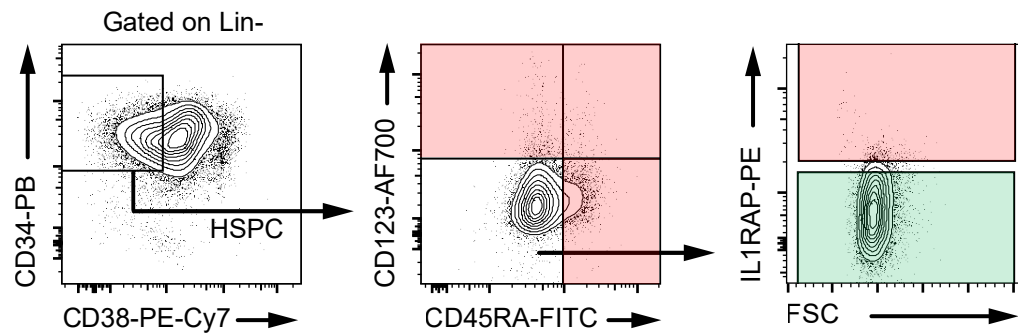


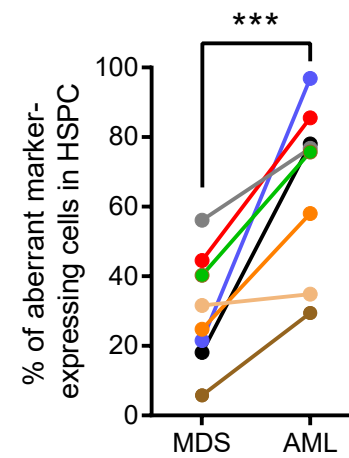
# Figure S1

**a**

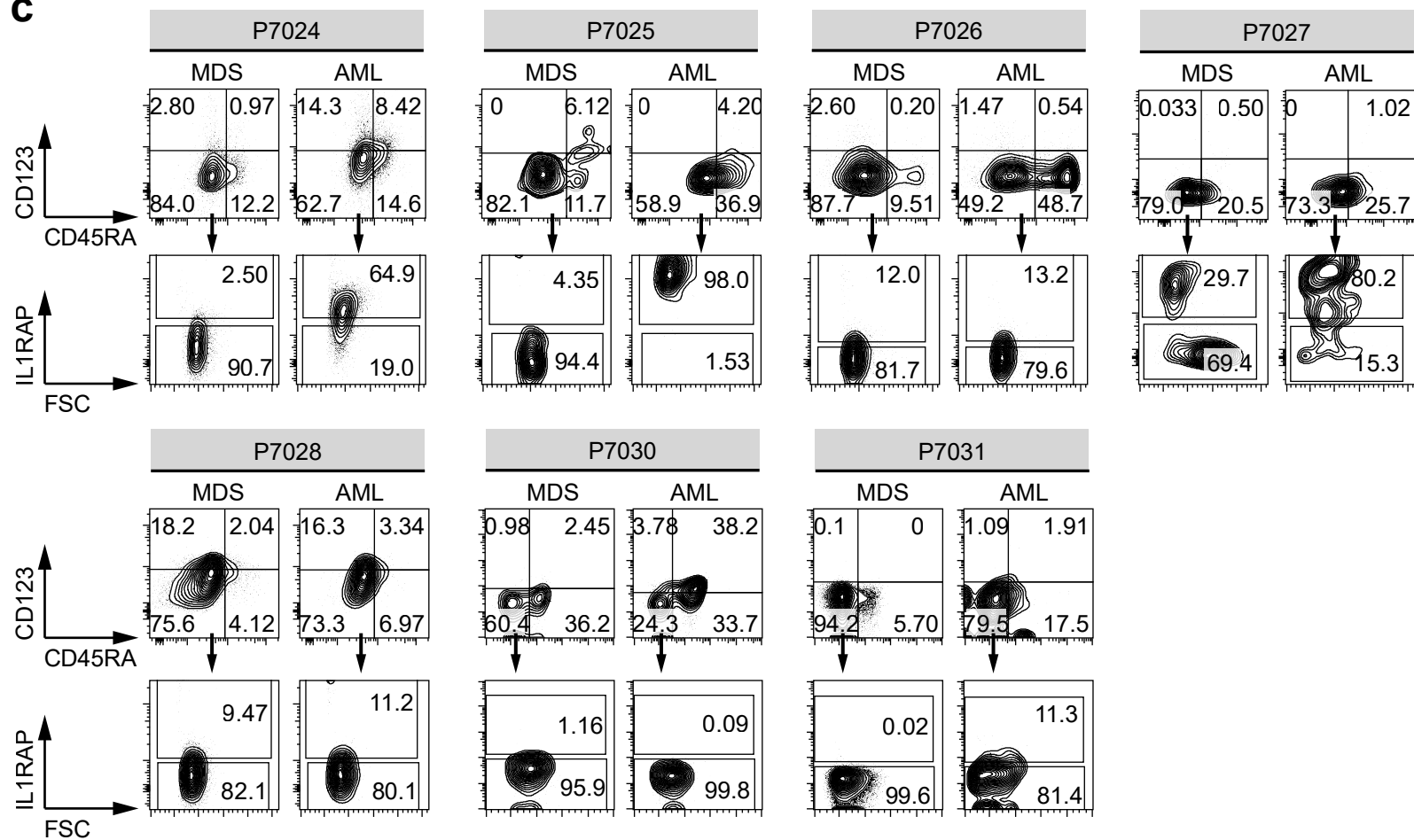


Population	Surface markers
PreMDS-SC or PreAML-SC	Lin-CD34+CD38-CD45RA-CD123-IL1RAP-
MDS-SC or AML-SC	Lin-CD34+CD38-CD45RA+/CD123+/IL1RAP+

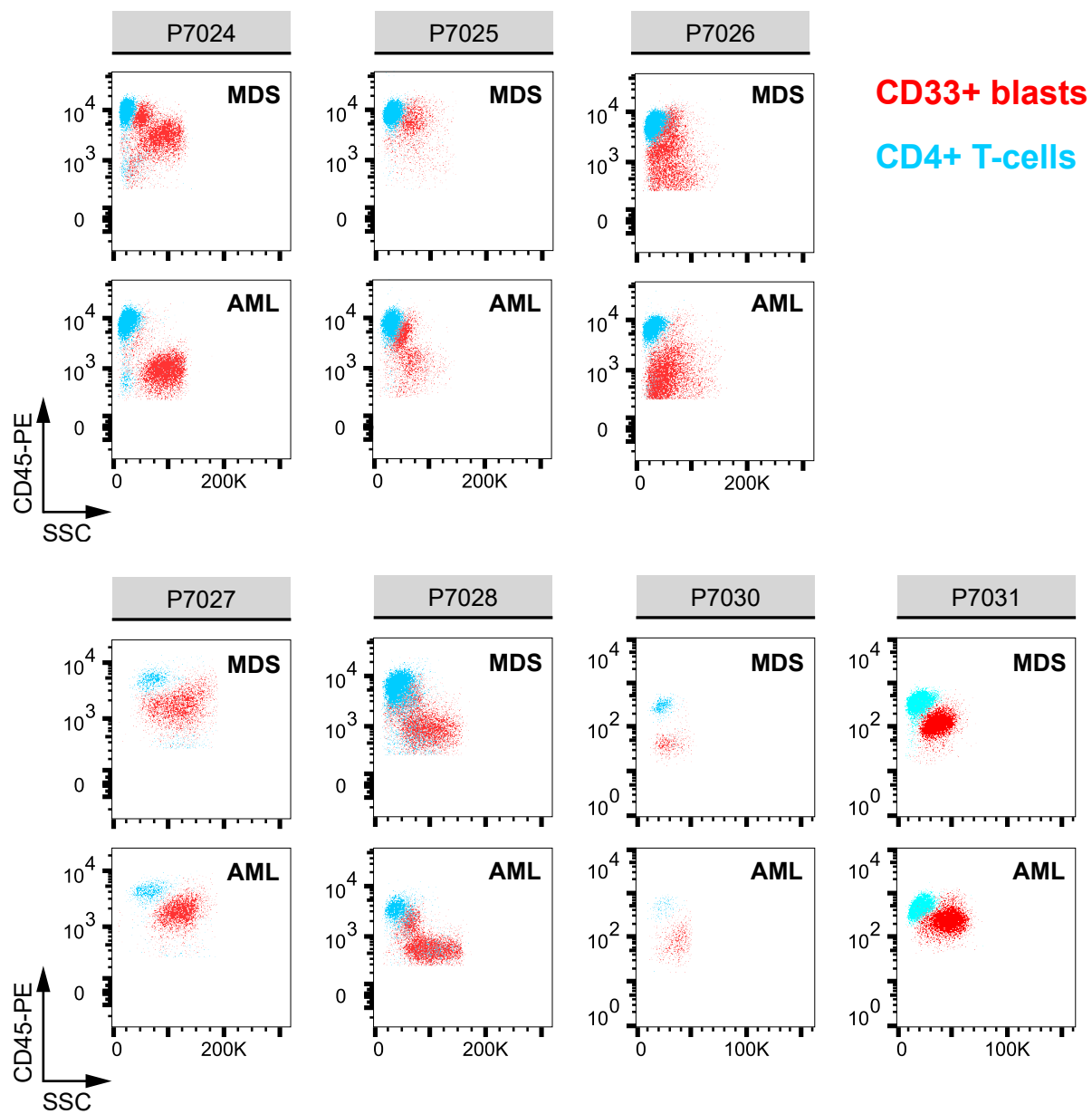
**b**

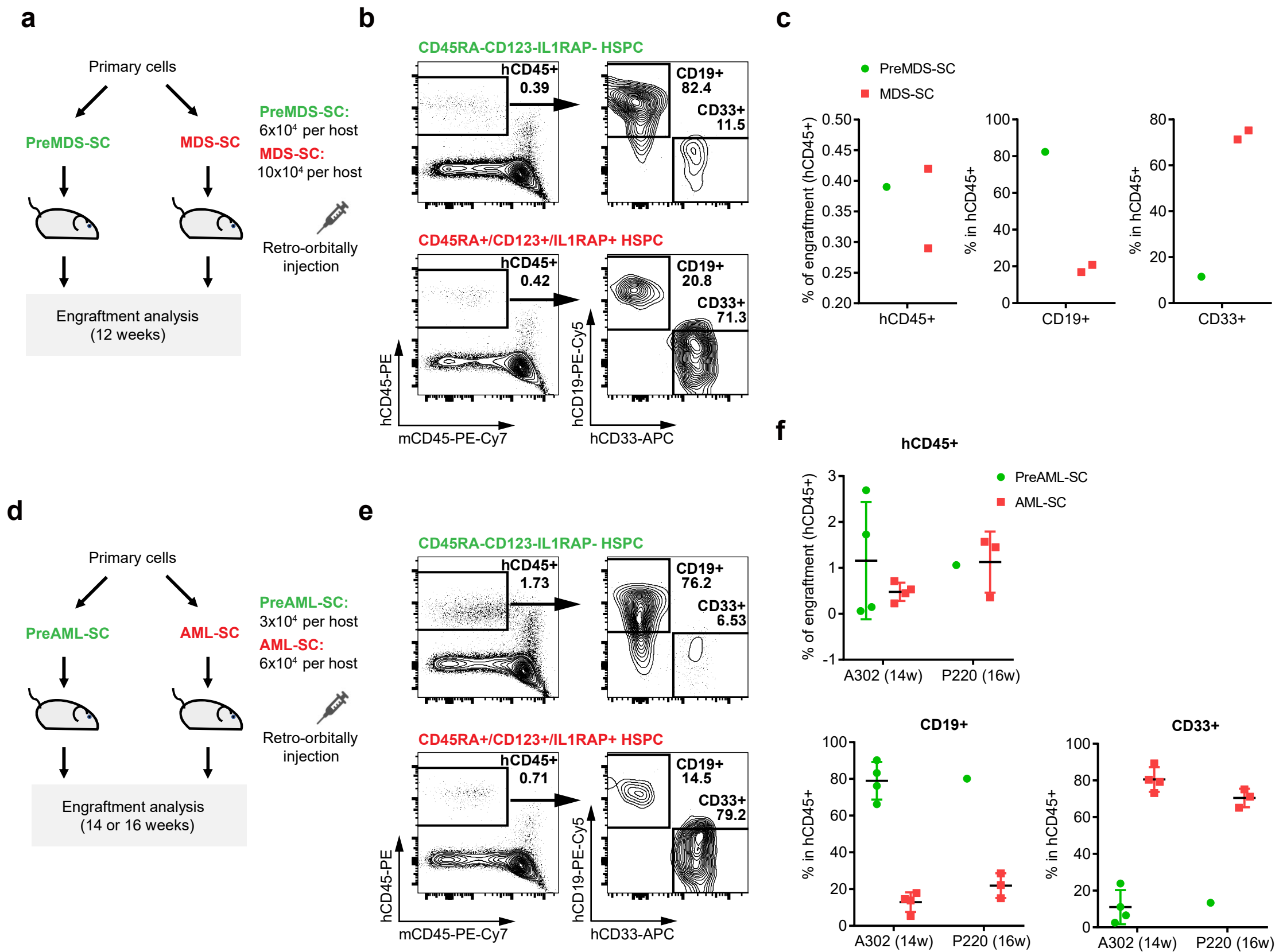


**c**



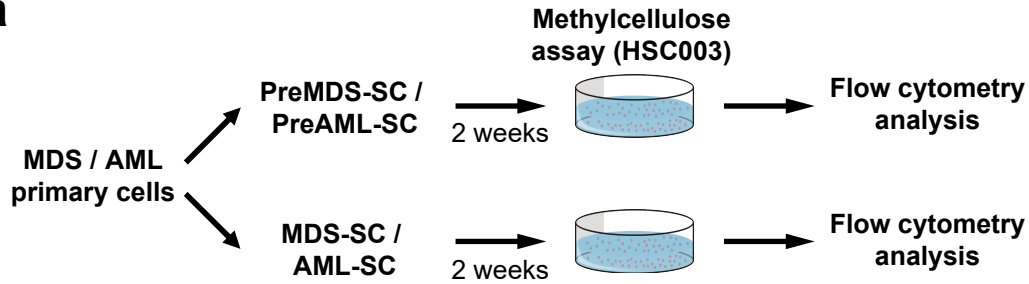
# Figure S2



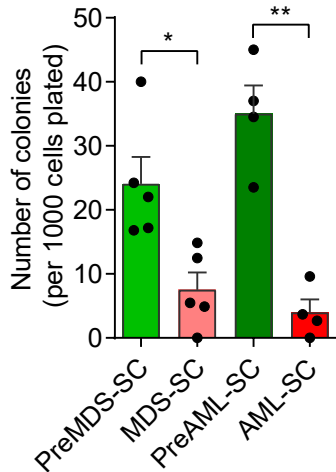


# Figure S4

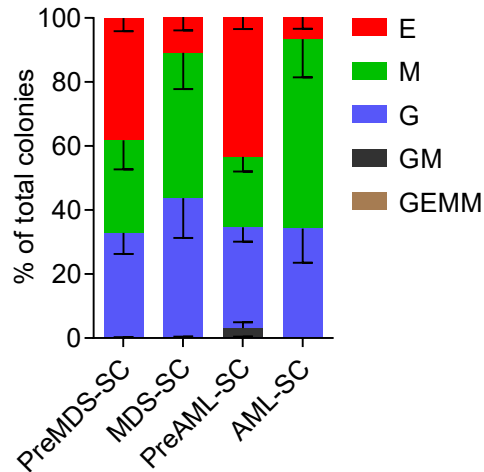
**a**



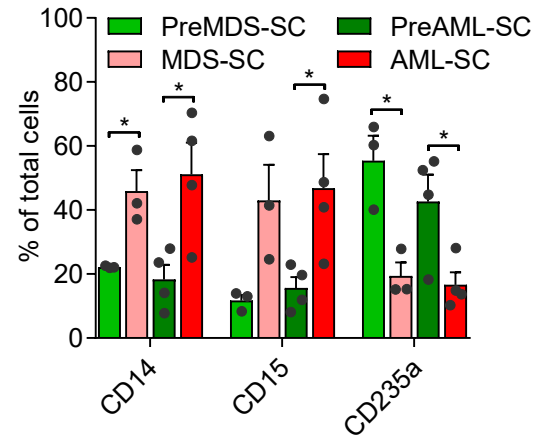
**b**

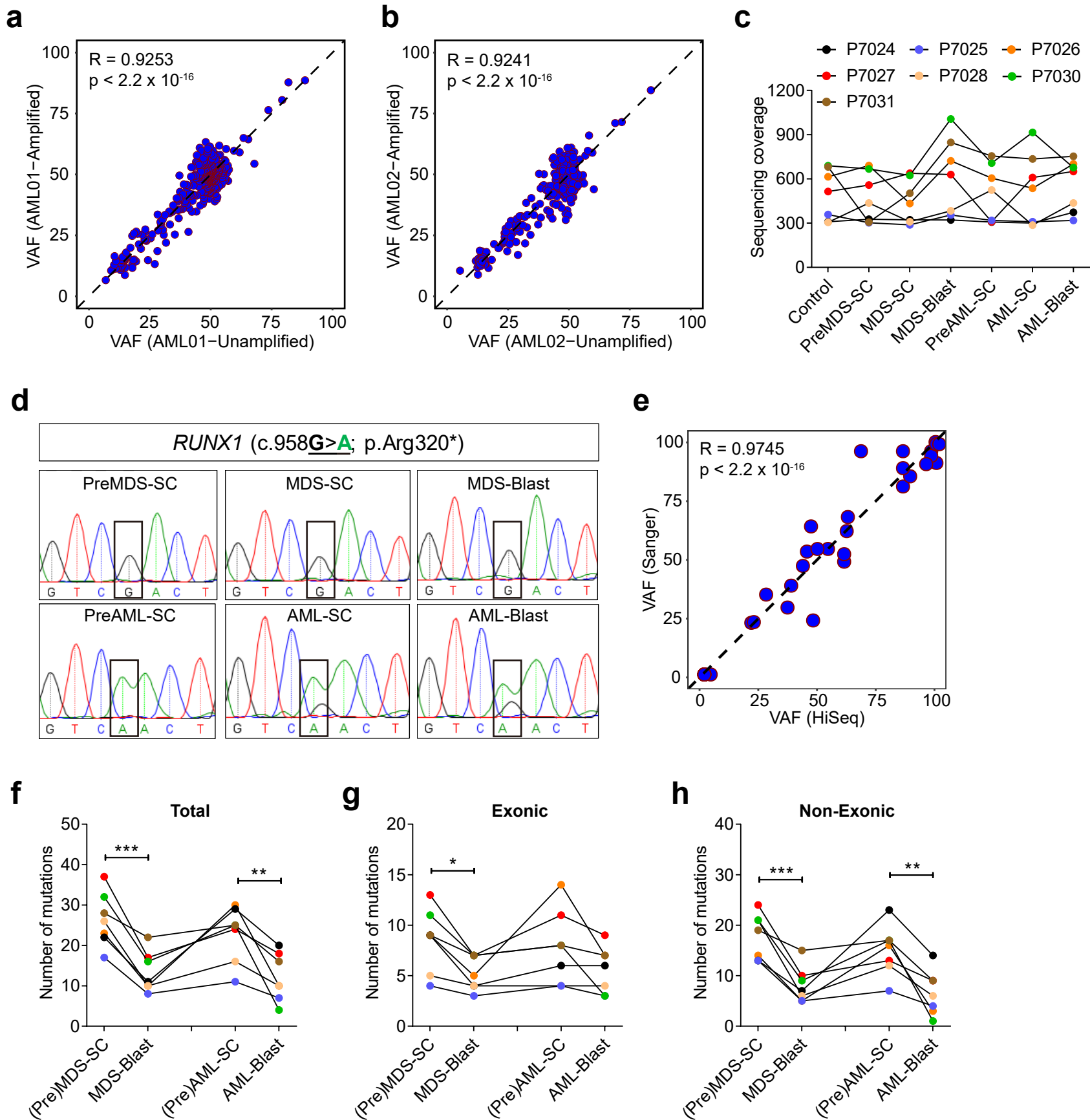


**c**

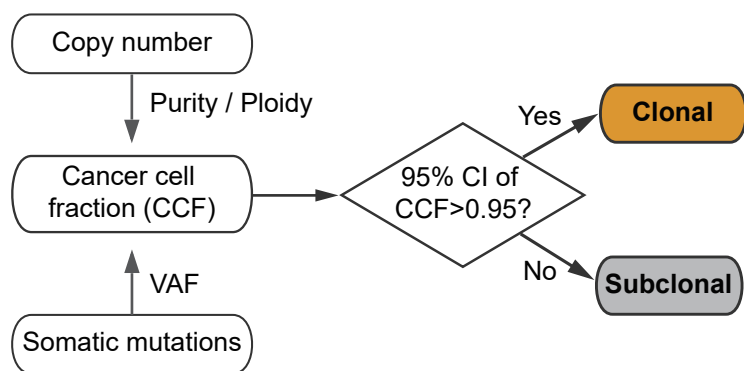


**d**

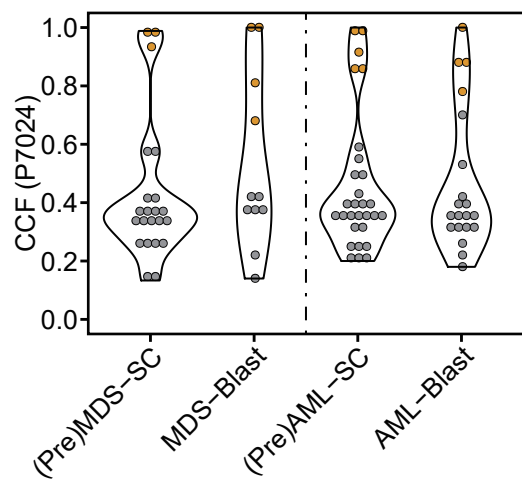


**Figure S5**

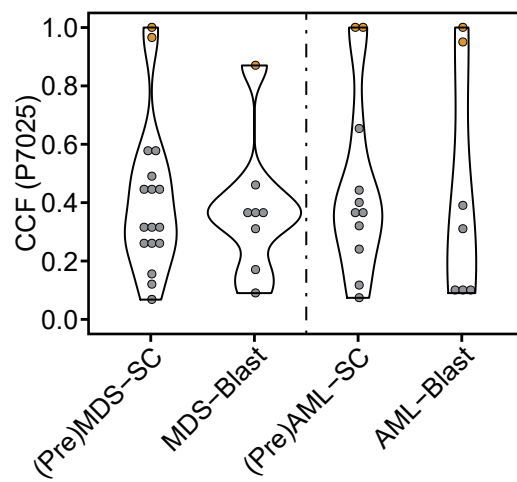
**a**



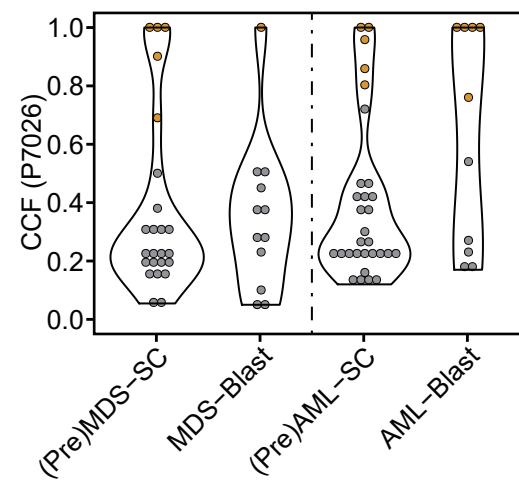
**b**



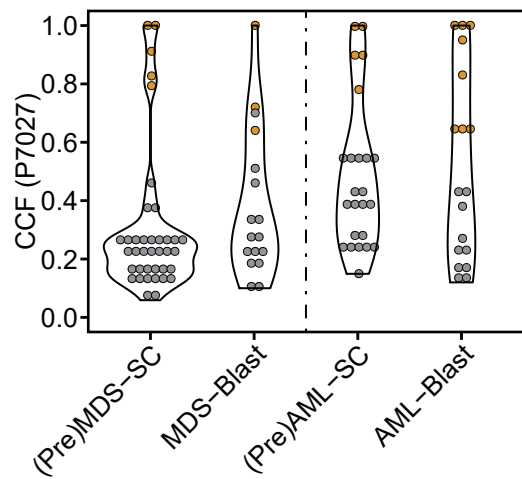
**c**



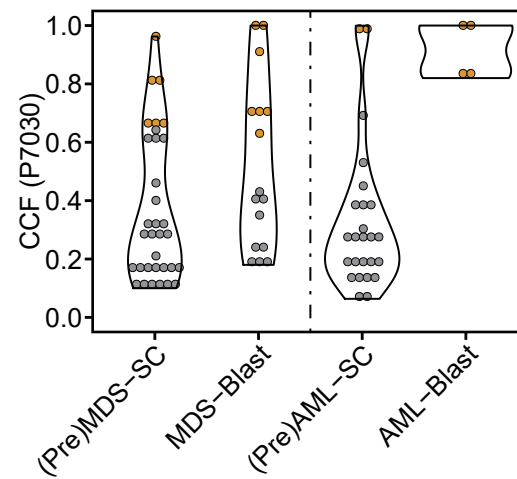
**d**



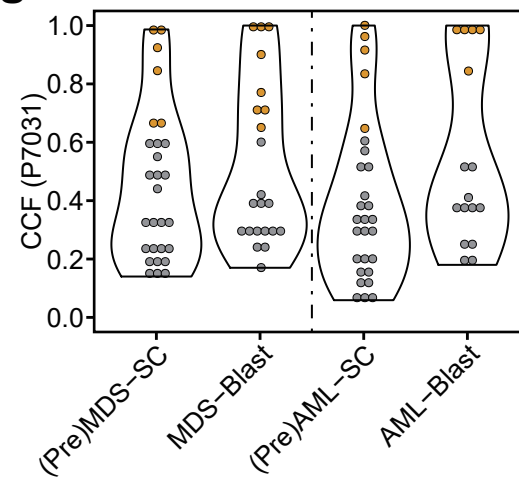
**e**



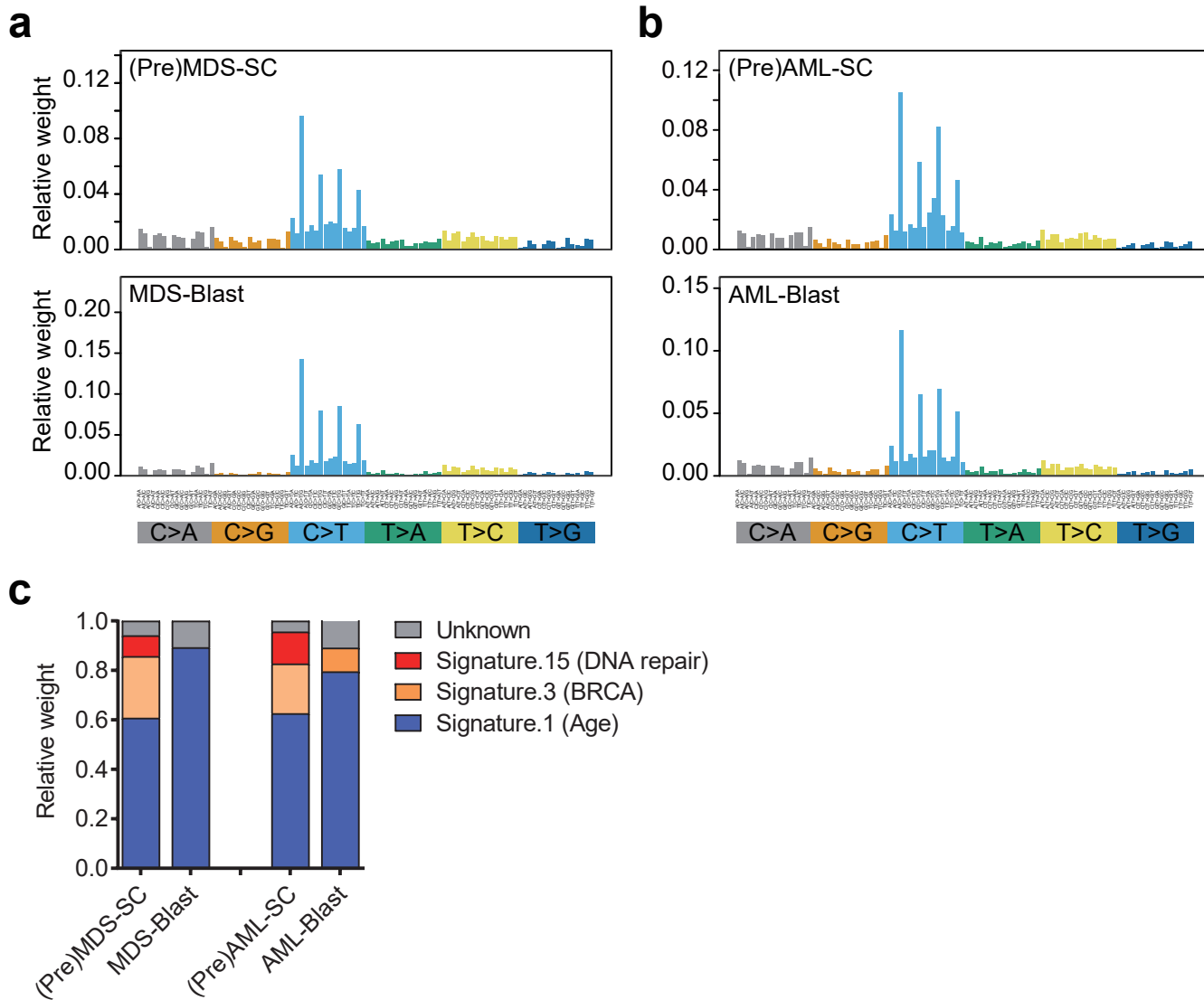
**f**



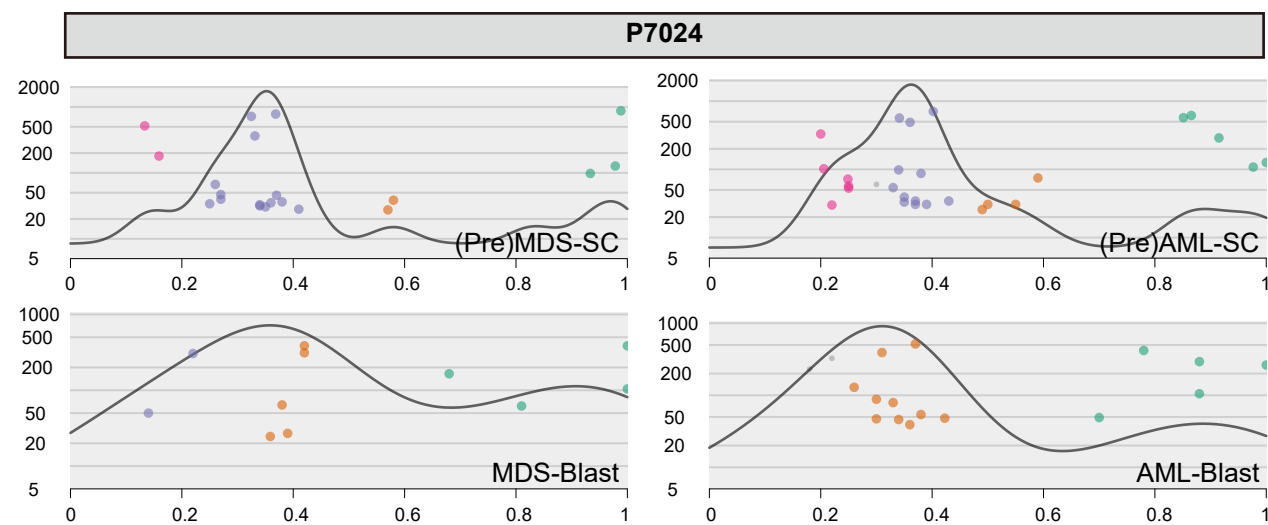
**g**



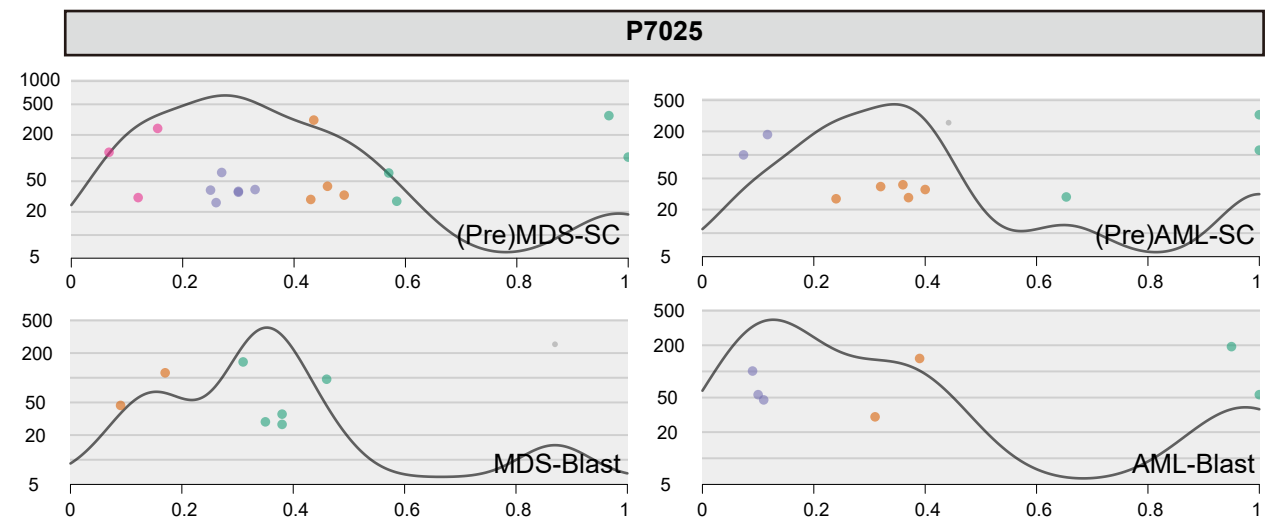
# Figure S7



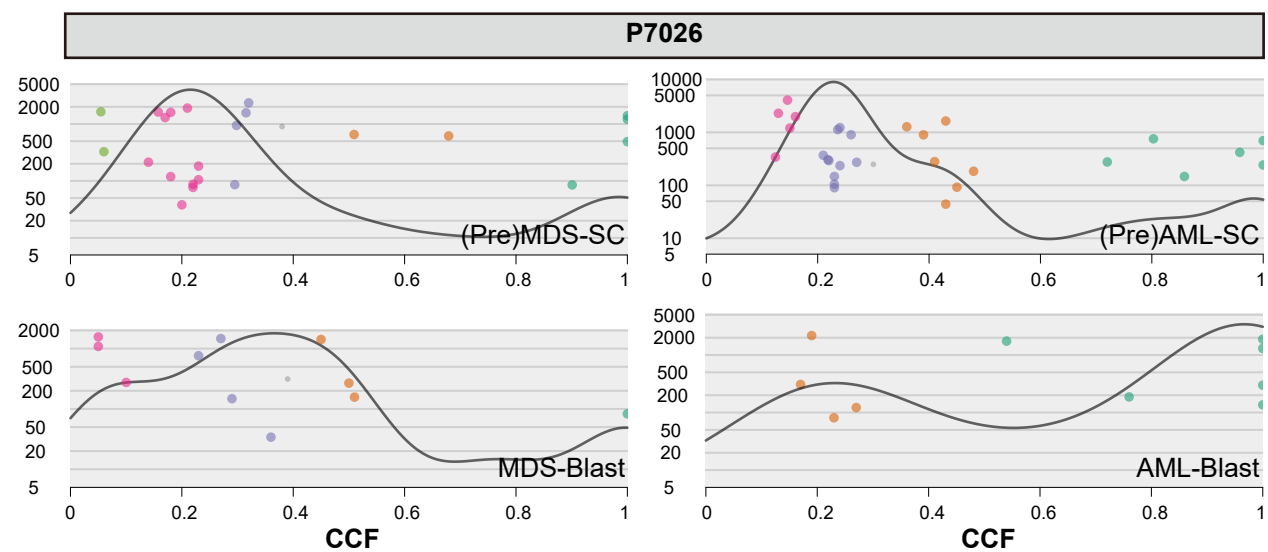
**a**



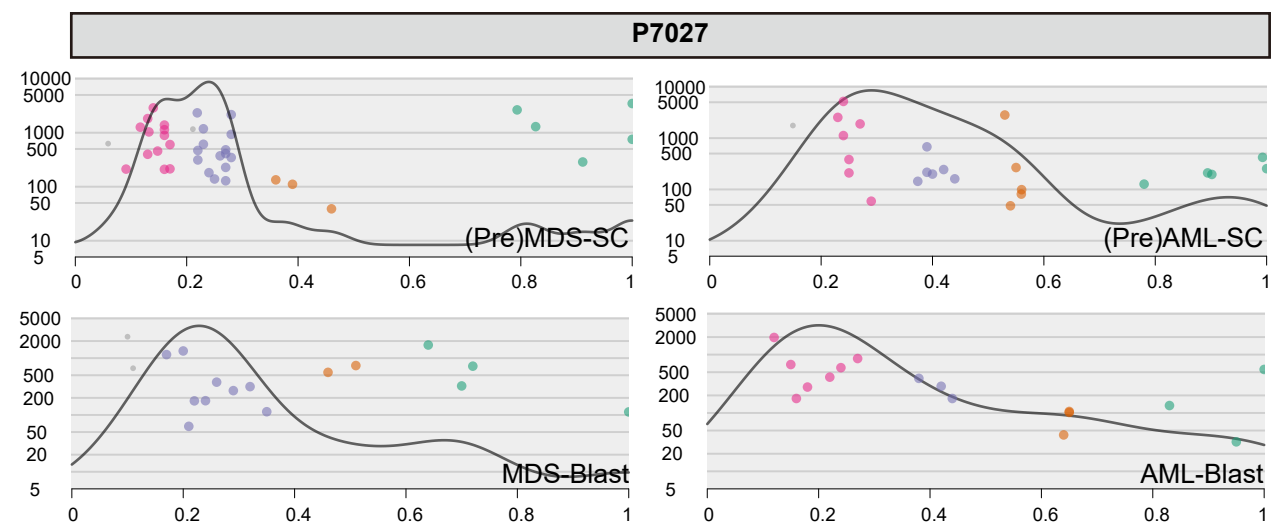
**b**



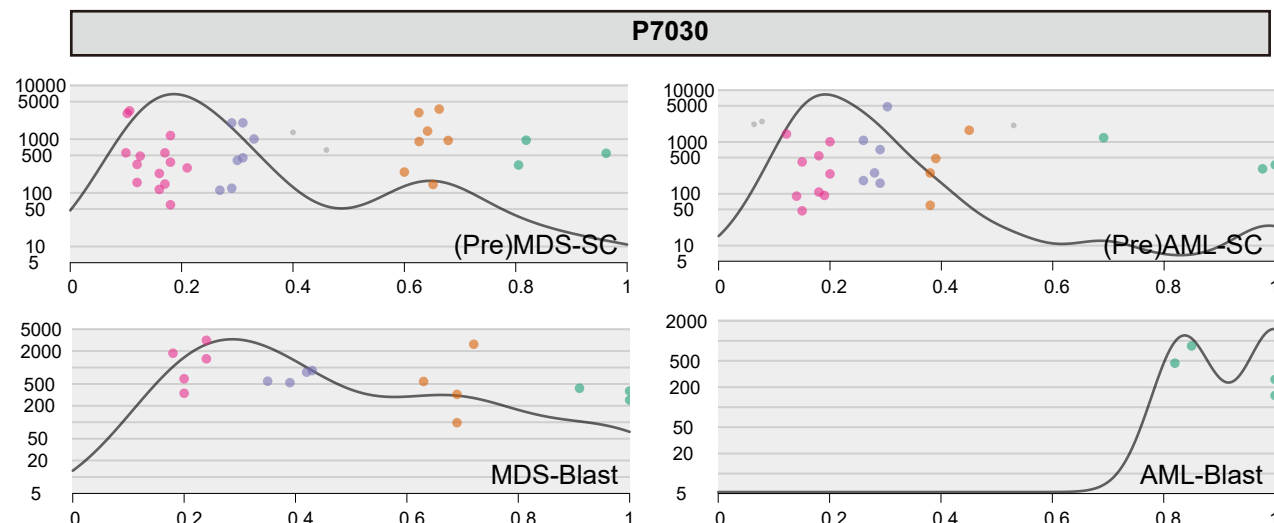
**c**



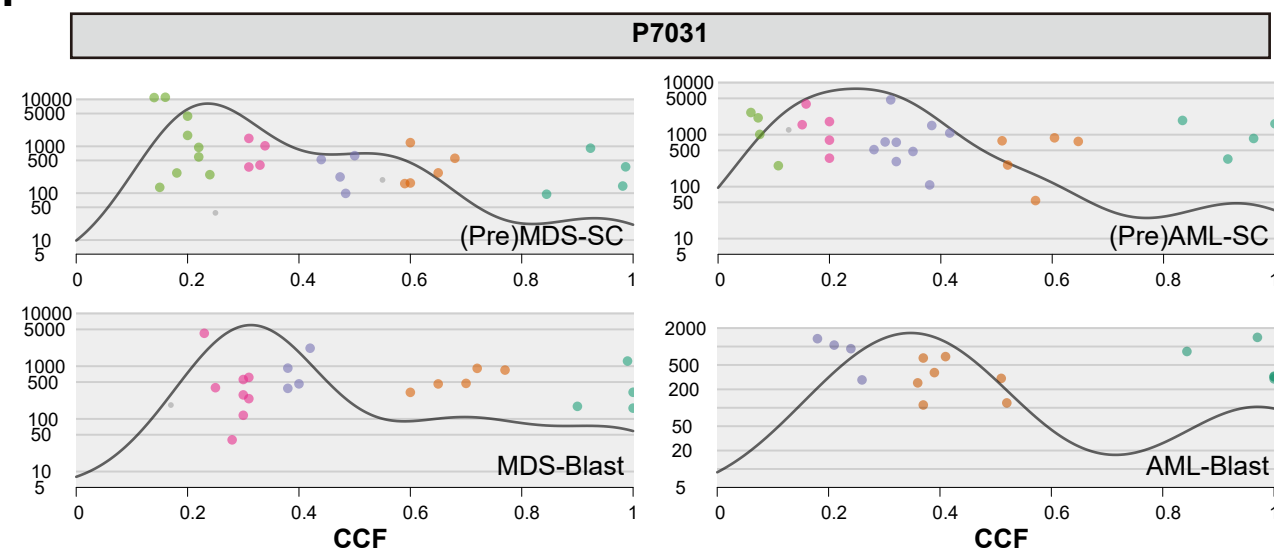
**d**



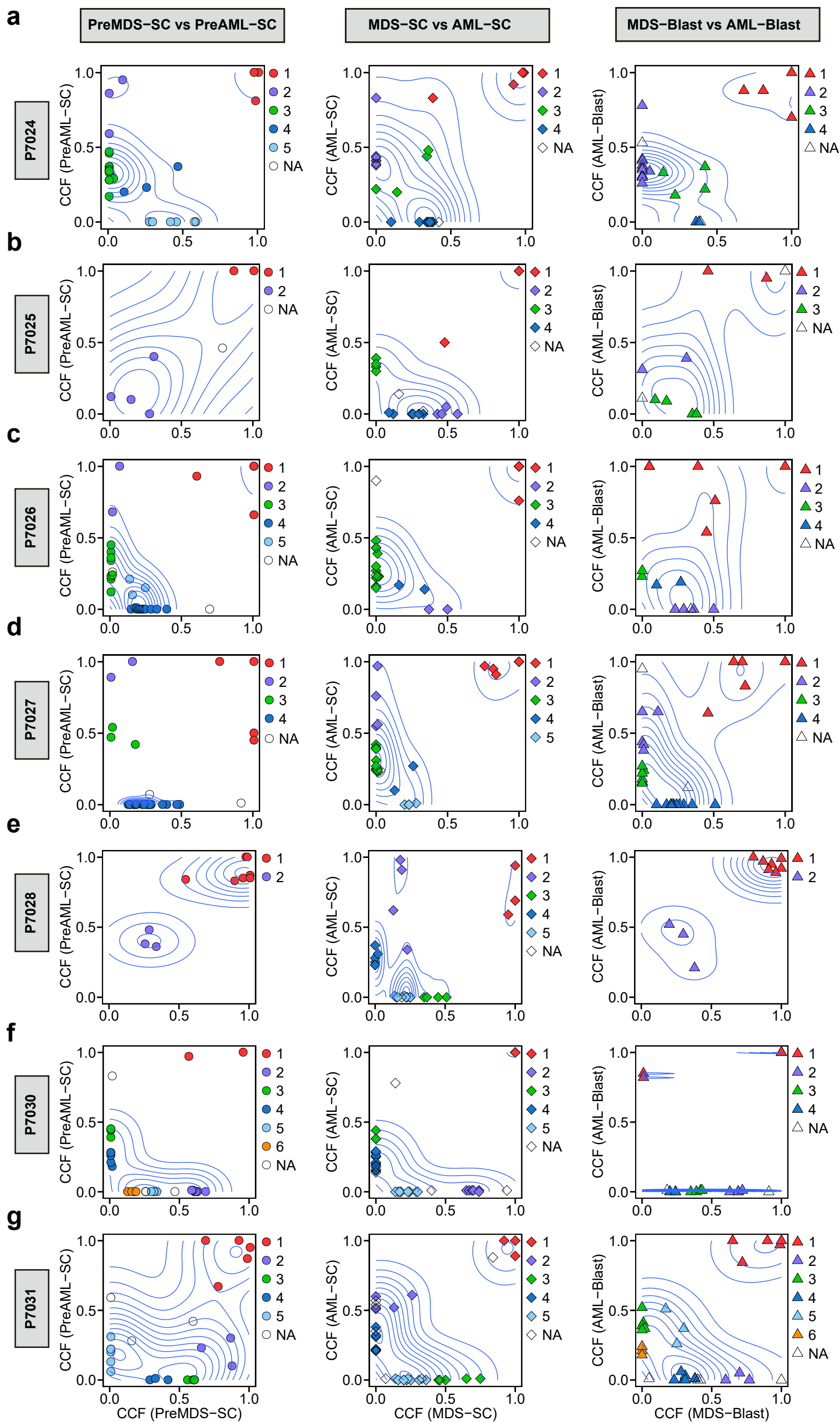
**e**



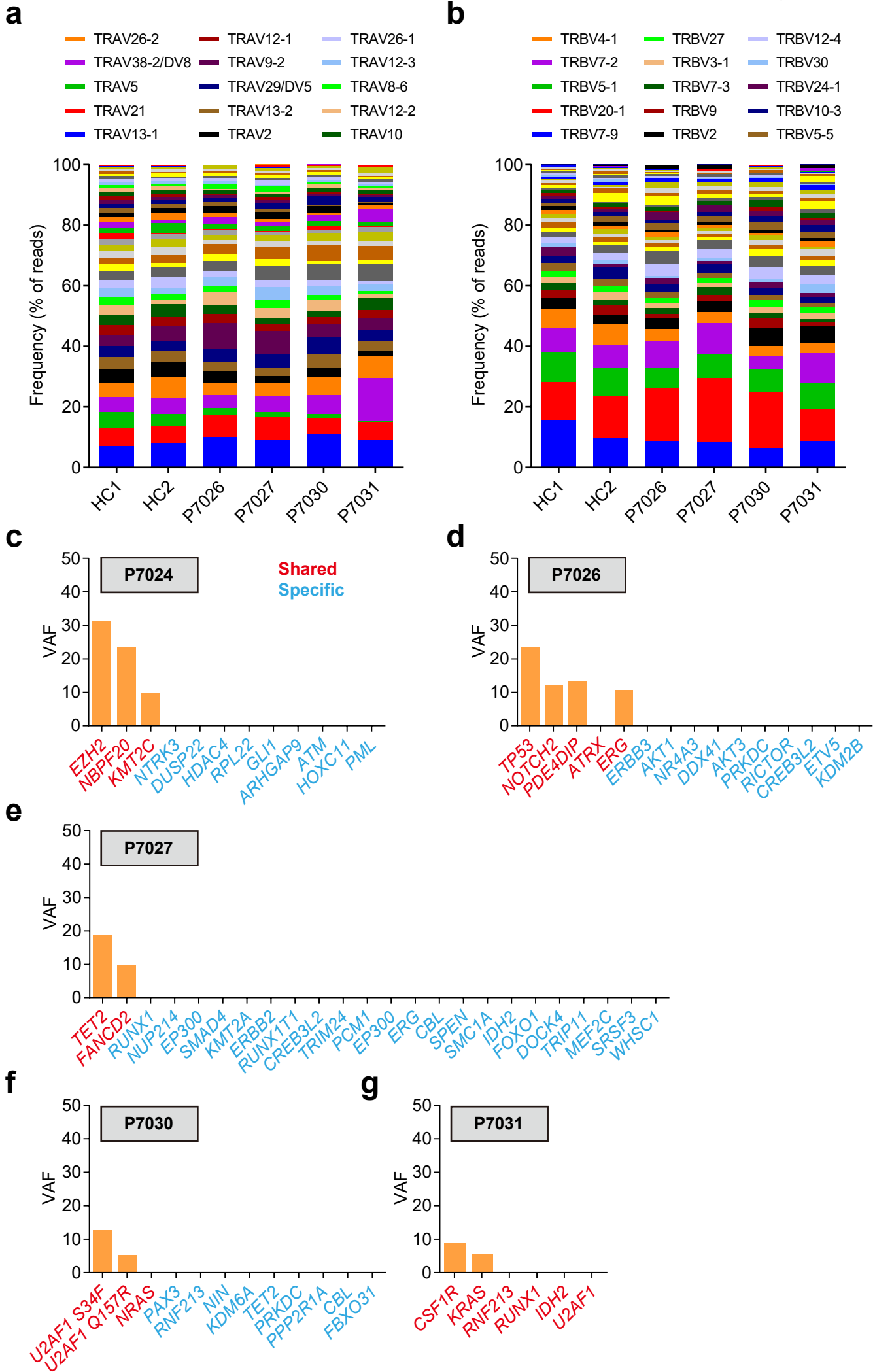
**f**





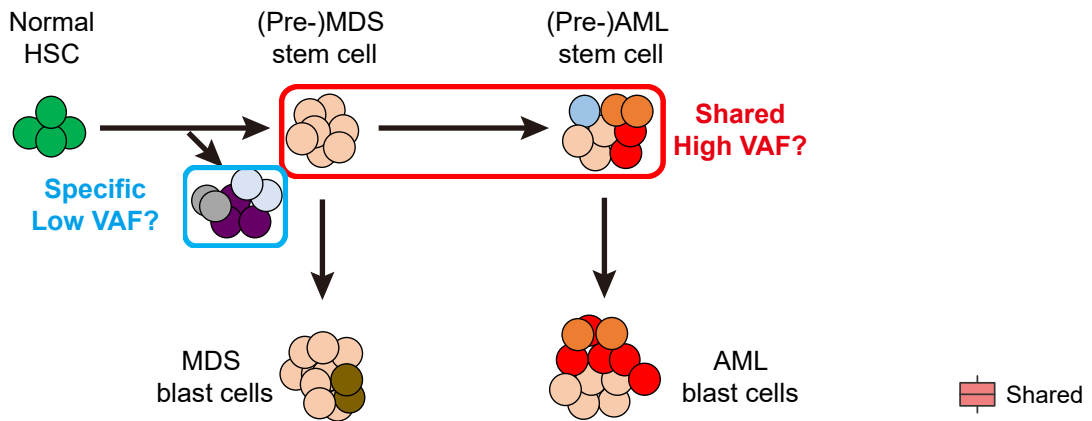


# Figure S10

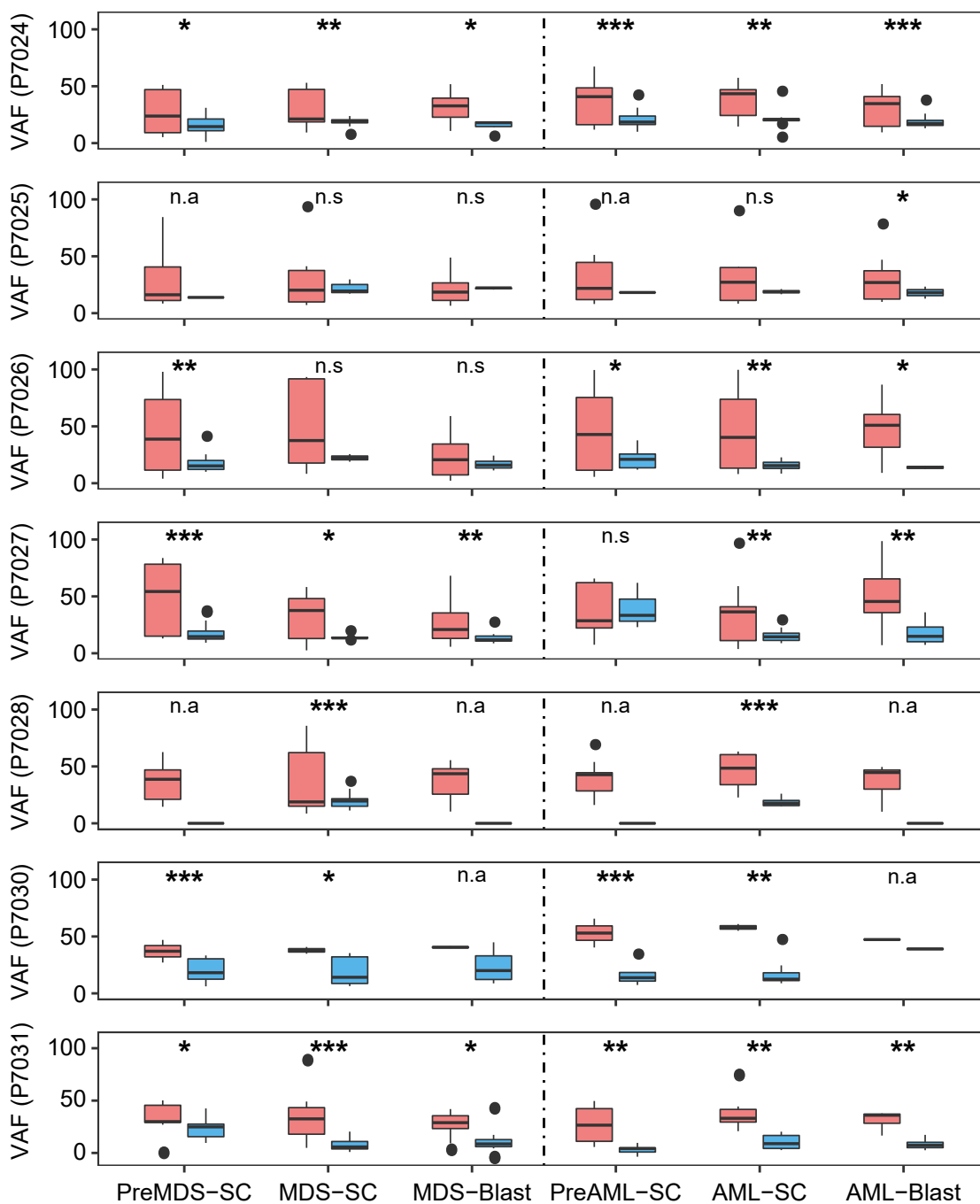


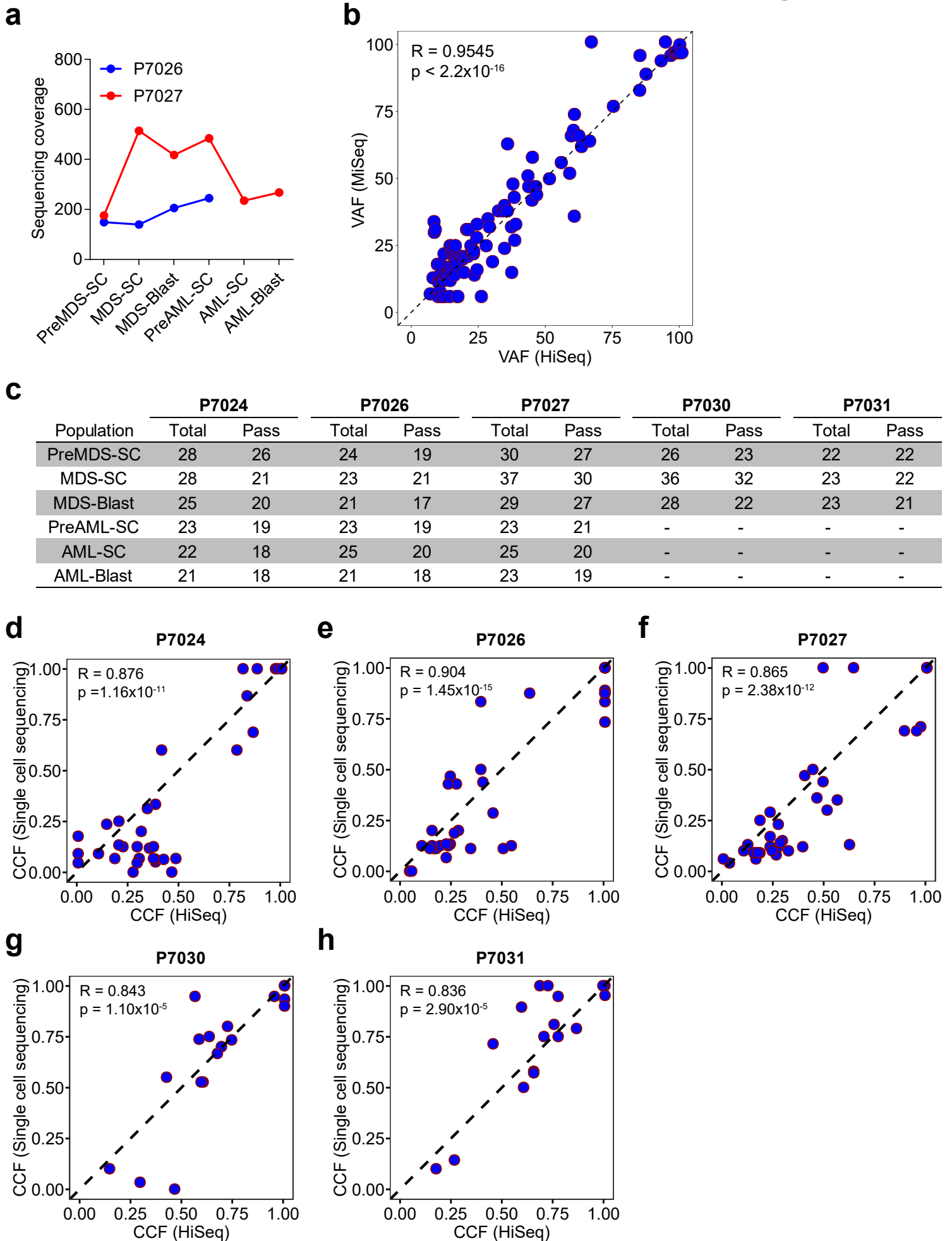
# Figure S11

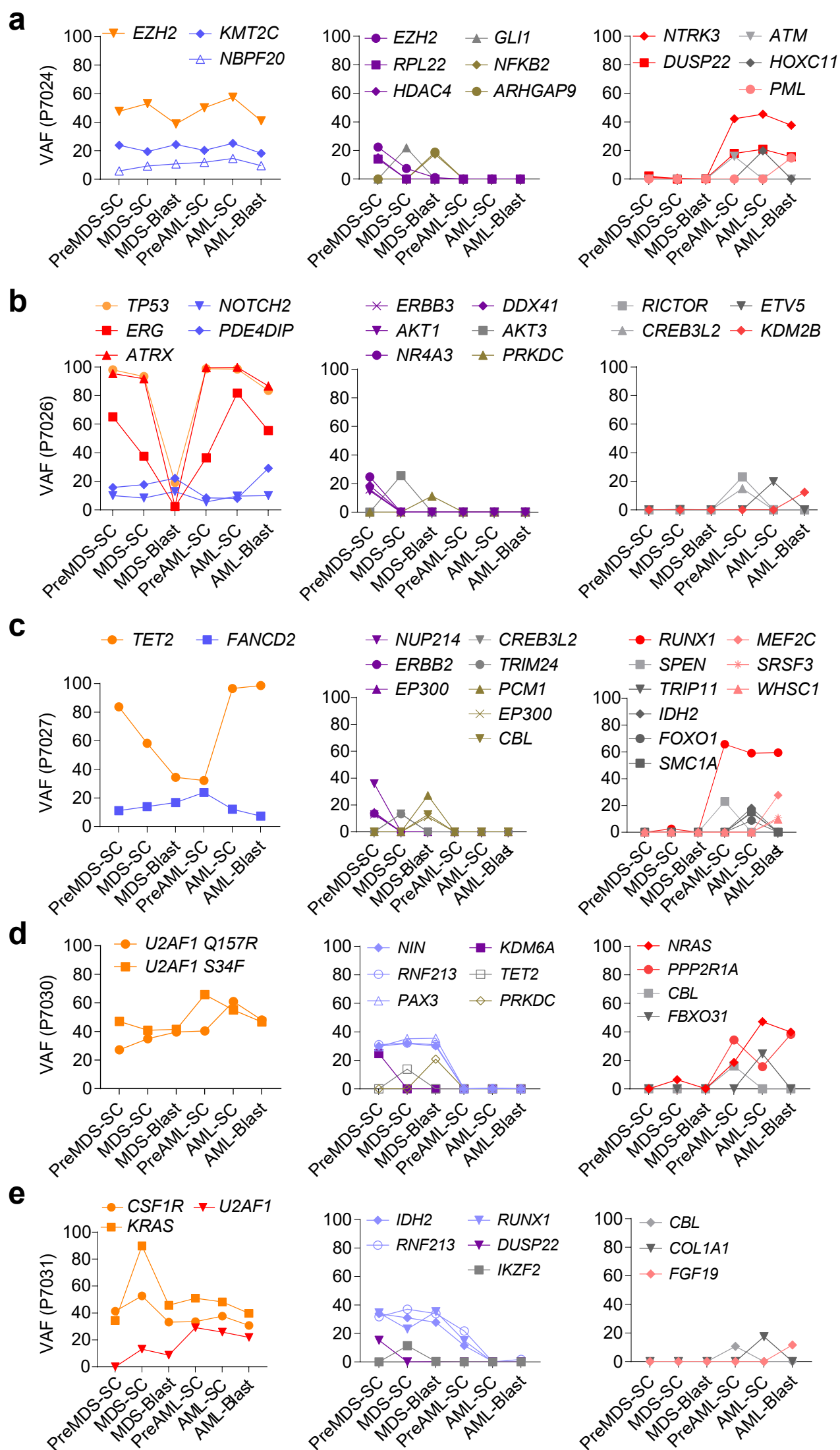
**a**



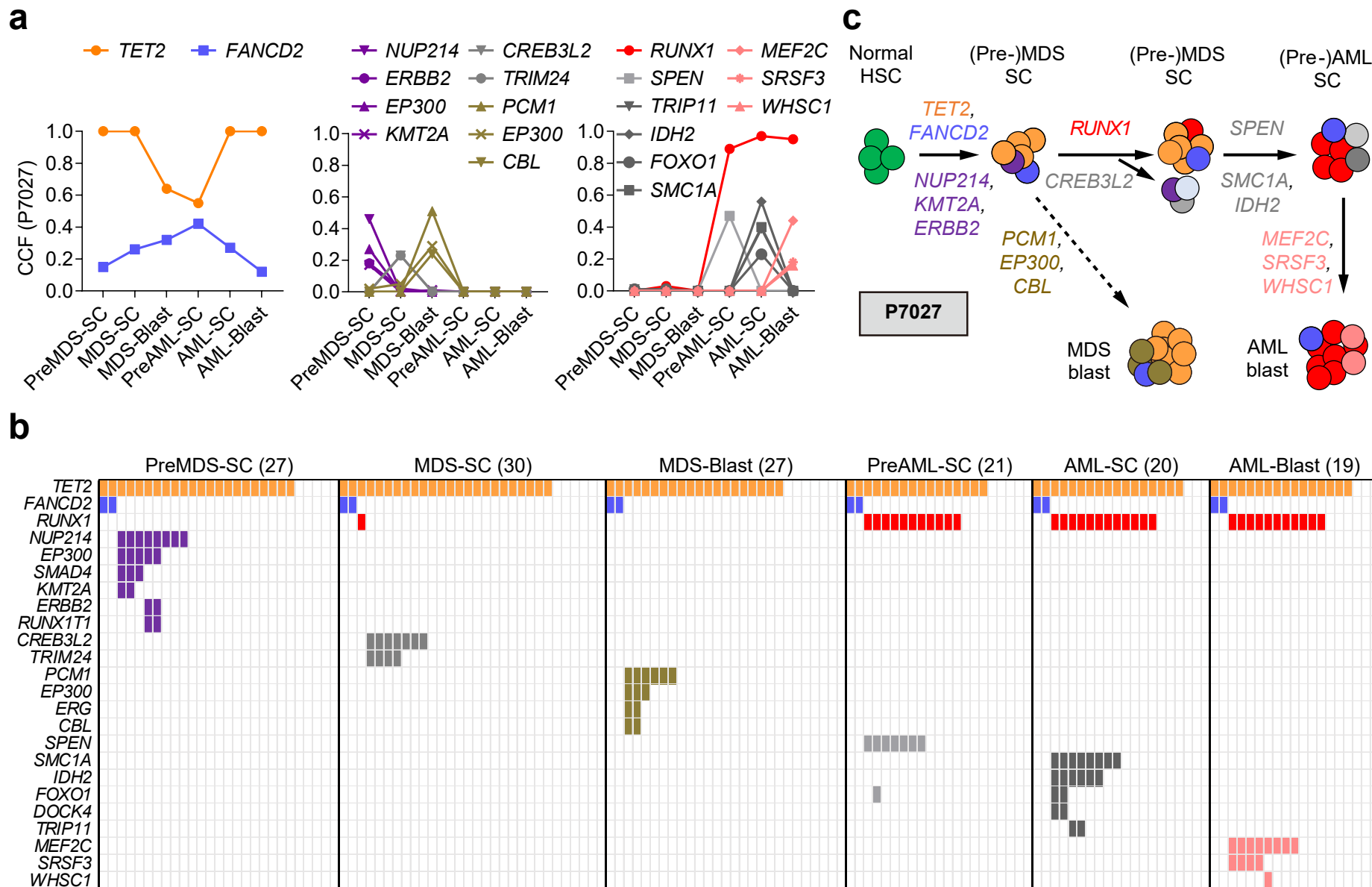
**b**

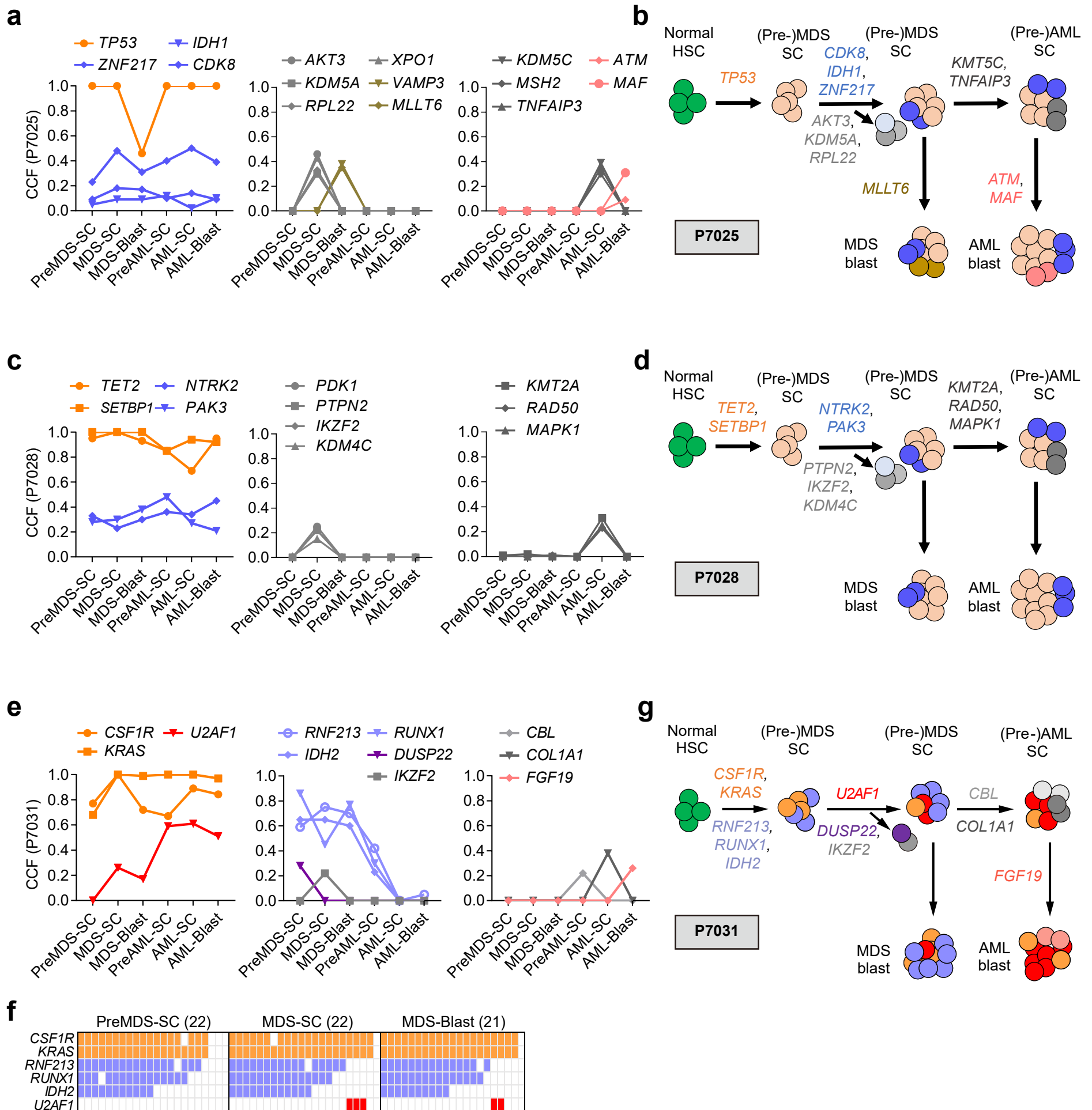






# Figure S14





## Supplementary Figure Legends

**Fig. S1. Isolation of stem cells from individual human patients with MDS and matched sAML for targeted deep sequencing and single cell sequencing.**

**a**, FACS strategy to isolate preMDS-SC and preAML-SC (green; CD45RA<sup>-</sup>CD123<sup>-</sup>IL1RAP<sup>-</sup> HSPC) and MDS-SC and AML-SC (red; CD45RA<sup>+</sup> and/or CD123<sup>+</sup> and/or IL1RAP<sup>+</sup> HSPC) from longitudinal bone marrow samples of individual patients with MDS who progressed to secondary AML (sAML). **b**, Frequency of stem cells expressing aberrant markers at the MDS and matched sAML stages of individual patients (n=8). \*\*\*p < 0.001 (p = 0.0006; two-tailed paired Student's t test). **c**, FACS plots comparing the expression of CD45RA/CD123/IL1RAP markers in MDS and matched sAML samples from each patient.

**Fig. S2. CD45/side scatter gating of blasts and T cells in MDS and matched sAML samples.**

Analysis of CD45 versus side scatter (SSC) profiles of blasts and T cells sorted in the studied MDS and matched sAML samples. CD33<sup>+</sup> blasts (red) and T cells (blue) were gated in a CD45/SSC bi-dimensional plot.

**Fig. S3. Xenotransplantation assays of sorted pre-malignant and malignant stem cells.**

**a**, Schema of xenotransplantation assay. Sorted preMDS-SC and MDS-SC were retro-orbitally transplanted into NBSGW immunocompromised mice (aged 6-8 weeks). PreMDS-SC and MDS-SC were sorted from fresh samples from independent MDS patient based on the same strategy used in the patients sequenced (PreMDS-SC: HSPC (Lin<sup>-</sup>CD34<sup>+</sup>CD38<sup>-</sup>) that were triple-negative for CD45RA, CD123, and IL1RAP; MDS-SC: HSPC expressing at least one (or more) of the markers). Engraftment analysis was performed within the bone marrow of recipient mice 12 weeks after transplantation. **b**, Representative FACS plots of engraftment analysis. The presence of human CD45<sup>+</sup> leukocyte engraftment was analyzed by flow cytometry, and the engrafted cells were further examined for the fraction of CD33<sup>+</sup> myeloid and



CD19<sup>+</sup> lymphoid cells. The number of recipient mice is shown in **c** (each data point represents an individual mouse). **c**, Quantification of engraftment analysis in **b**. **d**, Schema of xenotransplantation assay with sorted preAML-SC and AML-SC. PreAML-SC and AML-SC were sorted from fresh samples from independent AML patients based on the same strategy used in the patients sequenced (PreAML-SC: HSPC (Lin<sup>-</sup>CD34<sup>+</sup>CD38<sup>-</sup>) that were triple-negative for CD45RA, CD123, and IL1RAP; AML-SC: HSPCs expressing at least one (or more) of the markers). Engraftment analysis was performed within the bone marrow of recipient mice 14 or 16 weeks after transplantation. **e**, Representative FACS plots of engraftment analysis. The presence of human CD45<sup>+</sup> engraftment was analyzed by flow cytometry, and the engrafted cells were further examined for the fraction of CD33<sup>+</sup> myeloid and CD19<sup>+</sup> lymphoid cells. The number of recipient mice is shown in **f** (each data point represents an individual mouse). **f**, Quantification of engraftment analysis in **e**. Data are mean  $\pm$  SEM.

**Fig. S4. Methylcellulose differentiation assays of sorted pre-malignant and malignant stem cells.**

**a**, Schema of methylcellulose assay to examine the differentiation of immunophenotypic, sorted pre-malignant and malignant stem cells from independent MDS and AML samples (unpaired). Sorted HSPCs (Lin<sup>-</sup>CD34<sup>+</sup>CD38<sup>-</sup>) that were triple-negative for CD45RA, CD123, and IL1RAP (preMDS-SC and preAML-SC), and HSPCs expressing at least one (or more) of the markers (MDS-SC and AML-SC) were plated in semisolid methylcellulose medium. The colonies of different lineages (E: erythroid; M: macrophage; G: granulocyte; GM and GEMM: mixed lineages) were scored by microscopy 2 weeks after plating. Thereafter, the methylcellulose semisolid medium with colonies was dissolved and resuspend in PBS buffer to make single cell suspensions. Cells were then stained with cell surface markers against CD14, CD15 and CD235a, and analyzed by flow cytometry. **b**, Clonogenicities (per 1000 cells plated) of sorted preMDS-SC and MDS-SC (n=5), as well as preAML-SC and AML-SC (n=4). **c**, Frequency of colonies of different lineages (n=5 for MDS samples; n=4 for AML samples). **d**, Frequency of cells expressing specific lineage markers (n=3 for MDS samples; n=4 for AML samples). If not specified otherwise, data are mean  $\pm$  SEM. \*p < 0.05, \*\*p < 0.01, \*\*\*p < 0.001 (two-tailed paired Student's t test).

**Fig. S5. Targeted capture sequencing of sorted stem cells and blasts from MDS and matched sAML samples.**

**a,b**, Comparison of variant allele frequencies between whole genome amplified samples and unamplified samples of 2 different patients with AML. Prior to sequencing, we performed whole genome amplification (WGA) of sorted cell populations using the REPLI-g method, which utilizes the proofreading enzyme Phi 29 polymerase to achieve high-fidelity amplification of genomic DNA. WGA was performed with 10ng DNA of each patient, thereafter 500ng of purified WGA products were subjected to targeted sequencing along with 500ng of unamplified DNA. After pre-processing of the sequencing data, FreeBayes was used for variant detection (see methods). Then potential germline variants reported in dbSNP database with population frequency >1% were removed, before the comparison of VAFs between WGA sample and matched unamplified sample. 425 and 403 variants were obtained for AML01 and AML02, respectively. And we observed that the VAFs in WGA samples were highly consistent with the unamplified samples ( $p < 2.2 \times 10^{-16}$ ). In addition, we did not observe mutation artifacts that passed the cut-offs for detection of somatic mutations (see methods) when comparing WGA samples to the unamplified samples, indicating that WGA did not distort detection and quantitation of variants. Pearson coefficient R and p-value were calculated by R software. **c**, Mean sequencing coverages in target regions across different cell populations are shown. CD45- non-hematopoietic cells were used as germline control. **d**, Representative results of Sanger sequencing of a *RUNX1* mutation in different cell populations isolated from P7026. The Sanger sequencing was repeated twice. **e**, Correlation of variant frequency estimated by targeted sequencing and Sanger sequencing. A total of 36 pairs of mutations-samples were tested, and Pearson correlation coefficient R and significant p-value calculated by R software are shown. **f-h**, Number of mutations in all genomic regions (**f**), exonic regions (**g**), and non-exonic regions (**h**) across the different cell populations. On average, we observed 8.6 mutations in coding regions (range, 4 to 13) in MDS stem cells compared to 5.3 (range, 3 to 7) in the

MDS blasts ( $p = 0.012$ ); and 7.8 coding mutations (range, 4 to 14) in AML stem cells compared to 5.6 (range, 3 to 9) in AML blasts (**g**). Interestingly, we also observed high numbers of mutations in non-exonic regions, and these mutations in non-exonic regions were significantly more frequent in stem cell populations compared to blasts for both MDS (on average 17.8 vs. 8.1;  $p = 0.0008$ ) and AML (15.0 vs. 6.6;  $p = 0.003$ ) (**h**). \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$  (two-tailed paired Student's t test).

**Fig. S6. Distribution of clonal and subclonal mutations in stem cells and blast cells in each patient.**

**a**, Cancer cell fractions (CCFs) based on somatic mutations detected in targeted sequencing were estimated with variant allelic frequency (VAF), ploidy and purity as previous described (33). Mutations were defined as clonal mutation if the upper border of the 95% confidence interval of the CCF was greater than 0.95. **b-g**, Distribution of CCFs of clonal and subclonal mutations in stem cells and blasts of patients P7024 (**b**), P7025 (**c**), P7027 (**d**), P7028 (**e**), P7030 (**f**), and P7031 (**g**). Violin plot is showing frequency distribution (kernel density) of clonal mutations (orange) and subclonal mutations (grey).

**Fig. S7. Mutation spectrum in stem cells resembles age-associated and DNA repair-associated mutational signatures.**

**a,b**, The 96 trinucleotide mutational spectra of mutations detected in stem cells and blast cells at the stage of MDS (**a**) and sAML (**b**). X-axis is showing the 96 combination of nucleotide changes, and their relative weights inferred by deconstructSigs is shown on the Y-axis. **c**, Relative weight and underlying mechanisms of different mutation signatures inferred from the mutational spectra shown in panel **a, b**.

**Fig. S8. Clonal composition of stem cell and blast populations inferred with CCFs of mutations.**

Clonal composition of stem cell and blast populations in MDS and sAML, respectively, in patients P7024 (**a**), P7025 (**b**), P7026 (**c**), P7027 (**d**), P7030 (**e**), and P7031 (**f**). Based on the CCFs, mutations are clustered as clones with Sciclone and denoted with same color. Mutation was denoted with grey if the estimated possibility of the mutation to be clustered in the subclone was lower than 0.95.

**Fig. S9. Clonal evolution of stem cell and blast populations during progression from MDS to sAML, inferred by CCFs of mutations.**

**a-e**, Comparison of clonal architectures of pre-malignant stem cells (preMDS-SC vs. preAML-SC) (panels on the left), malignant stem cells (MDS-SC vs. AML-SC) (middle panels), and blasts (right panels) in individual patients progressing from MDS to sAML. **(a)** Patient P7024, **(b)** patient P7025, **(c)** patient P7026, **(d)** patient P7027, **(e)** patient P7028, **(f)** patient P7030, **(g)** patient P7031. Each dot represents one somatic mutation (covered by > 20x in both samples shown), and mutations are clustered as clones by SciClone based on the CCFs, indicated by differently colored symbols. Data was visualized with ggplot2 package of R with results from SciClone plus kernel density plot. “NA” indicates unclustered mutations with insufficient sequencing coverage, or mutations with possibility of less than 0.95 in the cluster. For 4 out of 7 patients (P7024, P7025, P7026, P7028) the blasts remained more stable during the MDS-to-sAML transition than the stem cell compartments. In all these patients, the MDS-SC showed massive subclonal differences compared to AML-SC, with subclones that were specific to either the MDS or the sAML stage of the same patient. In two patients (P7024 and P7026) this was additionally accompanied by subclonal differences at the pre-malignant stem cell level (preMDS-SC vs. preAML-SC). Interestingly, in patients P7028 and P7025, preMDS-SC, as well as the MDS blast population showed relatively stable clonal dynamics, whereas MDS-SC underwent extensive evolution during the progression from MDS to sAML. Two patients (P7027 and P7031) showed substantial subclonal evolution in each of the preMDS-SC/preAML-SC, the MDS-SC/AML-SC, and the MDS-blast/AML-blast transitions. Consistent patterns of clonal evolution were observed in the analyses with VAFs of mutations (data not shown).

**Fig. S10. Somatic mutations shared between samples at the MDS and sAML stages are present in T cells.**

**a,b**, Relative utilization of specific V $\alpha$  (**a**) and V $\beta$  (**b**) gene segments in T cells isolated from patients' samples, as well as T cells from two healthy controls (HC1 and HC2). The legends show the top 15 gene segments. We hypothesized that the shared mutations were acquired during early stages of disease initiation, thus shared across all cell populations. To test this, we performed targeted re-sequencing of these mutations in matched T cells isolated from MDS bone marrow samples of each patient (**c-g**), as the majority of T cells develop during childhood long before the diagnosis of MDS or AML. To rule out the possibility that T cells were derived from the expansion of a T cell clone arising from mutant HSCs, we performed T cell receptor (TCR) repertoire analysis of the T cells by sequencing the TCR mRNAs (see methods), and found that the TCR repertoire of the patients was highly diverse with polyclonal patterns indistinguishable from healthy controls. (**c-g**), Selected shared mutations, MDS-specific mutations, and AML-specific mutations were examined in sorted CD4<sup>+</sup> T cells isolated from the MDS bone marrow sample of the same patient with targeted re-sequencing by Fluidigm and MiSeq. (**c**) Patient P7024, (**d**) patient P7026, (**e**) patient P7027, (**f**) patient P7030, (**g**) patient P7031. Variant allele frequencies (VAFs) of mutations are shown. Data show presence in T cells of most shared mutations, but not of MDS- or AML-specific mutations.

**Fig. S11. Shared mutations between MDS and sAML have higher VAF compared to mutations specific to MDS or AML.**

**a**, Schematic non-linear clonal evolution of stem cells in MDS and progression to sAML. Accumulation of mutations gives rise to distinct subclones at the stem cell level, and different subclones contribute to the generation of MDS blasts and sAML progression. Based on this hypothesis, stage-specific (MDS or

AML) mutations would be expected to reside in smaller subset of cells, and thus have lower VAF values. We therefore examined the VAFs of “shared” (present in all sorted cellular subsets) and “specific” (only present in one cellular subset) mutations. **b**, Box plots of VAFs of shared (MDS and AML) or stage-specific (MDS or AML) mutations in different stem and blast cell populations of each individual patient (rows). Results were visualized with the ggplot2 package of R. Lower and upper hinges are showing the 25% and 75% quartiles, respectively. And the segment inside the rectangle shows the median of VAF values. We found that the stage specific mutations consistently had lower VAFs compared to the shared mutations in each patient \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$  (two-tailed unpaired Student’s t test). N.s: non-significant; n.a: p-value not available as there were less than 2 data points (mutations) in the either shared or specific group.

**Fig. S12. Targeted re-sequencing with Fluidigm Biomark HD followed by MiSeq platforms.**

**a**, Somatic mutations detected by HiSeq targeted sequencing were validated by MiSeq following targeted capture with the Fluidigm Biomark platform. We also re-sequenced the original DNA of patients P7026 and P7027 by targeted sequencing. Sequencing coverages of MiSeq data of the tested cell populations from patients P7026 and P7027 are shown. **b**, Comparison of VAFs determined by HiSeq vs MiSeq data. Each dot represents one somatic mutation. X- and Y-axis are showing the variant allele frequency (VAF) of mutations determined from the data of HiSeq and MiSeq, respectively. 82 pairs of mutations-samples were tested, and Pearson correlation coefficient  $R$  and significance  $p$ -value calculated by R software are shown. We found high concordance between the VAFs determined by HiSeq and MiSeq ( $p < 2.2 \times 10^{-16}$ ) **c**, Number of single cell sequencing samples tested and passing the quality control in each population. For single cell sequencing, we examined sorted, immunophenotypic preMDS-SC, MDS-SC, MDS blasts, preAML-SC, AML-SC, and AML blasts of the patients. **d-h**, CCFs of mutations estimated by HiSeq targeted sequencing are validated by single cell sequencing. Correlation between mutation CCFs estimated by HiSeq and CCFs determined by single cell sequencing in patient P7024 (**d**), P7026 (**e**), P7027 (**f**), P7030 (**g**), and P7031 (**h**). Examined mutations in each patient are shown in **Fig. 3**, and only

mutations detected by HiSeq and/or single cell sequencing are shown. CCFs in single cell sequencing studies were calculated as the fraction of single cells with the mutation within all single cells carrying early mutations (e.g. *TP53* and *U2AF1*). We found significant correlation between the CCFs estimated by HiSeq of sorted cell populations and CCFs determined by single cell sequencing in all patients, including P7024 ( $p = 1.16 \times 10^{-11}$ , n=40 pairs of mutations-samples), P7026 ( $p = 1.45 \times 10^{-15}$ , n=40), P7027 ( $p = 2.38 \times 10^{-12}$ , n=44), P7030 ( $p = 1.10 \times 10^{-5}$ , n=18), and P7031 ( $p = 2.90 \times 10^{-5}$ , n=17). Pearson coefficient R and p-value were calculated by *cor.test* function of R.

**Fig. S13. Variant allele frequencies of mutations across different stem and blast cell populations.**

**a-e**, Variant allele frequencies (VAFs) in sorted cell populations of the same mutations tested in single cell sequencing shown in **Fig. 3**. **(a)** Patient P7024, **(b)** patient P7026, **(c)** patient P7027, **(d)** patient P7030, **(e)** patient P7031.

**Fig. S14. Spatiotemporal subclonal evolution during the progression from MDS to sAML in patient P7027, determined by single cell sequencing of sorted stem and blast cells.**

**a**, CCFs of shared (left), MDS-specific (middle), AML-specific (right) mutations in different stem and blast populations at the MDS and sAML stage of patient P7027. **b**, Single cell targeted sequencing of mutations across different cell populations of patient P7027. **c**, Schematic model of clonal evolution in different stem and blast cell populations in patient P7027. Subclones of MDS stem cells with early founding mutations (i.e. *TET2*) remained present during MDS blast generation as well as AML progression, whereas other mutations, e.g. *PCMI* and *EP300*, only occurred in MDS blasts but not during progression to sAML. Interestingly, we had also identified a *RUNX1* mutation that had a low CCF in MDS-SC, but then expanded at the AML stage. Single cell sequencing confirmed that *TET2* was indeed present in almost all the cells across the different populations, whereas mutation of *FANCD2* resided in a subclone within *TET2*-mutated cells **(b)**. Importantly, the mutation of *RUNX1* also resided in a progeny subclone of *TET2* -mutated cells, however, one that was distinct from the subclone with *FANCD2*

mutation (**b**). In addition, within the *TET2*-mutated preMDS-SC, mutations of *NUP214* and *EP300* were acquired successively, however, these never gained dominance within the MDS or AML blast cells (**b**). Within the same parent subclone, *SMAD4* and *KMT2A* mutations were acquired in progeny different from subclones with *ERBB2* and *RUNX1T1* co-mutations (**b**). In this patient, the progression to sAML originated from a subclone of MDS stem cells and was triggered by a *RUNX1* mutation, whereas other mutations, e.g. *PCMI* and *EP300*, only occurred in MDS blasts but did not play a role in progression to sAML.

**Fig. S15. Spatiotemporal subclonal evolution in sorted stem and blast cells during the progression from MDS to sAML of patients P7025, P7028 and P7031.**

**a**, CCFs of shared (left), MDS-specific (middle), AML-specific (right) mutations in different stem and blast populations at the MDS and sAML stage of patient P7025. **b**, Schematic model of clonal evolution in different stem and blast cell populations in patient P7025. **c**, CCFs of shared (left), MDS-specific (middle), AML specific (right) mutations across all cell populations in patient P7028. **d**, Schematic model of clonal evolution in different stem and blast cell populations in patient P7028. Despite the larger number of stable mutations, results of P7025 and P7028 still indicate an overall model of parallel MDS and sAML evolution at the stem cell level (in these cases, with slightly later branching) rather than linear MDS to sAML progression. **e**, CCFs of shared (left), MDS-specific (middle), AML-specific (right) mutations across all cell populations in patient P7031. **f**, Single cell targeted sequencing of mutations across different MDS cells populations of patient P7031. **g**, Schematic model of clonal evolution in different stem and blast cell populations in patient P7031.