fatty acid composition and various  
439 low-molecular weight metabolites including amino acids, ketone bodies and glycolysis-related  
440 metabolites in molar concentration units. Details about the methods and applications of the  
441 NMR platform have been provided previously [22], [64]. The total amount of metabolic  
442 variables reported is 226 for EDTA plasma samples, including the ratios and derived  
443 measurements, but only 63 of these were considered for the current study, to prevent overfitting  
444 [21], [59]. The list comprises the total lipid concentrations, fatty acids composition and low-  
445 molecular-weight metabolites including ketone bodies, glycolysis-related metabolites, amino-

446 acids and metabolites related to immunity and fluid balance (See the Supplementary Materials  
447 for a full list).

448

### 449 **3. Data Pre-processing**

#### 450 **a. Pre-processing of metabolomics data**

451 We included in our analyses all the cohorts reporting on all the 63 metabolic biomarkers,  
452 therefore we omitted CODAM (N = 254) and VUNTR (N = 3,896), which are missing  
453 acetoacetate and glutamine, respectively. We also decided to not consider the metabolites with  
454 low detection rates in more than one cohort (3-hydroxybutyrate) or which frequently failed to  
455 reach the minimum detection threshold (XL\_VLDL\_L, XXL\_VLDL\_L, L\_VLDL\_L,  
456 XL\_HDL\_L, L\_HDL\_L). We removed outlier samples with 1 or more missing metabolic  
457 measure (232 removed samples), 1 or more zeroes per sample (74 removed samples) and  
458 samples with any metabolite concentration level more than 5 standard deviations away from  
459 the overall mean per metabolomic variable (604 removed samples). The remaining 551 missing  
460 values in the dataset were imputed using the function `nipals` of the R package `pcaMethods`, and  
461 we z-scaled the metabolic measures across all samples to have comparable concentration levels  
462 between metabolites. The final data matrix comprised 26,107 samples across 56 metabolic  
463 variables. For more details, see Supplementary Materials. The number of samples used to train  
464 a predictor for a clinical variable depended on the number of samples missing this phenotypic  
465 information (Table 1). More information about the range of each phenotype within each  
466 biobank can be found in the Supplementary Materials.

467

#### 468 **b. Binarization of the clinical variables**

469 To emphasize the relevant clinical conditions, we used clinical thresholds to obtain  
470 dichotomous variables out of the set of the available continuous risk factors, separating

471 between “normal” and “at risk” levels for each risk factor (in *Table1* and in the Supplementary  
472 Document 3).

### 473 **c. Composed clinical variables**

474 We chose to include some composed clinical variables: 1) LDL cholesterol, which was  
475 calculated using the Friedewald equation [65] with the measured hdl cholesterol, triglycerides  
476 levels and total cholesterol; 2) eGFR (estimated Glomerular Filtration Rate), which is a  
477 measure for the kidney filtration rate of an individual, was calculated using the creatinine-based  
478 CKD-EPI equation [66]; 3) obesity, which is a binary variable describing if a person is  
479 clinically obese or not variable that uses BMI, waist circumference and sex based on the finding  
480 of Flint et al. [26]; 4) high pressure, a binary variables which defines high blood pressure by  
481 using systolic and diastolic blood pressure [24]; 5) low\_hgb (low hemoglobin), which is a  
482 binary variables describing ‘at risk’ levels of hemoglobin by using hemoglobin and sex [27].

## 483 **4. Estimation of the metabolic surrogates**

### 484 **a. Method selection**

485 The models considered for each Risk Factor are logistic regression models:

$$487 \hat{c}_i \sim \beta_0 + \sum_{j=1}^{56} \beta_j m_j \beta_1 + \varepsilon_i$$

489 in which  $c_i$  represent one of clinical variables of interest,  $m_j$  one of the 56 measured  
490 metabolites,  $\beta_j$  the regression coefficient, and  $\varepsilon_i$  the normal distributed reconstruction error.

491 The regression coefficients are found by minimizing the ordinary least squares error, with the  
492 addition of an elastic net regularization term for the coefficients:

$$493 \boldsymbol{\beta}(\alpha, \lambda) = \operatorname{argmin}_{\boldsymbol{\beta}} ((c_i - \hat{c}_i)^2 + \lambda[\alpha \|\boldsymbol{\beta}\|_1 + (1 - \alpha) \|\boldsymbol{\beta}\|_2])$$

494 in which  $\lambda$  ( $\in (0, \infty)$ ) represents the “shrinkage parameter” and  $\alpha$  ( $\in (0,1)$ ) is the mixing  
495 parameter balancing the L1 and L2 norm regularizations. We fixed the mixing parameter  $\alpha$  at  
496



499 0.5 for the predictive models, like previously done by other authors and optimized the  
500 shrinkage parameter using an inner-fold cross-validation scheme.

501

## 502 **b. Training and validation procedure**

503 We employed two training-evaluation procedures to get an unbiased estimate of the  
504 models' possible performances (Figure 1). As a first scenario, we used a Double 5-Fold-Cross-  
505 Validation (5-Fold CV) with 5 repetitions. This procedure consists of two loops of 5FCV, one  
506 internal and one external, in which we first split the dataset in testing (20%) and training (80%)  
507 sets and then on the latter set we have another 5FCV repeated for 5 different times, which is  
508 done for an unbiased tuning of the model (setting the correct  $\lambda$  parameter) that is finally trained  
509 on the complete training dataset and tested on the left-out test data. Both 5-FoldCVs were done  
510 such that the original distribution of each clinical variable is maintained as much as possible  
511 (using the function *createFolds* from the R package *caret*). In the second training-testing  
512 procedure, we applied a Leave-One-Biobank-Out-Validation (LOBOV), which consists of  
513 holding out one of the biobanks with the considered variable available, which is then used as a  
514 test set, while training on the remaining biobanks [16]. Also, in this setting, we applied a 5FCV  
515 with 5 repetitions to tune the best model for each training set.

516

## 517 **5. Metabolome wide association studies**

518 We conducted Metabolome Wide Association Studies (MetaboWAS) using the middle-  
519 aged cohort of the Leiden Longevity Study (LLS-PARTOFFs, 2,307 individuals, median age  
520 at baseline = 59 years old). As metabolites distributions are often skewed, we first transformed  
521 all metabolite measurements using a rank inverse normalization (RIN). Applying a PCA on the  
522 LLS-PARTOFFs dataset revealed that the first 40 principal components explain 99% of the  
523 variance in the metabolites (**Figure S9A**). Hence, the *p-value* of the MetaboWASes were

524 Bonferroni corrected using 40 tests, i.e. a *p-value* designated significant when smaller than  
525 0.00125 (0.05/40) [60]. We performed 5 different MetaboWASs.

526

## 527 **6. Associations of the metabolic surrogates to all-cause mortality**

528 We used Cox proportional hazards models with follow-up time as the time scale, to test for  
529 associations between the metabolic surrogate measures and incident endpoints, i.e. all-cause-  
530 mortality in LLS-SIBS. We checked for associations adjusting for age and sex. To avoid bias  
531 due to familial correlations from pedigrees, we used robust standard errors (calculated with the  
532 Huber sandwich estimator) implemented in R *coxph* function. Considering that the population  
533 in LLS-SIBS has a different inclusion criterium for men (age > 89 years old) and women (age  
534 > 91 years old), we also evaluated associations separately in men and women. *P-values* were  
535 corrected using Benjamini Hochberg separately for each selection (all individuals, men and  
536 women) and considered significant the  $FDR < 0.05$ . To select potentially interesting metabolic  
537 surrogate, we used a stepwise procedure for the Cox regression models, corrected for sex and  
538 age. Starting from a model containing the full set of available variables, we removed or added  
539 an unselected metabolic surrogate at each round based on the improvement on the model  
540 calculated from the Akaike Information Criterion and considering the *p-value* of each variable  
541 included in the model.

542

## 543 **Acknowledgements**

544 This work was performed within the framework of the BBMRI Metabolomics Consortium  
545 funded by BBMRI-NL (a research infrastructure financed by the Dutch government, NWO  
546 184.021.007 and 184.033.111), by X-omics (NWO 184.034.019), VOILA (ZonMW  
547 457001001) and Medical Delta (scientific program METABODELTA: Metabolomics for  
548 clinical advances in the Medical Delta). EvdA is funded by a personal grant of the Dutch

549 Research Council (NWO; VENI: 09150161810095). A full list of acknowledgements for all  
550 the contributing studies can be found in the Supplementary Material table S1.

551

## 552 **Authors Contribution**

553 EbvDA, DB, MJTR and PES conceived and wrote the manuscript. DB performed the analyses.

554 EBvDA and MJTR verified and supervised the analyses. All authors discussed the results and

555 contributed to the final manuscript.

556

## 557 **Competing interests**

558 The authors declare that there are no competing interests.

559

## 560 **References**

561 [1] K. Strimbu and J. A. Tavel, ‘What are Biomarkers?’, *Curr. Opin. HIV AIDS*, vol. 5, no.  
562 6, pp. 463–466, Nov. 2010, doi: 10.1097/COH.0b013e32833ed177.

563 [2] R. Mayeux, ‘Biomarkers: Potential uses and limitations’, *NeuroRX*, vol. 1, no. 2, pp. 182–  
564 188, Apr. 2004, doi: 10.1602/neurorx.1.2.182.

565 [3] S. Naylor, ‘Biomarkers: current perspectives and future prospects’, *Expert Rev. Mol.*  
566 *Diagn.*, vol. 3, no. 5, pp. 525–529, Sep. 2003, doi: 10.1586/14737159.3.5.525.

567 [4] S. G. Liao *et al.*, ‘Missing value imputation in high-dimensional phenomic data:  
568 imputable or not, and how?’, *BMC Bioinformatics*, vol. 15, no. 1, p. 346, Nov. 2014, doi:  
569 10.1186/s12859-014-0346-6.

570 [5] A. Dahl *et al.*, ‘A multiple-phenotype imputation method for genetic studies’, *Nat. Genet.*,  
571 vol. 48, no. 4, Art. no. 4, Apr. 2016, doi: 10.1038/ng.3513.

572 [6] K. T. Do *et al.*, ‘Characterization of missing values in untargeted MS-based metabolomics  
573 data and evaluation of missing data handling strategies’, *Metabolomics*, vol. 14, no. 10,  
574 2018, doi: 10.1007/s11306-018-1420-2.

575 [7] M. J. Azur, E. A. Stuart, C. Frangakis, and P. J. Leaf, ‘Multiple imputation by chained  
576 equations: what is it and how does it work?’, *Int. J. Methods Psychiatr. Res.*, vol. 20, no.  
577 1, pp. 40–49, Feb. 2011, doi: 10.1002/mpr.329.

578 [8] M. R. Malarvizhi and D. A. S. Thanamani, *K-Nearest Neighbor in Missing Data*  
579 *Imputation*.

580 [9] S. Cheng *et al.*, ‘Potential Impact and Study Considerations of Metabolomics in  
581 Cardiovascular Health and Disease A Scientific Statement From the American Heart  
582 Association’, *Circ. Cardiovasc. Genet.*, vol. 10, no. 2, Apr. 2017, doi:  
583 10.1161/HCG.0000000000000032.

- 584 [10] J. Liu *et al.*, ‘Integration of epidemiologic, pharmacologic, genetic and gut microbiome  
585 data in a drug-metabolite atlas’, *Nat. Med.*, vol. 26, no. 1, pp. 110–117, Jan. 2020, doi:  
586 10.1038/s41591-019-0722-x.
- 587 [11] L. M. ‘t Hart *et al.*, ‘Blood Metabolomic Measures Associate With Present and Future  
588 Glycemic Control in Type 2 Diabetes’, *J. Clin. Endocrinol. Metab.*, vol. 103, no. 12, pp.  
589 4569–4579, Dec. 2018, doi: 10.1210/jc.2018-01165.
- 590 [12] G. L. J. Onderwater *et al.*, ‘Large-scale plasma metabolome analysis reveals alterations  
591 in HDL metabolism in migraine’, *Neurology*, vol. 92, no. 16, pp. e1899–e1911, Apr.  
592 2019, doi: 10.1212/WNL.0000000000007313.
- 593 [13] I. C. van den Munckhof *et al.*, ‘Microbial Impact on Plasma Metabolites is Linked to the  
594 Cardiovascular Risk and Phenotypes’, *Atheroscler. Suppl.*, vol. 32, pp. 118–119, Jun.  
595 2018, doi: 10.1016/j.atherosclerosis.2018.04.366.
- 596 [14] D. Vojinovic *et al.*, ‘Metabolic profiling of intra- and extracranial carotid artery  
597 atherosclerosis’, *Atherosclerosis*, vol. 272, pp. 60–65, May 2018, doi:  
598 10.1016/j.atherosclerosis.2018.03.015.
- 599 [15] J. Tynkkynen *et al.*, ‘Association of branched-chain amino acids and other circulating  
600 metabolites with risk of incident dementia and Alzheimer’s disease: A prospective study  
601 in eight cohorts’, *Alzheimers Dement.*, vol. 14, no. 6, pp. 723–733, 2018, doi:  
602 <https://doi.org/10.1016/j.jalz.2018.01.003>.
- 603 [16] van den Akker Erik B. *et al.*, ‘Metabolic Age Based on the BBMRI-NL 1H-NMR  
604 Metabolomics Repository as Biomarker of Age-related Disease’, *Circ. Genomic Precis.  
605 Med.*, vol. 0, no. 0, doi: 10.1161/CIRCGEN.119.002610.
- 606 [17] J. Deelen *et al.*, ‘A metabolic profile of all-cause mortality risk identified in an  
607 observational study of 44,168 individuals’, *Nat. Commun.*, vol. 10, no. 1, pp. 1–8, Aug.  
608 2019, doi: 10.1038/s41467-019-11311-9.
- 609 [18] J. E. Ho *et al.*, ‘Metabolomic Profiles of Body Mass Index in the Framingham Heart Study  
610 Reveal Distinct Cardiometabolic Phenotypes’, *PloS One*, vol. 11, no. 2, p. e0148361,  
611 2016, doi: 10.1371/journal.pone.0148361.
- 612 [19] M. J. Rist *et al.*, ‘Metabolite patterns predicting sex and age in participants of the  
613 Karlsruhe Metabolomics and Nutrition (KarMeN) study’, *PLoS ONE*, vol. 12, no. 8, Aug.  
614 2017, doi: 10.1371/journal.pone.0183228.
- 615 [20] M. Schoenmaker *et al.*, ‘Evidence of genetic enrichment for exceptional survival using a  
616 family approach: the Leiden Longevity Study’, *Eur. J. Hum. Genet.*, vol. 14, no. 1, Art.  
617 no. 1, Jan. 2006, doi: 10.1038/sj.ejhg.5201508.
- 618 [21] P. Würtz *et al.*, ‘Metabolic Signatures of Adiposity in Young Adults: Mendelian  
619 Randomization Analysis and Effects of Weight Change’, *PLoS Med.*, vol. 11, no. 12, Dec.  
620 2014, doi: 10.1371/journal.pmed.1001765.
- 621 [22] P. Soininen, A. J. Kangas, P. Würtz, T. Suna, and M. Ala-Korpela, ‘Quantitative serum  
622 nuclear magnetic resonance metabolomics in cardiovascular epidemiology and genetics’,  
623 *Circ. Cardiovasc. Genet.*, vol. 8, no. 1, pp. 192–206, Feb. 2015, doi:  
624 10.1161/CIRCGENETICS.114.000216.
- 625 [23] S. A. Williams *et al.*, ‘Plasma protein patterns as comprehensive indicators of health’,  
626 *Nat. Med.*, vol. 25, no. 12, pp. 1851–1857, Dec. 2019, doi: 10.1038/s41591-019-0665-2.
- 627 [24] E. Crimmins, S. Vasunilashorn, J. K. Kim, and D. Alley, ‘BIOMARKERS RELATED  
628 TO AGING IN HUMAN POPULATIONS’, *Adv. Clin. Chem.*, vol. 46, pp. 161–216,  
629 2008.
- 630 [25] L. M. Biasucci, CDC, and AHA, ‘CDC/AHA Workshop on Markers of Inflammation and  
631 Cardiovascular Disease: Application to Clinical and Public Health Practice: clinical use  
632 of inflammatory markers in patients with cardiovascular diseases: a background paper’,

- 633 *Circulation*, vol. 110, no. 25, pp. e560-567, Dec. 2004, doi:  
634 10.1161/01.CIR.0000148983.88334.80.
- 635 [26] A. J. Flint *et al.*, ‘Body mass index, waist circumference, and risk of coronary heart  
636 disease: a prospective study among men and women’, *Obes. Res. Clin. Pract.*, vol. 4, no.  
637 3, pp. e171–e181, 2010, doi: 10.1016/j.orcp.2010.01.001.
- 638 [27] L. Dean and L. Dean, *Blood Groups and Red Cell Antigens*. National Center for  
639 Biotechnology Information (US), 2005.
- 640 [28] R. J. Glassock and C. Winearls, ‘Ageing and the Glomerular Filtration Rate: Truths and  
641 Consequences’, *Trans. Am. Clin. Climatol. Assoc.*, vol. 120, pp. 419–428, 2009.
- 642 [29] E. Pinto, ‘Blood pressure and ageing’, *Postgrad. Med. J.*, vol. 83, no. 976, pp. 109–114,  
643 Feb. 2007, doi: 10.1136/pgmj.2006.048371.
- 644 [30] C. M. Schubert *et al.*, ‘Lipids, lipoproteins, lifestyle, adiposity and fat-free mass during  
645 middle age: the Fels Longitudinal Study’, *Int. J. Obes.*, vol. 30, no. 2, Art. no. 2, Feb.  
646 2006, doi: 10.1038/sj.ijo.0803129.
- 647 [31] V. Saini, ‘Molecular mechanisms of insulin resistance in type 2 diabetes mellitus’, *World*  
648 *J. Diabetes*, vol. 1, no. 3, pp. 68–75, Jul. 2010, doi: 10.4239/wjd.v1.i3.68.
- 649 [32] C. K. Roberts, A. L. Hevener, and R. J. Barnard, ‘Metabolic Syndrome and Insulin  
650 Resistance: Underlying Causes and Modification by Exercise Training’, *Compr. Physiol.*,  
651 vol. 3, no. 1, pp. 1–58, Jan. 2013, doi: 10.1002/cphy.c110062.
- 652 [33] ‘High White Blood Cell Count Is Associated With a Worsening of Insulin Sensitivity and  
653 Predicts the Development of Type 2 Diabetes | Diabetes’.  
654 <https://diabetes.diabetesjournals.org/content/51/2/455> (accessed Feb. 09, 2021).
- 655 [34] D. C. Muller, D. Elahi, J. D. Tobin, and R. Andres, ‘The effect of age on insulin resistance  
656 and secretion: a review’, *Semin. Nephrol.*, vol. 16, no. 4, pp. 289–298, Jul. 1996.
- 657 [35] P. J. Enriori, A. E. Evans, P. Sinnayah, and M. A. Cowley, ‘Leptin Resistance and  
658 Obesity’, *Obesity*, vol. 14, no. S8, pp. 254S-258S, 2006, doi:  
659 <https://doi.org/10.1038/oby.2006.319>.
- 660 [36] O. Ukkola and Y. A. Kesäniemi, ‘Leptin and high-sensitivity C-reactive protein and their  
661 interaction in the metabolic syndrome in middle-aged subjects’, *Metabolism*, vol. 56, no.  
662 9, pp. 1221–1227, Sep. 2007, doi: 10.1016/j.metabol.2007.04.019.
- 663 [37] P. W. Franks *et al.*, ‘Leptin Predicts a Worsening of the Features of the Metabolic  
664 Syndrome Independently of Obesity’, *Obes. Res.*, vol. 13, no. 8, pp. 1476–1484, 2005,  
665 doi: <https://doi.org/10.1038/oby.2005.178>.
- 666 [38] M. I. Schmidt *et al.*, ‘Leptin and incident type 2 diabetes: risk or protection?’,  
667 *Diabetologia*, vol. 49, no. 9, pp. 2086–2096, Sep. 2006, doi: 10.1007/s00125-006-0351-  
668 z.
- 669 [39] E. Nigro *et al.*, ‘New insight into adiponectin role in obesity and obesity-related diseases’,  
670 *BioMed Res. Int.*, vol. 2014, p. 658913, 2014, doi: 10.1155/2014/658913.
- 671 [40] M. Matsubara, K. Namioka, and S. Katayose, ‘Decreased plasma adiponectin  
672 concentrations in women with low-grade C-reactive protein elevation’, *Eur. J.*  
673 *Endocrinol.*, vol. 148, no. 6, pp. 657–662, Jun. 2003, doi: 10.1530/eje.0.1480657.
- 674 [41] T. B. Harris *et al.*, ‘Associations of elevated Interleukin-6 and C-Reactive protein levels  
675 with mortality in the elderly\*\*Access the “Journal Club” discussion of this paper at  
676 <http://www.elsevier.com/locate/ajmselect/>’, *Am. J. Med.*, vol. 106, no. 5, pp. 506–512,  
677 May 1999, doi: 10.1016/S0002-9343(99)00066-2.
- 678 [42] J. Helmersson, A. Larsson, B. Vessby, and S. Basu, ‘Active smoking and a history of  
679 smoking are associated with enhanced prostaglandin F<sub>2α</sub>, interleukin-6 and F<sub>2</sub>-  
680 isoprostane formation in elderly men’, *Atherosclerosis*, vol. 181, no. 1, pp. 201–207, Jul.  
681 2005, doi: 10.1016/j.atherosclerosis.2004.11.026.

- 682 [43] C. E. Byrne, A. Fitzgerald, C. P. Cannon, D. J. Fitzgerald, and D. C. Shields, 'Elevated  
683 white cell count in acute coronary syndromes: relationship to variants in inflammatory  
684 and thrombotic genes', *BMC Med. Genet.*, vol. 5, no. 1, p. 13, Jun. 2004, doi:  
685 10.1186/1471-2350-5-13.
- 686 [44] Y. Wang *et al.*, 'Sex differences in the association between diabetes and risk of  
687 cardiovascular disease, cancer, and all-cause and cause-specific mortality: a systematic  
688 review and meta-analysis of 5,162,654 participants', *BMC Med.*, vol. 17, no. 1, p. 136,  
689 Jul. 2019, doi: 10.1186/s12916-019-1355-0.
- 690 [45] K. M. Dale, C. I. Coleman, S. A. Shah, A. A. Patel, J. Kluger, and C. M. White, 'Impact  
691 of gender on statin efficacy', *Curr. Med. Res. Opin.*, vol. 23, no. 3, pp. 565–574, Mar.  
692 2007, doi: 10.1185/030079906X167516.
- 693 [46] Y. An, J. Jang, S. Lee, S. Moon, and S. K. Park, 'Sex-specific Associations Between  
694 Serum Hemoglobin Levels and the Risk of Cause-specific Death in Korea Using the  
695 National Health Insurance Service-National Health Screening Cohort (NHIS HEALS)',  
696 *J. Prev. Med. Public Health Yebang Uihakhoe Chi*, vol. 52, no. 6, pp. 393–404, Nov.  
697 2019, doi: 10.3961/jpmph.19.146.
- 698 [47] E. Prescott *et al.*, 'Mortality in women and men in relation to smoking', *Int. J. Epidemiol.*,  
699 vol. 27, no. 1, pp. 27–32, Feb. 1998, doi: 10.1093/ije/27.1.27.
- 700 [48] P. Muennig, E. Lubetkin, H. Jia, and P. Franks, 'Gender and the Burden of Disease  
701 Attributable to Obesity', *Am. J. Public Health*, vol. 96, no. 9, pp. 1662–1668, Sep. 2006,  
702 doi: 10.2105/AJPH.2005.068874.
- 703 [49] B. T. Palmisano, L. Zhu, R. H. Eckel, and J. M. Stafford, 'Sex differences in lipid and  
704 lipoprotein metabolism', *Mol. Metab.*, vol. 15, pp. 45–55, May 2018, doi:  
705 10.1016/j.molmet.2018.05.008.
- 706 [50] Y. Li *et al.*, 'Hs-CRP and all-cause, cardiovascular, and cancer mortality risk: A meta-  
707 analysis', *Atherosclerosis*, vol. 259, pp. 75–82, Apr. 2017, doi:  
708 10.1016/j.atherosclerosis.2017.02.003.
- 709 [51] J. Liang, J. M. Bennett, H. Sugisawa, E. Kobayashi, and T. Fukaya, 'Gender differences  
710 in old age mortality: roles of health behavior and baseline health status', *J. Clin.*  
711 *Epidemiol.*, vol. 56, no. 6, pp. 572–582, Jun. 2003, doi: 10.1016/s0895-4356(03)00060-  
712 x.
- 713 [52] M. A. Davis, J. M. Neuhaus, D. J. Moritz, D. Lein, J. D. Barclay, and S. P. Murphy,  
714 'Health behaviors and survival among middle-aged and older men and women in the  
715 NHANES I Epidemiologic Follow-up Study', *Prev. Med.*, vol. 23, no. 3, pp. 369–376,  
716 May 1994, doi: 10.1006/pmed.1994.1051.
- 717 [53] Y. Zhang *et al.*, 'Gender difference in cardiovascular risk factors in the elderly with  
718 cardiovascular disease in the last stage of lifespan: The PROTEGER study', *Int. J.*  
719 *Cardiol.*, vol. 155, no. 1, pp. 144–148, Feb. 2012, doi: 10.1016/j.ijcard.2011.09.073.
- 720 [54] F. Dudbridge, 'Power and Predictive Accuracy of Polygenic Risk Scores', *PLoS Genet.*,  
721 vol. 9, no. 3, Mar. 2013, doi: 10.1371/journal.pgen.1003348.
- 722 [55] C. M. Lewis and E. Vassos, 'Prospects for using risk scores in polygenic medicine',  
723 *Genome Med.*, vol. 9, Nov. 2017, doi: 10.1186/s13073-017-0489-y.
- 724 [56] A. V. Khera *et al.*, 'Genome-wide polygenic scores for common diseases identify  
725 individuals with risk equivalent to monogenic mutations', *Nat. Genet.*, vol. 50, no. 9, Art.  
726 no. 9, Sep. 2018, doi: 10.1038/s41588-018-0183-z.
- 727 [57] S. A. Lambert *et al.*, 'The Polygenic Score Catalog as an open database for reproducibility  
728 and systematic evaluation', *Nat. Genet.*, pp. 1–6, Mar. 2021, doi: 10.1038/s41588-021-  
729 00783-5.

- 730 [58] T. G. Richardson, S. Harrison, G. Hemani, and G. Davey Smith, ‘An atlas of polygenic  
731 risk score associations to highlight putative causal relationships across the human  
732 phenome’, *eLife*, vol. 8, doi: 10.7554/eLife.43657.
- 733 [59] P. Würtz *et al.*, ‘Metabolite profiling and cardiovascular event risk: a prospective study  
734 of 3 population-based cohorts’, *Circulation*, vol. 131, no. 9, pp. 774–785, Mar. 2015, doi:  
735 10.1161/CIRCULATIONAHA.114.013116.
- 736 [60] A. V. Ahola-Olli *et al.*, ‘Circulating metabolites and the risk of type 2 diabetes: a  
737 prospective study of 11,896 young adults from four Finnish cohorts’, *Diabetologia*, vol.  
738 62, no. 12, pp. 2298–2309, 2019, doi: 10.1007/s00125-019-05001-w.
- 739 [61] Nightingale Health UK Biobank Initiative, H. Julkunen, A. Cichońska, P. E. Slagboom,  
740 and P. Würtz, ‘Metabolic biomarker profiling for identification of susceptibility to severe  
741 pneumonia and COVID-19 in the general population’, *eLife*, vol. 10, p. e63033, May  
742 2021, doi: 10.7554/eLife.63033.
- 743 [62] W. E. Johnson, C. Li, and A. Rabinovic, ‘Adjusting batch effects in microarray expression  
744 data using empirical Bayes methods’, *Biostatistics*, vol. 8, no. 1, pp. 118–127, Jan. 2007,  
745 doi: 10.1093/biostatistics/kxj037.
- 746 [63] D. S. Wishart *et al.*, ‘HMDB 4.0: the human metabolome database for 2018’, *Nucleic  
747 Acids Res.*, vol. 46, no. Database issue, pp. D608–D617, Jan. 2018, doi:  
748 10.1093/nar/gkx1089.
- 749 [64] P. Würtz, A. J. Kangas, P. Soininen, D. A. Lawlor, G. Davey Smith, and M. Ala-Korpela,  
750 ‘Quantitative Serum Nuclear Magnetic Resonance Metabolomics in Large-Scale  
751 Epidemiology: A Primer on -Omic Technologies’, *Am. J. Epidemiol.*, vol. 186, no. 9, pp.  
752 1084–1096, Nov. 2017, doi: 10.1093/aje/kwx016.
- 753 [65] W. T. Friedewald, R. I. Levy, and D. S. Fredrickson, ‘Estimation of the Concentration of  
754 Low-Density Lipoprotein Cholesterol in Plasma, Without Use of the Preparative  
755 Ultracentrifuge’, *Clin. Chem.*, vol. 18, no. 6, pp. 499–502, Jun. 1972, doi:  
756 10.1093/clinchem/18.6.499.
- 757 [66] Carrero Juan Jesus, Andersson Franko Mikael, Obergfell Achim, Gabrielsen Anders, and  
758 Jernberg Tomas, ‘hsCRP Level and the Risk of Death or Recurrent Cardiovascular Events  
759 in Patients With Myocardial Infarction: a Healthcare-Based Study’, *J. Am. Heart Assoc.*,  
760 vol. 8, no. 11, p. e012638, Jun. 2019, doi: 10.1161/JAHA.119.012638.  
761