

## **Supplemental data**

Supplemental Methods.....	2
Supplemental Figures.....	11
Supplemental References.....	20

## 1 **Supplemental Methods**

### 2 **Cell lines**

3 NB4 cells were kindly provided by M. Lanotte (St. Louis Hospital, Paris, France).  
4 HEK-293T cells were obtained from the Cell Bank of the Chinese Academy of Science.  
5 NB4 cells were cultured in the RPMI-1640 medium supplemented with 10% fetal  
6 bovine serum (FBS) (Gibco, Carlsbad, CA, USA) . HEK-293T cells were cultured in  
7 DMEM supplemented with 10% FBS. All cell lines were authenticated by STR  
8 sequencing. Mycoplasma contamination was routinely detected by the one-step  
9 Quickcolor Mycoplasma detection Kit (Cat. # MD001; Shanghai Yise Medical  
10 Technology, Shanghai, China,).

### 11 **Antibodies**

12 The anti-H3K27ac (Cat. # 39133; Active Motif, Carlsbad, CA, USA) and anti-MYB  
13 (Cat. # ab45150; Abcam, Cambridge, UK) antibodies were used in chromatin  
14 immunoprecipitation (ChIP) assays. The anti-MYB (Cat. # sc-8558; Santa Cruz, CA,  
15 USA) and anti- $\beta$ -actin (Cat. # 66009; Proteintech, Wuhan, China) antibodies were used  
16 in Western blot experiments.

### 17 **Whole genome sequencing**

18 Whole genome sequencing (WGS) was performed on the paired tumor and normal  
19 samples from the same patient. High throughput sequencing was performed using the  
20 Illumina Hiseq X or NovaSeq 6000 platform. Raw reads were aligned to the reference  
21 human genome hg38 (Genome Reference Consortium GRCh38) using the BWA-MEM  
22 algorithm<sup>1</sup>. For datasets with multiple data lanes, the aligned files were merged using  
23 Picard MergeSamFiles (v25.0), and duplicate reads were marked with Picard  
24 MarkDuplicates (v25.0). In addition, the local realignment around indels and the base  
25 quality score recalibration functions were performed using GATK (v3.8). The quality  
26 of the alignments was assessed by metrics determined by Samtools stats (v1.6)<sup>2</sup>; the  
27 coverage and depth of alignments were assessed through each base depth calculated by  
28 Samtools depth (v1.6). The mean depth of all samples was 52 $\times$ . Detailed statistics on  
29 alignment and coverage are given in supplemental Table 3.

## 1 **ChIP-seq library construction**

2 APL blasts and NB4 cells were cross-linked with 1% formaldehyde (Cat. # F8775;  
3 Sigma, St Louis, MO, USA) for 10 min at room temperature, and the fixation reactions  
4 were quenched by adding glycine to a final concentration of 125 mM. Cells were  
5 washed twice in PBS, then incubated in lysis buffer (10 mM Tris-HCl pH 7.5, 10 mM  
6 NaCl, 3 mM MgCl<sub>2</sub>, 0.5% NP-40, 1 mM phenylmethylsulfonyl fluoride (PMSF, cat. #  
7 P7626, Sigma) and 1× protease inhibitor cocktail (Cat. # 11873580001; Roche,  
8 Mannheim, Germany) for 10 min. The pellets were re-suspended in the sonication  
9 buffer (20 mM Tris-HCl pH 8.0, 2 mM EDTA, 1% Triton X-100, 150 mM NaCl, 1 mM  
10 PMSF and 1× protease inhibitor cocktail). Chromatin was sonicated using the Biorupter  
11 Pico ultrasonicator device (Diagenode, Liege, Belgium) with 25 cycles of 30-second  
12 ON and 30-second OFF. Sheared chromatin was then incubated overnight at 4 °C with  
13 the indicated antibody. Protein G magnetic beads (Cat. # 10004D; Thermo Fisher  
14 Scientific, Norcross, GA, USA) were added to the immunoprecipitation reaction and  
15 incubated for another 2 hours at 4 °C. After finishing incubation, beads were washed  
16 for one time with low salt buffer (0.1% SDS, 1% TritonX-100, 2 mM EDTA, 20 mM  
17 Tris-HCl pH 8.0, 150 mM NaCl), high salt buffer (0.1% SDS, 1% TritonX-100, 2 mM  
18 EDTA, 20 mM Tris-HCl pH 8.0, 500 mM NaCl), LiCl buffer (250 mM LiCl, 1% NP-  
19 40, 1% sodium deoxycholate, 1 mM EDTA, 10 mM Tris-HCl pH 8.0) and two washes  
20 with TE buffer (10 mM Tris-HCl pH 8.0, 1 mM EDTA). The ChIP DNA was eluted  
21 using 100 μL of elution buffer (10 mM Tris, pH 8.0, 1 mM EDTA, pH 8.0, 1 % SDS,) at  
22 55 °C for 10min. Cross-linking was reversed by adding 5 μL of proteinase K (Cat. #  
23 P8107; New England Biolabs, Ipswich, MA, USA) and incubating overnight at 65 °C.  
24 Input and ChIPed DNA were purified using the DNA purification kit (Cat. # 28106;  
25 QIAGEN, Valencia, CA, USA). ChIP-qPCR was conducted in triplicate using the  
26 primers listed in supplemental Table 17. ChIP-seq libraries were prepared with the  
27 MicroPlex Library Preparation Kit v2 (Cat. # C05010014; Diagenode) according to the  
28 manufacturer's instructions. High throughput sequencing was performed on the  
29 Illumina Hiseq X or NovaSeq 6000 platform.

## 30 **Collection of ChIP-seq data**

31 ChIP-seq data for GFI1, IRF1, RUNX1, PML/RAR $\alpha$ , KLF13, MEF2D, NFE2 and

1 ETV6 were collected from the Gene Expression Omnibus (GEO) database under the  
2 accession numbers: GSM935505, GSM2026066, SRP103029, ENCSR608HVP\_2,  
3 ENCSR647ZXA\_1, GSM2527371 and GSM2527376. ChIP-seq data for MYB in  
4 MOLT-3, Jurkat and K562 were collected from the GEO database under the following  
5 accession numbers: GSM1519643, GSM1519641 and GSM2825506. ChIP-seq for  
6 H3K4me1, H3K4me3, H3K27ac and H3K27me3 in K562 were collected from the  
7 following accession code: GSM1782706, GSM2534289, GSM2877120 and  
8 GSM608166. ChIP-seq for H3K4me1, H3K4me3, H3K27ac and H3K27me3 in  
9 Kasumi-1 were collected from the GEO database under the following accession  
10 numbers: GSM3165518, GSM1534445, GSM3165517, GSM1534446. H3K27ac  
11 ChIP-seq for hematopoietic cell lines were collected from the GEO database under the  
12 following access numbers, i.e., GSM2836487 (HL-60), GSM2108046 (THP-1),  
13 GSM2136946 (MV4-11), GSM2136938 (MOLM-13), GSM3094374 (EOL-1),  
14 GSM3436213 (SET-2), GSM2445788 (SKNO-1), GSM1003462 (DND-41),  
15 GSM1246865 (Jurkat), and GSM3425377 (Nalm-6).

## 16 **ChIP-seq data processing**

17 The ChIP-seq raw data obtained from sequencing reactions or the above collections  
18 were aligned using Bowtie2<sup>3</sup> version 2.3.5.1 against the human genome hg38. Picard  
19 MarkDuplicates version 25.0 was used to remove the duplicate reads. The aligned reads  
20 were further normalized to the same library size by MACS2<sup>4</sup> randsample with  
21 NUMBER = 1e-7, SEED = 614, and default settings.

## 22 **Annotations of CREs in APL patients**

23 The H3K27ac ChIP-seq data from APL patients was used to define the CREs. The BED  
24 files of H3K27ac-positive peaks from patients, called using MACS2 as described above,  
25 were used to perform sample saturation analysis. Saturation analysis was based on the  
26 permutation process, randomly selecting the number of desired samples in 1,000  
27 iterations to combine and calculate the number of peaks. The permutation process was  
28 executed as a loop from 1 to the number of samples. We could evaluate the saturation  
29 of H3K27ac-positive regions based on a non-linear regression model fitted on the  
30 results of saturation analysis. Then, the H3K27ac-positive peaks of all patients were  
31 merged by Bedtools<sup>5</sup> v2.27.1 with the default setting. Considering that CREs are

1 usually flanked by H3K27ac positive regions<sup>6,7</sup>, we further obtained the APL-  
2 associated CREs by merging the H3K27ac-positive peaks through extending by  $\pm 500$   
3 bp (the average length of transcription factor (TF)-bound regions). The HOMER<sup>8</sup> peak  
4 annotate tool (annotatePeaks.pl) version 4.9.1 was used to annotate the relative genomic  
5 distribution of CREs, and the nearest neighbor gene of a given CRE was assigned as  
6 the CRE-associated gene.

7 To compare mutated CREs in APL versus those in other cancer types, we used WGS  
8 and H3K27ac ChIP-seq data to establish the landscapes of mutated CREs in other types  
9 of hematopoietic malignancies and solid cancers. Hematopoietic malignancies included  
10 chronic lymphocytic leukemia (CLLE) and malignant lymphoma (MALY). Solid  
11 cancers included bone cancer (BOCA), breast cancer (BRCA), liver cancer (LIRI),  
12 pancreatic cancer (PACA), pediatric brain cancer (PBCA), and prostate  
13 adenocarcinoma (PRAD). The variant call format files of WGS data for the above  
14 cancer types were downloaded from the Pan-Cancer Analysis of Whole Genomes  
15 project (PCAWG)<sup>9</sup>. The H3K27ac ChIP-seq data of these cancer types were  
16 downloaded from the Gene Expression Omnibus (GEO) database. Briefly, we  
17 constructed somatic mutation profiles for each cancer type by combining variant call  
18 data of multiple samples. We then established the CRE profiles for each cancer type  
19 using H3K27ac ChIP-seq data for the corresponding cell line. Then, we identified the  
20 mutated CREs for each cancer type by integrating the mutation profiles data and the  
21 CRE profiles (Detailed statistics in supplemental Table 7). We performed Gene  
22 Ontology (GO) enrichment analysis to identify the enriched pathways within mutated  
23 CRE-regulated genes in APL and other cancer types.

#### 24 **Identification of master TFs based on CRC analysis**

25 To construct the CRC model in APL, we first used the Rank Ordering of Super  
26 Enhancers (ROSE) algorithm<sup>10</sup> to define the super enhancers (SEs) of 16 APL samples  
27 based on the above described H3K27ac ChIP-seq data. The stitching distance was fixed  
28 at 12.5 kb to facilitate comparisons between samples. For other parameters, the default  
29 settings were used. Genes annotated by the ROSE2 ENHANCER\_TO\_TOP\_GENE.txt  
30 file were used for defining the target genes of SEs for subsequent analyses. Then, we  
31 applied the CRC mapper algorithm<sup>11</sup> to construct the CRC model in APL based on the

1 super-enhancer profiles of these 16 APL patients. A total of 25 transcription factors were  
2 identified using the 16 CRC models. The transcription factors included in more than  
3 25% of patients were considered as master transcription factors of APL. In the CRC  
4 analysis, the refGene.txt of hg38 downloaded from UCSC was used to annotate the  
5 genome files, and the motifs of transcription factors scanned by FIMO were obtained  
6 from the JASPAR database<sup>12</sup>. The previously reported RARE-half site bound by  
7 PML/RAR $\alpha$ <sup>13</sup> was included in the scanning process.

### 8 **Enrichment of binding regions for master transcription factors within mutated** 9 **CREs**

10 We collected the genomic binding regions of the master TFs in the hematopoietic cell  
11 lines. The genomic binding regions of MYB, IRF1, PU.1 and RAR $\alpha$  were obtained by  
12 analyzing their ChIP-seq data in NB4. The genomic binding regions of other TFs were  
13 obtained from the Cistrome database<sup>14</sup>. We used bedtools intersect v2.27.1 to calculate  
14 the overlap of the TF binding regions with mutated or unmutated CREs. The Fisher's  
15 exact test was used to evaluate the significance.

### 16 **Motif analysis**

17 To investigate whether mutated CREs were directly targeted by the identified  
18 transcription factors, we used the HOMER<sup>8</sup> motif discovery tool (findMotifsGenome.pl)  
19 version 4.9.1 to perform the motif enrichment analysis in the mutated CREs versus non-  
20 mutated CREs. We downloaded the position weight matrices of these transcription  
21 factors and their paralogs from the JASPAR database. The RARE-half motif for  
22 PML/RAR $\alpha$  was included in the motif enrichment analysis. Then, we calculated the  
23 enrichment of the above motifs within the mutated CREs (as the target regions) versus  
24 the non-mutated CREs (as the background regions). The number of the background  
25 regions was set to twice the target regions by randomly selecting from non-mutated  
26 CREs. This process was iterated 200 times to get the final result.

27 To define the MYB motif in APL, we used the findMotifsGenome.pl program to  
28 perform the motif enrichment analysis based on MYB-specific ChIP-seq data in NB4  
29 cells. The top 1000 intensity peaks were selected as the target regions. The most  
30 significant motif was defined as the specific binding motif of MYB in APL.

## 1 **Mutation enrichment analysis within the motifs and flanking regions**

2 To assess the significance of the mutation enrichment within the motifs and flanking  
3 regions, we referred to the previously published method design<sup>7</sup> for mutation  
4 enrichment analysis and realized it through self-developed code. First, bedtools v2.27.1  
5 was used to extract the overlap between the ChIP-seq peaks of the identified  
6 transcription factors in hematopoietic cell lines and the profile of APL CREs as  
7 previously established. We obtained the potential binding regions of each TF in APL as  
8 the regions of interest for subsequent analysis. Second, we scanned the motif positions  
9 in the regions of interest as the binding sites of the corresponding TFs. MOODS<sup>15</sup> v1.9.4  
10 was used to match the position frequency matrices of motifs from the database JASPAR  
11 against DNA sequences in the regions of interest. We set the parameter *P* value to  
12 0.0001. Third, we calculated the mutation frequency at the motif positions of the  
13 transcription factor combined with the APL somatic mutation profile as established  
14 above and further determined the mutation load by expanding the motif positions  
15 according to the specified base number, including 10bp, 20bp, 30bp, 40bp, 50bp, 100bp,  
16 200bp, 300bp, 400bp, 500bp. Fourth, the permutation test was performed to calculate  
17 the mutation frequency within randomly selected positions from the APL H3K27ac-  
18 positive regions. The randomly selected positions were consistent in number and length  
19 with the binding sites of the corresponding TFs in the regions of interest. This process  
20 was performed in 5000 iterations for each transcription factor. Finally, we transferred  
21 the mutation frequency within the motifs and around regions to the *Z* score according  
22 to the permutation test to evaluate the significance of the mutation enrichment.

## 23 **RNA extraction, reverse transcription and real time PCR**

24 Total RNA was extracted using the RNeasy mini kit (Cat. # 74106; QIAGEN)  
25 according to the manufacturer's instructions. RNA was reverse transcribed into cDNA  
26 using the PrimeScript™ RT reagent Kit with gDNA Eraser (Perfect Real Time) (Cat. #  
27 RR047B; Takara, Osaka, Japan). RT-qPCR was conducted using the SYBR Green  
28 Premix pro Taq HS qPCR Kit (Rox Plus) (Cat. # AG11719; Accurate Biology,  
29 Changsha, China) on the ABI ViiA 7 Real-Time PCR System. The relative expression  
30 level of each gene was calculated as  $2^{-\Delta\Delta C_t}$ . Primers for real time-qPCR are listed in  
31 supplemental Table 18.

## 1 **RNA sequencing and data processing**

2 RNA-seq libraries were constructed according to the manufacture's instruction using  
3 the TruSeq RNA Sample Preparation Kit v2 (Cat. # RS-122-2001 or RS-122-2002;  
4 Illumina, Hayward, CA, USA). The purified library was quantified using Qubit 4  
5 Fluorometer (Invitrogen, Carlsbad, CA, USA), and size distribution was analyzed by  
6 Bioanalyzer 2100 (Agilent Technologies, Palo Alto, CA). High throughput sequencing  
7 was performed on the Illumina HiSeq 2500 or NovaSeq 6000 platform. The software  
8 Hisat2<sup>16</sup> was used to align the reads from RNAseq to the reference human genome hg38  
9 (Genome Reference Consortium GRCh38). The gene counts of each sample were  
10 generated using the HTseq-count<sup>17</sup> and were normalized to transcripts per million  
11 (TPM).

## 12 **Identification of the functional regions harboring somatic non-coding mutations** 13 **within mutated CREs**

14 To identify the functional regions harboring somatic non-coding mutations, we  
15 performed a comprehensive analysis based on the recurrence occurring within a small  
16 region, the transcriptional influences, and the disease relevance. First, we identified the  
17 recurrently mutated loci by clustering all somatic non-coding mutations within 50 bp.  
18 Second, we assigned the target gene to each recurrently mutated locus within CREs  
19 according to the principle of proximity, and then used Phenolyzer<sup>18</sup> to assess the priority  
20 of those target genes related to leukemia. Third, we evaluated the transcriptional effects  
21 of the recurrently mutated loci within CREs by comparing the expression levels of the  
22 respective target gene in mutated and non-mutated samples. Then, each recurrently  
23 mutated locus within CREs was placed in a 3D feature space taking into account  
24 mutation recurrence, leukemia relevance, and gene expression. Finally, we calculated  
25 the Euclidean distances based on the three features described above and obtained the  
26 functional regions harboring somatic non-coding mutations within mutated CREs.

## 27 **PCR and Sanger sequencing**

28 Non-coding *WT1* somatic or germline variants were validated by PCR amplification  
29 and Sanger sequencing in tumor and paired normal samples. A somatic mutation refers  
30 to an alteration that occurs in the tumor sample but not in the paired normal sample. A



1 germline variant refers to an alteration that exists in both the tumor and the paired  
2 normal sample. The PCR primers were designed using the Primer3 software  
3 (supplemental Table 19). The PCR products were sequenced on the ABI 3730XL DNA  
4 sequencer (Applied Biosystems).

### 5 **Luciferase reporter assay**

6 The *WT1* promoter and enhancer region with strong H3K27ac signals were amplified  
7 using genomic DNA from APL patients with or without non-coding *WT1* variants. The  
8 primers used were described in supplemental Table 20. The amplified regions were  
9 cloned into the pGL3-basic vector. The luciferase constructs, in combination with the  
10 pRL-SV40 renilla plasmid and MYB overexpressing plasmids (pENTER-MYB (Cat. #  
11 CH806231; WZ Biosciences Inc., China) and empty vector) or knockdown siRNAs (si-  
12 MYB and si-NC), were delivered into HEK-293T or NB4 cells. After transfection for  
13 24 hours, both firefly luciferase activity and renilla luciferase activity were detected  
14 with the GloMax 20/20 Luminometer (Promega) using the Dual-Luciferase Reporter  
15 Assay System (Cat. # E1910; Promega, Madison, WI, USA) following the  
16 manufacturer's instructions. The firefly luciferase activity was normalized with the  
17 renilla luciferase activity to control the transfection efficiency.

### 18 **DNA pulldown assay**

19 The 5'-biotinylated double-stranded DNA probes (supplemental Table 21) were  
20 incubated with the whole cell lysate of HEK-293T cells overexpressing MYB overnight  
21 at 4 °C. The magnetic streptavidin beads were then added into the complexes of  
22 reactions, followed by rotation at 4 °C for another 4 hours. Beads were washed three  
23 times with cold lysis buffer (20 mM Tris-HCl pH 8.0; 2 mM EDTA; 1% Triton X100;  
24 150 mM NaCl) and resuspended in SDS loading buffer for western blot analysis with  
25 anti-MYB antibody.

### 26 **CRISPR/Cas9-mediated editing and knockout**

27 The Cas9-expressing NB4 cells were generated by lentiviral transduction of  
28 LentiCRISPR v2 GFP (Addgene plasmid #82416), followed by flow cytometry sorting  
29 of GFP-positive cells. The sorted cells were then sub-cloned and selected for the best  
30 CRISPR/Cas9 efficiency clone using a lentiviral reporter pKLV-sgGFP.

1 The sgRNAs oligonucleotides containing the BbsI restriction sites (supplemental Table  
2 22) were constructed into the lentiviral vector pKLV-U6gRNA(BbsI)-PGKpuro2ABFP  
3 (Addgene plasmid #50946). Supernatants containing lentivirus were generated by co-  
4 transfecting 16 µg of the pKLV construct, 9 µg of the psPAX2 package plasmid and 6  
5 µg of the pMD2G envelop plasmid into HEK-293T cells using the HilyMax reagent  
6 (Cat. # H357-10; Dojindo, Tokyo, Japan) and were harvested at 48 or 72 hours after  
7 transfection. Cas9-expressing NB4 cells were infected using the lentivirus with  
8 polybrene at 8 µg/mL (Cat. # H9268; Sigma) by centrifuging directly at  $1200 \times g$  for  
9 90 min at 37 °C. For cells transduced with sgRNAs targeting the MYB motif or control  
10 sgRNAs targeting regions surrounding the MYB motif within the intron 3, single cells  
11 were sorted into 96 cell plates. Once single cells had grown into colonies, genomic  
12 DNA was extracted and analyzed for mutations by Sanger sequencing.

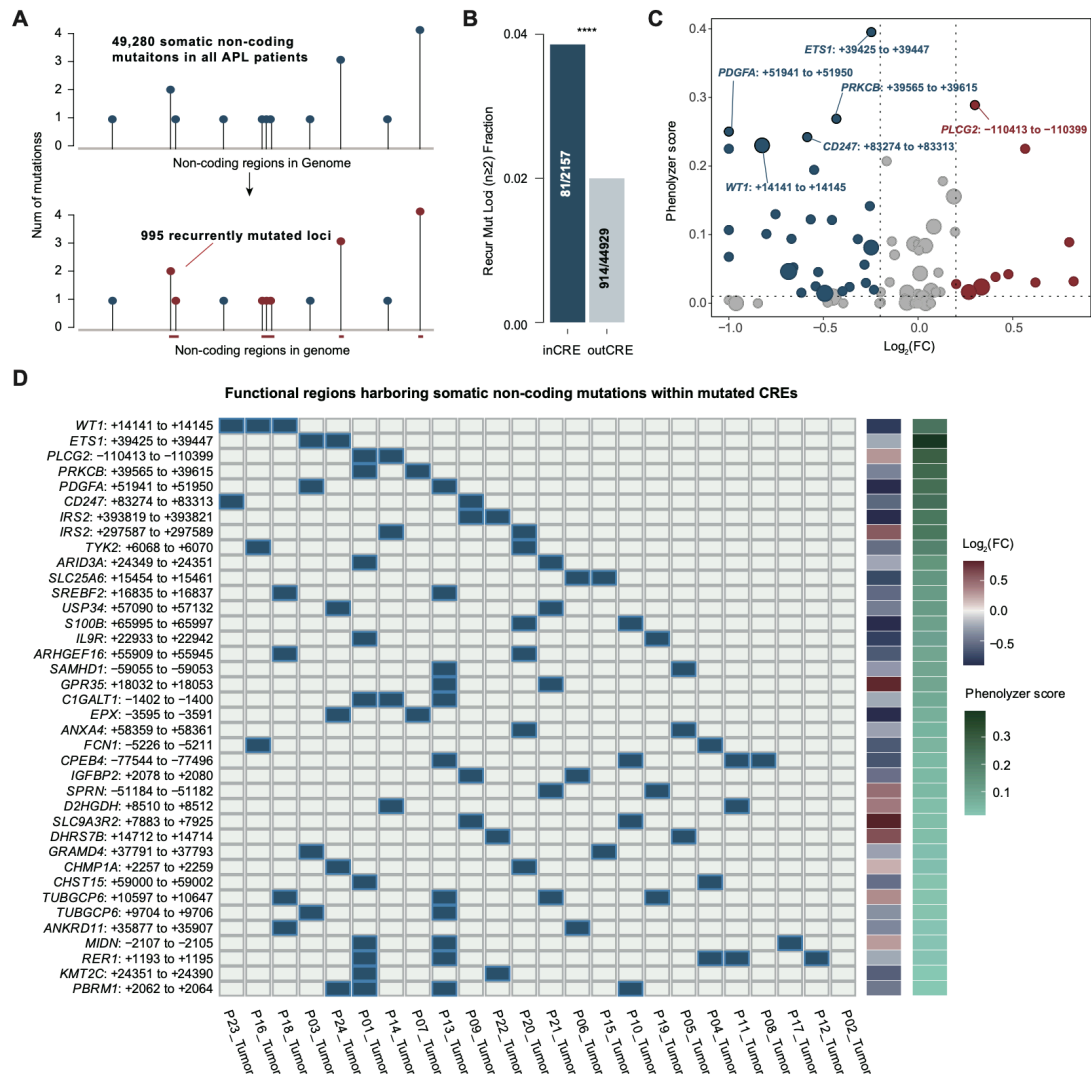
### 13 **Chromatin conformation capture (3C) experiment**

14 The 3C-qPCR experiments were performed as previously described<sup>19</sup> with some  
15 modifications. Briefly, cells were cross-linked with 1% formaldehyde and were  
16 quenched with 125 mM glycine. Cells were then lysed and resuspended in 0.3% SDS  
17 in  $1 \times$  NEBuffer 2.1 (10 µL of  $10 \times$  NEBuffer 2, 3 µL of 10% SDS, 87 µL of H<sub>2</sub>O) with  
18 shaking at 37 °C for 30 min. After adding Triton X-100 to a final concentration of 1%  
19 (v/v), the genomic DNA was digested overnight with XbaI (Cat. # R0145; New England  
20 Biolabs). DNA ligation was performed with the T4 DNA ligase (Cat. # M0202; New  
21 England Biolabs) at RT for 6 hours. The crosslinks were then reversed by overnight  
22 incubation at 65 °C with proteinase K. The digestion and ligation efficiencies were  
23 assessed using gel electrophoresis before proceeding to DNA purification by the  
24 phenol-chloroform extraction.

25 The purified 3C DNA was used to perform 3C-qPCR nearby the XbaI recognition sites  
26 among the regulatory regions of the *WT1* locus. The positive control sample was  
27 generated by amplifying, digesting and ligating fragments containing each of these  
28 seven XbaI sites with primers nearby the restriction sites, which was sequentially  
29 diluted to correct the PCR efficiency. The fragment within two adjacent XbaI sites in  
30 the *WT1* locus served as the loading control. The interaction frequency was calculated  
31 as the ratio relative to the region A in the vector. All the primers were designed using

- 1 Primer3 (supplemental Table 23). qPCR was performed using SsoFast EvaGreen
- 2 Supermix (Cat. #172-5211; Bio-Rad, Hercules, CA, USA).

# 1 Supplemental Figures

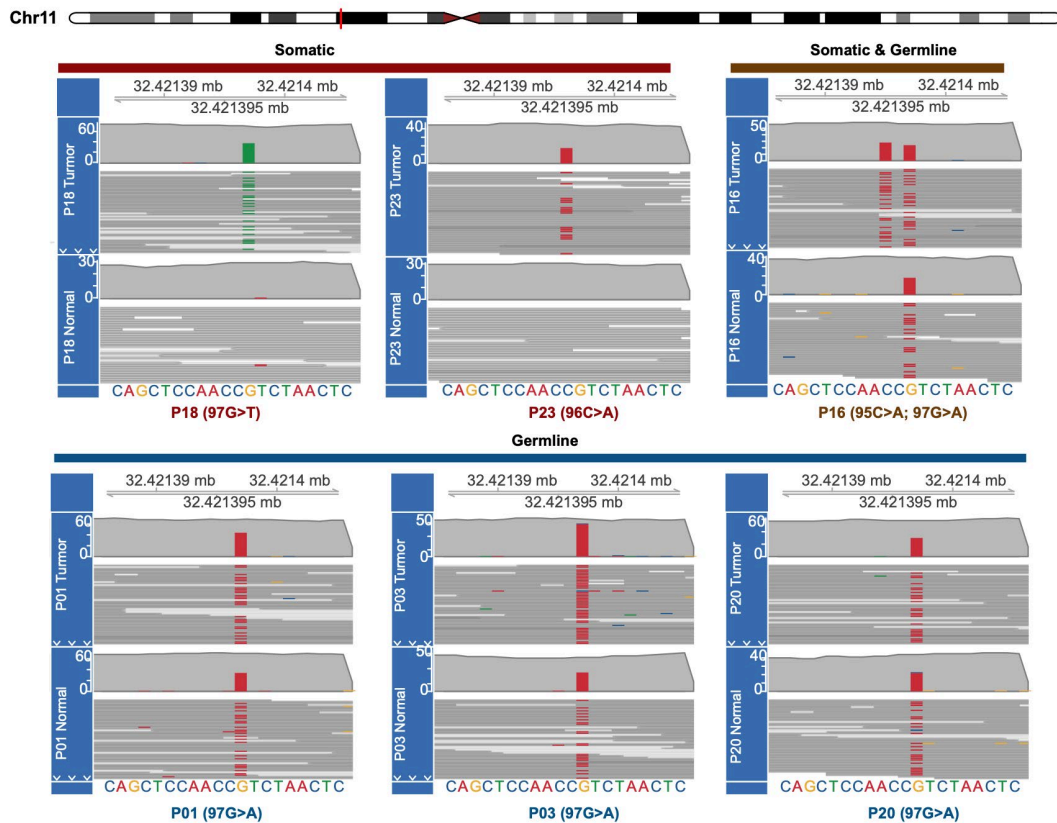


2

## 3 Supplemental Figure 1. Identification of the functional regions harboring somatic 4 non-coding mutations within mutated CREs.

5 (A) Diagram of the clustering of somatic non-coding mutations within 50 bp to identify  
6 the recurrently mutated loci. A total of 995 recurrently mutated loci were found from  
7 49,280 somatic non-coding mutations. (B) Distribution of the recurrently mutated loci  
8 according to whether they are located within CREs. The bar plot shows the distribution  
9 of the recurrently mutated loci within CREs. We identified a total of 81 recurrently  
10 mutated loci, which were significantly enriched within CREs relative to other non-  
11 coding regions. The statistical significance was calculated by the Fisher's exact test.  
12 \*\*\*\* $P < 0.0001$ . (C) Presentation of each recurrently mutated locus within CREs in a

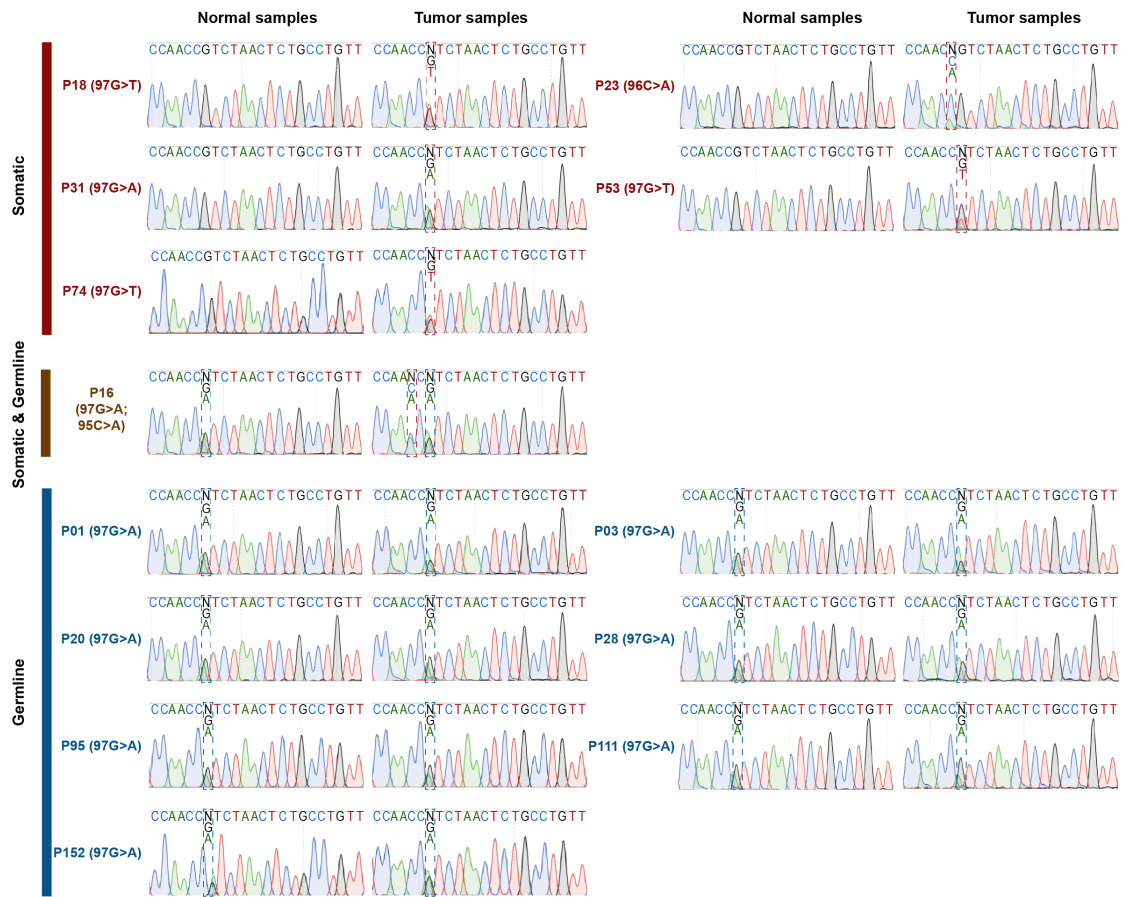
1 3D feature space, considering the mutation recurrence, leukemia relevance, and gene  
2 expression. The scatter plot shows the fold change and Phenolyzer score of 81  
3 recurrently mutated loci within CREs. The size of the point represents the number of  
4 the mutation recurrences. The blue and red points represent the candidate functional  
5 mutated loci that inhibit and promote the expression of the related target gene,  
6 respectively. **(D)** Displaying identified functional regions harboring somatic non-  
7 coding mutations within mutated CREs. The abscissa represents each patient, the  
8 ordinate represents the functional regions harboring somatic non-coding mutations for  
9 the indicated gene, the blue represents the sample with the specified mutation, and the  
10 gray represents the sample that does not contain the specified mutation. Also included  
11 is the log<sub>2</sub> fold change of associated genes by comparing the expression levels of  
12 respective target genes in mutated and non-mutated samples. Different shades of green  
13 squares represent the Phenolyzer score.



1

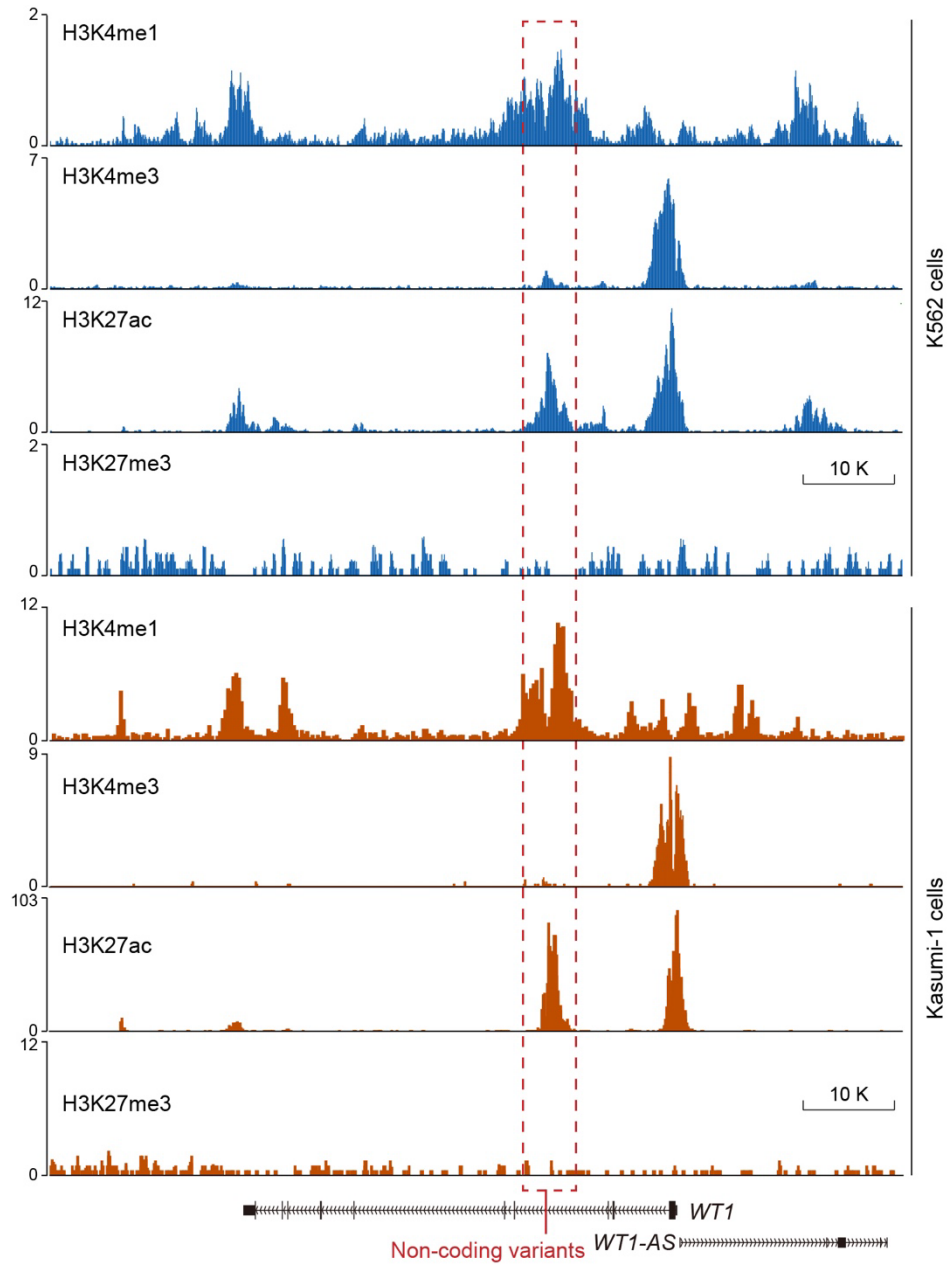
2 **Supplemental Figure 2. Visualization of the alignment pileup near the non-coding**  
 3 ***WT1* somatic mutations and/or germline variants in the tumor sample and the**  
 4 **paired normal sample from the same patient.**

5 P18 and P23 contained somatic *WT1* non-coding mutations. P16 harbored a somatic  
 6 *WT1* non-coding mutation and a germline *WT1* non-coding variant on two alleles,  
 7 respectively. P01, P03, and P20 contained germline *WT1* non-coding variants. A  
 8 somatic mutation refers to an alteration that occurs in the tumor sample but not in the  
 9 paired normal sample. A germline variant refers to an alteration that exists in both the  
 10 tumor and the paired normal sample.



1

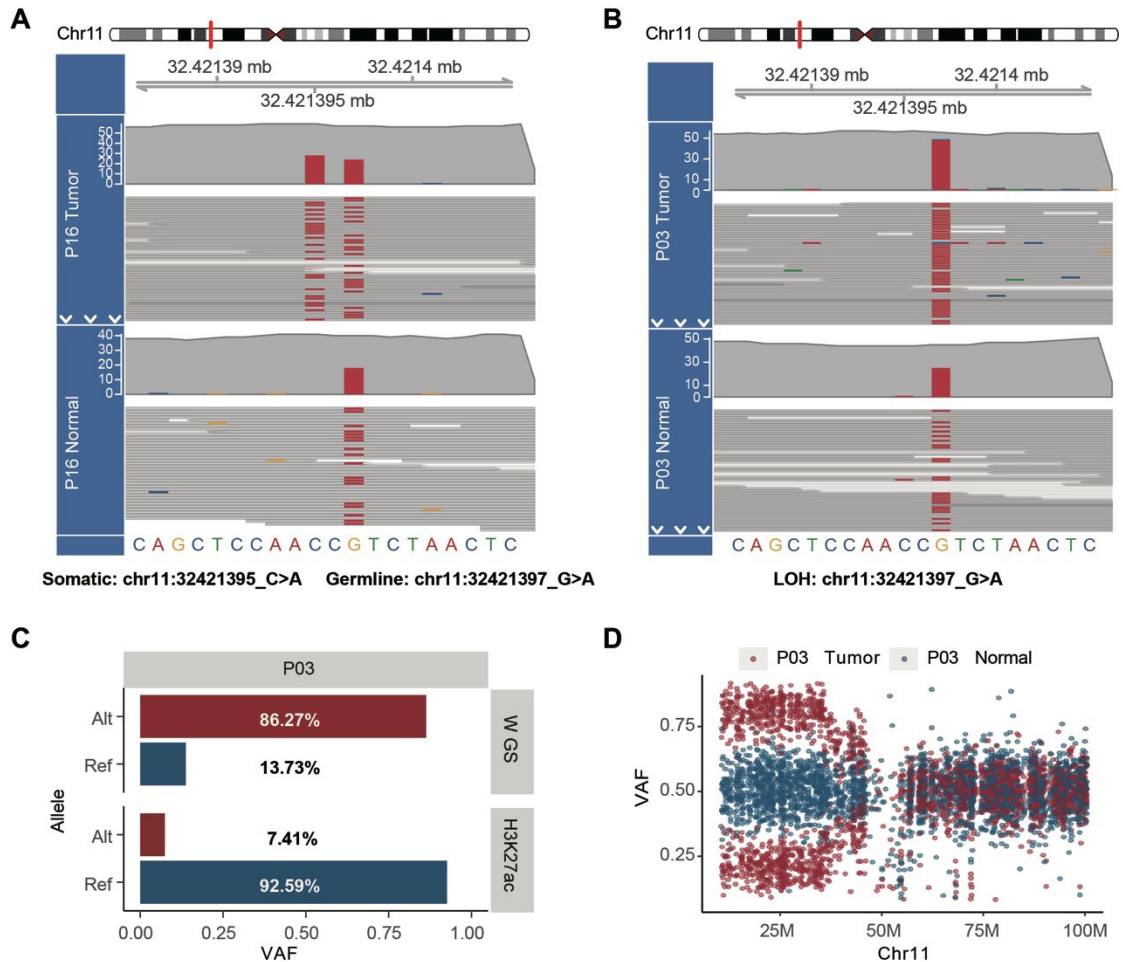
2 **Supplemental Figure 3.** The Sanger sequencing chromatograms of the third intron of  
 3 *WT1* in tumor and matched normal samples from APL patients with non-coding *WT1*  
 4 variants. Somatic and germline variants were identified by comparing the DNA  
 5 sequences of tumor to normal samples.



1

2 **Supplemental Figure 4.** ChIP-seq tracks of H3K4me1, H3K4me3, H3K27ac and  
 3 H3K27me3 occupancy at the *WT1* locus in K562 and Kasumi-1 cells without non-  
 4 coding *WT1* variants. The recurrently mutated site on the third intron of *WT1* is marked  
 5 with a red line.

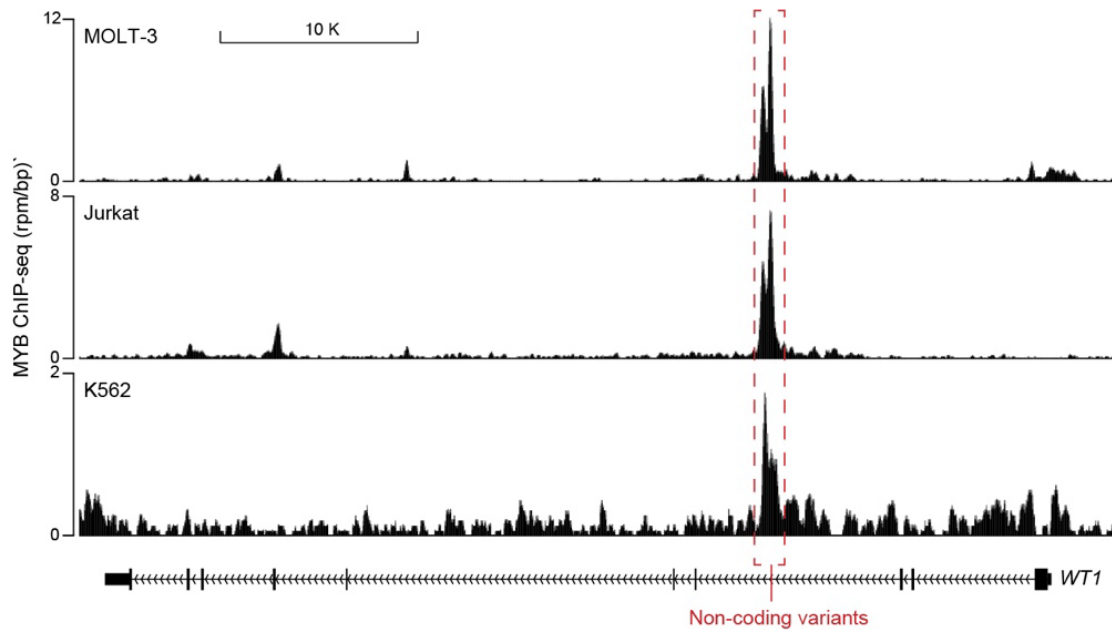




1

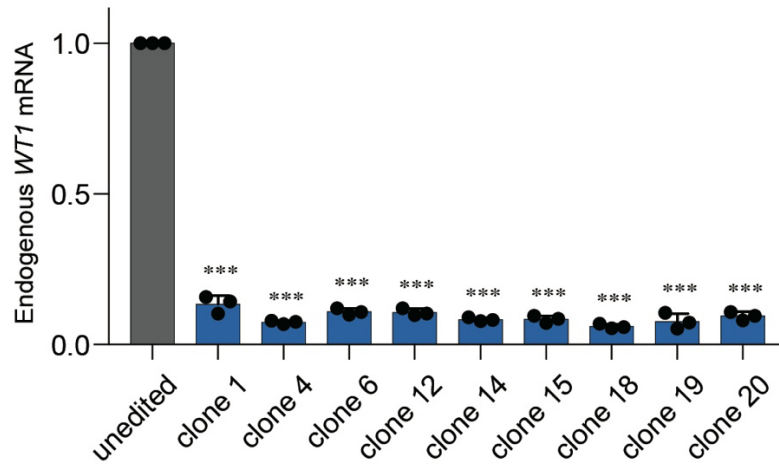
2 **Supplemental Figure 5. Two APL patients homozygous for the non-coding *WT1***  
 3 **variants.**

4 **(A-B)** Visualizing the alignment of WGS reads for non-coding *WT1* variants. **(A)** The  
 5 alignment track shows that P16 harbors a somatic mutation and a germline variant on  
 6 two alleles. **(B)** The alignment track shows that P03 contains the same germline variant  
 7 in a homozygous state due to copy-neutral LOH. **(C)** The variant allele frequency (VAF)  
 8 of the non-coding *WT1* variants in P03 based on the WGS data confirmed the  
 9 homozygous status. Still, the variant allele frequency shows extremely low levels in  
 10 H3K17ac ChIP-seq data, indicating the inhibitory effect of the non-coding *WT1* variant  
 11 on H3K27ac binding. **(D)** Tumor-specific loss of heterozygosity (LOH) for  
 12 chromosome 11p in P03 where the non-coding *WT1* variant is located. LOH was  
 13 assessed by examining variant allele frequencies (VAFs) for variants found to be  
 14 heterozygous in the normal sample. The red and blue points represent the VAFs of  
 15 single nucleotide polymorphisms (SNPs) in the tumor and the paired normal samples.



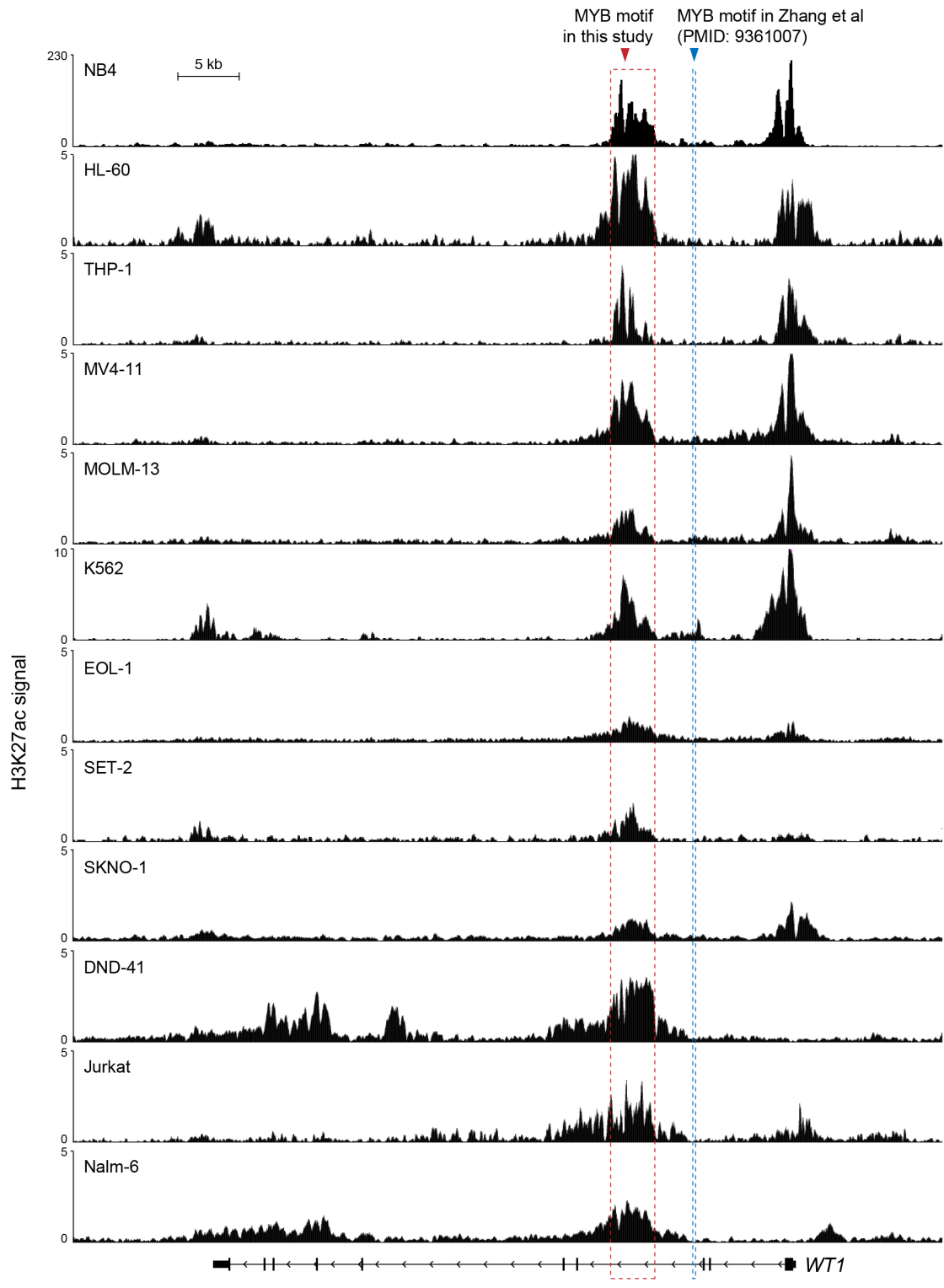
1

2 **Supplemental Figure 6.** ChIP-seq tracks showing MYB binding to the third intron of  
 3 the *WT1* gene in hematopoietic cell lines without non-coding *WT1* variants, i.e., MOLT-  
 4 3, Jurkat and K562. The recurrently mutated site on the third intron of *WT1* is marked  
 5 with a red line.



1

2 **Supplemental Figure 7.** The *WT1* mRNA levels in control and the MYB motif-mutated  
 3 clones. Data are represented as the fold change relative to the expression of the control  
 4 cells. \*\*\* $P < 0.001$ .



1  
 2 **Supplemental Figure 8. ChIP signals for H3K27ac at the WT1 locus in**  
 3 **hematopoietic cells.** Our identified enhancer and a previously reported MYB motif-  
 4 containing region<sup>20</sup> are shown by the red and blue dashed-line squares, respectively.  
 5 The MYB motif is indicated by the red arrow.

## 1 Supplemental References

- 2 1. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler  
3 transform. *Bioinformatics*. 2010;26(5):589-595.
- 4 2. Li H, Handsaker B, Wysoker A, et al. The Sequence Alignment/Map format and  
5 SAMtools. *Bioinformatics*. 2009;25(16):2078-2079.
- 6 3. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat*  
7 *Methods*. 2012;9(4):357-359.
- 8 4. Zhang Y, Liu T, Meyer CA, et al. Model-based analysis of ChIP-Seq (MACS).  
9 *Genome Biol*. 2008;9(9):R137.
- 10 5. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic  
11 features. *Bioinformatics*. 2010;26(6):841-842.
- 12 6. Kron KJ, Murison A, Zhou S, et al. TMPRSS2-ERG fusion co-opts master  
13 transcription factors and activates NOTCH signaling in primary prostate cancer. *Nat*  
14 *Genet*. 2017;49(9):1336-1345.
- 15 7. Mazrooei P, Kron KJ, Zhu Y, et al. Cistrome Partitioning Reveals Convergence of  
16 Somatic Mutations and Risk Variants on Master Transcription Regulators in Primary  
17 Prostate Tumors. *Cancer Cell*. 2019;36(6):674-689 e676.
- 18 8. Heinz S, Benner C, Spann N, et al. Simple combinations of lineage-determining  
19 transcription factors prime cis-regulatory elements required for macrophage and B cell  
20 identities. *Mol Cell*. 2010;38(4):576-589.
- 21 9. Consortium ITP-CAoWG. Pan-cancer analysis of whole genomes. *Nature*.  
22 2020;578(7793):82-93.
- 23 10. Loven J, Hoke HA, Lin CY, et al. Selective inhibition of tumor oncogenes by  
24 disruption of super-enhancers. *Cell*. 2013;153(2):320-334.
- 25 11. Saint-Andre V, Federation AJ, Lin CY, et al. Models of human core transcriptional  
26 regulatory circuitries. *Genome Res*. 2016;26(3):385-396.
- 27 12. Fornes O, Castro-Mondragon JA, Khan A, et al. JASPAR 2020: update of the open-  
28 access database of transcription factor binding profiles. *Nucleic Acids Res*.  
29 2020;48(D1):D87-D92.
- 30 13. Wang K, Wang P, Shi J, et al. PML/RARalpha targets promoter regions containing  
31 PU.1 consensus and RARE half sites in acute promyelocytic leukemia. *Cancer Cell*.  
32 2010;17(2):186-197.

- 1 14. Liu T, Ortiz JA, Taing L, et al. Cistrome: an integrative platform for transcriptional  
2 regulation studies. *Genome Biol.* 2011;12(8):R83.
- 3 15. Korhonen JH, Palin K, Taipale J, Ukkonen E. Fast motif matching revisited: high-  
4 order PWMs, SNPs and indels. *Bioinformatics.* 2017;33(4):514-521.
- 5 16. Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. Graph-based genome alignment  
6 and genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol.* 2019;37(8):907-  
7 915.
- 8 17. Anders S, Pyl PT, Huber W. HTSeq--a Python framework to work with high-  
9 throughput sequencing data. *Bioinformatics.* 2015;31(2):166-169.
- 10 18. Yang H, Robinson PN, Wang K. Phenolyzer: phenotype-based prioritization of  
11 candidate genes for human diseases. *Nat Methods.* 2015;12(9):841-843.
- 12 19. Tan Y, Wang X, Song H, et al. A PML/RARalpha direct target atlas redefines  
13 transcriptional deregulation in acute promyelocytic leukemia. *Blood.*  
14 2021;137(11):1503-1516.
- 15 20. Zhang X, Xing G, Fraizer GC, Saunders GF. Transactivation of an intronic  
16 hematopoietic-specific enhancer of the human Wilms' tumor 1 gene by GATA-1 and c-  
17 Myb. *J Biol Chem.* 1997;272(46):29272-29280.
- 18