1  **A single cell framework for multi-omic analysis of disease identifies malignant**

2  **regulatory signatures in mixed phenotype acute leukemia**

3

4  **AUTHORS**

5  Jeffrey M. Granja[1,2,3#], Sandy Klemm[3,#], Lisa M. McGinnis[3,4,#], Arwa S. Kathiria[3], Anja

6  Mezger[3], Benjamin Parks[3,5], Eric Gars[4], Michaela Liedtke[9], Grace X.Y. Zheng[6], Howard

7  Y. Chang[1,3,7,8], Ravindra Majeti[9], William J. Greenleaf[1,3,10,11]

8

9  **AFFILIATIONS**

10  [1] Center for Personal Dynamic Regulomes, Stanford University School of Medicine,

11  Stanford, CA 94305, USA.

12  [2] Biophysics Program, Stanford University School of Medicine, Stanford, CA 94305, USA.

13  [3] Department of Genetics, Stanford University School of Medicine, Stanford, CA 94305,

14  USA.

15  [4] Department of Pathology, Stanford University School of Medicine, Stanford, CA 94305,

16  USA.

17  [5] Department of Computer Science, Stanford University School of Engineering, Stanford,

18  CA 94305, USA.

19  [6] 10x Genomics, Inc., Pleasanton, CA 94566, USA.

20  [7] Department of Dermatology, Stanford University School of Medicine, Redwood City, CA

21  94063, USA.

22  [8] Howard Hughes Medical Institute, Stanford University School of Medicine, Stanford, CA

23  94305

24  [9] Department of Medicine, Division of Hematology, Stanford Cancer Institute, Stanford

25  University School of Medicine, Stanford, CA 94305

26  [10] Department of Applied Physics, Stanford University, Stanford, CA 94025

27  [11] Chan Zuckerberg Biohub, San Francisco, CA 94158, USA

28  [#] These authors contributed equally to this study.

29  Correspondence to: W.J.G. (wjg@stanford.edu), S.K. (klemm@stanford.edu), L.M.M

30  (lisa.mcginnis@stanford.edu)

31

2

32  **Abstract**

33  In order to identify the molecular determinants of human diseases, such as cancer, that

34  arise from a diverse range of tissue, it is necessary to accurately distinguish normal and

35  pathogenic cellular programs.[1–3] Here we present a novel approach for single-cell multi-

36  omic deconvolution of healthy and pathological molecular signatures within phenotypically

37  heterogeneous malignant cells. By first creating immunophenotypic, transcriptomic and

38  epigenetic single-cell maps of hematopoietic development from healthy peripheral blood

39  and bone marrow mononuclear cells, we identify cancer-specific transcriptional and

40  chromatin signatures from single cells in a cohort of mixed phenotype acute leukemia

41  (MPAL) clinical samples. MPALs are a high-risk subtype of acute leukemia characterized

42  by a heterogeneous malignant cell population expressing both myeloid and lymphoid

43  lineage-specific markers.[4,5] Our results reveal widespread heterogeneity in the

44  pathogenetic gene regulatory and expression programs across patients, yet relatively

45  consistent changes within patients even across malignant cells occupying diverse portions

46  of the hematopoietic lineage. An integrative analysis of transcriptomic and epigenetic

47  maps identifies 91,601 putative gene-regulatory interactions and classifies a number of

48  transcription factors that regulate leukemia specific genes, including *RUNX1*-linked

49  regulatory elements proximal to *CD69*. This work provides a template for integrative, multi-

50  omic analysis for the interpretation of pathogenic molecular signatures in the context of

51  developmental origin.

52

53

54

55

56

57

58    **Main**

59    To identify pathologic features within neoplastic cells, we first aimed to establish molecular

60    features of normal development for comparison. Since MPALs present with features of

61    multiple hematopoietic lineages, we first constructed independent immunophenotypic,

62    transcriptomic and epigenetic maps of normal blood development using droplet-based

63    CITE-seq[6] (single-cell antibody derived tag and RNA sequencing) and single-cell ATAC-

64    seq (scATAC-seq, single-cell chromatin accessibility profiling)[7] on bone marrow and

65    peripheral blood mononuclear cells (**Figure 1a**). For CITE-seq analyses, we

66    simultaneously generated 10x Genomics 3' single-cell RNA sequencing[8] (scRNA-seq) and

67    antibody derived tag sequencing[6] (scADT-seq) libraries from 35,882 bone marrow

68    mononuclear cells (BMMCs, n = 12,602), CD34$^+$ enriched BMMCs (n = 8,176), and

69    peripheral blood mononuclear cells (PBMC, n = 14,804). On average, 1,273 informative

70    genes (2,370 unique transcript molecules) were detected per cell and replicates were

71    highly correlated (Supplementary Figure 1a-b). We then selected a feature set of

72    transcripts to mitigate batch effects and linearly projected retained transcript counts into a

73    lower dimensional space using Latent Semantic Indexing (LSI, see Online Methods).[9,10]

74    Cells were clustered using Seurat's Shared Nearest Neighbor approach[11], annotated

75    using a manually curated maker gene list, and visualized using uniform manifold

76    approximation and projection (UMAP)[12] (**Figure 1b**, Supplementary Figure 1c-d).

77         We next established an epigenetic map of normal hematopoiesis by measuring

78    chromatin accessibility across 35,038 single BMMCs (n = 16,510), CD34$^+$ BMMCs (n =

79    10,160), and PBMCs (n = 8,368) using droplet scATAC-seq (10x Genomics)[7]. These cells

80    exhibited a canonical fragment size distribution with clearly resolved sub-, mono-, and

81    multi-nucleosomal modes, a high signal-to-noise ratio at transcription start sites, an

82    average of 11,597 uniquely accessible fragments per cell on average, and a majority

83    (61%) of Tn5 insertions aligning within peaks (Supplementary Figure 2a-c). After pooling

84    all scATAC-seq profiles from each experiment, we confirmed higher reproducibility across

85    replicates than across different samples, similar to the scRNA-seq analysis

86    (Supplementary Figure 2d). Using LSI, Seurat's Shared Nearest Neighbor clustering, and

87    UMAP, we generated a chromatin accessibility map of hematopoiesis that complements

88    the transcriptional map of hematopoiesis (**Figure 1c**, Supplementary 2e-f).

89         To validate the proposed transcriptomic and epigenetic single-cell maps of

90    hematopoiesis, we directly visualized lineage-restricted cell-surface marker and

91    transcription factor enrichment across each map. As anticipated, both scADT- and scRNA-

92    seq measurements of surface makers demonstrate *CD3* enrichment across bone marrow

93    and peripheral T cells; *CD14* enrichment within the monocytic lineage; broad up regulation

94    of *CD19* across the B cell lineage; and *CD8A* enrichment within cytotoxic T lymphocytes

95    (**Figure 1d**)[13]. Estimates of gene activity based on correlated variation in promoter and

96    distal peak accessibility (Cicero[14]) broadly recapitulates this pattern, confirming that

97    lineage specification is consistently reflected across the phenotypic, transcriptional and

98    epigenetic maps of hematopoietic development (**Figure 1d**). We then visualized our high

99    quality scADT-seq using UMAP and found that we could broadly recapitulate our

100   transcriptomic hematopoietic map (Supplementary Figure 3a-d). To further support these

101   cell type identifications and developmental mappings, we show concordance between

102   three separate single-cell measurements, including direct transcript measurements from

103   the scRNA-seq dataset, inferred gene activity scores from the scATAC-seq dataset, and

104   TF activity using chromVAR[15], for key developmental transcription factors, including

105   *CEBPB* in monocytic development, *GATA1* within the erythroid lineage, and *TBX21* in NK

106   and CD8[+] T memory cells, and *PAX5* in B cell and plasmacytoid dendritic cell development

107   (**Figure 1e**). High-resolution single cell multi-omic tracks for key marker genes in each of

108   the identified lineages further support these identifications (**Figure 1f-h,** Supplementary

109   Figure 4a-h). Collectively these results show that the proposed multi-omic maps of healthy

110    hematopoiesis are consistent and broadly capture essential phenotypic, transcriptomic

111    and epigenetic features of blood development.

112         Recent work has shown that immunophenotypically-distinct subpopulations of

113    MPAL blasts share similar genomic lesions within a patient, and that cells from one lineage

114    can reconstitute the alternate lineage in xenograft models[16], suggesting that MPAL lineage

115    plasticity may be epigenetically regulated. To explore the nature of this regulatory and

116    phenotypic dysfunction, we assayed six MPAL samples including three T-myeloid (T/M)

117    MPALs (MPAL1-3), 1 B-myeloid (B/M) MPAL (MPAL4), and one T/M MPAL sampled

118    before CALGB chemotherapy (MPAL5) and after post-treatment relapse (MPAL5R) (see

119    Supplementary Table 1). Across these samples, we observed extensive

120    immunophenotypic heterogeneity (via diagnostic flow cytometric analysis) including

121    bilineal patterns (multiple blast populations expressing both lymphoid and myeloid lineage

122    antigens), biphenotypic patterns (a dominant blast population that simultaneously

123    expresses both lymphoid and myeloid antigens), and both patterns (Supplementary

124    Figures 5a-c, 6a-c). We then performed Whole Exome Sequencing (WES) and found

125    mutational profiles similar to previous studies (Supplementary Figure 6d)[16,17]. To further

126    profile our MPAL samples, we performed CITE-seq (18,056 cells) and scATAC-seq

127    (35,423 cells) on either peripheral blood or bone marrow aspirates from these MPAL

128    patients, observing similar high data quality to that obtained for healthy samples

129    (Supplementary Figure 7a-f).

130         Using our transcriptomic and chromatin landscapes of healthy hematopoiesis, we

131    next sought to develop an analytical framework to identify the hematopoietic

132    developmental signature at single-cell resolution. First, the chromatin and gene

133    expression signatures of single cells are projected into our ATAC- and RNA-based healthy

134    hematopoietic map's LSI subspace, and the results are then visualized using UMAP

135    (**Figure 2a**, Supplementary Figure 8a). Next, by determining the closest hematopoietic

136    cells to the projected cells we can identify the hematopoietic developmental compartment.

137    This method does not require defining discrete cell type boundaries and uses a large

138    feature set to robustly position cells within the continuous landscape of hematopoiesis. To

139    validate this approach, we first projected downsampled, published bulk RNA-seq and

140    ATAC-seq data[18] from FACS-sorted subpopulations into our chromatin and transcription

141    hematopoietic maps and found high concordance with our healthy hematopoietic map and

142    cluster definitions (Supplementary Figure 8b). To further validate our approach, we

143    projected published scRNA-seq[19] and scATAC-seq[20–22] data from different platforms and

144    different genomes on our chromatin and transcription hematopoietic maps and found

145    striking agreement (Supplementary Figure 8c). These results confirm that this method can

146    accurately identify the hematopoietic signature for chromatin and gene expression at

147    single-cell resolution.

148        Using this LSI projection framework and landscapes of healthy hematopoiesis, we

149    next sought to deconvolve the normal and leukemic signatures of MPAL samples at single-

150    cell resolution. First, the leukemic single cells are projected into the hematopoietic linear

151    LSI subspace. Next we identify a non-redundant set of healthy hematopoietic cells that

152    were nearest neighbor normal cells to each leukemic cell, irrespective of their cell-type

153    boundaries. Lastly, we compute the differences between the leukemic cells and nearest

154    normal cells to identify the leukemic specific signature. We first tested our approach by

155    analyzing recently published scRNA-seq data from acute myeloid leukemia (AML) patient

156    samples[19]. By projecting the AMLs into our healthy hematopoietic map, we see general

157    agreement with previous classifications without need for the establishment of potential

158    arbitrary cell-type boundaries on normal hematopoiesis (Supplementary Figure 9a-c). We

159    next projected our phenotypically diverse MPAL patient samples onto our hematopoietic

160    maps and discovered broad epigenetic and gene expression diversity. To further resolve

161    this diversity, we grouped MPAL cells within individual patients into broad hematopoietic

162   developmental compartments: progenitors-like (purple) comprising human stem cell and

163   multipotent progenitor-like cells, lymphoid-like (blue) containing lymphoid-primed

164   multipotent progenitors, erythroid-like (red) which include megakaryocyte-erythroid

165   progenitors, myeloid-like (green) which include granulocyte-monocyte progenitors, and

166   T/NK-like (orange) which include differentiated T and NK cells[23] (**Figure 2a-b**). Strikingly,

167   we see that the scADT-seq data clearly resolves the dominant MPAL subpopulations in

168   MPAL1 and MPAL5; however it does not fully capture the transcriptional diversity of

169   MPALs 2-4 (Supplementary Figure 10a). We visualized these projected MPALs colored

170   by these broad hematopoietic compartments, observing the expected high concordance

171   between the scRNA- and scATAC-seq classifications (**Figure 2b**). Comparing MPAL gene

172   expression to this healthy nearest neighbor set allowed the identification of pathogenic

173   differential gene expression for MPALs from different compartments. In total, we identified

174   4,616 genes that were significantly up-regulated (LFC > 0.5 and FDR < 0.01) in at least

175   one MPAL subpopulation across the six patient samples, and grouped these genes with

176   k-means clustering (**Figure 2c**). We further categorized the most conserved differential

177   genes, TFs and KEGG pathways across the MPALs (Supplementary Figure 11a-c). Using

178   the same approach for the scATAC-seq data, we performed differential peak testing for

179   each MPAL subpopulation and found 72,196 significantly up-regulated peaks (LFC > 0.5

180   and FDR < 0.05) in at least one MPAL subpopulation (**Figure 2c**). Multi-omic differential

181   tracks for the cyclin dependent kinase *CDK11A* and cyclin dependent kinase inhibitor

182   *CDKN2A* , genes that are recurrently mutated in MPAL[16,24], demonstrate these leukemia-

183   specific ATAC- and RNA-seq differences (Supplementary Figure 11d-e). Additionally, we

184   calculated Pearson correlations of the differential genes and peaks; and found that

185   transcription and accessibility differs significantly *across* patients, but is relatively

186   conserved across subpopulations *within* patients. (**Figure 2d**).

187    To compare the MPAL hematopoietic compartments' leukemic programs to

188    previous studies, we downsampled bulk leukemia RNA-seq and projected onto our

189    transcriptomic hematopoietic UMAP for childhood AMLs, B-acute lymphoblastic

190    leukemias (B-ALLs), early T-cell precursor T-acute lymphoblastic leukemias (ETP T-

191    ALLs), non-ETP T-ALLs and MPALs[16] (Supplementary Figure 12a-b). We then calculated

192    differential expression with respect to the closest normal cell populations to identify their

193    respective leukemic programs. Next, we performed LSI on variable malignant genes

194    across all the leukemia subtypes, including MPALs 1-5, and then visualized these patients

195    with UMAP (**Figure 2e,** Supplementary Figure 12c-d). Interestingly, we found large

196    differences in the leukemic programs across various leukemias including T-ALLs, B-ALLs,

197    and across different cytogenetic subtypes. In addition, we found that the MPALs assayed

198    in this study were representative of previous characterized MPALs[16] (**Figure 2e**). Given

199    that we were insufficiently powered to detect unique leukemic differences between AML

200    and our MPAL samples when analyzing downsampled bulk data, we compared the

201    malignant transcriptomic profiles identified from re-analyzing AML scRNA-seq data[18] with

202    our MPALs in order to dissect further these unique malignant signatures (**Figure 2c**,

203    Supplementary Figure 9c). To this end, we identified genes that were more commonly

204    universally upregulated in AMLs, in MPALs, or jointly upregulated in both leukemias

205    (**Figure 2f**, Supplementary Figure 9c). These gene sets provide fine-grained phenotypic

206    resolution comparing the differences and similarities between AML and MPAL leukemic

207    programs and suggest possible insight into why MPALs respond poorly to AML

208    treatment[25,26].

209    Having compared our leukemic transcriptomic programs to other studies we

210    wanted to identify the key TFs that regulate these programs. First, we identified which TF

211    were differentially enriched in each k-means cluster of differentially accessible peaks

212    observed in Figure 2c. (**Figure 3a**). We found that *RUNX1* motifs were highly enriched in

213    both cluster 4 and 10 – the two clusters corresponding to the most commonly shared

214    accessible elements across MPAL subset populations. In addition, *RUNX1* is significantly

215    up-regulated in about half (7/17) of the MPAL subpopulations. *RUNX1* is one of the most

216    frequently mutated genes across hematologic malignancies acting as both a tumor

217    suppressor with loss-of-function mutations in AML[27], myelodysplastic syndrome (MDS)[28],

218    and ETP T-ALL[29,30], and as a putative oncogene in non-ETP T-ALL[31,32]. Furthermore,

219    wildtype *RUNX1* has been implicated as a potential driver of leukemogenesis in core-

220    binding factor (CBF) leukemia[33] and mixed lineage leukemia[34].

221          To link *RUNX1* and other putative regulatory TFs to their leukemic programs we

222    first developed an analytical framework that utilizes both our transcriptomic and chromatin

223    single-cell data to link putative regulator peaks to target genes. Using our matched

224    scATAC and scRNA data for all MPALs and concordant hematopoietic maps, and aligned

225    each single-cell into a common subspace using Canonical Correlation Analyses

226    (CCA)[10,11,35]. For each scATAC cell, we identified the nearest scRNA-seq neighbor

227    (**Figure 3b**, Supplementary Figure 13a-b**)**. We found that the mapping of scATAC cell

228    clusters to scRNA-defined cell clusters were highly consistent (single-cell overlap of 52%

229    across 26 clusters) (Supplementary Figure 14a-d). We then aggregated our scATAC cells

230    based on nearest neighbors in the LSI subspace using Cicero[14] and created a

231    corresponding scRNA aggregate for each cluster using the constructed CCA alignment.

232    We next identified 91,601 peak-to-gene links by correlating accessibility changes of ATAC

233    peaks within 250 kb of the gene promoter with the expression of the gene independently

234    for both healthy and MPAL aggregates (**Figure 3b**). This analysis revealed peak-to-gene

235    links that were specific to healthy hematopoiesis, others that were specific to MPALs, and

236    a conserved subset that was shared across both hematopoiesis and MPALs. We

237    hypothesize that the MPAL-specific peak-to-gene links may be important for leukemic

238    gene regulation. Overall, the identified set of peak-to-gene links had similar distributions

239   for peaks mapped per gene, genes mapped per peak, number of skipped genes and the

240   peak-to-gene as previously observed in a similar linkage analyses[2] (Supplementary Figure

241   14e). To further support these peak-to-gene links, we used previously published K27ac

242   HiChIP in primary T cells and a Human Coronary Artery Smooth Muscle Cells (HCASMC)

243   cell-line and found that the T/NK biased peak-to-gene links were more enriched in the T

244   cells than the HCASMC cell line[36] (Supplementary Figure 14f). We next examined GTEx

245   eQTL mappings within our inferred peak-to-gene links, finding enrichment of eQTLs in

246   several functionally related categories such as Whole Blood and Lymphocytes

247   (Supplementary Figure 14g). To demonstrate the utility of these peak-to-gene links, we

248   linked differentially accessible regions to known leukemic genes such as the surface

249   protein *CD96,* the leukemic stem cell marker *IL1RAP*, the cytokine receptor *FLT3*, and

250   apoptosis regulator *MCL1* (Supplementary Figure 15a-d). Overall, these analyses,

251   support that our peak-to-gene links are highly enriched in immune regulation and across

252   other previously published linkage data sets[2,36].

253   Having established a high-quality set of peak-to-gene links, we aimed to identify

254   the set of malignant genes putatively regulated by *RUNX1*. First, we utilized our peak-to-

255   gene links to identify differential peaks linked to a differential gene within at least 2 MPAL

256   subpopulations. Next, we selected all linked differential accessibility sites that contain the

257   *RUNX1* motif. Finally, for each linked gene we combined all linked peaks to create a

258   differential linkage score (see methods) and compared this score to the proportion of

259   MPAL subpopulations that exhibited differential expression and accessibility in at least

260   one linked peak and target gene (a measure of how common this *RUNX1*-driven

261   dysfunction is across MPAL subsets) (**Figure 3c**). Using this approach, we found 732

262   genes putatively regulated by a RUNX1-containing distal element in at least 2 MPAL

263   subsets, and found that *CD69*, gene implicated in lymphocyte activation through initiation

264   of JAK/STAT signaling[37] and lymphocyte retention in lymphoid organs[38], was both highly

265    enriched in the calculated differential linkage score and was observed to be differentially

266    up-regulated in almost every MPAL subpopulation (**Figure 3d**). To further support *RUNX1*

267    predicted regulation of *CD69*[39], we incorporated T cell K27ac HiChIP[36], CRISPRa

268    screens[40], and *RUNX1* ChIP-seq[41] onto our multi-omic differential track. These orthogonal

269    data sets show *RUNX1* binding to linked distal regulatory regions (**Figure 3e**). Finally, by

270    using the 732 identified *RUNX1* target genes to stratify TCGA AML[42] patients by

271    expression, we observe significantly decreased survival (p-value = 0.023) in donors with

272    a high *RUNX1* target gene signature[42] (**Figure 3f**). This analysis suggests that *RUNX1* is

273    an important TF that putatively up-regulates a portion of the leukemic signature in MPAL

274    and potentially AML.

275        Collectively, this work establishes an experimental and analytical approach for

276    deconstructing cancer-specific features using integrative analysis of multiple single-cell

277    technologies. We find that MPAL malignant programs are largely conserved across

278    phenotypically heterogenous cells within individual patients; this  observation is consistent

279    with a previous report[16] that MPAL cells likely originate from a multipotent progenitor cell,

280    thereby sharing a common mutational landscape while populating different regions of the

281    hematopoietic tree. We used integrative single-cell analyses to further define putative TF

282    regulation of these malignant programs. We inferred that *RUNX1* acts as a potential

283    oncogene in MPAL, regulating malignant genes associated with poor survival. We

284    anticipate that similar approaches will be used in future studies to both identify the

285    differentiation status of different tumor types (i.e. identify the closest "normal" cell type) as

286    well as enable molecular dissection of molecular dysfunction in pathogenic cellular sub-

287    types, with the ultimate goal of identifying personalized therapeutic targets through

288    integrative single-cell molecular characterization.

289

290

291    **METHODS**

292    **Experimental Methods**

293    **Description of Healthy Donors**

294    PBMCs, BMMCs, and CD34+ bone marrow cells were obtained from healthy donors

295    (AllCells).

296

297    **Description of Leukemic Patients/Donors**

298    Patient samples were collected with informed consent prospectively under a protocol

299    approved by the Institutional Review Board (IRB) at Stanford University Medical Center

300    (Stanford IRB, 42949, 18329, and 6453). Peripheral blood and bone marrow aspirate

301    samples were processed by Lymphoprep (STEMCELL Technologies) gradient

302    centrifugation and fresh frozen in Bambanker media. Diagnostic flow cytometric performed

303    on bone marrow aspirate samples were analyzed. In all cases, a retrospective review of

304    clinical parameters, hemogram data, peripheral blood smears, bone marrow aspirates,

305    trephine biopsies, results of karyotype and flow cytometry studies was performed. Clinical

306    follow-up information was obtained by retrospective review of the medical record charts.

307    Cases were classified using the 2016 WHO classification of hematopoietic and lymphoid

308    neoplasms[5].

309

310    **CITE-seq (combined single-cell antibody derived tag and RNA sequencing)**

311    Combined single-cell RNA and antibody derived tag sequencing (CITE-seq) was

312    performed as previously reported[6] using the (version 2) Chromium Single Cell 3' Library

313    and Gel Bead  Kit  (Cat # 120237, 10X Genomics). Six thousand cells were targeted for

314    each sample. Oligo-coupled antibodies were obtained from Biolegend indexed by PCR

315    (10 cycles) with custom barcodes (see Supplementary Table 3), quantified by PCR using

316    a PhiX Control v3 (Illumina, Cat #FC-110-3001) standard curve, and then sequenced on

317    an Illumina NextSeq 550 together with scRNA-seq at no more than 60% of the total library

318    composition (1.5pM loading concentration, 26 x 8 x 0 x 98 bp read configuration).

319

320    **Single-cell ATAC-seq (scATAC-seq).**

321    Single-cell ATAC-seq targeting four thousand cells per sample was performed using a

322    beta-version of Chromium Single Cell ATAC Library and Gel Bead Kit (Cat # 1000110,

323    10X Genomics). Each sample library was uniquely barcoded and  quantified by PCR using

324    a PhiX Control v3 (Illumina, Cat #FC-110-3001) standard curve. Libraries were then

325    pooled and loaded on a NextSeq 550 Illumina sequencer (1.4pM loading concentration,

326    33 x 8 x 16 x 33 bp read configuration) and sequenced to either  90% saturation or 30,000

327    unique reads per cell on average.

328

329    **Whole-Exome Sequencing of Leukemic Patients/Donors**

330    Genomic DNA was extracted from diagnostic peripheral blood mononuclear cells or bone

331    marrow samples using Zymo Clean and Concentrator Kit. Library construction (Agilent

332    SureSelect Human All Exon kit), quality assessment, and 150-bp paired-end sequencing

333    (HiSeq4000) were performed by Novogene (Beijing, China). Reads with adapter

334    contamination, uncertain nucleotides, and paired reads with >50% low-quality nucleotides

335    were discarded. Paired-end reads were then aligned to the reference genome (GRCh37)

336    using BWA software. Genome Analysis Toolkit (GATK) was used to ignore duplicates with

337    Picard-tool. Filtered variants (SNP, INDELs) were identified using GATK HaplotypeCaller

338    and variantFiltration. Variants obtained from initial analysis were further compared to

339    dbSNP and 1000 Genomes database. Finally, missense, stopgain and frameshift

340    mutations were compared against a custom panel of 300 genes that are recurrently

341    mutated in hematologic malignancies as described previously[16,17].

342

343 **Analytical Methods**

344 **FACS Analysis**

345 Flow cytometry was performed on a FACSCalibur or FACSCanto II (Becton Dickinson,

346 San Jose, Ca, USA) cytometer using commercially available antibodies (Supplementary

347 Table 2). Lymphocytes were identified by low side-scatter and bright CD45 expression.

348 The gate was validated by backgating on CD3-positive or CD19-positive events. Blasts

349 were identified by low side-scatter and dim CD45 expression. The gate was further

350 assessed by backgating on CD34-positive events. Gates were drawn by additionally using

351 isotype controls and internal positive and negative controls.

352

353 **scADT-seq Analysis**

354 Raw sequencing data were converted to fastq format using bcl2fastq (Illumina, version

355 v2.20.0.422). ADTs were then assigned to individual cells and antibodies (see reference

356 antibody barcodes in Supplementary Table 3) allowing for 2 and 3 barcode mismatches,

357 respectively. Unique molecular counts for each cell and antibody were then generated by

358 counting only barcodes with a unique molecular identifier. PBMC and BMMC ADT count

359 data were transformed using the centered log ratio (CLR) as previously described[6]. PBMC

360 and BMMC cells were visualized in two dimensions using uwot's implementation of

361 UMAP[43] in R (n_neighbors = 50, min_dist = 0.4).

362

363 ***scATAC-seq Analytical Methods***

364 **scATAC-seq Processing**

365 Raw sequencing data was converted to fastq format using cellranger atac mkfastq (10x

366 Genomics, version 1.0.0). Single-cell RNA-seq reads were aligned to the GRCh37 (hg19)

367 reference genome and quantified using cellranger count (10x Genomics, version 1.0.0).

368

369 **scATAC-seq Quality Control**

370 To ensure that each single-cell was both adequately sequenced and had high signal to

371 background, we filtered cells with less than 1000 unique fragments and enrichment at

372 transcription start sites (TSS) was below 8. To calculate a TSS enrichment[2], briefly Tn5

373 corrected insertions were aggregated +/- 2,000 bp relative (TSS strand-corrected) for each

374 unique TSS genome wide. This profile was normalized to the mean accessibility +/- 1,900-

375 2,000 bp from the TSS, smoothed every 51bp, and the maximum smoothed value was

376 reported as TSS enrichment in R. We estimate that the multiplet percentage for this study

377 was around 4%[7].

378

379 **scATAC-seq Counts Matrix**

380 To construct a counts matrix for each cell by each feature (window or peaks), we read

381 each fragment.tsv.gz fill into a Genomic Ranges object. For each Tn5 insertion, the "start"

382 and "end" of the ATAC-fragments, we used  findOverlaps" to find all overlaps with the

383 feature by insertions. Then we added a column with the unique id (integer) cell barcode to

384 the overlaps object and fed this into a sparseMatrix in R. To calculate the fraction of

385 reads/insertions in peaks, we used the colSums of the sparseMatrix and divided it by the

386 number of insertions for each cell id barcode using "table" in R.

387

388 **scATAC-seq Union Peak Set from Latent Semantic Indexing Clustering**

389 We adapted a previous workflow for generating a union peak set that will account for

390 diverse subpopulation structure[2,9,10]. First, we created 2.5kb windows genome wide using

391 "tile(hg19chromSizes, width = 2500)" in R. Next, a cell by 2.5kb window sparse matrix was

392 constructed as described above. The top 20,000 accessible windows were kept and the

393 binarized matrix was transformed with the term frequency-inverse document frequency

394 ("TF-IDF") transformation[8]. Briefly we divided each index by the colSums of the matrix to

395     compute the cell "term frequency". Next we multiplied these values by log(1 + ncol(matrix)

396     / rowSums(matrix)) which represents the "inverse document frequency". This

397     normalization resulted in a TF-IDF matrix that was then used as input to irlba's singular

398     value decomposition (SVD) implementation in R. The 2nd-25$^{th}$ SVD dimensions (1$^{st}$

399     dimension is correlated with cell read depth[15]) were used for creating a Seurat object and

400     identified clusters using Seurat's SNN graph clustering (v2.3.4) with "FindClusters" with a

401     default resolution of 0.8. If the minimum cluster size was below 200 cells, the resolution

402     was decreased until this criterion was reached leading to a final resolution of $0.8^{N}$ (where

403     N represents the iterations until the minimum cluster size is 200 cells). For each cluster,

404     peak calling was performed on Tn5-corrected insertions (each end of the Tn5-corrected

405     fragments) using the MACS2 callpeak command with parameters "--shift -75 --extsize 150

406     --nomodel --call-summits --nolambda --keep-dup all -q 0.05." The peak summits were then

407     extended by 250bp on either side to a final width of 501bp, filtered by the ENCODE hg19

408     blacklist (https://www.encodeproject.org/annotations/ENCSR636HFF/), and then filtered

409     to remove peaks that extend beyond the ends of chromosomes.

410

411     Overlapping peaks called were handled using an iterative removal procedure as

412     previously described[2]. First, the most significant (MACS2 score) extended peak summit is

413     kept and any peak that directly overlaps with that significant peak is removed. This process

414     re-iterates to the next most significant peak until all peaks have either been kept or

415     removed due to direct overlap with a more significant peak. The most significant 200,000

416     extend peak summits for each cluster were quantile normalized using

417     "trunc(rank(v))/length(v)" in R (where v represents the vector of MACS2 peaks scores).

418     These cluster peak sets were then merged and the previous iterative removal procedure

419     was used. Lastly, we removed any peaks whose nucleotide content had any "N"

420     nucleotides and any peaks mapping to chrY.

421

**422**    **scATAC-seq-centric Latent Semantic Indexing clustering and visualization**

**423**    scATAC-seq clustering was performed by adapting the strategy of Cusanovich et. al[9,10],

**424**    to compute the term frequency-inverse document frequency ("TF-IDF") transformation.

**425**    Briefly we divided each index by the colSums of the matrix to compute the cell "term

**426**    frequency." Next, we multiplied these values by log(1 + ncol(matrix) / rowSums(matrix)),

**427**    which represents the "inverse document frequency." This resulted in a TF-IDF matrix that

**428**    was used as input to irlba's singular value decomposition (SVD) implementation in R. The

**429**    first 50 SVD dimensions were used as input into a Seurat object and initial clustering was

**430**    performed using Seurat's (v2.3.4) SNN graph clustering "FindClusters" with a resolution

**431**    of 1.5 (25 SVD dimensions for Healthy Hematopoiesis and 50 for Healthy Hematopoiesis

**432**    and MPALs). We found that in some cases, that there was batch effect between

**433**    experiments. To minimize this effect, we identified the top 50,000 variable peaks across

**434**    the initial clusters (summed cell matrix for each cluster followed by edgeR logCPM

**435**    transformation[44]). These 50,000 variable peaks were then used to subset the sparse

**436**    binarized accessibility matrix and recomputed the "TF-IDF" transform. We used singular

**437**    value decomposition on the TF-IDF matrix to generate a lower dimensional representation

**438**    of the data by retaining the first 50 dimensions. We then used these reduced dimensions

**439**    as input into a Seurat object and then final clusters were identified by using Seurat's

**440**    (v2.3.4) SNN graph clustering "FindClusters" with a resolution of 1.5 (50 SVD dimensions

**441**    for Healthy Hematopoiesis and 50 for Healthy Hematopoiesis and MPALs). These same

**442**    reduced dimensions were used as input to uwots implementation of UMAP (n_neighbors

**443**    = 55, n_components = 2, min_dist = 0.45) and plotted in ggplot2 using R. We merged

**444**    scATAC-seq clusters from a total of 36 clusters for hematopoiesis to 26 final clusters that

**445**    best agreed with the scRNA-seq clusters (included in Supplemental Data). The objective

446    of this analysis is to optimize feature selection, that minimizes batch effects, and enable

447    projection of future data into the same manifold as described further below.

448

449    **scATAC-seq Visualization in Genomic Regions**

450    To visualize scATAC-seq data, we read the fragments into a GenomicRanges object in R.

451    We then computed sliding windows across each region we wanted to visualize every 100

452    bp "slidingWindows(region,100,100)". We computed a counts matrix for Tn5-corrected

453    insertions as described above and then binarized this matrix. We then returned all non-

454    zero indices (binarization) from the matrix (cell x 100bp intervals) and plotted them in

455    ggplot2 in R with "geom_tile". For visualizing aggregate scATAC-seq data, the binarized

456    matrix above was summed and normalized. Scale factors were computed by taking the

457    binarized sum in the global peakset and normalizing to 10,000,000. Tracks were then

458    plotted in ggplot in R.

459

460    **chromVAR**

461    We measured global TF activity using chromVAR[15]. We used the cell by peaks and the

462    CIS-BP motif (from chromVAR motifs "human_pwms_v1") matches within these peaks

463    from motifmatchr. We then computed the GC bias-corrected deviations using the

464    chromVAR "deviations" function. We then computed the GC bias-corrected deviation

465    scores using the chromVAR "deviationScores" function.

466

467    **Gene Activity Scores using Cicero and Co-Accessibility**

468    We calculated gene activities using the R package Cicero[14]. Briefly, we used the sparse

469    binary cell by peaks matrix and created a cellDataSet, detectedGenes, and

470    estimatedSizeFactors. We then created a "cicero_cds" with k=50 and the

471    "reduced_coordinates" being the latent semantic indexing singular value decompositions

472    coordinates (Hematopoiesis = 25, Hematopoiesis and MPALs = 50). This function returns

473    aggregated accessibility across groupings of cells based on nearest-neighbor rules from

474    FNN. We then identified all peak-peak linkages that were within 250 kb by resizing the

475    peaks to 250 kb and 1bp and using "findOverlaps" in R. We calculated the pearson

476    correlation for each unique peak-peak link and created a connections data.frame where

477    the first column is peak_i and the second column is peak_j and third coaccessibility

478    (pearson correlation). We then created a gene data.frame from the TxDb

479    "TxDb.Hsapiens.UCSC.hg19.knownGene" in R. We then resized each gene from its TSS

480    and created a window +/- 2.5 kb centered at the TSS and then annotated the "cicero_cds"

481    using "annotate_cds_by_site". We then calculated gene activities with

482    "build_gene_activity_matrix" (coaccess cutoff of 0.35). Lastly we normalized the gene

483    activities by using "normalize_gene_activities" and the read depth of the cells. We then

484    log normalized these gene activities scores for interpretability by computing

485    "log2(GA*1,000,000 +1)".

486

487    ***scRNA-seq Analytical Methods***

488    **scRNA-seq Processing**

489    Raw sequencing data was converted to fastq format using cellranger mkfastq (10x

490    Genomics, version 3.0.0). Single-cell RNA-seq reads were aligned to the GRCh37 (hg19)

491    reference genome and quantified using cellranger count (10x Genomics, version 3.0.0).

492    We kept genes that were present in both 10x gene transfer format (GTF) files v3.0.0 for

493    hg19 and hg38 (https://support.10xgenomics.com/single-cell-gene-

494    expression/software/release-notes/build). Mitochondrial and ribosomal genes were also

495    filtered prior to further analysis. Genes remaining after these filtering steps we refer to as

496    "informative" genes and enable cross genome comparison.

497

498 **scRNA-seq Quality Control**

499 We wanted to filter out cells whose transcripts were lowly captured and first plotted the

500 distribution of genes detected and UMIs for all experiments. Based on these plots we

501 chose to filter out cells that had less than 400 informative genes detected and 1000 UMIs.

502 In addition, to lower multiplet representation, we filtered cells with above 10,000 UMIs. We

503 estimate that the multiplet percentage for this study was around 6%[8]. We then plotted the

504 correlation for each replicate experiment and found high reproducibility.

505

506 **scRNA-seq-centric Latent Semantic Indexing clustering and visualization**

507 We initially tested out a few methods for clustering scRNA but settled on an approach that

508 enabled us to effectively capture the hematopoietic hierarchy without significant alteration

509 of transcripts expression. We first log-normalized the transcript counts by first depth

510 normalizing to 10,000 and adding a pseudo count prior to a log2 transform (log2(counts

511 per ten thousand transcripts + 1)). Next, we identified the top 3000 variable genes and

512 performed the TF-IDF transform on these 3000 genes. We then performed singular value

513 decomposition (SVD) on this transformed matrix keeping the first 25 dimensions and used

514 this as input to Seurat Shared Nearest Neighbor Clustering (v2.3.4) with an initial

515 resolution of 0.2. We then summed the individual clusters single cells and computed the

516 logCPM transformation, edgeR::cpm(mat,log=TRUE,prior.count=3), and then identified

517 the top 2500 variable genes across these initial clusters. These variable genes were then

518 used as input for a TF-IDF transform and then performed singular value decomposition

519 (SVD) on this transformed matrix keeping the first 25 dimensions and used this as input

520 to Seurat Shared Nearest Neighbor Clustering (v2.3.4) with an increased resolution of 0.6.

521 We then summed the individual clusters single cells and computed the logCPM

522 transformation, edgeR::cpm(mat,log=TRUE,prior.count=3), and then identified the top

523 2500 variable genes across these clusters. We then repeated this 1 more time (resolution

524   1.0) and then saved the final features and clusters. To align our clusters better with the

525   scATAC-seq data we merged a total of 26 clusters from 31 initial clusters (included in

526   Supplemental Data). These LSI dimensions were used as input to uwots implementation

527   of UMAP (n_neighbors = 35, n_components = 2, min_dist = 0.45) and plotted in ggplot2

528   using R. The objective of this analysis is to optimize feature selection, that minimizes batch

529   effects, and enable projection of future data into the same manifold as described further

530   below.

531

532   ***scATAC-seq and scRNA-seq Analytical Methods***

533   **LSI Projection for scATAC and scRNA-seq**

534   We designed the above analytical approach to clustering single cell data because it

535   optimized feature selection and enabled projection of new non-normalized data into low

536   dimension manifold. To enable this analyses, when computing the TF-IDF transformation

537   on the hematopoietic hierarchy, we kept the colSums, rowSums, and SVD from the

538   previous run and then when projecting new data into this subspace, we first identified

539   which row indices to zero out based on the initial TF-IDF rowSums. We then computed

540   the "term frequency" by dividing by the colSums in these features. Next, we computed the

541   "inverse document frequency" from the previous TF-IDF transform (diagonal(1+ncol(mat)/

542   rowSums(mat))) and computed the new TF-IDF transform. We then projected this TF-IDF

543   matrix into the SVD subspace previous generated. To do this calculation, we computed

544   the new coordinates by "t(TF_IDF) %*% SVD$u %*% diag(1/SVD$d)" where TF_IDF is

545   the transformed matrix and SVD is the previous SVD run using irlba in R (3.5.1). We then

546   computed the projected matrix by "SVD$u %*% diag(SVD$D) * t(V)" where V is the

547   projected coordinates above. For projecting bulk RNA-seq, we downsampled previously

548   published data to 5,000 reads in genes 100 times and then made a sparse matrix for

549   projection as single cell data. For projecting bulk scATAC-seq, we downsampled

550    previously published data to 10,000 reads in peaks 100 times and then made a binary

551    sparse matrix for projection as single cell data.

552

553    **Classification of MPAL single cells with scATAC and scRNA-seq**

554    We wanted to classify MPAL single cells based on their disease state and hematopoietic

555    progression. First, we determined which cells were healthy-like and disease-like. To do

556    this analysis, we clustered all of the healthy hematopoietic cells with the MPAL of interest

557    using our LSI workflow as described above (scRNA – 25 PCs, 1,000 variable genes and

558    Seurat SNN resolution of 0.2, 0.8 and 0.8; scATAC - 25 PCs, 25,000 variable peaks and

559    Seurat SNN resolution of 0.8 and 0.8). We then determined which clusters were "healthy-

560    like" if a high percentage (>80% for scRNA, >90% for scATAC) of the cells were from the

561    hematopoietic data. MPAL single cells belonging to these clusters were classified as

562    "healthy-like" and the remaining disease-like. We note that we did not detect significant

563    large-scale copy number amplifications with our previously described approach[7], and the

564    proportion of "disease-like" classified cells were consistent with our FACS estimation of

565    percent blast cells. In order to accurately characterize these MPAL "disease-like" by their

566    hematopoietic state, we established "hematopoietic compartments" across our scRNA

567    and scATAC-seq maps that broadly characterized the hematopoietic continuum. The

568    borders for these compartments were determined empirically using "fhs" in R, guided by

569    the initial clusters and agreement across the scRNA and scATAC-seq classifications. After

570    the hematopoietic continuum were classified, we then broadly classified the MPAL

571    "disease-like" cells based on their projected nearest neighbor in the UMAP subspace.

572    These classifications were used subsequently in differential analyses.

573

574    **Identifying differential features with scATAC and scRNA-seq**

575    To identify differential features for previously published AML data and MPALs, we

576    constructed a nearest neighbor healthy aggregate using the following approach. First, we

577    used FNN to identify the nearest 25 cells using "get.knnx(svdHealthy, svdProjected,

578    k=25)" based on Euclidean distance between the projected cells and hematopoietic cells

579    in LSI-SVD space. For each projected population, we used a minimum of 50 and maximum

580    of 500 cells (random sampling) as input. Next, we took the unique of all hematopoietic

581    single-cells and if this number was greater than 1.25 times the number of the projected

582    populations, we took the nearest 24 cells and repeated this procedure until this criterion

583    was met. Then the projected population and non-redundant hematopoietic cells were

584    downsampled to an equal number of cells (maximum 500). For scATAC-seq, we binarized

585    the matrix for both the projected populations and hematopoietic matrices. Next, we scaled

586    the sparse matrices to 10,000 total counts for scRNA and 5,000 total promoter counts for

587    scATAC-seq (promoter peaks defined as peaks within 500 bp of TSS from hg19 10x v3.0.0

588    gtf file). Next, we computed row-wise t-tests for each feature. We then calculated the FDR

589    using p.adjust(method="fdr"). We then computed the log2 mean and log2 fold changes for

590    each feature. We chose these parameters based on Soneson et al., study comparing

591    analytical methods for differential expression[45]. For scRNA-seq, differential expression

592    was determined by FDR < 0.01 and absolute log2 fold changes greater than 0.5. For

593    scRNA-seq, differential expression was determined by FDR < 0.05 and absolute log2 fold

594    changes greater than 0.05.

595

596    To identify differential genes for bulk leukemia RNA-seq, we downsampled the gene

597    counts to 10,000 counts randomly for 250 times. We then projected and used the above

598    framework to resolve differential genes with log2 fold change > 3 and FDR < 0.01. We

599    then removed genes that were differential in 33% or higher of the normal samples to

600    attempt to capture biased genes. In addition, we further removed genes differential in 50%

601    or higher of the leukemia samples. This filtering biases our identified malignant genes to

602    those variable across the leukemic types vs conserved across all leukemic types. We then

603    took the average malignancy for each remaining gene for each leukemic type and used

604    the top 300 variable malignant genes across the leukemic types for heatmap and LSI. For

605    computing differential LSI, we binarized each gene being malignant or not for the 300

606    variable malignant genes and computed the TF-IDF transform followed by SVD (LSI). We

607    then visualized this in 2 dimensions using uwot's implementation of UMAP (50 SVD

608    dimensions, n_neighbors = 50, min_dist = 0.005).

609

610    **Matching scATAC-scRNA-seq pairs using Seurat Canonical Correlation Analyses**

611    We wanted to be able to integrate our epigenetic and transcriptomic data and built off of

612    previous approaches for integration[10,35]. We found the approach that worked best for our

613    integrative analyses was using Seurat's Canonical Correlation Analysis. We performed

614    integration for each biological group separately because (1) it improved alignment

615    accuracy and (2) required much less memory. First, for both the Gene Activity Scores

616    matrix and scRNA matrix we created a Seurat Object "CreateSeuratObject", then

617    normalized with "NormalizeData", and found the top 2000 variable genes/activities ranked

618    by dispersion with "FindVariableGenes". We then defined the union of the top 2000

619    variable genes from scRNA-seq and gene scores from scATAC-seq and found this

620    increased the concordance downstream (defined by cluster to cluster mapping in

621    hematopoiesis and single cell spearman correlations). These genes were then used for

622    running Canonical Correlation analysis using "RunCCA" with the number of cc's to

623    compute as 25. We then calculated the explained variance using "CalcVarExpRatio"

624    grouping by each of the individual experimental protocols scATAC (Gene Activity Scores)

625    and scRNA. We then filtered cells where the variance explained by CCA is less than 2 fold

626    compared to PCA. We then Aligned the subspaces with "AlignSubspace" and 25

627     dimensions to align with reduction.type = "cca" and grouping.var = "protocol". We then

628     identified for each scATAC cell the nearest scRNA cell based on minimizing the euclidean

629     distance. We then created a UMAP using the aligned CCA coordinates as input into uwot's

630     UMAP implementation with n_neighbors = 50, min_dist = 0.5, metric = "euclidean" and

631     then plotting with ggplot2 in R. To enable more robust correlation based downstream

632     analyses, we used our initial KNN groupings (nGroups = 4998, KNN = 50) from Cicero[14]

633     to group scATAC accessibility, Gene Activity Scores, scRNA closest neighbor and

634     chromVAR[15] deviation scores.

635

636     **Peak-To-Gene Linkage**

637     Cicero[14] allows us to infer Gene Activity Scores by linking distal correlated ATAC peaks

638     to the promoter peak. While this measure is extremely useful, it does not actually mean it

639     is correlated to the gene expression. To circumvent this limitation, we used our grouped

640     scATAC and grouped linked scRNA-seq to identify peak-to-gene links. First we log-

641     normalized the accessibility and gene expression with log2(Counts Per 10,000 + 1) and

642     then we resized each of the gene GRanges to the start using resize(gr,1,"start") and then

643     resizing the start to a +- 250kb window using resize(gr, 2 * 250000 + 1, "center"). We then

644     overlapped all ATAC-seq peaks using "findOverlaps" to identify all putative peak-to-gene

645     links. We then split the aggregated ATAC and RNA matrices by whether majority of the

646     cells were from MPAL or Hematopoietic single cells. We then correlated the peaks and

647     genes for all putative peak-to-gene links. We used a previously described approach for

648     computing a null correlation based on *trans* correlations (correlating peaks and genes not

649     on the same chromosome)[2]. Briefly, for each chromosome 1000 peaks not on the same

650     chromosome are identified and correlated to every gene on that chromosome. Each

651     putative peak-to-gene correlation is converted into a z-score by using the mean and sd of

652     the null *trans* correlations. These are then converted to p-values and adjusted for multiple

653    hypothesis using the benjamini Hochberg correction "p.adjust" in R. We retained links

654    whose correlation (Pearson) was above 0.35 and FDR < 0.1, same correlation cutoff as

655    co-accessibility in Cicero[14], in either MPAL or Hematopoietic aggregations. We then kept

656    all peak-to-gene links that were greater than 2.5kb in distance. We identified peak-to-gene

657    links that are only present in hematopoiesis, MPALs or both. To visualize the peak-to-

658    gene links we plotted all of them as a heatmap with ComplexHeatmap. To determine the

659    column order we first computed PCA for the first 25 PCs using irlba. We then computed

660    Seurat[11] Shared Nearest Neighbor clustering with a resolution of 1 and then computed the

661    cluster means. We then computed the order of these clusters using hclust and the

662    dissimilarity 1-R as the distance. Next we then iterated through each cluster and

663    performed hclust with the dissimilarity calculations to get a final column order. The peak-

664    to-gene links were grouped by k-means clustering with 10 input centers 100 iterations and

665    10 random starts for healthy, disease and the overlapping links. We did this bi-clustering

666    because it enabled us to plot smaller rasterized chunks of the heatmap without

667    overwhelming the memory and put the individual rasterized k-means clusters together

668    post analysis.

669

670    **Peak-To-Gene links enrichment with GTEx eQTLs**

671    We adopted a previous approach for identifying the enrichment of our peak-to-gene links

672    in GTEx eQTL data. Briefly, we downloaded GTEx eQTL data (version 7) from

673    https://gtexportal.org/home/datasets and the *.signif_variant_gene_pairs.txt.gz files were

674    used. We in addition downloaded gencode v19 (matched to these eQTLs) and identified

675    all gene starts and identified all nearest gene starts to each peak and eQTL using

676    "distanceToNearest". We filtered all eQTLs that were further than 250kb from their

677    predicted gene to be consistent with our linkage approach. To calculate a conservative

678    overlap enrichment,  we further pruned all eQTL links that were to its nearest gene. We

679    then created a null set (n = 250) of peak-to-gene links by randomly selecting distal ATAC-

680    seq peak-to-gene links (within 250 kb) that are distance matched to the links tested at 5kb

681    resolution. We then calculated a z-score and enrichment for each peak-to-gene link set

682    compared to the null set and calculated an FDR using p.adjust(method = "fdr").

683

684    **Peak-To-Gene links enrichment with K27ac HiChIP metaV4C**

685    We wanted to determine the specify of our peak-to-gene links in published chromatin

686    conformation data as previously described. We downloaded previously published Naive T

687    cell and HCASMC K27ac HiChIP data. We then identified within each peak-to-gene links

688    subset the peaks that were most biased to T/NK cells. To do this analysis, we calculated

689    the z-score for each peak in the peak-to-gene links removed all links below 100kb and

690    floored each peak coordinate (start or end) to its nearest 10kb window. We then ranked

691    these links by the z-score for the peak, deduplicated the links at 10kb resolution and kept

692    the top 500 remaining peak-to-gene links. Next, we used juicer dump (no normalization

693    "NONE") at 10kb resolution for each chromosome in the ".hic" file. Then we read each

694    chromosomes into an individual "sparseMatrix" in R. We then scaled the sparse matrices

695    such that the total cis interactions summed up to 10 million PETs. Then, for each peak-to-

696    gene link, the upstream or downstream window (Column or Row) (whether the peak was

697    upstream or downstream of the gene promoter) was identified. To scale each interactions

698    distance for interpretability, we linearly interpolated the data to be on a -50-150% scale to

699    visualize the focal interaction. The mean interaction signal was reported and repeated for

700    both replicates. The mean and sd across both replicates were calculated and plotted with

701    ggplot in R.

702

703    **Identifying TF Malignant Target Genes and Survival Anlaysis**

704    We wanted to create a framework for identifying TFs that potentially directly regulate

705    malignant genes. To do this analysis, we first identified a set of transcription factors whose

706    hypergeometric enrichment in differential peaks were high across the MPAL

707    subpopulations (Comparing up-regulated peaks vs all peaks) and were identified as being

708    transcriptionally correlated with their motif's accessibility (see above). Next for a given TF

709    and all identified peak-to-gene links, we further subsetted these links by those containing

710    the TF motif. Then for each MPAL subpopulation, we determined for each peak-to-gene

711    link if both the peak and gene are up-regulated. Then for each gene, we gave a binary

712    score whether or not that MPAL subpopulation has at least one differential peak-to-gene

713    link (whose peak and gene are differentially up-regulated) and report the proportion of

714    subpopulations that were up-regulated. In addition, for each gene that has at least 1

715    differential peak-to-gene links we summed their squared correlation $R^2$ and report that as

716    the differential linkage score. We kept all genes that had least 1 MPAL subpopulation with

717    corresponding differential peak-to-gene links.

718

719    For survival analysis, we downloaded the RPKM TCGA-LAML data[42] (https://tcga-

720    data.nci.nih.gov/docs/publications/laml_2012/laml.rnaseq.179_v1.0_gaf2.0_rpkm_matrix

721    .txt.tcgaID.txt.gz). We downloaded the survival data from Bioconductor RTCGA.clinical

722    ("patient.vital_status") and matched using TCGA IDs the RPKM expression. Next, we took

723    all genes that were identified as target genes for *RUNX1* (n = 732), and computed row-

724    wise z-scores for each gene. Next, we took the column means of this matrix to get an

725    average z-score across all *RUNX1* target genes. We then identified the top 33% and

726    bottom 33% of donors based on this expression. We computed the p-value using the R

727    package survival "survfit(Surv(times,patient.vital_status)~Runx1_TG_Expression,

728    LAML_Survival)". We plotted the Kaplan-Meier curve using the R package survminer

729    "ggsurvplot" in R.

730

731 **FIGURE LEGENDS**

732

733 **Figure 1. Multi-omic epigenetic and phenotypic analysis of human hematopoiesis.**

734 **a**, Schematic of multi-omic profiling of chromatin accessibility, transcription, and cell

735 surface antibody abundance on healthy bone marrow and peripheral blood mononuclear

736 cells using scATAC-seq and CITE-seq (combined single-cell RNA and antibody derived

737 tag sequencing).

738 **b**, scRNA-seq LSI UMAP projection of 35,882 single cells of healthy hematopoiesis.

739 **c**, scATAC-seq LSI UMAP projection of 35,038 single cells of healthy hematopoiesis.

740 **d**, Surface marker overlay on single-cell RNA UMAP (**b**) of (Top) ADT antibody signal

741 (CLR normalized), (Middle) single-cell RNA, and (Bottom) log2 gene activity scores for

742 *CD3*, *CD14*, *CD19*, and *CD8A*.

743 **e**, Transcription factor overlay on single-cell ATAC UMAP (**c**) of (Top) TF deviations,

744 (Middle) gene activity scores, and (Bottom) single-cell RNA for *CEBPB*, *GATA1*, *TBX21*,

745 and *PAX5*.

746 **f-h**, Multi-omic tracks; (Top) average track of all clusters displayed, (Middle) binarized 100

747 random scATAC-seq tracks for each locus at 100bp resolution and (right) scRNA-seq log2

748 distribution of normalized expression for each cluster.

749 **f,** Multi-omic track of *CD14* (specific in these clusters for monocytes) across monocyte

750 development from HSC progenitor cells.

751 **g,** Multi-omic track of *CD19* (specific in these clusters for pre B cells) across B cell

752 development.

753 **h,** Multi-omic track of *PAX5* (specific in these clusters for pre B cells) across B cell

754 development.

755

756 **Figure 2. Multi-omic projection of MPALs into hematopoiesis identifies normal and**

757 **leukemic programs.**

758 **a,** Schematic for projection of MPAL single cells onto hematopoiesis for both scRNA-seq

759 and scATAC-seq classified into broad hematopoietic compartments.

760 **b,** (Left) MPAL single cell projections into hematopoiesis for both scRNA-seq and

761 scATAC-seq. (Right) The proportion of MPAL cells that were broadly classified as healthy

762 or disease and their respective hematopoietic compartment.

763 **c,** (Left) scRNA-seq heatmap of up-regulated genes log2 fold changes comparing MPAL

764 disease subpopulations to closest non-redundant normal cells. Differential genes were

765 clustered with k-means (k=10) based on their log2 fold changes. (Right) scATAC-seq

766 heatmap of differentially up-regulated accessible peaks log2 fold changes comparing

767 MPAL disease subpopulations to closest non-redundant normal cells. Differential peaks

768 were clustered with k-means (k=10) based on their log2 fold changes.

769 **d,** Pearson correlation of differentially up-regulated genes and peaks across all MPAL

770 subpopulations.

771 **e,** LSI UMAP of differentially up-regulated gene expression profiles across bulk

772 leukemias[16] and MPAL samples assayed in this study, colored by WHO 2016

773 classifications[5].

774 **f,** (Left) MA plot comparing the proportion of malignant (up-regulated) gene expression

775 profiles in AML and MPALs. The x-axis represents for each up-regulated gene, the

776 average proportion of AML and MPAL patient subpopulations broadly up-regulated (LFC

777 > 0.5). The y-axis represents for each up-regulated gene, the difference in the proportion

778 of MPAL and AML patient subpopulations up-regulated (LFC > 0.5). (Right) Genes that

779 are more malignant biased to AMLs, MPALs and conserved across both AMLs and

780 MPALs.

781

782 **Figure 3. Integrative scATAC and scRNA-seq analyses nominate putative**

783 **transcription factors that regulate leukemic programs**.

784 **a**, (Left) Hypergeometric TF motif enrichment FDR in differentially accessible peaks

785 across each k-means clusters identified in Figure 2c. TFs are also identified as being

786 differentially expressed and enriched in at least 3 MPAL hematopoietic compartments.

787 (Top) Number of accessible peaks in each k-means cluster. (Right) Proportion of

788 differentially up-regulated TF gene expression profiles across MPAL hematopoietic

789 compartments.

790 **b**, (Left) Schematic for alignment of scATAC and scRNA-seq data to link putative

791 regulatory regions to target genes. First, scATAC-seq data is converted from accessible

792 peaks to inferred gene activity scores using Cicero. Second, these gene activity scores

793 and scRNA-seq expression are aligned into a common subspace using Seurat's

794 Canonical Correlation Analyses. Third, each scATAC-seq cell is assigned its nearest

795 scRNA-seq neighbor. Fourth, ATAC-seq peaks within 2.5-250kb to a gene promoter are

796 correlated within the healthy hematopoietic and MPAL knn groupings. Lastly, significant

797 peak-to-gene links are identified by correlating peaks to genes on different chromosomes.

798 (Right) Heatmaps of 91,601 peak-to-gene links across hematopoiesis and MPALs. (Top)

799 peak-to-gene links that are identified only within hematopoiesis, (Middle) peak-to-gene

800 links that are unique to MPALs, and (Bottom) peak-to-gene links identified in both

801 hematopoiesis and MPALs.

802 **c**, Schematic for identifying genes that are putatively regulated by the transcription factor

803 of interest.

804 **d**, *RUNX1* putative target genes differentially up-regulated in at least 1 MPAL

805 subpopulations. The x-axis represents the proportion of MPAL subpopulations that are

806 differential in both scRNA-seq and a linked accessible peak. The y-axis represents the

807  cumulative linkage score between differentially up-regulated peaks linked to differentially

808  up-regulated genes.

809  **e,** *CD69* multi-omic differential track (Top) T cell Th17 K27ac HiChIP virtual4C of the *CD69*

810  locus, shading represents standard deviation between biological replicates (n = 2).

811  (Middle) Aggregated scATAC tracks showing MPAL disease subpopulations (red) and

812  aggregated nearest-neighbor healthy (grey). (Right) Distribution of log2 normalized

813  expression of *CD69* for MPAL disease subpopulations (red) and closest normal cells

814  (grey); black line represents the mean and asterisk denote significance (LFC > 0.5 and

815  FDR < 0.01). (Bottom) HL60 AML line ChIP-seq data across *CD69* locus, Jurkat CRISPRa

816  tiling screen across the *CD69* locus and *RUNX1* identified malignant peak-to-gene links.

817  **f,** Kaplan-Meier curve for TCGA AML patients (n=179) stratified by *RUNX1* putative target

818  genes top 33% vs bottom 33% (p-value = 0.023).

819

820  **SUPPLEMENTARY FIGURE LEGENDS**

821

822  **Supplementary Figure 1. Quality control of scRNA-seq data for hematopoiesis**

823  **samples.**

824  **a,** (Top) Number of cells passing filter for each experimental replicate (number of

825  informative genes > 400 and number of unique molecular identifiers (UMI) > 1000),

826  (Middle) number of informative genes detected per single cell and (Bottom) number of

827  unique molecular identified (UMI) transcripts.

828  **b,** Aggregated scRNA-seq one to one reproducibility plots for experimental replicates and

829  across experiments.

830  **c,** scRNA-seq biological cluster labels assigned to each cluster overlay on UMAP of

831  hematopoiesis.

832  **d**, scRNA-seq experimental sample labels overlay on UMAP of hematopoiesis.

833

834  **Supplementary Figure 2. Quality control of scATAC-seq data for hematopoiesis**

835  **samples.**

836  **a,** scATAC-seq cell filtering plot. The x-axis is the number of unique accessible fragments

837  and the y-axis is the enrichment of Tn5 insertions at transcription start sites, representing

838  the robust signal to background for each single cell.

839  **b,** Aggregated scATAC-seq fragment size distributions across individual experiments

840  demonstrating sub- , mono- and multi nucleosome spanning ATAC-seq fragments.

841  **c,** (Top) Number of cells passing filter for each experimental replicate (Unique fragments

842  > 1000 and TSS enrichment > 8), (Middle) log10 unique fragments, (Middle) fraction of

843  Tn5 insertions in the healthy hematopoietic union peak set, and (Bottom) enrichment at

844  transcription start sites.

845  **d,** Aggregated scATAC-seq one to one reproducibility plots for experimental replicates

846  and across experiments.

847  **e,** scATAC-seq biological cluster labels assigned to each cluster overlay on UMAP of

848  hematopoiesis.

849  **f**, scATAC-seq experimental sample labels overlay on UMAP of hematopoiesis.

850

851  **Supplementary Figure 3. Quality control of scADT-seq data for hematopoiesis.**

852  **a,** Proportion of scRNA-seq cells passing filter that were matched with corresponding

853  scADT data.

854  **b,** Aggregated scADT-seq one to one reproducibility plots for experimental replicates and

855  across experiments.

856  **c,** scADT-seq UMAP of bmmc and pbmc samples across 14 antibodies. scADT overlay of

857  experimental sample labels, *CD19*, *CD3*, *CD56*, *CD4*, *CD8A*, *CD14*, *CD16*, *CD45RA*,

858    *CD45RO, TIGIT* and *PD-1*. Color represents experimental labels or scADT-seq values

859    after CLR transformation.

860    **d,** Corresponding scRNA-seq biological cluster label overlay on the scADT-seq UMAP of

861    BMMC and PBMCs.

862

863    **Supplementary Figure 4. Validation of key marker genes for both scRNA-seq and**

864    **scATAC-seq for hematopoiesis.**

865    **a-h,** Multi-omic tracks; (Top) average track of all clusters displayed, (Middle) binarized 100

866    random scATAC-seq tracks for each locus at 100bp resolution and (right) scRNA-seq log2

867    distribution of normalized expression for each cluster, box-plot shows median and lower

868    and upper quartiles.

869    **a,** Multi-omic track of *GATA1* (specific in these clusters for Erythroid) for erythroid

870    development from HSC progenitor cells.

871    **b,** Multi-omic track of *GATA2* (specific in these clusters for Basophil) for erythroid

872    development from HSC progenitor cells.

873    **c,** Multi-omic track of *ELANE* (specific in these clusters for GMP/Neutrophil) for neutrophil

874    development from HSC progenitor cells.

875    **d,** Multi-omic track of *IRF8* (specific in these clusters for pDC) across pDC development

876    from HSC progenitor cells.

877    **e,** Multi-omic track of *SDC1* (specific in these clusters for Plasma cells) across B cell

878    development and plasma cells.

879    **f,** Multi-omic track of *CD1C* (specific in these clusters for cDC) across cDC development

880    from HSC progenitor cells.

881    **g,** Multi-omic track of *SELL* (specific in these clusters for Naive T cells vs memory, and

882    CD8 central memory vs CD8 effector memory) across NK and T cells.

883    **h,** Multi-omic track of *GZMB* (specific in these clusters for NK cells) across NK and T cells.

884

885    **Supplementary Figure 5. Diagnostic flow cytometry plots for MPALs 1-3.**

886    **a-c,** Diagnostic flow cytometry plots from three different MPAL cases gated on blasts area

887    (highlighted in red) and lymphocytes (highlighted in black) from CD45 and side scatter

888    area (SSC-A).

889    **a**, MPAL 1 shows classic bilineal phenotype with both T-lymphoblasts (cCD3-positive and

890    CD7-positve) and myeloid blasts (MPO-positive and CD33-positive).

891    **b**, MPAL 2 demonstrates a more complex phenotype with both biphenotypic (single

892    population expressing lymphoid marker CD7 and myeloid marker CD33) and bilineal T-

893    Myeloid patterns (subpopulation expressing monocytic markers CD64, CD33, and CD14).

894    **c**, MPAL 3 demonstrates a classic biphenotypic case with coexpression of both T-lineage

895    markers (cCD3-positive) and myeloid markers (MPO-positive).

896

897    **Supplementary Figure 6. Diagnostic flow cytometry plots for MPALs 4-5R.**

898    **a-c,** Diagnostic flow cytometry plots from three different MPAL cases gated on blasts area

899    (highlighted in red) and lymphocytes (highlighted in black) from CD45 and side scatter

900    area (SSC-A).

901    **a**, MPAL4 demonstrates a classic bilineal B/M phenotype expressing B-lineage markers

902    (CD79a and CD19-positive) and myeloid markers (MPO-positive and CD33-positive).

903    **b**, MPAL5 demonstrates a more complicated phenotype with a subpopulation of blasts

904    expressing T-lineage markers (cCD3-positive and CD7-positive) and a subpopulation

905    expressing myeloid marker MPO.

906    **c**, MPAL5R post-treatment relapse of MPAL5. Flow cytometry reveals expansion of the T-

907    lymphoblastic subpopulation (cCD3-positive, TdT-positive population) following

908    chemotherapy.

909   **d**, High-confidence mutations detected in 5 MPAL cases by whole exome sequencing.

910   Missense mutations are shown in blue, frameshift deletions are shown in yellow, stopgain

911   mutations are shown in purple, frameshift insertions are shown in orange, and

912   nonframeshift deletions are shown in dark gray.

913

914   **Supplementary Figure 7. Quality control of scRNA-seq and scATAC-seq data for**

915   **MPAL samples.**

916   **a,** (Top) Number of cells passing filter for each experimental replicate (number of

917   informative genes > 400 and number of unique molecular identifiers (UMI) > 1000),

918   (Middle) number of informative genes detected per single cell and (Bottom) number of

919   unique molecular identified (UMI) transcripts.

920   **b,** Aggregated scRNA-seq one to one reproducibility plots for experimental replicates and

921   across experiments.

922   **c,** scATAC-seq cell filtering plot. The x-axis is the number of unique accessible fragments

923   and the y-axis is the enrichment of Tn5 insertions at transcription start sites, representing

924   the robust signal to background for each single cell.

925   **d,** Aggregated scATAC-seq fragment size distributions across individual experiments

926   demonstrating sub- , mono- and multi nucleosome spanning ATAC-seq fragments.

927   **e,** (Top) Number of cells passing filter for each experimental replicate (Unique fragments

928   > 1000 and TSS enrichment > 8), (Middle) log10 unique fragments, (Middle) fraction of

929   Tn5 insertions in the MPAL union peak set, and (Bottom) enrichment at transcription start

930   sites

931   **f,** Aggregated scATAC-seq one to one reproducibility plots for experimental replicates and

932   across experiments.

933

934    **Supplementary Figure 8. Evaluation of LSI projection workflow for previously**

935    **published bulk and single-cell hematopoietic data sets across different platforms.**

936    **a,** Overview of LSI projection workflow. Briefly, using information from TF-IDF transform,

937    singular value decomposition and UMAP of hematopoiesis enables projection of new data

938    into the same subspace.

939    **b**, LSI projection of downsampled previously published bulk sorted hematopoietic data

940    sets[18,20]. (Left) RNA-seq downsampled bulk projections for 49 samples (n=250

941    downsampled cells). (Right) ATAC-seq downsampled bulk projections for 90 samples

942    (n=250 downsampled cells).

943    **c**, LSI projection of downsampled previously published single-cell hematopoietic data sets

944    labeled by previous classifications[20–22]. (Left) scRNA-seq projections of previous study

945    healthy bone marrow cells (different platform and different aligned genome) colored by

946    previous classifications. (Right) scATAC-seq projections for healthy bone marrow and

947    peripheral blood samples (2 different platforms across 3 studies), colored by ground truth

948    isolated populations.

949

950    **Supplementary Figure 9. LSI projection of previously published healthy and AML**

951    **scRNA-seq identifies malignant programs across AML subpopulations.**

952    **a**, (Left) Schematic of LSI projection. (Right) Initial projection of all AML malignant single-

953    cells colored by previous classifications[19].

954    **b**, Re-classification of scRNA-seq AML single-cells based on closest normal cells in

955    healthy hematopoiesis (See Methods). Broader re-classification increases the number of

956    cells per category for improved power in differential analyses. LSI projection for each

957    individual AML samples onto scRNA-seq healthy hematopoiesis colored by re-

958    classifications (denoted is the sample id and number of cells).

959    **c**, K-means differential scRNA-seq heatmap (k = 10), colored by log2 fold change,

960    comparing each AML sample subpopulations (classifications) vs their closest normal bone

961    marrow cells from the same study[19].

962

963    **Supplementary Figure 10. scADT-seq overlay of MPALs projected onto the**

964    **hematopoietic hierarchy**

965    **a**, (Left) Projected MPALs colored by hematopoietic compartments. (Right) scADT-seq

966    overlay of *CD7*, *CD33*, *CD14*, *CD4 and CD19* on MPAL single cells LSI projected onto

967    hematopoiesis.

968

969    **Supplementary Figure 11. Visualization of differential genes and accessible peak**

970    **regions.**

971    **a**, Top conserved differential genes across the MPAL hematopoietic compartments.

972    **b**, Top conserved differential transcription factors across the MPAL hematopoietic

973    compartments.

974    **c**, KEGG pathway enrichment in differential RNA k-means 2, 3, 4, and 10 (Figure 2c).

975    **d-e,** Multi-omic differential tracks (Left) scATAC tracks showing MPAL disease

976    subpopulations (red) closest normal cells (grey). (Right) Distribution of log2 normalized

977    expression for MPAL disease subpopulations (red) and closest normal cells (grey); black

978    line represents the mean and asterisk denote significance (LFC > 0.5 and FDR < 0.01).

979    **d,** Multi-omic differential track of *CDK11A*, up-regulated in MPALs 1, 2, 5 and 5R.

980    **e,** Multi-omic differential track of *CDKN2A*, up-regulated in MPALs 1, 2, 3, 4, and 5.

981

982    **Supplementary Figure 12. Seurat canonical correlation analysis alignment of**

983    **scRNA and scATAC-seq  hematopoietic and MPAL samples.**

984 **a**, Schematic of LSI projection of downsampled bulk leukemia RNA-seq onto healthy

985 hematopoiesis.

986 **b**, Representative downsampled LSI projections (n=250) for B-ALLs, non-ETP T-ALLs,

987 ETP T-ALLs, AMLs, T/M MPALs and B/M MPALs from previous studies[16].

988 **c**, LSI UMAP of differentially up-regulated gene expression profiles across bulk

989 leukemias[16] and MPAL samples assayed in this study, colored by cytogenetics.

990 **d**, Binary heatmap of variable malignant genes across leukemia classifications. Each cell

991 in the heatmap is colored whether the gene was identified as malignant for the leukemic

992 sample.

993

994 **Supplementary Figure 13. Seurat canonical correlation analysis alignment of**

995 **scRNA and scATAC-seq hematopoietic and MPAL samples.**

996 **a**, UMAP of CCA alignment of scATAC-seq using Cicero gene activity scores and scRNA-

997 seq for (Left) bone marrow, (Middle) CD34+ enriched bone marrow, (Right) peripheral

998 blood.

999 **b**, UMAP of CCA alignment of scATAC-seq using Cicero gene activity scores and scRNA-

1000 seq for MPAL samples.

1001

1002 **Supplementary Figure 14. Evaluation of scRNA and scATAC-seq alignment and**

1003 **peak-to-gene linkage across hematopoiesis and MPAL samples.**

1004 **a**, Spearman rank correlation between scATAC-seq Cicero gene activity scores to scRNA-

1005 seq for each mapped cell within across all biological experiments.

1006 **b**, Pearson correlation of CCA scRNA and scATAC-seq nearest-neighbors. The cutoff (R

1007 > 0.45) for high quality nearest neighbor mappings is shown.

1008 **c**, (Left) UMAP of scATAC-seq hematopoiesis colored by scATAC-seq clusters. (Right)

1009 UMAP of scATAC-seq hematopoiesis colored by mapped scRNA-seq clusters.

1010    **d**, Confusion matrix of initial clusters for mapped scRNA-seq to scATAC-seq clusters for

1011    hematopoiesis (Figure 1b-c).

1012    **e**, (Left) Distribution of peak-to-gene distances. (Left-Middle) Distribution of number of

1013    peaks mapped per gene (median = 6). (Right-Middle) Distribution of number of genes

1014    mapped per peak (median = 1). (Right) Distribution of number of genes skipped for peak-

1015    to-gene links (median = 2).

1016    **f**, MetaV4C plots of K27ac HiChIP in Naive T and HCASMC cells for top 500 biased T/NK

1017    (broad classification) peak-to-gene links that are identified only in  healthy hematopoiesis.

1018    Shading indicates standard deviation between replicate experiments (n = 2).

1019    **g**, Peak-to-genes enrichment in GTEx eQTLs over a permuted background distance-

1020    matched set (n=250) for the union set of peak-to-gene links.

1021

1022    **Supplementary Figure 15. Peak-to-gene links nominate putative regulatory regions**

1023    **that nominate key leukemic genes.**

1024    **a-d,** Multi-omic differential track; (Middle) Aggregated scATAC tracks showing MPAL

1025    disease subpopulations (red) and closest normal cells (grey). (Right) Distribution of log2

1026    normalized expression of gene of interest for MPAL disease subpopulations (red) and

1027    closest normal cells (grey); black line represents the mean and asterisk denote

1028    significance (LFC > 0.5 and FDR < 0.01). (Bottom) Peak-to-gene links for gene of interest.

1029    **a**, Multi-omic differential track for *IL1RAP*.

1030    **b**, Multi-omic differential track for *CD96*.

1031    **c**, Multi-omic differential track for *FLT3*.

1032    **d**, Multi-omic differential track for *MCL1*.

1033

1034    **Supplementary Figure 16. Analysis workflows for processing of scRNA-seq and**

1035    **scATAC-seq data.**

1036    **a**, scRNA-seq analysis workflow. Briefly cells are aligned using 10x cell ranger, quality

1037    filtered, and clustered using a feature optimization approach (see methods).

1038    **b**, scATAC-seq analysis workflow. Briefly cells are aligned using 10x cell ranger atac,

1039    quality filtered, clustered in large windows genome-wide, peak-calling on clusters, creation

1040    of a counts matrix and clustered using a feature optimization approach (see methods).

1041

1042

1043    **Supplementary Table 1. MPAL Patient Characteristics.**

1044    MPAL patient WHO Diagnosis, Age, Sex, Blast %, White Blood Cell Count, Cytogenetics,

1045    Prior Treatment.

1046

1047    **Supplementary Table 2. Antibodies used in flow cytometry of MPALs.**

1048

1049    **Supplementary Table 3. CITE-Seq Antibody List and Barcodes.**

1050    Antibody information for Hematopoietic and MPAL samples. Barcodes used for

1051    sequencing ADT libraries.

1052

1053    **Supplementary Table 4. Differential analyses for MPAL and AMLs.**

1054    MPAL differential RNA-seq k-means, MPAL differential ATAC-seq k-means, AML

1055    differential RNA-seq k-means and MPAL vs AML comparison.

1056

1057    **Supplementary Table 5. Motif enrichment and linkage to target genes.**

1058    MPAL differential ATAC-seq k-means enrichment for CIS-BP motifs shown in figure 3A,

1059    all motifs, significant peak-to-gene links, and RUNX1 target genes.

1060

1061    **References**

1062    1.  Hoadley, K. A. *et al.* Cell-of-Origin Patterns Dominate the Molecular Classification of

1063        10,000 Tumors from 33 Types of Cancer. *Cell* **173,** 291–304.e6 (2018).

1064    2.  Corces, M. R. *et al.* The chromatin accessibility landscape of primary human

1065        cancers. *Science* **362,** (2018).

1066    3.  Polak, P. *et al.* Cell-of-origin chromatin organization shapes the mutational

1067        landscape of cancer. *Nature* **518,** 360–364 (2015).

1068    4.  Weinberg, O. K. & Arber, D. A. Mixed-phenotype acute leukemia: historical overview

1069        and a new definition. *Leukemia* **24,** 1844–1851 (2010).

1070    5.  Arber, D. A. *et al.* The 2016 revision to the World Health Organization classification

1071        of myeloid neoplasms and acute leukemia. *Blood* **127,** 2391–2405 (2016).

1072    6.  Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single

1073        cells. *Nat. Methods* **14,** 865–868 (2017).

1074    7.  Satpathy, A. T. *et al.* Massively parallel single-cell chromatin landscapes of human

1075        immune cell development and intratumoral T cell exhaustion. *BioRxiv* (2019).

1076        doi:10.1101/610550

1077    8.  Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells.

1078        *Nat. Commun.* **8,** 14049 (2017).

1079    9.  Cusanovich, D. A. *et al.* The cis-regulatory dynamics of embryonic development at

1080        single-cell resolution. *Nature* **555,** 538–542 (2018).

1081    10. Cusanovich, D. A. *et al.* A Single-Cell Atlas of In Vivo Mammalian Chromatin

1082        Accessibility. *Cell* **174,** 1309–1324.e18 (2018).

1083    11. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell

1084        transcriptomic data across different conditions, technologies, and species. *Nat.*

1085        *Biotechnol.* **36,** 411–420 (2018).

1086    12. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and

1087        Projection for Dimension Reduction. *arXiv* (2018).

1088    13. Charles A Janeway, J., Travers, P., Walport, M. & Shlomchik, M. J. Immunobiology

1089        - NCBI Bookshelf. (2001).

1090    14. Pliner, H. A. *et al.* Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell

1091        Chromatin Accessibility Data. *Mol. Cell* **71,** 858–871.e8 (2018).

1092    15. Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. chromVAR: inferring

1093        transcription-factor-associated accessibility from single-cell epigenomic data. *Nat.*

1094        *Methods* **14,** 975–978 (2017).

1095    16. Alexander, T. B. *et al.* The genetic basis and cell of origin of mixed phenotype acute

1096        leukaemia. *Nature* **562,** 373–379 (2018).

1097    17. Takahashi, K. *et al.* Integrative genomic analysis of adult mixed phenotype acute

1098        leukemia delineates lineage associated molecular subtypes. *Nat. Commun.* **9,** 2670

1099        (2018).

1100    18. Corces, M. R. *et al.* Lineage-specific and single-cell chromatin accessibility charts

1101        human hematopoiesis and leukemia evolution. *Nat. Genet.* **48,** 1193–1203 (2016).

1102    19. van Galen, P. *et al.* Single-Cell RNA-Seq Reveals AML Hierarchies Relevant to

1103        Disease Progression and Immunity. *Cell* **176,** 1265–1281.e24 (2019).

1104    20. Satpathy, A. T. *et al.* Transcript-indexed ATAC-seq for precision immune profiling.

1105        *Nat. Med.* **24,** 580–590 (2018).

1106    21. Mezger, A. *et al.* High-throughput chromatin accessibility profiling at single-cell

1107        resolution. *Nat. Commun.* **9,** 3647 (2018).

1108    22. Buenrostro, J. D. *et al.* Integrated Single-Cell Analysis Maps the Continuous

1109        Regulatory Landscape of Human Hematopoietic Differentiation. *Cell* **173,** 1535–

1110        1548.e16 (2018).

1111    23. Mitchell, K. *et al.* IL1RAP potentiates multiple oncogenic signaling pathways in AML.

1112        *J. Exp. Med.* **215,** 1709–1727 (2018).

1113    24. Lim, S. & Kaldis, P. Cdks, cyclins and CKIs: roles beyond cell cycle regulation.

1114    *Development* **140,** 3079–3093 (2013).

1115  25.  Wolach, O. & Stone, R. M. How I treat mixed-phenotype acute leukemia. *Blood* **125,**

1116    2477–2485 (2015).

1117  26.  Zheng, C. *et al.* What is the optimal treatment for biphenotypic acute leukemia?

1118    *Haematologica* **94,** 1778–80; author reply 1780 (2009).

1119  27.  Osato, M. *et al.* Biallelic and heterozygous point mutations in the runt domain of the

1120    AML1/PEBP2alphaB gene associated with myeloblastic leukemias. *Blood* **93,** 1817–

1121    1824 (1999).

1122  28.  Haferlach, T. *et al.* Landscape of genetic lesions in 944 patients with

1123    myelodysplastic syndromes. *Leukemia* **28,** 241–247 (2014).

1124  29.  Zhang, J. *et al.* The genetic basis of early T-cell precursor acute lymphoblastic

1125    leukaemia. *Nature* **481,** 157–163 (2012).

1126  30.  Della Gatta, G. *et al.* Reverse engineering of TLX oncogenic transcriptional

1127    networks identifies RUNX1 as tumor suppressor in T-ALL. *Nat. Med.* **18,** 436–440

1128    (2012).

1129  31.  Wang, X. *et al.* Breast tumors educate the proteome of stromal tissue in an

1130    individualized but coordinated manner. *Sci. Signal.* **10,** (2017).

1131  32.  Sanda, T. *et al.* Core transcriptional regulatory circuit controlled by the TAL1

1132    complex in human T cell acute lymphoblastic leukemia. *Cancer Cell* **22,** 209–221

1133    (2012).

1134  33.  Ben-Ami, O. *et al.* Addiction of t(8;21) and inv(16) acute myeloid leukemia to native

1135    RUNX1. *Cell Rep.* **4,** 1131–1143 (2013).

1136  34.  Wilkinson, A. C. *et al.* RUNX1 is a key target in t(4;11) leukemias that contributes to

1137    gene activation through an AF4-MLL complex interaction. *Cell Rep.* **3,** 116–127

1138    (2013).

1139  35.  Stuart, T. *et al.* Comprehensive integration of single cell data. *BioRxiv* (2018).

1140    doi:10.1101/460147

1141    36. Mumbach, M. R. *et al.* Enhancer connectome in primary human cells identifies

1142        target genes of disease-associated DNA elements. *Nat. Genet.* **49,** 1602–1612

1143        (2017).

1144    37. Martín, P. *et al.* CD69 association with Jak3/Stat5 proteins regulates Th17 cell

1145        differentiation. *Mol. Cell. Biol.* **30,** 4877–4889 (2010).

1146    38. Shiow, L. R. *et al.* CD69 acts downstream of interferon-alpha/beta to inhibit S1P1

1147        and lymphocyte egress from lymphoid organs. *Nature* **440,** 540–544 (2006).

1148    39. Egawa, T., Tillman, R. E., Naoe, Y., Taniuchi, I. & Littman, D. R. The role of the

1149        Runx transcription factors in thymocyte differentiation and in homeostasis of naive T

1150        cells. *J. Exp. Med.* **204,** 1945–1957 (2007).

1151    40. Simeonov, D. R. *et al.* Discovery of stimulation-responsive immune enhancers with

1152        CRISPR activation. *Nature* **549,** 111–115 (2017).

1153    41. Feld, C. *et al.* Combined cistrome and transcriptome analysis of SKI in AML cells

1154        identifies SKI as a co-repressor for RUNX1. *Nucleic Acids Res.* **46,** 3412–3428

1155        (2018).

1156    42. Cancer Genome Atlas Research Network *et al.* Genomic and epigenomic

1157        landscapes of adult de novo acute myeloid leukemia. *N. Engl. J. Med.* **368,** 2059–

1158        2074 (2013).

1159    43. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and

1160        Projection for Dimension Reduction. *arXiv* (2018).

1161    44. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package

1162        for differential expression analysis of digital gene expression data. *Bioinformatics*

1163        **26,** 139–140 (2010).

1164    45. Soneson, C. & Robinson, M. D. Bias, robustness and scalability in single-cell

1165        differential expression analysis. *Nat. Methods* **15,** 255–261 (2018).

**ACKNOWLEDGEMENTS**

1177

**CODE AVAILABILITY**

1179    Code used in this study will be posted on GitHub for main analyses.

1180

**DATA AVAILABILITY**

1182    Sequencing data will be deposited in the Gene Expression Omnibus (GEO). There are no

1183    restrictions on data availability or use.

1184

**Author Contributions**

1186    L.M.M. and S.K. conceived the project and designed the experiments. L.M.M., M.L., E.G.,

1187    and R.M. curated patient samples. S.K. led data production and performed the

1188    experiments together with A.K., A.M., and L.M.M.. G.X.Y.Z. provided healthy bone marrow

1189    and peripheral blood CITE-seq data. S.K. analyzed the scADT-seq data with contribution

1190    from B.P.. J.M.G conceived the analytical workflows and performed the data analysis for

1191     scATAC-seq and scRNA-seq supervised by H.Y.C. and W.J.G.. J.M.G., S.K., L.M.M., and

1192     W.J.G wrote the manuscript with input from all authors.

1193

1194     **COMPETING FINANCIAL INTERESTS**

1195     R.M. is a founder, equity holder, and serves on the Board of Directors of Forty Seven Inc.

1196     H.Y.C. has affiliation with Accent Therapeutics (Founder, SAB), 10x Genomics (SAB), and

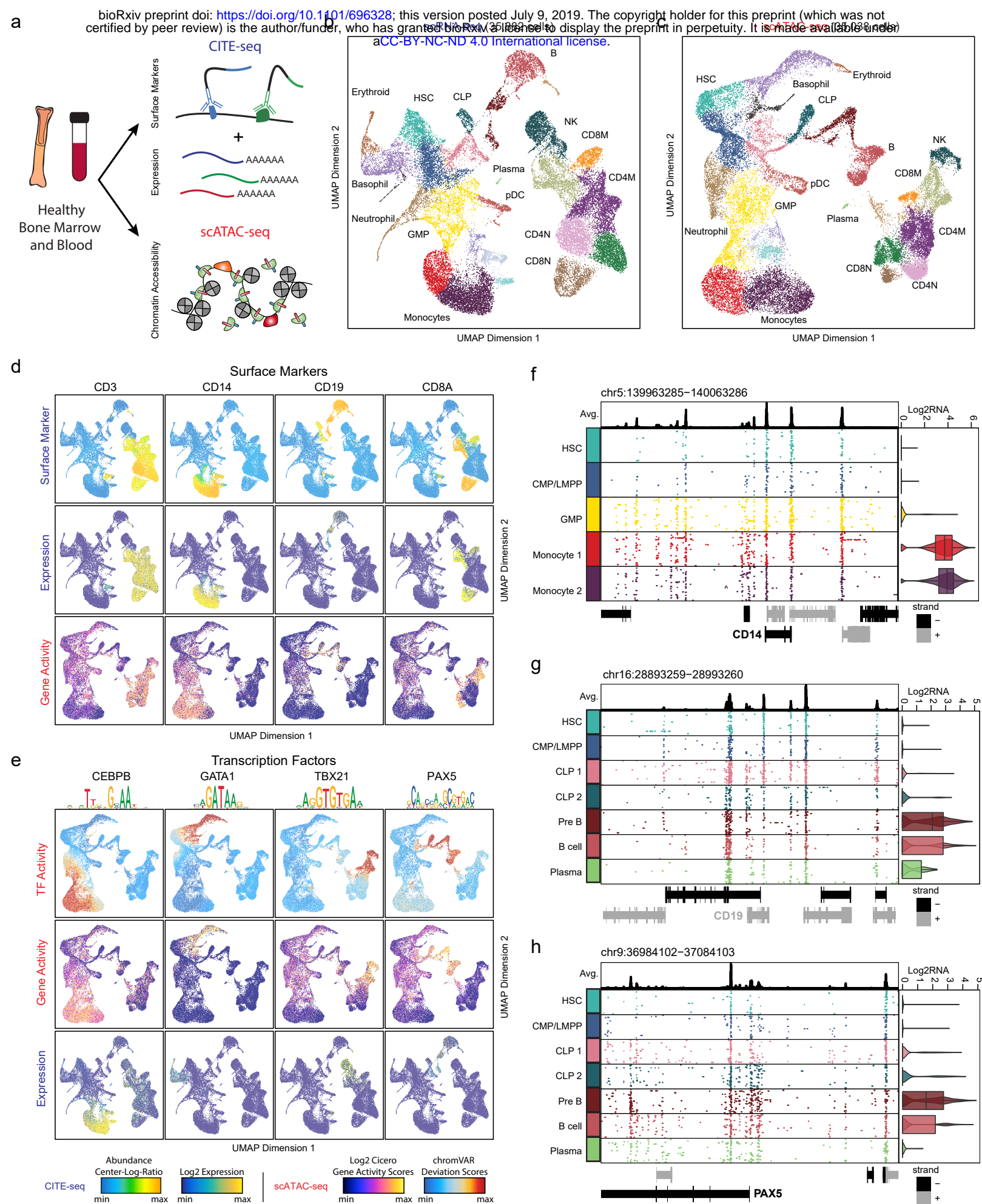1197     Spring Discovery (SAB). W.J.G. has affiliation with 10x Genomics (Consultant) and

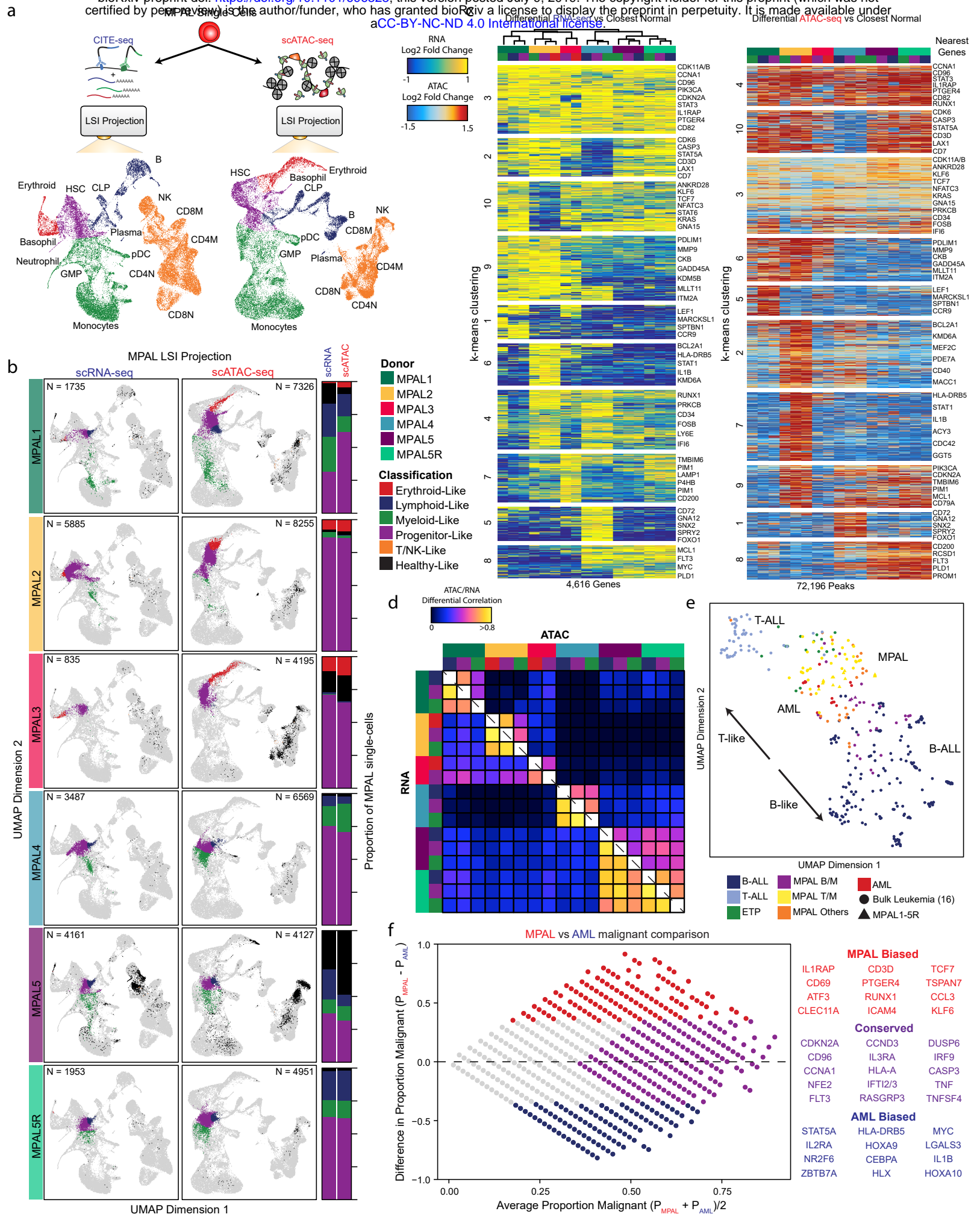1198     Guardant Health (Consultant).
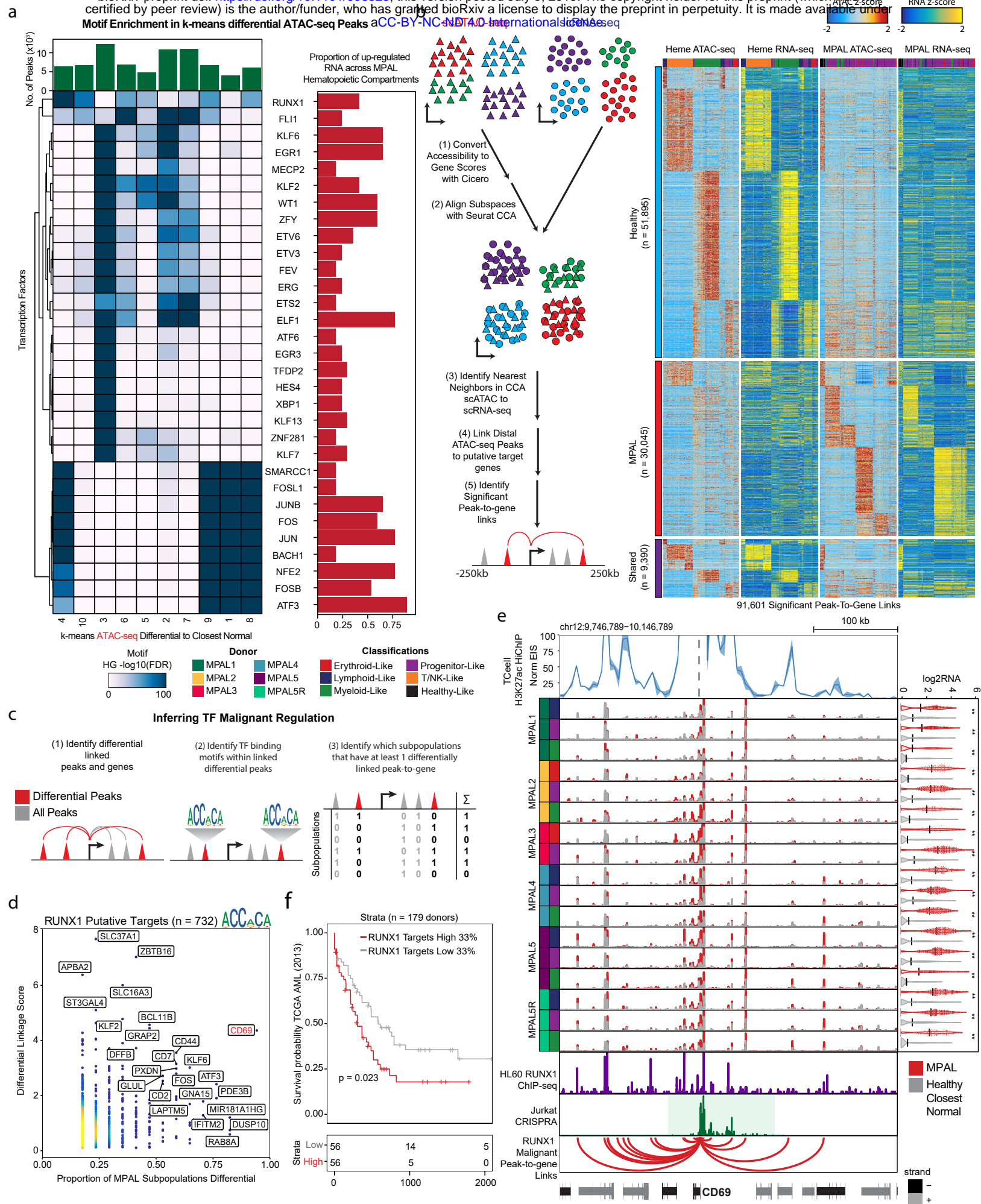
1199

Figure 1

Figure 2

Figure 3