**Stratified computational meta-analysis of 2213 acute myeloid leukemia patients reveals age- and sex-dependent gene expression signatures**

Raeuf Roushangar [1, 2], George I. Mias [1, 2*]

[1] Department of Biochemistry and Molecular Biology,

[2] Institute for Quantitative Health Science and Engineering,

Michigan State University, East Lansing MI 48824, USA


*Corresponding author

E-mail: gmias@msu.edu (GM)

In 2018 alone, an estimated 20,000 new acute myeloid leukemia (AML) patients were diagnosed, in the United States, and over 10,000 of them are expected to die from the disease. AML is primarily diagnosed among the elderly (median 68 years old at diagnosis). Prognoses have significantly improved for younger patients, but in patients older than 60 years old as much as 70% of patients will die within a year of diagnosis. In this study, we conducted stratified computational meta-analysis of 2,213 acute myeloid leukemia patients compared to 548 healthy individuals, using curated publicly available data. We carried out analysis of variance of normalized batch corrected data, including considerations for disease, age, tissue and sex. We identified 974 differentially expressed probe sets and 4 significant pathways associated with AML. Additionally, we identified 70 sex- and 375 age-related probe set expression signatures relevant to AML. Finally, we used a machine learning model (KNN model) to classify AML patients

24  compared to healthy individuals with 90+% achieved accuracy. Overall our

25  findings provide a new reanalysis of public datasets, that enabled the

26  identification of potential new gene sets relevant to AML that can potentially be

27  used in future experiments and possible stratified disease diagnostics.

28

29

30

31  **INTRODUCTION**

32  Acute myeloid leukemia (AML) is a heterogeneous malignant disease of the

33  hematopoietic system myeloid cell lineage. AML is best characterized by the

34  terminal differentiation in normal blood cells and excessive production and

35  release of cells at various stages of incomplete maturation (leukemia cells). As a

36  result of this faster than normal and uncontrolled growth of leukemia cells,

37  healthy myeloid precursors involved in hematopoiesis are suppressed, and

38  ultimately, can soar to death within months from diagnosis if untreated[1,2]. AML

39  accounts for 70% of myeloid leukemia and nearly 80% of acute leukemia cases,

40  making it the most common form of both myeloid and acute leukemia[2,3]. The

41  number of new AML cases is increasing each year – in 2018 alone, there have

42  been an estimated about 20,000 new diagnosed AML patients, over 10,000 of

43  them will die from the disease[4].

44

45  According to the 2016 World Health Organization (WHO) newly revised myeloid

46  neoplasms and acute leukemia classification system[5], AML prognosis criteria for

2

47  classification is highly dependent on the presence of chromosomal abnormalities,

48  including chromosomal deletions, duplications, translocations, inversions, and

49  gene fusion. Mostly, AML is diagnosed through microscopic, cytogenetics, and

50  molecular genetic analyses of patients' blood and/or bone marrow samples.

51  Microscopic examination is used to detect distinctive features (e.g. Auer rods) in

52  cell morphology, cytogenetic analysis to identify chromosomal structural

53  aberrations (e.g., t(8;21), inv(16), t(16;16), or t(9;11)), and molecular genetic

54  analysis to identify gene fusion (e.g., RUNX1-RUNX1T1 and CBFB-MYH11), and

55  mutations in genes frequently mutated in AML (e.g., NPM1, CEBPA, RUNX1,

56  FLT3)[6-8]. These cytogenetic and molecular genetic analyses are used to identify

57  prognosis markers that can be used to classify AML patients into three risk

58  categories: favorable, intermediate, and unfavorable. The largest group of AML

59  patients (almost 50%) however, present normal karyotype and lack genetic

60  abnormalities[7-10]. These patients are classified as intermediate risk, and often

61  have heterogeneous clinical outcome with standard therapy with risk of AML

62  relapse[11].

63

64  Additionally, AML prognosis worsens as age increases, and older patients

65  respond less to current treatments with poorer clinical outcomes than their

66  younger counterparts[12,13]. AML can occur in people of all ages but is primarily

67  diagnosed among the elderly (>60 years), with a median age of 68 year at

68  diagnosis[4]. Recent advances in AML biology expanded our understanding of its

69  complex genetic landscape and led to significant improvement in prognoses and

70  therapeutic strategy for younger patients[13,14]. However, in patients older than 60

71  years old, prognoses remain grim and therapeutic strategy has been nearly the

72  same for more than 30 years[2,6,13-15]. Approximately 70% of AML patients 65

73  years of age or older die within a year following diagnosis[16]. While it is apparent

74  that the nature of AML changes with age, still little is known about the extent of

75  these associations and how they vary with patient's age[14,17,18]. Taking into

76  consideration age considerations in the identification of changes in AML global

77  gene expression can lead to improved early diagnosis and improvement in

78  treatment approaches for elderly patients. Further complicating, AML has

79  multiple driver mutations and competing clones that evolve over time, making it a

80  very dynamic disease[19,20]

81

82  Multiple gene expression analyses of AML have been carried out, 25 of which

83  these have been systematically compared by Miller and Stamatoyannopoulos[21],

84  who analyzed information on 4918 genes, and identified 25 genes reported

85  across multiple, with potential prognostic features. In this study, we performed

86  comprehensive gene expression meta-analysis of 2213 acute myeloid leukemia

87  patients and 548 healthy subjects using 34 publicly available gene expression

88  microarray datasets (following strict inclusion criteria) to identify disease, sex-

89  and age-related gene expression changes associated with AML. We identified

90  sex- and age-related gene expression signatures that show similar alteration in

91  gene expression levels and associated signaling pathways in AML and have

92  used our results (gene sets) to predict AML or healthy status. We believe that our

4

93   results may lead to improved AML early detection and diagnostic testing with

94   target genes, which collectively can potentially serve as sex- and age-dependent

95   biomarkers for AML prognosis compared to healthy, as well as the identification

96   of new treatment targets with mechanisms of action different from those used in

97   conventional chemotherapy

98

99   **RESULTS**

100   **Data curation and gene expression preprocessing.**

101   We searched the Gene Expression Omnibus (GEO) public repository, based on

102   our systematic workflow and inclusion criteria, Fig. 1a-b. Overall, 2,132 datasets

103   were screened, 643 selected (577 were excluded as non-Affymetrix, various

104   platform arrays). From the 66 remaining, 34 studies were excluded due to lack of

105   metadata, non-peripheral blood and non-bone marrow tissues, cell line or cell-

106   type specific, treated subjects). After this curation we obtained 34 age-annotated

107   gene expression datasets from 32 different studies covering 2,213 AML patients

108   and 548 healthy individuals. The sets were re-analyzed, starting from raw data,

109   to perform a gene expression analysis of variance and functional pathway

110   enrichment analysis (see online Methods). Table 1 provides a description on

111   each dataset with a sub-table summary of all curated data used in our current

112   study. After pre-processing each individual data set separately, Fig. 1b, we

113   performed the statistical analysis on 44,754 probe sets which were common

114   across all samples (Affymetrix expression array data).

115

116    **Classification of missing metadata annotation.**

117    Following the data curation step, 805 arrays (802 AML and 3 healthy) of 2,761

118    curated data were found to be missing sex annotation, and 737 arrays (all AML

119    patients) were missing information regarding the sample source (i.e. tissue,

120    either bone marrow [BM] or peripheral blood [PB] annotation). To predict the

121    missing sex and sample source meta-data, we trained and validated various

122    machine learning supervised models, including logistic regression (LR)

123    classification models. The prediction of missing annotations for these arrays was

124    essential in our study, to increase the sample size, and statistical power[22]. The

125    models were trained and verified using our annotated preprocessed expression

126    data. Model training, parameters used in training, validation for this analysis are

127    discussed in the Methods. Results from model training and predictions, including

128    confusion matrix, model accuracy, and error can be viewed in Supplementary

129    Table S1 online and results from classification for missing annotation are

130    presented in Supplementary files 1 and 2 for sample source and sex annotations

131    respectively.

132

133    **Batch correction**

134    Our pre-processed data, AML and healthy, were processed using a "dataset-

135    wise" batch effect correction approach. The datasets used in this study did not

136    include within-study healthy controls, which would limit analysis of variance, and

137    particularly the ability to separate biological from batch effects. To address this,

138    we implemented an iterative batch effect correction approach, essentially

6

139 employing a weight-based method for correcting batch effects. Assuming the

140 batch effects due to each data set is a function of the number of samples in the

141 data set (weight), normalizing sets of unevenly sized datasets may lead to

142 unbalanced batch correction. We used 5 additional datasets as a reference set,

143 which we refer to as "covariate" hereafter. Each of the covariate reference

144 datasets included within study healthy controls. All 5 datasets together consisted

145 of a total 613 arrays (455 AML and 158 healthy) (Table 2), and pre-processed

146 exactly as our curated data sets. These were used together with each of the

147 remaining datasets to batch correct each dataset with respect the covariate

148 reference using ComBat[23]. After this dataset-wise correction, the 5 covariate

149 reference datasets were removed, and our expression data were clustered using

150 principal component analysis (PCA) to visually examine the effect of covariate

151 reference datasets on distributing the batch weight during batch correction. The

152 batch effect correction results were then compared to clustering results prior to

153 batch effect correction (Supplementary Fig. 1)

154

155 **Analysis 1: Gene expression meta-analysis and enrichment analysis of**

156 **AML disease state compared to healthy individuals**

157

158 **Gene expression meta-analysis of AML disease state.**

159 Following batch correction, we performed an analysis of differential expression

160 (DE) on 34 data sets including 2213 AML patients and 548 healthy controls.

161 Analysis of Variance (ANOVA)[24-26] was performed according to a linear model

162    (see method section **Meta-analysis**), including factors for age, sample source (

163    adjust for differences in tissue between AML and healthy), and sex, as well as

164    binary interactions thereof. In the analysis we used probe sets to avoid

165    assumptions on averaging over multiple probe sets corresponding the same

166    gene symbol. 974 Statistically significant differentially expressed probe sets

167    (DEPS) (with genes corresponding to 964 unique gene symbols) for AML versus

168    healthy were selected based on a Bonferroni[27] adjusted p-value < 0.01

169    (accounting for multiple hypothesis testing), in conjunction with a two-tailed 5%

170    quantile selection[28] based on the mean difference distribution between AML-

171    healthy group comparisons (post-hoc analyses using Tukey's Honestly

172    Significant Difference (HSD). The heatmap (Fig. 2a) shows the hierarchical

173    clustering of genee expression from the 974 DEPS, including 487 up- and 487

174    down-regulated with respect to AML as compared to healthy. From this analysis,

175    WT1 (Wilms tumor 1) with mean difference of 0.26 and adjusted p-value <

176    $4.11 \times 10^{-11}$ was the most DE up-regulated gene while CRISP3 (cysteine-rich

177    secretory protein 3) with mean difference of -0.52 and adjusted p-value <

178    $4.11 \times 10^{-11}$ was the least DE gene. Figure 2b shows the top 10 up- and down-

179    regulated DEPS with corresponding gene symbols, that resulted from this

180    analysis (also listed in Table 2, including mean difference and Bonferroni p-

181    adjusted values from post-hoc analysis using Tukey's HSD tests). The entire list

182    of all 974 DEPS can be found as Supplementary Table S2 online.

183

184

185 **(ii) Gene enrichment analysis AML disease state DEPS.**

186 To identify signaling pathways associated DEPS in AML, gene enrichment

187 analysis was performed on all 974 DEPS combined. Pathway over-

188 representation analysis in Kyoto Encyclopedia of Genes and Genomes

189 (KEGG)[29-31] signaling pathways, and Gene Ontology (GO) term[32,33] were carried

190 out using the Database for Annotation, Visualization and Integrated Discovery

191 (DAVID)[34,35]. Four KEGG signaling pathways were identified as enriched

192 (Benjamini and Hochberg[36] adjusted p-value < 0.05), including Hematopoietic

193 cell lineage, Cell cycle, p53 signaling pathway, and Transcriptional

194 misregulation in cancer. The 4 KEGG signaling pathways are summarized in

195 Table 3 (see also Supplementary Fig. 2a-d), including unadjusted p-values and

196 Benjamini and Hochberg[36] adjusted p-values. 56 DEPS including 27 up- and 29

197 down-regulated (Fig. 2c) were associated these signaling pathways, and the

198 heatmap of their mean differences is shown in Fig. 2d. From our gene

199 enrichment analysis for overrepresented biological GO terms, 21 GO terms were

200 statistically significant with 727 DE unique identities (335 up- and 392 down-

201 regulated). GO terms included protein and microtubule binding for the molecular

202 function (MF) category, inflammatory and immune responses, mitotic nuclear

203 division, and cell proliferation response for the biological process (BP) category,

204 and finally, cytoplasm, extracellular exosome, cytosol, extracellular space,

205 integral component of plasma membrane immune response, and others, for the

206 cellular component (CC) category (Fig. 2e). The entire list of our enrichment

207 analysis results (statistically significant over-representation in KEGG and GO

208 terms) can be found as Supplementary Table S3 online.

209

210 **Analysis 2. Gene expression meta-analysis and enrichment analysis of sex-**

211 **and age-related DEPS in AML.**

212 Further analysis of gene expression and pathways enrichment were conducted in

213 order to characterize sex- and age-specific gene expression changes in AML

214 patients compared to healthy individuals, (i) **Analysis 2a:** "**Sex-relevance**

215 **differential gene expression meta-analysis and associated signaling**

216 **pathways in AML",** and (ii) **Analysis 2b: "Age-dependent differential gene**

217 **expression meta-analysis and associated signaling pathways in AML"**. We

218 used the same filtering criteria in both analyses as those used in analysis 1 for

219 significant DEPS and signaling pathways between AML patients and healthy

220 controls. In addition, DEPS were regarded as statistically significantly (up- or

221 down-regulated) for each factor, sex and age, if they displayed Bonferroni

222 adjusted p-value from Tukey's HSD $< 2.2 \times 10^{-7}$ (=0.01/44,754 probe sets tested).

223

224 **Analysis 2a. Sex-relevance differential gene expression meta-analysis and**

225 **associated signaling pathways in AML.**

226 Gene expression meta-analysis was also used to identify DEPS that show sex

227 differences between male AML patients as compared to female AML patients.

228 266 DEPS were regarded statistically significant (p-value $< 2.2 \times 10^{-7}$). A list of all

229 266 DEPS (including whether higher in either males or females, gene title and

10

230   symbol, male-female mean difference, and Bonferroni corrected p-value) can be

231   found as Supplementary Table S3 online. 70 DEPS were found to overlap

232   between analysis 1 (AML disease state) and analysis 2 (Sex-relevance in AML).

233   Figure 3a shows these 70 DEPS with gene symbol annotations, and their mean

234   difference values in the heatmap, which displays differences in significance for a

235   common DEPS in both analyses 1 and 2. Figure 3b shows the hierarchical

236   clustering of the 70 DEPS (rows) on sex and disease state of all 2,213 AML and

237   548 healthy subjects (columns) indicated by color bars above the heatmap. The

238   top 10 DEPS higher in either males or females from this analysis are shown in

239   Figure 3c.

240

241   For enrichment analysis, we searched for common intersections in KEGG

242   pathways and GO terms between the sex meta-analysis and the 974 DE probe

243   sets from disease state in AML meta-analysis. Sex-relevant DEPS were found in

244   3 different signaling pathways, including, genes higher expressed in males FLT3

245   and CD34 in Hematopoietic cell lineage, FLT3 in Transcriptional misregulation in

246   cancer 1, and PMAIP1 in p53 signaling pathway 1, and MS4A1 was higher in

247   females and found in Hematopoietic cell lineage pathway (Table 3). Figure 3d

248   shows GO analysis results, where 15 overrepresented biological GO terms were

249   overlapped, including terms for extracellular space, immune response, protein

250   binding, spindle, and midbody. The entire list of our enrichment analysis

251   (statistically significant KEGG and GO terms) can be found as Supplementary

252   Table S4.

11

253

**Analysis 2b. Age-dependent differential gene expression meta-analysis and associated signaling pathways in AML.**

The subjects were binned in 8 age-groups: 0-19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, and 80-100 years old. From this meta-analysis, 1395 unique probe sets across all age-groups were identified as statistically significant (Bonferroni adjusted p-value < $2.2 \times 10^{-7}$) (Supplementary Table S5). From these 375 unique DEPS (372 unique gene symbols) were found to overlap with the 974 DEPS probe sets from our AML disease state meta-analysis, accounting for an overall 1400 binary comparisons between the multiple age groups deemed statistically significant, based on Tukey HSD tests between age-group pairs. The entire list of 1400 identified pairwise differences between age groups and associated probe set/gene information can be found as Supplementary Table S6 online. The top 10 up- and down- regulated DEPS (labeled with gene symbols) from this analysis are shown in Fig. 4a. Additionally, Fig. 4b shows 75 DEPS with gene symbols identified to have appeared specifically in one age-group comparison. Utilizing results for KEGG analysis for signaling pathways from analysis 1, Fig. 4c shows 17 DE genes identified in all 4 KEGG pathways according to age groups (see also Table 4).

272

To investigate further the progression with age, pairwise correlations between age-groups were computed. The 0-19 age-group was used as a common comparison reference with respect to other groups. Using this 0-19 group as a

12

276    baseline, Figure 4d shows the mean difference of 25 DEPS with respect to the 0-

277    19 baseline across all other groups. The mean difference values between AML

278    and healthy are shown in the right-most column of Fig. 4a, b and d for reference.

279

280    **AML Classification Machine Learning Model**

281    We used the 974 DEPS to train a k-nearest neighbor (KNN) algorithm in

282    ClassificaIO[37]. All 34 datasets (16 AML and 18 healthy) were used for training,

283    and testing was done on all 5 covariate reference datasets, include AML and

284    healthy subjects. The KNN algorithm trained was 98% accurate, and >90%

285    accurate in testing results (see online Methods for parameters and also details in

286    Supplementary File 3).

287

288    **DISCUSSION**

289    In the present study, we aimed to establish, disease sex-linked and age-

290    dependent biomarkers from genes with similar changes in gene expression

291    levels and associated signaling pathways relevant to AML. Utilizing microarray

292    gene expression data and combined with various machine learning models,

293    respectively, our biomarkers were indicative of prognostic signature for AML

294    prediction compared to healthy with 90+% achieved accuracy. We re-analyzed

295    data aggregated from our curation of 34 publicly available microarray gene

296    expression datasets covering 2213 AML patients and 548 healthy individuals to

297    identify changes in AML gene expression associated with disease state (AML

13

298    compared to healthy), sex-linked (male compared to female), and age-dependent

299    (across age-groups compared to baseline).

300        We performed 3 differential probe set (gene) expression and gene enrichment

301    analyses, as discussed below. We note here that our study identified multiple

302    potentially significant DEPS, with age and sex related differences associated with

303    AML. While our findings may generate further hypothesis-driven investigations,

304    we need to also identify the study's limitations: primarily the analysis of AML and

305    healthy subjects involved bone-marrow and blood samples respectively in each

306    disease group. We tried to account for this utilizing tissue as an effect in our

307    linear model, and including multiple interactions. Other limitations include an

308    unbalanced AML/healthy ratio, as well as the lack of in-study healthy controls. To

309    address these we attempted to account for batch effects using a dataset-wise

310    iterative batch correction transformation, as discussed in the methods. Finally,

311    we also included binary interactions between the factors in the analysis to

312    account for interaction-related confounding effects.

313

314     *i) Analysis 1: Gene expression meta-analysis and associated signaling*

315    *pathways of AML disease state compared to healthy individuals*, was carried out

316    to identify DEPS in AML disease state. The results from this analysis were then

317    used as baseline indicator for AML disease state. 974 DEPS (487 up- and 487

318    down-regulated) were identified as significantly differentially expressed between

319    AML patients and healthy individuals (Bonferroni adjusted p-value < 0.01).

320    Among these 6 genes are known to be involved in AML functional pathways,

14

321    including 4 up-regulated, JUP, CCNA1, FLT3, PIK3R1, and 2 down-regulated,

322    CD14, CEBPE. The top 10 up- and down-regulated genes from this analysis are

323    listed in Table 2 with their respected Tukey's HSD mean difference and

324    Bonferroni p-adjusted values. As shown in Figure 2b of the top 10 up- and down-

325    regulated DEPS and corresponding gene annotations -- WT1 (Wilms tumor 1)

326    was found to be the most expressed and CRISP3 (cysteine-rich secretory protein

327    3) was the most under-expressed gene. WT1 is a transcriptional regulatory

328    protein essential to cellular development and cell survival, and it has been known

329    to be highly expressed with an oncogenic role in AML[38,39], in agreement with our

330    findings. However, CRISP3's direct role in AML is still under investigation.

331    CRISP3 is a member of the cysteine-rich secretory protein CRISP family with

332    major role in female and male reproductive tract, and is mainly expressed in

333    salivary gland and bone marrow[40]. Recently, 80 genes were reported as

334    "extracellular matrix specific genes" in leukemia, and CRISP3 was among the

335    downregulated DE genes reported[41]. CRISP3 associations with AML merit further

336    investigation.

337

338    The enrichment analysis for GO terms of the 974 DE probe sets (Fig. 2c) results

339    showed 727 identifiers (335 up- and 392 down-regulated) enriched for 21 GO

340    terms. 592 of these (257 up- and 335 down-regulated) were enriched in the

341    cellular component (CC) categories mainly associated with cytoplasm,

342    extracellular exosome, cytosol, and extracellular space. These terms are rather

343    generic, but may still reflect relevance to AML development and progression[42,43].

344    Biological process (BP) category, GO terms included inflammatory and immune

345    responses, and cell proliferation, which are expected as AML is characterized by

346    terminal differentiation of normal blood cells and excessive proliferation and

347    release of abnormally differentiated myeloid cells, and likely affects many

348    biological processes associated to the immune system. The four statistically

349    significant KEGG pathways identified in the pathway enrichment analysis

350    encompassed 56 DEPS (Table 3). Transcriptional misregulation in cancer was

351    the most up-regulated pathway in AML (13 up-regulated DE genes, while

352    Hematopoietic cell lineage, and Cell cycle pathways were mostly down-

353    regulated, and the p53 signaling pathway was balanced in terms of

354    up/downregulated DE genes (Fig. 2c). The enriched pathways Fig. 2d shows the

355    mean difference values of the 56 DE pathway-associated genes, including 27

356    genes up- and 29 down-regulated. These KEGG pathways are known to be

357    involved in tumorigenesis. Additionally, the majority of the associated DE genes

358    from AML meta-analysis with the identified signaling pathways are known to be

359    abnormally expressed in AML. These findings are consistent with findings from

360    other studies and our current understanding of AML pathogenesis.

361

362    The DEPS overlap with the 25 genes reported by Miller and

363    Stamatoyannopoulos that were reported in at least 8 studies[21], namely HOXA10,

364    CD34, MEIS1,VCAN, RBPMS and MN1. In terms of the genes reported in the

365    same study for poor progression we also consistently identified as upregulated

366    HOXA10, RBPMS, CD34, GNAI1, CLIP2, DAPK1, GUCY1A3, ANGPT1 and

16

367    FLT3, and as downregulated UGCG. While these are known markers, with

368    consistent expression differences, our additional results need to be investigated

369    further and experimentally validated, including mechanistic considerations.

370

371    *ii) Analysis 2a: Sex-dependent gene expression meta-analysis and associated*

372    *signaling pathways in AML compared to healthy individuals*, was performed to

373    explore the relevance of patients' sex on gene expression and to identify sex-

374    linked genes and associated signaling pathways in AML. A total of 266 DEPS

375    were found statistically significant in this analysis, with 70 found to overlap with

376    the DEPS from Analysis 1 (Fig 3a-b). The top10 up- and down-regulated DE

377    genes with respect to females include (Fig. 3c) – DDX3Y (DEAD-Box Helicase 3

378    Y-Linked), EIF1AY (Eukaryotic Translation Initiation Factor 1A Y-Linked),

379    KDM5D (Lysine Demethylase 5D), RPS4Y1 (Ribosomal Protein S4 Y-Linked 1)

380    with higher expression in males compared to females, and XIST (X Inactive

381    Specific Transcript), TSIX (TSIX Transcript, XIST Antisense RNA), and PRKX

382    (Protein Kinase X-Linked) were as higher in females. These genes are known to

383    be sex-specific and show such differences and sex separation within the AML

384    and the healthy groups respectively (Fig. 3d). The role of these genes as positive

385    controls in studies with AML needs to be investigated further. We also reported

386    sex and AML known genes that were statistically significant in our analysis,

387    including FLT3 and MAL.

388

17

389  iii) *Analysis 2b: Age-dependent gene expression meta-analysis and associated*

390  *signaling pathways in AML compared to healthy individuals*, was carried out to

391  identify common set of age-dependent genes and associated signaling pathways

392  and to explore age-dependent trends in gene expression in AML. The age-

393  dependent meta-analysis in AML using ANOVA, identified 1,395 DEPS

394  (Bonferroni adjusted p-value <0.01). To identify age-related DEPS in AML we

395  overlapped the 1,395 DEPS to our findings of 974 DEPS in AML disease state

396  (Analysis 1) (Fig. 4a), and identified an overlap of 375 DEPS (Bonferroni

397  adjusted p.value <0.01). As shown in Figure 4b, the top 10 most and least DE

398  age-associate genes in AML according to the mean difference values in seven

399  age-groups, including their corresponding values from AML disease state in

400  column "AML - healthy" for comparisons. Interestingly, CRISP3 was among the

401  down regulated genes specifically and involved in this analysis as well,

402  specifically associated with differences in younger age groups, 20 to 49 years of

403  age as compared to 0 to 19 age group. Other genes showing age-specific

404  differences included HOXA3, HOXA5 and HOXA10-HOXA9, which belong to the

405  homeobox genes (HOX) family of transcription factors, essential to embryonic

406  development and hematopoiesis, and associated with chromosomal

407  abnormalities translocation and over-expression in AML[44,45]. Also identified with

408  age-specific DE, was ORM1, which in Analysis 1 was among the top-10 most

409  under-expressed genes, and was also among the 70 DE genes in analysis 2a.

410  ORM1's direct role in AML also merits further investigation, given ORM1

411  involvement in immunosuppression and inflammation[46]. Finally, we have

18

412    identified 75 DEPS that show association with only one age-group, exclusively

413    from all other age-groups, suggestive of potential age-specific differential gene

414    expression signature.

415

416    In summary, our study successfully integrated multiple datasets to perform a

417    study of gene expression in AML, across multiple factors that included disease,

418    sex and age considerations, and identified interesting genes, both known and not

419    previously reported as differentially expressed in each factor. We identified 974

420    DEPS and 4 associated significant pathways involved in AML, and 70 sex- and

421    375 age-related DE signatures. Using the 974 DEPS, a KNN model allowed AML

422    with 91.7% accuracy. We hope that these findings may provide additional

423    relevant targets for further experimental mechanistic studies, and to help identify

424    new markers and therapeutic targets for AML.

425

426 **METHODS**

427 The generalized workflow consisted of five main steps: i) Curation of microarray

428 gene expression data, ii) Preprocessing of raw data files followed by batch effect

429 correction, iii) Predictions of missing annotation data using supervised machine

430 learning, iv) Differential gene expression analysis, and v) Gene enrichment for

431 pathway analysis that includes gene annotation, and finally gene expression-

432 based prediction of AML (Fig. 1a).

433

434 **Gene expression data curation and screening criteria.**

435 Datasets used in this study were selected from the GEO public repository,

436 maintained by the National Center for Biotechnology Information (NCBI)[47]

437 (https://www.ncbi.nlm.nih.gov/geo/). To facilitate speed of search and keep up-to-

438 date with possible new and relevant datasets, as soon as they were released, a

439 Python script was used that utilized functions from the Entrez Utilities from

440 Biopython[48]. We used the script to navigate the GEO records, and download

441 microarray gene expression datasets up to 10/18. We additionally utilized Python

442 packages, including Pandas, NumPy, and Matplotlib for data structure, numerical

443 computing for data processing, and data visualization respectively. We used

444 strict inclusion criteria to maintain consistency in each dataset selection, screen

445 for availability of both raw and meta-data annotation files provided, human

446 samples used from untreated subjects, and that the sample source was from

447 either bone marrow (BM) and/or peripheral blood (PB). Array platform was

448 restricted to Affymetrix, which was found to have the most available data, and to

449  avoid cross-platform normalization issues. Inclusion criteria and the data curation

450  workflow are illustrated in Fig. 1 a-b.

451

452  **Gene expression data sets used in our analysis.**

453  The curation method is summarized in the Supplementary File 4 flowchart and in

454  the Results section. For our analysis we included 34 age-dependent datasets

455  from 32 different studies, 16 included AML and 18 healthy subjects respectively.

456  From the 34 datasets, 32 were produced from Affymetrix GeneChip Human

457  Genome U133 Plus 2.0 (GPL570) and 2 conducted on Affymetrix GeneChip

458  Human Genome U133 Array Set (GPL96 & GPL97) arrays. Table 1 provides

459  detailed information about each data set, including the number of samples used

460  from each dataset, sample tissue source, as well as the total number of AML

461  patients and healthy subjects. Two studies, GSE12417[49] and GSE37642[50-53],

462  were originally conducted on two different Affymetrix array types (GPL570, and

463  GPL96 & GPL97), so each was separated into two subgroups and each

464  subgroup was considered as individual dataset in our meta-analysis, data set

465  GSE12417: (i) subgroup 1 included 73 BM and 5 PB samples, and (ii) subgroup

466  2 included 160 BM and 2PB. For dataset GSE37642 (i) subgroup 1 included 140

467  BM and (ii) subgroup 2 422 BM samples (Table 1).

468

469  **Dataset annotation and preprocessing.**

470  Figure 1b outlines the workflow of our preliminary data analysis including

471  preprocessing. For each dataset used in our analysis, raw microarray CEL files

21

472    were downloaded from GEO, metadata was reviewed, and the data was

473    manually curated to guarantee that and each array, which corresponded to either

474    an AML patient or healthy individual, was verified and correctly annotated for

475    sample source (BM or PB), platform technology used, age, sex, and disease

476    state (AML or healthy). Raw CEL files from individual datasets were individually

477    pre-processed using the RMA (Robust Multi-Array Average) algorithm[54-56].

478    Datasets with mixed sample source, i.e both BM and PB, were pre-processed

479    together irrespective of sample source. Preprocessing consisted of correction for

480    background noise using RMA background correction on perfect match (PM) raw

481    intensities, quantile normalization to obtain the same empirical distribution of

482    intensities for each array, median polish summarization of probes into probe sets

483    to estimate gene-level expression value, and logarithm base-2 transformations of

484    gene expression values to facilitate data interpretation (normal distributions) and

485    comparisons between arrays. Additionally, our expression data were first

486    reduced to 44,754 probe sets that are common to and appeared in all data. Data

487    sets were z-score standardized across all probe sets and arrays.

488

489    **Prediction of missing sex- and sample source annotations from curated**

490    **data sets.**

491    805 arrays (802 from AML patients and 3 were healthy subjects) of curated data

492    were not annotated for sex, while 737 arrays (all AML patients) were missing

493    sample source information. Without these metadata, we would have to discard

494    the data, which in turn would limit the statistical power for the study, and our

22

495     ability to correct for biases stemming from individual datasets[22]. To address this,

496     we used supervised machine learning classifiers to predict metadata. For all

497     prediction, we used ClassificaIO[37], a machine learning for classification user

498     interface, which we recently developed, to carry out the machine learning

499     classification analyses utilizing the sklearn package in Python[57]

500

501     To predict sex pre-processed data sets, 1956 arrays (including both healthy and

502     AML), that include 44,754 probe sets and their annotated sex information were

503     used to train logistic regression (LR) classification models, and to predict 805 sex

504     annotations. Additionally, 2024 arrays were used to train for sample source, and

505     the prediction was performed on 737 arrays.

506

507     The supervised machine learning LR classifier we used with the following

508     parameters:

509

510     *random_state = None, shuffle = True, penalty = l2, multi_class = ovr, solver =*

511     *liblinear, max_iter= 100, tol = 0.0001, intercept_scaling = 1.0, verbose = 0,*

512     *n_jobs = 1, C = 1.0, fit_intercept = True, dual = False, warm_start = False,*

513     *class_weight = None*

514

515     The trained models for classification of missing sex and sample source

516     annotation from curated data achieved > 95% classification accuracy with ~ 3-5%

517     classification errors. Confusion matrix details, model accuracy and error for

518    training and testing are presented in Supplementary Table S1 online, and results

519    in Supplementary files 1 and 2. To account for training overfitting, we used 10-

520    fold cross-validation on all 1,956 gene expression data arrays for training and

521    validation.

522

523    **Dataset-wise correction approach for batch effects correction.**

524    Batch correction was done using a dataset-wise correction. Here we refer to the

525    term "dataset-wise correction," to indicate performing batch correction iteratively

526    on one dataset at a time, against a reference set of datasets chosen to account

527    for variability. We used this approach to account for the lack within-study healthy

528    controls in the curated gene expression datasets. To address this issue, we used

529    5 additional datasets the included within-study controls, GEO accessions:

530    GSE107968, GSE68172[58], GSE17054[59], GSE33223[60], and GSE15061[61] (Table

531    1B). We refer to the latter datasets hereafter as "covariate" reference datasets,

532    as they were as the reference datasets in the batch correction. Our approach

533    aimed to balance/distribute the weight of batch effects exerted by each dataset,

534    as this is dependent on the number of observations within a given dataset.

535    Combined, the covariate reference datasets included 613 total arrays, totaling

536    455 AML and 158 healthy controls. We used ComBat[23] to correct for study batch

537    effects, as its empirical Bayes-based algorithm uses both scale and mean center

538    based methods, providing an appropriate algorithm[23]. Covariate reference

539    datasets were treated as the covariate for batch during batch correction, to

540    improve performance in correcting for batch effects rather than biological

541    variation. After batch correction, we used principal component analysis (PCA),

542    visualizing components in both 2 and 3 dimensions, to compare the clustering

543    results for corrections. Covariate reference datasets were removed after the

544    batch correction step and were not part of our downstream meta-analysis.

545    (Supplementary Fig. S1).

546

547    **Gene expression meta-analysis.**

548    After batch correction step, we performed gene expression meta-analysis for

549    differential expression on the merged datasets (34 data sets, 16 AML and 18

550    healthy), where the expression values for all 44,754 common probe sets were

551    aggregated. The effects of patients' age, sex, and sample source, including their

552    pairwise interactions were investigated using an analysis of variance (ANOVA)[8,62]

553    . For each gene $i$, where $i$=[1,…44,754], the gene expression probe set $Y_i$ was

554    modeled computationally as a linear model:

555    $Y_i \sim (a + s + d + t) + (a{:}s + a{:}d + a{:}t) + (s{:}d + s{:}t) + (d{:}t) + \varepsilon,$

556    where $d$ is the disease state (AML or healthy), $a$ is age (between 0 to 100 years),

557    $s$ is sex (female or male), $t$ is sample source (BM or PB), and $\varepsilon$ is a random error

558    term. We note that the model includes sample source and its interactions to

559    address comparisons involving different tissues in AML and healthy subjects (BM

560    or PB respectively).

561

562    From the ANOVA analysis, genes were deemed to be disease state statistically

563    significant (differentially expressed) if they displayed ANOVA Bonferroni-adjusted

25

564    p-value < 0.01. Post-hoc analysis for significant genes was conducted for

565    comparisons (between groups) using Tukey's Honestly Significant Difference

566    (HSD) tests. Additionally, we performed a quantile-based effect filter, were genes

567    were deemed to show biological effects in our analysis if they displayed mean

568    difference values in the <5% and/or > 95% quantiles of the mean difference

569    distributions of the binary group comparisons. Based on the post-hoc analysis,

570    genes were deemed to be statistically significantly (up- or down-regulated) if they

571    displayed Tukey HSD using a Bonferroni adjusted cutoff for p-value <

572    0.01/44,754.

573

574    **Functional and pathway enrichment analysis**

575    We carried our enrichment analysis for DEPS using the Database DAVID[34,35], the

576    KEGG database[29-31] for signaling pathways, GO terms functional annotation for

577    over representation of biological function [32,33] were utilized and signaling

578    pathways were deemed significant based on Benjamini-Hochberg adjusted p-

579    value < 0.05.

580

581    **Using a k-nearest neighbor model to predict AML**

582    Before gene expression data passed to the k-nearest neighbor (KNN) algorithm

583    to train, gene expression signatures resulted from our meta-analysis were used

584    to extract expression values. KNN in ClassificaIO[37] was used to carry out this

585    analysis. All 34 data sets (16 AML and 18 healthy) were used for training, and

586    testing was done on all 5 covariate data sets, include AML and healthy subjects.

26

587   Dependent, target , and testing data files were prepared in accordance with

588   ClassificaIO[37] user guide. The KNN model used the following parameters

589   (Supplementary File 3):

590

591   *random_state = None, shuffle = True, metric = minkowski, weights = uniform,*

592   *algorithm = auto, n_neighbors = 5, leaf_size = 30, n_jobs = 1, p = 2,*

593   *metric_params = None*

594

595   The trained model was 98% accurate, while testing was 91.7% accurate (details

596   of training and testing are given in Supplementary File 3.

597

598   **DATA AVAILABILITY STATEMENT**

599   The datasets generated in the study, supplementary data, tables, figures and

600   files are available online at http://doi.org/10.5281/zenodo.1492796

601   Datasets re-analyzed in the study are publicly available on the Gene Expression

602   Omnibus repository, at https://www.ncbi.nlm.nih.gov/geo/ under accessions

603   summarized in Table 1.

604

605

27

## REFERENCES

1       Kumar, C. C. Genetic abnormalities and challenges in the treatment of acute myeloid leukemia. *Genes Cancer* **2**, 95-107, doi:10.1177/1947601911408076 (2011).

2       De Kouchkovsky, I. & Abdul-Hay, M. 'Acute myeloid leukemia: a comprehensive review and 2016 update'. *Blood Cancer J* **6**, e441, doi:10.1038/bcj.2016.50 (2016).

3       Siegel, R. L., Miller, K. D. & Jemal, A. Cancer statistics, 2018. *CA Cancer J Clin* **68**, 7-30, doi:10.3322/caac.21442 (2018).

4       Institute, N. C. SEER Cancer Stat Facts: Acute Myeloid Leukemia (Percent of New Cases by Age Group). [https://seer.cancer.gov/statfacts/html/amyl.html].  ((accessed 11.30.18), 2011-2015).

5       Arber, D. A. *et al.* The 2016 revision to the World Health Organization classification of myeloid neoplasms and acute leukemia. *Blood* **127**, 2391-2405, doi:10.1182/blood-2016-03-643544 (2016).

6       Dohner, H. *et al.* Diagnosis and management of acute myeloid leukemia in adults: recommendations from an international expert panel, on behalf of the European LeukemiaNet. *Blood* **115**, 453-474, doi:10.1182/blood-2009-07-235358 (2010).

7       Grimwade, D. & Hills, R. K. Independent prognostic factors for AML outcome. *Hematology Am Soc Hematol Educ Program*, 385-395, doi:10.1182/asheducation-2009.1.385 (2009).

8       Dohner, H. Implication of the molecular characterization of acute myeloid leukemia. *Hematology Am Soc Hematol Educ Program*, 412-419, doi:10.1182/asheducation-2007.1.412 (2007).

9       Walter, M. J. *et al.* Acquired copy number alterations in adult acute myeloid leukemia genomes. *Proc Natl Acad Sci U S A* **106**, 12950-12955, doi:10.1073/pnas.0903091106 (2009).

10      Suela, J., Alvarez, S. & Cigudosa, J. C. DNA profiling by arrayCGH in acute myeloid leukemia and myelodysplastic syndromes. *Cytogenet Genome Res* **118**, 304-309, doi:10.1159/000108314 (2007).

11      Martelli, M. P., Sportoletti, P., Tiacci, E., Martelli, M. F. & Falini, B. Mutational landscape of AML with normal cytogenetics: biological and clinical implications. *Blood Rev* **27**, 13-22, doi:10.1016/j.blre.2012.11.001 (2013).

12      Klepin, H. D., Rao, A. V. & Pardee, T. S. Acute myeloid leukemia and myelodysplastic syndromes in older adults. *J Clin Oncol* **32**, 2541-2552, doi:10.1200/JCO.2014.55.1564 (2014).

13      Short, N. J., Rytting, M. E. & Cortes, J. E. Acute myeloid leukaemia. *Lancet* **392**, 593-606, doi:10.1016/S0140-6736(18)31041-9 (2018).

14      Dohner, H., Weisdorf, D. J. & Bloomfield, C. D. Acute Myeloid Leukemia. *N Engl J Med* **373**, 1136-1152, doi:10.1056/NEJMra1406184 (2015).

15      Reese, N. D. & Schiller, G. J. High-dose cytarabine (HD araC) in the treatment of leukemias: a review. *Curr Hematol Malig Rep* **8**, 141-148, doi:10.1007/s11899-013-0156-3 (2013).

651    16    Meyers, J., Yu, Y., Kaye, J. A. & Davis, K. L. Medicare fee-for-service enrollees
652          with primary acute myeloid leukemia: an analysis of treatment patterns,
653          survival, and healthcare resource utilization and costs. *Appl Health Econ*
654          *Health Policy* **11**, 275-286, doi:10.1007/s40258-013-0032-2 (2013).
655    17    Ferrara, F. & Schiffer, C. A. Acute myeloid leukaemia in adults. *Lancet* **381**,
656          484-495, doi:10.1016/S0140-6736(12)61727-9 (2013).
657    18    Appelbaum, F. R. *et al.* Age and acute myeloid leukemia. *Blood* **107**, 3481-
658          3485, doi:10.1182/blood-2005-09-3724 (2006).
659    19    Cancer Genome Atlas Research, N. *et al.* Genomic and epigenomic landscapes
660          of adult de novo acute myeloid leukemia. *N Engl J Med* **368**, 2059-2074,
661          doi:10.1056/NEJMoa1301689 (2013).
662    20    Walter, M. J. *et al.* Clonal architecture of secondary acute myeloid leukemia. *N*
663          *Engl J Med* **366**, 1090-1098, doi:10.1056/NEJMoa1106968 (2012).
664    21    Miller, B. G. & Stamatoyannopoulos, J. A. Integrative meta-analysis of
665          differential gene expression in acute myeloid leukemia. *PLoS One* **5**, e9466,
666          doi:10.1371/journal.pone.0009466 (2010).
667    22    Ramasamy, A., Mondry, A., Holmes, C. C. & Altman, D. G. Key issues in
668          conducting a meta-analysis of gene expression microarray datasets. *PLoS*
669          *Med* **5**, e184, doi:10.1371/journal.pmed.0050184 (2008).
670    23    Chen, C. *et al.* Removing batch effects in analysis of expression microarray
671          data: an evaluation of six batch adjustment methods. *PLoS One* **6**, e17238,
672          doi:10.1371/journal.pone.0017238 (2011).
673    24    Pavlidis, P. Using ANOVA for gene selection from microarray studies of the
674          nervous system. *Methods* **31**, 282-289, doi:10.1016/S1046-2023(03)00157-
675          9 (2003).
676    25    Pavlidis, P. & Noble, W. S. Matrix2png: a utility for visualizing matrix data.
677          *Bioinformatics* **19**, 295-296, doi:DOI 10.1093/bioinformatics/19.2.295
678          (2003).
679    26    Mias, G. in *Mathematica for Bioinformatics: A Wolfram Language Approach to*
680          *Omics*    193-226 (Springer International Publishing, 2018).
681    27    Neyman, J. & Pearson, E. S. On the use and interpretation of certain test
682          criteria for purposes of statistical inference. Part II. *Biometrika* **20a**, 263-294,
683          doi:DOI 10.1093/biomet/20A.3-4.263 (1928).
684    28    Waltman, L. & Schreiber, M. On the calculation of percentile-based
685          bibliometric indicators. *J Am Soc Inf Sci Tec* **64**, 372-379,
686          doi:10.1002/asi.22775 (2013).
687    29    Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new
688          perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res* **45**,
689          D353-D361, doi:10.1093/nar/gkw1092 (2017).
690    30    Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a
691          reference resource for gene and protein annotation. *Nucleic Acids Research*
692          **44**, D457-D462, doi:10.1093/nar/gkv1070 (2016).
693    31    Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes.
694          *Nucleic Acids Res* **28**, 27-30 (2000).
695    32    Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The
696          Gene Ontology Consortium. *Nat Genet* **25**, 25-29, doi:10.1038/75556 (2000).

697   33   Carbon, S. *et al.* Expansion of the Gene Ontology knowledgebase and
698        resources. *Nucleic Acids Research* **45**, D331-D338, doi:10.1093/nar/gkw1108
699        (2017).
700   34   Huang, D. W., Sherman, B. T. & Lempicki, R. A. Bioinformatics enrichment
701        tools: paths toward the comprehensive functional analysis of large gene lists.
702        *Nucleic Acids Research* **37**, 1-13, doi:10.1093/nar/gkn923 (2009).
703   35   Huang, D. W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative
704        analysis of large gene lists using DAVID bioinformatics resources. *Nat Protoc*
705        **4**, 44-57, doi:10.1038/nprot.2008.211 (2009).
706   36   Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate - a Practical
707        and Powerful Approach to Multiple Testing. *J Roy Stat Soc B Met* **57**, 289-300
708        (1995).
709   37   Roushangar, R. & Mias, G. I. ClassificaIO: machine learning for classification
710        graphical user interface. *bioRxiv*, doi:10.1101/240184 (2017).
711   38   Hou, H. A. *et al.* WT1 mutation in 470 adult patients with acute myeloid
712        leukemia: stability during disease evolution and implication of its
713        incorporation into a survival scoring system. *Blood* **115**, 5222-5231,
714        doi:10.1182/blood-2009-12-259390 (2010).
715   39   Ho, P. A. *et al.* Prevalence and prognostic implications of WT1 mutations in
716        pediatric acute myeloid leukemia (AML): a report from the Children's
717        Oncology Group. *Blood* **116**, 702-710, doi:10.1182/blood-2010-02-268953
718        (2010).
719   40   Udby, L., Calafat, J., Sorensen, O. E., Borregaard, N. & Kjeldsen, L. Identification
720        of human cysteine-rich secretory protein 3 (CRISP-3) as a matrix protein in a
721        subset of peroxidase-negative granules of neutrophils and in the granules of
722        eosinophils. *J Leukocyte Biol* **72**, 462-469 (2002).
723   41   Izzi, V. *et al.* An extracellular matrix signature in leukemia precursor cells and
724        acute myeloid leukemia. *Haematologica* **102**, E245-E248,
725        doi:10.3324/haematol.2017.167304 (2017).
726   42   Buggins, A. G. *et al.* Microenvironment produced by acute myeloid leukemia
727        cells prevents T cell activation and proliferation by inhibition of NF-kappaB,
728        c-Myc, and pRb pathways. *J Immunol* **167**, 6021-6030 (2001).
729   43   Rashidi, A. & Uy, G. L. Targeting the Microenvironment in Acute Myeloid
730        Leukemia. *Curr Hematol Malig R* **10**, 126-131, doi:10.1007/s11899-015-
731        0255-4 (2015).
732   44   Borrow, J. *et al.* The t(7;11)(p15;p15) translocation in acute myeloid
733        leukaemia fuses the genes for nucleoporin NUP98 and class I homeoprotein
734        HOXA9. *Nature Genetics* **12**, 159-167, doi:DOI 10.1038/ng0296-159 (1996).
735   45   Andreeff, M. *et al.* HOX expression patterns identify a common signature for
736        favorable AML. *Leukemia* **22**, 2041-2047, doi:10.1038/leu.2008.198 (2008).
737   46   Fan, C., Stendahl, U., Stjernberg, N. & Beckman, L. Association between
738        Orosomucoid Types and Cancer. *Oncology* **52**, 498-500 (1995).
739   47   Barrett, T. *et al.* NCBI GEO: archive for functional genomics data sets--update.
740        *Nucleic Acids Res* **41**, D991-995, doi:10.1093/nar/gks1193 (2013).

741    48    Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational
742         molecular biology and bioinformatics. *Bioinformatics* **25**, 1422-1423,
743         doi:10.1093/bioinformatics/btp163 (2009).
744    49    Metzeler, K. H. *et al.* An 86-probe-set gene-expression signature predicts
745         survival in cytogenetically normal acute myeloid leukemia. *Blood* **112**, 4193-
746         4201, doi:10.1182/blood-2008-02-134411 (2008).
747    50    Li, Z. *et al.* Identification of a 24-gene prognostic signature that improves the
748         European LeukemiaNet risk classification of acute myeloid leukemia: an
749         international collaborative study. *J Clin Oncol* **31**, 1172-1181,
750         doi:10.1200/JCO.2012.44.3184 (2013).
751    51    Herold, T. *et al.* Isolated trisomy 13 defines a homogeneous AML subgroup
752         with high frequency of mutations in spliceosome genes and poor prognosis.
753         *Blood* **124**, 1304-1311, doi:10.1182/blood-2013-12-540716 (2014).
754    52    Janke, H. *et al.* Activating FLT3 Mutants Show Distinct Gain-of-Function
755         Phenotypes In Vitro and a Characteristic Signaling Pathway Profile
756         Associated with Prognosis in Acute Myeloid Leukemia. *Plos One* **9**, doi:ARTN
757         e89560
758    10.1371/journal.pone.0089560 (2014).
759    53    Jiang, X. *et al.* Eradication of Acute Myeloid Leukemia with FLT3 Ligand-
760         Targeted miR-150 Nanoparticles. *Cancer Res* **76**, 4470-4480,
761         doi:10.1158/0008-5472.CAN-15-2949 (2016).
762    54    Bolstad, B. M., Irizarry, R. A., Astrand, M. & Speed, T. P. A comparison of
763         normalization methods for high density oligonucleotide array data based on
764         variance and bias. *Bioinformatics* **19**, 185-193 (2003).
765    55    Irizarry, R. A. *et al.* Summaries of Affymetrix GeneChip probe level data.
766         *Nucleic Acids Res* **31**, e15 (2003).
767    56    Irizarry, R. A. *et al.* Exploration, normalization, and summaries of high
768         density oligonucleotide array probe level data. *Biostatistics* **4**, 249-264,
769         doi:10.1093/biostatistics/4.2.249 (2003).
770    57    Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J Mach Learn Res*
771         **12**, 2825-2830 (2011).
772    58    Schneider, V. Z., L.; Markus, R.; Fekete, N.; Schrezenmeier, H.; Erle, A. ; Lars,
773         B.; Hofmann, S.; Götz, M.; Döhner, K.; Ihme, S.; Döhner, H.; Buske, C.; Feuring-
774         Buske, M.; Greiner, J. Leukemic progenitor cells are susceptible to targeting
775         by stimulated cytotoxic T cells against immunogenic leukemia-associated
776         antigens.  (2015).
777    59    Majeti, R. *et al.* Dysregulated gene expression networks in human acute
778         myelogenous leukemia stem cells. *Proc Natl Acad Sci U S A* **106**, 3396-3401,
779         doi:10.1073/pnas.0900089106 (2009).
780    60    Bacher, U. *et al.* Multilineage dysplasia does not influence prognosis in
781         CEBPA-mutated AML, supporting the WHO proposal to classify these patients
782         as a unique entity. *Blood* **119**, 4719-4722, doi:10.1182/blood-2011-12-
783         395574 (2012).
784    61    Mills, K. I. *et al.* Microarray-based classifiers and prognosis models identify
785         subgroups with distinct clinical outcomes and high risk of AML

786          transformation of myelodysplastic syndrome. *Blood* **114**, 1063-1072,
787          doi:10.1182/blood-2008-10-187203 (2009).
788   62   Tasaki, S. *et al.* Multi-omics monitoring of drug response in rheumatoid
789          arthritis in pursuit of molecular remission. *Nat Commun* **9**, 2755,
790          doi:10.1038/s41467-018-05044-4 (2018).
791   63   Zatkova, A. *et al.* AML/MDS with 11q/MLL amplification show characteristic
792          gene expression signature and interplay of DNA copy number changes. *Genes*
793          *Chromosomes Cancer* **48**, 510-520, doi:10.1002/gcc.20658 (2009).
794   64   Tomasson, M. H. *et al.* Somatic mutations and germline sequence variants in
795          the expressed tyrosine kinase genes of patients with de novo acute myeloid
796          leukemia. *Blood* **111**, 4797-4808, doi:10.1182/blood-2007-09-113027
797          (2008).
798   65   Taskesen, E. *et al.* Prognostic impact, concurrent genetic mutations, and gene
799          expression features of AML with CEBPA mutations in a cohort of 1182
800          cytogenetically normal AML patients: further evidence for CEBPA double
801          mutant AML as a distinctive disease entity. *Blood* **117**, 2469-2475,
802          doi:10.1182/blood-2010-09-307280 (2011).
803   66   Wouters, B. J. *et al.* Double CEBPA mutations, but not single CEBPA
804          mutations, define a subgroup of acute myeloid leukemia with a distinctive
805          gene expression profile that is uniquely associated with a favorable outcome.
806          *Blood* **113**, 3088-3091, doi:10.1182/blood-2008-09-179895 (2009).
807   67   Figueroa, M. E. *et al.* Genome-wide epigenetic analysis delineates a
808          biologically distinct immature acute leukemia with myeloid/T-lymphoid
809          features. *Blood* **113**, 2795-2804, doi:10.1182/blood-2008-08-172387
810          (2009).
811   68   Klein, H. U. *et al.* Quantitative comparison of microarray experiments with
812          published leukemia related gene expression signatures. *BMC Bioinformatics*
813          **10**, 422, doi:10.1186/1471-2105-10-422 (2009).
814   69   Luck, S. C. *et al.* Deregulated apoptosis signaling in core-binding factor
815          leukemia differentiates clinically relevant, molecular marker-independent
816          subgroups. *Leukemia* **25**, 1728-1738, doi:10.1038/leu.2011.154 (2011).
817   70   Opel, D. *et al.* Targeting inhibitor of apoptosis proteins by Smac mimetic
818          elicits cell death in poor prognostic subgroups of chronic lymphocytic
819          leukemia. *Int J Cancer* **137**, 2959-2970, doi:10.1002/ijc.29650 (2015).
820   71   Cao, Q. *et al.* BCOR regulates myeloid cell proliferation and differentiation.
821          *Leukemia* **30**, 1155-1165, doi:10.1038/leu.2016.2 (2016).
822   72   Li, L. *et al.* Altered hematopoietic cell gene expression precedes development
823          of therapy-related myelodysplasia/acute myeloid leukemia and identifies
824          patients at risk. *Cancer Cell* **20**, 591-605, doi:10.1016/j.ccr.2011.09.011
825          (2011).
826   73   Warren, H. S. *et al.* A genomic score prognostic of outcome in trauma
827          patients. *Mol Med* **15**, 220-227, doi:10.2119/molmed.2009.00027 (2009).
828   74   Karlovich, C. *et al.* A longitudinal study of gene expression in healthy
829          individuals. *BMC Med Genomics* **2**, 33, doi:10.1186/1755-8794-2-33 (2009).

830  75    Kong, S. W. *et al.* Characteristics and predictive value of blood transcriptome
831        signature in males with autism spectrum disorders. *PLoS One* **7**, e49475,
832        doi:10.1371/journal.pone.0049475 (2012).
833  76    Sharma, S. M. *et al.* Insights in to the pathogenesis of axial
834        spondyloarthropathy based on gene expression profiles. *Arthritis Res Ther*
835        **11**, R168, doi:10.1186/ar2855 (2009).
836  77    Rosell, A. *et al.* Brain perihematoma genomic profile following spontaneous
837        human intracerebral hemorrhage. *PLoS One* **6**, e16750,
838        doi:10.1371/journal.pone.0016750 (2011).
839  78    Schmidt, S. *et al.* Identification of glucocorticoid-response genes in children
840        with acute lymphoblastic leukemia. *Blood* **107**, 2061-2069,
841        doi:10.1182/blood-2005-07-2853 (2006).
842  79    Tasaki, S. *et al.* Multiomic disease signatures converge to cytotoxic CD8 T
843        cells in primary Sjogren's syndrome. *Ann Rheum Dis* **76**, 1458-1466,
844        doi:10.1136/annrheumdis-2016-210788 (2017).
845  80    Leday, G. G. R. *et al.* Replicable and Coupled Changes in Innate and Adaptive
846        Immune Gene Expression in Two Case-Control Studies of Blood Microarrays
847        in Major Depressive Disorder. *Biol Psychiatry* **83**, 70-80,
848        doi:10.1016/j.biopsych.2017.01.021 (2018).
849  81    Shamir, R. *et al.* Analysis of blood-based gene expression in idiopathic
850        Parkinson disease. *Neurology* **89**, 1676-1683,
851        doi:10.1212/WNL.0000000000004516 (2017).
852  82    Clelland, C. L. *et al.* Utilization of never-medicated bipolar disorder patients
853        towards development and validation of a peripheral biomarker profile. *PLoS*
854        *One* **8**, e69082, doi:10.1371/journal.pone.0069082 (2013).
855  83    Ducreux, J. *et al.* Interferon alpha kinoid induces neutralizing anti-interferon
856        alpha antibodies that decrease the expression of interferon-induced and B
857        cell activation associated transcripts: analysis of extended follow-up data
858        from the interferon alpha kinoid phase I/II study. *Rheumatology (Oxford)* **55**,
859        1901-1905, doi:10.1093/rheumatology/kew262 (2016).
860  84    Lauwerys, B. R. *et al.* Down-regulation of interferon signature in systemic
861        lupus erythematosus patients by active immunization with interferon alpha-
862        kinoid. *Arthritis Rheum* **65**, 447-456, doi:10.1002/art.37785 (2013).
863  85    Xiao, W. *et al.* A genomic storm in critically injured humans. *J Exp Med* **208**,
864        2581-2590, doi:10.1084/jem.20111354 (2011).
865  86    Zhou, B. *et al.* Analysis of factorial time-course microarrays with application
866        to a clinical study of burn injury. *Proc Natl Acad Sci U S A* **107**, 9923-9928,
867        doi:10.1073/pnas.1002757107 (2010).
868
869

870 **ACKNOWLEDGEMENTS**

877

878 **AUTHOR CONTRIBUTIONS STATEMENT**

879 R.R. and G.I.M. wrote the main manuscript text and prepared the figures. All

880 authors reviewed the manuscript.

881

882 **ADDITIONAL INFORMATION**

883 **Competing interests.** G.I.M. has consulted for Colgate-Palmolive North America

884 and received compensation. R.R. declares no potential conflict of interest.

885

886 **FIGURE LEGENDS**

887 **Figure 1. General approach, data curation, and analysis workflow summary.**

888 The flowchart shows in **(a)** the five main steps that summarize our method of

889 approach for our study, and in (**b**) the curation and screening criteria for raw

890 gene expression and annotation data files curation, data pre-processing,

891 supervised machine learning for missing metadata prediction, and batch effects

892 correction. (**c**) The meta-analysis included a linear model analysis of variance

893   (ANOVA) coupled Tukey's Honestly Significant Difference (HSD) post-hoc tests,

894   and KEGG pathway and GO enrichment. Finally, we performed a machine

895   learning classification of AML based on our findings.

896

897   **Figure 2: Functional classification of DEPS from AML meta-analysis and**

898   **associated KEGG and GO enrichment analysis.** For all panels, normalized

899   values are represented in with blue for down-regulation and red for up-regulation,

900   while light red/gray represents no reported specific direction. (**a**) Heatmap of 974

901   DEPS (rows) on 2,761 arrays (columns) including 2213 AML patients and 548

902   healthy individuals from AML meta-analysis, using unsupervised hierarchical

903   clustering and Euclidean distance for clustering. The age of each individual is

904   displayed at the bottom and illustrated in the color bar on the top (dark green for

905   young and yellow for old). The disease state (AML vs healthy), sex of each

906   subject and age-groups are represented in color bars on the top. (**b**) Horizontal

907   barplot of the top 10 DEPS (gene symbols on vertical axis) from AML meta-

908   analysis with mean difference values between AML and healthy (horizontal axis).

909   Enrichment analysis identified 4 KEGG signaling pathways (**c**) for our AML

910   DEPS, also visualized as a heatmap (**d**) of DEPS mean difference values

911   between AML and healthy DEPS (rows) identified in these 4 KEGG signaling

912   pathways (columns). The GO enrichment analysis results are summarized in (**e**).

913

914   **Figure 3: Sex-related gene expression meta-analysis in AML. (a).** The

915   heatmap of mean difference values comparison between the 70 DE overlapping

916   genes between Analysis 1 and Analysis 2a. (**b**) Heatmap the 70 DEPS

917    expression (rows) on 2,761 arrays (columns) including 2213 AML patients and

918    548 healthy individuals from Analysis 2a of sex-relevance in AML (using

919    unsupervised hierarchical clustering and Euclidean distance for clustering). The

920    disease state (AML vs healthy) and sex of each subject are indicated in color

921    bars at the top.  **(c).** Horizontal barplot of the top 10 DEPS (gene symbols on

922    vertical axis), with the mean difference values between male-female (horizontal

923    axis). (**d).** Enrichment analysis for statistically significant overrepresented

924    biological GO terms on the 70 DE genes.

925

926    **Figure 4: Age-related gene expression meta-analysis in AML. (a)** The top 10

927    up- and down- regulated DEPS overlapping AML and age-related analyses. 75

928    DEPS specific to a single age-group comparison, **(b)**. **(c)** The mean difference of

929    25 DEPS with respect to the 0-19 baseline across all other groups are plotted to

930    illustrate changes with aging. We note that the mean difference values between

931    AML and healthy cohorts are shown in the right-most column of panes (**a**)-(**c**) for

932    reference comparisons. (**d**) Overlaps over KEGG pathways of 17 DE genes

933    identified in 4 KEGG pathways according to age groups.

934 **Table 1: Summary table of all 34 gene expression datasets used in this**

935 **study.**

| Author, Year | GEO accession | Disease Status* | Affymetrix platform id: Number of samples used & Sample source* | Refs. |
|---|---|---|---|---|
| Zatkova et al, 2009 | GSE10258 | AML | GPL570: 8 BM | 63 |
| Tomasson et al, 2008 | GSE10358 | AML | GPL570: 300 BM | 64 |
| Metzeler et al, 2008 | GSE12417 | AML | GPL570: 73 BM & 5 PB<br>GPL96/97: 160 BM & 2PB | 49 |
| Wouters et al, 2009, Taskesen et al, 2011 | GSE14468 | AML | GPL570: 482 BM & 43 PB | 65,66 |
| Figueroa et al, 2009 | GSE14479 | AML | GPL570: 16 BM | 67 |
| Klein et al, 2009 | GSE15434 | AML | GPL570: 231 BM & 20 PB | 68 |
| Lück et al, 2011 | GSE29883 | AML | GPL570: 10 BM & 2 PB | 69 |
| Li et al, 2013, Herold et al, 2014, Janke et al, 2014, Jiang et al, 2016 | GSE37642 | AML | GPL570: 140 BM<br>GPL96/97: 422 BM | 50-53 |
| Bullinger et al, 2014 | GSE39363 | AML | GPL570: 11 BM & 2 PB | NYP |
| Opel et al, 2015 | GSE46819 | AML | GPL570: 8 BM & 4 PB | 70 |
| TCGA et al, 2015 | GSE68833 | AML | GPL570: 183 BM | NYP |
| Cao et al, 2016 | GSE69565 | AML | GPL570: 12 PB | 71 |
| Bohl et al, 2016 | GSE84334 | AML | GPL570: 25 BM & 20 PB | NYP |
| Li et al, 2011 | GSE23025 | AML | GPL570: 21 BM & 13 PB | 72 |
| Warren et al, 2009 | GSE11375 | Healthy | GPL570: 26 PB | 73 |
| Green et al, 2009 | GSE14845 | Healthy | GPL570: 1 PB | NYP |
| Wu et al, 2012 | GSE15932 | Healthy | GPL570: 8 PB | NYP |
| Karlovich et al, 2009 | GSE16028 | Healthy | GPL570: 22 PB | 74 |
| Krug et al, 2011 | GSE17114 | Healthy | GPL570: 14 PB | NYP |
| Kong et al, 2012 | GSE18123 | Healthy | GPL570: 17 PB | 75 |
| Sharma et al, 2009 | GSE18781 | Healthy | GPL570: 25 PB | 76 |
| Rosell et al, 2011 | GSE25414 | Healthy | GPL570: 12 PB | 77 |
| Schmidt et al, 2006 | GSE2842 | Healthy | GPL570: 2 PB | 78 |
| Meng et al, 2015 | GSE71226 | Healthy | GPL570: 3 PB | NYP |
| Tasaki et al, 2017 | GSE84844 | Healthy | GPL570: 30 PB | 79 |
| Leday et al, 2018 | GSE98793 | Healthy | GPL570: 64 PB | 80 |
| Shamir et al, 2017 | GSE99039 | Healthy | GPL570: 121 PB | 81 |
| Tasaki et al, 2018 | GSE93272 | Healthy | GPL570: 35 PB | 62 |
| Clelland et al, 2013 | GSE46449 | Healthy | GPL570: 24 PB | 82 |
| Lauwerys et al, 2013 Ducreux et al, 2016 | GSE39088 | Healthy | GPL570: 46 PB | 83,84 |
| Xiao et al, 2011 | GSE36809 | Healthy | GPL570: 35 PB | 85 |
| Zhou et al, 2010 | GSE19743 | Healthy | GPL570: 63 PB | 86 |
| Jiang et al, 2018[#] | GSE107968[*] | 2 AML,<br>1 Healthy | GPL570: 3 BM | NYP |
| Greiner et al, 2015[#] | GSE68172[*] | 20 AML,<br>5 Healthy | GPL570: 25 PB | 58 |
| Majeti et al, 2009[#] | GSE17054[*] | 9 AML,<br>4 Healthy | GPL570: 13 BM | 59 |
| Bacher et al, 2012[#] | GSE33223[*] | 20 AML,<br>10 Healthy | GPL570: 30 PB | 60 |
| Mills et al, 2009[#] | GSE15061[*] | 404 AML,<br>138 Healthy | GPL570: 542 BM | 61 |

**Meta-analysis data sets summary**

| Disease state | | Sample source | | Affymetrix platform id | | Unique probe sets | |
|---|---|---|---|---|---|---|---|
| AML | Healthy | BM | PB | GPL570 | GPL96/97 | GPL570 | GPL96/97 |
| 2213 | 548 | 2090 | 671 | 2177 | 584 | 54,675 | 44,760 |

**Table 1.** A summary table of all our data sets using in our meta-analysis and disease classification.
[#]"Covariate reference data sets," 5 data sets that were used during the batch correction step., datasets were used only during the batch effect correction steps.
*GEO, Gene Expression Omnibus; AML, acute myeloid leukemia; Ref. reference; NYP, not yet published, GPL570, Affymetrix Human Genome U133 Plus 2.0 Array; GPL96, Affymetrix Human Genome U133A Array; GPL97, Affymetrix Human Genome U133B Array; BM, Bone Marrow; PB, Peripheral Blood.

936 **Table 2. Top 10 up- and down-regulated of DEPS in AML from disease state**
937

| Up-regulated* | | | |
|---|---|---|---|
| DEG name | DEPS Gene Symbol | Tukey's HSD Mean difference | Bonferroni (p-adjusted) |
| Wilms tumor 1 | WT1 | 0.255353 | < 4.11E-11 |
| MAM domain containing 2 | MAMDC2 | 0.248983 | 5.47E-09 |
| X inactive specific transcript (non-protein coding) | XIST | 0.230331 | < 4.11E-11 |
| homeobox A3 | HOXA3 | 0.195790 | 1.1E-06 |
| fms-related tyrosine kinase 3 | FLT3 | 0.193420 | < 4.11E-11 |
| cyclin A1 | CCNA1 | 0.185050 | 1.35E-07 |
| mex-3 RNA binding family member B | MEX3B | 0.181068 | < 4.11E-11 |
| collagen, type IV, alpha 5 | COL4A5 | 0.177721 | 1.7E-05 |
| neurexin 2 | NRXN2 | 0.166598 | < 4.11E-11 |
| ATPase, Na+/K+ transporting, beta 1 polypeptide | ATP1B1 | 0.165197 | 5.47E-09 |
| **Down-regulated** | | | |
| cysteine-rich secretory protein 3 | CRISP3 | -0.51965625 | < 4.11E-11 |
| olfactomedin 4 | OLFM4 | -0.489845396 | < 4.11E-11 |
| orosomucoid 1 | ORM1 | -0.465232864 | < 4.11E-11 |
| cytochrome P450, family 4, subfamily F, polypeptide 3 | CYP4F3 | -0.453467442 | < 4.11E-11 |
| chitinase 3-like 1 (cartilage glycoprotein-39) | CHI3L1 | -0.421520435 | < 4.11E-11 |
| annexin A3 | ANXA3 | -0.390688999 | < 4.11E-11 |
| oxidized low density lipoprotein (lectin-like) receptor 1 | OLR1 | -0.35525472 | < 4.11E-11 |
| carcinoembryonic antigen-related cell adhesion molecule 8 | CEACAM8 | -0.351181264 | < 4.11E-11 |
| orosomucoid 1 | ORM1 | -0.336303304 | < 4.11E-11 |
| tumor-associated calcium signal transducer 2 | TACSTD2 | -0.323939961 | < 4.11E-11 |

**Table 2.** From the Post-hoc Tukey's test, gene expression means difference value < 5% or > 95% between AML and healthy (AML - healthy) were deemed statistically significant for AML. Genes were considered disease state statistically significant from the analysis of all 2761 cases (2213 AML patients and 548 healthy controls) using. The p-values were adjusted based on Bonferroni correction for false discovery rate (FDR). Significant DEPS (gene symbols) are listed in descending order of the mean difference value comparisons for disease state.

938

**Table 3. KEGG pathway analysis of DEPS from meta-analysis of 34 gene expression datasets.**

**AML Vs Healthy DEPS and associated signaling pathways**

| Pathway | No. of genes* | Down-regulated | Up-regulated | p-value | p-value Benjamini adjusted |
|---|---|---|---|---|---|
| **Hematopoietic cell lineage** | 11, 6 | IL1R2, CD59, GYPA, MS4A1, EPOR, CD24, CD14, EPOR, IL1R1, MME, CR1 | ITGA4, FLT3, CD34, IL3RA, ITGA5, CD44 | 2.3E-5 | 5.8E-3 |
| **Cell cycle** | 12, 6 | CDC7, CDC6, CCNB1, CDC20, CCNA2, CCNE2, TTK, CDC14B', CDK1, BUB1, CCNB2, BUB1B | RB1, CCNA1, CDK6, ATM, TFDP2, CDKN2A | 1.4E-4 | 1.2E-2 |
| **p53 signaling pathway** | 6, 7 | THBS1, CCNB1, CCNE2, CDK1, RRM2, CCNB2 | SIAH1, CDK6, ATM, SERPINE1, CDKN2A, PMAIP1, ZMAT3 | 1.0E-4 | 1.3E-2 |
| **Transcriptional misregulation in cancer** | 7, 13 | IL1R2, GZMB, CD14, ELANE, MMP9, CEBPE, PBX1 | WT1, RUNX2, ETV5, MEIS1, JUP, EWSR1, ATM, HOXA10, MLF1, FLT3, CCNT2, MEF2C, SLC45A3 | 6.5E-4 | 4.1E-2 |

**AML sex relevant (male - female) DEPS & associated signaling pathways**

| Pathway | No. of genes* | High in Females | High in Males |
|---|---|---|---|
| **Hematopoietic cell lineage** | 1, 2 | – | FLT3, CD34 |
| **p53 signaling pathway** | –, 1 | – | PMAIP1 |
| **Transcriptional misregulation in cancer** | –, 1 | MS4A1 | FLT3 |

**Table 3:** Enrichment analysis was done using 974 DEPS, including KEGG enrichment analysis identified 4 statistically significant pathways from AML Vs Healthy meta-analysis, shown with overlaps with sex-specific analysis.
* up and down regulated genes displayed

942 **Table 4. KEGG pathway analysis of DEPS from meta-analysis of 34 gene**
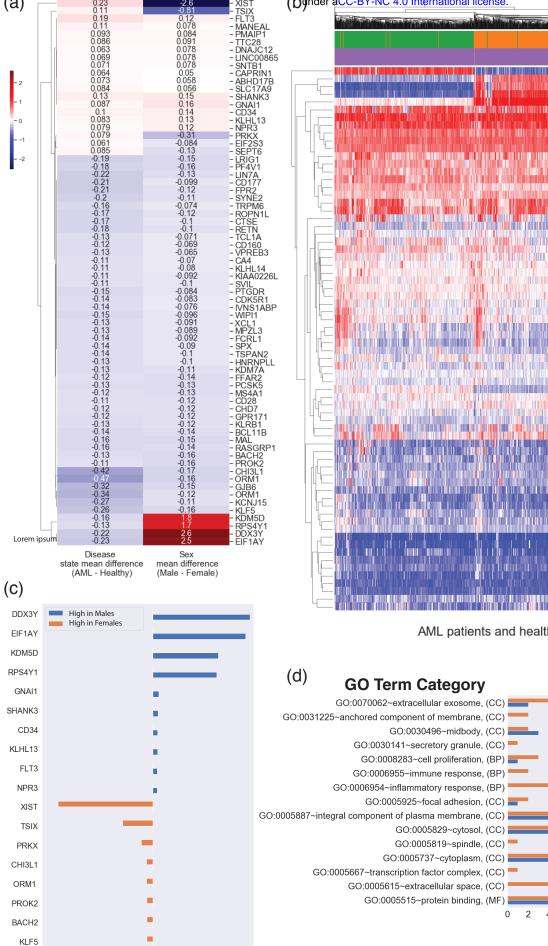943 **expression datasets overlap with age-specific findings.**

| AML age-dependent (AML - healthy) DEPS & associated signaling pathways | | | |
|---|---|---|---|
| Pathway | No. of genes* | Down-regulated Age-group | Up-regulated Age-group |
| Hematopoietic cell lineage | 4, 1 | CD14 (30 to 39) - (0 to 19) / MME (30 to 39) - (0 to 19), (40 to 49) - (0 to 19), (50 to 59) - (0 to 19) / CD24 (30 to 39) - (0 to 19), (40 to 49) - (0 to 19), (50 to 59) - (0 to 19) / MS4A1 (40 to 49) - (0 to 19), (50 to 59) - (0 to 19), (60 to 69) - (0 to 19), (70 to 79) - (0 to 19), (80 to 100) - (0 to 19) | FLT3 (20 to 29) - (0 to 19), (30 to 39) - (0 to 19), (40 to 49) - (0 to 19), (50 to 59) - (0 to 19), (60 to 69) - (0 to 19), (70 to 79) - (0 to 19), (80 to 100) - (0 to 19) |
| Cell cycle | 3, 2 | CCNA2 (50 to 59) - (0 to 19) / CDK6 (60 to 69) - (30 to 39) / CDC14B (30 to 39) - (0 to 19), (40 to 49) - (0 to 19), (50 to 59) - (0 to 19), (60 to 69) - (0 to 19), (70 to 79) - (0 to 19) | CCNA1 (30 to 39) - (0 to 19), (40 to 49) - (0 to 19), (50 to 59) - (0 to 19), (60 to 69) - (0 to 19) / CDKN2A (40 to 49) - (0 to 19) |
| p53 signaling pathway | 1, 1 | CDK6 (60 to 69) - (30 to 39) | CDKN2A (40 to 49) - (0 to 19) |
| Transcriptional misregulation in cancer | 5, 4 | CD14 (30 to 39) - (0 to 19) / MMP9 (20 to 29) - (0 to 19), (30 to 39) - (0 to 19), (40 to 49) - (0 to 19), (50 to 59) - (0 to 19), (60 to 69) - (0 to 19), (70 to 79) - (0 to 19) / EWSR1 (60 to 69) - (50 to 59), (70 to 79) - (50 to 59) / CEBPE (20 to 29) - (0 to 19), (30 to 39) - (0 to 19), (40 to 49) - (0 to 19), (50 to 59) - (0 to 19), (50 to 59) - (20 to 29), (60 to 69) - (0 to19), (70 to 79) - (0 to 19), (70 to 79) - (20 to29), (80 to 100) - (0 to 19) / CCNT2 (60 to 69) - (30 to 39), (70 to 79) - (30 to 39), (60 to 69) - (50 to 59) | MEIS1 (50 to 59) - (0 to 19), (50 to 59) - (20 to 29), (60 to 69) - (0 to 19), (60 to 69) - (20 to 29), (70 to 79) - (0 to 19) / WT1 (20 to 29) - (0 to 19), (30 to 39) - (0 to 19), (40 to 49) - (0 to 19), (50 to 59) - (0 to 19), (60 to 69) - (0 to 19), (70 to 79) - (0 to 19) / FLT3 (20 to 29) - (0 to 19), (30 to 39) - (0 to 19), (40 to 49) - (0 to 19), (50 to 59) - (0 to 19), (60 to 69) - (0 to 19), (70 to 79) - (0 to 19), (80 to 100) - (0 to 19) / HOXA10 (40 to 49) - (0 to 19), (50 to 59) - (0 to 19), (50 to 59) - (20 to 29), (60 to 69) - (0 to 19), (60 to 69) - (20 to 29), (70 to 79) - (0 to 19) |

**Table 4:** Enrichment analysis was done using 974 DEPS overlapped with age-specific analysis
* up and down regulated genes displayed

944

## a. Overview of approach and meta-analysis methods

**i) Gene expression data curation**
- Curate raw data and meta data for AML and healthy

**ii) Data pre-processing & batch correction**
- Remove non-biological variation in each dataset. Remove effects caused by grouping of different datasets

**iii) Supervised learning (classification)**
- Predict missing meta data required for downstream analyses (e.g. sex, sample source, etc.)

**v) Gene annotation, pathways enrichment, & AML prediction**

**iv) ANOVA with Tukey HSD for differential gene expression**
- Changes in AML vs Healthy with respect to age, sex, and sample source

## b. Pre-processing & Dataset-wise correction for batch effects

**Microarray data**
Group based on technology

**Data screening inclusion criteria**
1: Human & untreated?
2: Blood/bone marrow source?
3: Data normally distributed?

**Probe set level pre-processing of each data set separately**
1: Background correction
2: Normalization (quantile)
3: Summarization (median polish)
4: Log(base)2 transformation

**z-score standardization to account for variance within each subject**

**Supervised machine learning**
Classification of missing meta-data annotation. e.g. sample source, sex, etc.

**Batch effect correction**
1: Covariant data sets (5 AML datasets) each includes healthy controls
2: Used ComBat
3: Remove effects due to batch

## c. Meta-analysis

**Linear model**
*For each gene i, where i=[1,...44,754] the gene expression*
$Y_i$ *is modeled computationally as a linear model:*
$$Y_i \sim (a + s + d + t) + (a{:}s + a{:}d + a{:}t) + (s{:}d + s{:}t) + (d{:}t) + \varepsilon$$

**ANOVA**
**Test disease state for statistical significance**
ANOVA Bonferroni adjusted p.value <0.01

**Tukey's HSD**
**Test disease state for significant biological effect**
5%> Mean Group Difference >95%

**Test age, sex, & sample source statistical significance**
Tukey's Bonferroni adjusted p.value <(0.01/44,754)

**Gene Annotation**
Probe set annotation file, release 36 ThermoFisher

**Gene enrichment analysis**
KEGG & GO terms

**Classification of AML!**
Based on our gene expression results

(a)

(b)

AML patients and healthy individuals

(c)

(d) GO Term Category