# Human blood lipoprotein predictions from [1]H NMR spectra: protocol, model performances and cage of covariance

Bekzod Khakimov[a,1], Huub C.J. Hoefsloot[b], Nabiollah Mobaraki[c], Violetta Aru[a], Mette Kristensen[d], Mads V. Lind[d], Lars Holm[e], Josué L Castro-Mejía[a], Dennis S Nielsen[a], Doris M. Jacobs[f], Age K. Smilde[a,b], Søren Balling Engelsen[a,2]

[a]Department of Food Science, University of Copenhagen, Rolighedsvej 26, DK-1958 Frederiksberg C, Denmark
[b]Swammerdam Institute for Life Sciences, University of Amsterdam, Postbus 94215, Amsterdam 1090 GE, The Netherlands
[c]Department of Chemistry, Faculty of Science, Shiraz University, Shiraz, 7194684795, Iran
[d]Department of Nutrition, Exercise and Sports, University of Copenhagen, Rolighedsvej 26, DK-1958 Frederiksberg C, Denmark
[e] School of Sport, Exercise and Rehabilitation Sciences, University of Birmingham, Edgbaston, Birmingham, B15 2TT, UK
[f]Unilever Global Food Innovation Centre, 6708 WH Wageningen, The Netherlands

[1]Corresponding author: Bekzod Khakimov, Department of Food Science, University of Copenhagen, Rolighedsvej 26, DK-1958 Frederiksberg C, Denmark. Tel.; +45-2887-4454, Email: bzo@food.ku.dk, ORCID: 0000-0002-6580-2034

[2]Corresponding author: Søren Balling Engelsen, Department of Food Science, University of Copenhagen, Rolighedsvej 26, DK-1958 Frederiksberg C, Denmark. Tel.: +45-2020-0064, Email: se@food.ku.dk, ORCID: 0000-0003-4124-4338

**Author Contributions:** S.B.E., D.S.N., A.K.S. designed research, B.K., H.C.J.H., N.M., V.A. performed research, B.K., V.A., M.K., M.V.L., L.H., J.L.C.M., D.M.J. acquired data, B.K. and N.M. developed software, B.K. analysed data, B.K. and S.B.E. wrote the paper, all authors read and approved the paper.

**Competing Interest Statement:** The authors declare no competing interest.

**Keywords:** Lipoprotein, biomarker, metabolomics, NMR, PLS, cage of covariance

35 **Abstract**

36 Lipoprotein subfractions are biomarkers for early diagnosis of cardiovascular diseases. The reference method,

37 ultracentrifugation, for measuring lipoproteins is time consuming and there is a need to develop a rapid method

38 for cohort screenings. Here we present partial least squares regression models developed using $^1$H-NMR spectra

39 and concentrations of lipoproteins as measured by ultracentrifugation on 316 healthy Danes. Different regions of

40 the $^1$H-NMR spectra representing signals of the lipoproteins and different lipid species were investigated to

41 develop parsimonious, reliable and best performing prediction models. 65 LP main and subfractions were

42 predictable with an accuracy $Q^2$ of $> 0.6$ on test set samples. The models were tested on an independent cohort of

43 290 healthy Swedes with predicted and reference values matching by up to 85-95%. The software was developed

44 to predict lipoproteins in human blood using $^1$H-NMR spectra and made freely available to be applied for future

45 cohort screenings.

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

## Introduction

66

67    Cardiovascular diseases (CVD) are still the leading cause of mortality and morbidity [1]. Although the trend in CVD
68    mortality has plateaued in many countries [1], ~18 million people die annually worldwide due to CVD [2]. Blood
69    lipids, including lipoproteins (LP), play an important role in the development of this pathology and serve as
70    diagnostic markers. Total cholesterol (*chol*) in blood and the level of low density lipoprotein particles (LDL*chol*)
71    have long been used as markers of risk of CVD [3-5]. Previous studies have shown that ratios of *chol* or LDL*chol* to
72    high density lipoprotein particles (HDL*chol*, "good" LP) are strongly associated with CVD risk [6,7]. A recent study
73    has shown that very low density lipoprotein (VLDL) and intermediate density lipoprotein (IDL) are also
74    associated with CVD risk and should be integrated into clinical practice as secondary targets of lipid-lowering
75    therapy [8]. Moreover, an increase in LDL*chol* to HDL*chol* ratio has been suggested to be a sign of developing
76    CVD, such as myocardial infarction, as well as being a predictor of CVD mortality [8,9].

77    Lipoproteins are micellar-like particles, with heterogeneous density and size, made up primarily of lipids and
78    proteins [10]. The inside of the typical LP particle is composed of triglycerides (*tg*) and cholesterol esters (*chole*).
79    The outer shell is composed of free cholesterols (*fchol*), phospholipids (*phosl*) and apolipoproteins (*apoA* and
80    *apoB*). During fasting, LP in human blood can be divided into four main fractions based on density and size:
81    VLDL, IDL, LDL and HDL. Very low density lipoproteins are the largest particles with the lowest density, and
82    HDL are the smallest particles with the highest density. There are many analytical methods, including
83    ultracentrifugation (UC), gel-electrophoresis, high-performance liquid chromatography (HPLC) and numerous
84    assays, to determine these main fractions of LP particles and their subfractions in human blood plasma and serum.
85    However, the definition of subfractions differs between the various separation techniques. The LP particles are
86    not clearly distinct groups of particles, but rather form a distribution of particles differing in size and density, and
87    in lipid and protein composition [11]. For example, UC can separate seven subfractions of LDL, while HPLC can
88    separate five subfractions, and these subfractions are not directly comparable. Ultracentrifugation has become the
89    clinical reference method thanks to its capacity to better separate LP subfractions [10]. However, not all aspects of
90    density and/or size of all subfractions are fully determined, and their biochemical function and physiological roles
91    in human metabolism are only partly understood. The usefulness of LP as biomarkers is a strong impetus for the
92    development of a rapid and accurate quantification method.

93    A promising analytical method for the rapid quantification of LP in human blood plasma is proton ([1]H) nuclear
94    magnetic resonance (NMR) spectroscopy. In 1991-1992 Otvos *et al.* demonstrated that [1]H NMR spectra and curve
95    fitting of the plasma lipid methyl signal envelope (0.87-0.67 ppm) can be used to predict the distribution of the
96    main LP fractions and subfractions [12,13]. This seminal work was greatly expanded by Ala-Korpela *et al.*, who
97    investigated multiple NMR effects that can influence the quantification of LP, including spectral regions and time-
98    domain versus frequency domain solutions [14-16]. In 2005 Petersen *et al.* applied a more pragmatic multivariate
99    data analytical approach to this work. They applied more robust partial least squares (PLS) regression [17] for
100   quantifying LP main and subfractions by regressing the reference values measured by UC against the lipid region
101   (5.7-0.2 ppm) of the [1]H NMR spectra (600 MHz) [18]. This approach was soon applied by others and there are now
102   more than 20 studies demonstrating LP quantification by combining [1]H NMR spectroscopy with a reference
103   measurement method, predominantly UC [10]. More recently, it has become clear that optimization and
104   standardization of the [1]H NMR measurement protocols is of paramount importance if reliable, accurate and

3

105   comparable LP quantifications are to be obtained across different laboratories [19]. Consensus regarding
106   standardization is that high resolution, typically 600 MHz, NMR spectrometer be used for the measurement of
107   one-dimensional (1D) nuclear Overhauser effect spectroscopy (NOESY) $^1$H NMR spectra of blood plasma or
108   serum samples for quantification of LP.

109   The main advantage of using the 1D NOESY pulse sequence is that it is fast, robust and simple, and thus suitable
110   for standardization. Other more sophisticated NMR techniques, such as 2D Diffusion edited NMR (DOSY), which
111   is sensitive to diffusion and can even recover the pure underlying NMR spectra of LP fractions using parallel
112   factor analysis (PARAFAC) [20], are less suitable for standardization because of their longer acquisition times, and
113   a more complex pulse sequence using field gradients. Reproducible and standardized measurement protocols
114   allow data fusion of multiple cohorts for continued improvements of the calibration models, which continues
115   accumulation of knowledge.

116   This study shows the development and extensive validation of PLS models to predict concentrations of LP main
117   and subfractions using blood plasma $^1$H NMR spectra and corresponding LP data measured by UC for 300+
118   volunteers. Furthermore, it provides an integrated software to predict LP from NMR spectra in future studies. The
119   datasets and the software have been made freely available to the public. The $^1$H NMR spectra were measured
120   using the most recent standard operating procedures (SOP) covering blood collection, sample handling, and NMR
121   data acquisition, which have been published previously [19,21]. To the best of our knowledge, this study is the first
122   to illustrate the prediction performances of PLS models for a wide range of LP subfractions, prediction model
123   parameters and complexity, and provide open access datasets of $^1$H NMR spectra and corresponding LP data. The
124   study investigated regression models using different parts of the entire range of $^1$H NMR spectra in order to
125   identify the best regions representing signals of LP and different lipid species with an aim to develop the simplest,
126   most robust and best performing PLS prediction models. The prediction of nearly 100 parameters from a single
127   NMR spectrum raises concerns about the covariations amongst the independent variables (the reference
128   parameters), and about the rank of the dependent variables (the spectra). This study therefore investigated the rank
129   of the NMR and LP datasets to describe and understand the level of inter-correlations *"cage of covariance"*
130   between the individual LP and the information content in the NMR spectra. Finally, the PLS models developed
131   were validated by prediction of LP concentrations in an independent cohort of 290 healthy Swedes [22].

132

133   **Results**

134   **Overview of the implemented workflow**

135   Figure 1 shows an overview of the workflow implemented to predict LP in human blood using $^1$H NMR spectra
136   and UC data as response variables by applying PLS modelling. As described previously [19], the NMR spectra were
137   scaled to the electronic reference to access *in vivo* concentrations (ERETIC) signal positioned at 15 ppm,
138   equivalent to 10 mmol L$^{-1}$ protons and aligned towards the doublet of alanine's methyl group (1.507-1.494 ppm)
139   using *icoshift* [23]. Scaling the NMR spectra according to the area of the ERETIC signal minimizes variations
140   originating from the NMR experiment and allows for inter-laboratory data comparisons. Unlike urine NMR
141   spectra [24], shifting of the entire $^1$H NMR spectra of blood plasma samples towards the doublet of alanine is
142   sufficient to eliminate minor misalignments present in the spectra due to small differences in pH and/or

4

143  experimental error (e.g. plasma to buffer ratio, small fluctuations on spectrum acquisition temperature). It should

144  be noted that any additional spectral alignment shifting affecting the methyl (0.92-0.8 ppm) and methylene (1.4-

145  1.2 ppm) regions will impair the quantification of LP, as it is the small shifts of the lipid signals that are modelled

146  by PLS. 1D $^1$H NMR spectra of human blood plasma contain ~13 unique spectral regions representing signals of

147  the blood lipids. A total of 33 NMR datasets (Figure S1) were constructed using the 13 regions containing signals

148  of lipids, either alone or in various combinations, and were subsequently used to develop PLS models with an

149  optimal prediction performance (*vide supra*) (Figure 1). Briefly, an individual PLS model was developed for each

150  LP variable. Firstly, a training model was developed using randomly selected 70% of samples, and was then tested

151  on the remaining 30% of samples. The optimal number of latent variables (LV) for the training models was

152  selected by developing one to twenty LV PLS models using 10 fold cross validation for each model, and the model

153  with the smallest root mean square error of cross validation (RMSECV) value was chosen.

154  Ultracentrifugation determined LP concentrations (milligram per decilitre, mg/dL) in fresh fasting plasma

155  samples. In total, 97 variables were used as response variables for the development of the PLS models. Samples

156  with LP concentrations below limit of detection (LOD) or missing value for a given LP variable were removed

157  prior to PLS modelling. In addition, a few samples with an extraordinary underperforming PLS prediction of a

158  LP variable were removed prior to the development of the final models (Figure 1). These extreme samples could

159  primarily be related to the relatively high uncertainty of the UC measurements of a few individuals due to faulty

160  results from lipids assays used in the UC. The uncertainty of the UC method for LP quantification ranged from

161  5% to 40%, depending mainly on the molecular type [25]. The lowest uncertainties were observed for apolipoprotein

162  A (*ApoA*) (5-15%), apolipoprotein B (*ApoB*) (6-20%), and cholesterol (5-18%) molecules, while the highest

163  uncertainty was observed for free cholesterol (8-40%). More details on reproducibility and source of uncertainty

164  on reference method for LP measurement are published elsewhere [25].

165

166  **Optimal regions of the $^1$H NMR spectra for lipoprotein quantification**

167  When performing regression on spectra, a selection of an optimal spectral region is important for developing

168  parsimonious models with less noise and fewer interferences, especially when spectra are large (e.g. >16k

169  variables). For this reason, the $^1$H NMR spectra of blood plasma acquired from 316 Danish volunteers were

170  divided into 13 spectral regions representing signals from different lipid functional groups. These regions include

171  the signal of the methyl group (C18) of cholesterol (0.75-0.65 ppm), the signals of methyl (0.92-0.8 ppm),

172  methylene (1.4-1.2 ppm), and methine (5.40-5.24 ppm) protons of different LP molecules, including triglycerides,

173  phospholipids, apolipoproteins (*apoA* and *apoB*), as well as spectral regions containing signals from other lipid

174  functional groups (Figure S1). Each spectral region was used either alone or in combination with other regions,

175  resulting in 33 NMR datasets (including the entire NMR spectrum), to develop the PLS regression models for

176  predicting LP concentrations in human blood plasma from UC reference values.

177  The concentrations of 65 of the 97 measured LP were found to be predictable from the majority of the NMR

178  datasets (Figure 2a). One criterion set for an acceptable prediction performance of the PLS model was to have a

179  prediction accuracy ($Q^2$) > 0.6 for test set samples. $Q^2$ is a statistical measure of prediction accuracy often used in

180  PLS modelling [26] and defined as $Q^2 = (1-PRESS/SS)$, where PRESS is predictive residual sum of squares and SS

5

181     is sum of squares of actual values (LP concentrations). Total triglyceride and total cholesterol in blood plasma

182     and in main fractions were predictable ($Q^2 > 0.6$) from all 33 spectral datasets apart from spectral regions 26 and

183     31 (Figure S1), which represent only the choline head group in lipids and the aromatic protons, respectively. Not

184     surprisingly, these two regions were also the worst performing for predicting the other LP since they do not carry

185     signals of LP molecules. They were therefore removed before further analysis. Spectral region 18, containing a

186     weak signal of lysine residue in albumin and more abundant signals of creatine, creatinine, was also not suitable

187     for prediction of most LP, and only total concentrations of *tg, chol, apoA,* and *chole* were weakly predictable. The

188     most up field spectral region, representing the protons in the methyl group (C18) of cholesterol (region 1 in Figure

189     S1), was able to predict all LP molecular classes in plasma and main fractions, as well as *tg, chol, fchol, phosl,*

190     and *apoB* in some subfractions. The spectral regions representing protons of the methylene groups, located either

191     one or two bonds away from carbonyl group or double bond of lipids (regions 11, 12, 13, 17, 27 in Figure S1) are

192     primarily able to predict total concentrations of LP in blood plasma and the main fractions. The spectral regions

193     23 and 25 (Figure S1), corresponding to glyceryl protons of lipids, C*H₂*OCOR and C*H*OCOR, respectively,

194     exhibit relatively high prediction power for total *tg* concentration in plasma and main fractions, as well as *tg* of

195     VLDL and IDL. The performance of these regions in predicting LP subfractions is sub-optimal, and none of them

196     are able to predict all subfractions. Thus, all these spectral regions could only selectively predict some LP, and

197     were mostly limited to prediction of total concentrations of LP in plasma and/or their cumulative concentrations

198     of main fractions. However, 20 of the 33 investigated NMR regions displayed relatively high prediction

199     performances for all 65 LP (Figure 2a).These 20 NMR datasets were therefore used for further modelling. The

200     prediction performance parameters of the PLS models developed for all 65 predictable LP using the 20 NMR

201     regions are shown in Table S1.

202     The five spectral regions (regions 7-10,14 in Figure 2b) representing the main signals of the LP, including methyl

203     and methylene protons as well as methyl of cholesterol, showed the best prediction results, with $Q^2 > 0.6$ for all

204     65 predictable LP. Overall, the prediction performances of all five spectral regions were similar, with a mean $Q^2$

205     of ~0.84 and coefficients of variation (CV) of ~13% for the test set prediction models. However, for a few LP,

206     the prediction performances differed significantly between the five regions. Unlike the other spectral regions,

207     region 8 (Figure 2b) does not contain the methyl signal of cholesterol, though this did not influence the prediction

208     performance for *chol* content in different LP particles. In contrast, a significantly lower prediction performance

209     was observed for *phosl* in the LDL1 subfraction using region 8. The $Q^2$ of the region 8 (Figure 2a) based model

210     for LDL1*phosl* was 0.7, while the other LP regions exhibit a $Q^2$ of at least 0.77. Regions 9 and 10 (Figure 2b),

211     representing the entire LP region of 1.4-0.6 ppm, displayed a better prediction performance for *tg, chol,* and *chole*

212     of LDL1, *chol* of LDL4, *chole* of HDL2b, and of *chol* and *apoA* of HDL3 subfractions, compared to the other

213     three regions. Despite interfering signals from non-LP related molecules such as lactic acid, valine, leucine,

214     isoleucine, the spectral regions 9 and 10 showed an overall better performance in prediction of all 65 LP than

215     regions 7, 8 and 14, where the signals of the non-LP molecules were removed. Spectral regions containing only

216     the methyl (region 2 and 6) or the methylene (region 3 and 4) protons of LP performed significantly worse

217     compared to regions 8 and 9, where the two proton populations are combined. This was especially pronounced in

218     the prediction performances of the PLS models developed for LDL and HDL subfractions (Table S1).

219     Interestingly, region 15, which contains only signals of methylene protons located either one or two bonds away

220     from a carbonyl group or double bond of the fatty acid chain, is also able to predict 60 out of 65 LP, with a $Q^2$ of

6

221    at least 0.6 for test set samples. However, the prediction performances of region 15 for many LP are significantly

222    lower than those of the spectral regions 9 or 10.Similarly, region 24, which represents olefinic protons of lipids,

223    enabled predictions of 58 of 65 LP, with a $Q^2 > 0.6$ for test set samples. Using a larger part of the spectral regions

224    "as is" (region 16, 22, 28) or after selecting only LP related signals (region 20, 21, 29, 30) gave prediction

225    performance similar to that of regions 2-4, with $Q^2 > 0.6$ for the majority of LP in test set prediction. Finally, the

226    use of the entire spectral range of 9.8-0.6 ppm (region 32) showed significantly lower LP prediction performances

227    than regions 9 or 10.

228    Despite comparable LP prediction performances of several spectral regions, especially regions 7-10 and 14

229    (Figure 2a), region 9 (1.4-0.6 ppm, later referred to as LP region) proved to be an optimal spectral region, showing

230    consistently high prediction performances for all 65 LP. This LP region is also the simplest one to extract from

231    the entire spectra, and unlike other regions it does not require the removal of interferences (non-LP molecules)

232    from narrow ranges of spectral intervals. Thus, the use of LP region minimizes complications that may arise when

233    dealing with a large number of samples, such as possible chemical shifts, occurrence of unforeseen signals, or a

234    high spectral complexity, which may impair reliable LP predictions across laboratories. Comparison of the relative

235    standard deviation (RSD) of the 65 LP predicted from NMR spectra of the 40 quality control (QC) pooled blood

236    plasma samples using the LP region and other spectral regions showed up to 10% lower RSD values in favour to

237    the LP region (data not shown). Thus, it was decided to use the LP region for the development of optimal PLS

238    models for NMR-based LP prediction.

239

240    **Prediction performance of PLS models using the LP region (1.4-0.6 ppm)**

241    Concentrations of 65 LP including main and subfractions were predictable from the [1]H NMR spectral region of

242    1.4-0.6 ppm (region 9; Figure 2b) with a $Q^2$ of 0.6 or higher (Table S2). The PLS regression models varied in

243    terms of their prediction performance and complexity. From 4 to 17 latent variables (LV) are required in order to

244    obtain optimal PLS regression models. More than 40 LP predictions require 10 or less LV. Only one LP

245    (LDL4*chol*) needed 15 LV, and three LP (*apoB* in LDL3, LDL4, and LDL6) needed 17 LV. The PLS regression

246    models for total concentrations of LP in plasma and of the main LP fractions were found to be less complex than

247    the corresponding models for the subfractions. Models for *phoslp* were generally less complex than the models

248    for *chol* or *apoB*. For example, in the same subfraction, LDL1, the *phoslp* model required 7 LV while the *chol*

249    and *apoB* models required 15 and 17 LV respectively. No clear trend explaining model complexity by LP

250    molecular type, particle type or particle size was observed. However, it is assumed that the number of LV in PLS

251    regression models is mainly determined by three factors: the chemical complexity of the LP (number of

252    representative NMR signals and their resolution), the variation range present in a cohort (concentration span of

253    the individual LP) and the presence of spectral interferences (overlap of signals of LP and non-LP molecules).

254    Figure 3 shows the PLS predictions of selected LPs representing *tg, chol, phoslp,* and *apoA* (see Figure S2 for all

255    65 LP). Concentrations of all seven LP molecules, *tg, chol, fchol, chole, phoslp, apoA,* and *apoB* in blood plasma

256    were predictable and showed test set $Q^2$ and CV values of 0.87-0.97 and 3-6%, respectively (Table S2). A

257    cumulative concentration of VLDL, IDL, LDL, and HDL particles (e.g., *chol*_main_fraction = VLDL*chol* +

258    IDL*chol* + LDL*chol* + HDL*chol*) was also predicted with high accuracy. For test set samples, $Q^2$ and CV values

259    ranged between 0.70-0.97 and 4-8% respectively, for all main fractions. Partial least squares regression models

260    performed slightly worse in predicting individual LP molecules in each main fraction separately (e.g. VLDL*chol*)

261    than for cumulative amounts across all main fractions or in plasma. For example, $Q^2$ and CV of models predicting

262    concentration of *chol* in VLDL, IDL, LDL, and HDL ranged between 0.87-0.94 and 7-20% respectively, while

263    the $Q^2$ and CV of the model predicting total level of *chol* in main fractions was 0.95 and 5% respectively. Similar

264    trends were observed for *tg, chole, phoslp, apoA,* and *apoB* molecules. Consistent for all seven LP molecules, the

265    prediction performances of the PLS models developed for the main fractions were better than for their

266    corresponding subfractions. For example, the $Q^2$ and CV of the LDL*phoslp* model were 0.83 and 13%

267    respectively, while the $Q^2$ and CV values obtained from the PLS models developed for subfractions of this particle,

268    including LDL1*phoslp*, LDL2*phoslp*, LDL3*phoslp*, and LDL4*phoslp*, had a range of 0.67-0.76 and 17-19%

269    respectively. Overall, the PLS models of all 65 LP showed a $Q^2$ of at least 0.6 for the test set samples with a mean,

270    quartile (Q) 50%, Q 75%, and Q 90% values of 0.84, 0.86, 0.93, and 0.95 respectively. Similarly, the coefficients

271    of variations (CV) of predicted LP concentrations in the test set samples were relatively low for all 65 LP and

272    ranged between 3-30%, with mean, Q 50%, Q 75%, and Q 90% values of 13, 13, 17, and 21 respectively.

273    Generally, the prediction performances of the PLS models developed on training set samples were similar to the

274    test set models (differences between the training and test set RMSE values were <5%). Overall, the prediction

275    performances of PLS regression models for different LP molecules decreased with specificity and can be ordered

276    as follows: total concentration in blood plasma > cumulative in all main fractions > individual main fractions >

277    subfractions.

278    In summary, a total of 65 of 97 LP measured using UC were predictable with reasonable performance ($Q^2$ >0.6

279    for the test set). The prediction of the remaining 32 LP variables was sub-optimal, with a relatively low test set

280    $Q^2$ of 0.3-0.6. These models were therefore deemed unreliable and, as a minimum, require additional data (NMR

281    and corresponding UC data) in order to be improved. It is assumed that there are two main reasons behind these

282    suboptimal predictions: 1) a high uncertainty of the reference method (UC) related to assay limitations and freeze-

283    thaw cycle of plasma samples, and 2) a lack of variability and/or close to LOD levels of those 32 LP concentrations

284    in the investigated cohort. A previous study found that the lowest repeatability in UC based LP quantification was

285    observed for *fchol*, *phoslp* and *tg* molecules [25]. This is in agreement with our PLS modelling results, where these

286    LP molecules were not well predicted using the [1]H NMR spectra. Monsonis-Centelles *et al.* [25] reported an average

287    within-individual coefficient of variation (WCV) as high as 12 to 16% for LDL2*tg* - LDL6*tg* using fresh blood

288    plasma samples duplicated from the same volunteer. In addition, their study also showed that the repeatability of

289    the UC based quantification of *tg* molecules is significantly less using frozen plasma compared to fresh plasma

290    samples. Given this, we speculate that the main reasons behind the relatively low prediction performances of the

291    PLS models for LDL2*tg*-LDL6*tg* are twofold: sample matrix disruption due to freeze-thaw cycles of the plasma

292    samples, and a relatively high uncertainty of the UC measurements. A similar trend was observed for *fchol*

293    molecules in the LDL subfractions [25]. Within-individual coefficient of variation values for LDL1*fchol*- LDL6*fchol*

294    vary by 12-24%  dependent on the two different types of assays. The PLS models developed in the present study

295    for *fchol* in LDL subfractions showed a $Q^2$ of 0.2-0.4 for test set samples. In contrast, *chol* and *phoslp* molecules

296    in LDL subfractions (LDL1-LDL6) were predicted with moderate to high prediction performance (test set $Q^2$ of

297    0.67-0.76 for *phoslp* and 0.70-0.81 for *chol*), with the exception of LDL6*chol*, LDL5*phoslp*, and LDL6*phoslp*,

298 which were not well predicted. Cholesterol esters (*chole*) were also predictable in all LDL subfractions with a

299 moderate prediction performance ($Q^2$ of 0.60-0.81), except for LDL3*chole* and LDL6*chole*.

300

**Validation of the final LP prediction models in an independent cohort**

302 In order to perform an external validation, PLS prediction models developed in this study were applied to an

303 independent Swedish cohort [22] using externally measured [1]H NMR spectra as input data. The predicted LP

304 concentrations were subsequently compared to the actual concentrations. The Swedish cohort included [1]H NMR

305 spectra of blood serum from 290 healthy subjects (sex: 210 females/80 males; age: $57.8 \pm 11$ years old; BMI: 25.0

306 $\pm 2.6$ kg/m$^2$) measured using a protocol similar to that of the present study. Unlike the present study, the LP

307 concentrations of the Swedish cohort were quantified using the HPLC method [27] at LipoSEARCH (Skylight

308 Biotech Inc., Akita, Japan). Lipoprotein subfractions were thus not directly comparable, as the HPLC method is

309 based on size distribution in contrast to the UC method, which separates LP particles based on their density. A

310 direct comparison was therefore only possible for total concentrations of *chol* and *tg* in blood. While the [1]H NMR

311 spectra of the Swedish cohort were acquired using the same [1]H NMR pulse sequence, temperature, and similar

312 acquisition parameters as in the Danish cohort, the blood sample preparation procedure for NMR measurements

313 differed slightly, resulting in noticeable differences in spectral intensities. Accordingly, prior to predictions, the

314 spectra of the Swedish cohort were aligned and scaled towards the Danish cohort as described previously [19]. Figure

315 4 shows scatter plots comparing actual concentrations (mg/dL) of total *chol* and *tg* as measured by HPLC, with

316 predicted concentrations from the corresponding PLS models developed using the Danish cohort in this study.

317 The total concentrations of *chol* and *tg* were predicted well with a Pearson's correlation coefficient ($r^2$) of 0.94

318 and 0.97 respectively. In the case of total *tg*, the root mean square error (RMSE) calculated between HPLC and

319 PLS based predicted values was as low as 8.2 mg/dL, and mean standard deviation (STD) was 4.6, and mean

320 relative standard deviation (RSD) was 5.1%. However, the PLS based predicted concentrations of total *chol* were

321 underestimated in the majority of samples. Despite a high correlation between the actual and predicted values, the

322 presence of the offset resulted in a relatively high RMSE (23.9). This may be related to systematic differences

323 between the two reference measurement methods for *chol* quantification. The UC method has been shown to result

324 in denaturation or degradation of some LP particles [28,29], which might be the reason for the systematic

325 underestimation of *chol* particles when using UC calibrated PLS models.

326 Despite the fact that subfractions are not directly comparable between HPLC and UC, correlations of *chol* and *tg*

327 values in the some subfractions were found to be high. For example, $r^2$ between the actual HPLC concentration

328 of *chol* in the G17 subfraction (which is defined as medium HDL with a diameter of 10.9 nm) and *chol* of HDL2a

329 subfraction predicted from UC calibrated PLS model was 0.83. As a consequence, RMSE and RSD values were

330 also low, 2.4 and 6.3 respectively, suggesting that G17*chol* quantified by HPLC may in fact represent HDL2a*chol*

331 measured by UC. However, concentrations of *tg* in the same subfractions, quantified by HPLC or predicted using

332 the PLS model calibrated by UC, were not comparable and resulted in a low $r^2$ (0.3). Concentrations of *tg* in the

333 G08 subfraction of the HPLC method, which represents *tg* in large LDL subfraction with the diameter of 28.6 nm,

334 correspond to the *tg* of IDL fraction quantified using UC and showed $r^2$, RMSE and RSD values of 0.72, 1.5, 8.4,

335 respectively. However, concentrations of *chol* in the same subfractions, G08 from HPLC and IDL from UC, were

336  not comparable ($r^2$ = 0.1 and RSD = 90%). Instead, the IDL*chol* content predicted by the PLS model showed a

337  relative high correlation with the *chol* of G06 subfraction (medium VLDL with a particle diameter of 36.8 nm)

338  measured by HPLC ($r^2$ = 0.64 and RSD = 18%). Interestingly, a similar correlation was observed for the *tg* content

339  in the same fractions, IDL*tg* versus *tg* in G06 ($r^2$ = 0.68 and RSD = 22%) (Figure S3). Furthermore, the PLS based

340  predicted *chol* concentration in the HDL2b subfraction was highly correlated to the *chol* of G16 subfraction ($r^2$ =

341  0.89) measured by HPLC, which represents a large HDL with a diameter of 12.1 nm. Although a high correlation

342  coefficient was observed between the UC based predicted and HPLC values, a significant offset was present (LP

343  concentrations were systematically overestimated by the PLS model or underestimated by the HPLC method),

344  and predictions were not accurate (RSD = 46%). Concentrations of *tg* in HDL2b particle predicted by PLS and in

345  the G16 subfraction of the HPLC were not comparable and resulted in an $r^2$ of 0.34. Whereas, *tg* concentrations

346  in G07 seem to represent *tg* in LDL1 subfraction measured by UC and showed relatively high correlation ($r^2$ =

347  0.65) and low RSD (15%).

348

### Rank of LP region of the $^1$H NMR spectra

350  In order to better understand the feasibility of predicting the concentrations of 65 LP from a relatively small $^1$H

351  NMR spectral region of human blood plasma, the LP region (1.4-0.6 ppm), the rank estimation of the NMR data

352  was performed. It is known that concentrations of many LP particles co-vary in blood and that biology makes it

353  extremely difficult to break this covariation. This phenomenon is known as *cage of covariance* [30]. In practice, this

354  means that some LP prediction models rely on information related to highly co-varying LP particles, and causal

355  relationships between the LP particles thus remain largely unknown. Nevertheless, a rank estimation was

356  performed in order to better understand the level of inter-correlations, i.e. the cage of covariance, between the

357  individual LP, and compare it with the possible information content in the NMR spectra. In order to estimate the

358  rank of the NMR spectra and UC data, a principal component analysis (PCA) based iterative approach [31,32] was

359  employed (SI). It is assumed that the rank estimated in this way approximates the true chemical variation present

360  in the data. This rank estimation revealed that the LP data used in this study (316 subjects by 65 LP variables) has

361  a rank of 33 (Figure 5), which indicates that there are at least 33 independent systematic variations present in the

362  LP data, and thus a medium level of cage of covariance. To further investigate this, the LP data was subjected to

363  a correlation analysis where each LP variable of the original UC data ($Y_{ACTUAL}$) was correlated to all other LP

364  variables individually, resulting in 65-by-65 diagonal matrix consisting of Pearson's correlation coefficients. The

365  same correlation analysis was then performed on the predicted UC data obtained from PLS models developed on

366  the Danish cohort in this study ($Y_{HAT}$) [30]. The symmetric heat map shows Pearson's correlation coefficients

367  between the 65 LP in $Y_{ACTUAL}$ and $Y_{HAT}$ (Figure S4), where LP variables are ordered using K-Nearest Neighbour

368  based clustering of LP in $Y_{ACTUAL}$. The heat map and the distribution of the correlation coefficients suggest that

369  an inter-correlation of LP variables is similar (symmetric) in $Y_{ACTUAL}$ and $Y_{HAT}$ and ranged between -0.55 till 0.99.

370  The fact that the correlations are not increased in the predicted concentrations ($Y_{HAT}$) compared to the actual

371  ($Y_{ACTUAL}$) is an indication of a reliable prediction network which is not over-fitted by inter-correlations between

372  LP. Clustering of the correlation matrices further revealed three major clusters of positive correlations representing

373  HDL, LDL and VLDL main and subfractions. Correlation between the same molecules but in different particles

374    is weak or insignificant. The most significant negative correlation trends were observed between VLDL and HDL
375    particles, and to a lesser extent between LDL and HDL particles.

376    In contrast to the UC data, the LP region of the $^1$H NMR spectra is much more complex. Despite the global
377    alignment performed, minor signal misalignments are still present across samples, which complicates an unbiased
378    rank estimation using the iterative PCA approach. Therefore, the rank of the LP region of the $^1$H NMR spectra
379    obtained from the Danish cohort is estimated without and after binning with different bin sizes of 2-32 (Figure 5).
380    The estimated rank of the LP region gradually decreases as bin size increases, with a notable decline after a bin
381    size of 11. Interestingly, the bin size of 11 was the largest bin size that could be used before losing spectral
382    resolution to the point that characteristic shape of the methyl protons between 0.95-0.8 ppm was kept intact, and
383    before the loss of spin coupling information from signals of amino acids. Thus bin size of 11 was found to be
384    optimal for reducing spectral misalignments while keeping the shape and resolution of signals intact. The rank of
385    the LP region without any binning was found to be as high as 93. After binning with bin size of 11, the rank of
386    the system reduced to 83. Further increase of the bin size caused rapid decline of the rank. For comparison, the
387    rank of the NMR spectra of the Swedish cohort [22], using bin size of 11, was 92. This could be explained by the
388    greater heterogeneity of the Swedish cohort compared to the Danish cohort. In conclusion, it is indeed possible
389    that the LP region of the NMR spectra is able to predict a large number of independently varying LP in human
390    blood plasma or serum.

391

**Spectral signatures of LP depend on particle size**

393    A few studies have characterised signature signals of the LP main fractions and subfractions. This has been done
394    mathematically using either curve fitting [16] or by calculating selectivity ratio from PLS models developed to
395    predict the LP [18,19], and experimentally by measuring pure fractions after UC [33]. The concentrations of the three
396    main populations of protons, cholesterol (0.75-0.6 ppm), methyl (0.95-0.80 ppm) and methylene (1.4-1.2 ppm),
397    differ significantly among the $^1$H NMR spectra of the four main fractions. The chemical shifts of these signals
398    also differ between the four main fractions, and it is mostly pronounced in the signals of the methyl and methylene
399    groups. The spectral differences between the subclasses belonging to the same main fraction are much smaller,
400    especially among the LDL subfractions. This section discusses the unique spectral signatures of LP main and
401    subfractions, within and between molecular types, using a selectivity ratio matrix obtained from the PLS models.
402    Selectivity ratio (SR) can be regarded as a spectral signature responsible for the prediction of a given LP particle.
403    The SR matrix consisting of all the 65 predictable LP was analysed using PCA and ANOVA-simultaneous
404    component analysis (ASCA) [34] (Figure 6). The first three principal components (PC) of the PCA model developed
405    on the mean centred SR explained 95% of the data variation, and clearly distinguished the SR of VLDL, IDL,
406    LDL, and HDL. By far the largest variation was captured by PC1 (68%), which fully separates LDL from HDL,
407    and IDL particles were in between them, while VLDL fractions slightly overlapped with the LDLs. However,
408    PC2 (20%) separates LDL from VLDL and IDL. Principal components analysis did not reveal any variation
409    related to LP molecular types. The ASCA analysis (partition of variations) of the mean centred SR data showed
410    that only particle type was significant and explained 59.5% of the total data variation, and that molecular type was
411    not significant. This was confirmed by superimposing plots of SR that revealed unique and consistent SR

412 specificity for the majority of fractions within a particle type. Figure 6c shows the mean SR of the four main

413 fractions across the molecular types (e.g. VLDL*mean* = mean of SR of VLDL*chol*, VLDL*tg*, VLDL*fchol* etc.)

414 and shows that the SRs are notably different, both in terms of relative ratio of the three main signals (cholesterol,

415 methyl and methylene) and in their chemical shift profiles. The SR of main fractions across different molecular

416 types were very similar, with $r^2$ ranging from 0.89 to 0.97 (Figure S5). The main difference between the SR of the

417 four main fractions can be quantified as relative ratios of the three major proton populations (Figure 6c). For HDL

418 and LDL the relative proportion of cholesterol methyl is significantly higher than in VLDL and IDL, whereas the

419 relative ratio of methylene protons is similarly higher in SR of VLDL and IDL compared to HDL and LDL. Such

420 relative proportions of the three major signals were consistent across all subfractions per particle type. When

421 comparing different particle types, but for the same molecular type, the relative ratios and the chemical shifts of

422 the signals were different (Figure S6a). However, the SR of subfractions belonging to the same main fraction and

423 the same molecular type were similar, though they possessed clear trends of chemical shift changes in cholesterol

424 methyl, LP methyl and methylene signals according to the particle size. For example, SR of the *chol* prediction

425 models for all five subfractions of LDL exhibited a clear shift of cholesterol methyl and LP methyl signals towards

426 a lower field from LDL1 to LDL5 subfractions (Figure S6b). The same trend was observed in SR of *phoslp* and

427 *apoB* models of LDL subfractions - a clear shift of the LP methyl signal towards lower field was observed for

428 both types of models. Similar shifts of LP methyl and methylene signals were observed for the *chol*, *phoslp*, and

429 *apoA* models of HDL subfractions. In all cases, signals shifted towards a lower field from HDL3 to HDL2a and

430 HDL2b subfractions.

431

**Implementation of lipoprotein prediction PLS models into the Signature Mapping (SigMa) software**

433 Based on the [1]H NMR spectra and LP concentrations obtained from UC of blood plasma samples from 316 Danes,

434 and subsequent PLS regression models, an open access software was developed for future use in research and

435 biomarker discovery. It was developed as an extension of the SigMa software, originally developed to process

436 complex [1]H NMR metabolomics datasets [24,35]. The software is based on the latest developments in the NMR

437 spectral data processing methods and through the regression models developed in this study, and able to predict

438 LP concentrations from human blood plasma using [1]H NMR NOESY spectra. SigMa takes the user [1]H NMR

439 spectra as an input and returns the predicted concentrations (mg/dL) of 65 LP main and subclasses. This requires

440 that the input spectra are compatible and acquired using a similar experimental protocol including blood sample

441 preparation and [1]H NMR spectral data acquisition. The SigMa software initializes by scaling the user spectra by

442 the ERETIC signal and aligning the spectra towards the doublet of alanine's methyl group at 1.49 ppm using the

443 *icoshift* [23] algorithm, as described previously [19]. Then the for LP predictions the input spectra are constrained to

444 the LP region by selecting spectral range of 1.4-0.6 ppm. Prior to LP prediction SigMa ensures compatibility of

445 the user spectra by spectral length correction and intensity normalization. Then 65 LP main and subclasses are

446 predicted including a "traffic light" marker that keep track of the individual LP predictions quality. This quality

447 check is based on an "*X-Y relation*" test and a "*Y-predicted*" test, which evaluates if the predicted LP

448 concentrations are within the cohort calibration range. The "traffic light" marker classifies each predicted LP

449 value as green (both input spectrum and predicted LP values are within the cohort calibration model), yellow (at

450 least one parameter, either input spectrum or predicted LP concentration, is outside the cohort calibration model),

12

451 or red (both input spectrum and predicted LP values are outside the cohort calibration model). More details of the

452 SigMa LP prediction method are explained in Supplementary Information. SigMa LP software can be freely

453 downloaded from www.food.ku.dk/foodomics.

454

455 **Discussion**

456 The feasibility of LP prediction in human blood using $^1$H NMR spectra and PLS based regression was

457 demonstrated some time ago, but prediction performances, parameters and complexity of PLS models, as well as

458 their transferability to new cohorts, remains unclear. The present study describes, for the first time, the entire

459 workflow of LP prediction using $^1$H NMR spectra, including spectral processing steps, the optimization and

460 validation of PLS regression models, and their prediction performances on test set and independent cohort

461 samples. Coherent datasets of $^1$H NMR spectra of human blood plasma and corresponding UC data acquired on

462 300+ volunteers in the Danish cohort were made publicly available for future research, with an aim to improving

463 PLS models for LP prediction. $^1$H NMR spectra were comprehensively investigated to find the best, simplest, and

464 most robust spectral signatures able to predict concentration of LP. A total of 13 distinct spectral intervals

465 representing LP signals and other lipid species were identified and tested in different combinations for their

466 performance to develop PLS LP prediction models. We found that a relatively small interval of the $^1$H NMR

467 spectra, namely the LP region (1.4-0.6 ppm), was the optimal region for LP prediction in terms of model

468 performance, robustness and simplicity. The LP region contained not only the three most important proton

469 populations representing LP signals, including methyl groups of cholesterol (0.75-0.65 ppm), and methyl (0.92-

470 0.8 ppm) and methylene (1.4-1.2 ppm) protons of different LP molecules, but also signals of a few amino acids

471 (e.g., valine, leucine and isoleucine), and of lactic acid. Spectral regions representing only the methyl (0.92-0.8

472 ppm) or the methylene (1.4-1.2 ppm) protons performed significantly worse, especially for LDL and HDL

473 particles, than their counterparts, where the two proton populations were combined. This suggests that the entire

474 LP profile information is better preserved in the LP region, which includes all major signature signals of LP

475 molecules., This region was therefore the final NMR dataset used to develop the PLS models.

476 Using the LP region, all main fractions in plasma were predicted with high accuracy. Consistently for all LP

477 molecular types, the prediction performance of PLS models was best for total concentrations in blood plasma

478 followed by main fractions, whilst relatively lower performances were observed for subfractions. The complexity

479 of PLS models depends on LP particle type, the smaller the particle size the more LVs were needed to develop an

480 optimal model. Overall, 4 to 17 LVs were required to predict the different LP variables, and the least number of

481 LVs were required for models of LP molecules in plasma (4-7 LV) and in main fractions (4-8 LV). A greater

482 number of LVs were required for subfractions (4-11 LV). The greatest number of LVs was needed for LDL*apoB*

483 subfractions (8-17 LV). It can be assumed that the optimal number of LVs depends not only on the complexity of

484 the signals, but also on the cohort, number of subjects, heterogeneity of the volunteers, as well as on the LP

485 concentration span and spectral interferences. An average difference in RMSE of the training and test set models

486 was <5%, which indicates robustness of the PLS models developed in this study. The models were further

487 validated externally on an independent Swedish cohort (290 volunteers) and showed high accuracy in prediction

488 of the LP variables that were comparable between the two cohorts, plasma concentrations of *tg* and *chol*. The PLS

13

489     models developed in this study were implemented in the SigMa LP software, which is freely available. In the

490     Danish cohort we were able to predict concentrations of 65 of the 97 measured LP variables using UC. Sub-

491     optimal PLS models for the remaining 32 LP can be explained by relatively high uncertainty of the reference

492     measurements and/or limitation of the variability in the cohort. The PLS models implemented in SigMa LP can

493     be improved/upgraded in the future when a new datasets with coherent $^1$H NMR spectra and UC data become

494     available from other cohorts. This will significantly increase the coverage of the LP prediction models.

495     Analysis of selectivity ratios (SR) mostly shows characteristic spectral pattern for LP particle types, and to a less

496     extent reflects particle size or molecular type. However, chemical shifts of a few signals in SR were dependent on

497     the particle size of subfractions. This was most pronounced for LDL subfractions (Figure S5). The SR developed

498     in this study can be used for PLS model comparison across laboratories to validate the spectral signatures of LP

499     from $^1$H NMR spectra. Rank estimation performed on the LP region of the NMR spectra suggest the rank of 83

500     for the Danish cohort. This is a surprisingly high number considering the relatively small NMR spectral interval,

501     which only contains signals of a handful of blood metabolites and three major proton populations representing

502     methyl group of cholesterol, and methyl and methylene protons of other LP molecules. Interestingly, a similarly

503     high rank of 92 was observed for the LP region of the spectra from the Swedish cohort. These high *"chemical"*

504     ranks reflect the complexity of the LP region due to the small but distinct spectral signatures of the LP main and

505     subfractions. As observed from the SR, the same LP molecules possessed significantly different spectral

506     signatures. Even more strikingly the same LP molecules in the same LP particle's subfraction had significantly

507     different SR spectral profiles (Figure S6a,b). The prediction of as many as 65 LP variables from a relatively small

508     NMR spectral interval is therefore feasible, and the PLS models developed do not appear to be assisted by the co-

509     variation of LP alone.

510     In conclusion, this study describes a protocol and open access data to build PLS models to predict LP concentration

511     from standard $^1$H NMR spectra acquired on human blood plasma or serum using the most advanced/recent SOPs

512     applied in all NMR phenotyping laboratories around the world. The models are optimized, use the most

513     informative and reproducible region of the spectra, and are based on a relatively large and heterogeneous cohort.

514     Most importantly, it is possible to enrich and maintain the calibration models when new datasets from different

515     laboratories become available.

516     **Materials and Methods**

517     The study was approved by the Research Ethics Committees of the Capital Region of Denmark in accordance

518     with the Helsinki declaration (H-15008313) and the Danish Data Protection Agency (2013-54-0522). The Danish

519     cohort included 316 subjects recruited from the COUNTERSTRIKE cohort. Subjects included 206 females (51.1

520     ± 19.8 years old) and 110 males (57.4 ± 19.7 years old) with the mean body mass index (BMI (kg/m$^2$)) of 24.9 (±

521     4.4) for females and 25.3 (± 3.6) for males. The mean values for systolic and diastolic blood pressure (mm Hg)

522     were 124.4 (± 14.5), and 76.7 (± 9.4), respectively, for females, and 130.6 (± 16.8) and 77.5 (± 11.0), respectively,

523     for males. All subjects were apparently healthy and without diagnosis of any form of cardiovascular disease or

524     diabetes, reporting no chronic gastrointestinal disorders, and not receiving antibiotic treatment within three

525     months of starting the study, or using pre- or probiotic supplements within one month of starting the study. All

526     subjects visited The Department of Nutrition, Exercise and Sports, where blood samples were taken, and

527    anthropometric and clinical parameters were recorded. Fasting blood samples were collected in vacutainers

528    containing an anticoagulant reagent ethylenediaminetetraacetic acid (EDTA). Plasma was separated after blood

529    sample collection and stored in 500 ul aliquot cryovials at −80 °C until measurement.

530    Details of the UC based quantification of LP particles, measurement of One-dimensional (1D) proton ([1]H) NMR

531    spectra on human blood plasma samples, the chemical and reagents used, and details of the spectral data

532    processing and PLS model development are given in Supplementary Information. Briefly, quantification of LP

533    particles was performed using UC as described previously [25]. One-dimensional [1]H NMR spectra were measured

534    on fasting stage EDTA plasma samples as described previously [19] at the Department of Food Science (University

535    of Copenhagen) using a Bruker Avance III 600 MHz NMR spectrometer equipped with a 5-mm broadband inverse

536    RT (BBI) probe, automated tuning, and matching accessory (ATMA) and cooling unit BCU-05. The spectrometer

537    was equipped with an automated sample changer (SampleJet, Bruker BioSpin) with sample cooling (278 K) and

538    preheating stations (298 K), where samples were stored at 278 K and measured within 72h. Phase and baseline

539    corrected 1D [1]H NMR spectra were then imported to the SigMa software [24], scaled to the ERETIC signal [36], and

540    aligned towards alanine's doublet corresponding to its methyl group (1.507−1.494 ppm) using *icoshift* [23]. Prior to

541    PLS regression analysis subjects with a LP concentration below limit of detection (LOD) or with missing values

542    were removed. This led to slightly different numbers of subjects for different PLS models. An optimal number of

543    LVs was selected by fitting one to twenty LV models to the training samples (70%) using 10 fold cross validation

544    and 10 times Monte-Carlo repetitions. Thus, a total of 640,200 PLS models (33 NMR datasets × 97 LP variables

545    × 20 LVs × 10 cross validations) were developed in this study. It should be noted that a few subjects (zero to

546    seven) whose LP values were predicted with a large error were regarded as "X-Y relation" outliers and were

547    removed from the training models. A PLS calibration model was then recalculated and tested on independent

548    subjects (30%) that were not used in the training model optimization. The final PLS calibration models were also

549    tested to predict LP variables, plasma *tg* and *chol*, in the independent Swedish cohort [22]. All data analysis including

550    PLS, PCA, and ASCA, were performed in MATLAB (version R2016b, The Mathworks, Inc., U.S.A.) using

551    customised MATLAB scripts written by the authors.

552    **Data Availability**

553    The [1]H NMR data acquired on the human blood plasma of 316 healthy subjects from the Danish cohort and the

554    corresponding lipoprotein concentration data, measured by ultracentrifugation, are available upon request by

555    contacting the corresponding authors. Signature Mapping for Lipoprotein Quantification (SigMa LP) software

556    can be freely downloaded from www.food.ku.dk/foodomics. All other data supporting the findings of this study

557    are included in the article text and supporting information.

558    **ACKNOWLEDGEMENTS**

**References**

1    Roth, G. A. *et al.* Global, Regional, and National Burden of Cardiovascular Diseases for 10 Causes, 1990 to 2015. *J Am Coll Cardiol* **70**, 1-25, doi:10.1016/j.jacc.2017.04.052 (2017).

2    *World Health Organization (WHO): Cardiovascular diseases*, https://www.who.int/health-topics/cardiovascular-diseases/#tab=tab_1

3    Ference, B. A. *et al.* Low-density lipoproteins cause atherosclerotic cardiovascular disease. 1. Evidence from genetic, epidemiologic, and clinical studies. A consensus statement from the European Atherosclerosis Society Consensus Panel. *Eur Heart J* **38**, 2459-2472, doi:10.1093/eurheartj/ehx144 (2017).

4    D'Agostino, R. B. *et al.* General Cardiovascular Risk Profile for Use in Primary Care. *Circulation* **117**, 743-753, doi:10.1161/CIRCULATIONAHA.107.699579 (2008).

5    Wilson, P. W. F. *et al.* Prediction of Coronary Heart Disease Using Risk Factor Categories. *Circulation* **97**, 1837-1847, doi:10.1161/01.CIR.97.18.1837 (1998).

6    Rosenson, R. S. *et al.* HDL and atherosclerotic cardiovascular disease: genetic insights into complex biology. *Nat Rev Cardiol* **15**, 9-19, doi:10.1038/nrcardio.2017.115 (2018).

7    Rader, D. J. & Hovingh, G. K. HDL and cardiovascular disease. *Lancet* **384**, 618-625, doi:10.1016/s0140-6736(14)61217-4 (2014).

8    Chapman, M. J. & Caslake, M. Non-high-density lipoprotein cholesterol as a risk factor: addressing risk associated with apolipoprotein B-containing lipoproteins. *European Heart Journal Supplements* **6**, A43-A48, doi:10.1016/j.ehjsup.2004.01.010 (2004).

9    Fernandez, M. L. & Webb, D. The LDL to HDL cholesterol ratio as a valuable tool to evaluate coronary heart disease risk. *J Am Coll Nutr* **27**, 1-5, doi:10.1080/07315724.2008.10719668 (2008).

10   Aru, V. *et al.* Quantification of lipoprotein profiles by nuclear magnetic resonance spectroscopy and multivariate data analysis. *TrAC Trends in Analytical Chemistry* **94**, 210-219, doi:10.1016/j.trac.2017.07.009 (2017).

11   Musliner, T. A. & Krauss, R. M. Lipoprotein subspecies and risk of coronary disease. *Clin Chem* **34**, B78-83 (1988).

12   Otvos, J. D., Jeyarajah, E. J., Bennett, D. W. & Krauss, R. M. Development of a Proton Nuclear Magnetic Resonance Spectroscopic Method for Determining Plasma Lipoprotein Concentrations and Subspecies Distributions from a Single, Rapid Measurement. *Clinical Chemistry* **38**, 1632-1638, doi:10.1093/clinchem/38.9.1632 (1992).

13   Otvos, J. D., Jeyarajah, E. J. & Bennett, D. W. Quantification of plasma lipoproteins by proton nuclear magnetic resonance spectroscopy. *Clinical Chemistry* **37**, 377-386 (1991).

14   van den Boogaart, A., Ala-Korpela, M., Jokisaari, J. & Griffiths, J. R. Time and frequency domain analysis of NMR data compared: an application to 1D 1H spectra of lipoproteins. *Magn Reson Med* **31**, 347-358, doi:10.1002/mrm.1910310402 (1994).

15   Ala-Korpela, M. 1H NMR spectroscopy of human blood plasma. *Progress in Nuclear Magnetic Resonance Spectroscopy* **27**, 475-554, doi:10.1016/0079-6565(95)01013-0 (1995).

16   Ala-Korpela, M. *et al.* 1H NMR-based absolute quantitation of human lipoproteins and their lipid contents directly from plasma. *J Lipid Res* **35**, 2292-2304 (1994).

17   Wold, S., Ruhe, A., Wold, H. & W. J. Dunn, I. The Collinearity Problem in Linear Regression. The Partial Least Squares (PLS) Approach to Generalized Inverses. *SIAM Journal on Scientific and Statistical Computing* **5**, 735-743, doi:10.1137/0905052 (1984).

18   Petersen, M. *et al.* Quantification of Lipoprotein Subclasses by Proton Nuclear Magnetic Resonance–Based Partial Least-Squares Regression Models. *Clinical Chemistry* **51**, 1457-1461, doi:10.1373/clinchem.2004.046748 (2005).

19   Monsonis Centelles, S. *et al.* Toward Reliable Lipoprotein Particle Predictions from NMR Spectra of Human Blood: An Interlaboratory Ring Test. *Analytical Chemistry* **89**, 8004-8012, doi:10.1021/acs.analchem.7b01329 (2017).

20   Dyrby, M. *et al.* Analysis of lipoproteins using 2D diffusion-edited NMR spectroscopy and multi-way chemometrics. *Analytica Chimica Acta* **531**, 209-216, doi:10.1016/j.aca.2004.10.052 (2005).

21   Dona, A. C. *et al.* Precision High-Throughput Proton NMR Spectroscopy of Human Urine, Serum, and Plasma for Large-Scale Metabolic Phenotyping. *Analytical Chemistry* **86**, 9887-9894, doi:10.1021/ac5025039 (2014).

22   Mihaleva, V. V. *et al.* A Systematic Approach to Obtain Validated Partial Least Square Models for Predicting Lipoprotein Subclasses from Serum NMR Spectra. *Analytical Chemistry* **86**, 543-550, doi:10.1021/ac402571z (2014).

23   Savorani, F., Tomasi, G. & Engelsen, S. B. icoshift: A versatile tool for the rapid alignment of 1D NMR spectra. *Journal of Magnetic Resonance* **202**, 190-202, doi:10.1016/j.jmr.2009.11.012 (2010).

24   Khakimov, B., Mobaraki, N., Trimigno, A., Aru, V. & Engelsen, S. B. Signature Mapping (SigMa): An efficient approach for processing complex human urine [1]H NMR metabolomics data. *Analytica Chimica Acta* **1108**, 142-151, doi:10.1016/j.aca.2020.02.025 (2020).

25   Monsonis-Centelles, S., Hoefsloot, H. C. J., Engelsen, S. B., Smilde, A. K. & Lind, M. V. Repeatability and reproducibility of lipoprotein particle profile measurements in plasma samples by ultracentrifugation. *Clinical Chemistry and Laboratory Medicine (CCLM)* **58**, 103, doi:10.1515/cclm-2019-0729 (2020).

26   Wold, S., Sjöström, M. & Eriksson, L. PLS-regression: a basic tool of chemometrics. *Chemometrics and Intelligent Laboratory Systems* **58**, 109-130, doi:10.1016/S0169-7439(01)00155-1 (2001).

27   Okazaki, M. *et al.* Identification of unique lipoprotein subclasses for visceral obesity by component analysis of cholesterol profile in high-performance liquid chromatography. *Arterioscler Thromb Vasc Biol* **25**, 578-584, doi:10.1161/01.Atv.0000155017.60171.88 (2005).

28   Kunitake, S. T. & Kane, J. P. Factors affecting the integrity of high density lipoproteins in the ultracentrifuge. *J Lipid Res* **23**, 936-940 (1982).

29  Murdoch, S. J. & Breckenridge, W. C. Development of a Density Gradient Ultracentrifugation Technique for the Resolution of Plasma Lipoproteins which Avoids Apo E Dissociation. *Analytical Biochemistry* **222**, 427-434, doi:10.1006/abio.1994.1512 (1994).

30  Berhe, D. T. *et al.* Prediction of total fatty acid parameters and individual fatty acids in pork backfat using Raman spectroscopy and chemometrics: Understanding the cage of covariance between highly correlated fat parameters. *Meat Science* **111**, 18-26, doi:10.1016/j.meatsci.2015.08.009 (2016).

31  Vitale, R. *et al.* Selecting the number of factors in principal component analysis by permutation testing—Numerical and practical aspects. *Journal of Chemometrics* **31**, e2937, doi:10.1002/cem.2937 (2017).

32  Hotelling, H. Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology* **24**, 417-441, doi:10.1037/h0071325 (1933).

33  Baumstark, D. *et al.* (1)H NMR spectroscopy quantifies visibility of lipoproteins, subclasses, and lipids at varied temperatures and pressures. *J Lipid Res* **60**, 1516-1534, doi:10.1194/jlr.M092643 (2019).

34  Smilde, A. K. *et al.* ANOVA-simultaneous component analysis (ASCA): a new tool for analyzing designed metabolomics data. *Bioinformatics* **21**, 3043-3048, doi:10.1093/bioinformatics/bti476 (2005).

35  Cui, M., Trimigno, A., Aru, V., Khakimov, B. & Engelsen, S. B. Human Faecal 1H NMR Metabolomics: Evaluation of Solvent and Sample Processing on Coverage and Reproducibility of Signature Metabolites. *Analytical Chemistry* **92**, 9546-9555, doi:10.1021/acs.analchem.0c00606 (2020).

36  Akoka, S., Barantin, L. & Trierweiler, M. Concentration measurement by proton NMR using the ERETIC method. *Analytical Chemistry* **71**, 2554-2557, doi:10.1021/ac981422i (1999).

**Figure captions**

**Figure 1.** An overview of the implemented workflow to predict concentrations of lipoproteins in human blood plasma using 1D $^1$H NMR spectra. An outer loop indicated with a grey line represents selection of NMR spectra regions used for PLS modelling. A total of 33 NMR datasets were constructed using 13 spectral regions, representing LP signals and signals of other lipid species, either alone or in different combinations. An inner lop indicated with a black line represents PLS modelling of each lipoprotein particle individually using a selected NMR spectral region. * **m** corresponds to a number of spectral data points in NMR data, **k** corresponds to a number of lipoprotein variables, **n** corresponds to a number of subjects recruited in cohorts.

**Figure 2.** Lipoprotein prediction performance of PLS models developed on 20 NMR spectral regions, of 33 investigated, that showed relatively high prediction performances for at least 65 of 97 modelled lipoproteins. **a)** $Q^2$ obtained from test set prediction of 65 LP using PLS models developed on 20 NMR spectral regions. Black bars on the left side of the heat map show overall prediction performance of each LP variable (normalized cumulative value of $Q^2$ obtained from 20 PLS models). Black bars at the bottom of the heat map show overall prediction performance of each NMR spectral regions used for PLS modelling (normalized cumulative value of $Q^2$ obtained from 65 LP variables). The NMR spectral region shown inside the dashed line corresponds to the LP region (1.4-0.6 ppm) with one of the highest performances for predicting concentrations of LP. For more details see Table S1. **b)** 20 NMR spectral regions corresponding to the heat map on the left. *Regions of spectra highlighted in grey were excluded from the PLS modelling.

**Figure 3.** Training PLS model and test set prediction performances of selected LP variables, included triglycerides (*tg*), cholesterol (*chol*), phospholipid (*phoslp*), and apolipoprotein A (*apoA*) molecules in different fraction or sub-fraction of LP particles using the LP region (1.4-0.6 ppm) of the 1D $^1$H NMR spectra and ultracentrifugation as a reference method.

**Figure 4.** Validation of the PLS based LP prediction models developed in this study in an independent cohort, the Swedish cohort (290 healthy subjects).

**Figure 5.** Rank estimation of the LP region of the 1D $^1$H NMR spectra and ultracentrifugation data obtained from the Danish cohort (316 subjects) was performed separately using an iterative permutation based PCA modelling. Rank of the LP region of 1D $^1$H NMR spectra of the Swedish cohort (290 healthy subjects) was evaluated in the same way. Correlation coefficients ($R^2$), relative standard deviation (RSD), standard deviation (STD), and root mean square error (RMSE) were calculated between the predicted value of LP variable in this study using UC and the HPLC based measured value obtained from the Swedish cohort. RMSE-M, $Q^2$-M, and CV-M correspond to the RMSE of test set prediction (RMSEP), $Q^2$ and coefficients of variance, respectively, obtained from the test set prediction of an LP variable obtained in this study.

**Figure 6.** Comparison of selectivity ratios (SR) obtained from the PLS models developed for 65 LP variables. **a)** Score plots of PCA model developed on mean centred selectivity ratio data. **b)** Loadings of the corresponding PCA model. **c)** Mean of 1D $^1$H NMR spectra and mean selectivity ratios of four main fractions of LP across all molecular types. *Main difference in SR is between particles, VLDL, IDL, LDL, and HDL. ASCA analysis of the SR data show that only particle type was significant (p-val = 1.2e-10, Exp.Var. = 59.5%), while LP molecular type or the two-factor interaction terms were not significant.
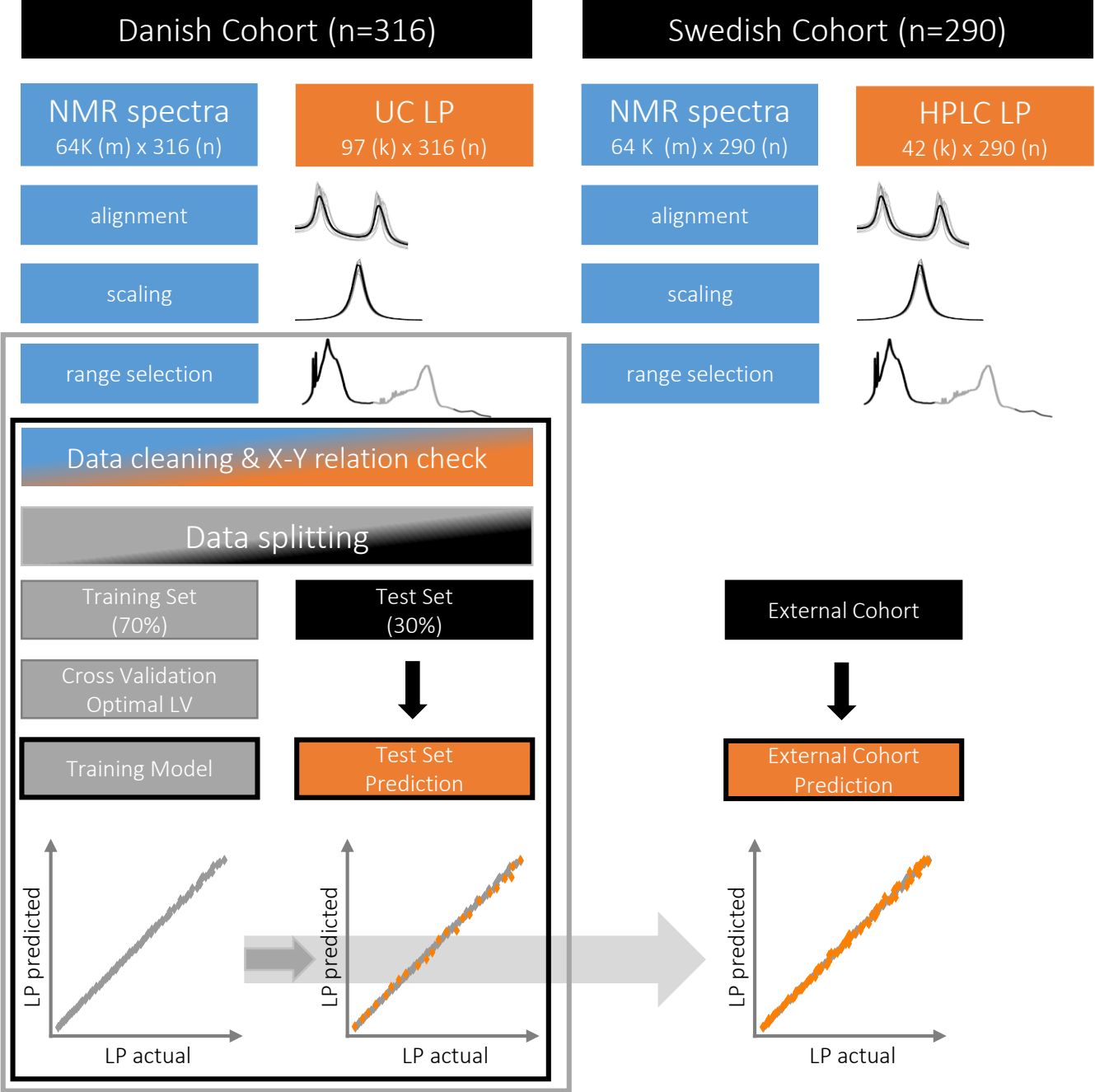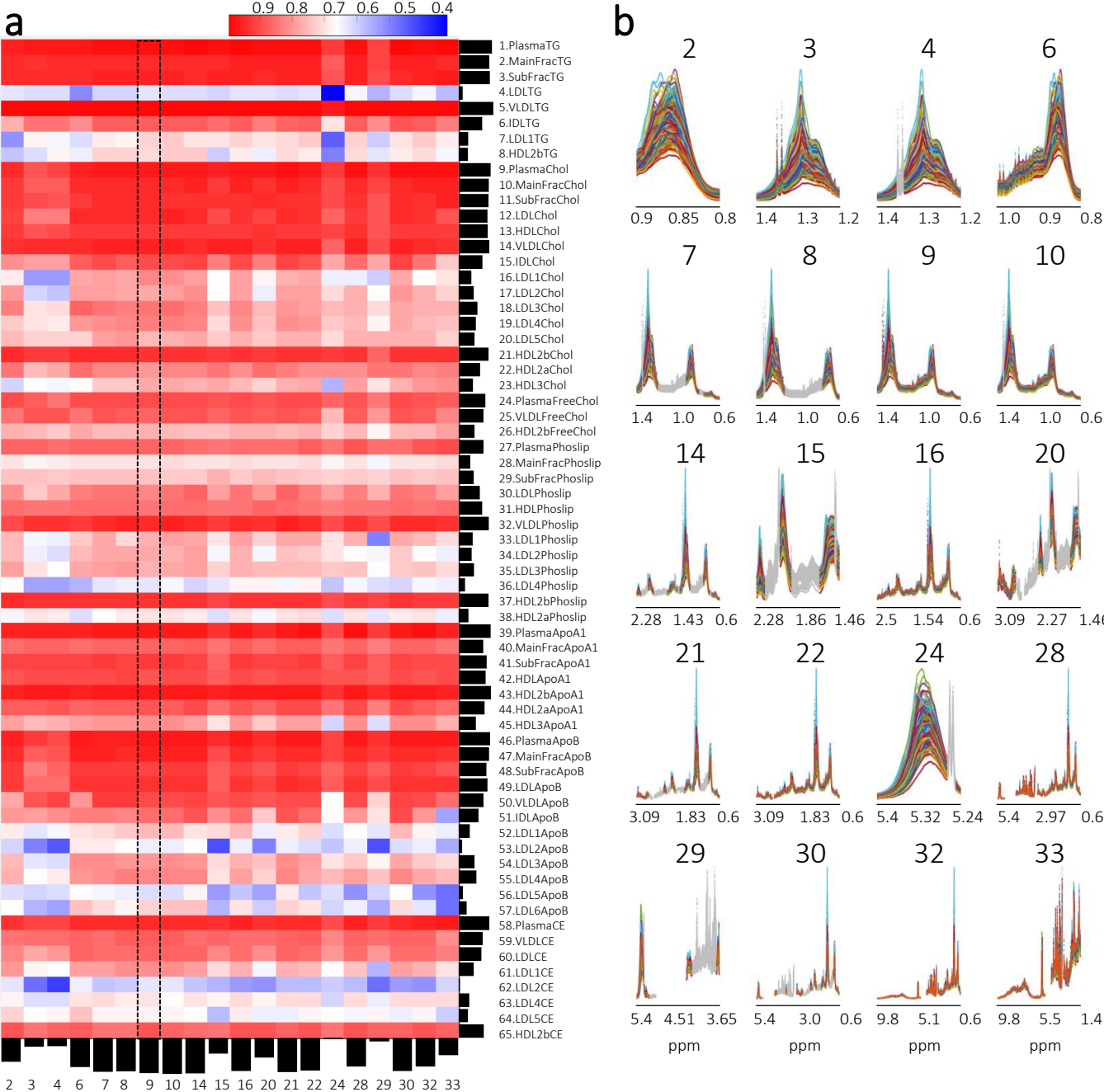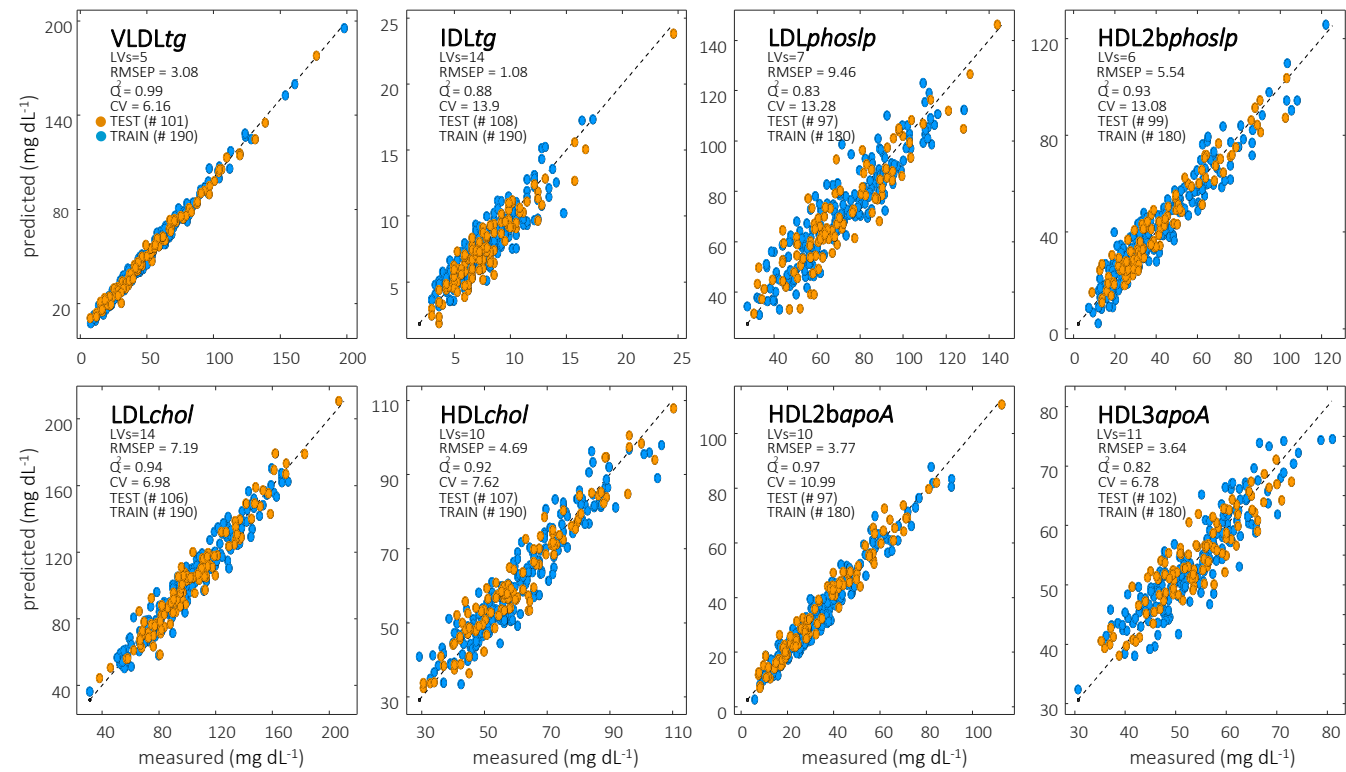
# Figure 1



**Danish Cohort (n=316)**

| NMR spectra 64K (m) x 316 (n) | UC LP 97 (k) x 316 (n) |

alignment

scaling

range selection

Data cleaning & X-Y relation check

Data splitting

| Training Set (70%) | Test Set (30%) |

Cross Validation Optimal LV

| Training Model | Test Set Prediction |

LP predicted / LP actual

LP predicted / LP actual

**Swedish Cohort (n=290)**

| NMR spectra 64 K (m) x 290 (n) | HPLC LP 42 (k) x 290 (n) |

alignment

scaling

range selection

External Cohort

External Cohort Prediction

LP predicted / LP actual

**Figure 2**



a

0.9 0.8 0.7 0.6 0.5 0.4

1.PlasmaTG
2.MainFracTG
3.SubFracTG
4.LDLTG
5.VLDLTG
6.IDLTG
7.LDL1TG
8.HDL2bTG
9.PlasmaChol
10.MainFracChol
11.SubFracChol
12.LDLChol
13.HDLChol
14.VLDLChol
15.IDLChol
16.LDL1Chol
17.LDL2Chol
18.LDL3Chol
19.LDL4Chol
20.LDL5Chol
21.HDL2bChol
22.HDL2aChol
23.HDL3Chol
24.PlasmaFreeChol
25.VLDLFreeChol
26.HDL2bFreeChol
27.PlasmaPhoslip
28.MainFracPhoslip
29.SubFracPhoslip
30.LDLPhoslip
31.HDLPhoslip
32.VLDLPhoslip
33.LDL1Phoslip
34.LDL2Phoslip
35.LDL3Phoslip
36.LDL4Phoslip
37.HDL2bPhoslip
38.HDL2aPhoslip
39.PlasmaApoA1
40.MainFracApoA1
41.SubFracApoA1
42.HDLApoA1
43.HDL2bApoA1
44.HDL2aApoA1
45.HDL3ApoA1
46.PlasmaApoB
47.MainFracApoB
48.SubFracApoB
49.LDLApoB
50.VLDLApoB
51.IDLApoB
52.LDL1ApoB
53.LDL2ApoB
54.LDL3ApoB
55.LDL4ApoB
56.LDL5ApoB
57.LDL6ApoB
58.PlasmaCE
59.VLDLCE
60.LDLCE
61.LDL1CE
62.LDL2CE
63.LDL4CE
64.LDL5CE
65.HDL2bCE

2 3 4 6 7 8 9 10 14 15 16 20 21 22 24 28 29 30 32 33

b

**Figure 3**



VLDL*tg*
LVs=5
RMSEP = 3.08
$q^2$ = 0.99
CV = 6.16
● TEST (# 101)
● TRAIN (# 190)

IDL*tg*
LVs=14
RMSEP = 1.08
$q^2$ = 0.88
CV = 13.9
TEST (# 108)
TRAIN (# 190)

LDL*phoslp*
LVs=7
RMSEP = 9.46
$q^2$ = 0.83
CV = 13.28
TEST (# 97)
TRAIN (# 180)

HDL2b*phoslp*
LVs=6
RMSEP = 5.54
$q^2$ = 0.93
CV = 13.08
TEST (# 99)
TRAIN (# 180)

LDL*chol*
LVs=14
RMSEP = 7.19
$q^2$ = 0.94
CV = 6.98
TEST (# 106)
TRAIN (# 190)

HDL*chol*
LVs=10
RMSEP = 4.69
$q^2$ = 0.92
CV = 7.62
TEST (# 107)
TRAIN (# 190)

HDL2b*apoA*
LVs=10
RMSEP = 3.77
$q^2$ = 0.97
CV = 10.99
TEST (# 97)
TRAIN (# 180)

HDL3*apoA*
LVs=11
RMSEP = 3.64
$q^2$ = 0.82
CV = 6.78
TEST (# 102)
TRAIN (# 180)

predicted (mg dL$^{-1}$)

measured (mg dL$^{-1}$)

# Figure 4

# Figure 5

# Figure 6

# Supplementary Information

# Human blood lipoprotein predictions from $^1$H NMR spectra: protocol, model performances and cage of covariance

Bekzod Khakimov[a,1], Huub C.J. Hoefsloot[b], Nabiollah Mobaraki[c], Violetta Aru[a], Mette Kristensen[d], Mads V. Lind[d], Lars Holm[e], Josué L Castro-Mejía[a], Dennis S Nielsen[a], Doris M. Jacobs[f], Age K. Smilde[a,b], Søren Balling Engelsen[a,2]


[a]Department of Food Science, University of Copenhagen, Rolighedsvej 26, DK-1958 Frederiksberg C, Denmark

[b]Swammerdam Institute for Life Sciences, University of Amsterdam, Postbus 94215, Amsterdam 1090 GE, The Netherlands

[c]Department of Chemistry, Faculty of Science, Shiraz University, Shiraz, 7194684795, Iran

[d]Department of Nutrition, Exercise and Sports, University of Copenhagen, Rolighedsvej 26, DK-1958 Frederiksberg C, Denmark

[e] School of Sport, Exercise and Rehabilitation Sciences, University of Birmingham, Edgbaston, Birmingham, B15 2TT, UK

[f]Unilever Global Food Innovation Centre, 6708 WH Wageningen, The Netherlands


[1]Corresponding author: Bekzod Khakimov, Department of Food Science, University of Copenhagen, Rolighedsvej 26, DK-1958 Frederiksberg C, Denmark. Tel.; +45-2887-4454, Email: bzo@food.ku.dk, ORCID: 0000-0002-6580-2034

[2]Corresponding author: Søren Balling Engelsen, Department of Food Science, University of Copenhagen, Rolighedsvej 26, DK-1958 Frederiksberg C, Denmark. Tel.: +45-2020-0064, Email: se@food.ku.dk, ORCID: 0000-0003-4124-4338

## Materials and Methods

### *Chemicals and reagents*

Unless otherwise stated, all chemicals and reagents were purchased from Sigma-Aldrich (Søborg, Denmark). These include deuterium oxide (D2O, 99.9 atom % D), monobasic sodium phosphate ($NaH_2PO_4$, ≥ 99.0%), and dibasic sodium phosphate ($Na_2HPO_4$, ≥ 98.0%), sodium salt of 3-(trimethylsilyl) propionic-2,2,3,3-d4 acid (TSP, 98 atom % D, ≥ 98.0%), and sodium azide (NaN3, ≥ 99.5%). The water used throughout the study was purified using a Millipore lab water system (Merck KGaA, Darmstadt, Germany) equipped with a 0.22 μm filter membrane. For the stock solutions used during ultracentrifugation NaCl (VWR Chemicals, US), NaN3 (Riedel-de Haën, Germany), EDTA (Merck, Germany), and NaBr (Alfa Aesar, US) were used.

### *Ultracentrifugation based quantification of lipoprotein particles*

Quantification of lipoprotein (LP) particles was performed using ultracentrifugation (UC) as previously described (1). Seven different lipoprotein molecules including cholesterol (*chol*), triglycerides (*tg*), cholesterol ester (*chole*), free cholesterol (*fchol*), phospholipids (*phosl*), apolipoprotein A (*apoA*) and apolipoprotein B (*apoB*) were quantified in all or in some of the LP main fractions (VLDL, IDL, HDL, LDL) and/or in their subfractions (HDL2a, HDL2b, HDL3, LDL1, LDL2, LDL3, LDL4, LDL5, LDL6). Fractionation was done by sequential centrifugation of 3 mL EDTA plasma using an Optima L-80 XP ultracentrifuge with a fixed angle rotor type 50.4 Ti (Beckman Coulter, Inc., US). The UC process was initiated immediately after the fasting plasma sample was collected, and it was completed over a period of 8 days. A detailed description of the separation steps can be found in (1). Immediately after the fractionation step, all subfractions were frozen and stored at -80⁰C for later analysis. Colorimetric and turbidimetric assays were performed on an ABX Pentra 400 analyzer (ABX Pentra; Horiba ABX, Montpellier, France) to determine the plasma, main class and subclass concentrations of total *chol*, *tg*, *apoA* and *apoB* (ABX Pentra; Horiba Medical, France). Free cholesterol and *phoslp* were determined using colorimetric and turbidimetric assays (MTI Diagnostics, Germany).

### *Measurement of ¹H NMR spectra on human blood plasma*

Fasting EDTA-plasma samples were measured using one dimensional (1D) nuclear Overhauser effect spectroscopy (NOESY) proton (¹H) NMR spectra as previously described (2). Briefly, 350 µL of plasma was carefully mixed with the same volume of phosphate buffer into a 2.0 mL Eppendorf tube and 600 µL of the mixture was transferred into SampleJet NMR tube (103.5 mm length and 5.0 mm diameter). The phosphate buffer was prepared as previously described (3). Sample preparation and measurements were randomized. Pooled control human blood plasma samples were measured at regular intervals throughout the whole measurement sequence. The ¹H NMR spectra of blood plasma samples were acquired at the Department of Food Science (University of Copenhagen) using a Bruker Avance III 600 MHz NMR spectrometer equipped with a 5-mm broadband inverse RT (BBI) probe, automated tuning and matching accessory (ATMA) and cooling unit BCU-05. The spectrometer was equipped with an automated sample changer (SampleJet, Bruker BioSpin, Rheinstetten, Germany) with sample cooling (278 K) and preheating stations (298 K) where samples were stored at 278 K and measured within

72h. Data acquisition and processing were carried out using TOPSPIN 3.5 PL6 (Bruker BioSpin, Rheinstetten, Germany) and automation of the overall measurement procedure was controlled by Icon NMR (Bruker BioSpin, Rheinstetten, Germany). Each sample was pre-heated at 298 K for 60 sec in SampleJet and kept inside the NMR probe head for 5 minutes to reach temperature equilibrium at 310 ±0.1 K. Before each measurement automated tuning and matching, automated locking, and automated shimming (TOPSHIM routine) were performed. Automation included also the 90° hard pulse calibration, and optimized presaturation power for each sample. The $^{1}$H NMR spectra were acquired using the standard pulse sequence with water suppression (*noesygppr1d*) from the Bruker pulse program library. A total of 32 scans were acquired after 4 dummy scans, and the generated free induction decays (FIDs) were collected into 96k data points using a spectral width of 30 ppm. The relaxation delay and mixing time were set to 4.0 and 0.01 sec, respectively. The receiver gain was set to 90.5 for all samples. Automated data processing, including Fourier transform of FID (free induction decay), apodization with a 0.3 Hz line-broadening, automated phasing, and baseline correction was carried out, for each $^{1}$H NMR spectrum, in the TOPSPIN software.

### *Data analysis*

The $^{1}$H NMR spectra were imported into the SigMa software (4) and scaled towads the Electronic REference To access *In vivo* Concentrations (ERETIC) signal (5) positioned at 15 ppm, which is equivalent to 10 mmol L$^{-1}$ protons. The scaled spectra were then aligned towards the doublet of alanine's methyl group (1.507−1.494 ppm) using *icoshift* (6). Subsequently, the spectra were divided into 13 different regions representing LP signals. Each region alone or in various combinations were used for partial least squares (PLS) regression analysis (7) making a total of 33 different NMR datasets with different lengths (Figure S1). Prior to the PLS analysis, subjects with LP concentrations below the limit of detection (LOD) or with missing values were removed from the datasets leading to a small difference in the number of subjects included in the PLS models. Each NMR dataset was separately used to predict 97 LP variables obtained from UC. NMR spectra and LP data were mean centred prior to PLS. First a training model was developed using 70% of randomly selected subjects (e.g., 210 subjects out of 300). An optimal number of latent variables (LVs) was selected by fitting one to twenty LV models to the training samples using a 10-fold cross validation and 10 times Monte-Carlo repetitions. After a PLS calibration model was developed and optimized it was tested on 30% subjects (independent set) that were not used in training model optimization. In addition, the final PLS calibration models were also tested to predict LP variables, plasma *tg* and *chol*, of the independent Swedish cohort (8). In order to estimate the rank of the NMR and UC data, a principal component analysis (PCA) (9) based iterative approach (10) was employed. This method fits one component at a time by deflation of the original matrix by the corresponding modelled data and the procedure is continued on the resulting residual matrix. In parallel, the same iterative PCA procedure is repeated with the residual matrix after its columns are independently permuted. The method then compares the "so-called" F-ratio, which is a ratio of an eigenvalue obtained from the PCA analysis of the original matrix (or unpermuted residual matrix) to the value obtained from the PCA of the permuted residual matrix. After deflating a certain number of principal components (PCs), the F-ratio of the unpermuted residual matrix will become equal to or lower than the F-ratio obtained on the permuted residual matrix which in turn is a sign that there is no more sys

tematic variation left in the residuals. All data analysis including PLS, PCA, and ANOVA-simultaneous component analysis (ASCA) (11) were performed in MATLAB (version R2016b, The Mathworks, Inc., U.S.A.) using customised MATLAB scripts written by the authors.

**Captions of Supporting Information**

**Figure S1.** A total of 33 regions of the $^1$H NMR spectra were used to develop the PLS models for predicting concentrations of lipoproteins in human blood plasma. These regions represented 13 unique NMR spectral intervals corresponding to protons derived from different lipoprotein molecules and other lipids: **Region 1** - represented the methyl group (C18) of cholesterol (δ 0.75-0.65); **Region 2** – methyl group of lipoprotein molecules (δ 0.92-0.8); **Region 3** – methylene group of lipoprotein molecules (δ 1.4-1.2); **Region 11** - the methylene groups of lipids (C*H₂*CH2CO) located two bonds away from carbonyl group (δ 1.48-1.46, 1.65-1.51); **Region 12** - the methylene groups of lipids (C*H₂*C=C) located one bond away from double bond (δ 2.04-1.94); **Region 13** - the methylene groups of lipids (C*H₂*CO) located one bond away from carbonyl group (δ 2.28-2.2); **Region 17** - the methylene groups of lipids (C=CHC*H₂*CH=C) located one bond away from two double bonds (δ 2.84-2.74); **Region 18** – signals derived from lysine residue in albumin (δ 2.84-3.09); **Region 23** – the methylene group (C*H₂*OCOR) of glyceryl of lipids and the methylene group (C*H*2OH) of choline (δ 4.35-4.24); **Region 24** – the methine protons (C*H*) of unsaturated lipids (δ 5.4-5.24); **Region 25** – the methine group (C*H*OCOR) of glyceryl of lipids (δ 5.23-5.14); **Region 26** – the methylene group (NC*H₂*) of choline (δ 3.71-3.65); **Region 27** – the methylene groups of lipids (CH₂C*H₂*C=C) located one bond away from double bond (δ 1.88-1.65). Chemical shifts range of the remaining 20 regions were as follows: **Region 4**, δ 1.34-1.2, 1.4-1.36; **Region 5**, δ 1.07-0.92; **Region 6**, δ 1.07-0.8; **Region 7**, δ 0.75-0.6, 0.92-0.8, 1.34-1.2, 1.4-1.36; **Region 8**, δ 0.92-0.8, 1.34-1.2, 1.4-1.36; **Region 9**, δ 1.4-0.6; **Region 10**, δ 1.34-0.6, 1.4-1.36; **Region 14**, δ 0.75-0.6, 0.92-0.8, 1.34-1.2, 1.4-1.36, 1.48-1.46, 1.65-1.51, 2.04-1.94, 2.28-2.2; **Region 15**, δ 1.48-1.46, 1.65-1.51, 2.04-1.94, 2.28-2.2; **Region 16**, δ 2.5-0.6; **Region 19**, δ 3.09-2.74; **Region 20**, δ 1.48-1.46, 1.65-1.51, 2.04-1.94, 2.28-2.2, 3.09-2.74; **Region 21**, δ 0.75-0.6, 0.92-0.8, 1.34-1.2, 1.4-1.36, 1.48-1.46, 1.65-1.51, 2.04-1.94, 2.28-2.2, 3.09-2.74; **Region 22**, δ 2.56-0.6, 2.7-2.62, 3.09-2.73; **Region 28**, δ 2.56-0.6, 2.7-2.62, 3.09-2.73, 3.59-3.25, 4.35-3.65, 5.4-5.0, **Region 29**, δ 3.71-3.65, 4.35-4.24, 5.23-5.14, 5.25-5.24, 5.4-5.26; **Region 30**, δ 0.75-0.6, 0.92-0.8, 1.34-1.2, 1.4-1.36, 1.48-1.46, 1.65-1.51, 2.04-1.94, 2.28-2.2, 3.09-2.74, 3.71-3.65, 4.35-4.24, 5.23-5.14, 5.25-5.24, 5.4-5.26; **Region 31**, δ 9.8-5.45; **Region 32**, δ 2.56-0.6, 2.7-2.62, 3.09-2.73, 3.59-3.25, 4.35-3.65, 9.8-5.0; **Region 33**, δ 2.56-1.4, 2.7-2.62, 3.09-2.73, 3.59-3.25, 4.35-3.65, 9.8-5.0. *Part of the spectra deemed in grey colour corresponds to the excluded range of the spectra from PLS modelling.

**Figure S2A.** Lipoprotein prediction performance of test set PLS models developed on 20 NMR spectral regions (see Figure 2), of 33 investigated, that showed relatively high prediction performances for at least 65 of 97 modelled lipoproteins. **Q²** is a statistical measure of prediction accuracy often used in PLS modelling (12) and defined as Q²=(1-PRESS/SS), where PRESS is predictive residual sum of squares and SS is sum of squares of actual values (LP concentrations). **RMSEP**, root mean square error of prediction, is a statistical measure of absolute prediction error estimated from test set PLS models. **CV**, coefficient of variation, is a statistical measure of relative (%) prediction error. Higher the Q², and lower the CV and RMSE values indicate high prediction performance of PLS models. *Four best performing regions are highlighted in colour: LP (Region 9) in blue, -CH₃ (Region 2) in yellow, -CH₂- (Region 3) in green, and –CH₃ and –CH₂- (Region 8) in magenta.

**Figure S2B.** Lipoprotein prediction performance of training set PLS models developed on 20 NMR spectral regions (see Figure 2), of 33 investigated, that showed relatively high prediction performances for at least 65 of 97 modelled lipoproteins. **Q²** is a statistical measure of prediction accuracy often used in PLS modelling (12) and defined as Q²=(1-PRESS/SS), where PRESS is predictive residual sum of squares and SS is sum of squares of actual values (LP concentrations). **RMSECV**, root mean square error of cross validation, is a statistical measure of absolute prediction error estimated from test set PLS models. **CV**, coefficient of variation, is a statistical measure of relative (%) prediction error. Higher the Q², and lower the CV and RMSE values indicate high

prediction performance of PLS models. *Four best performing regions are highlighted in colour: LP (Region 9) in blue, -CH₃ (Region 2) in yellow, -CH₂- (Region 3) in green, and –CH₃ and –CH₂- (Region 8) in magenta.

**Figure S3.** Validation of the PLS based LP prediction models developed in this study in an independent cohort, the Swedish cohort (290 healthy subjects) (see Figure 4).

**Figure S4.** Correlations between lipoprotein particles. Heat map demonstrates clustered Pearson correlation coefficients calculated between lipoprotein concentrations, separately for LP measured using ultracentrifugation ($Y_{ACTUAL}$) and predicted ($Y_{HAT}$) using the PLS models developed in this study. A distribution plot demonstrates an overview of positive and negative correlations between LP in $Y_{ACTUAL}$ and $Y_{HAT}$.

**Figure S5.** Selectivity ratios (SR) calculated from PLS models. SR of different LP molecular are compared within the same main class, including all four main classes, VLDL, IDL, LDL, and HDL.

**Figure S6A.** SR of the same LP molecular are compared across main classes, VLDL, IDL, LDL, and HDL.

**Figure S6B.** SR of the same LP molecular are compared within the same main and subfractions.

**Table S1.** Lipoprotein prediction performance of PLS models developed on 20 NMR spectral regions that showed relatively high prediction performances for 65 lipoproteins. $Q^2$ (prediction accuracy), $R^2$ (Pearson correlation coefficient between actual and predicted values), RMSE (root mean square error of cross validation (for training set) or prediction (for test set)), CV (coefficients of variation), and P (p-value) values are given separately for training and test set prediction models.

**Table S2.** Lipoprotein prediction performance of PLS models developed using the LP region (δ 1.4-0.6) of the ${}^1$H NMR spectra and LP measured using ultracentrifugation on the Danish cohort subjects. $Q^2$ (prediction accuracy), $R^2$ (Pearson correlation coefficient between actual and predicted values), RMSE (root mean square error of cross validation (for training set) or prediction (for test set)), CV (coefficients of variation), and P (p-value) values are given separately for training and test set prediction models. The table also contains the number of subjects included in training and test set PLS models, as well as mean, median, min, maximum, standard deviation, relative standard deviation, and quartile 0.25, 0.5, and 0.75 are given.

# References

1. S. Monsonis-Centelles, H. C. J. Hoefsloot, S. B. Engelsen, A. K. Smilde, M. V. Lind, Repeatability and reproducibility of lipoprotein particle profile measurements in plasma samples by ultracentrifugation. *Clinical Chemistry and Laboratory Medicine (CCLM)* **58**, 103 (2020).

2. S. Monsonis Centelles *et al.*, Toward Reliable Lipoprotein Particle Predictions from NMR Spectra of Human Blood: An Interlaboratory Ring Test. *Analytical Chemistry* **89**, 8004-8012 (2017).

3. O. Beckonert *et al.*, Metabolic profiling, metabolomic and metabonomic procedures for NMR spectroscopy of urine, plasma, serum and tissue extracts. *Nature Protocols* **2**, 2692-2703 (2007).

4. B. Khakimov, N. Mobaraki, A. Trimigno, V. Aru, S. B. Engelsen, Signature Mapping (SigMa): An efficient approach for processing complex human urine $^1$H NMR metabolomics data. *Analytica Chimica Acta* **1108**, 142-151 (2020).

5. S. Akoka, L. Barantin, M. Trierweiler, Concentration measurement by proton NMR using the ERETIC method. *Analytical Chemistry* **71**, 2554-2557 (1999).

6. F. Savorani, G. Tomasi, S. B. Engelsen, icoshift: A versatile tool for the rapid alignment of 1D NMR spectra. *Journal of Magnetic Resonance* **202**, 190-202 (2010).

7. S. Wold, A. Ruhe, H. Wold, I. W. J. Dunn, The Collinearity Problem in Linear Regression. The Partial Least Squares (PLS) Approach to Generalized Inverses. *SIAM Journal on Scientific and Statistical Computing* **5**, 735-743 (1984).

8. V. V. Mihaleva *et al.*, A Systematic Approach to Obtain Validated Partial Least Square Models for Predicting Lipoprotein Subclasses from Serum NMR Spectra. *Analytical Chemistry* **86**, 543-550 (2014).

9. H. Hotelling, Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology* **24**, 417-441 (1933).

10. R. Vitale *et al.*, Selecting the number of factors in principal component analysis by permutation testing—Numerical and practical aspects. *Journal of Chemometrics* **31**, e2937 (2017).

11. A. K. Smilde *et al.*, ANOVA-simultaneous component analysis (ASCA): a new tool for analyzing designed metabolomics data. *Bioinformatics* **21**, 3043-3048 (2005).

12. S. Wold, M. Sjöström, L. Eriksson, PLS-regression: a basic tool of chemometrics. *Chemometrics and Intelligent Laboratory Systems* **58**, 109-130 (2001).
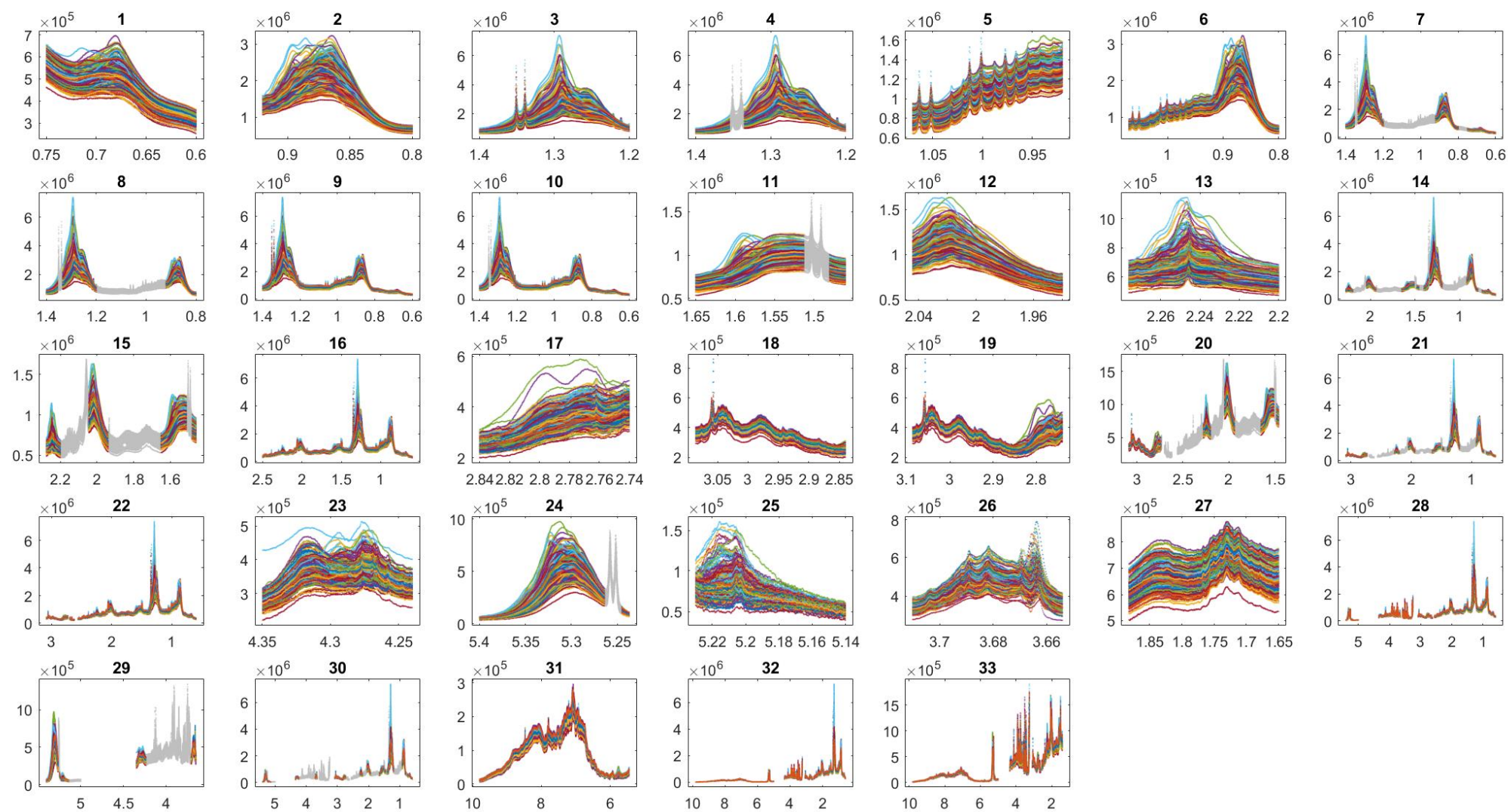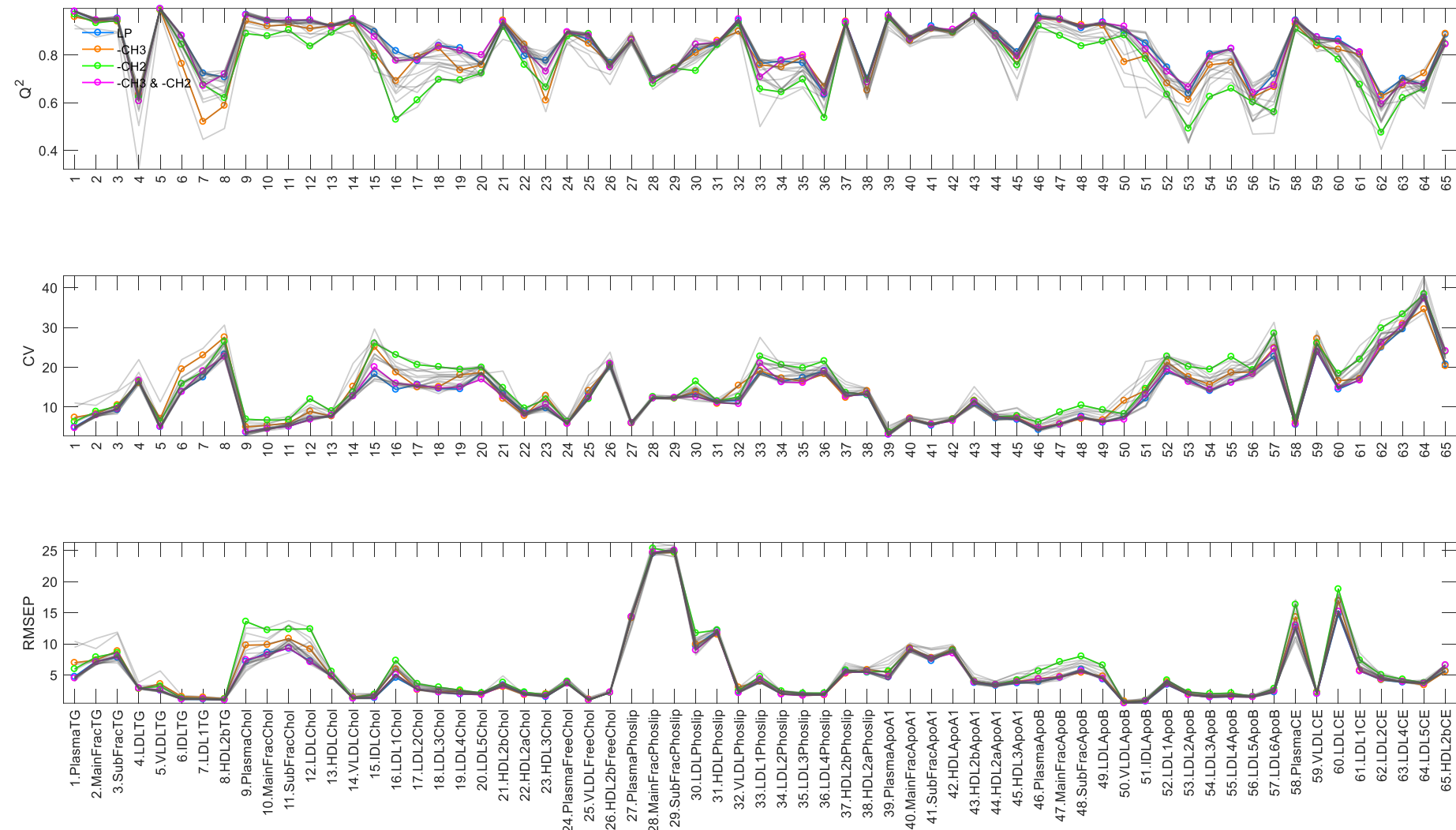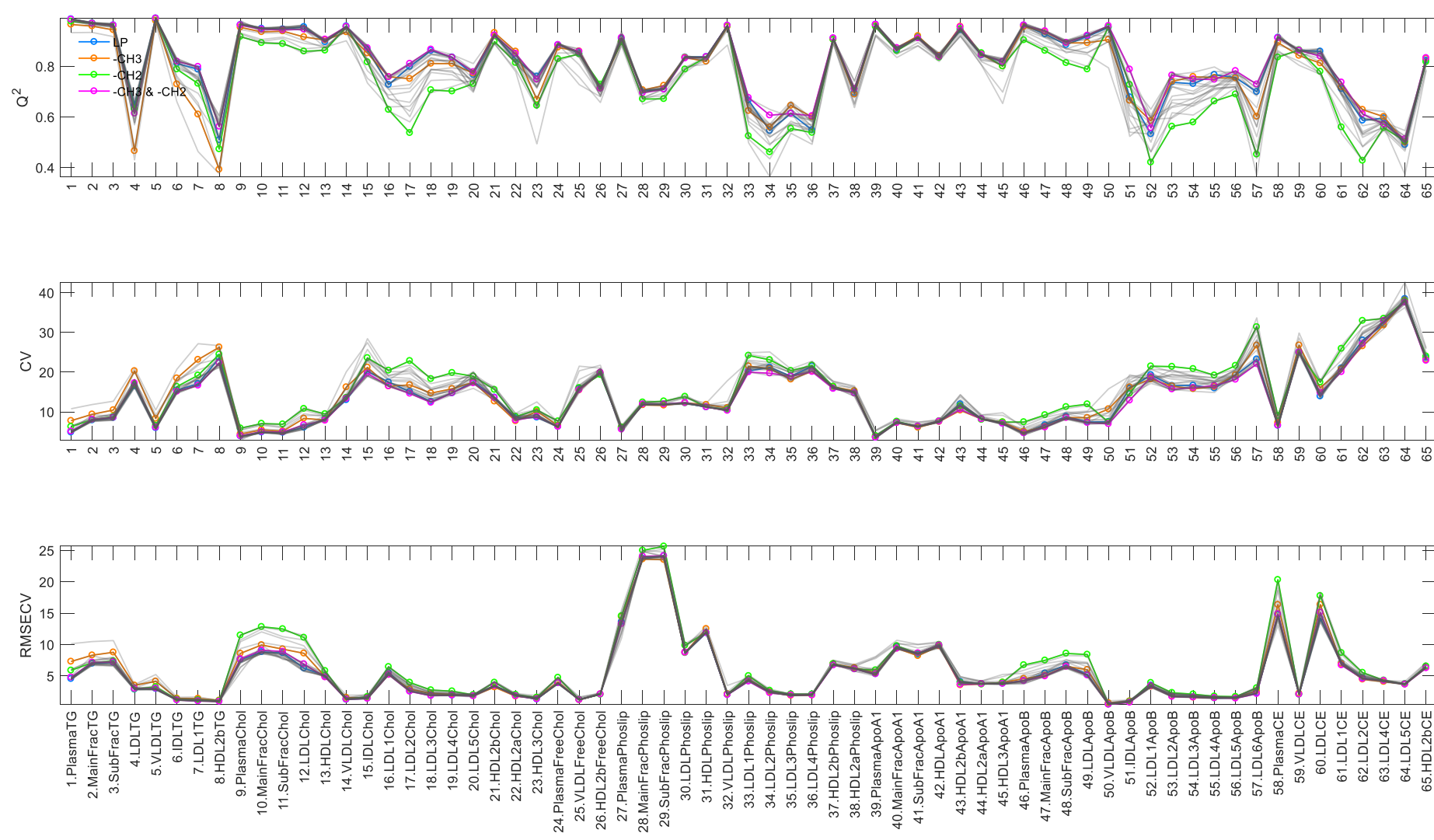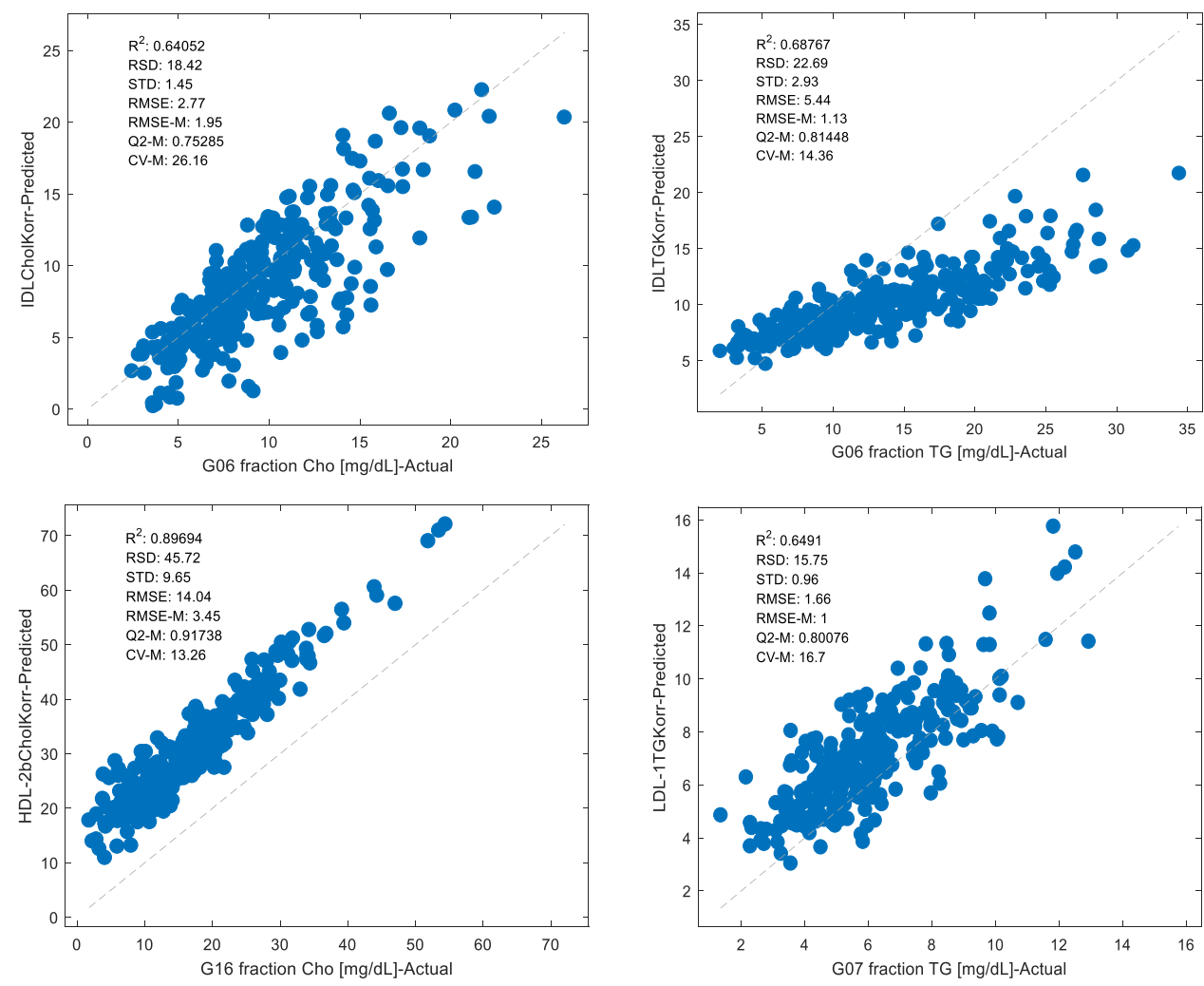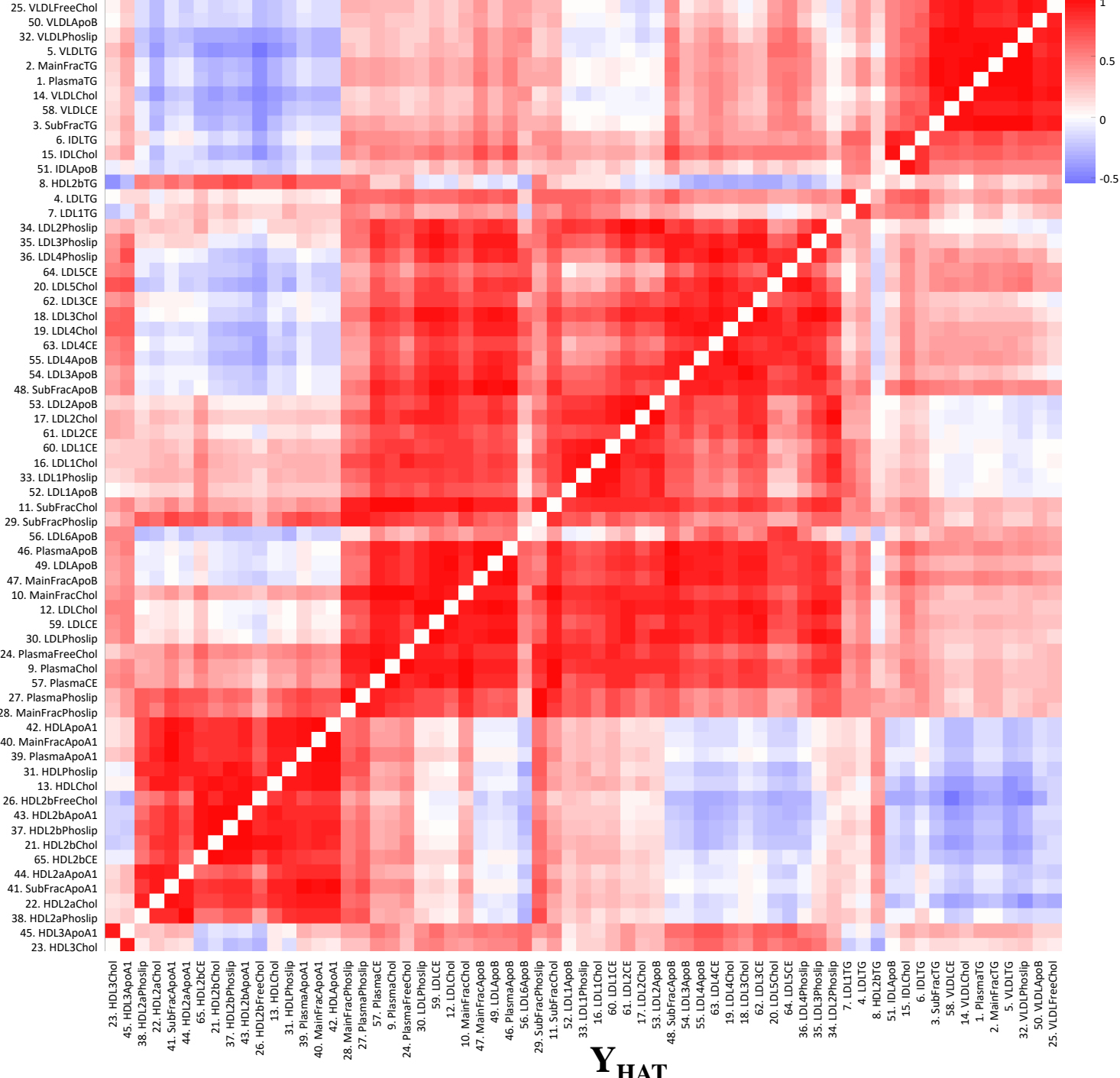
# Figure S1

# Figure S2A

# Figure S2B

# Figure S3

**Figure S4**

# Figure S5

# Figure S6A

# Figure S6B