

In the format provided by the authors and unedited.

# Somatic mutations precede acute myeloid leukemia years before diagnosis

**Pinkal Desai<sup>1,7\*</sup>, Nuria Mencia-Trinchant<sup>1,7</sup>, Oleksandr Savenkov<sup>2</sup>, Michael S. Simon<sup>3</sup>, Gloria Cheang<sup>4</sup>, Sangmin Lee<sup>1</sup>, Michael Samuel<sup>1</sup>, Ellen K. Ritchie<sup>1</sup>, Monica L. Guzman<sup>1</sup>, Karla V. Ballman<sup>2</sup>, Gail J. Roboz<sup>1,8</sup> and Duane C. Hassane<sup>1,5,6,8\*</sup>**

---

<sup>1</sup>Division of Hematology and Oncology, Weill Cornell Medical College, New York, NY, USA. <sup>2</sup>Health Care Policy and Research, Weill Cornell Medical College, New York, NY, USA. <sup>3</sup>Barbara Ann Karmanos Cancer Institute, Detroit, MI, USA. <sup>4</sup>Department of Pathology and Laboratory Medicine, Weill Cornell Medical College, New York, NY, USA. <sup>5</sup>Institute for Computational Biomedicine, Weill Cornell Medical College, New York, NY, USA. <sup>6</sup>Caryl and Israel Englander Institute for Precision Medicine, Weill Cornell Medical College, New York, NY, USA. <sup>7</sup>These authors contributed equally: Pinkal Desai, Nuria Mencia-Trinchant. <sup>8</sup>These authors jointly supervised this work: Gail J. Roboz, Duane C. Hassane. \*e-mail: [pid9006@med.cornell.edu](mailto:pid9006@med.cornell.edu); [dhassane@med.cornell.edu](mailto:dhassane@med.cornell.edu)

## SUPPLEMENTAL APPENDIX

### Somatic Mutations Precede Acute Myeloid Leukemia Years Before Diagnosis

Pinkal Desai, MD, MPH<sup>1†\*</sup>, Nuria Mencia-Trinchant, PhD<sup>1†</sup>, Oleksandr Savenkov, PhD<sup>2</sup>, Michael S. Simon, MD MPH<sup>3</sup>, Gloria Cheang<sup>4</sup>, Sangmin Lee, MD<sup>1</sup>, Michael Samuel<sup>1</sup>, MD, Ellen K Ritchie, MD<sup>1</sup>, Monica L. Guzman, PhD<sup>1</sup>, Karla V Ballman, PhD<sup>2</sup>, Gail J. Roboz, MD<sup>1‡</sup> and Duane C. Hassane, PhD<sup>1,5,6‡\*</sup>

#### Table of Contents

9	<b>SUPPLEMENTAL METHODS</b> .....	<b>3</b>
10	WOMEN’S HEALTH INITIATIVE COHORT .....	3
11	CASE-CONTROL MATCHING .....	3
12	TARGETED EXOME SEQUENCING .....	3
13	STATISTICAL ANALYSIS .....	3
14	STATISTICAL TOOLS .....	4
15	NEXT GENERATION SEQUENCING .....	4
16	ANALYZED GENES .....	5
17	PATHOGENICITY OF VARIANTS.....	5
18	VISUALIZATION AND ASSESSMENT OF MUTATION SPECTRUM, CO-MUTATION SPECTRUM, AND CLONAL	
19	EVOLUTION.....	6
20	DETERMINATION OF NUCLEOTIDE COMPOSITION AND CONTEXT OF MUTATIONS.....	6
21	FUNCTIONAL DOMAIN ANALYSIS.....	6
22	DATA AVAILABILITY STATEMENT .....	6
23	<b>SUPPLEMENTAL TABLES</b> .....	<b>7</b>
24	TABLE S1. SUMMARY OF CLINICAL PARAMETERS.....	7
25	TABLE S2: MUTATED GENES IN THE CASE AND CONTROL GROUP.....	8
26	TABLE S3: MUTATED CASES IN AML CASE GROUP WITH TIME TO AML LESS OR GREATER THAN 5	
27	YEARS .....	9
28	TABLE S4. EFFECT OF ALLELIC FRACTION CUTOFF AND VARIANT CLASSIFICATION ON MUTATION RATES	
29	IN CASES VS. CONTROLS.....	10
30	TABLE S5. MUTATION FREQUENCY IN AML CASES AND CONTROLS AT VARYING SNP EXCLUSION	
31	CRITERIA .....	11
32	<b>SUPPLEMENTAL FIGURES</b> .....	<b>12</b>
33	FIGURE S1. COVERAGE ACROSS GENE REGIONS FOR TOP MUTATED GENES – PART 1 .....	12
34	FIGURE S1. COVERAGE ACROSS GENE REGIONS FOR TOP MUTATED GENES – PART 2 .....	13
35	FIGURE S2. SPECTRUM AND CO-MUTATION PATTERN OF PRE-LEUKEMIC MUTATIONS – PART 1.....	14
36	FIGURE S2. SPECTRUM AND CO-MUTATION PATTERN OF PRE-LEUKEMIC MUTATIONS – PART 2.....	15
37	FIGURE S3. CO-OCCURRENCE AND MUTUAL EXCLUSIVITY PATTERNS.....	16
38	FIGURE S4. CLONAL HEMATOPOIESIS RATE IN THE CONTROL GROUP WHEN APPLYING DIFFERENT	
39	VARIANT CUTOFFS AND PATHOGENICITY CLASSIFICATION CRITERIA.....	17
40	FIGURE S5. MUTATION VARIANT ANALYSIS SUMMARY IN THE AML CASE GROUP. ....	18
41	FIGURE S6. MUTATION VARIANT ANALYSIS SUMMARY IN THE CONTROL GROUP. ....	19
42	FIGURE S7. DISTRIBUTION OF POINT MUTATION TYPES IN THE AML AND CONTROL GROUP.....	20
43	FIGURE S8. SNV SIGNATURE ANALYSIS AND TRANSCRIBED STRAND BIAS IN THE CASE AND CONTROL	
44	GROUP.....	21

1	FIGURE S9. VAF CUTOFF TABLES FOR SPECIFICITY AND SENSITIVITY .....	22
2	FIGURE S10. VARIATION IN VAF OVER TIME IN RECURRENTLY MUTATED GENES.....	23
3	FIGURE S11. PERSISTENCE OF MUTATIONS DETECTED AT BASELINE IN LONGITUDINALLY EVALUATED	
4	PARTICIPANTS. ....	24
5	FIGURE S12. ABSOLUTE RISK ESTIMATION FOR RECURRENTLY MUTATED GENES. ....	25
6	FIGURE S13. MAPPING OF CODING ALTERATIONS TO PROTEIN DOMAINS – PART 1 .....	26
7	FIGURE S13. MAPPING OF CODING ALTERATIONS TO PROTEIN DOMAINS – PART 2 .....	27
8	FIGURE S13. MAPPING OF CODING ALTERATIONS TO PROTEIN DOMAINS – PART 3 .....	28
9	FIGURE S13. MAPPING OF CODING ALTERATIONS TO PROTEIN DOMAINS – PART 4 .....	29
10	FIGURE S14. SPATIAL CLUSTERING OF TP53 POINT MUTATIONS IN 3-D STRUCTURE.....	30
11	FIGURE S15. DISTRIBUTION OF MUTATIONS IN PROTEIN DOMAINS (PFAM) IN THE AML CASE AND	
12	CONTROL GROUP. ....	31
13	FIGURE S16. MUTATION FREQUENCIES AND ASSOCIATION OF PROBABLE PATHOGENIC SOMATIC	
14	VARIANTS WITH AML. ....	32
15	FIGURE S16. MUTATION FREQUENCIES AND ASSOCIATION OF PROBABLE PATHOGENIC SOMATIC	
16	VARIANTS WITH AML – CONTINUED. ....	33
17	FIGURE S17. ODDS OF AML ADJUSTED FOR THE PRESENCE OF ADDITIONAL MUTATIONS. ....	34
18	FIGURE S18. ROC ANALYSIS OF NUMBER OF MUTATIONS IN SIGNIFICANT HIGH-RISK GENES. ....	35
19	FIGURE S19. ASSOCIATIONS BETWEEN AML AND MUTATIONS ARE ROBUST TO DIFFERENT VARIANT	
20	CLASSIFICATION METHODS AND VAF CUTOFFS. ....	36
21	FIGURE S20. ASSOCIATION BETWEEN NUMBER OF SOMATIC VARIANTS AND CLONE SIZE. ....	37
22	FIGURE S21. THE PATTERN OF SUBCLONAL SOMATIC VARIATIONS IN CASES AND CONTROLS REVEALS	
23	SIMILAR DRIVER GENES AS A DE NOVO AML COHORT.....	38
24	FIGURE S22. ODDS OF AML ARE ELEVATED WHEN MUTATIONS ARE PRESENT AT HIGHER VAF. ....	39
25	FIGURE S23. ABSENCE OF COLLINEARITY BETWEEN PREDICTORS.....	40
26	<b>SHORT LIST OF WHI INVESTIGATORS .....</b>	<b>41</b>
27	<b>REFERENCES .....</b>	<b>42</b>

28

29

## SUPPLEMENTAL METHODS

### ***Women's Health Initiative Cohort***

Study population: The WHI enrolled more than 160,000 women in one or more of three clinical trials (CT group) or an observational study cohort (OS group) in 40 U.S. clinical centers from October 1, 1993 through December 31, 1998 with data collection updated through September 2012 and an average follow up of 10.8 years (SD 3.3 years). The participants in the CT group were followed at baseline and years 1, 3 and 6, with samples collected at WHI baseline, year 1 and 3 during these follow up visits. The participants in the OS group were followed at baseline and year 3 with samples collected similarly during these visits. Available clinical parameters for this study are white blood cell count (WBC), hemoglobin, platelet count, and hematocrit.

### ***Case-control matching***

Controls were matched to cases by age at baseline within 2 years, WHI component (clinical trial vs. observational study), history of non-myeloid cancers at baseline as well as the exact timing of blood draw and follow up. In addition to these criteria, controls were also matched by type and timing of any cancers that occurred in cases after WHI baseline, but before the diagnosis of AML. For example, if a clinical trial participant with AML had prior breast cancer and was age 60 at baseline, a corresponding 60-year-old control participant with prior breast cancer history was selected who never developed subsequent AML and had exactly the same timing of follow up visits and blood draws in follow up. Similarly, if a participant with AML (case) had breast cancer diagnosed after her year 1 PB draw, she was matched to a control who also had breast cancer diagnosed after her year 1 blood draw, but did not develop AML. Matching was done in a time forward manner to ensure that each control had as much control time as its matched case. For example, a participant who developed AML five years after randomization on the WHI protocol would be matched with a control with at least five years of follow-up.

### ***Targeted Exome Sequencing***

Genomic DNA was provided by WHI in a blinded manner, in which case-control status and clinical covariates were revealed only after variant calling was completed. Library generation and amplification were performed using a low error rate Hi-Fi DNA polymerase according to the Kapa HyperPrep protocol (Kapa Biosystems). Dual sample indexing, rather than single indexing, of libraries was performed to minimize signal spread errors arising from misidentification of multiplexed samples<sup>35</sup>. Targeted sequencing using a panel of 68 recurrently mutated genes in hematological malignancies was performed using a custom capture probes (Nimblegen) to a median coverage of 2000x for both AML cases and controls (Figure S1). Variant analysis was performed following rigorous quality control and filtration of low quality sequence information. To identify somatic variants, filtration based on population allele frequency data was applied so as to enrich somatic variants that are not likely inherited. To this end, variants were classified as probable somatic if exhibiting a dbSNP v142 or ExAC adjusted population allele frequency  $\leq 0.25\%$  or a median VAF in the cohort  $< 40\%$ . Only mutations present at  $> 1\%$  VAF were evaluated for association with AML development and time-to-AML.

### ***Statistical Analysis***

Baseline characteristics of AML cases and matched controls were compared with the use of the two-sample t-test for continuous variables and the Fisher's exact test for categorical variables. Among the 188 cases, participants with baseline precursor mutations were compared to participants without precursor mutations with regard to demographic characteristics and baseline hematological characteristics (i.e. WBC count and differential counts, hemoglobin value and platelet count). The

1 relationship between specific precursor mutations and AML development were estimated by exact odds  
2 ratios (OR) and adjusted ORs were obtained from penalized-likelihood logistic regression<sup>36</sup>. Age was  
3 added as a continuous variable in all analysis and OR values were adjusted accordingly. OR and adjusted  
4 ORs are presented with their associated 95% confidence interval (95% CI). Multivariable penalized-  
5 likelihood logistic regression analysis was performed to assess the independent effect of demographic  
6 and prognostic factors of interest on precursor mutation status. Collinearity between predictors in the  
7 models was evaluated prior to the formulation of the final multivariable models. Time to development  
8 of AML was estimated with a Kaplan-Meier estimator. Differences between groups based on mutational  
9 status were evaluated with a log-rank test. Significant differences in variant allele fraction were  
10 determined in serial sampling using Fisher's Exact Test based on the count of supporting alternate and  
11 reference reads for each sample at a mutated site. For the absolute risk estimates we used a weighted  
12 partial likelihood approach and weighted baseline hazard<sup>37,38</sup>. We used sampling weights to adjust the  
13 contribution from controls, since the ratio of cases to controls is much higher in a case-control study than  
14 in the general population of interest. The weight for a control is proportional to a probability of being  
15 selected from a risk set (pool of potential controls for each case). This probability depends both on  
16 censoring time and on the matching variable. All p-values were two-sided with statistical significance  
17 set a priori at the 0.05 level. Ninety-five percent confidence intervals (95 % C.I.) were calculated to  
18 assess the precision of the obtained estimates.

19

#### 20 **Statistical tools.**

21 All statistical analyses were performed with the use of R software v3.4.0<sup>39</sup>. Multivariable odds ratios and  
22 p-values were computed using Firth's Bias-Reduced Logistic Regression implemented in the *logistf*  
23 v1.22 package<sup>40</sup>. Plots were produced using *ggplot2* v2.2.1 (<https://github.com/tidyverse/ggplot2>). Data  
24 summarization and reshaping were performed using *plyr* v1.8.4 (<https://github.com/hadley/plyr>), *dplyr*  
25 v0.7.4 (<https://github.com/hadley/dplyr>) and *reshape2* v1.4.3 (<https://github.com/hadley/reshape>)

26

#### 27 **Next generation sequencing.**

28 Following targeted enrichment according to Nimblegen protocols, libraries were sequenced on the  
29 Illumina HiSeq 4000 using dual-indexed sample adapters (Integrated DNA Technologies). To reduce  
30 errors arising from misalignment, reads were trimmed of contaminating adapter sequences and low-  
31 quality bases using Trimmomatic v0.32 (trimmed when median Illumina base quality score < 20 over 6  
32 nt sliding window). To further improve sequence quality, overlapping paired end reads were merged into  
33 a single long consensus read using AdapterRemoval v2<sup>41</sup> when at least 12 bp overlap was present. The  
34 remaining high-quality reads were mapped against the 1000 genomes phase 2 human reference genome  
35 + decoy contigs (hs37d5) using BWA MEM<sup>42</sup>. Duplicate marking was performed using SamBlaster  
36 v0.1.21<sup>43</sup> and MarkDupsByStartEnd v0.2.1<sup>44</sup>. Single nucleotide variants (SNVs) and insertions/deletions  
37 (indels) were detected using VarDictJava v1.4.6<sup>45</sup> in single sample mode with indel realignment. Marked  
38 duplicates were excluded. Copy number variations (CNV) were detected using CNVkit v0.8.6<sup>46</sup>.  
39 Annotation of variants and their functional impact was performed using Variant Effect Predictor (VEP)  
40 v85<sup>47</sup> and snpEff v4.1g<sup>48</sup>. Indel representations in both the call set and annotation data were left-aligned  
41 and harmonized using vt analysis toolkit v0.5 to maximize concordance between variants and  
42 annotations<sup>49</sup>.

43

44 The resulting call set was filtered with guidance regarding artifact removal as described<sup>50</sup>. Variants were  
45 classified as probable artifacts if any of the following conditions were met: occurs within a low  
46 complexity region subject to high alignment error defined in the hs37d5 reference genome according to  
47 the mdust-LC algorithm; exhibits strand bias; present within or immediately flanked by repetitive or low  
48 entropy regions (10 bp window); < 4 supporting reads per strand except for *NPM1* exon 12 (Ensembl

1 transcript ENST00000517671.5) where any supporting reads indicative of insertions longer than 4  
2 nucleotides were considered; mean BWA MEM mapping quality of supporting reads < 45; mean  
3 Illumina base quality of variant-supporting bases < 30; exhibits read position bias toward the 5' or 3'  
4 end of reads; < 100x unique depth of high quality coverage at site of mutation; coverage depth x VAF <  
5 8; exhibits high recurrence suggestive of artifact except for mutations known to be present  $\geq 10$  times  
6 in COSMIC v74<sup>51</sup> or at least once in the TCGA AML study<sup>52</sup> when recurrence is defined by an identical  
7 mutation occurring in >20% of samples evaluated or the position is mutated in >40% of evaluated  
8 samples or identical mutation is present in >1 sequenced instance of NA12878 or the identical position  
9 is mutated in >2 sequenced instances of NA12878. Next, for each sample, we applied binomial  
10 probability filters based on the verifyBamID mix estimate that determined the lowest allelic fraction at  
11 which a somatic variant can be distinguished from sample cross-talk as described<sup>53,54</sup> and that ascertained  
12 the probability of being a heterozygous SNP based on the number of supporting alternate and reference  
13 reads controlling the false discovery rate at 1%. Subsequent filtering was performed for false variants  
14 that may arise from reads mapping to simple tandem repeats not flagged in previous steps. Variants  
15 above  $\geq 40\%$  VAF that remained unfiltered following these steps were required to be present at <40%  
16 in same sample on serial evaluation or in the COSMIC or AML TCGA database. Finally, mutations  
17 classified by VEP were considered when categorized as missense, stop gain, splice acceptor, splice  
18 donor, frameshift insertion, frameshift deletion, in-frame insertion, and in-frame deletion mutations.  
19 Variants above 40% VAF that remained after these filtration steps included 4 instances JAK2  
20 p.Val617Phe serially maintained in 2 study participants, 2 DNMT3A variants, and an isolated 18  
21 nucleotide CREBBP in-frame insertion that is present in serial evaluation of the same subject and that  
22 may be a germline polymorphism exhibiting VAF under-representation as a result of capture bias.

### 23 **Analyzed genes.**

24 Coding exons and flanking DNA +/- 5 bp were evaluated for mutations in 68 genes including *ASXL1*,  
25 *ASXL2*, *ATRX*, *BCOR*, *BCORL1*, *BRAF*, *CALR*, *CARD11*, *CBL*, *CBLB*, *CBLC*, *CD70*, *CD79B*,  
26 *CDKN2A*, *CEBPA*, *CREBBP*, *CSF3R*, *CUX1*, *DIS3*, *DNMT3A*, *EPPK1*, *ETV6*, *EZH2*, *FAM46C*,  
27 *FBXW7*, *FLT1*, *FLT3*, *GATA1*, *GATA2*, *GNAS*, *HRAS*, *IDH1*, *IDH2*, *IKZF1*, *JAK1*, *JAK2*, *JAK3*,  
28 *KDM6A*, *KIT*, *KRAS*, *MPL*, *MYD88*, *NOTCH1*, *NPM1*, *NRAS*, *PAX5*, *PDGFRA*, *PHF6*, *PTEN*,  
29 *PTPN11*, *RAD21*, *RUNX1*, *SETBP1*, *SF3B1*, *SMC1A*, *SMC3*, *SRSF2*, *STAG1*, *STAG2*, *STAT6*, *TET1*,  
30 *TET2*, *TNF*, *TNFRSF14*, *TP53*, *U2AF1*, *WT1*, *ZRSR2*. Depth of coverage statistics were determined  
31 using Picard tools v2.6.0<sup>55</sup>.

### 32 **Pathogenicity of variants.**

33 Detected variants were classified as “probable pathogenic” if any of the following conditions were met:

- 34 • Known hotspot in genes of known pathogenic significance in myeloid malignancies including  
35 *DNMT3A*, *IDH1*, *IDH2*, *SRSF2*, *SF3B1*, *TP53*, *JAK2*, *FLT3*, *NRAS*, *KRAS*, *KIT*, *MPL*, *CBL*.
- 36 • Disruptive mutations that produce frameshifts, duplications, stop codons, or affect splice acceptor or  
37 donor sites in genes where such alterations are known to be of pathogenic impact in myeloid  
38 malignancies including *DNMT3A*, *TET2*, *TP53*, *NPM1* exon 12, *FLT3* exon 13/14, *RUNX1*, *ASXL1*  
39 exon 12/13, *CALR* exon 9, *CEBPA*, *ATRX*, *CBL*, *EZH2*, *BCOR*, *CREBBP*, *KDM6A*, *NOTCH1*,  
40 *RAD21*, *PHF6*, *CUX1*.
- 41 • Mutations demonstrating  $\geq 10$  instances in COSMIC v74 or previously identified in the AML TCGA  
42 study.
- 43 • SNVs classified as pathogenic using the Mendelian Clinically Applicable Pathogenicity (M-CAP)  
44 score classifier using the suggested cutoff of including variants scoring  $>0.025$ .<sup>56,57</sup>

1 **Visualization and assessment of mutation spectrum, co-mutation spectrum, and clonal evolution.**  
2 Variants and clinical annotations were visualized as an OncoPrint using the *ComplexHeatmap v1.12.0*  
3 package<sup>58</sup>. Chord diagrams indicating co-mutational patterns of myeloid malignancy genes were  
4 producing using the *circlize v0.4.3* package<sup>59</sup>. Mutual exclusivity of and co-occurrence of mutated genes  
5 were determined using the *maftools v1.4.25* R package.  
6

7 **Determination of nucleotide composition and context of mutations.**  
8 Analysis of transitions, transversions, and relative proportion of mutation types was performed in  
9 R/BioConductor using the *maftools v1.4.25* package. Relative frequencies of transitions and  
10 transversions within trinucleotide and transcribed strand context was performed using the  
11 R/BioConductor *MutationalPatterns v1.0* package<sup>60</sup> using transcript strand notations derived from  
12 known transcript data present in the BioConductor *TxDb.Hsapiens.ucsc hg19.knowngene* database with  
13 hg19 coordinates converted to hs37d5.  
14

15 **Functional domain analysis.**  
16 To determine the presence of mutations within known functional domains, individual mutations were  
17 mapped to Pfam functional domains based on the gene and protein-level HGVS description of each  
18 alteration (e.g. *DNMT3A* p.Arg882Cys) with the *maftools v1.4.25* package. Quantification and plotting  
19 of most frequently altered Pfam functional domains and the number of involved genes was performed  
20 using *maftools v1.4.25* package. 3-D special clusters were identified using mutation3D<sup>61</sup>. 3-D molecule  
21 visualizations were generated using PyMol 2.0 (Schrödinger, LLC; New York, NY).  
22

23 **Data availability statement**

24 The datasets generated/or analyzed during the current study are available from the corresponding  
25 author on reasonable request.  
26

1 SUPPLEMENTAL TABLES

2

3 **Table S1. Summary of clinical parameters**

4 Summary of clinical parameters for the entire cohort, AML cases or control participants (mutated, non-  
 5 mutated and in total). Median and range values are provided for each variable, time to AML or last  
 6 follow up (years), age, hematocrit, white blood cell (WBC), hemoglobin and platelet counts.

7

8

	All participants (n =369)			AML cases (n =188)			Controls (n =181)		
	Mutated n = 185	Non- mutated n = 184	Total n = 369	Mutated n = 129	Non- mutated n = 59	Total n = 188	Mutated n = 56	Non- mutated n =125	Total n = 181
<b>Years to AML (cases) / follow up (controls)</b>									
median	11.1	16.9	13.9	8.2	11.9	9.6	17.2	17.9	17.9
range	0.5-20.9	0.4-21	0.4-21	0.5-19.2	0.4-18.7	0.4-19.2	5-20.9	5.9-21	5-21
<b>Age (years)</b>									
median	68.1	64.4	66.1	67.6	62.9	66.2	69.2	65.3	66.1
range	51.5-79.7	51.4-79.7	51.4-79.7	51.5-79.7	52.6-78.4	51.5-79.7	52-79.4	51.4-79.7	51.4-79.7
<b>Hematocrit (%)</b>									
median	39.8	40	39.9	40	40	40	39.6	40	39.8
range	31.8-50.2	33.2-46.7	31.8-50.2	31.8-50.2	33.2-44.5	31.8-50.2	33.2-45.5	33.2-46.7	33.2-46.7
<b>WBC (%)</b>									
median	5.6	5.8	5.7	5.57	5.7	5.7	5.68	5.8	5.7
range	1.8-17.9	2.5-13	1.8-17.9	1.8-17.9	2.5-9.9	1.8-17.9	3.5-8.7	3.2-13	3.2-13
<b>Hemoglobin (g/dL)</b>									
median	13.5	13.5	13.5	13.5	13.3	13.5	13.3	13.6	13.5
range	10.4-16.4	10.8-16	10.4-16.4	10.4-16.4	10.8-15	10.4-16.4	10.8-15.3	11.7-16	10.8-16
<b>Platelet (1x10<sup>9</sup>/L)</b>									
median	242	239.5	240	234	242	241.5	249	235	239
range	38-874	73-410	38-874	38-874	101-387	38-874	140-438	73-410	73-438

9

10

11



1 **Table S2: Mutated genes in the case and control group**

2  
3 Number of participants with mutations in selected genes within the AML case or control groups overall  
4 and for patients < 65 or ≥ 65 years of age.  
5  
6

Gene		AML cases (N = 188) no. of participants (% of group)		Controls (N = 181) no. of participants (% of group)	
		Mutated	Non-mutated	Mutated	Non-mutated
<b>DNMT3A</b>	< 65	23 (28.75)	57 (71.25)	9 (11.69)	68 (88.31)
	≥ 65	46 (42.59)	62 (57.41)	25 (24.04)	79 (75.96)
	Total	69 (36.7)	119 (63.3)	34 (18.78)	147 (81.22)
<b>TET2</b>	< 65	8 (10)	72 (90)	2 (2.6)	75 (97.4)
	≥ 65	39 (36.11)	69 (63.89)	8 (7.69)	96 (92.31)
	Total	47 (25)	141 (75)	10 (5.52)	171 (94.48)
<b>TP53</b>	< 65	6 (7.5)	74 (92.5)	0 (0)	77 (100)
	≥ 65	15 (13.89)	93 (86.11)	0 (0)	104 (100)
	Total	21 (11.17)	167 (88.83)	0 (0)	181 (100)
<b>SRSF2</b>	< 65	2 (2.5)	78 (97.5)	0 (0)	77 (100)
	≥ 65	11 (10.19)	97 (89.81)	0 (0)	104 (100)
	Total	13 (6.91)	175 (93.09)	0 (0)	181 (100)
<b>IDH2</b>	< 65	2 (2.5)	78 (97.5)	0 (0)	77 (100)
	≥ 65	10 (9.26)	98 (90.74)	0 (0)	104 (100)
	Total	12 (6.38)	176 (93.62)	0 (0)	181 (100)
<b>JAK2</b>	< 65	5 (6.25)	75 (93.75)	0 (0)	77 (100)
	≥ 65	5 (4.63)	103 (95.37)	1 (0.96)	103 (99.04)
	Total	10 (5.32)	178 (94.68)	1 (0.55)	180 (99.45)
<b>SF3B1</b>	< 65	2 (2.5)	78 (97.5)	2 (2.6)	75 (97.4)
	≥ 65	9 (8.33)	99 (91.67)	0 (0)	104 (100)
	Total	11 (5.85)	177 (94.15)	2 (1.1)	179 (98.9)
<b>ASXL1</b>	< 65	1 (1.25)	79 (98.75)	2 (2.6)	75 (97.4)
	≥ 65	5 (4.63)	103 (95.37)	4 (3.85)	100 (96.15)
	Total	6 (3.19)	182 (96.81)	6 (3.31)	175 (96.69)
<b>U2AF1</b>	< 65	1 (1.25)	79 (98.75)	0 (0)	77 (100)
	≥ 65	5 (4.63)	103 (95.37)	0 (0)	104 (100)
	Total	6 (3.19)	182 (96.81)	0 (0)	181 (100)
<b>RUNX1</b>	< 65	1 (1.25)	79 (98.75)	0 (0)	77 (100)
	≥ 65	2 (1.85)	106 (98.15)	0 (0)	104 (100)
	Total	3 (1.6)	185 (98.4)	0 (0)	181 (100)
<b>IDH1</b>	< 65	1 (1.25)	79 (98.75)	0 (0)	77 (100)
	≥ 65	2 (1.85)	106 (98.15)	0 (0)	104 (100)
	Total	3 (1.6)	185 (98.4)	0 (0)	181 (100)

7  
8  
9

1 **Table S3: Mutated cases in AML case group with time to AML less or greater than 5 years**  
 2  
 3 Number of participants with mutations in selected genes within the AML case group with <5 or ≥ 5 years  
 4 to the diagnosis of AML  
 5  
 6

Gene	Time to AML < 5 years (n = 48)	Time to AML ≥ 5 years (n = 140)
	Mutated cases (% of group)	Mutated cases (% of group)
<b>TP53</b>	11 (22.92)	10 (7.14)
<b>RUNX1</b>	3 (6.25)	0 (0)
<b>DNMT3A</b>	22 (45.83)	47 (33.57)
<b>TET2</b>	16 (33.33)	31 (22.14)
<b>ASXL1</b>	2 (4.17)	4 (2.86)
<b>SRSF2</b>	5 (10.42)	8 (5.71)
<b>IDH2</b>	4 (8.33)	8 (5.71)
<b>JAK2</b>	2 (4.17)	8 (5.71)
<b>SF3B1</b>	5 (10.42)	6 (4.29)
<b>U2AF1</b>	2 (4.17)	4 (2.86)
<b>IDH1</b>	1 (2.08)	2 (1.43)

7  
 8  
 9

**Table S4. Effect of allelic fraction cutoff and variant classification on mutation rates in cases vs. controls.**

Retained mutations are compared using different allelic fraction cutoffs and pathogenicity classification criteria. Approaches compared include (top to bottom on table below): **(i)** variants exhibiting >1% VAF are retained after filtering is performed as described in Supplemental Methods; **(ii)** SNVs present at >3.5% VAF and indels present >7% VAF are retained after filtering is performed as described in Supplemental Methods; **(iii)** same as (i) except only probable pathogenic mutations are retained; **(iv)** same as (ii) except only probable pathogenic mutations are retained; **(v)** variants present >1% VAF are retained when satisfying inclusion criteria described by Jaiswal et al 2014<sup>62</sup>; **(vi)** same as (v) using the same SNV and indel VAF cutoffs of >3.5% and >7% described by Jaiswal et al<sup>62</sup>.

VAF cutoffs have a greater impact on clonal hematopoiesis (CH) rate in the control group than the variant classification criteria.

		AML cases (N=188)		Controls (N=181)		Odds Ratio	
		# Mutated (%)	# Non-mutated (%)	# Mutated (%)	# Non-mutated (%)	OR (95% CI)	P value
Desai classification, VAF cutoff 0.01	< 65	43 (53.75)	37 (46.25)	16 (20.78)	61 (79.22)	4.39 (2.08-9.61)	3.1 x 10 <sup>-5</sup>
	≥ 65	86 (79.63)	22 (20.37)	40 (38.46)	64 (61.54)	6.19 (3.25-12.14)	1.0 x 10 <sup>-9</sup>
	Total	129 (68.62)	59 (31.38)	56 (30.94)	125 (69.06)	4.86 (3.07-7.77)	3.8 x 10 <sup>-13</sup>
Desai classification, VAF cutoff 0.035 (SNV) or 0.07 (indels)	< 65	25 (31.25)	55 (68.75)	5 (6.49)	72 (93.51)	6.47 (2.24-23.05)	7.9 x 10 <sup>-5</sup>
	≥ 65	60 (55.56)	48 (44.44)	13 (12.5)	91 (87.5)	8.65 (4.2-18.98)	2.5 x 10 <sup>-11</sup>
	Total	85 (45.21)	103 (54.79)	18 (9.94)	163 (90.06)	7.43 (4.14-13.93)	1.1 x 10 <sup>-14</sup>
Desai classification, only Pathogenic variants, VAF cutoff 0.01	< 65	42 (52.5)	38 (47.5)	15 (19.48)	62 (80.52)	4.52 (2.12-10.05)	2.7 x 10 <sup>-5</sup>
	≥ 65	85 (78.7)	23 (21.3)	38 (36.54)	66 (63.46)	6.35 (3.35-12.4)	2.5 x 10 <sup>-11</sup>
	Total	127 (67.55)	61 (32.45)	53 (29.28)	128 (70.72)	5 (3.16-8.02)	1.1 x 10 <sup>-14</sup>
Desai classification, only Pathogenic variants, VAF cutoff 0.035 (SNV) or 0.07 (indels)	< 65	25 (31.25)	55 (68.75)	5 (6.49)	72 (93.51)	6.47 (2.24-23.05)	7.9 x 10 <sup>-5</sup>
	≥ 65	60 (55.56)	48 (44.44)	13 (12.5)	91 (87.5)	8.65 (4.2-18.98)	4.9 x 10 <sup>-10</sup>
	Total	85 (45.21)	103 (54.79)	18 (9.94)	163 (90.06)	7.43 (4.14-13.93)	1.6 x 10 <sup>-13</sup>
Jaiswal et al., 2014 classification, VAF cutoff 0.01	< 65	39 (48.75)	41 (51.25)	11 (14.29)	66 (85.71)	5.64 (2.5-13.65)	3.2 x 10 <sup>-6</sup>
	≥ 65	78 (72.22)	30 (27.78)	29 (27.88)	75 (72.12)	6.65 (3.54-12.84)	1.3 x 10 <sup>-10</sup>
	Total	117 (62.23)	71 (37.77)	40 (22.1)	141 (77.9)	5.78 (3.59-9.45)	3.3 x 10 <sup>-15</sup>
Jaiswal et al., 2014 classification, VAF cutoff 0.035 (SNV) or 0.07 (indels)	< 65	23 (28.75)	57 (71.25)	2 (2.6)	75 (97.4)	14.92 (3.45-136.04)	4.4 x 10 <sup>-6</sup>
	≥ 65	50 (46.3)	58 (53.7)	10 (9.62)	94 (90.38)	8.02 (3.67-19.16)	2.5 x 10 <sup>-9</sup>
	Total	73 (38.83)	115 (61.17)	12 (6.63)	169 (93.37)	8.89 (4.54-18.83)	3.9 x 10 <sup>-14</sup>

1 **Table S5. Mutation frequency in AML cases and controls at varying SNP exclusion criteria**  
 2 Comparison of the effect of varying population allele frequency cutoffs on the rate of mutations in cases  
 3 and controls. (a) Population frequency cutoffs of 0.001%, 0.01%, 0.25% (applied in this study) and and  
 4 1% (recommendation of Association for Molecular Pathology (AMP) / American Society for Clinical  
 5 Oncology (ASCO) / College of American Pathologists (CAP) for selecting somatic mutations)<sup>63</sup> (b)  
 6 Forest plot indicating odds ratio of mutations in cases vs. controls using a population allele frequency  
 7 cutoff 0.001%. Genes or gene categories significantly associated with AML include *TP53* ( $P = 3.0 \times 10^{-6}$ ),  
 8 *IDH* ( $P = 3.0 \times 10^{-4}$ ), spliceosome, *TET2* ( $P = 2.4 \times 10^{-6}$ ), and *DNMT3A* ( $P = 3.4 \times 10^{-4}$ ). *IDH* category  
 9 includes *IDH1* and *IDH2*. The spliceosome category includes *SRSF2*, *SF3B1*, and *U2AF1*. OR per gene  
 10 are adjusted by age (years) as a continuous variable. Abbreviations: CI, confidence interval; N, number  
 11 affected. P-values are shown for penalized likelihood multivariable logistic regression.  
 12

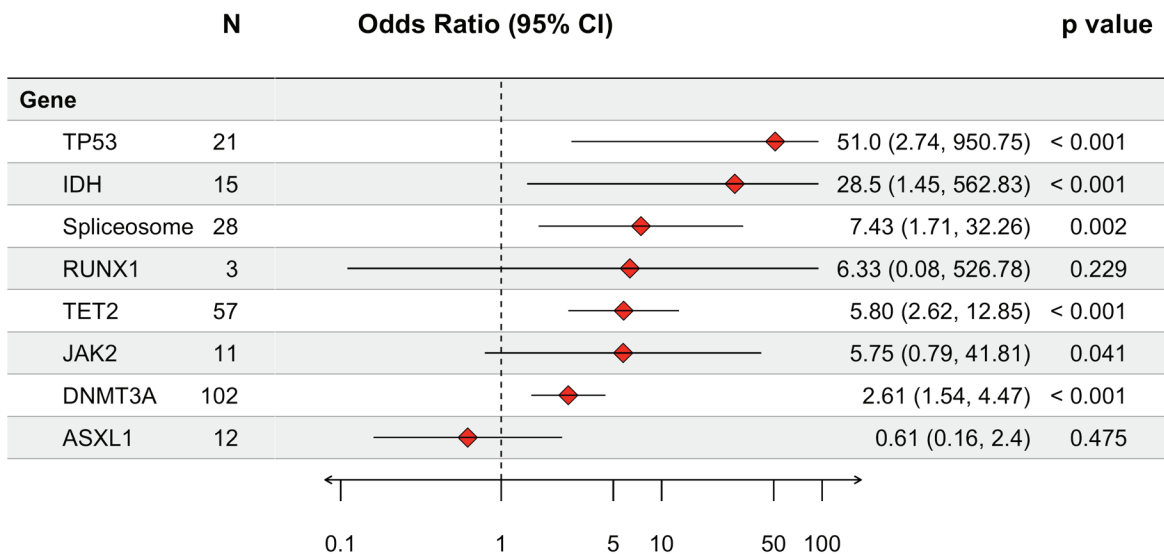
13 **a**

	AML cases (N=188)		Controls (N=181)		Odds Ratio		
	# Mutated (%)	# Non-mutated (%)	# Mutated (%)	# Non-mutated (%)	OR (95% CI)	P value	
dbSNP & ExAC AF < 0.001%	< 65	43 (53.75)	37 (46.25)	16 (20.78)	61 (79.22)	4.39 (2.08-9.61)	$3.1 \times 10^{-5}$
	≥ 65	86 (79.63)	22 (20.37)	40 (38.46)	64 (61.54)	6.19 (3.25-12.14)	$1.0 \times 10^{-9}$
	Total	129 (68.62)	59 (31.38)	56 (30.94)	125 (69.06)	4.86 (3.07-7.77)	$3.8 \times 10^{-13}$
dbSNP & ExAC AF < 0.01%	< 65	43 (53.75)	37 (46.25)	16 (20.78)	61 (79.22)	4.39 (2.08-9.61)	$3.1 \times 10^{-5}$
	≥ 65	86 (79.63)	22 (20.37)	40 (38.46)	64 (61.54)	6.19 (3.25-12.14)	$1.0 \times 10^{-9}$
	Total	129 (68.62)	59 (31.38)	56 (30.94)	125 (69.06)	4.86 (3.07-7.77)	$3.8 \times 10^{-13}$
dbSNP & ExAC AF < 0.25%	< 65	43 (53.75)	37 (46.25)	16 (20.78)	61 (79.22)	4.39 (2.08-9.61)	$3.1 \times 10^{-5}$
	≥ 65	86 (79.63)	22 (20.37)	40 (38.46)	64 (61.54)	6.19 (3.25-12.14)	$1.0 \times 10^{-9}$
	Total	129 (68.62)	59 (31.38)	56 (30.94)	125 (69.06)	4.86 (3.07-7.77)	$3.8 \times 10^{-13}$
dbSNP & ExAC AF < 1%	< 65	43 (53.75)	37 (46.25)	17 (22.08)	60 (77.92)	4.06 (1.95-8.78)	$6.9 \times 10^{-5}$
	≥ 65	86 (79.63)	22 (20.37)	41 (39.42)	63 (60.58)	5.95 (3.13-11.65)	$2.3 \times 10^{-9}$
	Total	129 (68.62)	59 (31.38)	58 (32.04)	123 (67.96)	4.62 (2.92-7.37)	$1.9 \times 10^{-12}$

14

15 **b**

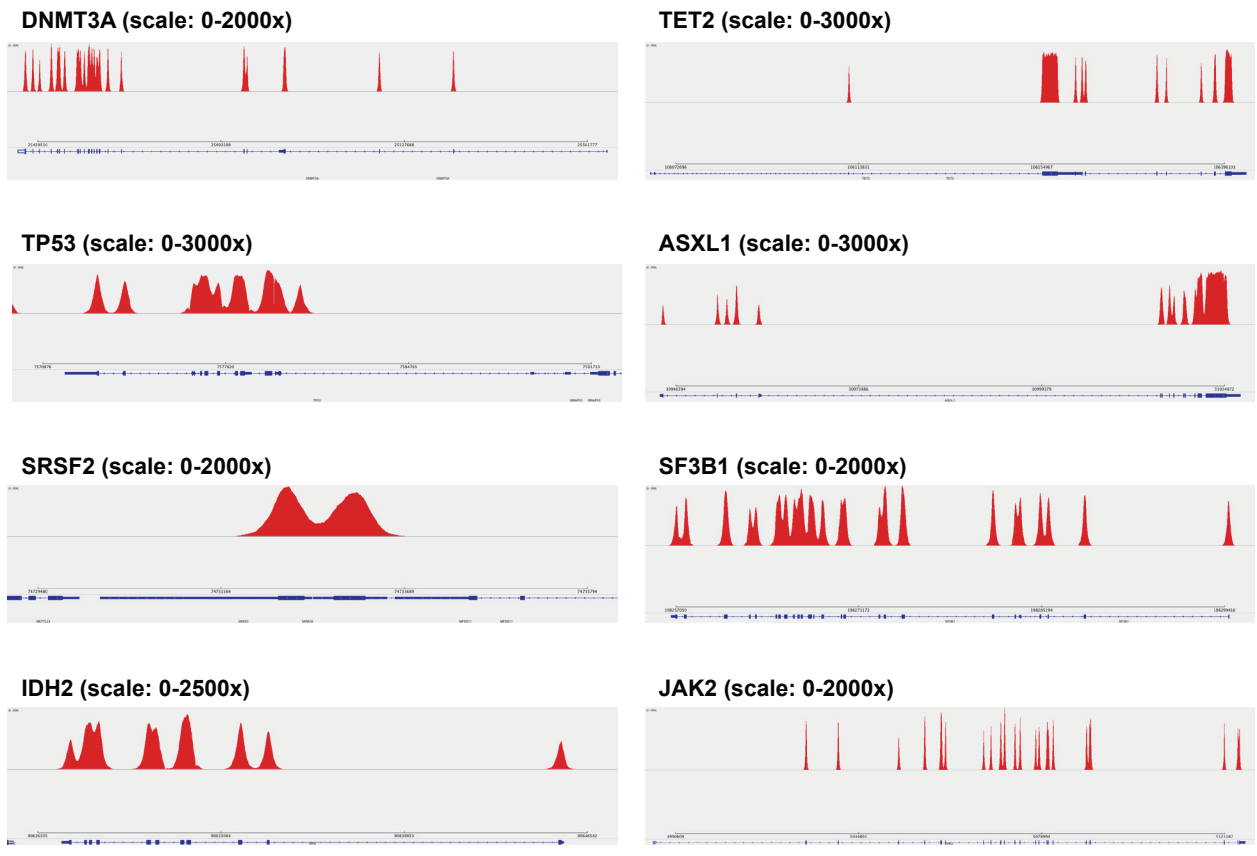
16



1 SUPPLEMENTAL FIGURES

2  
3 **Figure S1. Coverage across gene regions for top mutated genes – Part 1**

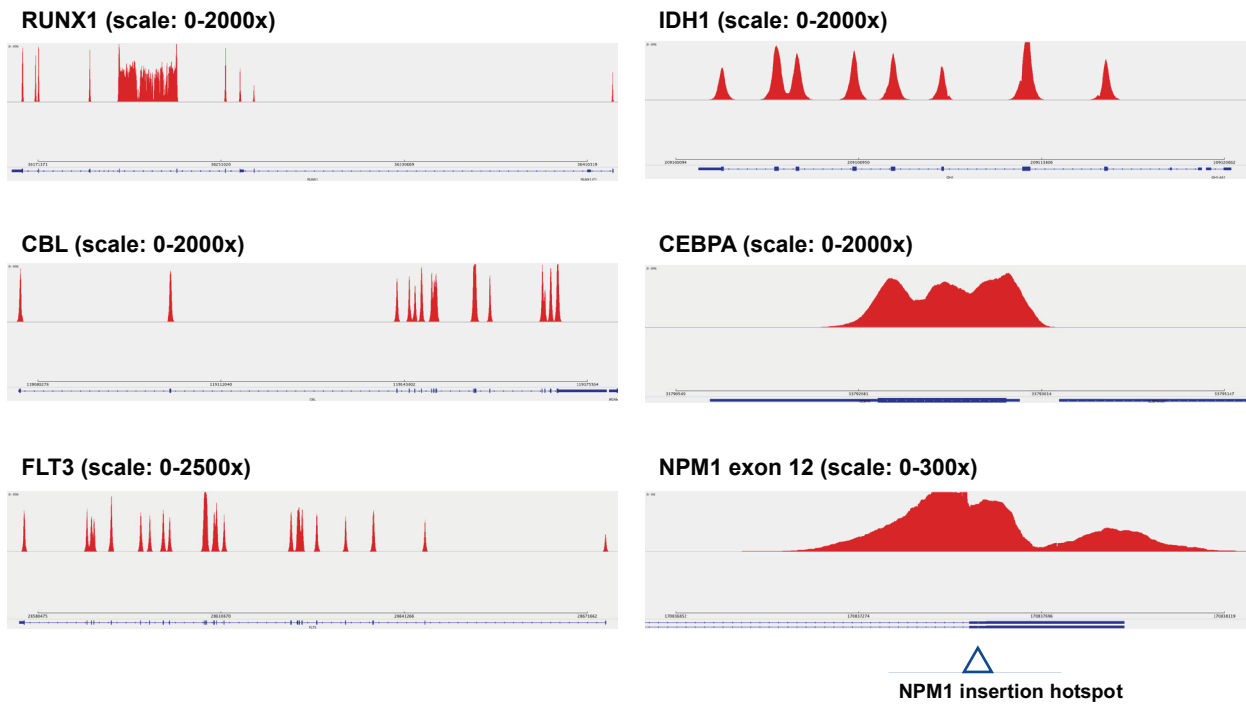
4  
5 Representative images for coverage of recurrently mutated genes in the cohort at the median depth of  
6 coverage of 2000x. For each gene, coverage across pertinent exons approaches or exceeds 2000x  
7 including important regions of driver genes such as *FLT3*. *CEBPA* also achieved >500x median coverage  
8 across its single coding exon. Median *NPM1* exon 12 coverage was relatively lower (~280x) potentially  
9 resulting in more false negatives. However, given the zero background rate of 4 nucleotide insertions in  
10 the *NPM1* insertion hotspot, this lower coverage maintains >80% power to detect *NPM1* insertions to a  
11 VAF > 1% (1-sample binomial power calculation; background mutation rate of 0.1%, alpha = 0.01).



13  
14

1 **Figure S1. Coverage across gene regions for top mutated genes – Part 2**

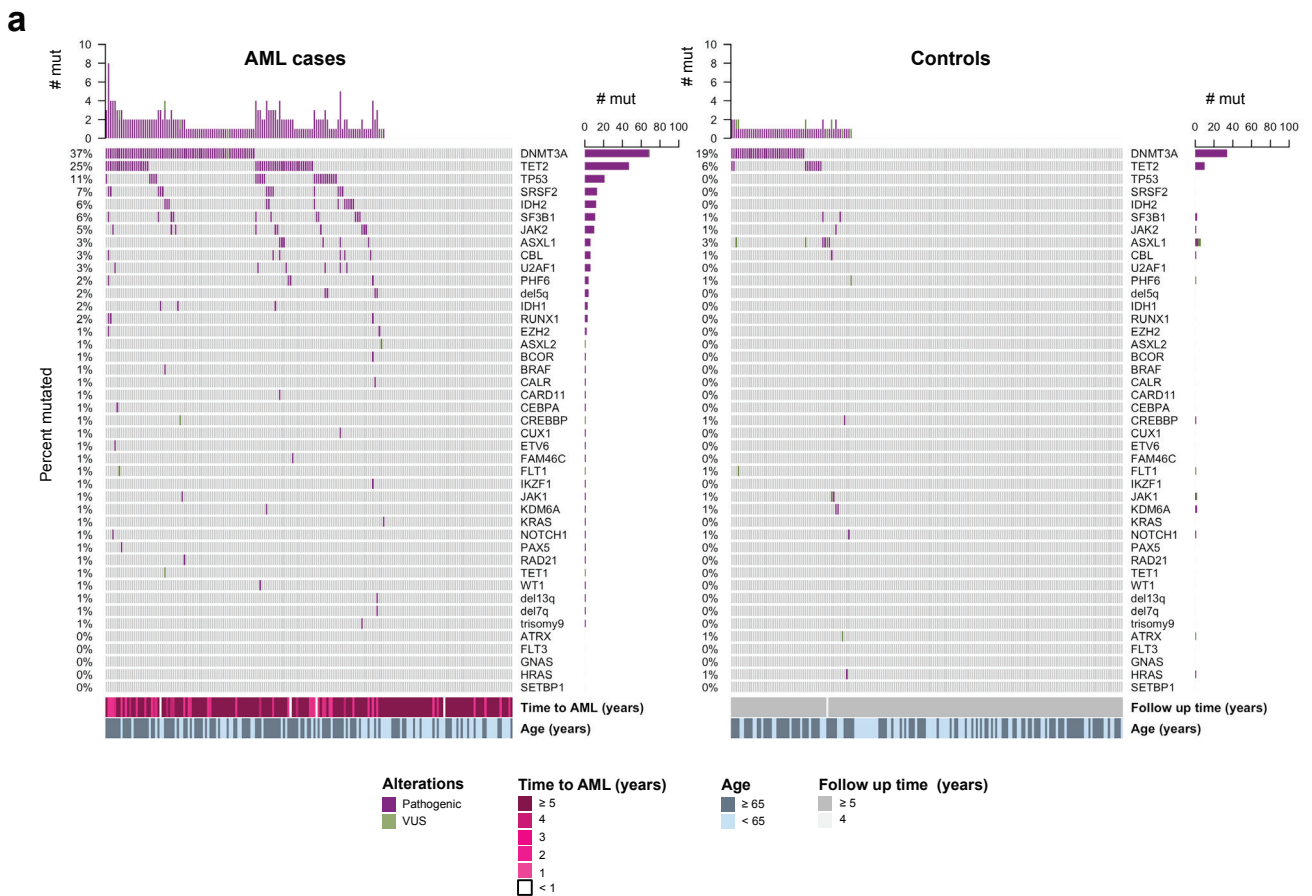
2



3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

1 **Figure S2. Spectrum and co-mutation pattern of pre-leukemic mutations – Part 1**

2  
 3 (a) OncoPrint for all participants in the study at baseline. AML cases (N = 188) are represented in the  
 4 left panel and controls (N = 181) in the right panel. Each row represents a gene and each column  
 5 corresponds to a participant in the study. Bar plots indicate the number of mutations per patient (top bar  
 6 plot), and the number of patients with mutations in each gene (side bar plot). For each patient, bottom  
 7 panels show: time to AML diagnosis (Time to AML) for the cases or last follow up for the controls  
 8 (Follow up time) and age at diagnosis (Age). Dark grey, patients older or equal than 65 years old; light  
 9 grey, patients younger than 65 years old. Alterations classified as variant of unknown significance (VUS,  
 10 green) or pathogenic (Pathogenic, magenta) according to the criteria specified in Supplemental Methods.  
 11 All analyzed genes are included.  
 12

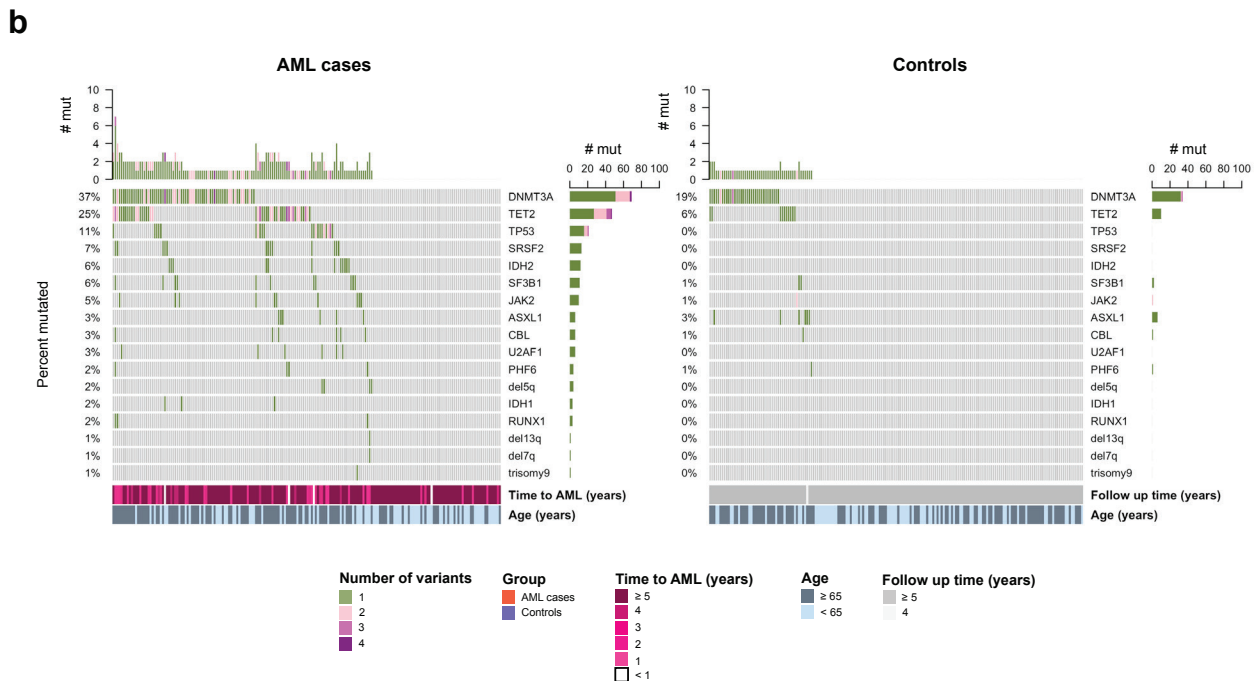


13  
 14

**Figure S2. Spectrum and co-mutation pattern of pre-leukemic mutations – Part 2**

(b) Oncoprint showing the number of variants per gene detected for each participant. Each row represents a gene and each column corresponds to a participant in the study. Bar plots indicate the number of mutations per patient (top bar plot), and the number of patients with mutations in each gene (side bar plot). For each patient, bottom panels indicate: time to AML diagnosis (Time to AML) for the cases or last follow up for the controls (Follow up time) and age at diagnosis (Age). Dark grey, patients older or equal than 65 years old; light grey, patients younger than 65 years old. Number of variants ranges from 1 to 4 (1, green; 2, light pink; 3, pink and 4, magenta).

Multiple variants per participant were mainly found in DNMT3A, TET2 and TP53 genes.

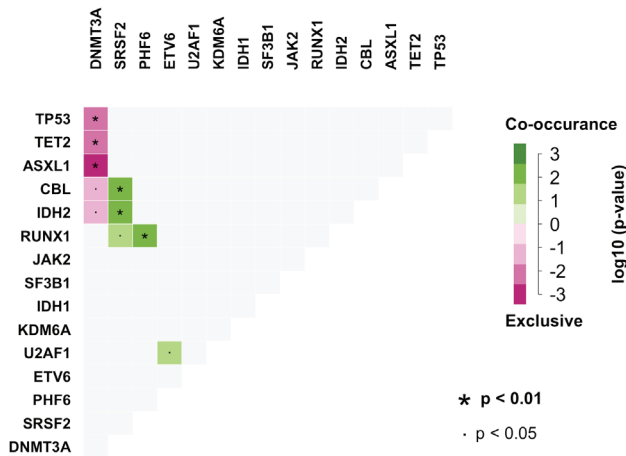


13  
14



1 **Figure S3. Co-occurrence and mutual exclusivity patterns**  
 2 Tendency of co-occurrence and mutual exclusivity in mutated genes are shown for the (a) overall cohort,  
 3 (b) AML case group, and (c) control group. (d) Table showing P values each significant gene pair  
 4 association as well as tendency toward mutually exclusivity or co-occurrence. Pairwise significance of  
 5 associations are determined using the *maftools* R package (\* p < 0.01; • p < 0.05; two-sided Fisher's  
 6 Exact Test). Mutations present >=3 times per group were considered. n = number mutated  
 7  
 8

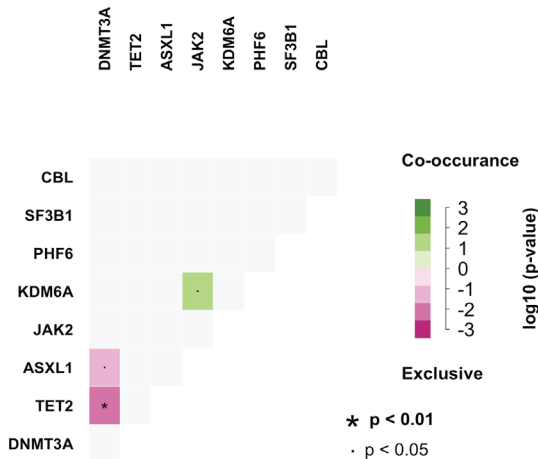
**a Overall (n = 185)**



**b AML cases (n = 129)**



**c Controls (n = 56)**



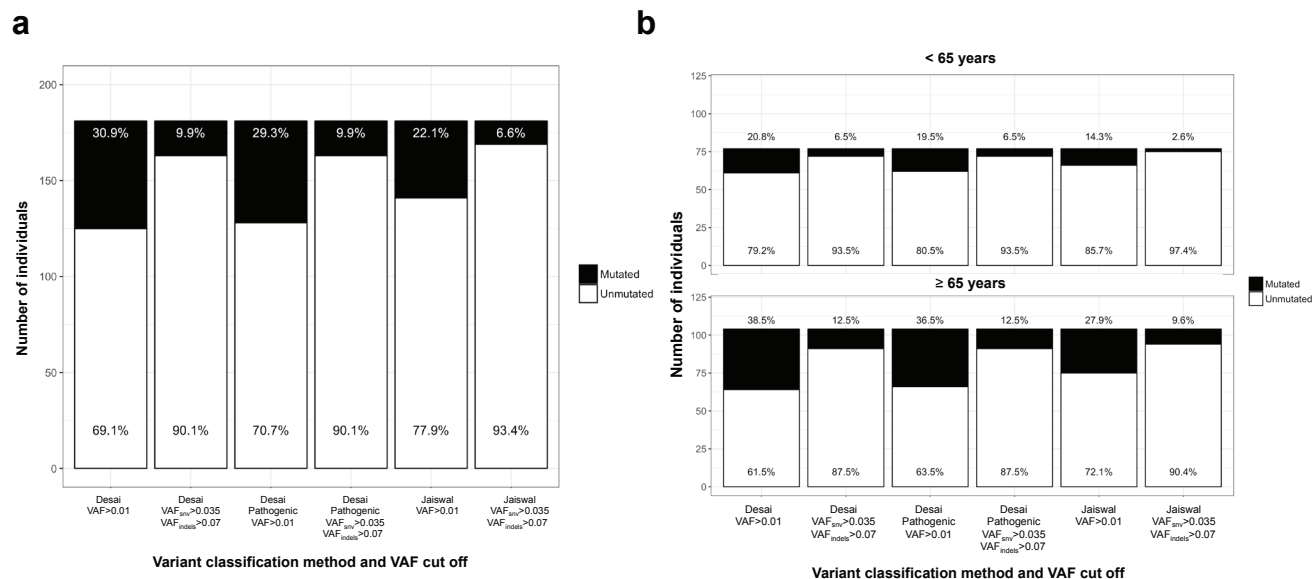
**d**

	Gene 1	Gene 2	Pvalue	Event
<b>Overall</b>	DNMT3A	ASXL1	0.0006	Mutually exclusive
	PHF6	RUNX1	0.0018	Co-occurrence
	TP53	DNMT3A	0.0021	Mutually exclusive
	TET2	DNMT3A	0.0022	Mutually exclusive
	IDH2	SRSF2	0.0055	Co-occurrence
	CBL	SRSF2	0.0083	Co-occurrence
	SRSF2	RUNX1	0.0133	Co-occurrence
	ETV6	U2AF1	0.0326	Co-occurrence
	DNMT3A	IDH2	0.0344	Mutually exclusive
	CBL	DNMT3A	0.0449	Mutually exclusive
<b>Cases</b>	RUNX1	PHF6	0.0022	Co-occurrence
	TP53	DNMT3A	0.0035	Mutually exclusive
	DNMT3A	ASXL1	0.0083	Mutually exclusive
	CBL	SRSF2	0.0139	Co-occurrence
	IDH2	SRSF2	0.0204	Co-occurrence
	CBL	ASXL1	0.0254	Co-occurrence
	RUNX1	SRSF2	0.0271	Co-occurrence
ETV6	U2AF1	0.0469	Co-occurrence	
<b>Controls</b>	TET2	DNMT3A	0.0092	Mutually exclusive
	ASXL1	DNMT3A	0.0299	Mutually exclusive
	JAK2	KDM6A	0.0357	Co-occurrence

1 **Figure S4. Clonal hematopoiesis rate in the control group when applying different variant cutoffs**  
 2 **and pathogenicity classification criteria.**

3  
 4 (a) Percent of control participants exhibiting clonal hematopoiesis when defined as (left to right): any  
 5 variant present at VAF > 1% irrespective of probable driver status (Desai, VAF > 0.01); any variant  
 6 present irrespective of probable driver status using higher VAF cutoffs of 3.5% for SNVs and 7% for  
 7 indels that were used by previous low depth whole exome sequencing studies<sup>62</sup> (Desai VAF<sub>SNV</sub> > 0.035,  
 8 VAF<sub>indel</sub> > 0.07); pathogenic variants classified as in Supplemental Methods present at >1% VAF (Desai  
 9 Pathogenic VAF > 0.01); pathogenic variants classified as in Supplemental Methods using higher VAF  
 10 cutoffs of >3.5% for SNVs and >7% for indels (Desai Pathogenic VAF<sub>SNV</sub> > 0.035, VAF<sub>indel</sub> > 0.07);  
 11 any variant present at VAF > 1% passing the pathogenicity criteria of Jaiswal and colleagues (Jaiswal  
 12 VAF > 0.01)<sup>62</sup>; any variant passing both the pathogenicity criteria and VAF cutoffs used by Jaiswal and  
 13 colleagues (Jaiswal VAF<sub>SNV</sub> > 0.035, VAF<sub>indel</sub> > 0.07)<sup>62</sup>. (b) Same classification and cutoff criteria as in  
 14 (a) are applied stratified by age group (< 65 years; top) vs. (≥ 65 years; bottom). Each bar represents the  
 15 fraction of mutated participants (black) vs. non-mutated participants (white). Percentages are shown.

16  
 17 VAF cutoffs have a greater impact on incidence of mutations than the variant classification criteria.  
 18

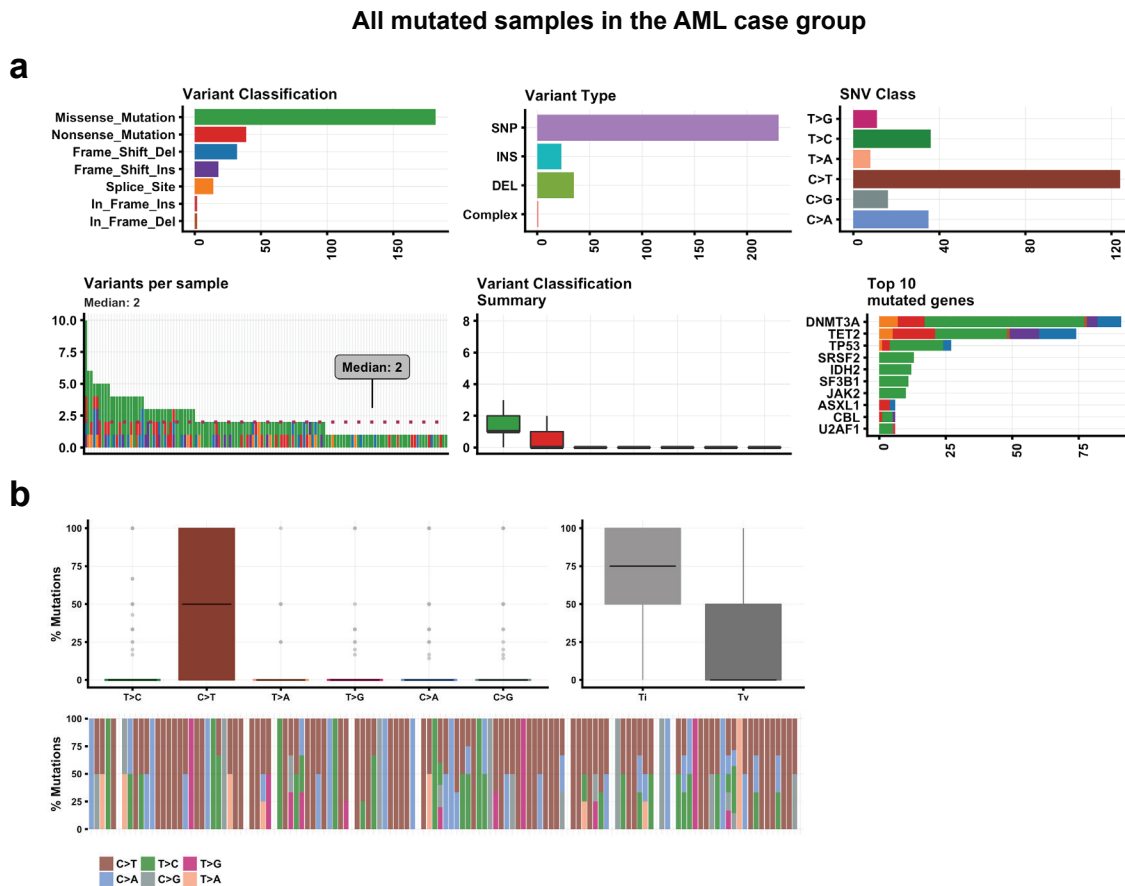


19  
 20  
 21  
 22  
 23  
 24

1 **Figure S5. Mutation variant analysis summary in the AML case group.**

2  
3 A total of 289 variants were identified in AML case group (n = 129). Missense mutations accounted for  
4 63% (181 mutations) followed by deletions (11.7%; 34 mutations), insertions (6.9%; 20 mutations), and  
5 1.7% CNVs (5 total found exclusively in AML cases, denoted as “Complex” or “Other”).

6  
7 (a) Top row: bar plots enumerate variant classification (left); variant types (middle); SNV substitution  
8 (right). SNV, single nucleotide variant; INS, insertion; DEL, deletion. Bottom row: number of variants  
9 per sample with colors indicating variant classification (left); variant classification summary;  
10 and frequently mutated genes (right). N indicates number of mutations per sample. Colors indicate variant  
11 classifications. (b) Top row: boxplots indicating overall distribution of six different types of conversions  
12 (left) or transitions (Ti) and transversions (Tv) (right). Bottom row: Stacked bar plot shows the fraction  
13 of conversions per sample (bottom panel). Box and whiskers plots: box indicates the 1<sup>st</sup> quartile, median,  
14 and 3<sup>rd</sup> quartile whereas whiskers represent 1.5x the interquartile range. Plots are generated using  
15 *maftools*.  
16



17  
18  
19

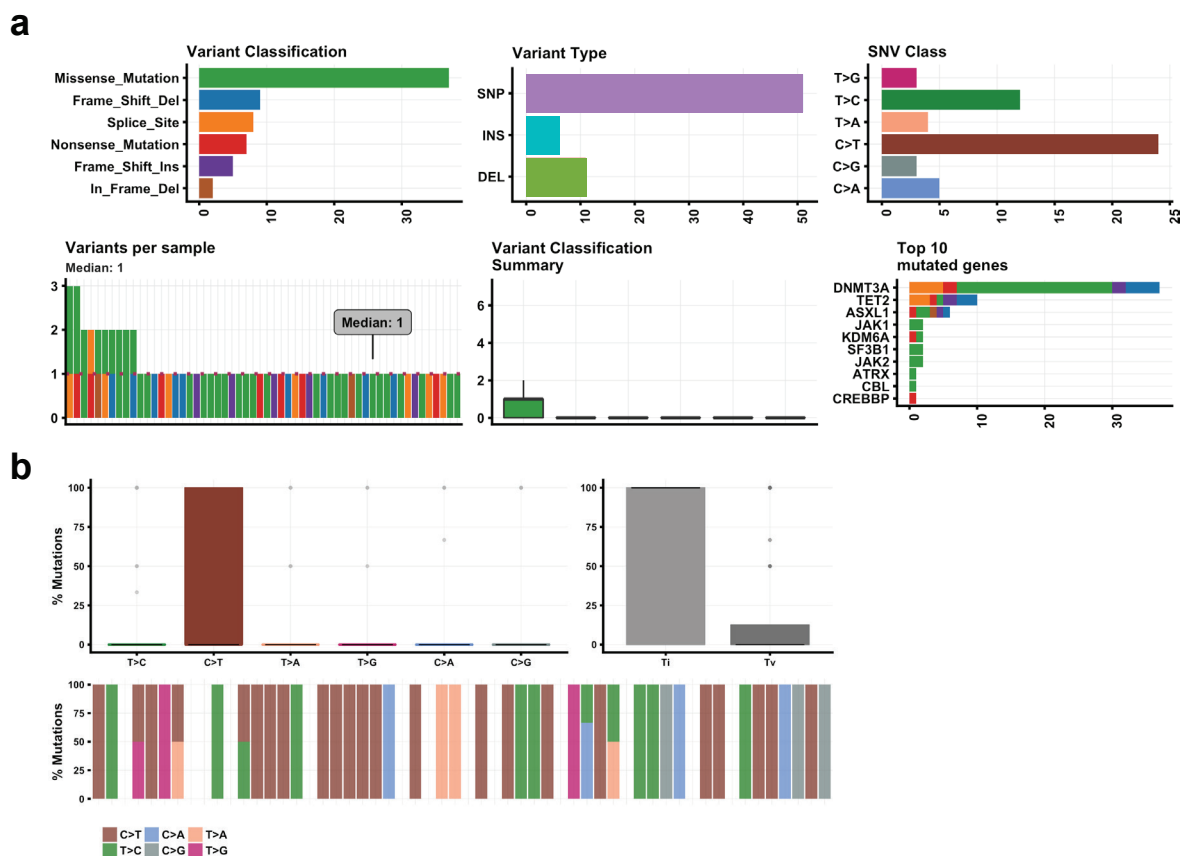
1 **Figure S6. Mutation variant analysis summary in the control group.**

2  
3 A total of 68 variants were identified in the control group (n = 56). Missense mutations accounted for  
4 54.4% of mutations (37 mutations total), followed by deletions (16.1%; 11 mutations), insertions (7.3%;  
5 5 mutations).

6  
7 As with the AML cases, the most common single-nucleotide change was a cytosine-to-thymine (C>T)  
8 transition occurring in CpG context for both groups (supplemental figures 3,4,5), a lesion and context  
9 associated with age-related mutagenesis<sup>64</sup> and consistent with other reports<sup>62</sup>.

10  
11 (a) Top row: bar plots enumerate variant classification (left); variant types (middle); SNV substitution  
12 (right). SNV, single nucleotide variant; INS, insertion; DEL, deletion. Bottom row: number of variants  
13 per sample with colors indicating variant classification (left); variant classification summary;  
14 and frequently mutated genes (right). N indicates number of mutations per sample. Colors indicate variant  
15 classifications. (b) Top row: boxplots indicating overall distribution of six different types of conversions  
16 (left) or transitions (Ti) and transversions (Tv) (right). Bottom row: Stacked bar plot shows the fraction  
17 of conversions per sample (bottom panel). Box and whiskers plots: box indicates the 1<sup>st</sup> quartile,  
18 and 3<sup>rd</sup> quartile whereas whiskers represent 1.5x the interquartile range. Plot is generated using *maftools*.

19 **All mutated samples in the control group**

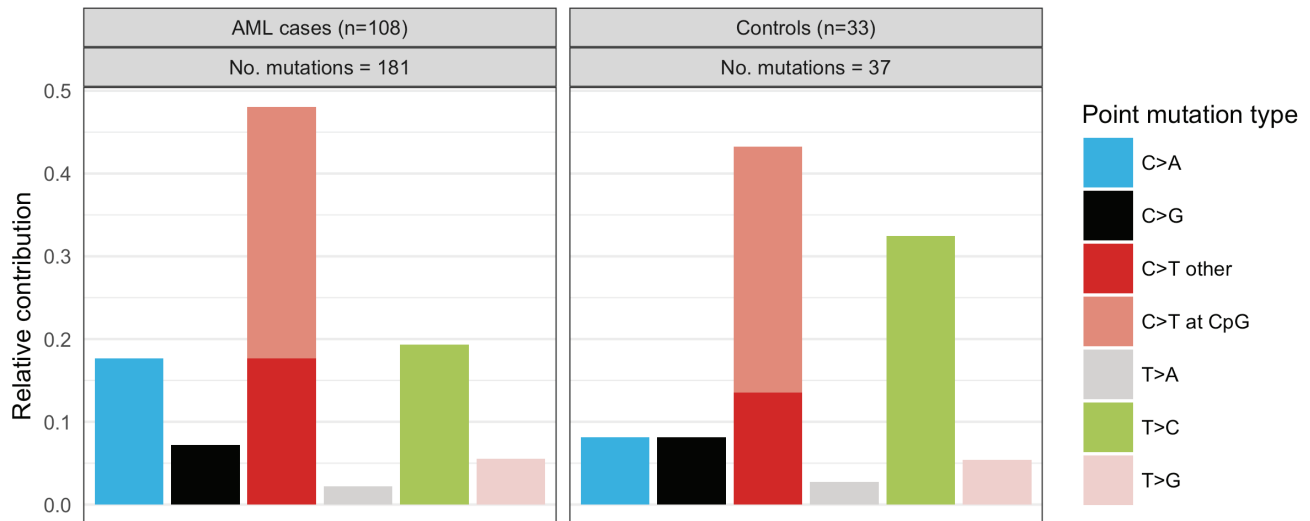


20  
21  
22  
23

1 **Figure S7. Distribution of point mutation types in the AML and control group.**

2  
3 Relative contribution of base substitutions when focusing on missense mutations enumerated in the AML  
4 case group (n = 181; Figure S1) and the control group (n = 37; Figure S2). The majority of C>T  
5 transitions occur in CpG context for both groups, suggesting acquisition from missrepair deamination of  
6 5-methylcytosine. The most common single-nucleotide change was a cytosine-to-thymine (C>T)  
7 transition occurring in CpG context for both groups, a lesion and context associated with age-related  
8 mutagenesis<sup>64</sup> and consistent with other findings<sup>62</sup>.

9

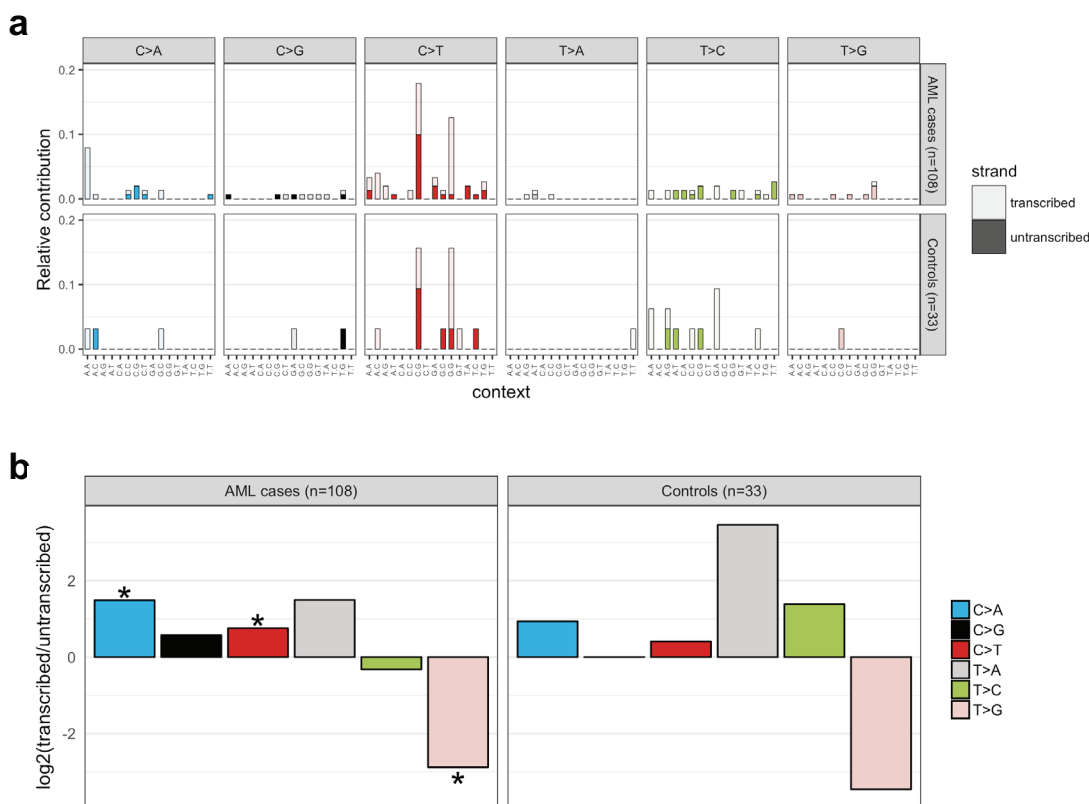


10  
11  
12

**Figure S8. SNV signature analysis and transcribed strand bias in the case and control group.**

A broader analysis of base substitution was performed taking into account the bases immediately upstream and downstream of the mutated base providing mutation context information (panel A below). An elevated rate of spontaneous deamination of 5-methyl-cytosine occurring predominantly NpCpG trinucleotide is consistent with reports of mutation patterns in AML<sup>65</sup>. Of note, a similar pattern is observed for AML cases and controls, which suggests common underlying mutation processes in the two groups most likely driven by aging<sup>66</sup>. Further analysis of the mutation pattern shows preference for certain substitutions in the transcribed strand over the untranscribed strand suggesting additional mutational processes driven by transcription-coupled repair (TCR)<sup>67</sup>, a nucleotide excision repair (NER) process that has been shown to decrease in efficiency with normal aging<sup>68</sup>. While strand bias suggesting TCR approached significance ( $P < 0.05$ ) for the AML group. The control group demonstrated the same trend toward strand bias also suggesting TCR but did not approach significance because of lower number of cases.

(a) Trinucleotide context of C>A, C>G, C>T, T>A, T>C, and T>G point mutations is shown among non-synonymous point mutations in the AML case and control group. For each context, the stacked bar chart indicates mutations occurring on the untranscribed strand (filled with color) vs the transcribed strand (not filled). (b) Transcribed vs. untranscribed strand bias (log2 scale) is indicated for each type of point mutation type in the AML case (n = 181 variants, n = 108 participants) and control group (n = 37 variants, n = 33 participants). Significant differences are indicated by asterisks (\*) using the Poisson test. \*  $P < 0.05$  as implemented in the *MutationalPatterns* R package.



24  
25  
26  
27

1 **Figure S9. VAF cutoff tables for specificity and sensitivity**

2  
3 Analysis of sensitivity and specificity showed that the false positive rate for individuals bearing  
4 mutations in *TP53*, *SRSF2*, *IDH2*, *SF3B1* or *U2AF1* in the 1-2% VAF range is less than 1%.

5  
6 Table of VAF cut offs for the significantly mutated genes producing no greater than a (a) 1% false  
7 positive rate, (b) 5% false positive rate, or (c) 10% false positive rate while maximizing sensitivity. False  
8 positives represent controls misclassified as AML cases.  
9

**a**

**< 1% false positive rate**

Gene	VAF cutoff	True positive rate	False positive rate
DNMT3A	34.09	6.2	0
TET2	20.11	13.18	0
SRSF2	1.29	10.08	0
IDH2	1.18	9.3	0
TP53	1.02	16.28	0
SF3B1	1.95	7.75	0
U2AF1	1.79	4.65	0
Any of the above genes	31.21	11.63	0
Any of the above genes except DNMT3A	20.11	22.48	0

**b**

**< 5% false positive rate**

Gene	VAF cutoff	True positive rate	False positive rate
DNMT3A	21.65	10.08	3.57
TET2	5.84	19.38	3.57
SRSF2	1.29	10.08	0
IDH2	1.18	9.3	0
TP53	1.02	16.28	0
SF3B1	1.49	8.53	1.79
U2AF1	1.79	4.65	0
Any of the above genes	20.11	27.91	3.57
Any of the above genes except DNMT3A	5.84	37.98	3.57

**c**

**< 10% false positive rate**

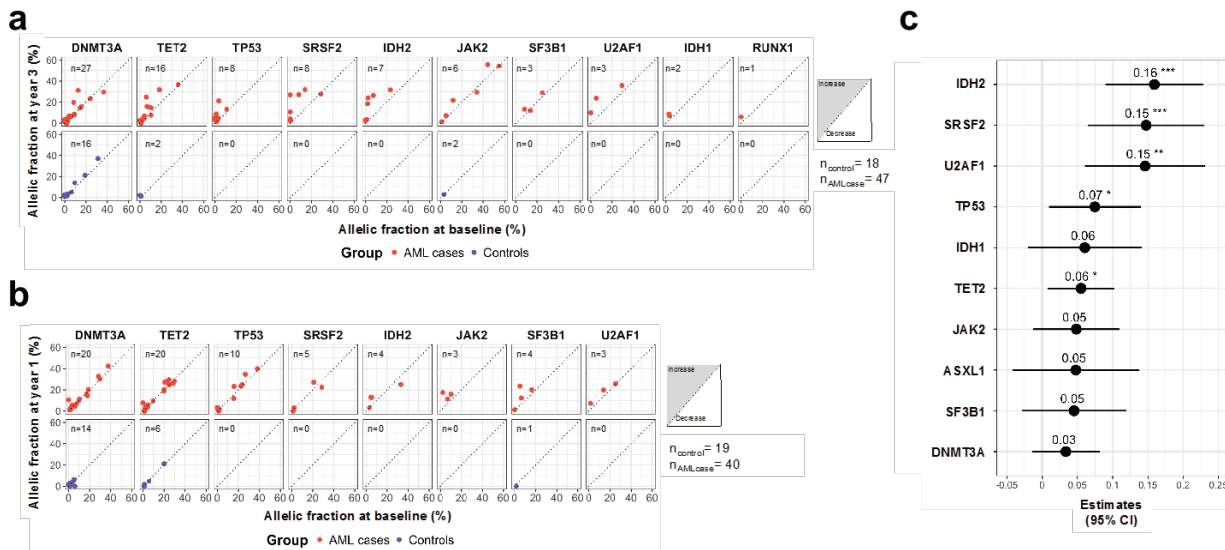
Gene	VAF cutoff	True positive rate	False positive rate
DNMT3A	8.88	22.48	8.93
TET2	1.37	31.78	8.93
SRSF2	1.29	10.08	0
IDH2	1.18	9.3	0
TP53	1.02	16.28	0
SF3B1	1.49	8.53	1.79
U2AF1	1.79	4.65	0
Any of the above genes	9.66	44.96	8.93
Any of the above genes except DNMT3A	1.47	55.04	8.93

1 **Figure S10. Variation in VAF over time in recurrently mutated genes.**

2  
3 Mutations in recurrently mutated genes in patients with serial samples were tracked over time. Generally,  
4 no changes were observed within the first year in AML cases or controls. At 3 years, however, the AML  
5 cases demonstrated elevations in VAF. In contrast, the VAF in the controls group remained mostly stable  
6 up to 3 years. Mutations in *IDH2*, *SRSF2*, *U2AF1*, *TP53* and *TET2* showed significant increase in mean  
7 VAF between baseline and year 3 follow up.

8  
9 Comparison of allelic fractions per mutated gene at baseline (horizontal axis) and (a) after 3 years of  
10 follow up (vertical axis) or (b) after 1 year of follow up (vertical axis) irrespective of whether mutation  
11 is present at baseline evaluation. The diagonal represents no change in VAF. (c) Forest plot indicating  
12 linear model estimates of mean allelic fraction changes when a mutation is present at baseline evaluation  
13 for specific genes: *IDH2* ( $P = 1.9 \times 10^{-5}$ ,  $n=5$ ); *SRSF2* ( $P = 6.4 \times 10^{-4}$ ,  $n = 3$ ); *U2AF1* ( $P = 1.0 \times 10^{-3}$ ,  $n =$   
14  $2$ ); *TP53* ( $P = 2.5 \times 10^{-2}$ ,  $n = 6$ ); *TET2* ( $P = 2.4 \times 10^{-2}$ ,  $n = 20$ ); *IDH1* ( $P = 0.140$ ,  $n = 2$ ); *JAK2* ( $P = 0.120$ ,  
15  $n = 8$ ); *ASXL1* ( $P = 0.294$ ,  $n = 2$ ); *SF3B1* ( $P = 0.227$ ,  $n = 3$ ); *DNMT3A* ( $n = 38$ ,  $P = 0.169$ ). \*  $P < 0.05$ ,  
16 \*\* $P < 0.01$ , \*\*\*  $P < 0.001$ ;  $n$ , number of pairwise evaluations.

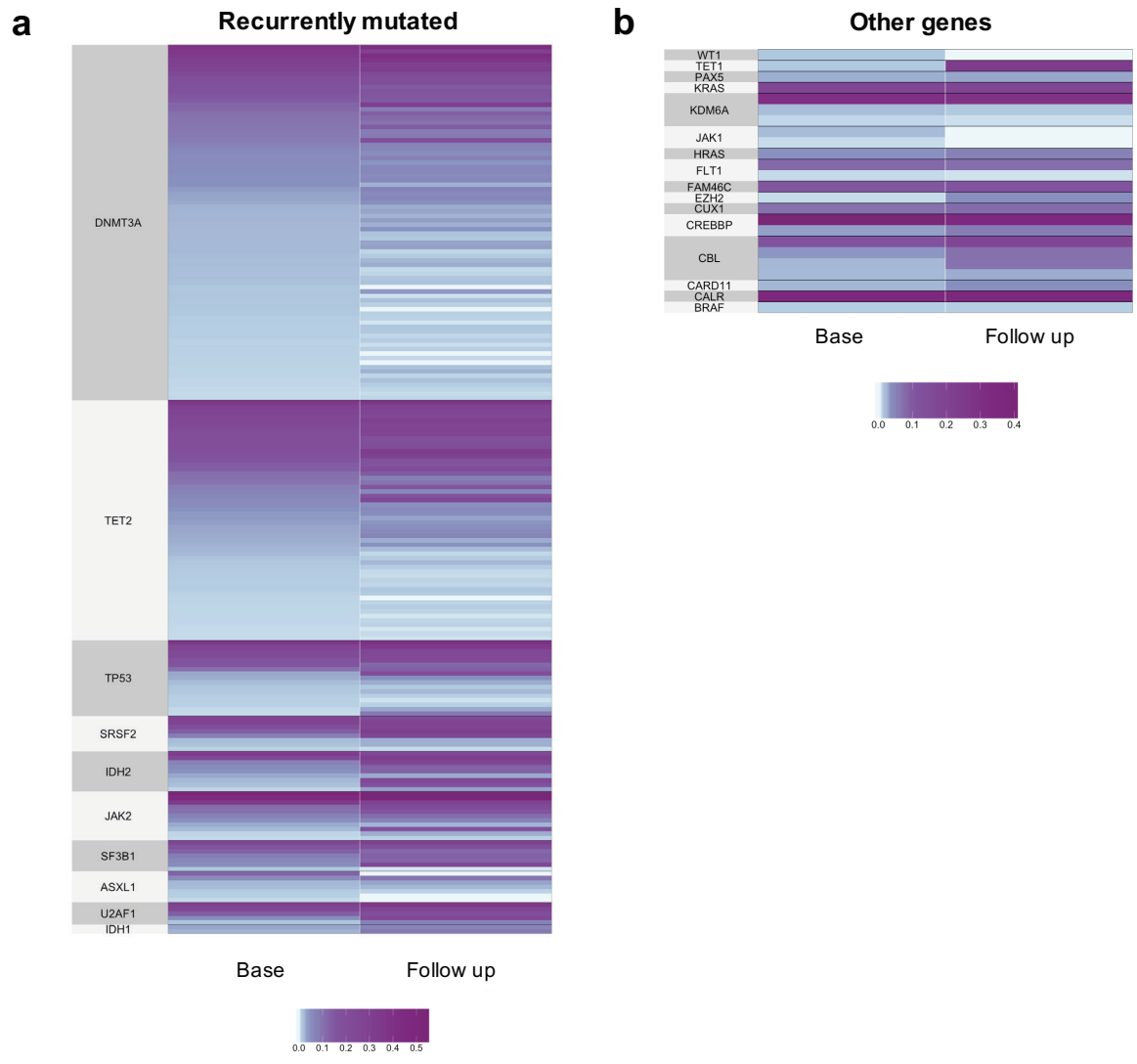
17  
18 The diagonal represents no change in VAF. Individual AML cases mutated in the gene indicated are  
19 indicated with red dot (AML case) or blue dot (control). The maximum VAF was selected when  
20 participants harbored more than 1 mutation.  
21



22



1 **Figure S11. Persistence of mutations detected at baseline in longitudinally evaluated participants.**  
 2  
 3 Heatmap indicating persistence of mutations in evaluable cases and controls with serial samples  
 4 available at baseline and 1-year or 3-year follow-up. Blue to purple color gradient in heatmap indicates  
 5 VAF ranging from 0% to 40%. >95% of variants present at baseline VAF > 1% are stably maintained at  
 6 year 1 or year 3 (N=213/224 serially evaluable variants in 121 individuals). SNP signatures were verified  
 7 in all longitudinally monitored participants and did not explain the non-persistent variants. Mutations  
 8 were force-called down to 0.5% VAF to detect persistence of variants near the 1% VAF cutoff. Non-  
 9 canonical variants in genes such as *FLT1* and *CARD11* are maintained similarly as canonical drivers  
 10 such as *IDH2*, *TET2*, *TP53*, and *SRSF2*. Variants present at lower VAF demonstrated a tendency to drop  
 11 out.  
 12

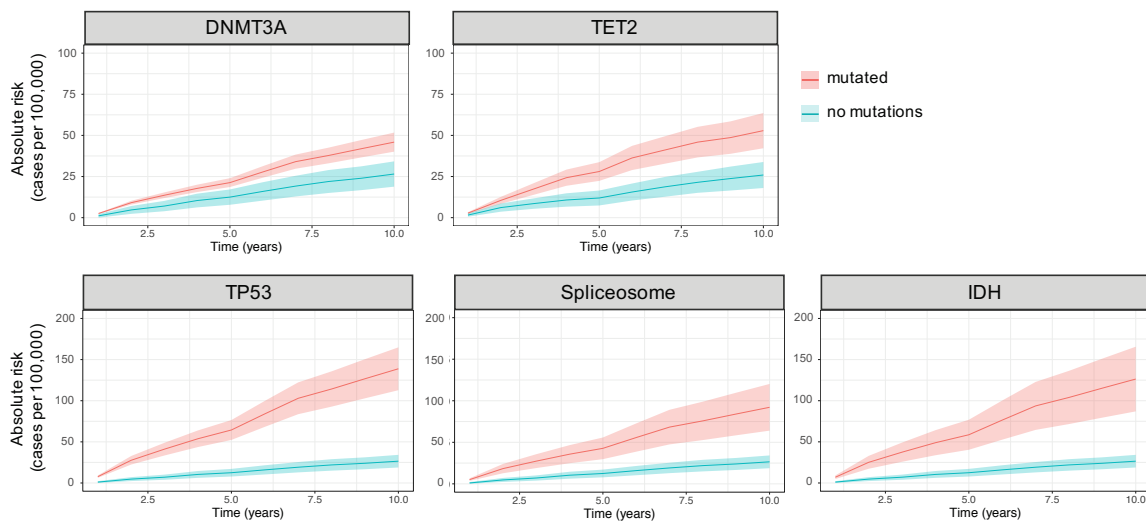


13  
 14  
 15

1 **Figure S12. Absolute risk estimation for recurrently mutated genes.**

2  
 3 Odds ratios from comparison of mutations in cases and controls are converted into estimated absolute  
 4 risk estimates using a weighted partial likelihood method and weighted baseline hazard approach<sup>38</sup>  
 5 described in Supplemental Methods. (a) Analysis shows estimated absolute risk over a 10-year period  
 6 using demographic data from the WHI cohort for individual mutated in the genes indicated (mutated)  
 7 alongside mutation-free participants (no mutations) (n = 59). Shaded region indicates the standard error  
 8 (N=200 bootstrap estimates). Baseline incidence per year is 2.6 AML cases per 100,000 women in the  
 9 absence of mutations. *DNMT3A* (n = 103), *TET2* (n = 57), *TP53* (n = 21), spliceosome (n = 28) and  
 10 IDH genes (n= 15) (b) Incidence rate values per gene as well as for mutation-free participants  
 11 expressed as cases per 100,000 persons per year.  
 12

**a**



**b**

Mutation	Incidence rate (per 100,000 persons per year)
No mutations	2.650
DNMT3A	4.591
TET2	5.462
TP53	13.884
Spliceosome	9.213
IDH	12.634

13  
 14

1 **Figure S13. Mapping of coding alterations to protein domains – Part 1**

2  
3 Mapping of coding alterations to protein domains for recurrently mutated genes in the cohort: *DNMT3A*,  
4 *TET2*, *TP53*, *SRSF2*, *IDH2*, *JAK2*, *SF3B1*, *U2AF1*, *ASXL1*, *IDH1* and *RUNX1*. For each gene, mutations  
5 identified in the AML case and control groups are plotted except for genes that were not mutated in the  
6 control group: *TP53*, *U2AF1* and *RUNX1*. Mutations are classified as Missense (green), In-frame  
7 (brown) and Truncating (grey).

8  
9 Mutations in *IDH2*, *SRSF2*, *JAK2*, *SF3B1* and *U2AF1* occurred in positions R140, P95, V617, K700 and  
10 Q157, respectively. Other point mutations were detected in the HEAT domain of *SF3B1* in close  
11 proximity to K700 or in the zinc finger domain of *U2AF1* in close proximity to Q157. All of these  
12 positions are known hotspots for the aforementioned genes and highly associated with hematological  
13 malignancies and especially AML<sup>69</sup>.

14  
15 Mutations in *DNMT3A* were localized in exons 8-23. 83% of the variants detected in the gene  
16 corresponded to SNV with missense mutations in the R882 position accounting for 26% of the SNVs.  
17 The second most common alteration were truncating variants affecting all the functional domains. SNVs  
18 demonstrated an overall tendency to occur in functional domains whereas truncating mutations occurred  
19 in the N-terminal half of the protein.

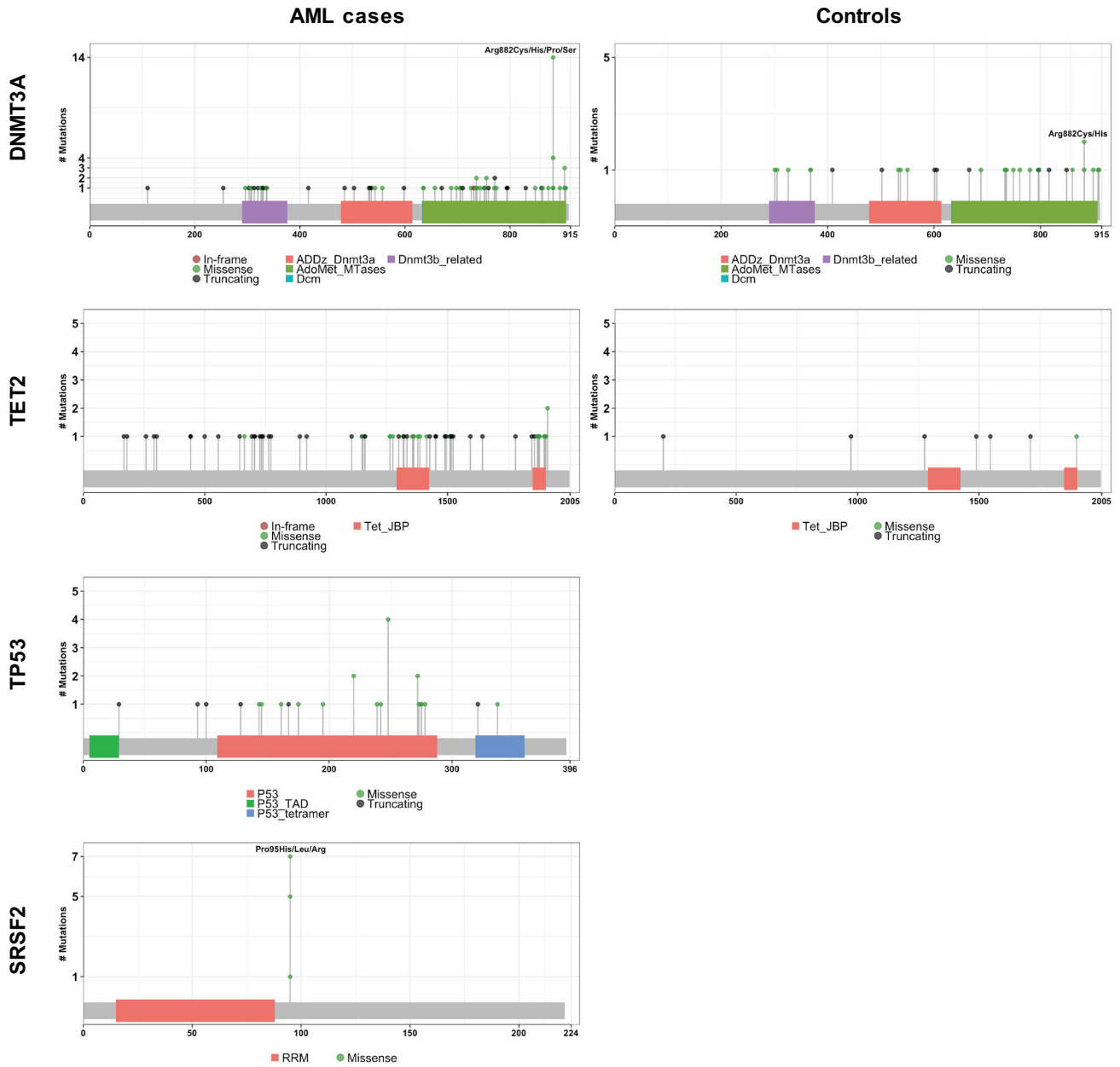
20  
21 Missense mutations comprised 63% of the observed mutations in *TET2* with the majority of the missense  
22 SNVs being confined to oxygenase domain of *TET2* (Tet2\_JDP). Truncating mutations were distributed  
23 across coding exons.

24  
25 Mutations in *ASXL1* were predominantly found in exon 13, with the most common type of alteration  
26 being non-sense SNV. Truncations in the carboxy-terminus or premature stop variants have a disrupting  
27 effect whereas missense variants have an unknown significance<sup>69</sup>.

28  
29 Mutations in *TP53* were also found distributed along the gene. SNVs were concentrated in the DNA-  
30 binding domain with additional variations in tetramerization and transactivating domains. Truncating  
31 mutations occurred throughout. The vast majority of mutations were missense SNV (75%) including  
32 known hotspots with few frameshift deletions, consistent with the mutation pattern observed for  
33 mutations in *TP53* in human cancer<sup>70</sup>.

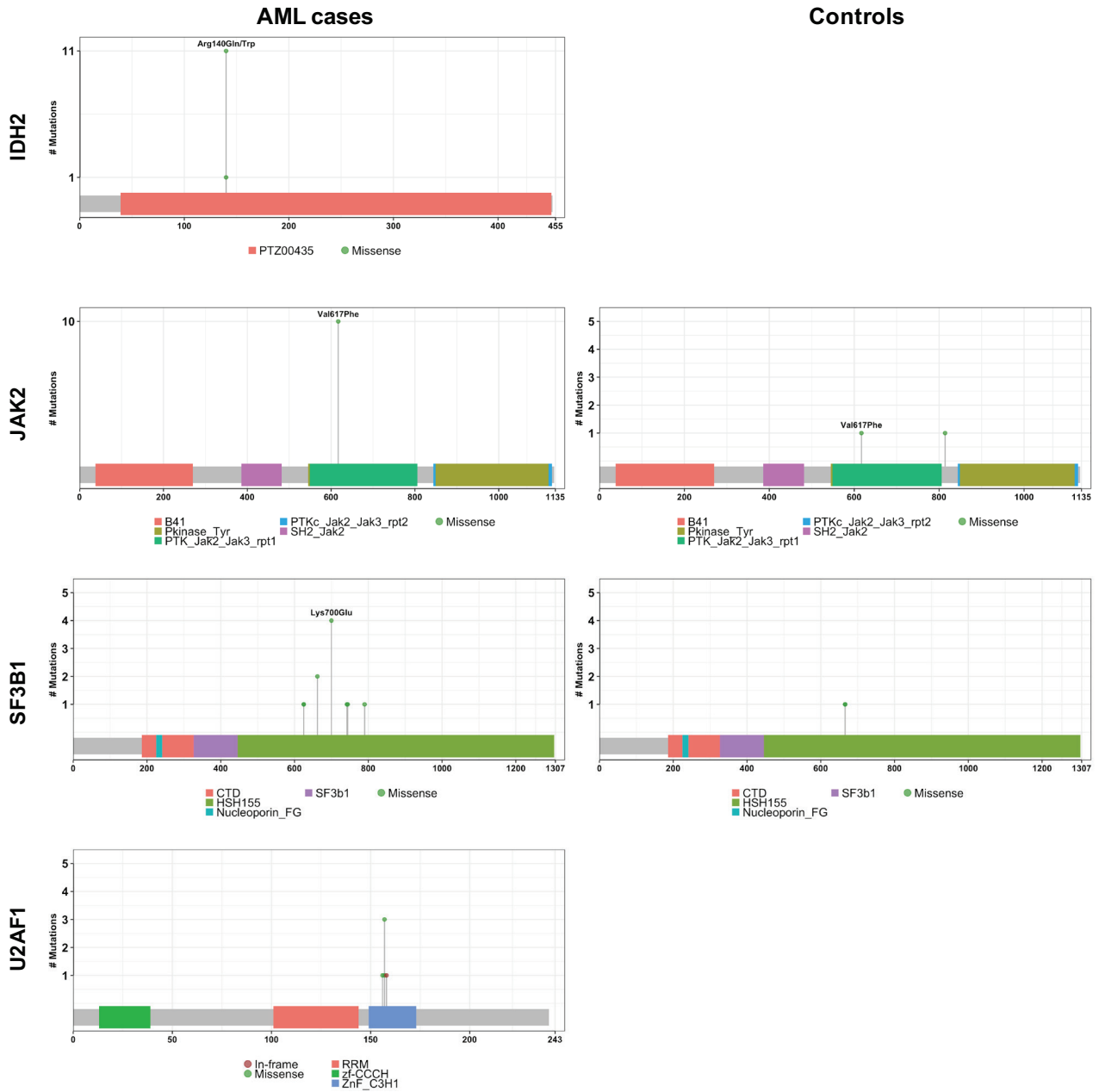
34  
35 A complete list of mutations is available as a supplemental spreadsheet.  
36  
37  
38  
39  
40

1 **Figure S13. Mapping of coding alterations to protein domains – Part 2**  
 2



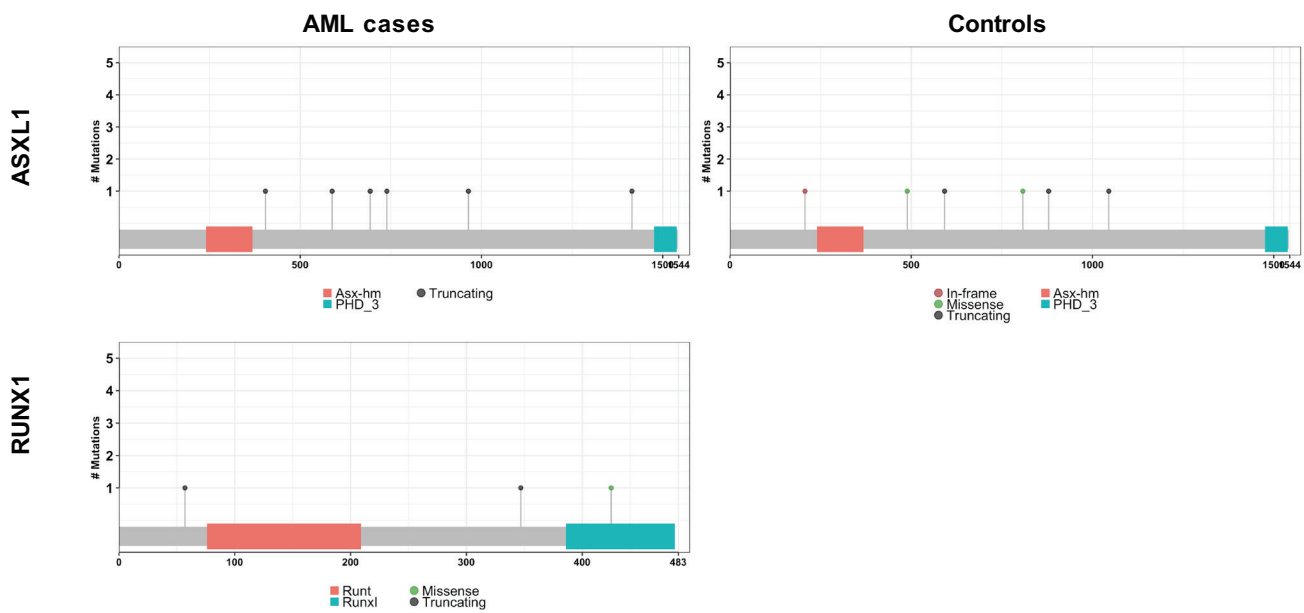
3  
4

1 **Figure S13. Mapping of coding alterations to protein domains – Part 3**  
 2



3  
 4  
 5

1 **Figure S13. Mapping of coding alterations to protein domains – Part 4**  
 2



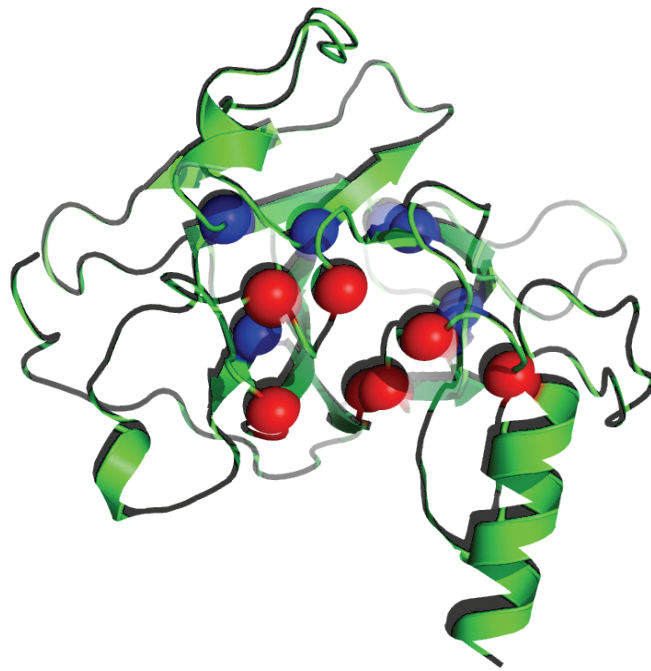
3  
4

1 **Figure S14. Spatial clustering of TP53 point mutations in 3-D structure.**

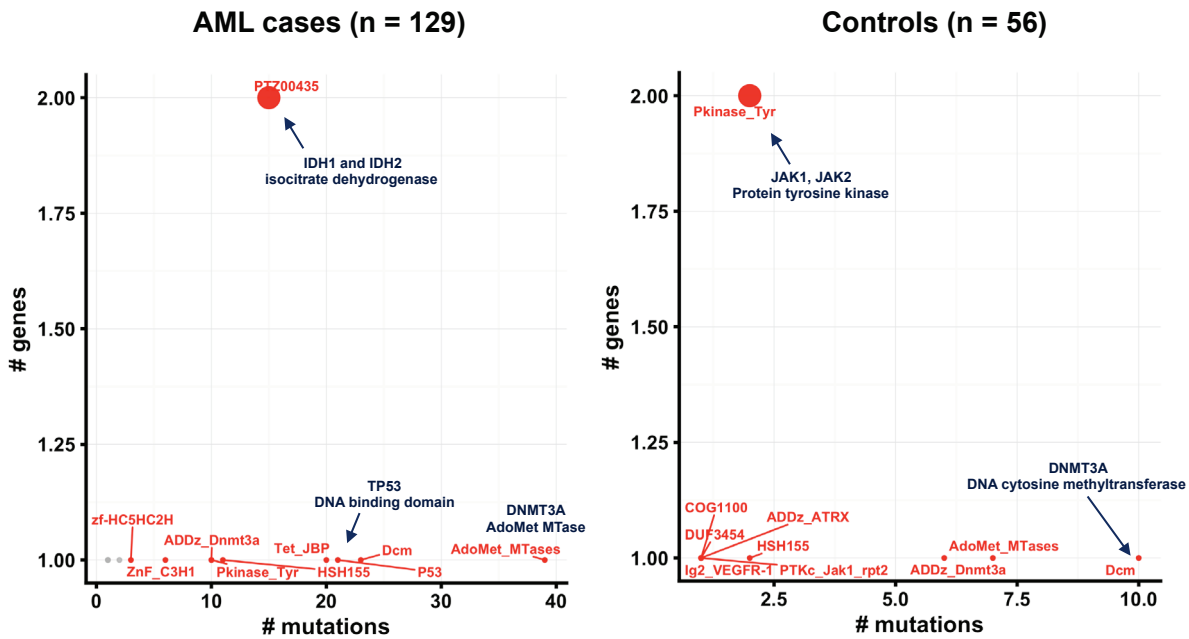
2  
3 Despite having distant amino acid positions, mutations within *TP53* across the AML cases were spatially  
4 localized when mapped to the tertiary structure of the protein. Positions of SNVs are indicated by  
5 spheres. Clustering was performed using mutation3D<sup>61</sup>. Two distinct clusters were identified each  
6 encompassing known regions involved in DNA contact (e.g. R248 and R273) and structural support (e.g.  
7 R175)<sup>71</sup>.

8  
9 Cluster 1 (red spheres) indicates alterations in amino acid positions 239, 242, 248 (4 participants), 272  
10 (2 participants), 273, 275 (2 participants), and 278 (P = 0.02, non-parametric bootstrap). Cluster 2 (blue  
11 spheres) indicates mutations in positions 161, 175, 195, and 234 (P = 0.06, non-parametric bootstrap).  
12 PDB accession number is 2XWR.  
13

14  
15



1 **Figure S15. Distribution of mutations in protein domains (Pfam) in the AML case and control group.**  
 2  
 3 Number of genes (nGenes; vertical axis) is plotted vs. number of mutations (nMuts; horizontal axis).  
 4 AML case samples: Top mutated Pfam domains include PTZ00435 (isocitrate dehydrogenase *IDH1*  
 5 Arg132 and *IDH2* Arg140), AdoMet\_MTases (AdoMet methyltransferase; *DNMT3A*), and P53 (DNA  
 6 binding domain, *TP53*). Control samples: Top mutated Pfam domains include Dcm (DNA cytosine  
 7 methyltransferase; *DNMT3A*), and Pkinase\_Tyr (tyrosine kinase domain; *JAK1* Lys696 and *JAK2*  
 8 Val617).  
 9



10  
 11  
 12  
 13



1 **Figure S16. Mutation frequencies and association of probable pathogenic somatic variants with AML.**

2  
3 (a) Number and frequency of pathogenic mutations in AML cases vs. controls overall and for participants  
4 younger than 65 years vs.  $\geq 65$  years. (b) Forest plot indicating odds ratio of pathogenic mutations in  
5 each gene occurring in the AML cases vs. controls. Genes or gene categories significantly associated  
6 with AML include *TP53* ( $P = 7.9 \times 10^{-6}$ ), *IDH* ( $P = 2.6 \times 10^{-4}$ ), spliceosome, *TET2* ( $P = 5.1 \times 10^{-6}$ ), and  
7 *DNMT3A* ( $P = 4.6 \times 10^{-4}$ ). (c) Forest plot indicating odds of developing AML within 5 years from  
8 baseline, depicted as odds ratios for the specific pathogenic mutations. Mutations in *TP53* and *DNMT3A*  
9 are significantly associated with rapid development of AML. (d) Forest plot indicating odds ratio of  
10 pathogenic mutations in each gene occurring in the AML cases vs. controls adjusted for presence of  
11 mutations in other genes (Others).  $P < 0.001$ : *TET2* ( $P = 4.8 \times 10^{-5}$ ) and *DNMT3A* ( $P = 2.9 \times 10^{-4}$ ). OR  
12 per gene are adjusted by age (years) as a continuous variable. *IDH* category includes *IDH1* and *IDH2*.  
13 The spliceosome category includes *SRSF2*, *SF3B1*, and *U2AF1*. Abbreviations: CI, confidence interval;  
14 N, number affected. P-values are shown for penalized likelihood multivariable logistic regression.

15  
16 All participants bearing mutations in significant genes presented mutations classified as pathogenic, with  
17 the exception of a *DNMT3A* (see Figure 1A and Figure 1B).

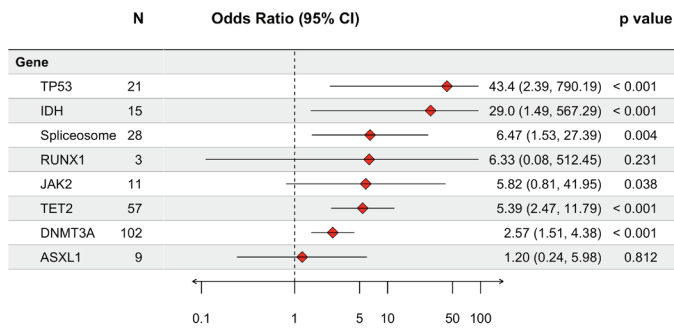
1 **Figure S16. Mutation frequencies and association of probable pathogenic somatic variants with AML**  
 2 **– continued.**  
 3  
 4

**a**

Age	AML cases (N=188)		Controls (N=181)		Odds Ratio	
	# Mutated (%)	# Non-mutated (%)	# Mutated (%)	# Non-mutated (%)	OR (95% CI)	P value
< 65	42 (52.5)	38 (47.5)	15 (19.48)	62 (80.52)	4.52 (2.12-10.05)	2.7 x 10 <sup>-5</sup>
≥ 65	85 (78.7)	23 (21.3)	38 (36.54)	66 (63.46)	6.35 (3.35-12.4)	2.5 x 10 <sup>-11</sup>
Total	127 (67.55)	61 (32.45)	53 (29.28)	128 (70.72)	5.00 (3.16-8.02)	1.1 x 10 <sup>-14</sup>

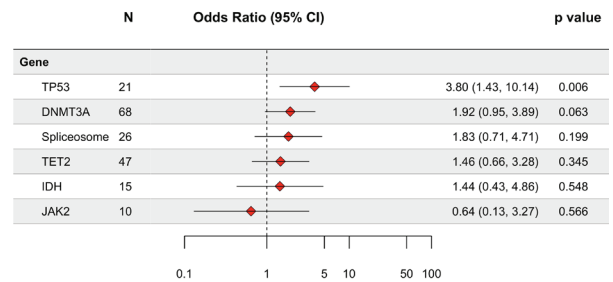
**b**

**Odds ratio: AML cases vs. controls**



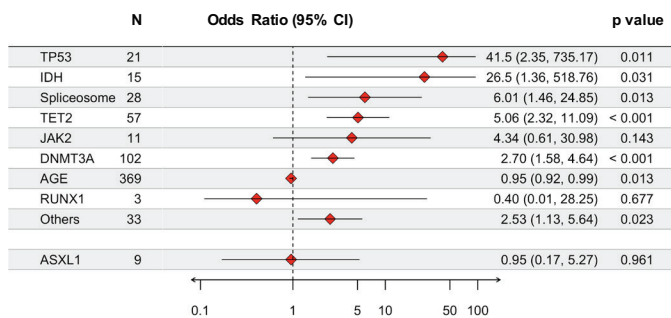
**c**

**Odds ratio: AML diagnosis < 5 years from baseline**



**d**

**Odds ratio: AML cases vs. controls, pathogenic variants with others**

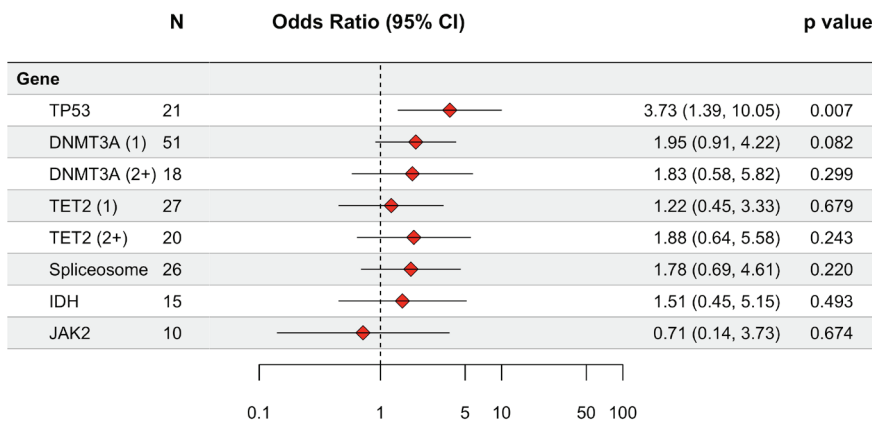


5  
6

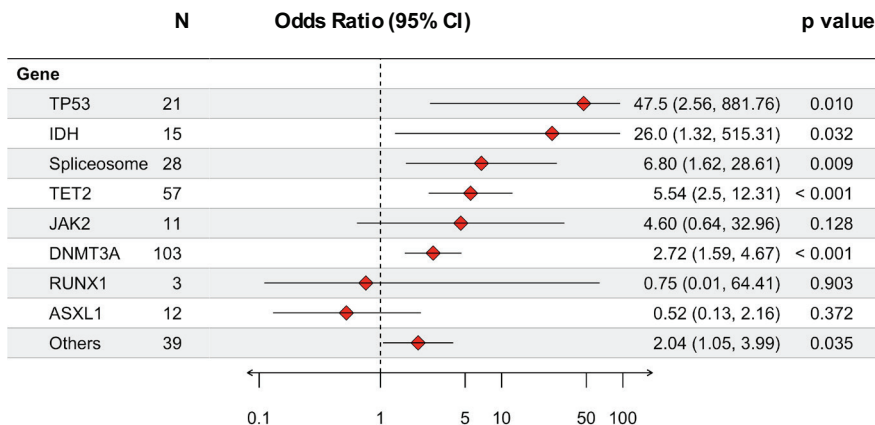
**Figure S17. Odds of AML adjusted for the presence of additional mutations.**

(a) Forest plot indicating odds of developing AML within 5 years from baseline, depicted as odds ratios for mutations in significant genes accounting for presence of 1 or multiple mutations in *DNMT3A* and *TET2*. Mutations in *TP53* are significantly associated with rapid development of AML. The number of mutations per participant in *DNMT3A* or *TET2* does not present significant differences in the odds to develop AML. (b) Forest plot indicating odds ratio of mutations in each gene occurring in the AML cases vs. controls. Odds ratios are corrected by the presence of mutations in other genes not represented in Table 1b.  $P < 0.001$ : *TET2* ( $P = 2.6 \times 10^{-5}$ ) and *DNMT3A* ( $P = 2.6 \times 10^{-4}$ ). OR per gene are adjusted by age (years) as a continuous variable. IDH category includes *IDH1* and *IDH2*. The spliceosome category includes *SRSF2*, *SF3B1*, and *U2AF1*. Abbreviations: CI, confidence interval; N, number affected. P-values are shown for penalized likelihood multivariable logistic regression.

**a Odds ratio: AML diagnosis < 5 years from baseline, including number of variants for DNMT3A and TET2**



**b Odds ratio: AML vs Controls, adjusted with other mutations**

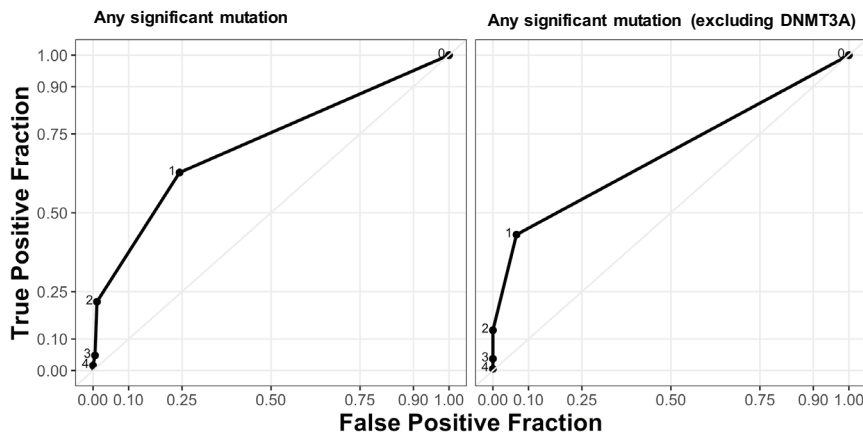


1 **Figure S18. ROC analysis of number of mutations in significant high-risk genes.**

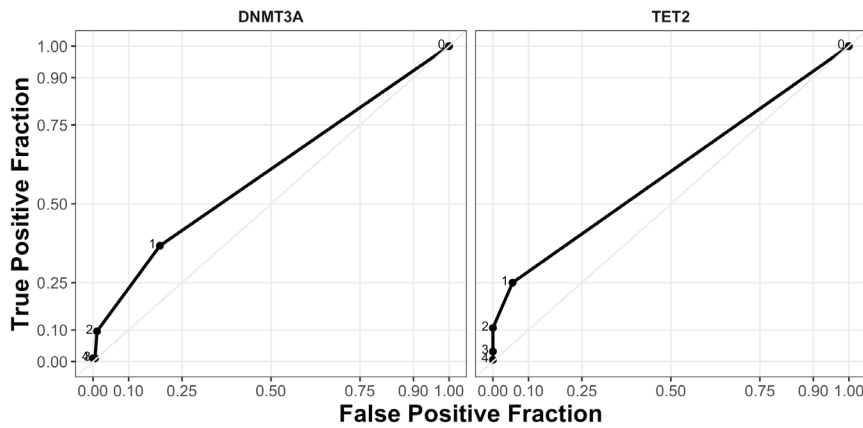
2  
3 Receiver operating characteristic (ROC) curves indicating the % true positive rate (vertical axis) vs. the  
4 % false positive rate (horizontal axis) of the number of mutations to detect AML cases. The curves  
5 indicate performance at decreasing number of variants per gene for (a) significant genes (left plot;  
6 *DNMT3A*, *TET2*, *IDH1*, *IDH2*, *SRSF2*, *SF3B1*, *U2AF1*, *TP53*; n = 164 [118 AML cases, 44 controls])  
7 or the same set of genes excluding *DNMT3A*; n = 94 [81 AML cases, 12 controls] (right plot). (b)  
8 performance is shown for *DNMT3A* and *TET2* genes; left and right plot, respectively.

9  
10 Presence of 2 or more mutations in selected genes was able to detect AML cases with less than 5% false  
11 positive fraction.  
12

**a**



**b**



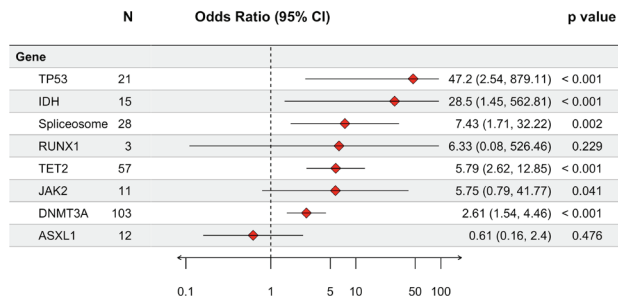
13

**Figure S19. Associations between AML and mutations are robust to different variant classification methods and VAF cutoffs.**

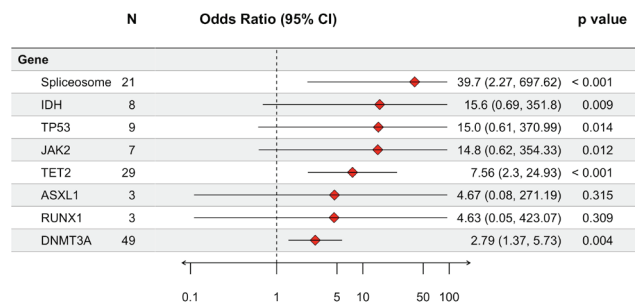
Forest plot indicating odds ratio of mutations per gene occurring in AML cases vs. controls applying different variant classification and VAF cutoff criteria. (a) Somatic variants the indicated genes at VAF > 1%. P < 0.001: *TP53* (P = 5.5 x 10<sup>-6</sup>); *IDH* (P = 3.0 x 10<sup>-4</sup>); *TET2* (P = 2.4 x 10<sup>-6</sup>) and *DNMT3A* (P = 3.4 x 10<sup>-4</sup>). (b) Somatic variants selected for the indicated genes at VAF > 3.5% for SNVs and VAF > 7% for indels. P < 0.001: spliceosome (P = 1.5 x 10<sup>-5</sup>) and *TET2* (P = 9.8 x 10<sup>-5</sup>) (c) Somatic variants selected according to the method of Jaiswal and colleagues<sup>62</sup> using a cutoff of VAF > 1%. P < 0.001: *TP53* (P = 4.8 x 10<sup>-6</sup>); *IDH* (P = 5.0 x 10<sup>-4</sup>); *TET2* (P = 5.0 x 10<sup>-4</sup>) and *DNMT3A* (P = 7.5 x 10<sup>-5</sup>) (d) Somatic variants selected according to the method of Jaiswal and colleagues<sup>62</sup> using a cutoff of VAF > 3.5% for SNVs and VAF > 7% for indels. P < 0.001: spliceosome (P = 9.8 x 10<sup>-6</sup>) and *DNMT3A* (P = 9.5 x 10<sup>-4</sup>). OR per gene are adjusted by age (years) as a continuous variable. *IDH* category includes *IDH1* and *IDH2*. The spliceosome category includes *SRSF2*, *SF3B1*, and *U2AF1*. Abbreviations: CI, confidence interval; N, number affected. P-values are shown for penalized likelihood multivariable logistic regression.

VAF cutoffs has a greater impact on the incidence of mutations than the variant classification criteria. See Table 1b, Table S4 and figure S4.

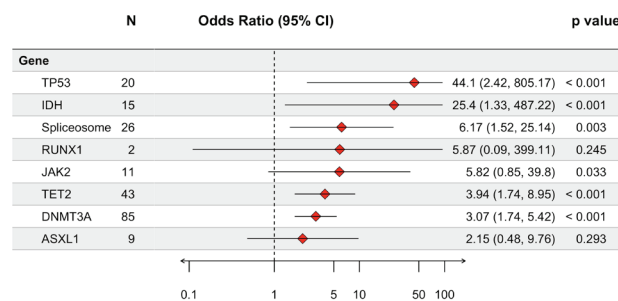
**a Variants Desai VAF > 0.01**



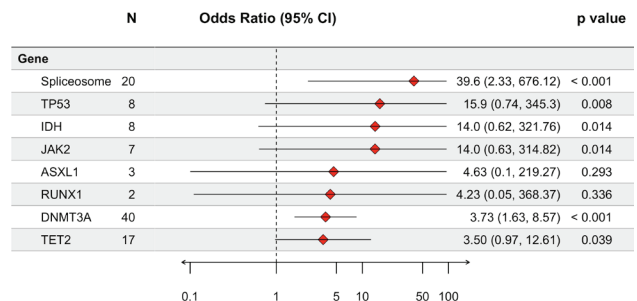
**b VAF<sub>SNV</sub> 0.035 & VAF<sub>indel</sub> 0.07**



**c Variants in Jaiswal VAF > 0.01**



**d Variants in Jaiswal VAF<sub>SNV</sub> 0.035 & VAF<sub>indel</sub> 0.07**

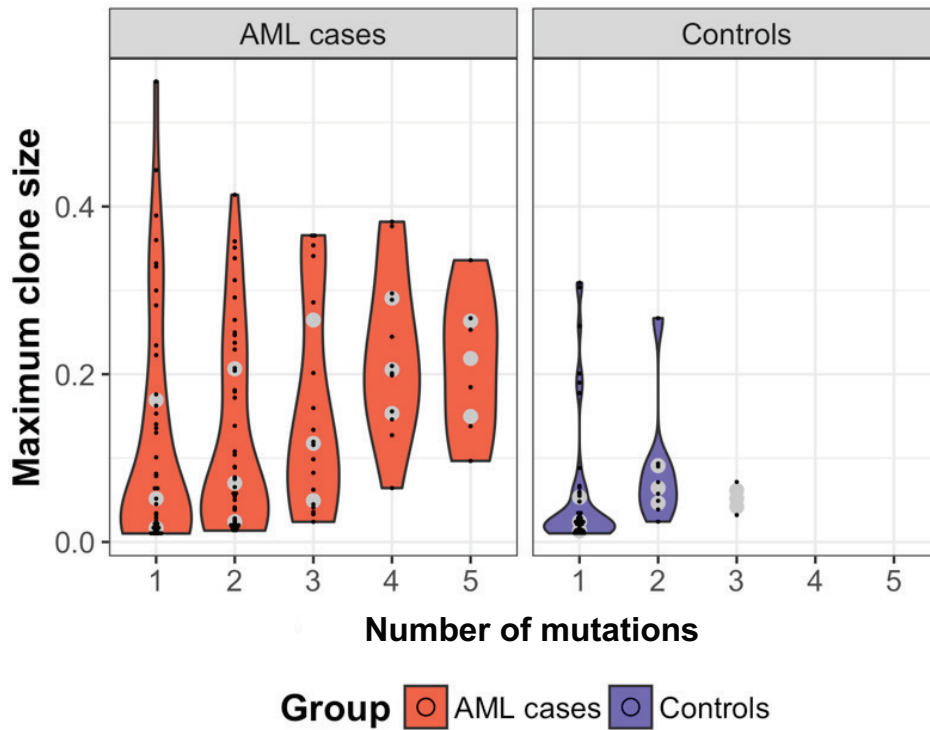


21  
22  
23  
24  
25

1 **Figure S20. Association between number of somatic variants and clone size.**

2  
3 Significant positive correlation was found between the number of somatic variants and the clone size  
4 defined as the maximum VAF of all somatic mutations detected per participant ( $\rho = 0.43$ ,  $P = 1.5 \times 10^{-9}$ ,  
5 Spearman's correlation). Median, 1<sup>st</sup> quantile and 3<sup>rd</sup> quantile maximum clone VAF is shown for each  
6 number of variants (middle, lower and upper grey dots, respectively) for AML cases (n = 125) (red, left  
7 panel) and controls (n=56) (blue, right panel). AML cases: 1 (n = 43, VAF [0.01-0.55]), 2 (n = 46, VAF  
8 [0.01-0.41]), 3 (n = 18, VAF [0.02-0.37]), 4 (n = 12, VAF [0.06-0.38]) and 5 variants (n = 6, VAF  
9 [0.10-0.33]); Controls: 1 (n = 46, VAF [0.01-0.31]), 2 (n = 8, VAF [0.02-0.27]) and 3 variants (n = 2,  
10 VAF [0.03-0.07]). n, number of participants.

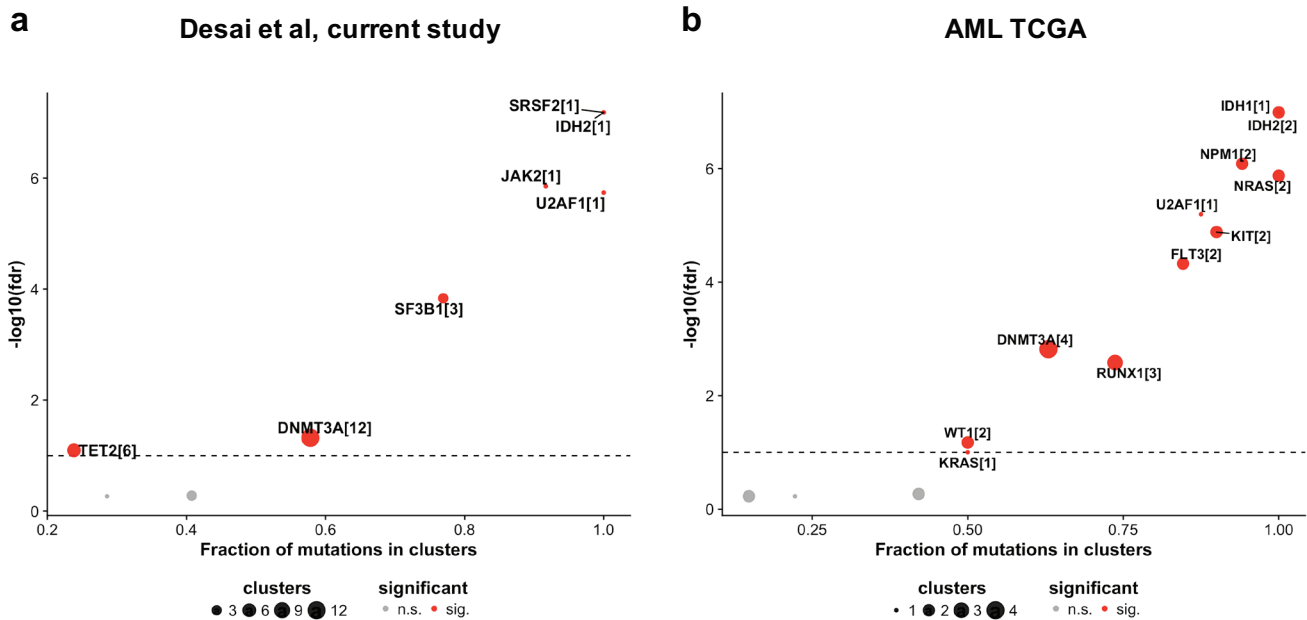
11  
12



13  
14

1 **Figure S21. The pattern of subclonal somatic variations in cases and controls reveals similar driver**  
 2 **genes as a de novo AML cohort.**

3  
 4 Identification of driver genes was performed using the oncoDriveCLUST algorithm implemented in  
 5 *maftools*<sup>72</sup>. The horizontal axis shows the fraction of mutations within clusters while the vertical axis  
 6 indicates the  $-\log_{10}(\text{false discovery rate})$ . The FDR cutoff was set to 1%. Each red dot represents a  
 7 probable driver gene informed by the variant call set. (a) Driver genes identified using all somatic  
 8 variants at baseline evaluation for genes with  $\geq 5$  mutations (N = 302 variants). (b) Driver genes  
 9 identified using all somatic variants reported in the AML TCGA study for genes with  $\geq 5$  mutations (N  
 10 = 305 variants). The size of the dot is related to the number of clusters per gene with the number of  
 11 clusters indicated by the number in brackets. Significant clusters were identified for *SRSF2*, *IDH2*, *JAK2*,  
 12 *U2AF1*, *SF3B1*, *DNMT3A* and *TET2*. Similarly, the TCGA AML cohort reveals *IDH1*, *IDH2*, *NRAS*,  
 13 *U2AF1*, *KIT*, *FLT3*, *DNMT3A*, *RUNX1*, *WT1*, and *KRAS*.  
 14

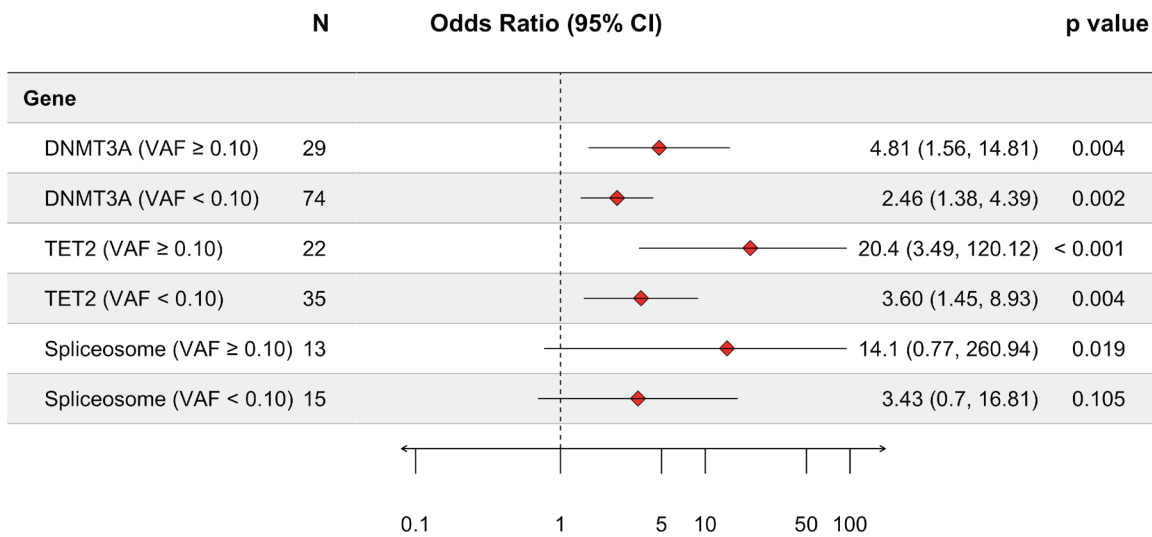


15

1 **Figure S22. Odds of AML are elevated when mutations are present at higher VAF.**

2  
3 Forest plot indicating odds ratio of mutations in *DNMT3A*, *TET2* and spliceosome genes at high VAF (>  
4 10%) vs low VAF (< 10%). The spliceosome category includes *SRSF2*, *SF3B1*, and *U2AF1*. Odds ratio  
5 (OR) per gene are adjusted by age (years) as a continuous variable. Abbreviations: CI, confidence  
6 interval; N, number affected. P-values are shown for penalized likelihood multivariable logistic  
7 regression. Exact p-values: TET2 (P = 7.5 x 10<sup>-6</sup>).

8  
9 Participants with mutations in these genes and VAF > 10% have increased odds for AML development  
10



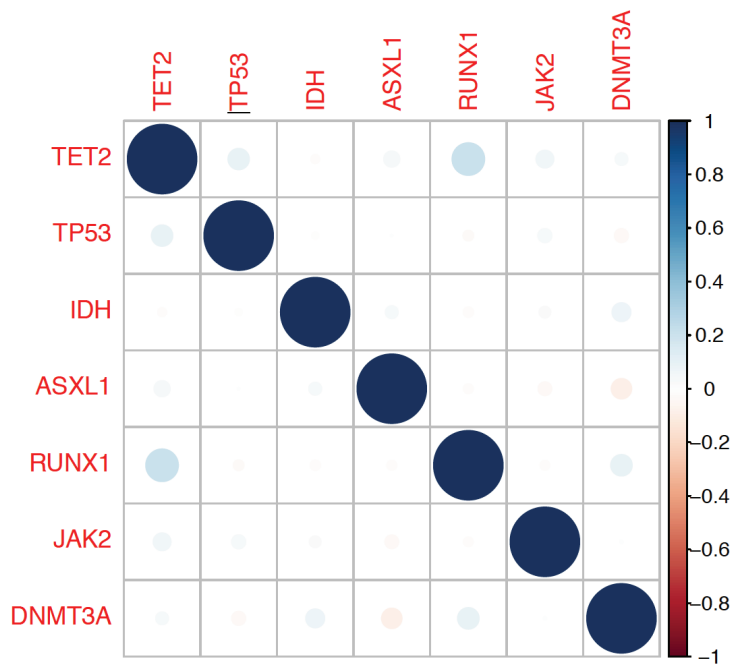


1 **Figure S23. Absence of collinearity between predictors.**

2

3 Correlation (Pearson's) matrix between recurrently mutated genes demonstrates the absence of  
4 collinear variables in multivariable models (N=185 ).

5



6

7

1 **Short list of WHI investigators**  
2

3 **Program Office:** (National Heart, Lung, and Blood Institute, Bethesda, Maryland) Jacques Rossouw,  
4 Shari Ludlam, Joan McGowan, Leslie Ford, and Nancy Geller  
5

6 **Clinical Coordinating Center:** (Fred Hutchinson Cancer Research Center, Seattle, WA) Garnet  
7 Anderson, Ross Prentice, Andrea LaCroix, and Charles Kooperberg I  
8

9 **Investigators and Academic Centers:** (Brigham and Women's Hospital, Harvard Medical School,  
10 Boston, MA) JoAnn E. Manson; (MedStar Health Research Institute/Howard University, Washington,  
11 DC) Barbara V. Howard; (Stanford Prevention Research Center, Stanford, CA) Marcia L. Stefanick;  
12 (The Ohio State University, Columbus, OH) Rebecca Jackson; (University of Arizona,  
13 Tucson/Phoenix, AZ) Cynthia A. Thomson; (University at Buffalo, Buffalo, NY) Jean Wactawski-  
14 Wende; (University of Florida, Gainesville/Jacksonville, FL) Marian Limacher; (University of Iowa,  
15 Iowa City/Davenport, IA) Jennifer Robinson; (University of Pittsburgh, Pittsburgh, PA) Lewis Kuller;  
16 (Wake Forest University School of Medicine, Winston-Salem, NC) Sally Shumaker; (University of  
17 Nevada, Reno, NV) Robert Brunner; (University of Minnesota, Minneapolis, MN) Karen L. Margolis  
18

19 **Women's Health Initiative Memory Study:** (Wake Forest University School of Medicine, Winston-  
20 Salem, NC) Mark Espeland For a list of all the investigators who have contributed to WHI science,  
21 please visit:  
22

23 [https://www.whi.org/researchers/Documents%20%20Write%20a%20Paper/WHI%20Investigator](https://www.whi.org/researchers/Documents%20%20Write%20a%20Paper/WHI%20Investigator%20Long%20List.pdf)  
24 [%20Long%20List.pdf](https://www.whi.org/researchers/Documents%20%20Write%20a%20Paper/WHI%20Investigator%20Long%20List.pdf)  
25  
26

## References

- 35 Kircher, M., Sawyer, S. & Meyer, M. Double indexing overcomes inaccuracies in multiplex sequencing on the Illumina platform. *Nucleic Acids Res* **40**, e3, doi:10.1093/nar/gkr771 (2012).
- 36 Firth, D. Bias Reduction of Maximum Likelihood Estimates. *Biometrika* **80**, 27-38, doi:10.2307/2336755 (1993).
- 37 Cai, T. & Zheng, Y. Non-parametric Evaluation of Biomarker Accuracy under Nested Case-control Studies. *J Am Stat Assoc* **106**, 569-580, doi:10.1198/jasa.2011.tm09807 (2011).
- 38 Samuelsen, S. O. A pseudolikelihood approach to analysis of nested case-control studies. *Biometrika* **84**, 379-394, doi:10.1093/biomet/84.2.379 (1997).
- 39 Ihaka, R. & Gentleman, R. R. A language for data analysis and graphics. *Journal of Computational and Graphical Statistics* **5**, 299-314, doi:10.2307/1390807 (1996).
- 40 Heinze, G. *Logistic regression using Firth's bias reduction: a solution to the problem of separation in logistic regression*, <<https://cemsis.meduniwien.ac.at/en/kb/science-research/software/statistical-software/fllogistf/>> (
- 41 Lindgreen, S. AdapterRemoval: easy cleaning of next-generation sequencing reads. *BMC Res Notes* **5**, 337, doi:10.1186/1756-0500-5-337 (2012).
- 42 Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv preprint arXiv* **00**, 3-3, doi:arXiv:1303.3997 [q-bio.GN] (2013).
- 43 Faust, G. G. & Hall, I. M. SAMBLASTER: fast duplicate marking and structural variant read extraction. *Bioinformatics* **30**, 2503-2505, doi:10.1093/bioinformatics/btu314 (2014).
- 44 Beraldi, D. *MarkDupsByStartEnd*, <<https://github.com/dariober/Java-cafe/tree/master/MarkDupsByStartEnd>> (
- 45 Lai, Z. *et al.* VarDict: a novel and versatile variant caller for next-generation sequencing in cancer research. *Nucleic Acids Res* **44**, e108, doi:10.1093/nar/gkw227 (2016).
- 46 Talevich, E., Shain, A. H., Botton, T. & Bastian, B. C. CNVkit: Genome-Wide Copy Number Detection and Visualization from Targeted DNA Sequencing. *PLoS Comput Biol* **12**, e1004873, doi:10.1371/journal.pcbi.1004873 (2016).
- 47 McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* **17**, 122, doi:10.1186/s13059-016-0974-4 (2016).
- 48 Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* **6**, 80-92, doi:10.4161/fly.19695 (2012).
- 49 *Vt analysis toolkit*, <<https://github.com/atks/vt>> (2017).
- 50 Li, H. Toward better understanding of artifacts in variant calling from high-coverage samples. *Bioinformatics* **30**, 2843-2851, doi:10.1093/bioinformatics/btu356 (2014).
- 51 Forbes, S. A. *et al.* COSMIC: somatic cancer genetics at high-resolution. *Nucleic Acids Res* **45**, D777-D783, doi:10.1093/nar/gkw1121 (2017).
- 52 *The Cancer Genome Atlas. National Cancer Institute*, <<https://cancergenome.nih.gov/>> (
- 53 Jun, G. *et al.* Detecting and estimating contamination of human DNA samples in sequencing and array-based genotype data. *Am J Hum Genet* **91**, 839-848, doi:10.1016/j.ajhg.2012.09.004 (2012).
- 54 Landau, D. A. *et al.* The evolutionary landscape of chronic lymphocytic leukemia treated with ibrutinib targeted therapy. *Nat Commun* **8**, 2185, doi:10.1038/s41467-017-02329-y (2017).
- 55 *Picard Tools*, <<http://broadinstitute.github.io/picard/>> (
- 56 Jagadeesh, K. A. *et al.* M-CAP eliminates a majority of variants of uncertain significance in clinical exomes at high sensitivity. *Nat Genet* **48**, 1581-1586, doi:10.1038/ng.3703 (2016).

1 57 Jagadeesh, K. A., Wenger, A., Berger, M., Guturu, H., Stenson, P., Cooper, D., Bernstein, J.,  
2 and Bejerano, G. *Mendelian Clinically Applicable Pathogenicity (M-CAP) Score*,  
3 <<http://bejerano.stanford.edu/mcap/>> (2016).

4 58 Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in  
5 multidimensional genomic data. *Bioinformatics* **32**, 2847-2849,  
6 doi:10.1093/bioinformatics/btw313 (2016).

7 59 Gu, Z., Gu, L., Eils, R., Schlesner, M. & Brors, B. circlize Implements and enhances circular  
8 visualization in R. *Bioinformatics* **30**, 2811-2812, doi:10.1093/bioinformatics/btu393 (2014).

9 60 Janssen, R. *MutationalPatterns*, <<https://github.com/UMCUGenetics/MutationalPatterns>> (  
10 61 Meyer, M. J. *et al.* mutation3D: Cancer Gene Prediction Through Atomic Clustering of Coding  
11 Variants in the Structural Proteome. *Hum Mutat* **37**, 447-456, doi:10.1002/humu.22963 (2016).

12 62 Jaiswal, S. *et al.* Age-related clonal hematopoiesis associated with adverse outcomes. *N Engl J*  
13 *Med* **371**, 2488-2498, doi:10.1056/NEJMoa1408617 (2014).

14 63 Jennings, L. J. *et al.* Guidelines for Validation of Next-Generation Sequencing-Based  
15 Oncology Panels: A Joint Consensus Recommendation of the Association for Molecular  
16 Pathology and College of American Pathologists. *J Mol Diagn* **19**, 341-365,  
17 doi:10.1016/j.jmoldx.2017.01.011 (2017).

18 64 Ryan, S. L. *et al.* The role of the RAS pathway in iAMP21-ALL. *Leukemia* **30**, 1824-1831,  
19 doi:10.1038/leu.2016.80 (2016).

20 65 Roberts, S. A. *et al.* An APOBEC cytidine deaminase mutagenesis pattern is widespread in  
21 human cancers. *Nat Genet* **45**, 970-976, doi:10.1038/ng.2702 (2013).

22 66 Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415-  
23 421, doi:10.1038/nature12477 (2013).

24 67 Hanawalt, P. C. & Spivak, G. Transcription-coupled DNA repair: two decades of progress and  
25 surprises. *Nat Rev Mol Cell Biol* **9**, 958-970, doi:10.1038/nrm2549 (2008).

26 68 Gorbunova, V., Seluanov, A., Mao, Z. & Hine, C. Changes in DNA repair during aging.  
27 *Nucleic Acids Res* **35**, 7466-7474, doi:10.1093/nar/gkm756 (2007).

28 69 Khwaja, A. *et al.* Acute myeloid leukaemia. *Nat Rev Dis Primers* **2**, 16010,  
29 doi:10.1038/nrdp.2016.10 (2016).

30 70 Kadia, T. M. *et al.* TP53 mutations in newly diagnosed acute myeloid leukemia:  
31 Clinicomolecular characteristics, response to therapy, and outcomes. *Cancer*,  
32 doi:10.1002/encr.30203 (2016).

33 71 Olivier, M., Hollstein, M. & Hainaut, P. TP53 mutations in human cancers: origins,  
34 consequences, and clinical use. *Cold Spring Harb Perspect Biol* **2**, a001008,  
35 doi:10.1101/cshperspect.a001008 (2010).

36 72 Tamborero, D., Gonzalez-Perez, A. & Lopez-Bigas, N. OncodriveCLUST: exploiting the  
37 positional clustering of somatic mutations to identify cancer genes. *Bioinformatics* **29**, 2238-  
38 2244, doi:10.1093/bioinformatics/btt395 (2013).

39