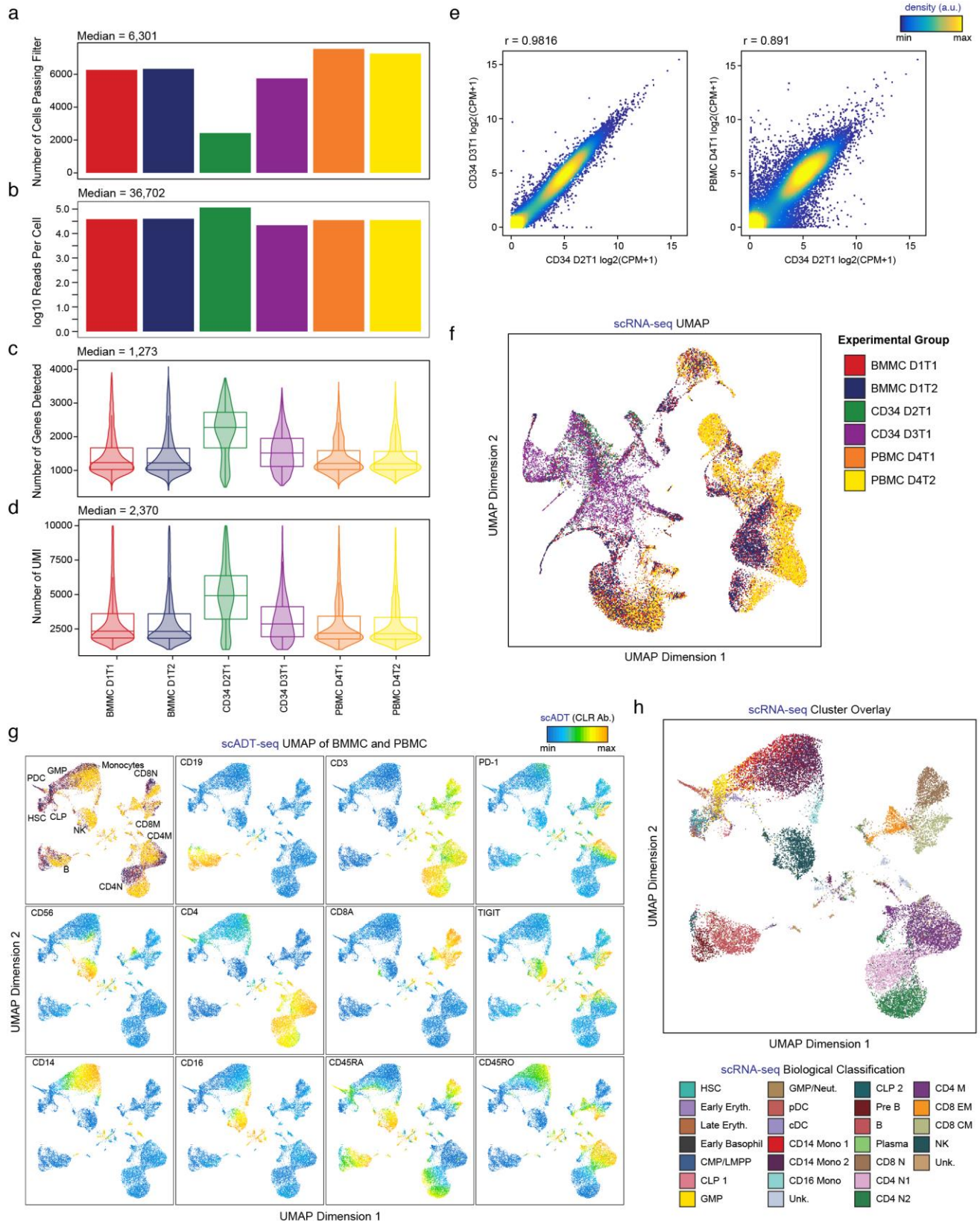


In the format provided by the authors and unedited.

Single-cell multiomic analysis identifies regulatory programs in mixed-phenotype acute leukemia

Jeffrey M. Granja^{1,2,3,13}, Sandy Klemm^{3,13*}, Lisa M. McGinnis^{3,4,13*}, Arwa S. Kathiria³, Anja Mezger^{3,5}, M. Ryan Corces^{1,4}, Benjamin Parks^{3,6}, Eric Gars⁴, Michaela Liedtke⁷, Grace X. Y. Zheng⁸, Howard Y. Chang^{1,3,9,10}, Ravindra Majeti⁷ and William J. Greenleaf^{1,3,11,12*}

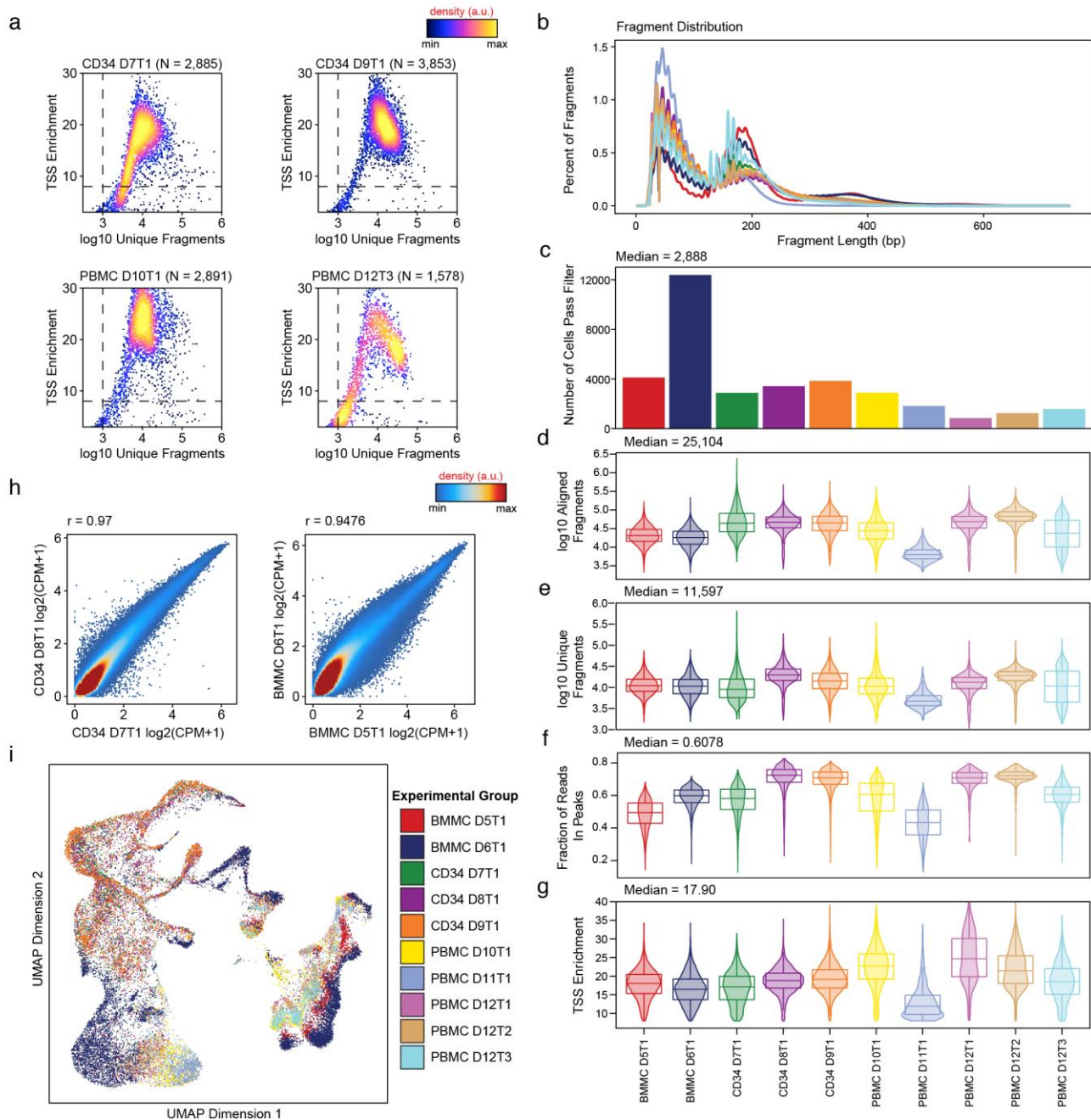
¹Center for Personal Dynamic Regulomes, Stanford University School of Medicine, Stanford, CA, USA. ²Biophysics Program, Stanford University School of Medicine, Stanford, CA, USA. ³Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA. ⁴Department of Pathology, Stanford University School of Medicine, Stanford, CA, USA. ⁵Department of Medical Biochemistry and Biophysics, Karolinska Institute, Stockholm, Sweden. ⁶Department of Computer Science, Stanford University School of Engineering, Stanford, CA, USA. ⁷Department of Medicine, Division of Hematology, Stanford Cancer Institute, Stanford University School of Medicine, Stanford, CA, USA. ⁸10x Genomics, Pleasanton, CA, USA. ⁹Department of Dermatology, Stanford University School of Medicine, Redwood City, CA, USA. ¹⁰Howard Hughes Medical Institute, Stanford University School of Medicine, Stanford, CA, USA. ¹¹Department of Applied Physics, Stanford University, Stanford, CA, USA. ¹²Chan-Zuckerberg Biohub, San Francisco, CA, USA. ¹³These authors contributed equally: Jeffrey M. Granja, Sandy Klemm, Lisa M. McGinnis. *e-mail: klemm@stanford.edu; lisa.mcginnis@stanford.edu; wjg@stanford.edu



Supplementary Figure 1

Quality control of CITE-seq data for hematopoiesis samples.

(a) Number of cells passing filter for each experimental replicate (number of informative genes > 400 and number of unique molecular identifiers (UMI) > 1000). **(b)** Number of aligned reads on average per cell passing filter for each experimental sample. **(c)** Violin and box-whisker plot of the number of informative genes detected per single cell passing filter per experimental sample (n = 2,424 – 7,544). **(d)** Violin and box-whisker plot of the number of unique molecular identified (UMI) transcripts per cell passing filter per experimental sample (n = 2,424 – 7,544). **(e)** Aggregated scRNA-seq (n = 20,287) one to one reproducibility plots for biological replicates (left) and across sample types (right) colored by the density. The correlation (r) represents the Pearson correlation across all genes. **(f)** scRNA-seq experimental sample labels overlay on UMAP of hematopoiesis (n = 35,582). **(g)** scADT-seq UMAP of BMMC and PBMC samples (n = 4) across 14 antibodies. scADT overlay of experimental sample labels, CD19, CD3, CD56, CD4, CD8A, CD14, CD16, CD45RA, CD45RO, TIGIT and PD-1 respectively. Color represents experimental labels or scADT-seq values after CLR transformation. **(h)** scADT-seq UMAP of BMMC and PBMC samples (n = 4) across 14 antibodies colored by scRNA-seq clusters with biological classification. Box-whisker plot; lower whisker is the lowest value greater than the 25% quantile minus 1.5 times the interquartile range (IQR), the lower hinge is the 25% quantile, the middle is the median, the upper hinge is the 75% quantile and the upper whisker is the largest value less than the 75% quantile plus 1.5 times the IQR.

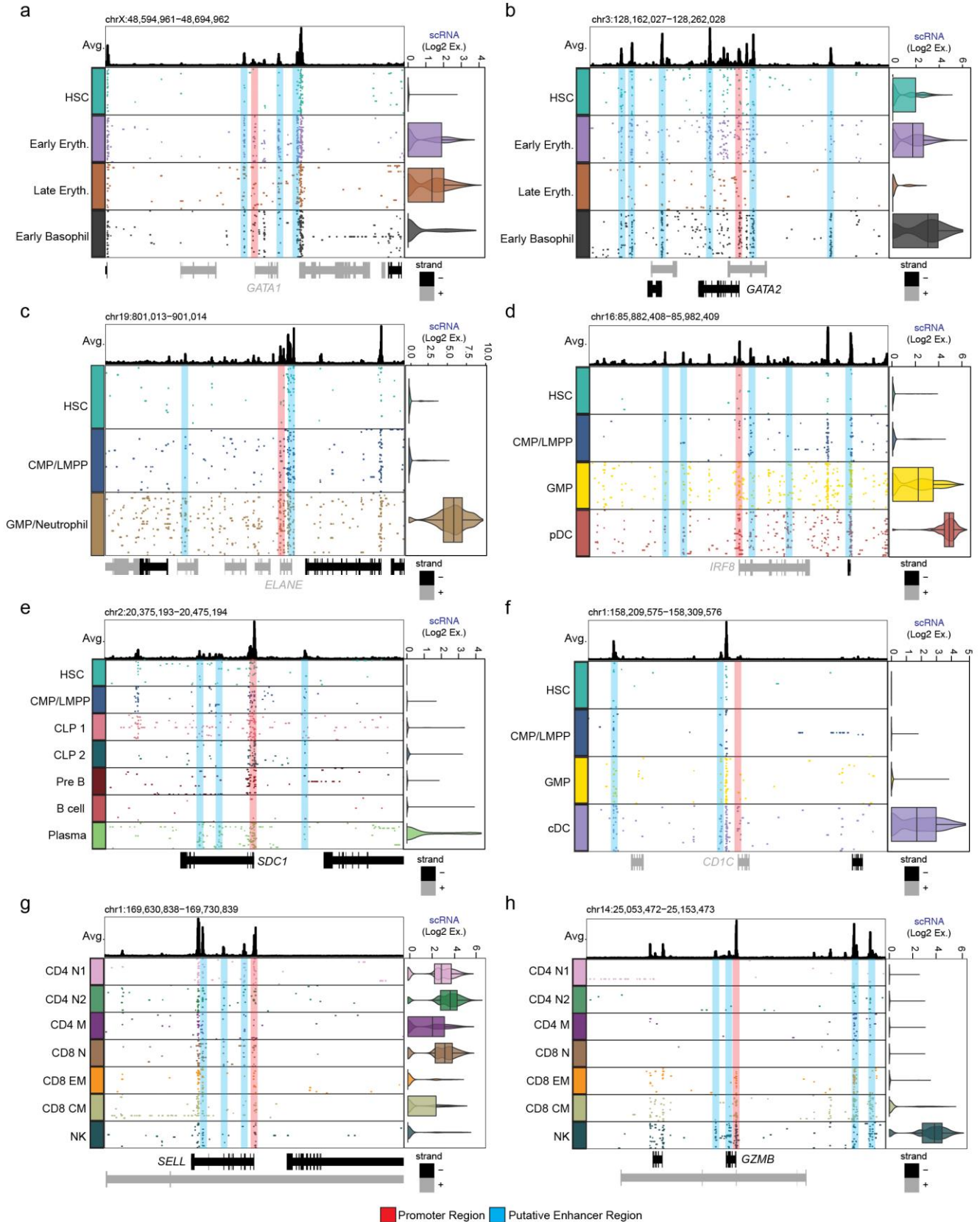


Supplementary Figure 2

Quality control of scATAC-seq data for hematopoiesis samples.

(a) scATAC-seq cell filtering plot of 4 representative scATAC-seq hematopoietic samples. The x-axis is the number of unique accessible fragments and the y-axis is the enrichment of Tn5 insertions at transcription start sites, representing the robust signal to background for each single cell. **(b)** Aggregated scATAC-seq fragment size distributions across individual experiments demonstrating sub-, mono- and multi nucleosome spanning ATAC-seq fragments. **(c)** Number of cells passing filter for each experimental replicate (Unique nuclear fragments > 1000 and TSS enrichment > 8). **(d)** Violin and box-whisker plot of the number of total aligned fragments for each single cell passing filter per experimental sample ($n = 836 - 12,394$). **(e)** Violin and box-whisker plot of the number of unique aligned nuclear fragments for each single cell passing filter per experimental sample ($n = 836 - 12,394$). **(f)** Violin and box-whisker plot of the fraction of the total number Tn5 insertions (Reads) that are within the healthy hematopoietic union peak set ($n = 452,004$) for each single cell passing filter. **(g)** Violin and box-whisker plot of the normalized transcription start site (TSS) enrichment for each single cell passing filter per experimental sample. **(h)** Aggregated scATAC-seq ($n = 452,004$) one to one reproducibility plots for biological

replicates colored by the density. The correlation (r) represents the Pearson correlation across all genes. **(i)** scATAC-seq experimental sample labels overlay on UMAP of hematopoiesis ($n = 35,038$). Box-whisker plot; lower whisker is the lowest value greater than the 25% quantile minus 1.5 times the interquartile range (IQR), the lower hinge is the 25% quantile, the middle is the median, the upper hinge is the 75% quantile and the upper whisker is the largest value less than the 75% quantile plus 1.5 times the IQR.



Supplementary Figure 3

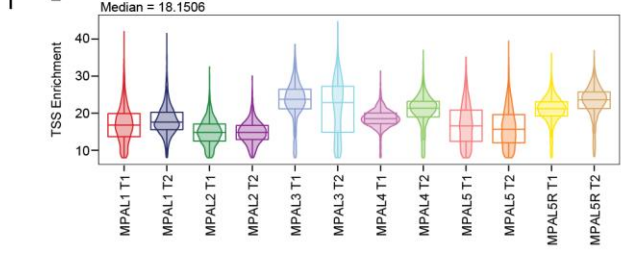
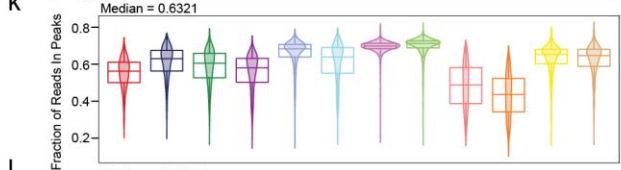
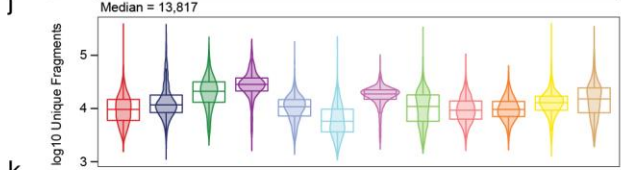
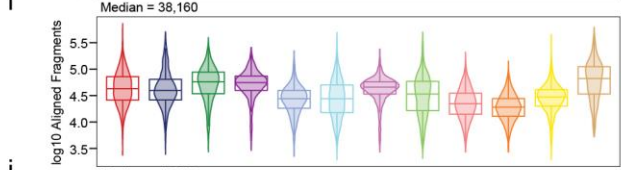
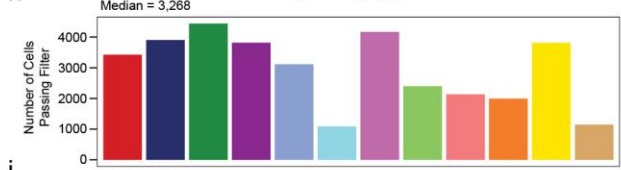
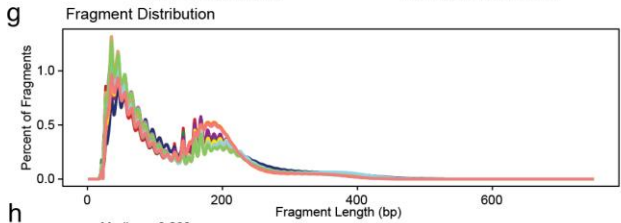
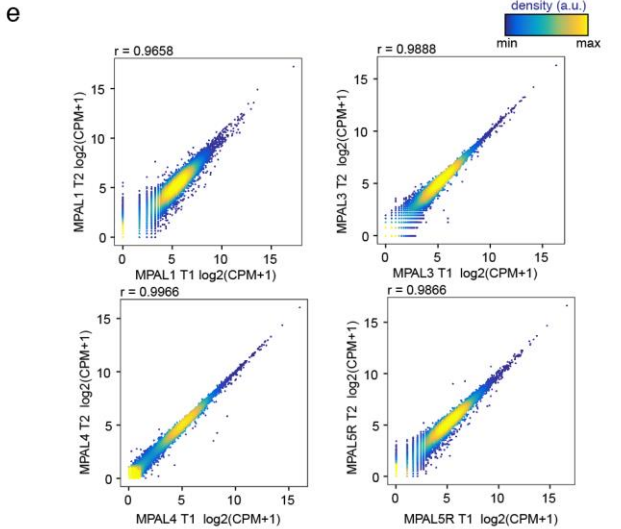
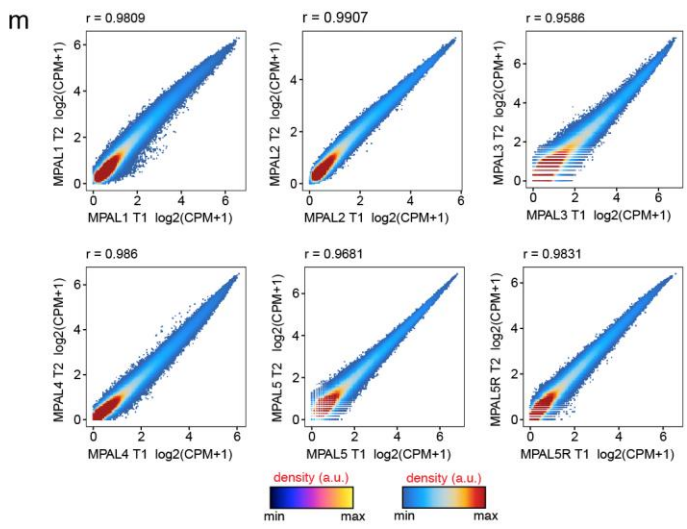
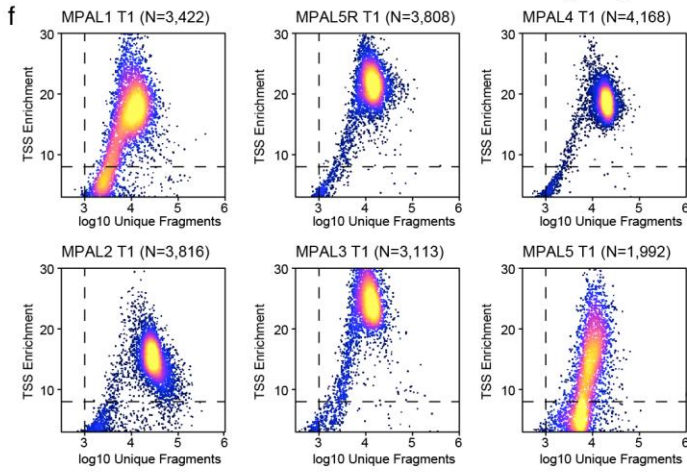
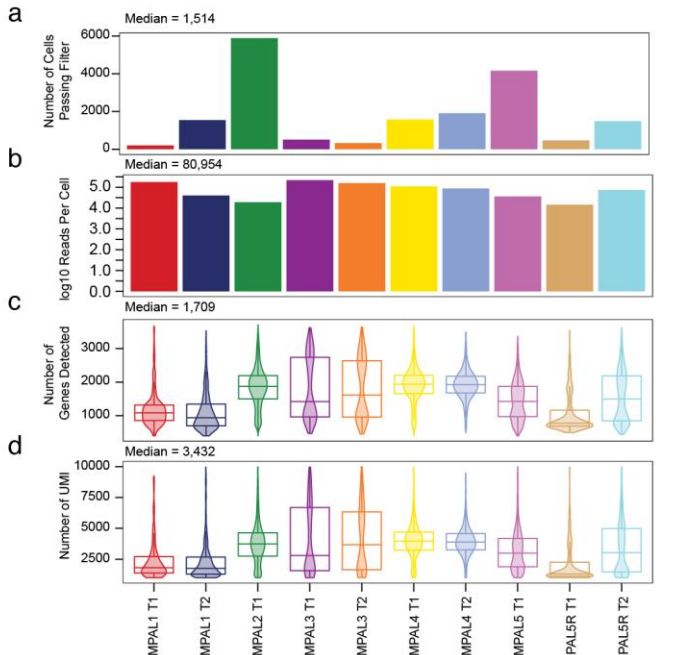
Validation of key marker genes for both scRNA-seq and scATAC-seq for hematopoiesis.

(a-h) Multi-omic tracks; (Top) average track of all clusters displayed, (Middle) binarized 100 random scATAC-seq tracks for each locus at 100bp resolution and (right) violin and box-whisker plot of the scRNA-seq log₂ normalized expression for each cluster. Box-whisker plot; lower whisker is the lowest value greater than the 25% quantile minus 1.5 times the interquartile range (IQR), the lower hinge is the 25% quantile, the middle is the median, the upper hinge is the 75% quantile and the upper whisker is the largest value less than the 75% quantile plus 1.5 times the IQR. **(a)** Multi-omic track of *GATA1* (specific in these clusters for Erythroid) for erythroid development from HSC progenitor cells (n = 111 - 1,653). **(b)** Multi-omic track of *GATA2* (specific in these clusters for Basophil) for erythroid development from HSC progenitor cells (n = 111 - 1,653). **(c)** Multi-omic track of *ELANE* (specific in these clusters for GMP/Neutrophil) for neutrophil development from HSC progenitor cells (n = 1,050 - 2,260). **(d)** Multi-omic track of *IRF8* (specific in these clusters for pDC) across pDC development from HSC progenitor cells (n = 544 - 2,260). **(e)** Multi-omic track of *SDC1* (specific in these clusters for Plasma cells) across B cell development and plasma cells (n = 62 - 2,260). **(f)** Multi-omic track of *CD1C* (specific in these clusters for cDC) across cDC development from HSC progenitor cells (n = 325 - 2,260). **(g)** Multi-omic track of *SELL* (specific in these clusters for Naive T cells vs memory, and CD8 central memory vs CD8 effector memory) across NK and T cells (n = 796 - 3,539). **(h)** Multi-omic track of *GZMB* (specific in these clusters for NK cells) across NK and T cells (n = 796 - 3,539).

Supplementary Figure 4

Diagnostic flow cytometry plots for MPALs 1-5R.

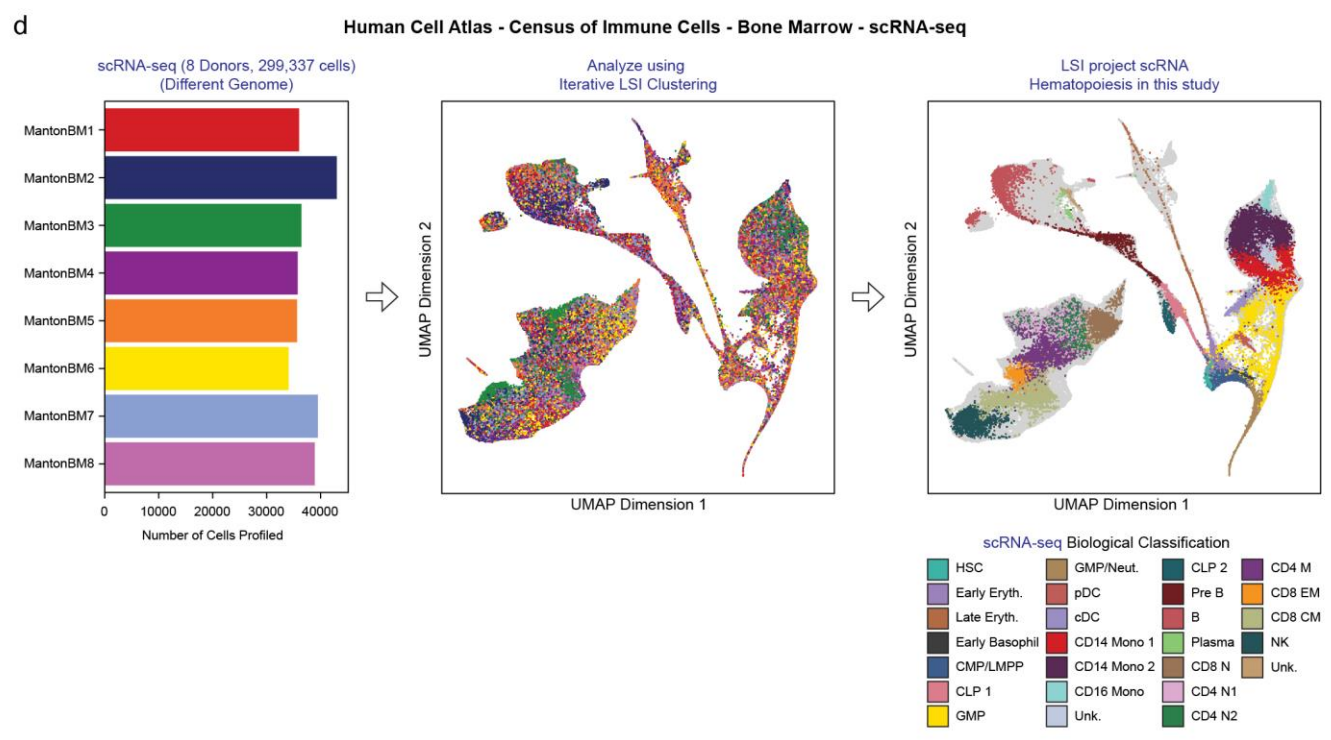
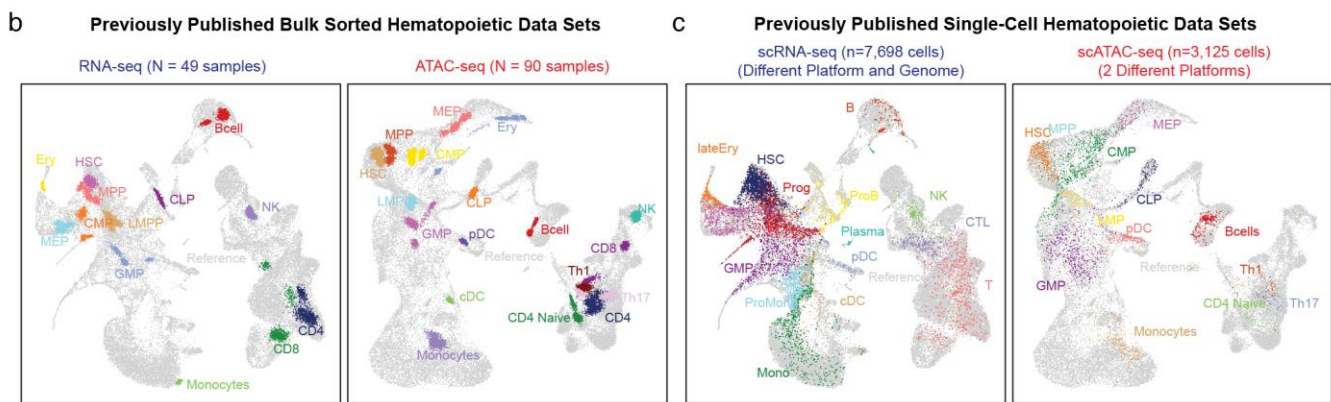
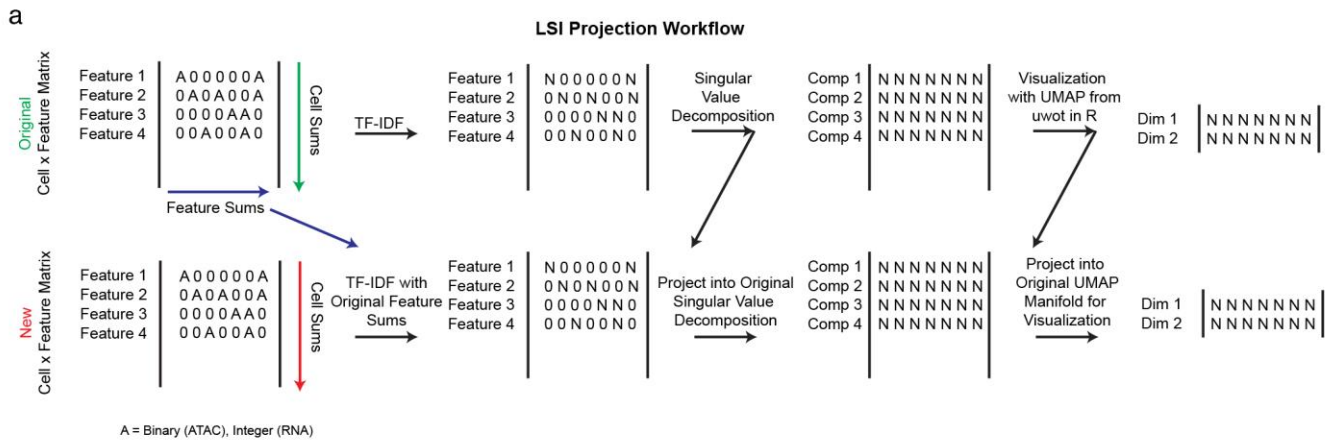
(a-f) Diagnostic flow cytometry plots from the 5 different MPAL cases (MPAL1-5R) gated on blasts area (highlighted in red) and lymphocytes (highlighted in black) from CD45 and side scatter area (SSC-A). **(a)** MPAL 1 shows classic bilineal phenotype with both T-lymphoblasts (cCD3-positive and CD7-positive) and myeloid blasts (MPO-positive and CD33-positive). **(b)** MPAL 2 demonstrates a more complex phenotype with both biphenotypic (single population expressing lymphoid marker CD7 and myeloid marker CD33) and bilineal T-Myeloid patterns (subpopulation expressing monocytic markers CD64, CD33, and CD14). **(c)** MPAL 3 demonstrates a classic biphenotypic case with coexpression of both T-lineage markers (cCD3-positive) and myeloid markers (MPO-positive). **(d)** MPAL4 demonstrates a classic bilineal B/M phenotype expressing B-lineage markers (CD79a and CD19-positive) and myeloid markers (MPO-positive and CD33-positive). **(e)** MPAL5 demonstrates a more complicated phenotype with a subpopulation of blasts expressing T-lineage markers (cCD3-positive and CD7-positive) and a subpopulation expressing myeloid marker MPO. **(f)** MPAL5R post-treatment relapse of MPAL5. Flow cytometry reveals expansion of the T-lymphoblastic subpopulation (cCD3-positive, TdT-positive population) following chemotherapy. **(g)** High-confidence mutations detected in 5 MPAL cases by whole exome sequencing. Missense mutations are shown in blue, frameshift deletions are shown in yellow, stop-gain mutations are shown in purple, frameshift insertions are shown in orange, and nonframeshift deletions are shown in dark gray.



Supplementary Figure 5

Quality control of scRNA-seq and scATAC-seq data for MPAL samples.

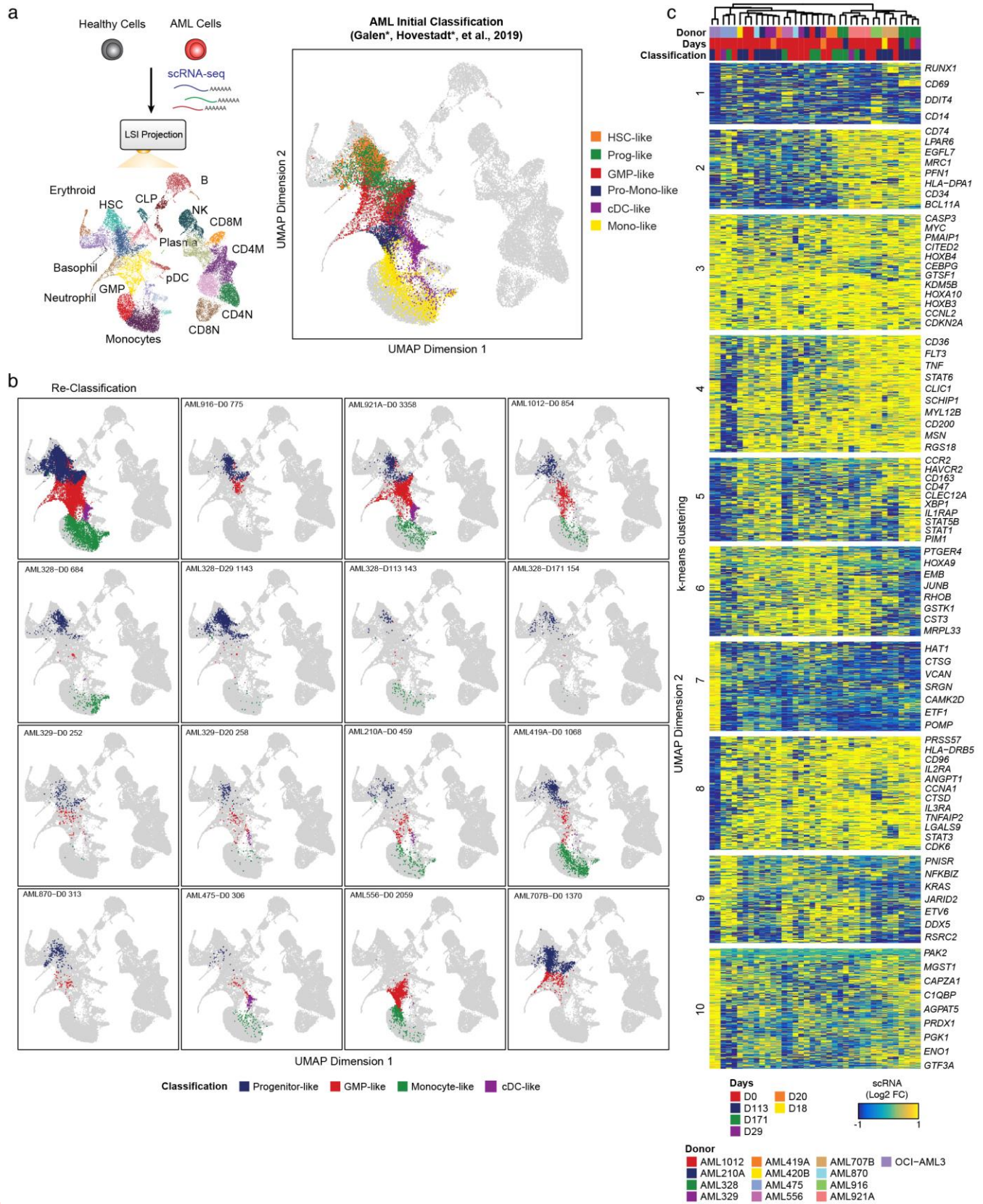
(a) Number of cells passing filter for each experimental replicate (number of informative genes > 400 and number of unique molecular identifiers (UMI) > 1000). **(b)** Number of aligned reads on average per cell passing filter for each experimental sample. **(c)** Violin and box-whisker plot of the number of informative genes detected per single cell passing filter per experimental sample. **(d)** Violin and box-whisker plot of the number of unique molecular identified (UMI) transcripts per cell passing filter per experimental sample. **(e)** Aggregated scRNA-seq ($n = 20,287$) one to one reproducibility plots for technical replicates colored by the density. The correlation (r) represents the Pearson correlation across all genes. **(f)** scATAC-seq cell filtering plot of 6 representative scATAC-seq MPAL samples. The x-axis is the number of unique accessible fragments and the y-axis is the enrichment of Tn5 insertions at transcription start sites, representing the robust signal to background for each single cell. **(g)** Aggregated scATAC-seq fragment size distributions across individual experiments demonstrating sub-, mono- and multi nucleosome spanning ATAC-seq fragments. **(h)** Number of cells passing filter for each experimental replicate (Unique nuclear fragments > 1000 and TSS enrichment > 8). **(i)** Violin and box-whisker plot of the number of total aligned fragments for each single cell passing filter per experimental sample. **(j)** Violin and box-whisker plot of the number of unique aligned nuclear fragments for each single cell passing filter per experimental sample. **(k)** Violin and box-whisker plot of the fraction of the total number Tn5 insertions (Reads) that are within the MPAL union peak set ($n = 346,274$) for each single cell passing filter. **(l)** Violin and box-whisker plot of the normalized transcription start site (TSS) enrichment for each single cell passing filter per experimental sample. **(m)** Aggregated scATAC-seq ($n = 346,274$) one to one reproducibility plots for technical replicates colored by the density. The correlation (r) represents the Pearson correlation across all genes. Box-whisker plot; lower whisker is the lowest value greater than the 25% quantile minus 1.5 times the interquartile range (IQR), the lower hinge is the 25% quantile, the middle is the median, the upper hinge is the 75% quantile and the upper whisker is the largest value less than the 75% quantile plus 1.5 times the IQR.



Supplementary Figure 6

Evaluation of LSI projection workflow for previously published bulk and single-cell hematopoietic data sets across different platforms.

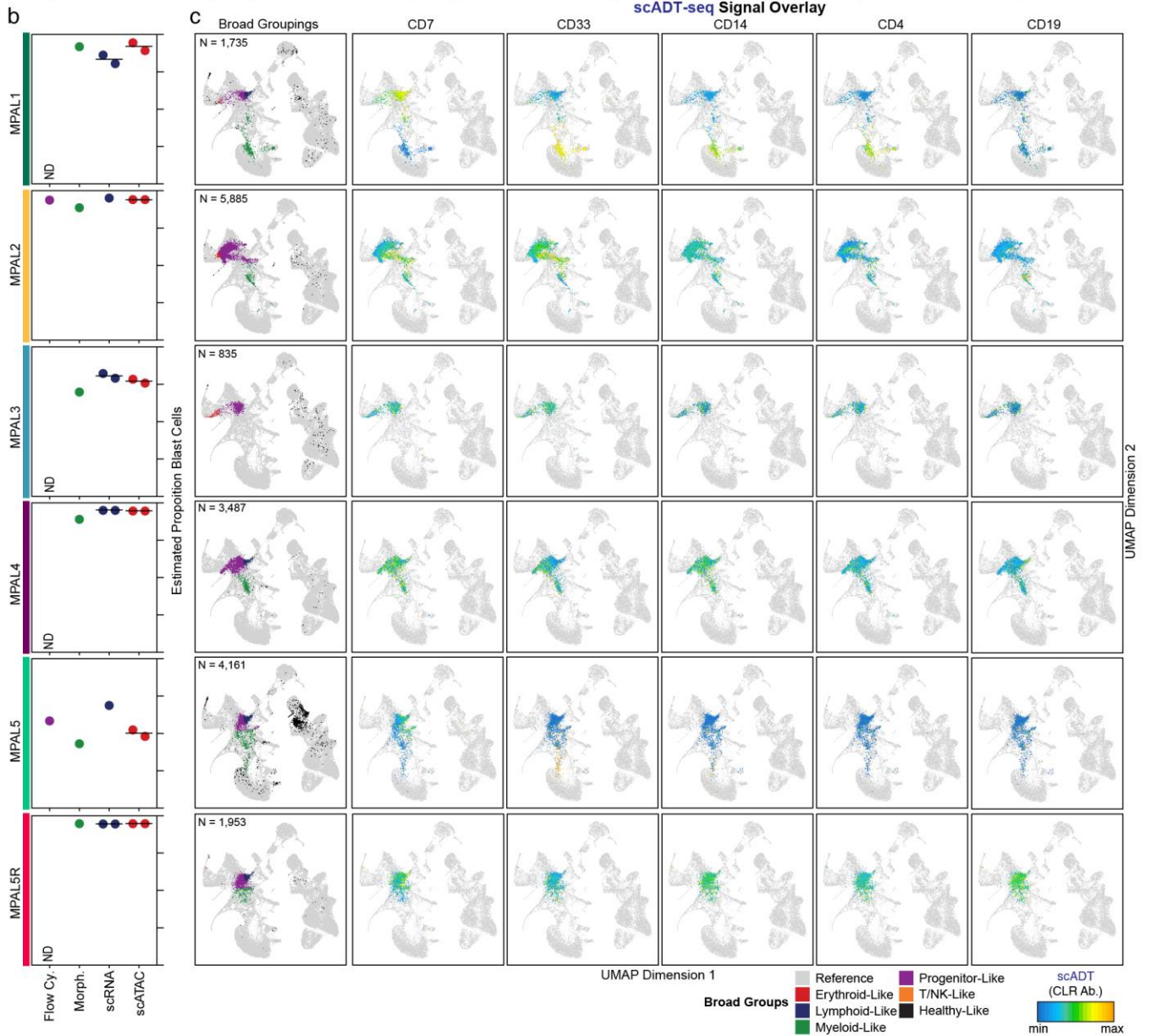
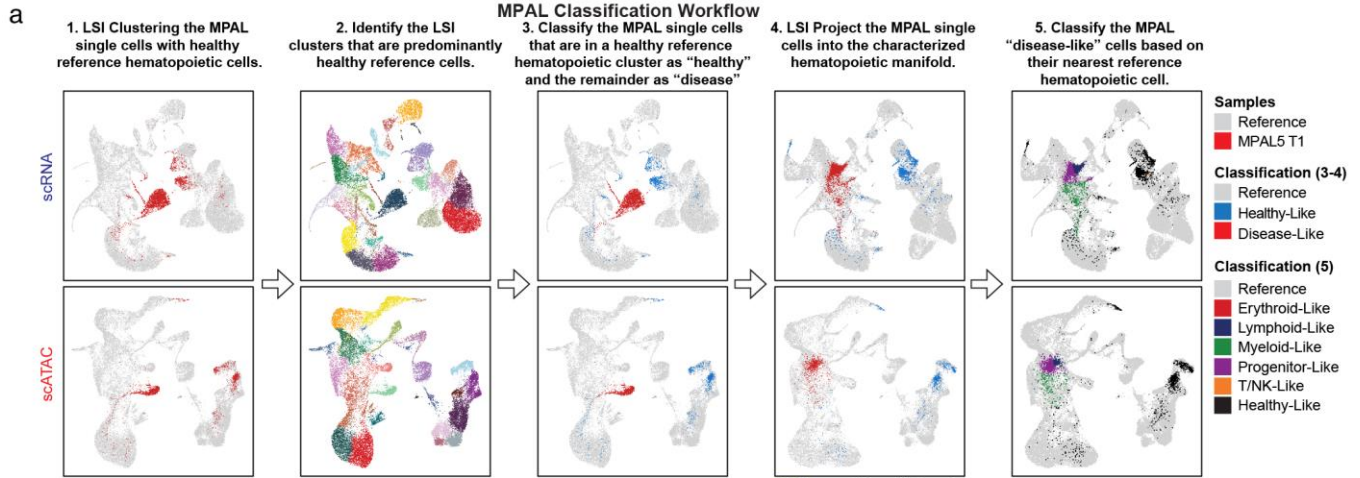
(a) Overview of LSI projection workflow. Briefly, when previously computing the TF-IDF transform for the original hematopoietic manifold, we store the feature sums (document frequency) and use this information to compute the TF-IDF transform for the new matrix. We then use the hematopoietic singular value decomposition loadings on the new TF-IDF matrix to project the new matrix into a common subspace. To visualize the hematopoietic subspace, we constructed a UMAP projection using uwot in R. We then used this learned manifold to project the new matrix subspace into the original UMAP projection. **(b)** LSI projection of downsampled previously published bulk sorted hematopoietic data sets^{18,20}. (Left) RNA-seq downsampled bulk projections for 49 samples (n = 250 downsampled cells). (Right) ATAC-seq downsampled bulk projections for 90 samples (n = 250 downsampled cells). **(c)** LSI projection of downsampled previously published single-cell hematopoietic data sets labeled by previous classifications²⁰⁻²². (Left) scRNA-seq projections of previous study healthy bone marrow cells (different platform and different aligned genome) colored by previous classifications. (Right) scATAC-seq projections for healthy bone marrow and peripheral blood samples (2 different platforms across 3 studies), colored by ground truth isolated populations. **(d)** Projection of hematopoietic scRNA-seq into of Human Cell Atlas (HCA) Census of Immune Cells. (Left) Number of cells per each of 8 bone marrow donors. (Middle) UMAP projection of LSI iterative clustering of HCA bone marrow scRNA-seq. (Right) LSI projection of our scRNA-seq hematopoietic single cells into HCA bone marrow UMAP colored by cluster definitions.



Supplementary Figure 7

LSI projection of previously published healthy and AML scRNA-seq identifies malignant programs across AML subpopulations.

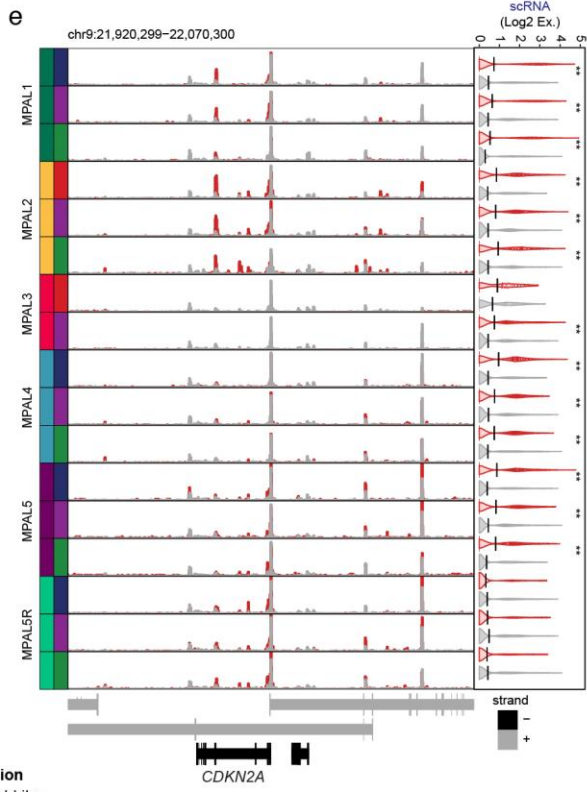
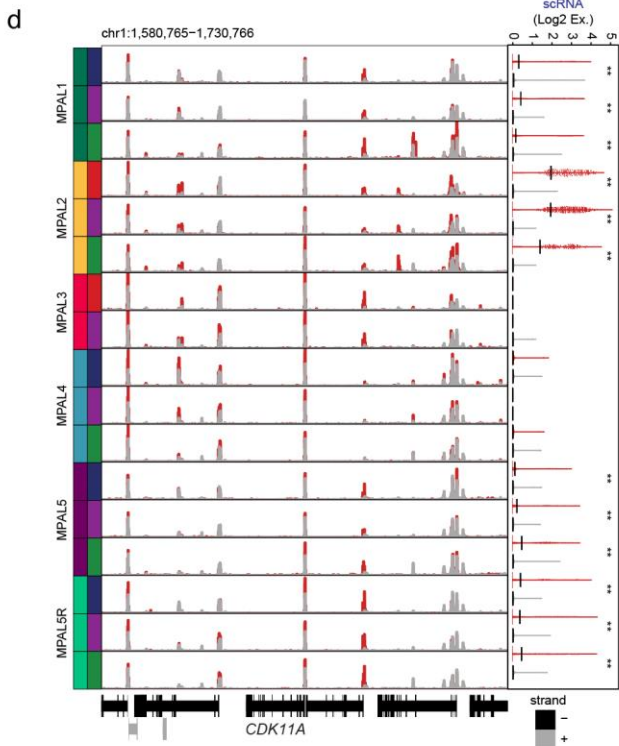
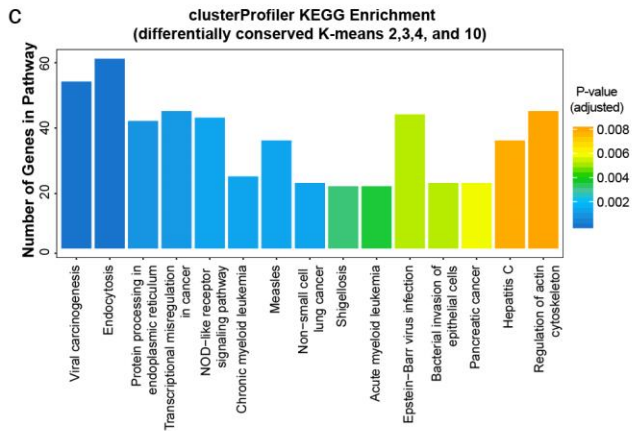
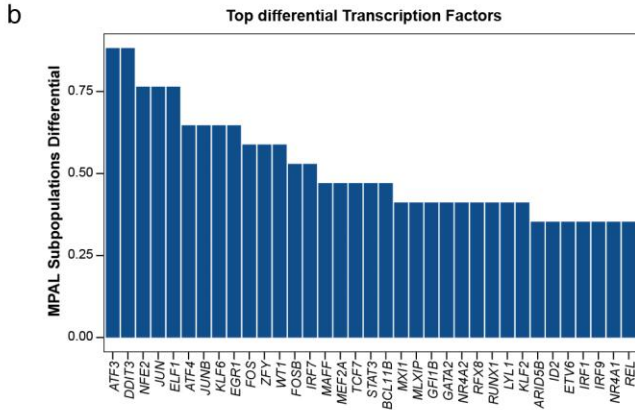
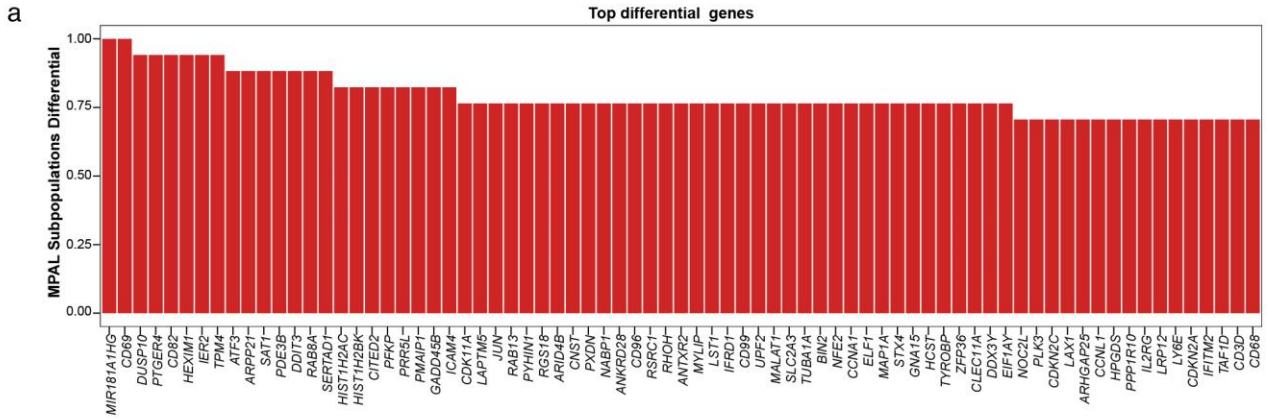
(a) (Left) Schematic of LSI projection. (Right) Initial projection of all AML malignant single-cells colored by previous classifications¹⁹ (n = 31,890). **(b)** Re-classification of scRNA-seq AML single-cells based on closest normal cells in healthy hematopoiesis (See Methods). Broader re-classification increases the number of cells per category for improved power in differential analyses. LSI projection for each individual AML samples onto scRNA-seq healthy hematopoiesis colored by re-classifications (denoted is the sample id and number of cells, n = 143 – 3,358). **(c)** K-means differential scRNA-seq heatmap (k = 10), colored by log₂ fold change, comparing each AML sample subpopulations (classifications) vs their closest normal bone marrow cells from the same study¹⁹.



Supplementary Figure 8

Classification of MPALs by projection onto the hematopoietic hierarchy.

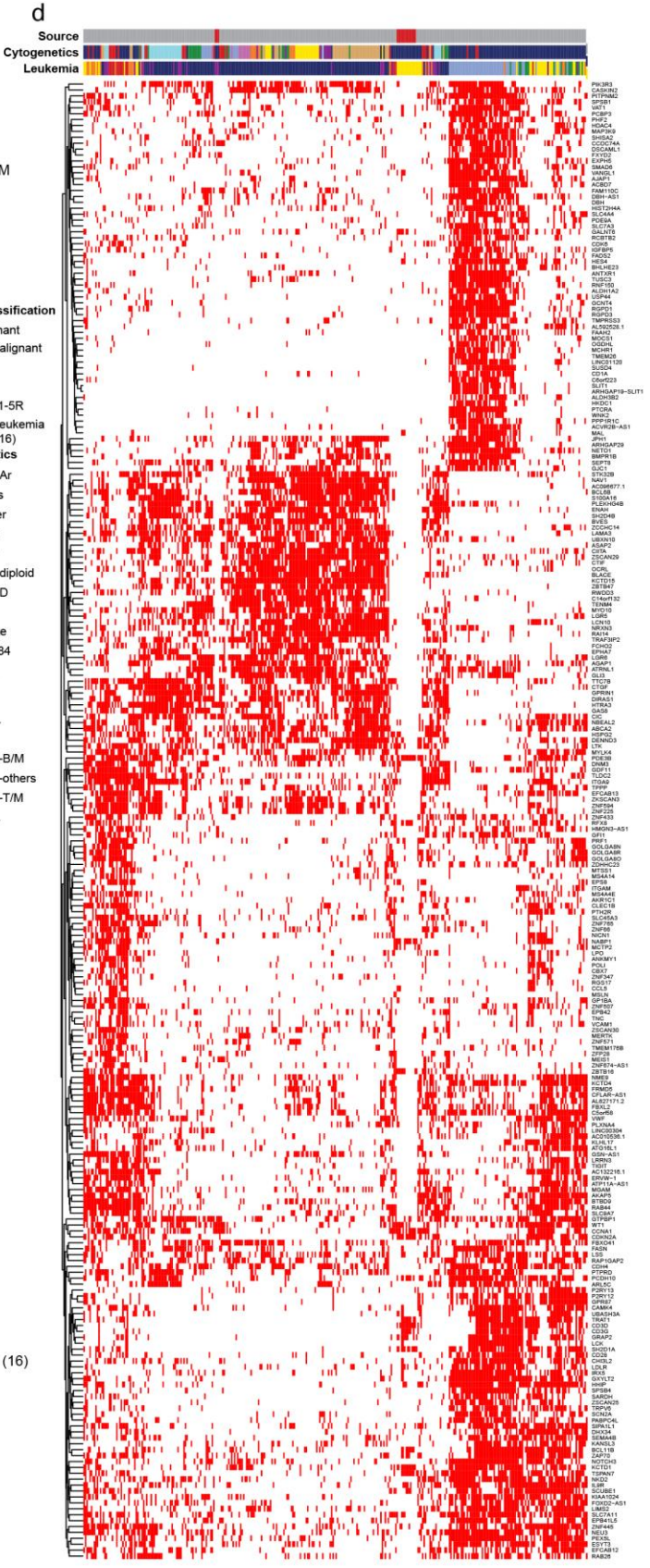
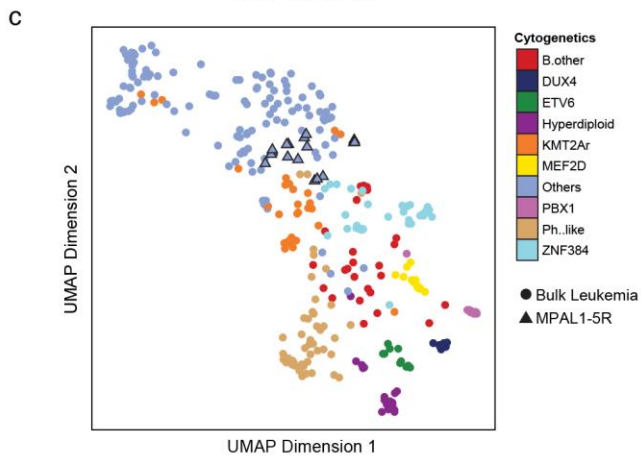
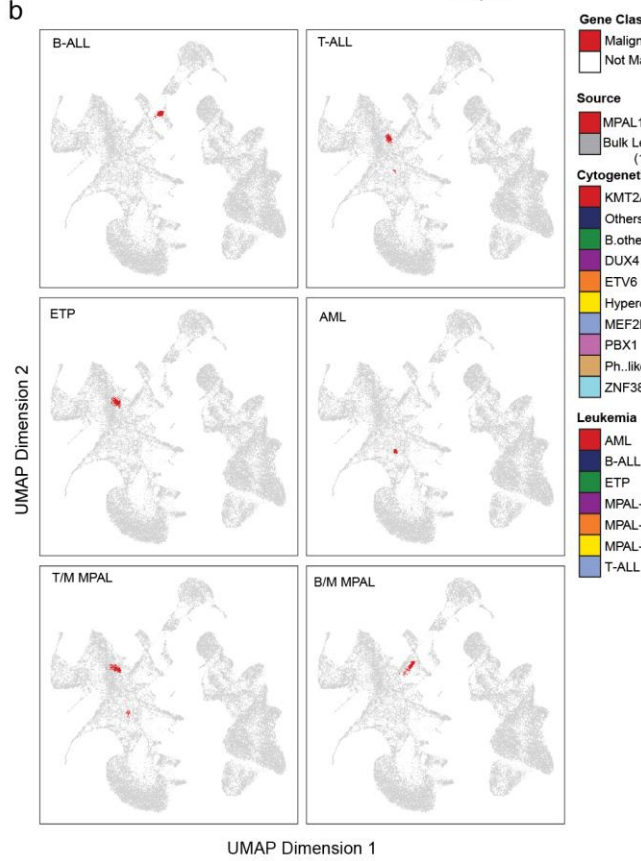
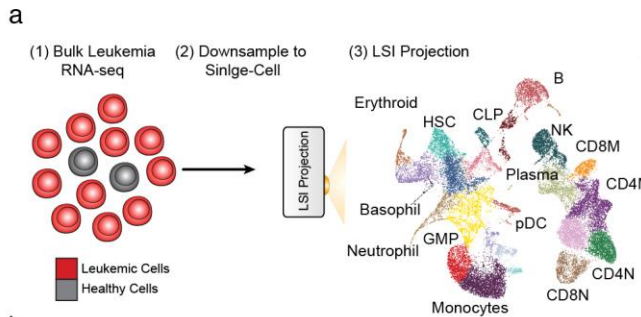
(a) MPAL single cell classification workflow for scRNA and scATAC-seq. First cells were clustered with reference hematopoietic cells (1) and classified based on healthy hematopoietic clusters (2-3). Cells were then LSI projected into the reference hematopoietic manifold (4) and then classified based on the nearest reference cell hematopoietic compartment (5). MPAL 5 replicate 1 is shown as an example for scRNA (Top) and scATAC-seq (Bottom). **(b)** Proportion of estimated blast cells for each MPAL with Flow Cytometry, Morphology, scRNA and scATAC-seq (Range is from 0 to 1). **(c)** (Left) Projected MPALs colored by hematopoietic compartments as described in **a**. (Right) scADT-seq overlay of CD7, CD33, CD14, CD4 and CD19 on MPAL single cells LSI projected onto hematopoiesis.



Supplementary Figure 9

Visualization of top differential genes and accessible peak regions.

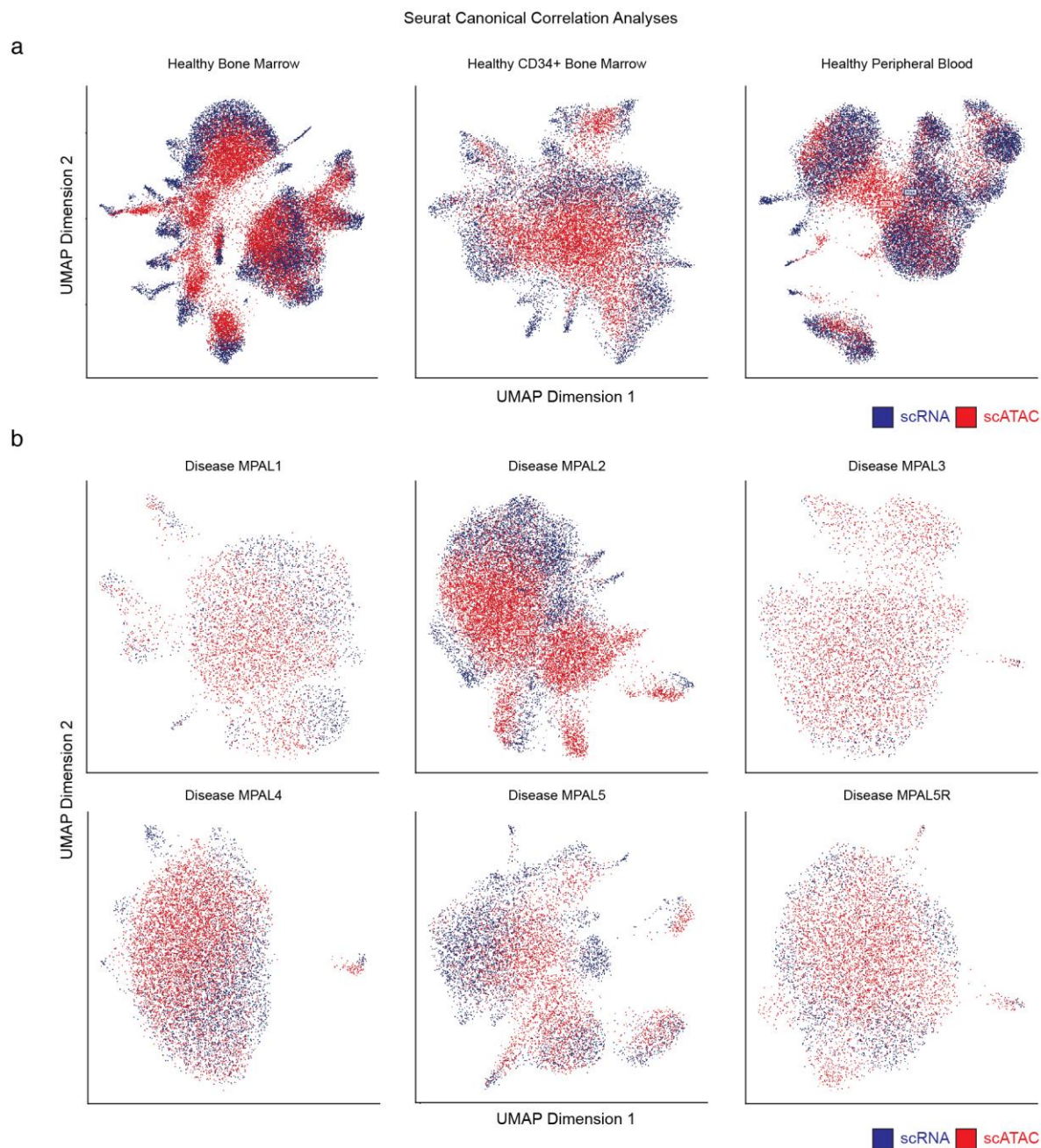
(a) Top conserved differential genes across the MPAL hematopoietic compartments. **(b)** Top conserved differential transcription factors across the MPAL hematopoietic compartments. **(c)** KEGG pathway enrichment of genes (scRNA) differentially conserved k-means 2, 3, 4, and 10 (Figure 2c, $n = 2,117$ genes). Color represents the significance (hypergeometric test adjusted p-value with the Benjamini-Hochberg correction) of the KEGG path way generated using clusterProfiler (Yu, G., Wang, L.-G., Han, Y. & He, Q.-Y. clusterProfiler: an R package for comparing biological themes among gene clusters. *OMICS* **16**, 284–287 (2012)). **(d-e)** Multi-omic differential tracks (Left) scATAC tracks showing MPAL disease subpopulations (red) closest normal cells (grey). (Right) Violin plot of the log₂ normalized expression for MPAL disease subpopulations (red) and closest normal cells (grey); black line represents the mean and asterisk denote significance ($LFC > 0.5$ and $FDR < 0.01$ from Figure 2c). **(d)** Multi-omic differential track of *CDK11A*, up-regulated in MPALs 1, 2, 5 and 5R ($n = 89 - 500$). **(e)** Multi-omic differential track of *CDKN2A*, up-regulated in MPALs 1, 2, 3, 4, and 5 ($n = 89 - 500$).



Supplementary Figure 10

LSI projection of bulk leukemia RNA-seq onto hematopoietic hierarchy.

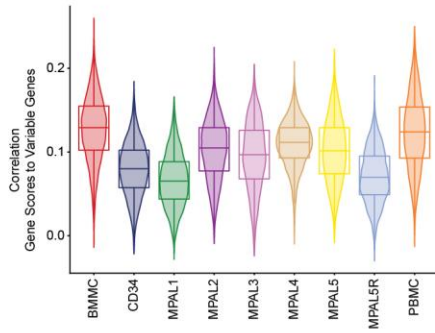
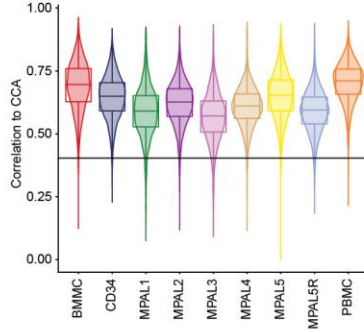
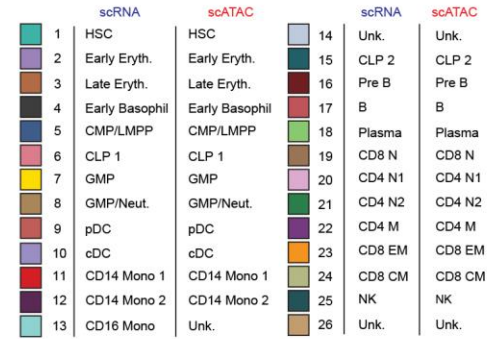
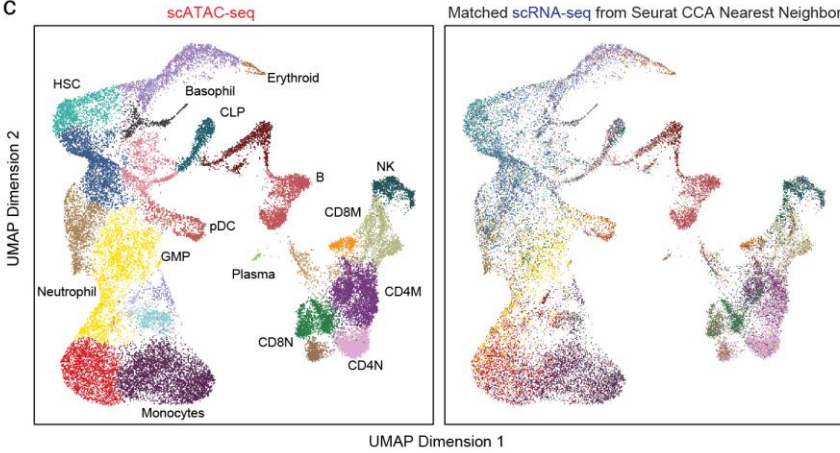
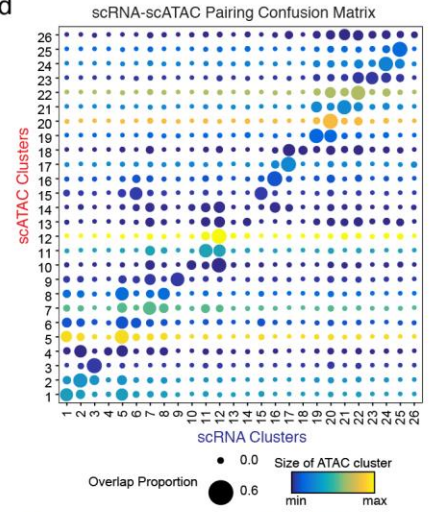
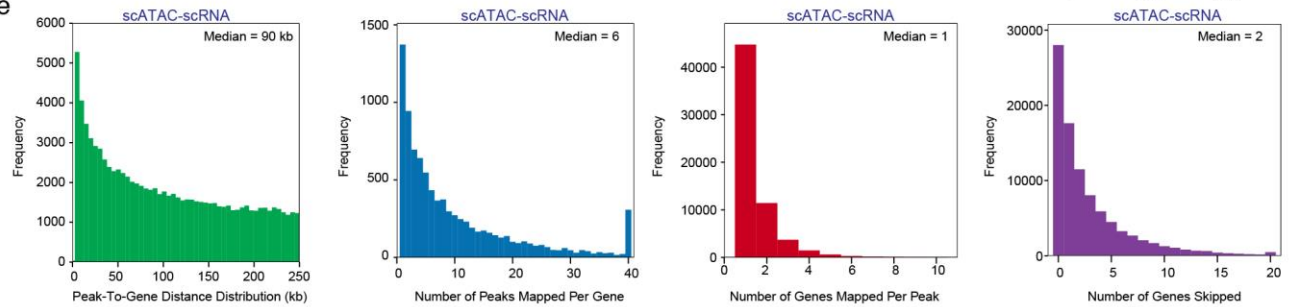
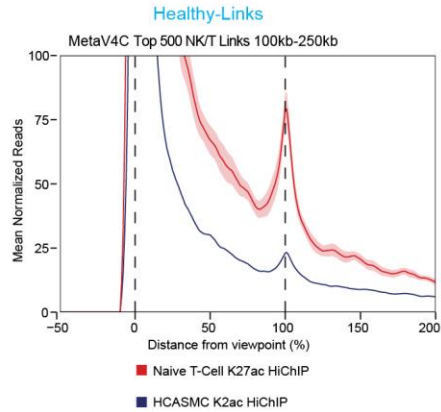
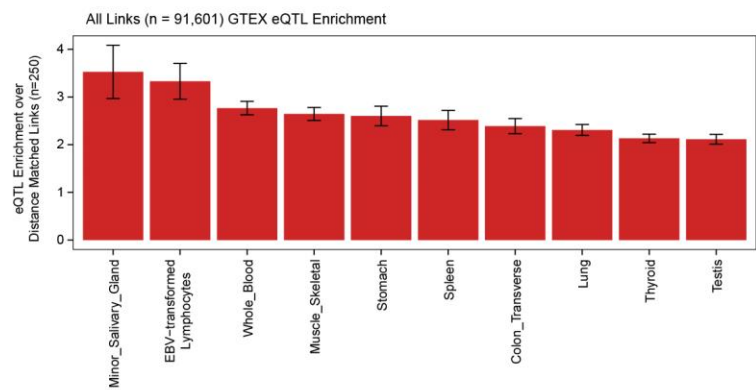
(a) Schematic of LSI projection of downsampled bulk leukemia RNA-seq onto healthy hematopoiesis. **(b)** Representative downsampled LSI projections (n = 250) for B-ALLs, non-ETP T-ALLs, ETP T-ALLs, AMLs, T/M MPALs and B/M MPALs from previous studies¹⁶. **(c)** LSI UMAP of differentially up-regulated gene expression profiles across bulk leukemias¹⁶ (n = 321) and MPAL subpopulations assayed in this study (n = 17), colored by cytogenetics. **(d)** Binary heatmap of variable malignant genes across leukemia classifications. Each cell in the heatmap is colored whether the gene was identified as malignant for the leukemic sample.



Supplementary Figure 11

Seurat canonical correlation analysis alignment of scRNA and scATAC-seq hematopoietic and MPAL samples.

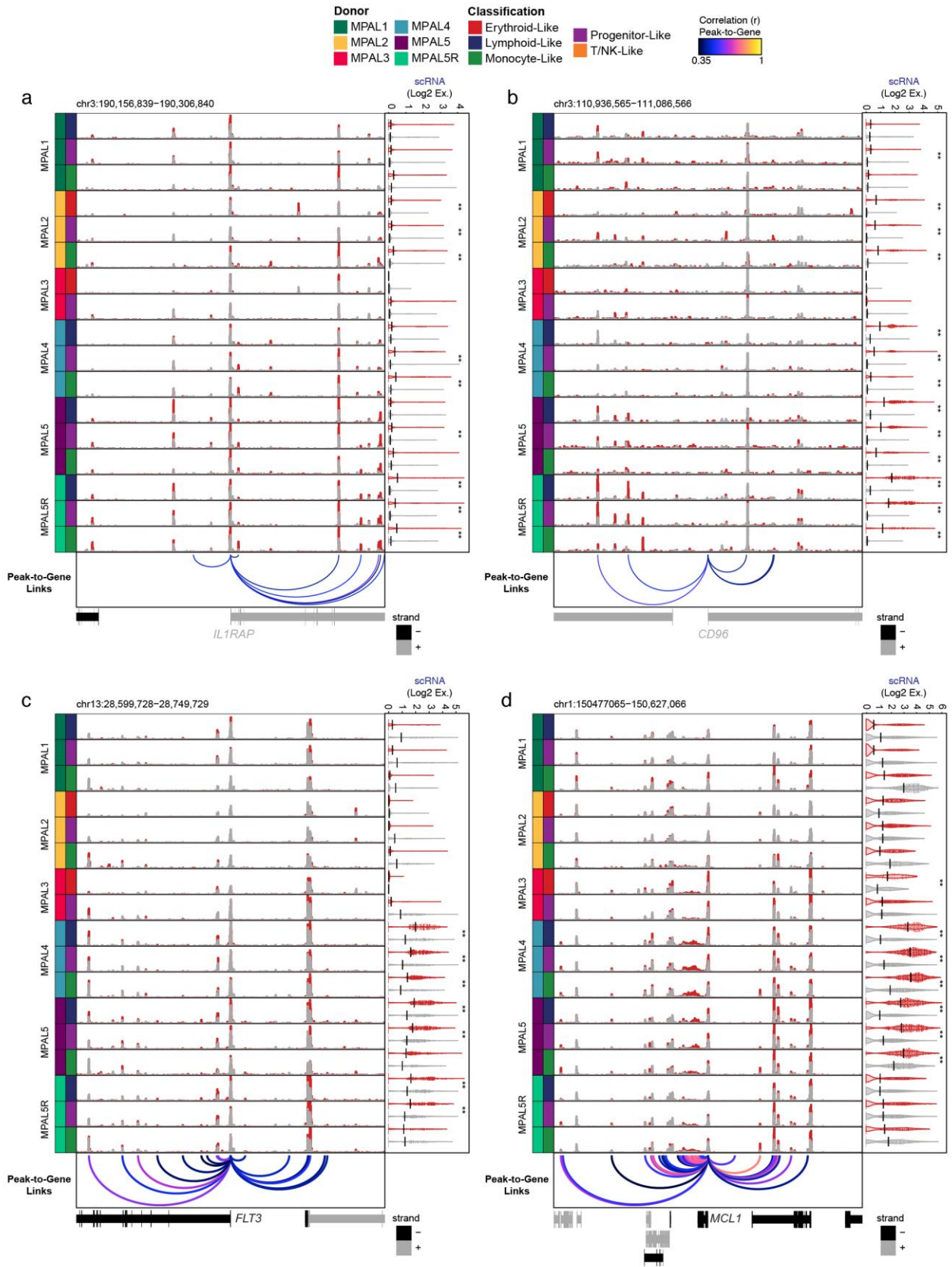
(a) UMAP of CCA alignment of scATAC-seq using Cicero gene activity scores and scRNA-seq for (Left) bone marrow (nATAC = 12,602; nRNA = 16,510), (Middle) CD34+ enriched bone marrow (nATAC = 8,176; nRNA = 10,160), (Right) peripheral blood (nATAC = 14,804; nRNA = 8,368). **(b)** UMAP of CCA alignment of scATAC-seq using Cicero gene activity scores and scRNA-seq for MPAL samples 1-5R (nATAC = 4,127 – 8,255; nRNA = 835 – 5,885).

a**b****Classification****c****d****e****f****g**

Supplementary Figure 12

Evaluation of scRNA and scATAC-seq alignment and peak-to-gene linkage across hematopoiesis and MPAL samples.

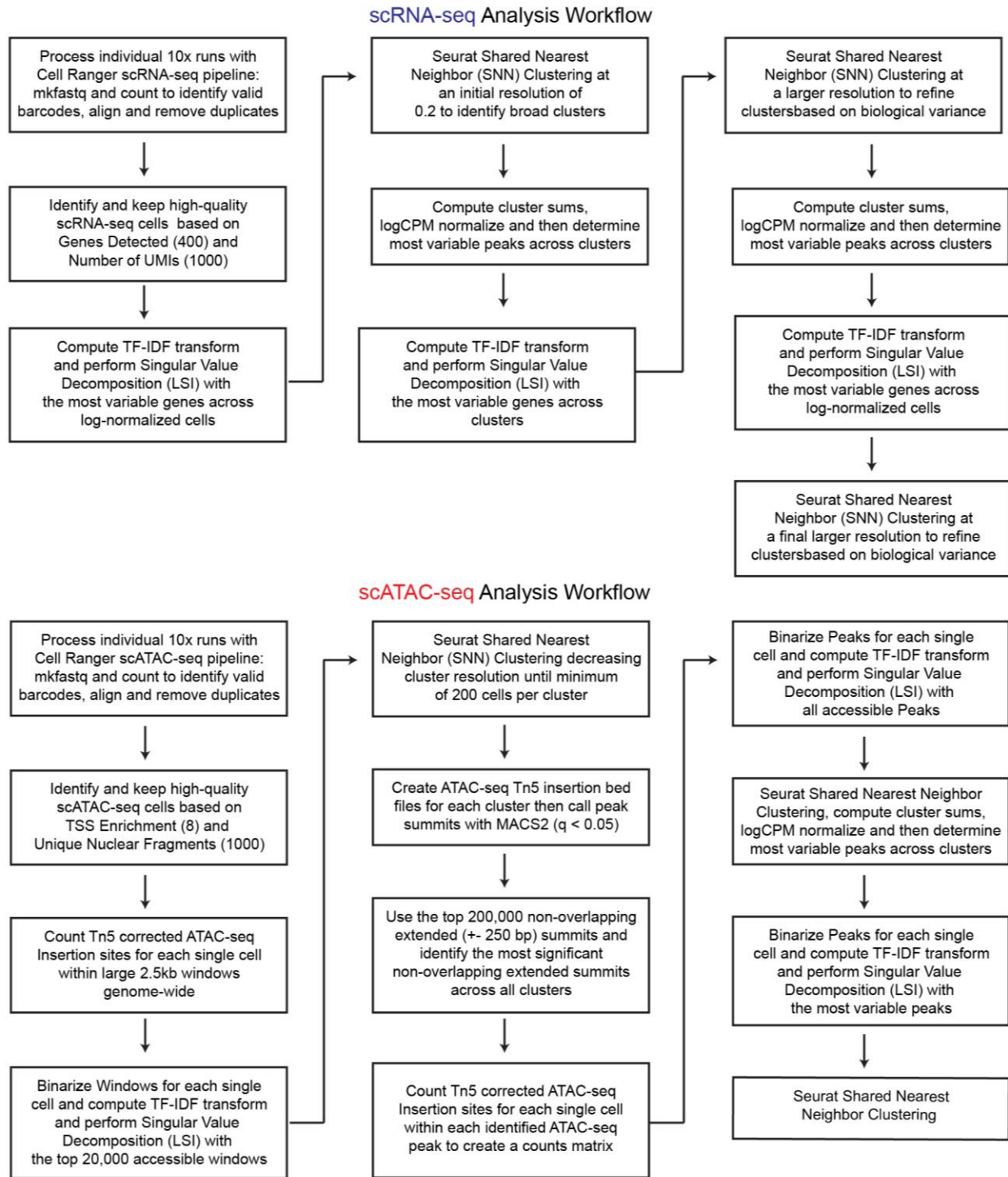
(a) Spearman rank correlation between scATAC-seq Cicero gene activity scores to scRNA-seq for each mapped cell within across all biological experiments ($n = 4,127 - 16,510$). **(b)** Pearson correlation of CCA scRNA and scATAC-seq nearest-neighbors. The cutoff ($R > 0.45$) for high quality nearest neighbor mappings is shown ($n = 4,127 - 16,510$). **(c)** (Left) UMAP of scATAC-seq hematopoiesis colored by scATAC-seq clusters ($n = 35,038$). (Right) UMAP of scATAC-seq hematopoiesis colored by filtered mapped scRNA-seq clusters ($n = 34,507$). **(d)** Confusion matrix of initial clusters for mapped scRNA-seq to scATAC-seq clusters for hematopoiesis (Figure 1b-c). Above shows the assigned biological classifications for each cluster across both scRNA and scATAC-seq. **(e)** (Left) Distribution of peak-to-gene distances. (Left-Middle) Distribution of number of peaks mapped per gene (median = 6). (Right-Middle) Distribution of number of genes mapped per peak (median = 1). (Right) Distribution of number of genes skipped for peak-to-gene links (median = 2). **(f)** MetaV4C plots of K27ac HiChIP in Naive T and HCASMC cells for top 500 biased T/NK (broad classification) peak-to-gene links that are identified only in healthy hematopoiesis. The line represents the average signal and the shading represents the range of the signal times the square root of 2 between biological replicates ($n = 2$). **(g)** Peak-to-gene links ($n = 91,601$) enrichment in GTEx eQTLs over a permuted background distance-matched set (permutations = 250) for the union set of peak-to-gene links. The mean enrichment is shown, and the error bars indicate 1 standard deviation. Box-whisker plot; lower whisker is the lowest value greater than the 25% quantile minus 1.5 times the interquartile range (IQR), the lower hinge is the 25% quantile, the middle is the median, the upper hinge is the 75% quantile and the upper whisker is the largest value less than the 75% quantile plus 1.5 times the IQR.



Supplementary Figure 13

Peak-to-gene links nominate putative regulatory regions that nominate key leukemic genes.

(a-d) Multi-omic differential track; (Middle) Aggregated scATAC tracks showing MPAL disease subpopulations (red) and closest normal cells (grey). (Right) Distribution of log₂ normalized expression of gene of interest for MPAL disease subpopulations (red) and closest normal cells (grey); black line represents the mean and asterisk denote significance (LFC > 0.5 and FDR < 0.01 from Figure 2c). Violin plot represents the smoothed density of the distribution of the log₂ normalized expression and the black line represents the mean log₂ normalized expression. (Bottom) Peak-to-gene links for gene of interest colored by Pearson correlation of the peak accessibility and gene expression (see methods). **(a)** Multi-omic differential track for *IL1RAP* (n = 89 - 500). **(b)** Multi-omic differential track for *CD96* (n = 89 - 500). **(c)** Multi-omic differential track for *FLT3* (n = 89 - 500). **(d)** Multi-omic differential track for *MCL1* (n = 89 - 500).



Supplementary Figure 14

Analysis workflows for processing of scRNA-seq and scATAC-seq data.

(Top) scRNA-seq analysis workflow. Briefly cells are aligned using 10x cell ranger, quality filtered, and clustered using a feature optimization approach (see methods). (Bottom) scATAC-seq analysis workflow. Briefly cells are aligned using 10x cell ranger atac, quality filtered, clustered in large windows genome-wide, peak-calling on clusters, creation of a counts matrix and clustered using a feature optimization approach (see methods). (see methods).

