

1 **Multi-omic approach to identify phenotypic modifiers underlying** 2 **cerebral demyelination in X-linked adrenoleukodystrophy**

3 **Phillip A. Richmond**^{†1}, **Frans van der Kloet**^{†2,3}, **Frederic M. Vaz**⁴, **David Lin**¹, **Anuli Uzozie**^{5,6},
4 **Emma Graham**¹, **Michael Kobor**¹, **Sara Mostafavi**¹, **Perry D. Moerland**², **Philipp F. Lange**^{5,6},
5 **Antoine H. C. van Kampen**^{2,7}, **Wyeth Wasserman**^{#1}, **Marc Engelen**^{#*8,9}, **Stephan Kemp**^{#*7,8},
6 **Clara van Karnebeek**^{#*8}

7 ¹Center for Molecular Medicine and Therapeutics, BC Children's Hospital, University of British
8 Columbia, Vancouver, BC, Canada

9 ²Bioinformatics Laboratory, Department of Clinical Epidemiology, Biostatistics and Bioinformatics,
10 Amsterdam Public Health research institute, Amsterdam University Medical Centers, University of
11 Amsterdam, Amsterdam, The Netherlands

12 ³Department of Pediatrics, Amsterdam University Medical Centers, Amsterdam, The Netherlands

13 ⁴Laboratory Genetic Metabolic Diseases, Amsterdam UMC, University of Amsterdam, Department
14 of Clinical Chemistry, Amsterdam Gastroenterology & Metabolism, Amsterdam, The Netherlands

15 ⁵Department of Pathology and Laboratory Medicine, University of British Columbia, BC, Canada

16 ⁶Michael Cuccione Childhood Cancer Research Program, BC Children's Hospital, Vancouver, BC,
17 Canada

18 ⁷Biosystems Data Analysis, Swammerdam Institute for Life Sciences, University of Amsterdam,
19 Amsterdam, The Netherlands

20 ⁸Department of Pediatric Neurology, Emma Children's Hospital, Amsterdam University Medical
21 Centers, Amsterdam, The Netherlands

22 ⁹Department of Neurology, Amsterdam University Medical Centers, Amsterdam, The Netherlands

23 ¹⁰Department of Pediatrics, Emma Children's Hospital, Amsterdam University Medical Centres,
24 Amsterdam NL

25 ¹¹Department of Pediatrics, Amalia Children's Hospital, Radboud University Medical Center,
26 Nijmegen NL

27

28 **† Co-first authors, # Co-last authors, * Corresponding authors**

29

30 *** Correspondence:**

31 Marc Engelen <m.engelen@amsterdamumc.nl>, Stephan Kemp <s.kemp@amsterdamumc.nl>
32 >, Clara van Karnebeek <c.d.vankarnebeek@amsterdamumc.nl>

33

34 **Keywords: Multi-omics, LCMS, Methylation, Lipidomics, Proteomics, Bioinformatics, X-ALD,**
35 **cerebral demyelination**

36

37 Abstract

38 X-linked adrenoleukodystrophy (ALD) is a peroxisomal metabolic disorder with a highly
39 complex clinical presentation. ALD is caused by mutations in the *ABCD1* gene, and is characterized
40 by the accumulation of very long-chain fatty acids in plasma and tissues. Disease-causing mutations
41 are ‘loss of function’ mutations, with no prognostic value with respect to the clinical outcome of an
42 individual. All male patients with ALD develop spinal cord disease and a peripheral neuropathy in
43 adulthood, although age of onset is highly variable. However, the lifetime prevalence to develop
44 progressive white matter lesions, termed cerebral ALD (CALD), is only about 60%. Early
45 identification of transition to CALD is critical since it can be halted by allogeneic hematopoietic stem
46 cell therapy only in an early stage. The primary goal of this study is to identify molecular markers
47 which may be prognostic of cerebral demyelination from a simple blood sample, with the hope that
48 blood-based assays can replace the current protocols for diagnosis. We collected six well-
49 characterized brother pairs affected by ALD and discordant for the presence of CALD and performed
50 multi-omic profiling of blood samples including genome, epigenome, transcriptome,
51 metabolome/lipidome, and proteome profiling. In our analysis we identify discordant genomic alleles
52 present across all families as well as differentially abundant molecular features across the omics
53 technologies. The analysis was focused on univariate modeling to discriminate the two phenotypic
54 groups, but was unable to identify statistically significant candidate molecular markers. Our study
55 highlights the issues caused by a large amount of inter-individual variation, and supports the
56 emerging hypothesis that cerebral demyelination is a complex mix of environmental factors and/or
57 heterogeneous genomic alleles. We confirm previous observations about the role of immune
58 response, specifically auto-immunity and the potential role of PFN1 protein overabundance in CALD
59 in a subset of the families. We envision our methodology as well as dataset has utility to the field for
60 reproducing previous or enabling future modifier investigations.

61
62

63 1 Introduction

64 Adrenoleukodystrophy (ALD) is a rare peroxisomal X-linked degenerative disease (MIM
65 300100), caused by deficiency of the ABC half-transporter encoded by the *ABCD1* gene. Over 800
66 different disease-causing loss-of-function *ABCD1* mutations have been reported
67 (www.adrenoleukodystrophy.info). Mutations lead to a defect in the import of very long-chain fatty
68 acids (VLCFA) into peroxisomes for further degradation and a subsequent accumulation of VLCFA
69 in plasma and tissues. The overall incidence is 1:17,000. In males, ALD often manifests with
70 adrenocortical insufficiency in childhood (50% before 10 years)(Huffnagel, Laheji, et al. 2019).
71 During adulthood virtually all male and, eventually, female patients develop a progressive
72 myelopathy termed adrenomyeloneuropathy (Engelen et al. 2014, 2012). Additionally, during
73 childhood or sometimes through adulthood male patients can develop cerebral demyelination, termed
74 cerebral ALD (CALD). It is estimated that eventually more than 60% of male patients develop
75 CALD(de Beer and Scheltens 2016)(Kemp, Berger, and Aubourg 2012). Untreated CALD is often
76 progressive, but can spontaneously arrest in 10 - 20% of patients. It causes vegetative state and death
77 2-3 years after onset, so early identification as well as careful and frequent monitoring of all male
78 ALD patients is necessary. If diagnosed early, hematopoietic stem cell therapy can be used to halt
79 further progression of cerebral ALD. To ensure timely stem cell therapy for males with CALD,
80 affected individuals are subjected to rigorous neurological and MRI follow-ups that pose
81 considerable physical, emotional and financial burden. As such, the unresolved and unpredictable
82 phenotypic variability of ALD is a crucial roadblock for patient care.

83 As newborn screening for ALD has recently been implemented, there is an urgent need for
84 identification of markers which may be prognostic of cerebral demyelination in many newly
85 diagnosed patients around the world. Our research focuses on delineating the enormous phenotypic
86 variability in ALD, with the overarching goal of identifying biomarkers prognostic of the
87 advancement to CALD. If successful in identifying biomarkers with prognostic power, then the
88 biomarkers could replace existing expensive monitoring protocols and potentially highlight
89 therapeutic targets, as is the case with other rare genetic disorders. For example, in Spinal Muscular
90 Atrophy, the genes *PLS3* and *CORO1C* were identified as protective modifiers, unravelling impaired
91 endocytosis as a rescue mechanism for the phenotype (Hosseinibarkooie et al. 2016). These
92 modifiers were identified from studies focusing on siblings with discordant disease severity, and are
93 opening novel therapeutic targets for treatment. Patients with ALD may benefit from similar research
94 advances.

95 Phenotypic discordance in individuals with the same *ABCD1* genotype, including siblings and
96 even monozygotic twins (Korenke et al. 1996), strongly supports the hypothesis that other modifying
97 factors play a role in the progression of the disease (Kemp et al. 2016; Wiesinger, Eichler, and Berger
98 2015). As yet, however, modifier studies using candidate gene approaches have had little success and
99 resulted in the identification of only a single modifier gene (*CYP4F2*) with limited prognostic power
100 (van Engen et al. 2016). Other candidate variants have been proposed, including a candidate cis-
101 regulatory SNP in the promoter region of *ELOVLI*—a gene involved in VLCFA synthesis (Ofman et
102 al. 2010). The functional consequences of this SNP with respect to the expression of *ELOVLI* in the
103 brain is still under investigation (Kemp, Berger, and Aubourg 2012). The lack of modifier
104 identification could be due to the limited genomic search space that was explored, which to date has
105 focused only on candidate gene approaches. Owing to the small sample size inherent to rare disease
106 cohorts, traditional genome wide association studies (GWAS) approaches are not feasible.
107 Employing a strategy which utilizes family structure may allow for a narrower search space
108 compared to GWAS, while allowing a broader interrogation of the genome than candidate gene
109 approaches. Beyond a search space which involves genetic mapping, other high throughput “omics”
110 technologies allow the exploration of complex biological systems at many levels. It is now possible
111 to identify differences between individuals or phenotypic states at the DNA, methylated DNA, RNA,
112 lipid, and protein levels. Our goal is to delineate personal molecular characteristics that contribute to
113 phenotypic variability in male ALD siblings enabling the identification of biomarkers that
114 prognosticate onset and progression of CALD. Because modifying factors (Génin, Feingold, and
115 Clerget-Darpoux 2008) could also include environmental, epigenetic and microbiome factors (Génin,
116 Feingold, and Clerget-Darpoux 2008; Argmann et al. 2016), multi-omics approaches are key.

117 In this study, we carefully selected a set of six well characterized brother pairs who have the
118 same *ABCD1* pathogenic allele but are discordant for cerebral ALD: one brother has CALD and the
119 other has no white matter lesions on MRI (non-CALD). The brother pairs are close in age (no more
120 than 2 years apart), and range in age from 6-38 years at sample collection (Table 1). Blood samples
121 were obtained from each of these patients and underwent profiling through five omics technologies
122 (Figure 1) including whole genome sequencing (WGS), RNA sequencing (RNA-seq), EPIC DNA
123 methylation (DNAm) microarray, lipidomic profiling via liquid chromatography mass spectrometry
124 (LCMS), and protein profiling by LCMS. Each omics dataset was processed to quantify/map
125 features, undergo quality control analysis, and then used for group-wise comparisons between CALD
126 and non-CALD phenotype groups using univariate analysis. We first investigated the potential for a
127 single, shared modifier allele which could discriminate the two groups from the WGS data. Next, we
128 systematically compared the groups for each of these omics data sets to find potential markers
129 specific to the phenotype. We aggregated the datasets together after performing pairwise comparisons
130 and identified heterogeneous signals within sub-groups of the 6 families. To the best of our

131 knowledge this is the most comprehensive study to date in terms of systems biology characterization
132 of human ALD using a unique collection of samples.

133

134 **2 Materials and Methods**

135 **2.1 Project Overview**

136 An overview of the project can be found in Figure 1, which depicts the project phases
137 including patient phenotyping/sample collection, multi-omic data collection, feature
138 quantification/processing, quality control, and group-wise comparisons between phenotype groups.
139 In this project, six brother pairs affected by ALD but discordant for the presence of cerebral ALD
140 were included. Patients were selected from the Dutch cohort, an ongoing prospective natural history
141 study (Huffnagel, van Ballegoij, et al. 2019). Blood was drawn from the brother pairs and
142 lymphocyte pellets or plasma was isolated. Lymphocyte pellets were used for whole WGS, RNA-seq,
143 and DNAm. Fasted plasma was used for downstream LCMS analysis identifying either lipid or
144 protein abundances. Data was then processed independently for each of the platforms including
145 feature quantification/mapping, followed by platform specific quality control and group-wise
146 comparisons. Details regarding sample collection, platform specifications, and specific methodology
147 for each analysis performed in this project can be found in the Supplemental Methods section.

148 **2.2 Patient selection and phenotyping**

149 All patients were selected from the Dutch cohort, an ongoing prospective cohort study. All
150 patients are examined yearly (by ME) and undergo an MRI of the brain at the time of examination.
151 Samples are collected in the PEROX biobank. The presence of cerebral ALD is defined as the
152 presence of white matter lesions in a distribution consistent with ALD. The classification of the sibs
153 (CALD versus non-CALD) is valid at this time, but the non-affected individuals can theoretically
154 convert to cerebral ALD.

155 All samples were collected and stored in the PEROX Biobank according to a protocol
156 (METC2015_066) approved by the biobank review board of the Amsterdam UMC (BioBank
157 Toetsingscommissie AMC). All patients provided written informed consent for storage and use of
158 materials for medical research.

159

160

161 **2.3 Feature quantification and data processing**

162 For each platform, the data was processed independently following best-practices guidelines
163 from the groups generating the datasets. Details regarding feature quantification and assignment at
164 the gene, lipid, protein, and differentially methylated region (DMR) level can be found in the
165 Supplemental Methods section.

166

167 **2.4 Univariate modeling of CALD vs. non cerebral ALD**

168 Using univariate modeling techniques the prognostic power of each lipid, transcript (RNA) or protein
169 is calculated as:

170

$$y = \rho\beta + fam + \epsilon \quad (\text{eq1})$$

171

172

173 in which y is the observed value and ρ the phenotype (0 or 1), β is the weight and fam the family
174 cofactor. ϵ is the remaining error. Because methylation of DNA changes with age (McEwen et al.
175 2018), age is also included as a cofactor:

176

$$y = \rho\beta + fam + age + \epsilon \quad (\text{eq2})$$

177

178

179

180 when analysing the DNAm results. The significance (p-value) of the discriminating phenotype
181 (fixed) and family (random) effects are determined by ordinary least squares modelling (OLS) of the
182 data using the model from Eq. 1 in case of lipid and proteomic data (Harrison et al. 2018). In the case
183 of methylation and RNA sequencing data the p-values are determined by maximum likelihood
184 estimates (MLE) of the fixed and random effects using Limma (Ritchie et al. 2015) and edgeR
185 respectively (Ritchie et al. 2015; Robinson, McCarthy, and Smyth 2010).

186

187 **2.5 Allele comparisons in whole genome sequencing data**

188 Details on data processing, including variant calling and comparing across samples can be
189 found in the Supplemental Methods section. Briefly, allele comparisons were performed in whole
190 genome sequencing data on jointly genotyped variant datasets. For SNVs and indels, variants were
191 jointly genotyped and converted into a GEMINI database (Paila et al. 2013). This database was then
192 queried to identify subsets of discriminating alleles. For structural variants and mobile element
193 insertions, custom scripts were used to identify discordant genotypes from annotated jointly
194 genotyped variant tables. Discordant genotypes, stored as unique variant identifiers, were then placed
195 into Intervene for intersection analysis (Khan and Mathelier 2017).

196

197 **2.6 Aggregation of signal across platforms**

198 To assess the added value of combining the different platforms, significant signals prior to
199 multiple testing correction were collected for each omics platform and intersected at the annotated
200 gene level (hg19). Because of the lack of a clear mapping of lipids to genes the lipidomics platform
201 was excluded from this intersection allowing 4 possible intersections; DNAm-RNA, DNAm-Protein,
202 RNA-Protein and the overall intersection of DNAm-RNA-Protein. Further investigation of a shared
203 signal was performed by clustering the first 3 principal components (i.e. capturing the most variance)
204 of the log fold changes (top 10 and $p < 0.05$) of the combined platform data (including lipids).

205

206 **2.7 Assessing contribution of family effect per feature**

207 The data were modelled using the equations (above) in which both phenotypic and family effects are
208 estimated. We partitioned the variance for the lipids, proteins, and RNA datasets to identify the
209 contribution of the family effect, the phenotype effect, or the residual variance using the
210 variancePartition package (Hoffman and Schadt 2016). The same was repeated for the DNA
211 methylation dataset with the addition of the age, and phenotype:age variance terms. Next, we plotted
212 the top two principal components for each omics dataset before and after the removal of the variance

213 contributed from the family effect with the limma:removeBatchEffect tool (Ritchie et al. 2015).
214 Lastly, to determine the sensitivity/specificity of the findings for leaving out a one or two families all
215 the analyses (excluding methylation data analysis) that were run for the case of all families were
216 repeated with a one or two families left out (e.g. without fam 1, without fam 2, without fam 1 and
217 fam 2, etc.). We encapsulated this information in separate upset plots for each platform.
218

219

220 **3 Results**

221 **3.1 Lipidomics analysis of a fatty acid storage disorder**

222 Patients affected by ALD have a buildup of VLCFAs within cells in the body. Recent mass
223 spectrometry advances allow for broad, untargeted profiling of lipids ([Huffnagel, Dijkgraaf, et al.](#)
224 [2019](#)). We applied LCMS from plasma samples of each of the patients within this cohort as well as
225 matched control samples.

226 First, we identified differential lipid abundances between ALD (both non-CALD and CALD)
227 samples and control samples, with 139 lipids passing the threshold of p -value < 0.05 (Eq. 1, OLS),
228 and 17 lipids remaining significant after multiple testing correction (Bonferroni) (Figure S1, Table
229 S1) (Methods). The measured lipids are plotted as a volcano plot, that is the \log_2 fold change of ALD
230 over control versus corrected p -value (Figure 2A). We confirm that untargeted lipidomic profiling
231 can distinguish ALD from control samples via principal component analysis, and also capture the
232 expected differentially abundant lipids between control and ALD samples including the known ALD
233 biomarker LPC(26:0) (Figure 2B,C).

234 Next, we compared CALD and non-CALD groups for differences in lipid abundance which
235 could act as markers of cerebral demyelination. Of note, the principal component analysis which
236 separates ALD from control did not separate CALD from non-CALD, i.e. the differences in lipid
237 profiles between these two phenotypes are much less pronounced than the differences separating
238 ALD patients from controls (Figure 2B). The measured lipids are plotted as a volcano plot, that is the
239 \log_2 fold change of CALD over non-CALD versus transformed p -value (Figure 2D). In total 22
240 lipids were found to have different abundances between the two groups with p -value < 0.05 , however
241 none of the lipids remained significant after correcting for multiple testing (Table S2). The observed
242 differences are much smaller between CALD and non-CALD compared to ALD and control, as
243 highlighted by the differences in fold change axes (Figure 2A, D). Interestingly, there was a higher
244 abundance in the non-CALD group for several key VLCFAs involved in ALD including PC(44:4)
245 and Cer(d42:3), the latter reaching p -value < 0.05 (Figure 2C,E). While some lipids show a relatively
246 large fold change between CALD and non-CALD groups as a whole, the signal is not consistent for
247 every family. An example of this can be seen in SM(d36:2) or PS(43:3) (Figure S2). This limits the
248 prognostic power of these lipids as consistent markers delineating the phenotype. Lastly, we
249 observed a large range of lipid abundances within the control group for several of the differential
250 lipids between CALD and non-CALD, which could indicate that these lipids are variable within
251 healthy individuals and the signal we observe between CALD and non-CALD could be due to noise
252 or variation in the healthy population (Figure 2E, Figure S1).

253

254 **3.2 Discordant genotype analysis for the identification of a modifier allele**

255 Using whole genome sequencing, we investigated a range of variant classes for discordant
256 alleles between siblings. These discordant alleles are then intersected across multiple families under
257 the hypothesis that polymorphic differences contribute to cerebral demyelination.

258 We first focused on alleles which emerged from previous modifier studies to see if they are
259 confirmed. Proposed modifier alleles from target gene studies have identified two candidates within
260 *ELOVL1* (rs839765) and *CYP4F2* (rs2108622) (van Engen et al. 2016; Kemp et al. 2012). Within
261 this cohort, those modifier alleles do not segregate with ALD phenotype (Table 1), nor are the
262 genotypes shared or lacking in the confidently phenotyped CALD patients. Furthermore, it has been
263 suggested that *APOE* genotypes--which are a combination between two SNP sites to produce *APOE2*
264 ($\epsilon 2$), *APOE3* ($\epsilon 3$), and *APOE4* ($\epsilon 4$) alleles--may be markers of disease severity and cerebral
265 progression (Orchard, Markowski, et al. 2019). These *APOE* alleles do not segregate with disease nor
266 are they shared by all CALD patients. Together, these results suggest limited prognostic power of
267 these alleles, and perhaps supports heterogeneous contributions of genetic background to disease
268 progression.

269 Next, for several variant classes, we performed a discordant analysis between siblings and
270 intersected these alleles across families (Figure 3). We considered four genotypic categories termed
271 dominant protective, recessive protective, dominant damaging, or recessive damaging based on the
272 genotype (heterozygous: dominant or homozygous: recessive) and the sibling which carries the
273 genotype (CALD: damaging or non-CALD: protective). Performing this genotypic analysis on SNVs
274 and indels, we identified $\sim 6.0 \times 10^5$ discordant candidate variants in the dominant categories from each
275 family, and $\sim 3.0 \times 10^5$ discordant candidate variants from the recessive categories (Figure 3A, Table
276 S3). Despite the large number of discordant candidates per family, intersecting these sets across
277 families reduces the candidates dramatically, resulting in only two candidate variants at the
278 intersection of all 6 families (Figure 3B) (Table S4). A recessive damaging variant downstream of the
279 *WIBG/PYMI* gene (rs7980776) and recessive protective allele (rs55639747/rs61327784) within a
280 *CCDC67/DEU1* intronic region (Figure S3 & S4). We validated our approach with a parallel
281 pipeline utilizing the new DeepVariant tool (Poplin et al. 2018), which claims higher accuracy than
282 GATK HaplotypeCaller (Table S4). There is high concordance between the two variant call sets, and
283 they produced the same two variants within the intersection. A single additional variant was reported
284 using DeepVariant under the recessive damaging model, however the variant did not pass the manual
285 inspection quality assessment. In silico analysis of both variants suggests these variants have little
286 functional effect, and the associated genes did not link to the cerebral demyelination phenotype
287 (Supplemental Results).

288 For the other variant classes, including structural variants (SVs) and mobile element
289 insertions (MEIs), we performed joint genotyping to identify shared and discordant alleles in the
290 same manner as SNVs and indels. We identified ~ 1500 and ~ 400 SVs in the dominant and recessive
291 categories respectively, and ~ 400 and ~ 100 MEIs (Table S3). Unsurprisingly, these discordant events
292 were not shared across more than 4 families (Table S4). We further manually inspected the regions
293 around the discordant SNVs/indels identified above, and did not find any other segregating SVs or
294 MEIs.

295 Lastly, we extended our discordance analysis to the mitochondrial genome to examine
296 candidate alleles which may show evidence of heteroplasmy which are not shared between two
297 siblings. We identified that between 129 and 547 mitochondrial variants per sample, of which 52 to
298 476 are heteroplasmic, and none are consistent discriminating variants between phenotypes shared
299 across all families (Table S5). Further, if we aggregated at the gene level, we did not find any
300 heteroplasmic variants consistent across the same gene.

301 In recognition of a study limitation--the fact that some non-CALD patients may progress to
302 CALD--we further intersected alleles shared by all CALD patients. These variants were annotated by
303 impact or as eQTLs defined in GTEx (Supplemental Methods). There were 48 variants present in the

304 heterozygous state across all CALD patients, where no non-CALD patients were heterozygous (Table
305 S6). Of these, a haploblock containing 20 variants was identified overlapping the gene *TPCN2*,
306 including a missense variant (rs3750965) (Figure S5). Interestingly, the only patient which was
307 homozygous for this variant is the youngest non-CALD patient within the cohort, suggesting that this
308 gene could be of significance should the patient develop the cerebral demyelination phenotype.
309

310 **3.3 Univariate modeling of phenotype differences across omics platforms**

311 Beyond identifying a single genetic modifier allele, the omics platforms allow for the
312 identification of candidate molecular signatures which can discriminate between the CALD and non-
313 CALD phenotypes. Using univariate analysis we identify differences across each platform at the
314 feature level, to search for a signal which can be used as a marker for transition to CALD. Further,
315 we leverage these molecular signatures to provide insight into the pathogenesis of cerebral
316 demyelination.
317

318 **3.3.1 Transcriptomics**

319 Examining RNA expression using RNA-seq provides a measurement for nearly all expressed
320 protein coding genes in the genome. Differential gene expression was calculated between the two
321 phenotype groups using the univariate model accounting for family effect (equation 1). There were
322 199 genes found with a p-value < 0.05, although none remained significant after multiple testing
323 correction (Bonferroni) (Figure 4A,B, Figure S6). This is likely due to the low number of samples
324 and relatively small differences that were observed between the two groups. Furthermore, many of
325 the genes identified as significant were inconsistent in one or more of the sibling pairs, limiting the
326 diagnostic utility as a marker (Figure 4B). Despite not having significant genes after multiple testing
327 correction, we performed enrichment analysis using GO (gene annotation) and KEGG (pathway
328 annotation) to derive insights based on the 199 genes passing a threshold of p-value < 0.05 (Figure
329 S7). Of note, elevated interferon related processes suggest that the host may be reacting to pathogens
330 activating the immune system (Hoffmann, Schneider, and Rice 2015). It is therefore no surprise that
331 3 chemokines (*CXCL6*, *CXCL8* and *IFI27*) were found in the top 10 differentially expressed genes.
332 Amongst the remainder of the proteins encoded by the top 10 differentially expressed genes, the D-
333 Xylulokinase gene (*XYLB*) encodes for the protein that catalyzes the ATP-dependent phosphorylation
334 of D-xylulose to produce xylulose-5-phosphate (Xu5P) therefore *XYLB* may play an important role in
335 metabolic disease given that Xu5P is a key regulator of glucose metabolism and lipogenesis (Bunker
336 et al. 2013). *GATM* has been associated with statin intolerance (V Willrich et al. 2018) and its
337 function to catalyze creatine and possibly affect the production of ceramides (Turer et al. 2017).
338 *MYOB1B* is a protein that may participate in a process critical to neuronal development and function
339 such as cell migration, neurite outgrowth and vesicular transport (Sittaramane and Chandrasekhar
340 2008).

341
342

343 **3.3.2 Epigenomics**

344 DNA methylation has been linked to changes in gene expression, and is an important readout
345 of some environmental impacts upon the cell. Measuring DNA methylation is typically done at
346 specific methylation sites (CpGs), and then aggregated across regions where several sites have
347 similar trends of methylation levels to find differentially methylated regions (DMRs). Here, we used
348 the MethylationEPIC BeadChip which targets over 850,000 CpGs. Using LIMMA modeling
349 including age as cofactor (equation 2) 264 CpGs had a nominal p-value < 0.0005. Of these 264 CpGs,

350 16 passed the delta beta (i.e. difference between methylation levels of CALD vs non-CALD) of >5%
351 (Table S8). When aggregating these loci into a DMR analysis, we identified 22 regions passing
352 thresholds of $FDR < 0.05$ and >10% methylation change (Figure 4C). Multiple CpGs map to the same
353 gene and show a large delta beta, which we identified in the genes *PTPRN2* and *RGS14* (Figure 4D).
354 *RGS14* may alter calcium levels to enhance long term potentiation and learning (Lee et al. 2010).
355 Due to its presence in neurosecretory vesicles, *PTPRN2* has been implicated in insulin and
356 neurotransmitter exocytosis (Sengelaub et al. 2016). Furthermore, *PTPRN2* hypermethylation has
357 been identified within a separate study which compared DNA methylation between CALD and non-
358 CALD patients (Schlüter et al. 2018).

359

360 **3.3.3 Proteomics**

361 In addition to profiling lipids, LCMS can be used for high throughput profiling of proteins thus
362 enabling the identification of differential protein abundances between samples. Applying proteomics
363 to these 12 patients yielded a quantification of 5,862 peptides which were matched against 351
364 protein groups. Comparing CALD and non-CALD groups, we found 16 proteins with differential
365 abundances ($p < 0.05$) (Figure 4E,F; Figure S8; Table S9). Investigating the top hits we find 4/16
366 proteins associated with immunoglobulin heavy chain (IGHV4-34, IGHV3-30, IGHV3-7, and
367 PODOX6), 2/16 are associated with immunoglobulin kappa variables (IGKV6D-21 and IGKV1D-
368 33), and with PODOX8 also being related to immunoglobulin, half of these proteins are related to the
369 immune system (Parra et al. 2016). All of these immunoglobulin proteins were up-regulated in the
370 CALD samples. Also related to the immune system is CD5L, a secreted glycoprotein that participates
371 in host response to bacterial infection (Sanjurjo et al. 2015) and is also known to regulate lipid
372 biosynthesis (Wang et al. 2015). Beyond immune system proteins, we identified proteins associated
373 with the brain or with involvement in lipid metabolism. ECM1 has been associated with lipid
374 proteinosis in which brain damage develops over time and is associated with the development of
375 cognitive disabilities and epileptic seizures (Zhang et al. 2014). The role of APOL1 is not yet clear
376 but it has been associated with the lipid biology in the podocyte (Fornoni, Merscher, and Kopp
377 2014). Copy number variants of MINPP1 have been associated with varying IP6 levels (Waugh
378 2016) and IP6 has been reported to suppress lipid peroxidation (Foster et al. 2017). APOC3 is a key
379 player in triglyceride-rich lipoprotein metabolism (Ramms and Gordts 2018) and regulated by the
380 peroxisome proliferator-activated receptor- α (Liu et al. 2015). PFN1 has recently been reported in a
381 CALD study which looked at markers of autoreactivity, identifying anti-PFN1 antibodies present in a
382 large proportion of CALD patients (Orchard, Nascene, et al. 2019). Together, these protein signals
383 could have significance with respect to the pathophysiology of cerebral demyelination, by
384 highlighting differences around proteins involved in lipid metabolism as well as immune response.

385

386 **3.3.4 Estimating variance of family effect**

387 The univariate modeling of CALD vs. non-CALD for each of the individual omics platforms
388 was unsuccessful in identifying significant hits after multiple testing correction. While traditional
389 multiple testing correction methods may be too strict for the omics technologies, we still cannot rule
390 out the possibility that our top hits arise by chance due to variability. Furthermore, our top hits per
391 platform still exhibited a high amount of variance between families, and a lack of consistent signal in
392 molecular features across the entire cohort (Figure 4 B, D, F). Within our model we included the
393 effect of the family on the level of the measured signal, and thus we are able to capture the
394 contribution of family structure to a feature's abundance (equation 1, equation 2). To illustrate the
395 contribution of these effects, we partitioned the variance contribution within our linear models
396 (Methods). The phenotype effect, total family effect, and residual variance were extracted from our

397 model for each of the features within the RNA-seq, proteomics, and lipidomics platforms (Figure
398 S9). As DNAm varies with age we additionally extracted the variance contributed from the age or
399 phenotype-by-age effects. Clearly, the contribution of variance from the phenotype is small in the
400 majority of features across all omics datasets, and a large residual variance indicates a high level of
401 noise present in these high dimensional assays (Figure S9). We further demonstrated the
402 heterogeneity in the data by subsetting the families and then repeating comparisons between CALD
403 and non-CALD phenotypes. By leaving out one or two families, the β in equations 1-2 are re-
404 evaluated for the RNA, protein, and lipid datasets. The number of candidates increased with removal
405 of each family, which could be interpreted as potential modifier signatures present in a subset of
406 families, but absent from others (Figure S10).
407

408 **3.4 Integrating multi-omic datasets**

409 As it was our intention to identify molecular marker features underlying cerebral
410 demyelination, we investigated the omics datasets independently to identify a consistent signal.
411 However, owing to a large amount of inter-family variance, we are limited in our ability to identify a
412 statistically significant feature which separates the two phenotypes. As the multi-omic assays should
413 be complementary to each other, we searched for genes which showed differences between the
414 groups in multiple assays. We searched the phenotype comparison between all families, as well as
415 the results from the leave-one-out analysis, wherein we withheld a family and repeated the modeling
416 between the two phenotype groups (Methods). Intersections showed overlapping evidence at the
417 DNA methylation and RNA levels, as well as overlap between RNA and protein levels, for eight
418 genes. Focusing only on the intersection of all families, only *PTPRN2* has differential signal from
419 both DNA methylation and RNA levels (Figure 4 B,D). Additional genes were identified in the
420 leave-one-out subsets (Table 2).

421 In the multi-omic data we observed that several of the molecular features have trends of
422 differential abundance/expression in a subset of the families. To illustrate this, and attempt to identify
423 clusters within the data, we gathered per-family log₂-fold-change of CALD over non-CALD for the
424 top hits from the lipid, protein, and RNA datasets. We took this approach because it removes the
425 differences in absolute levels of expression between families. Noticeably, the fold-change values are
426 not consistent for each family, as evidenced by a lack of consistent colouring for each of the features
427 (rows) within the heatmap (Figure S11 A). Family 2 and family 6 were more similar in their
428 CALD/non-CALD ratios for these features. This is further supported by a principal component
429 analysis, wherein family 2 and family 6 are separated from the other four families on the first
430 principal component (Figure S11 C). However, this trend does not hold when the set of features is
431 increased to all hits with p-value < 0.05 across the three platforms, as family 1 and 5 cluster together
432 with the other four families as an outer group (Figure S11 B). Thus, clustering these families based
433 on top differential features does not reveal confident sub-groupings within the small cohort.
434

435 **3.5 Specific modifier hypothesis testing**

436 Finding molecular markers which delineate cerebral demyelination in patients with ALD is an
437 ongoing research problem. Additionally, understanding the pathophysiology of cerebral
438 demyelination and potential disruption of the blood brain barrier has implications for diseases beyond
439 ALD. Different hypotheses have been suggested, including involvement of the immune system in
440 autoreactivity or as a response to severe viral infections. Using the multi-omics dataset, which gives
441 us insight into the complexities of the underlying complex biological system, we tested recently

442 proposed modifiers of cerebral demyelination to see if there is evidence of their discriminatory power
443 within the blood samples profiled in our dataset.

444 It has recently been demonstrated that autoreactivity to profilin (PFN1) occurs in patients
445 affected by CALD, and may be a discriminating marker of cerebral demyelination (Orchard,
446 Nascene, et al. 2019). We investigated differences in PFN1 methylation, RNA, and protein levels
447 between CALD and non-CALD patients to see if this observation is confirmed in our dataset. At the
448 methylation and RNA level, we did not see a consistent signal differentiating the CALD and non-
449 CALD groups, but at the protein level we observe an increased amount of PFN1 in the CALD group
450 for four out of six families (Figure 5 A, B, Figure 4F). This is consistent with the observation from
451 the previous study that not all patients exhibit PFN1 autoreactivity, and the dramatically increased
452 protein levels could precede or act as biomarkers of the autoimmune response within the subset of
453 patients who exhibit this trend.

454 Another study focused on DNA methylation (DNAm) as a marker of CALD, and investigated
455 the intact white matter of brains from patients affected by ALD with and without the cerebral
456 demyelination phenotype (Schlüter et al. 2018). Whether or not the signals they identify confirm
457 within the blood within a separate cohort is important if these proposed marker genes are to be used
458 within newborn screening. Within their analysis they identified differential methylation signals at
459 several genes, two of which are *LPIN1* and *UNC45A*. Within this cohort, we see no differential
460 methylation signal in the blood for *LPIN1*, and a slight hypermethylation (although not significant) in
461 *UNC45A* (Figure 5C,E). Investigating the RNA shows that while both these genes are highly
462 expressed, there are no consistent differences between the two phenotype groups.

463 Lastly, it is possible that a viral infection causing an immune response is the phenotypic trigger
464 for progression to CALD, as this is suggested to be a candidate environmental modifier from other
465 cerebral demyelination diseases including multiple sclerosis (Libbey, Lane, and Fujinami 2014). As
466 is the case in several cancers, RNA-seq can capture actively expressing viral RNA within a sample.
467 To test the hypothesis of whether or not we could observe different expressing viruses within the
468 RNA-seq of these patients, we used the tool Centrifuge to identify traces of viral (or bacterial)
469 sequences (Table S10) (Kim et al. 2016). Aside from identifying human, synthetic construct, and
470 endogenous retrovirus, no significant viral or bacterial sequences were identified.

471

472 **4 Discussion**

473

474 In this study we took a systems biology approach to identify personal molecular
475 characteristics, either genetic or molecular markers, which may prognosticate the onset of cerebral
476 demyelination in patients affected by ALD. Identifying a single modifier consistent across all
477 individuals has importance because of its potential utility as prognosticator or biomarker heralding
478 the transition to cerebral demyelination, and this carries tremendous treatment implications.

479 Our cohort was comprised of carefully phenotyped brothers affected by ALD who were
480 discordant for the severe cerebral demyelination phenotype. We collected blood and performed high
481 throughput experiments to profile the DNA, methylated DNA, RNA, lipids, and proteins. In
482 summary, we did not find a strong, convincing, univariate marker which can differentiate all of the
483 CALD and non-cerebral patients in this small cohort. There are several explanations for this negative
484 result: the small cohort with only six discordant sibling pairs of different ethnic background, the
485 possibility that one or more non-CALD patients may still develop cerebral demyelination, high inter-
486 individual variability, and finally the possibility of multiple modifiers and/or an exogenous or non
487 genetic modifier such as infection or physical trauma. In spite of these limitations, we still emerged
488 with interesting results from each of the omics platforms from this pilot study including discordant
489 genotypes separating all CALD and non-CALD patients, confirmations of recently proposed CALD

490 modifiers, and a suspected involvement of differential activity within the immune system in patients
491 with cerebral demyelination.

492 In our genetic approach, we identified two discordant genotypes shared between all six
493 brother-pairs: an intronic SNV in *DEU1* and an SNV downstream of *WIBG*. Although in silico
494 analysis of the variants and the function of the associated genes did not link these alleles to the
495 cerebral demyelination phenotype, it is of interest to see if they replicate in a larger cohort.
496 Examining variants shared by all CALD patients led to the identification of a missense
497 polymorphism in *TPCN2*, a gene which localizes to lysosomal membranes. This exists in a
498 segregating haplotype block, and is absent from all non-CALD patients except for family 4--the
499 youngest patient with the highest chance to develop cerebral demyelination--where the haploblock is
500 homozygous. How this variant segregates in a larger patient cohort could be of interest. None of the
501 previously proposed modifier alleles, emerging from GWAS or target-gene studies, confirmed within
502 our cohort.

503 Although our analysis was burdened by high inter-individual variability, we were able to
504 identify univariate molecular markers with increased confidence due to replication--by multiple
505 omics levels and/or by confirming previously proposed modifier markers. A recent study (Schlüter et
506 al. 2018) showed CALD patients with DNA hypermethylation within *PTPRN2*, which we confirm in
507 our study and support with decreased mRNA expression in CALD patients (both platforms reaching
508 p-value <0.05 before multiple testing correction). The same study showed hypermethylation of
509 *LPIN1* and *UNC45A*, the latter of which we confirm (although not statistically significant) as slightly
510 hypermethylated in CALD samples. Of note, that study used brain tissue to derive their signal
511 whereas we use blood samples. Another study utilized CSF and blood plasma, including longitudinal
512 data from ALD patients pre- and post cerebral demyelination, to identify autoreactivity to Profilin 1
513 (*PFN1*) within CALD patients (Orchard, Nascene, et al. 2019). They observed auto-antigens to *PFN1*
514 in the blood, and increased *PFN1* levels in CSF, in ~50% of CALD patients. In our cohort, four out
515 of six patients exhibit increased *PFN1* protein levels, in-line with the observation that *PFN1*
516 phenotype is not ubiquitous across all CALD patients. We further contribute to this observation by
517 showing no differences at the DNAm or mRNA levels, pointing towards a separate mechanism of
518 upregulation/overabundance of *PFN1*.

519 As ALD is a peroxisomal disorder, the lipidomic analysis presented here is of interest. The
520 lipid profiling data confirmed previous observations regarding VLCFA abundance differences in
521 ALD samples when compared to controls. Specifically, the phosphatidylcholines (PC) species
522 containing very long-chain fatty acids are more abundant in the ALD group compared to the control
523 group. Furthermore, the suitability of *LPC(C26:0)* to function as a marker for ALD in newborn
524 screening was confirmed. Differences in lipid abundance between CALD and non-CALD groups did
525 not reach significance after multiple testing correction, likely due to a lack of consistent lipid
526 differences between all brother pairs. Nevertheless, the differential lipids between CALD and non-
527 CALD provide insight into the pathophysiology of CALD as CALD patients had lower levels of
528 sphingomyelin and its precursor ceramide, in line with disease progression. This could support the
529 findings that the sphingolipid systems hold important roles in CNS disorders like Alzheimer's,
530 Parkinson's and Huntington's (Assi et al. 2013).

531 Beyond identifying phenotype-stratifying molecular features, we investigated the top hits at
532 the gene-level from each omics platform for any relation to the pathophysiology of CALD. Literature
533 searches highlighted genes involved in lipid metabolism, the nervous system, and the immune
534 system. Gene Ontology and KEGG pathways further supported these observations. Larger datasets
535 are needed to draw conclusions from differentially abundant molecular features.

536 Throughout this work we have identified certain limitations of our approach which should be
537 considered in future work focused on modifiers of rare disease, especially for other inborn errors of
538 metabolism (e.g. Gaucher disease). First, we suffered from having a small number of samples and a

539 high number of observed features. For future univariate marker investigations we recommend
540 focusing only on protein or mRNA and increasing the number of samples. Second, our genetic
541 analysis was limited by the possibility of future transition to the CALD state for any of our non-
542 CALD patients, especially those patients who have not reached maturity. Recent epidemiological
543 analysis shows that cerebral demyelination can occur throughout the lifetime of an ALD patient
544 (Huffnagel, Laheji, et al. 2019), so genetic studies should focus on older (60-70 years old) patients
545 who have not developed the cerebral demyelination phenotype. While discordant brother pairs
546 reaching old age are challenging to find, a collection of genotyped non-CALD patients older than 60-
547 70 years of age could serve as a good control. Third, we are limited in capturing relevant biological
548 insights because we are profiling blood not CSF/brain tissue. Lastly, while we profile DNA
549 methylation, we don't capture other components of the environment which could have an impact
550 including microbiome and pathogen exposure history.

551 With newborn screening now a reality for ALD, prognostication and timing of therapy
552 becomes more relevant than ever before; thus modifier studies to decipher a protector or marker for
553 cerebral demyelination will continue (Moser and Fatemi 2018). We believe that this dataset can
554 continue to be mined and used for testing the replication of proposed phenotypic markers. Further,
555 the data within this study could be used as part of a larger dataset examining multivariate signals
556 differentiating the two classes. Whether it is a collection of genetic markers or a pattern of multiple
557 molecular features, it is clear that there is a need for a larger sample size. As such, we make the
558 measurements through this study available for future use to the community, with the hopes that the
559 data can serve as a secondary confirmation of new modifier hypotheses, or as part of a larger dataset
560 for investigating the complex nature of cerebral demyelination.

561

562 **5 References**

- 563 Argmann, Carmen A., Sander M. Houten, Jun Zhu, and Eric E. Schadt. 2016. "A Next Generation
564 Multiscale View of Inborn Errors of Metabolism." *Cell Metabolism* 23 (1): 13–26.
- 565 Assi, Emma, Denise Cazzato, Clara De Palma, Cristiana Perrotta, Emilio Clementi, and Davide
566 Cervia. 2013. "Sphingolipids and Brain Resident Macrophages in Neuroinflammation: An
567 Emerging Aspect of Nervous System Pathology." *Clinical & Developmental Immunology* 2013
568 (September): 309302.
- 569 Beer, Marlijn H. de, and Philip Scheltens. 2016. "Cognitive Decline in Patients with Chronic
570 Hydrocephalus and Normal Aging: 'Growing into Deficits.'" *Dementia and Geriatric Cognitive
571 Disorders Extra* 6 (3): 500–507.
- 572 Bunker, Richard D., Esther M. M. Bulloch, James M. J. Dickson, Kerry M. Loomes, and Edward N.
573 Baker. 2013. "Structure and Function of Human Xylulokinase, an Enzyme with Important Roles
574 in Carbohydrate Metabolism." *The Journal of Biological Chemistry* 288 (3): 1643–52.
- 575 Engelen, Marc, Mathieu Barbier, Inge M. E. Dijkstra, Remmelt Schür, Rob M. A. de Bie, Camiel
576 Verhamme, Marcel G. W. Dijkgraaf, et al. 2014. "X-Linked Adrenoleukodystrophy in Women:
577 A Cross-Sectional Cohort Study." *Brain: A Journal of Neurology* 137 (Pt 3): 693–706.
- 578 Engelen, Marc, Stephan Kemp, Marianne de Visser, Björn M. van Geel, Ronald J. A. Wanders,
579 Patrick Aubourg, and Bwee Poll-The. 2012. "X-Linked Adrenoleukodystrophy (X-ALD):
580 Clinical Presentation and Guidelines for Diagnosis, Follow-up and Management." *Orphanet
581 Journal of Rare Diseases*. <https://doi.org/10.1186/1750-1172-7-51>.
- 582 Engen, Catherine E. van, Rob Ofman, Inge M. E. Dijkstra, Tessa Jacobs van Goethem, Eveline
583 Verheij, Jennifer Varin, Michel Vidaud, et al. 2016. "CYP4F2 Affects Phenotypic Outcome in
584 Adrenoleukodystrophy by Modulating the Clearance of Very Long-Chain Fatty Acids."
585 *Biochimica et Biophysica Acta* 1862 (10): 1861–70.

- 586 Fornoni, Alessia, Sandra Merscher, and Jeffrey B. Kopp. 2014. "Lipid Biology of the Podocyte—
587 new Perspectives Offer New Opportunities." *Nature Reviews Nephrology*.
588 <https://doi.org/10.1038/nrneph.2014.87>.
- 589 Foster, Shadae R., Lowell L. Dilworth, Felix O. Omoruyi, Rory Thompson, and Ruby L. Alexander-
590 Lindo. 2017. "Pancreatic and Renal Function in Streptozotocin-Induced Type 2 Diabetic Rats
591 Administered Combined Inositol Hexakisphosphate and Inositol Supplement." *Biomedicine &
592 Pharmacotherapy = Biomedecine & Pharmacotherapie* 96 (December): 72–77.
- 593 Génin, Emmanuelle, Josué Feingold, and Françoise Clerget-Darpoux. 2008. "Identifying Modifier
594 Genes of Monogenic Disease: Strategies and Difficulties." *Human Genetics* 124 (4): 357–68.
- 595 Harrison, Xavier A., Lynda Donaldson, Maria Eugenia Correa-Cano, Julian Evans, David N. Fisher,
596 Cecily E. D. Goodwin, Beth S. Robinson, David J. Hodgson, and Richard Inger. 2018. "A Brief
597 Introduction to Mixed Effects Modelling and Multi-Model Inference in Ecology." *PeerJ* 6
598 (May): e4794.
- 599 Hoffman, Gabriel E., and Eric E. Schadt. 2016. "variancePartition: Interpreting Drivers of Variation
600 in Complex Gene Expression Studies." *BMC Bioinformatics* 17 (1): 483.
- 601 Hoffmann, Hans-Heinrich, William M. Schneider, and Charles M. Rice. 2015. "Interferons and
602 Viruses: An Evolutionary Arms Race of Molecular Interactions." *Trends in Immunology* 36 (3):
603 124–38.
- 604 Hosseinbarkooie, Seyyedmohsen, Miriam Peters, Laura Torres-Benito, Raphael H. Rastetter,
605 Kristina Hupperich, Andrea Hoffmann, Natalia Mendoza-Ferreira, et al. 2016. "The Power of
606 Human Protective Modifiers: PLS3 and CORO1C Unravel Impaired Endocytosis in Spinal
607 Muscular Atrophy and Rescue SMA Phenotype." *American Journal of Human Genetics* 99 (3):
608 647–65.
- 609 Huffnagel, Irene C., Wouter J. C. van Ballegoij, Björn M. van Geel, Johanna M. B. W. Vos, Stephan
610 Kemp, and Marc Engelen. 2019. "Progression of Myelopathy in Males with
611 Adrenoleukodystrophy: Towards Clinical Trial Readiness." *Brain: A Journal of Neurology* 142
612 (2): 334–43.
- 613 Huffnagel, Irene C., Marcel G. W. Dijkgraaf, Georges E. Janssens, Michel van Weeghel, Björn M.
614 van Geel, Bwee Tien Poll-The, Stephan Kemp, and Marc Engelen. 2019. "Disease Progression
615 in Women with X-Linked Adrenoleukodystrophy Is Slow." *Orphanet Journal of Rare Diseases*
616 14 (1): 30.
- 617 Huffnagel, Irene C., Fiza K. Laheji, Razina Aziz-Bose, Nicholas A. Tritos, Rose Marino, Gabor E.
618 Linthorst, Stephan Kemp, Marc Engelen, and Florian Eichler. 2019. "The Natural History of
619 Adrenal Insufficiency in X-Linked Adrenoleukodystrophy: An International Collaboration." *The
620 Journal of Clinical Endocrinology and Metabolism* 104 (1): 118–26.
- 621 Kemp, Stephan, Johannes Berger, and Patrick Aubourg. 2012. "X-Linked Adrenoleukodystrophy:
622 Clinical, Metabolic, Genetic and Pathophysiological Aspects." *Biochimica et Biophysica Acta*
623 1822 (9): 1465–74.
- 624 Kemp, Stephan, Irene C. Huffnagel, Gabor E. Linthorst, Ronald J. Wanders, and Marc Engelen.
625 2016. "Adrenoleukodystrophy – Neuroendocrine Pathogenesis and Redefinition of Natural
626 History." *Nature Reviews Endocrinology*. <https://doi.org/10.1038/nrendo.2016.90>.
- 627 Khan, Aziz, and Anthony Mathelier. 2017. "Intervene: A Tool for Intersection and Visualization of
628 Multiple Gene or Genomic Region Sets." *BMC Bioinformatics* 18 (1): 287.
- 629 Kim, Daehwan, Li Song, Florian P. Breitwieser, and Steven L. Salzberg. 2016. "Centrifuge: Rapid
630 and Sensitive Classification of Metagenomic Sequences." *Genome Research* 26 (12): 1721–29.
- 631 Korenke, G. C., S. Fuchs, E. Krasemann, H. G. Doerr, E. Wilichowski, D. H. Hunneman, and F.
632 Hanefeld. 1996. "Cerebral Adrenoleukodystrophy (ALD) in Only One of Monozygotic Twins
633 with an Identical ALD Genotype." *Annals of Neurology* 40 (2): 254–57.

- 634 Lee, Sarah Emerson, Stephen B. Simons, Scott A. Heldt, Meilan Zhao, Jason P. Schroeder,
635 Christopher P. Vellano, D. Patrick Cowan, et al. 2010. “RGS14 Is a Natural Suppressor of Both
636 Synaptic Plasticity in CA2 Neurons and Hippocampal-Based Learning and Memory.”
637 *Proceedings of the National Academy of Sciences of the United States of America* 107 (39):
638 16994–98.
- 639 Libbey, Jane E., Thomas E. Lane, and Robert S. Fujinami. 2014. “Axonal Pathology and
640 Demyelination in Viral Models of Multiple Sclerosis.” *Discovery Medicine* 18 (97): 79–89.
- 641 Liu, Chao, Qianqian Guo, Mengchen Lu, and Yunman Li. 2015. “An Experimental Study on
642 Amelioration of Dyslipidemia-Induced Atherosclerosis by Clematichinenoside through Regulating
643 Peroxisome Proliferator-Activated Receptor- α Mediated Apolipoprotein A-I, A-II and C-III.”
644 *European Journal of Pharmacology*. <https://doi.org/10.1016/j.ejphar.2015.04.015>.
- 645 McEwen, Lisa M., Meaghan J. Jones, David Tse Shen Lin, Rachel D. Edgar, Lucas T. Husquin, Julia
646 L. MacIsaac, Katia E. Ramadori, et al. 2018. “Systematic Evaluation of DNA Methylation Age
647 Estimation with Common Preprocessing Methods and the Infinium MethylationEPIC BeadChip
648 Array.” *Clinical Epigenetics* 10 (1): 123.
- 649 Moser, Ann B., and Ali Fatemi. 2018. “Newborn Screening and Emerging Therapies for X-Linked
650 Adrenoleukodystrophy.” *JAMA Neurology* 75 (10): 1175–76.
- 651 Moser, Ann, Richard Jones, Walter Hubbard, Silvia Tortorelli, Joseph Orsini, Michele Caggana, Beth
652 Vogel, and Gerald Raymond. 2016. “Newborn Screening for X-Linked Adrenoleukodystrophy.”
653 *International Journal of Neonatal Screening*. <https://doi.org/10.3390/ijns2040015>.
- 654 Ofman, Rob, Inge M. E. Dijkstra, Carlo W. T. van Roermund, Nena Burger, Marjolein Turkenburg,
655 Arno van Cruchten, Catherine E. van Engen, Ronald J. A. Wanders, and Stephan Kemp. 2010.
656 “The Role of ELOVL1 in Very Long-Chain Fatty Acid Homeostasis and X-Linked
657 Adrenoleukodystrophy.” *EMBO Molecular Medicine* 2 (3): 90–97.
- 658 Orchard, Paul J., Todd W. Markowski, Leeann Higgins, Gerald V. Raymond, David R. Nascene,
659 Weston P. Miller, Elizabeth I. Pierpont, and Troy C. Lund. 2019. “Association between APOE4
660 and Biomarkers in Cerebral Adrenoleukodystrophy.” *Scientific Reports* 9 (1): 7858.
- 661 Orchard, Paul J., David R. Nascene, Ashish Gupta, Mandy E. Taisto, Leeann Higgins, Todd W.
662 Markowski, and Troy C. Lund. 2019. “Cerebral Adrenoleukodystrophy Is Associated with Loss
663 of Tolerance to Profilin.” *European Journal of Immunology* 49 (6): 947–53.
- 664 Paila, Umadevi, Brad A. Chapman, Rory Kirchner, and Aaron R. Quinlan. 2013. “GEMINI:
665 Integrative Exploration of Genetic Variation and Genome Annotations.” *PLoS Computational
666 Biology* 9 (7): e1003153.
- 667 Parra, David, Tomáš Korytář, Fumio Takizawa, and J. Oriol Sunyer. 2016. “B Cells and Their Role
668 in the Teleost Gut.” *Developmental and Comparative Immunology* 64 (November): 150–66.
- 669 Pierpont, Elizabeth I., Julie B. Eisengart, Ryan Shanley, David Nascene, Gerald V. Raymond, Elsa
670 G. Shapiro, Rich S. Ziegler, Paul J. Orchard, and Weston P. Miller. 2017. “Neurocognitive
671 Trajectory of Boys Who Received a Hematopoietic Stem Cell Transplant at an Early Stage of
672 Childhood Cerebral Adrenoleukodystrophy.” *JAMA Neurology* 74 (6): 710–17.
- 673 Poplin, Ryan, Pi-Chuan Chang, David Alexander, Scott Schwartz, Thomas Colthurst, Alexander Ku,
674 Dan Newburger, et al. 2018. “A Universal SNP and Small-Indel Variant Caller Using Deep
675 Neural Networks.” *Nature Biotechnology* 36 (10): 983–87.
- 676 Ramms, Bastian, and Philip L. S. M. Gordts. 2018. “Apolipoprotein C-III in Triglyceride-Rich
677 Lipoprotein Metabolism.” *Current Opinion in Lipidology* 29 (3): 171–79.
- 678 Ritchie, Matthew E., Belinda Phipson, Di Wu, Yifang Hu, Charity W. Law, Wei Shi, and Gordon K.
679 Smyth. 2015. “Limma Powers Differential Expression Analyses for RNA-Sequencing and
680 Microarray Studies.” *Nucleic Acids Research* 43 (7): e47.

- 681 Robinson, Mark D., Davis J. McCarthy, and Gordon K. Smyth. 2010. “edgeR: A Bioconductor
682 Package for Differential Expression Analysis of Digital Gene Expression Data.” *Bioinformatics*
683 26 (1): 139–40.
- 684 Sanjurjo, Lucía, Núria Amézaga, Gemma Aran, Mar Naranjo-Gómez, Lilibeth Arias, Carolina
685 Armengol, Francesc E. Borràs, and Maria-Rosa Sarrias. 2015. “The Human CD5L/AIM-CD36
686 Axis: A Novel Autophagy Inducer in Macrophages That Modulates Inflammatory Responses.”
687 *Autophagy*. <https://doi.org/10.1080/15548627.2015.1017183>.
- 688 Schlüter, Agatha, Juan Sandoval, Stéphane Fourcade, Angel Díaz-Lagares, Montserrat Ruiz, Patrizia
689 Casaccia, Manel Esteller, and Aurora Pujol. 2018. “Epigenomic Signature of
690 Adrenoleukodystrophy Predicts Compromised Oligodendrocyte Differentiation.” *Brain*
691 *Pathology* 28 (6): 902–19.
- 692 Sengelaub, Caitlin A., Kristina Navrazhina, Jason B. Ross, Nils Halberg, and Sohail F. Tavazoie.
693 2016. “PTPRN 2 and PLC β 1 Promote Metastatic Breast Cancer Cell Migration through PI
694 (4,5)P 2 -dependent Actin Remodeling.” *The EMBO Journal*.
695 <https://doi.org/10.15252/embj.201591973>.
- 696 Sittaramane, Vinoth, and Anand Chandrasekhar. 2008. “Expression of Unconventional Myosin
697 Genes during Neuronal Development in Zebrafish.” *Gene Expression Patterns: GEP* 8 (3): 161–
698 70.
- 699 Turer, Emre, William McAlpine, Kuan-Wen Wang, Tianshi Lu, Xiaohong Li, Miao Tang, Xiaoming
700 Zhan, et al. 2017. “Creatine Maintains Intestinal Homeostasis and Protects against Colitis.”
701 *Proceedings of the National Academy of Sciences of the United States of America* 114 (7):
702 E1273–81.
- 703 V Willrich, Maria Alice, Erin J. Kaleta, Sandra C. Bryant, Grant M. Spears, Laura J. Train, Sandra E.
704 Peterson, Vanda A. Lennon, Stephen L. Kopecky, and Linnea M. Baudhuin. 2018. “Genetic
705 Variation in Statin Intolerance and a Possible Protective Role for UGT1A1.” *Pharmacogenomics*
706 19 (2): 83–94.
- 707 Wang, Chao, Nir Yosef, Jellert Gaublumme, Chuan Wu, Youjin Lee, Clary B. Clish, Jim Kaminski,
708 et al. 2015. “CD5L/AIM Regulates Lipid Biosynthesis and Restrains Th17 Cell Pathogenicity.”
709 *Cell* 163 (6): 1413–27.
- 710 Waugh, Mark G. 2016. “Chromosomal Instability and Phosphoinositide Pathway Gene Signatures in
711 Glioblastoma Multiforme.” *Molecular Neurobiology* 53 (1): 621–30.
- 712 Wiesinger, Christoph, Florian S. Eichler, and Johannes Berger. 2015. “The Genetic Landscape of X-
713 Linked Adrenoleukodystrophy: Inheritance, Mutations, Modifier Genes, and Diagnosis.” *The*
714 *Application of Clinical Genetics* 8 (May): 109–21.
- 715 Zhang, Rong, Yang Liu, Yang Xue, Yinan Wang, Xinwen Wang, Songtao Shi, Tao Cai, and
716 Qintao Wang. 2014. “Treatment of Lipoid Proteinosis due to the p.C220G Mutation in ECM1, a
717 Major Allele in Chinese Patients.” *Journal of Translational Medicine*.
718 <https://doi.org/10.1186/1479-5876-12-85>.
- 719
720

721 **6 Conflict of Interest**

722 *The authors declare that the research was conducted in the absence of any commercial or financial*
723 *relationships that could be construed as a potential conflict of interest.*

724 **7 Author Contributions**

725 Project was conceived by SK, ME, and CvK. Patients were recruited by ME. Paper written by PAR,
726 FvdK, WW, CvK, SK, and ME. Data analysis for this work includes: WGS data analysis by PAR;

727 RNA-seq analysis by FvdK; Lipidomics analysis by FV; Proteomics analysis by AU and PL; DNA
728 methylation analysis by DL and MK; statistical analysis by EG, SM, PDM, and AHCvK.

729 **8 Funding**

730

731 PAR was supported by BC Children's Hospital Research Institute Graduate Studentship. ME is
732 supported by an NWO/ZonMW Vidi grant (016.196.310). CvK is supported by a Stichting Metakids
733 grant. The research project was funded by Stichting Steun Emma Kinderziekenhuis (Amsterdam;
734 project number WAR 2016-014).

735

736 **9 Acknowledgments**

737

738 We gratefully acknowledge the patients and families for their participation in this study, and the
739 clinicians and colleagues for their expert management.

740 **10 Data Availability Statement**

741

742 The datasets generated and analyzed for this study can be found in the Zenodo repository: DOI-
743 10.5281/zenodo.3698292 ; URL-<https://zenodo.org/record/3698292#.XmrVHi0ZNTY>.

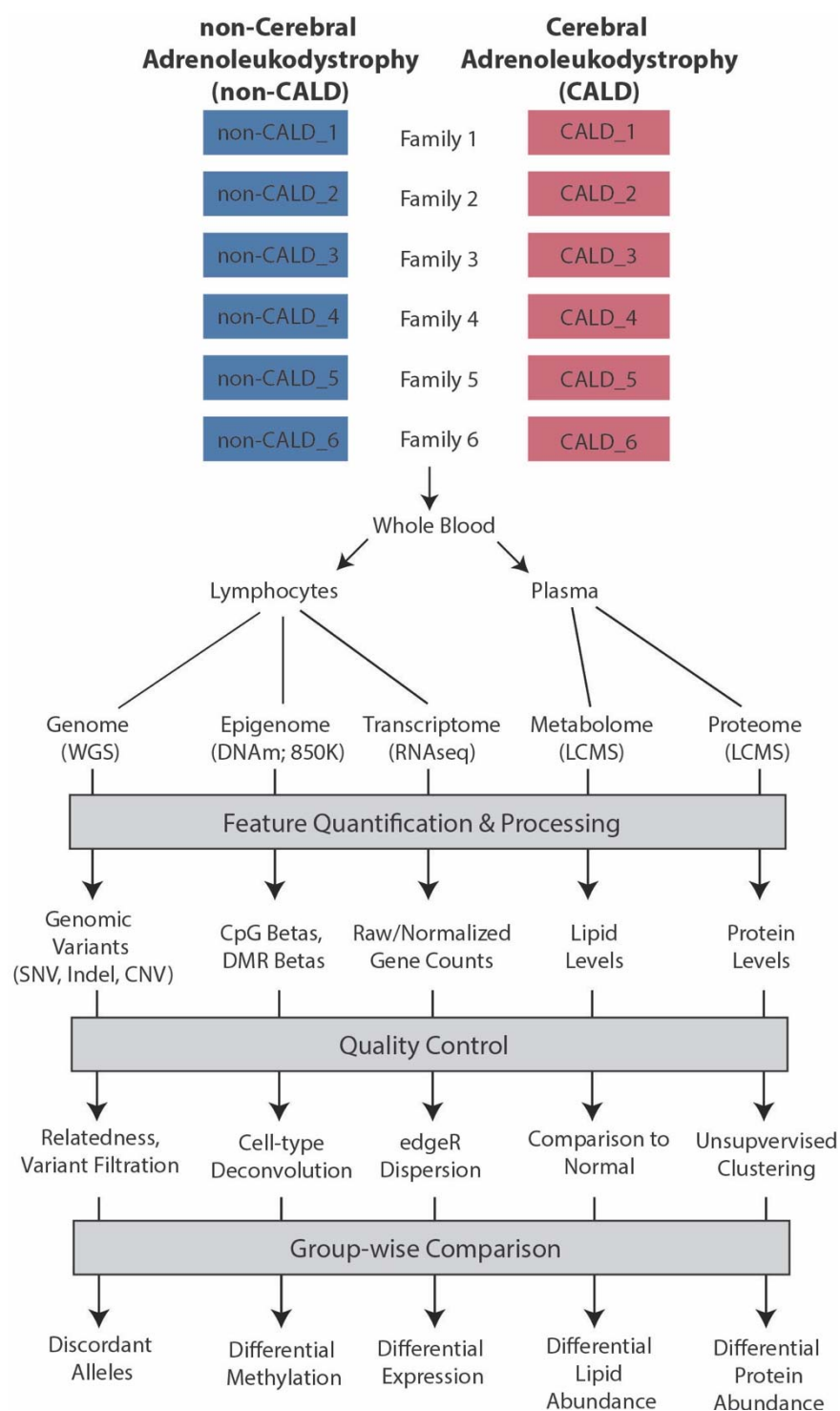
744 Analysis code for this project can be found in the github repository: [https://github.com/Phillip-a-
745 richmond/ALD_Modifier_Project](https://github.com/Phillip-a-richmond/ALD_Modifier_Project).

746

747 **11 Figures**

748

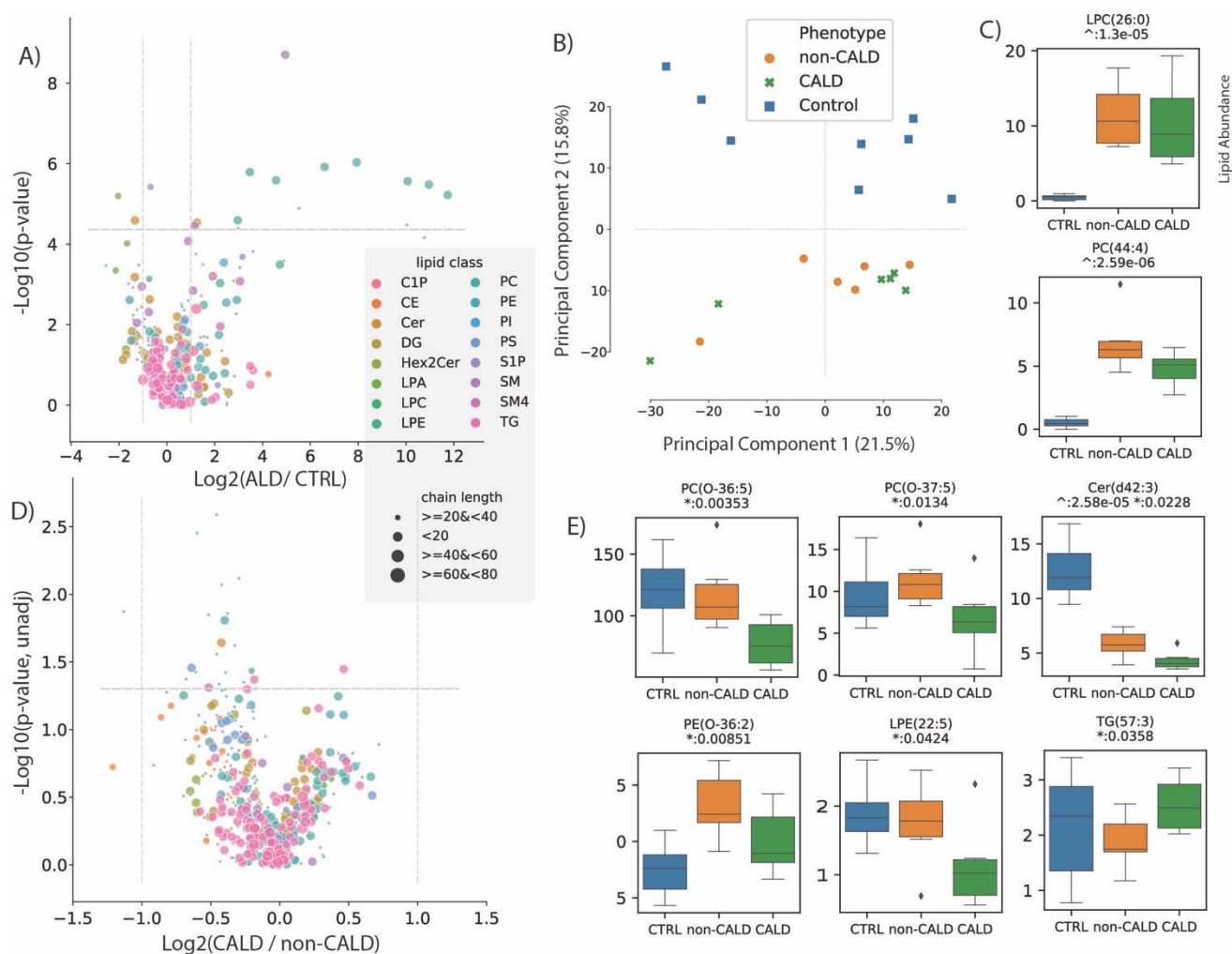
749



750
751
752
753
754
755
756
757
758
759

Figure 1 – Project overview

An overview of the data and processes involved in the project including samples from two brothers across unrelated families, blood isolated into lymphocytes and plasma, and then profiling with five omics technologies including WGS for the genome, DNA methylation (DNAm) via the 850K EPIC microarray, transcriptome profiling with RNA sequencing (RNA-seq), metabolome profiling with liquid chromatography mass spectrometry (LCMS), and protein profiling with LCMS. These data are then taken through feature quantification/processing, quality control metrics, and group-wise comparison through univariate modeling.



760

761

Figure 2 – Lipidomic analysis of ALD

762

The univariate analysis comparing the lipid abundances between control vs ALD, and CALD vs non-

763

CALD is depicted. A) Volcano plot showing the log₂ fold change between ALD and control (CTRL)

764

samples for all of the measured metabolites within the LCMS assay, versus the -log₁₀ transformed

765

adjusted p-value. B) Principal component analysis plot showing the first two principal components

766

which can discriminate between control (blue) and ALD (orange:non-CALD, green:CALD) samples.

767

C) Boxplots showing the abundances of a known marker for ALD, LPC(26:0), and another lipid

768

differentially abundant between ALD and control samples. Values are lipid abundances measured on

769

LCMS. D) Volcano plot showing the log₂ fold change between CALD and non-CALD samples

770

versus the -log₁₀ transformed p-value. E) Boxplots for lipids different between CALD and non-

771

CALD before p-value correction. For A) and D), the lipids are coloured according to their assigned

772

class and their size corresponds to the lipid chain length. For boxplots: \wedge represents unadjusted p-

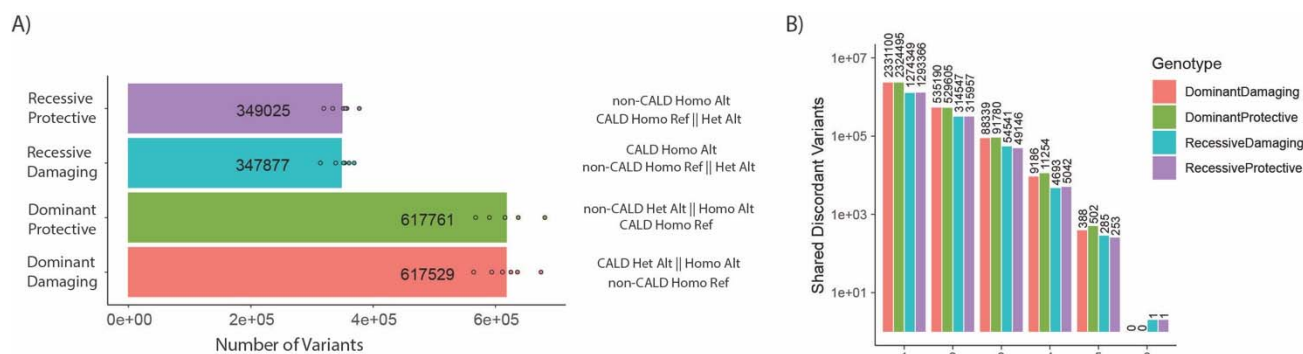
773

values of comparison between ALD and control, * represents unadjusted p-values of comparison

774

between CALD and non-CALD.

775

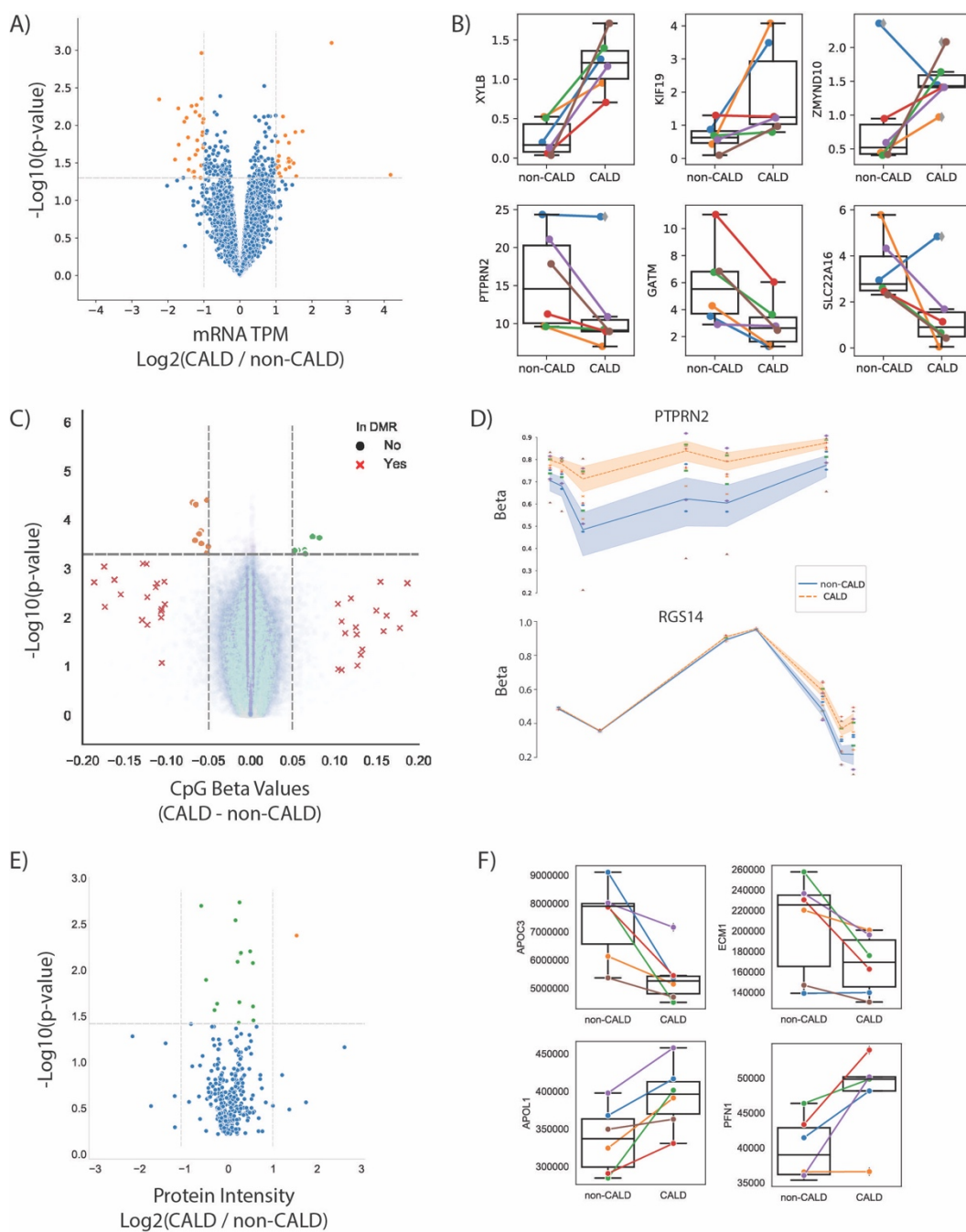


776
777

778 **Figure 3 – Discordant genotype analysis**

779 A) Number of discordant genotypes in each category for each of the 6 families, with description of
 780 genotypes for non-CALD and CALD pairs per genotypic category (per-category means are
 781 displayed). B) Upon intersection of discordant genotypes, the number of variants which exist within
 782 any intersection with set sizes of 1-6, meaning the set size of 6 is the intersection of all families, and
 783 a set of 1 are discordant variants only found in one family.

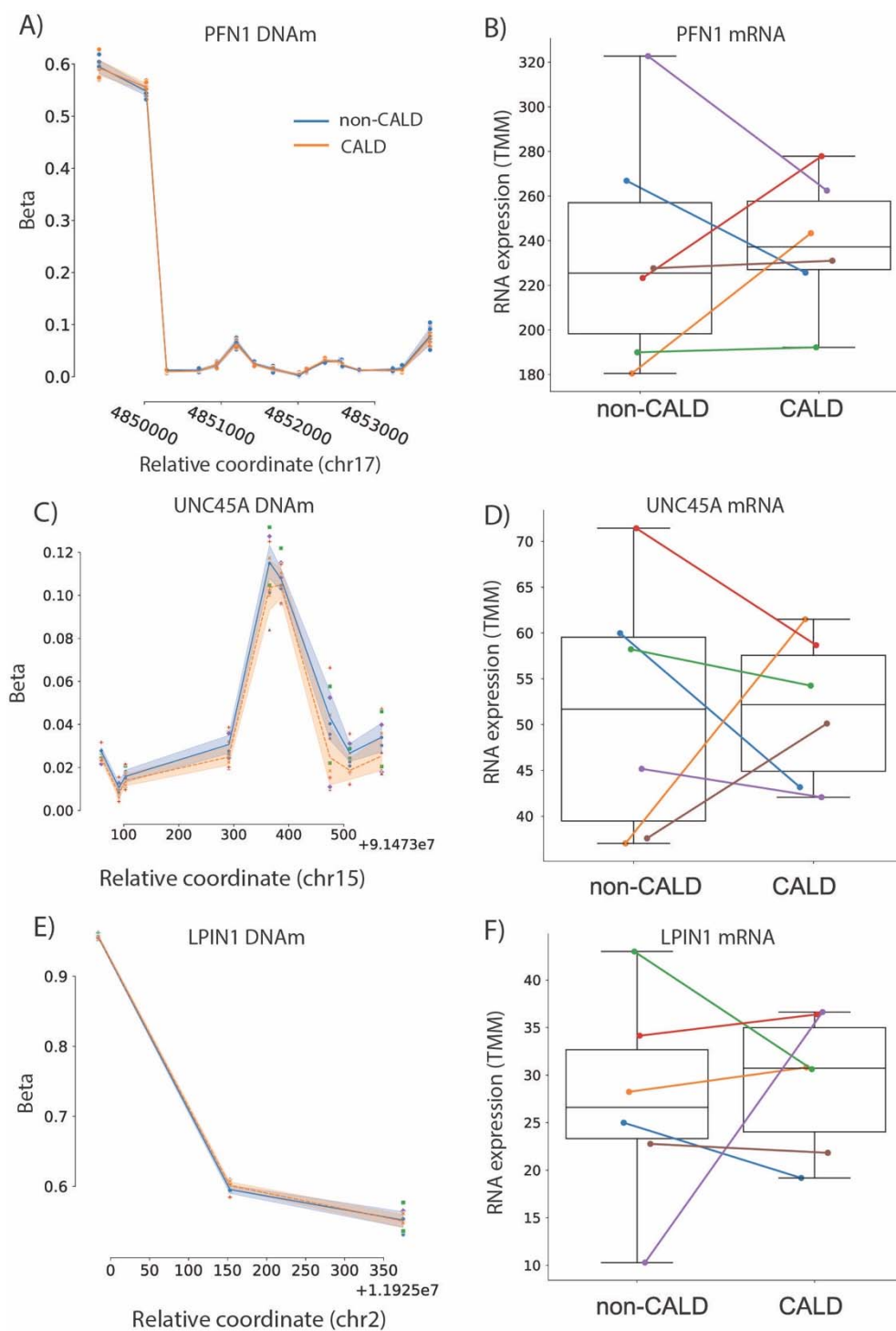
784



785

786 **Figure 4 – Multi-omic analysis**

787 A) Volcano plot showing p-value and log2 fold change of gene expression from RNA-seq.
 788 Significant genes at $p < 0.05$ (orange dots), non-significant genes (blue dots). B) Selected genes
 789 plotted as normalized RNA-seq values with boxplots for each group where each line/point is
 790 coloured by family. C) Volcano plot of DNA methylation over CpG probes from EPIC array, with
 791 non-significant ($p > 0.05$) DMR probes (blue dots), significant CpGs at the DMR level (red Xs),
 792 higher methylated non-CALD probes (orange dots), and higher methylated CALD probes (green
 793 dots). D) DNAm over two significant DMRs within *PTPRN2* and *RGS14*, points coloured by family
 794 and lines coloured by phenotype, with shading denoting inner quartile range. E) Volcano plot of
 795 protein levels from LCMS with non significant proteins (blue), significant proteins with log2 fold-
 796 change (CALD / non-CALD) of -1 to 1 (green), and log2 fold-change greater than 1 (orange). F)
 797 Selected proteins which passed the p-value threshold of 0.05.



798
799
800
801
802
803
804
805
806
807

Figure 5 – Testing previously suggested markers

DNA methylation and mRNA abundance for PFN1 (A,B), UNC45A (C,D), and LPIN1 (E,F). DNA methylation is shown for all CpGs associated to the listed genes, with the non-CALD mean methylation shown as a blue line with standard error shading, and the orange dashed line showing the mean methylation of CALD. Individual points are shown and colored by family. RNA expression is shown as a boxplot for non-CALD and CALD phenotype groups, with individual families labeled with family 1 as blue, family 2 as orange, family 3 as green, family 4 as red, family 5 as purple, and family 6 as brown.

808

809 **12 Tables**

810

811

Family	Age	Phenotype	ABCD1 Mutation	ELOVL1 (A>G)	CYP4F2 (C>T)	APOE rs429358	APOE rs7412	APOE Genotype
1	28	CALD	c.1390C>T	G/G	T/T	T/T	C/C	ε3 / ε3
	28	non-CALD		G/G	C/T	T/T	C/C	ε3 / ε3
2	30	CALD	c.1899delC	A/G	C/C	T/T	C/C	ε3 / ε3
	30	non-CALD		A/G	C/C	T/T	C/C	ε3 / ε3
3	38	CALD	c.1992-2A>G	A/G	C/C	T/T	C/C	ε3 / ε3
3	36	non-CALD		A/G	C/C	T/T	C/C	ε3 / ε3
4	6	CALD	c.659T>C	A/A	C/T	T/T	T/T	ε2 / ε2
	8	non-CALD		A/A	C/T	T/T	T/T	ε2 / ε2
5	16	CALD	c.1866-2A>T	A/G	C/C	T/C	C/C	ε3 / ε4
	18	non-CALD		G/G	C/C	T/C	C/C	ε3 / ε4
6	27	CALD	c.892G>A	A/A	C/C	T/C	C/C	ε3 / ε4
	25	non-CALD		A/G	C/C	T/C	C/C	ε3 / ε4

812

813 **Table 1 - Summary of patients within ALD cohort**

814 The family number, patient ID, age at sample collection, ALD phenotype, *ABCD1* variant, and
815 genotypes for previously associated modifier alleles for all patients within the cohort.
816

817

818

Comparison	DNAm & RNA	DNAm & Protein	RNA & Protein	DNAm & RNA & Protein
all_families	PTPRN2(↑ - ↓)	□	□	□
wo_fam_1	□	□	□	□
wo_fam_2	□	□	□	□
wo_fam_3	HLA-DQB1(↓-↓), IL5RA(↓ - ↑),KIF19(↑ - ↓)	□	□	□
wo_fam_4	□	□	□	□
wo_fam_5	□	□	ICAM1(↑ - ↓), APOL1(↑ - ↑), CD14(↑ - ↑)	□
wo_fam_6	□	□	JCHAIN(↑ - ↑)	□

819

820 **Table 2 - Intersections of significant hits from multiple platforms**

821 For each comparison including all families, and each possible 5x5 comparison between CALD and
822 non-CALD, the significant hits (p-value < 0.05 before multiple testing correction) from DNA
823 methylation (DNAm), RNA-sequencing (RNA), and protein LCMS (Protein) were intersected. (↑
824 means up-regulated/higher for CALD, ↓ means lower in CALD).
825

826

827

828

829
830
831
832
833