# Joint profiling of proteins and DNA in single cells reveals extensive proteogenomic decoupling in leukemia

Benjamin Demaree[1,2,†], Cyrille L. Delley[1,†], Harish N. Vasudevan[1,3], Cheryl A.C. Peretz[4], David Ruff[5], Catherine C. Smith[4], Adam R. Abate[1,2,6,*]

[1]Department of Bioengineering and Therapeutic Sciences, University of California, San Francisco, San Francisco, California, USA
[2]UC Berkeley-UCSF Graduate Program in Bioengineering, University of California, San Francisco, California, USA
[3]Department of Radiation Oncology, University of California, San Francisco, California, USA
[4]Division of Hematology/Oncology, Department of Medicine, University of California, San Francisco, California, USA
[5]Mission Bio, Inc., San Francisco, California, USA
[6]Chan Zuckerberg Biohub, San Francisco, California, USA

*Corresponding author.
†These authors contributed equally to this work.

Correspondence to:
Adam R. Abate, Ph.D.
Department of Bioengineering and Therapeutic Sciences, University of California, San Francisco, 1700 4th St, San Francisco, California, USA; Email: adam@abatelab.org

## Summary

Current leukemia therapies target cancer cells with specific phenotypes or genotypes, but this assumes that either genomic mutations or immunophenotypes alone serve as faithful proxies for treatment response[1]. Moreover, the heterogeneity inherent to all cancers, including leukemias, makes direct mapping of genotype-phenotype relationships challenging[2,3]. Here, we present a method to genotype and phenotype single cells at high throughput, allowing direct characterization of proteogenomic states on tens of thousands of cancer cells rapidly and cost efficiently. Using the approach, we analyze the disease of three leukemia patients over multiple treatment timepoints and recurrences. We observe complex genotype-phenotype dynamics and extensive decoupling of the relationships over disease progression and response to therapy, illustrating the subtlety of the disease process and the inability to use genotypes as direct proxies for phenotypes. Our technology has enabled the first rigorous test of the prevailing paradigm that treatment of a disease phenotype is equivalent to treatment of its underlying genotype. More broadly, our results highlight the power of single-cell multiomic measurements to resolve complex biology in heterogeneous populations, and illustrate how this information can be used to inform treatment. We thus expect that our methodology will find broad application to study proteogenomic tumor landscapes across cancers and will support the next generation of immunotherapy.

## Main

Acute myeloid leukemia (AML) is an aggressive hematologic malignancy prone to relapse that often manifests as a polyclonal ensemble of cells with distinctive genotypes but diverse immunophenotypes[4,5]. Because of this disparity, it is difficult to directly link genotypes to immunophenotypes beyond circumstantial evidence from epidemiologic studies. Moreover, while AML blasts often exhibit immunophenotypes distinct from normal cells, with some surface markers even serving as therapeutic targets[6], genotypes are the strongest prognostic factors, suggesting a weak correspondence between these domains[7,8]. Cellular heterogeneity is an intrinsic aspect of essentially all cancers, including leukemias. Because cancer cells are heterogeneous in genotype and phenotype, single-cell analysis provides a powerful tool for characterizing this complexity and thereby advancing our understanding of different cancers. The value of single-cell analysis is its ability to correlate co-occurrence of different features in individual cells, with high-throughput technologies permitting analysis of thousands of cells to generate rich and intricate feature maps. For example, single-cell genotyping of AML-relevant loci has revealed co-occurrence of mutations and mapping of the clonal relationships between blasts[9–12]. These studies, however, have yet to map DNA genotypes and phenotypes in the same cells, precluding direct linkage of phenotypes to the genetic mutations that drive them.

To obtain simultaneous genotype and immunophenotype information, single cells can be sorted based on multi-parametric antibody analysis, and sequenced. While severely limited in throughput, these studies have uncovered important insights into the genetics of AML, identifying relevant aberrations such as single nucleotide polymorphisms (SNPs) and gene fusions[13]. Single-cell RNA sequencing (scRNA-seq) has emerged as a potentially valuable approach for genotype-phenotype linkage because it is cost effective and scalable[3,10,14,15]. The mRNA sequences provide genotype information[15,16] while their counts relate phenotype[17–21]. Moreover, modern approaches are extremely high throughput, allowing characterization of thousands of cells. Nevertheless,

48 genotyping from mRNA remains a challenging and error-prone procedure that, even in the best
49 case, provides incomplete information. For example, stochastic gene expression, biological
50 biases[22], and limited coverage of essential genes combine to make assigning a genotype more
51 difficult than can be achieved by direct analysis of DNA. Moreover, since RNA methods analyze
52 only the expressed portion of the genome, mutations in intronic and other non-transcribed
53 elements, like transcription factor binding sites, are omitted[23,24]. Thus, while several technologies
54 have highlighted the importance of high-throughput single cell genotype-phenotype
55 measurements, none provide the scalability and precision for comprehensive and accurate mapping
56 of these clinically valuable biomarkers.
57
58 In this paper, we describe DAb-seq, a novel approach for joint profiling of DNA and surface
59 proteins in single cells at high throughput. While existing methods attempt to obtain this
60 information from the transcriptome alone, ours directly characterizes DNA for genotype and
61 surface proteins for phenotype – both the gold standards in AML for these annotations. Our
62 approach is thus complementary to scRNA-seq methods and, as we show, provides novel and
63 important information for characterizing the disease. To illustrate the power of DAb-seq, we
64 characterize the immunophenotypic and genotypic diversity underpinning AML in three patients
65 at multiple timepoints, exploiting its throughput to characterize 50 DNA targets and 23
66 hematopoietic markers in a total of 54,717 cells. This analysis allows tracking of proteogenomic
67 dynamics for multiple patients over multiple treatments and recurrences. We identify extensive
68 genotype-phenotype decoupling, observing immunophenotypic heterogeneity among cells with a
69 shared pathogenic mutation and genotypically diverse cells with a convergent malignant
70 immunophenotype. These findings indicate substantial variability of blast fate upon treatment in
71 AML, and that independent phenotype or genotype measurements do not adequately capture the
72 proteogenomic heterogeneity. More broadly, our work demonstrates how single-cell technologies
73 can inform the diagnosis and treatment of AML by elucidating the complex interplay between
74 DNA mutations and their effects on protein expression.
75
76
77 **Results**
78
79 *Combined single-cell DNA sequencing and antibody profiling (DAb-seq) robustly delineates*
80 *single-cell genotypes and immunophenotypic diversity*
81
82 The commercially available Mission Bio Tapestri supports highly multiplexed targeted sequencing
83 of thousands of single cells and is being used across cancers for genotype and lineage mapping[11].
84 While the instrument runs a flexible workflow, it does not natively support Abseq, a separate
85 method we developed[25] that allows characterization of single-cell surface proteins by sequencing,
86 and is analogous to flow cytometry in its ability to provide immunophenotype information. Thus,
87 our major technical innovation is to adapt Tapestri to enable simultaneous DNA and Abseq
88 analysis. As in our published Abseq approach, DAb-seq begins with immunostaining of a cell
89 suspension using a mixture of antibody-oligo conjugates (Figure 1A). Each antibody is associated
90 with a known oligo tag; thus, when cells are stained with a pool of tagged antibodies, each cell is
91 bound with a combination of antibodies and their tags based on surface protein profile. To
92 characterize the profile, the tags must be sequenced and counted which, in flow cytometry, is
93 analogous to measuring fluorescence of the dyes associated with each antibody, except that photon
94 counting is replaced with tag counting.

95   The stained cells are processed through a modified Tapestri workflow to amplify and barcode
96   genomic targets and surface-bound antibody tags. The workflow follows a two-step protocol to
97   lyse cells and digest chromatin, making the genome accessible to amplification; the droplets are
98   then subjected to a multiplex PCR to simultaneously amplify the genomic targets and capture
99   antibody tags, labeling them with a droplet barcode relating sequences from the same cell (Figure
100  1B). For genotype, we target recurrently mutated genomic DNA loci in AML with primers
101  containing a unique cell barcode against 50 amplicons spanning 19 genes. The primers and PCR
102  conditions are tuned to enable uniform and quantitative amplification of all targets, since count
103  information is necessary for accurate genotype and immunophenotype characterization. These
104  primers also capture antibody tags from a 23-plex immunophenotyping panel based on those used
105  in clinical minimal residual disease studies[26,27] (Figure 1C; Supplementary Table 3). Sequencing
106  yields a multiomic data set where each cell is represented by two vectors and which can be
107  visualized as a low-dimensional embedding (Figure 1D).
108
109  Peripheral blood mononuclear cells (PBMCs) comprise a diverse and well-understood population
110  easily obtained from a blood draw, and thus provide an excellent sample by which to assess the
111  effectiveness of DAb-seq for mapping hematopoietic immunophenotypes. When applied to
112  PBMCs from a healthy donor, we obtain expected cell subsets across blood compartments,
113  identifying both rare and abundant cells in peripheral blood (Figures 2A, 2B). To test single-cell
114  genotyping capability, we also perform DAb-seq on a mixture of three cell lines derived from
115  distinct hematopoietic lineages (Jurkat, Raji, K562) with documented mutations in the targeted
116  genomic regions covered by our single-cell DNA sequencing panel[28] (Supplemental Table 2). For
117  all genetic variants, we assign genotype calls to each individual cell: homozygous wildtype,
118  heterozygous alternate, or homozygous alternate. We observe the expected correspondence
119  between single-cell genotypes and phenotypes, as cells of the same genotype segregate within a
120  common immunophenotypic cluster (Figures 2C, 2D). Notably, we also find that DAb-seq's
121  genotyping is sufficiently sensitive to differentiate the cells based on zygosity of a given mutation
122  (Figure 2D). These results show that DAb-seq can simultaneously profile genotype from direct
123  analysis of genomic DNA and immunophenotype from barcoded antibodies.
124
125  *NPM1-mutated cells persist across therapy timepoints with a static immunophenotype*
126
127  AML therapies targeted to cell surface proteins require ubiquitous expression of the target marker
128  on the malignant cells. We therefore reason that mutated cells should robustly associate with a
129  common targeted phenotype in patients responsive to this therapy. To investigate this, we perform
130  DAb-seq on 21,952 total cells from bone marrow aspirates of a patient with AML receiving
131  gemtuzumab, a CD33-targeted therapy, across four treatment timepoints (Figure 3A). This patient
132  received multiple rounds of chemotherapy, including a stem cell transplantation, prior to the first
133  timepoint sampled in this study (Supplementary Table 1). In the single-cell DNA genotyping data,
134  we identify a recurrent frameshift mutation in the *NPM1* gene (*NPM1^{mut}*) across relapse, salvage
135  therapy, and progression timepoints. In addition, the *NPM1* mutation is found to always co-occur
136  with a mutation at the *DNMT3A* locus (Figure 3A). Gemtuzumab targets CD33$^+$ cells, which are
137  extinguished at the remission timepoint[29]. To examine the immunophenotypic profile of the
138  *NPM1^{mut}* blast population, we plot single-cell CD33 and CD34 values with *NPM1* mutation status
139  across timepoints (Figure 3B). The proportion of *NPM1^{mut}* cells in the CD34$^-$ and CD34$^+$
140  compartments does not vary extensively across treatments, suggesting the lack of a therapeutic

141  response in the blast immunophenotype. $CD33^+$ myeloid cells targeted by the drug are absent at
142  remission.

144  In all timepoints for this patient, our analysis suggests an equivalence between the dominant blast
145  genotype and corresponding phenotype. To further explore this relationship between genotype and
146  phenotype, we visualize the high-dimensional single-cell immunophenotype as a Uniform
147  Manifold Approximation and Projection[30] (UMAP) embedding of the antibody data (Figure 3C).
148  Cells within single immunophenotypic clusters originate from different timepoints, highlighting
149  the stability of normal and malignant immunophenotypes over time. When we overlay *NPM1*
150  genotype on the immunophenotype UMAP space, we find a clear association between a single
151  malignant immunophenotype composed of $CD33^+$ cells with *NPM1* mutation status, with variable
152  expression of CD34, CD38, and CD117 in this population (Figure 3D). Indeed, this is in agreement
153  with previous observations in flow cytometric studies where blast cells have been found to
154  uniformly express CD33 and variably express CD34, CD38, and CD117[31]. Among the *NPM1$^{wt}$*
155  cells, we identify classical blood cell markers including CD3 and CD5 (lymphocyte), CD15
156  (monocyte), and CD56 (natural killer). Taken together, in this patient, DAb-seq confirms
157  elimination of $CD33^+$ cells by gemtuzumab treatment and reveals a strong correspondence
158  between genotype and phenotype across timepoints.

160  *Genotypic subclones form overlapping subsets across an immunophenotypic continuum*

162  To investigate whether such tight genotype-phenotype association is a universal feature of AML,
163  we apply DAb-seq to a pediatric patient who underwent induction and consolidation
164  chemotherapy, but ultimately relapsed (Supplementary Table 1). We identify two mutually
165  exclusive *KRAS* and *FLT3*-mutated clones at diagnosis and relapse (*KRAS$^{mut}$*, *FLT3$^{mut}$*). The
166  *FLT3$^{mut}$* population, although the minor subclone at diagnosis comprising just 43 of 4,563 cells
167  (0.94%) compared to 1,539 cells (33.7%) for the *KRAS$^{mut}$* variant, dominates at relapse (6,800 of
168  7,516 cells, 90.5%) (Figure 4A). Immunophenotypically, we also identify a third subset
169  comprising *KRAS$^{WT}$*/*FLT3$^{WT}$* cells expressing a blast-like $CD33^+CD38^+$ immunophenotype with
170  no identifiable DNA mutations in the targeted loci. When we group cells from all timepoints by
171  genotype, pathogenic blasts display variable patterns in immunophenotype, with no clear mapping
172  between the two (Figure 4B).

174  In the absence of an obvious genotype-phenotype mapping for this patient, we sought to investigate
175  the underlying relationship between these domains. Using UMAP, we project the antibody data
176  into two dimensions, coloring the points according to genotype (Figure 4C). We observe a single
177  immunophenotypic compartment with incomplete separation between genotypes. To estimate
178  antibody profile expression within the blast compartment continuum, we identify the dominant
179  gradient in the phenotypic space, ordering all points along the gradient. We then calculate the local
180  average antibody and genotypic composition for neighboring cells (Figure 4C, D) (Methods). As
181  expected, many markers are anticorrelated (CD11b, CD33, CD56) or correlated (CD15) with the
182  principal immunophenotypic gradient. Less trivially, genotypic compositions vary along the
183  gradient, with *KRAS$^{mut}$* clone frequencies anticorrelated and *FLT3$^{mut}$* correlated (Figure 4D).
184  Nevertheless, genotype composition never completely separates into individual clonal
185  populations, making it impossible to define distinct genotype-phenotype clusters; consequently,

186 technologies profiling one modality, such as genotyping or immunophenotyping, cannot
187 adequately capture the heterogeneity inherent to this case of AML.
188
189 *FLT3 inhibitor therapy induces erythroid differentiation in a case of AML*
190
191 Our first two cases feature either a strong genotype-phenotype correlation (Patient 1) or mixed
192 genotyping comprising a single immunophenotype (Patient 2). Thus, for our final case, we analyze
193 a patient treated with gilteritinib, a FLT3 inhibitor therapy reported to promote *in vivo*
194 differentiation of myeloid blasts. This treatment is thought to disperse distinct genotypes into
195 multiple immunophenotypes, although the terminal lineage of the cells remains poorly
196 understood[32–34]. Accordingly, we hypothesize DAb-seq should allow tracking of
197 immunophenotypic dispersal and confirmation of their terminal hematopoietic lineage. We
198 analyze 18,287 cells across treatment timepoints, beginning at diagnosis, discovering a subclone
199 with co-mutated *DNMT3A* and *NPM1* (Figure 5A; Supplementary Table 1). Following
200 cytarabine/daunorubicin induction therapy, a fraction of *DNMT3A^{mut}* cells remain at remission. At
201 relapse and after treatment with the FLT3 inhibitor gilteritinib ("FLT3 Inhibitor"), most cells
202 contain a 24-bp *FLT3* internal tandem duplication (ITD), in addition to the initial *DNMT3A* and
203 *NPM1* mutations. The genotypic structure inferred from the single-cell data indicates a linear,
204 branching hierarchy of sequentially acquired mutations in response to therapy. To explore the
205 immunophenotypic features of this patient's disease, we integrate cells from all timepoints and
206 construct a UMAP representation using the antibody data (Figure 5B). We cluster this data using
207 the Leiden method for cluster detection, an improved algorithm over Louvain modularity[35,36], and
208 manually annotate with phenotypic labels corresponding to hematopoietic lineage from the
209 antibody data (Figure 5C). We identify three blast populations expressing high levels of CD33 and
210 CD38, a monocytic population expressing CD15 and CD16, and erythroid and lymphoid clusters
211 with elevated CD71 and CD3. As expected, samples across treatment timepoints comprise a
212 mixture of immunophenotypically normal and blast-like cells.
213
214 Hypothesizing that different therapies should yield different genotype-phenotype coupling
215 patterns, we sought to characterize how mutated and normal cells distribute across
216 immunophenotypic clusters. For each timepoint, we thus label cells in UMAP space according to
217 DNA genotype and generate density distributions of CD33 signal, a pan-myeloid marker (Figure
218 5D). We also evaluate counts of phenotype cluster membership in each timepoint, subdivided by
219 DNA genotype. At diagnosis, cells mutated at both the *DNMT3A* and *NPM1* locus reside primarily
220 in the Blast 1 cluster (81.8% of *DNMT3A^{mut}/NPM1^{mut}* cells) and express high levels of CD33. A
221 secondary clone mutated exclusively at the *DNMT3A* locus exhibits comparable CD33 expression
222 and resides mainly in the Blast 1 and monocytic clusters (62.5% and 27.7% of *DNMT3A^{mut}* cells,
223 respectively). At remission, the same *DNMT3A^{mut}* clone is identified but with decreased CD33
224 expression and a primarily monocytic immunophenotype (92.7% of *DNMT3A^{mut}* cells) co-
225 localizes with cells of normal genotype, consistent with clonal hematopoiesis of a pre-leukemic
226 clone[37,38]. A newly acquired *FLT3*-ITD clone emerges in high numbers at relapse (99.8% of
227 genotyped cells), coinciding with a phenotypic shift of cells to the CD33^+ Blast 2 cluster.
228 Following FLT3 inhibitor treatment, the same *FLT3*-ITD clone persists but exhibits a transformed
229 immunophenotype, as evidenced by membership of the *FLT3* clone in multiple immunophenotypic
230 clusters. The new *FLT3*-ITD immunophenotype is primarily erythroid (82.2% of *FLT3*-ITD cells),
231 with minor fractions in the Blast 3 and monocytic compartments (11.1% and 4.84% of *FLT3*-ITD

232    cells, respectively). Furthermore, the *FLT3*-ITD clone at relapse lacks uniform CD33 expression,
233    indicating that this clone is no longer restricted to the myeloid compartment. Taken together, these
234    findings support the model of terminal erythroid differentiation of blasts in a case of leukemia
235    treated with gilteritinib. In agreement with a recent study[34], proteogenomic analysis by DAb-seq
236    challenges a prior report of gilteritinib-induced terminal differentiation towards a myeloid fate[33].
237    DAb-seq elucidates the rich and complex dynamics of this process and illustrates how distinct
238    DNA genotypes can fractionate into multiple phenotypic identities in response to treatment.
239
240
241    **Discussion**
242
243    Through its ability to jointly profile DNA and immunophenotype, DAb-seq captures the
244    complexity of proteogenomic states underlying AML. Analysis of multiple patients over
245    timepoints and treatments demonstrates the plasticity of the disease and the complex and
246    unpredictable way it progresses in different contexts. In a patient with extensive clinical history
247    including multiple rounds of chemotherapy, we found a robust relationship between mutant *NPM1*
248    cells and a malignant phenotype; this suggested that a single CD33-targeted therapy would
249    eradicate the blast population, as indeed it did. By contrast, in a separate case of pediatric AML,
250    we observed that genetically distinct populations shared overlapping immunophenotype,
251    demonstrating that this domain alone is insufficient for characterizing how cells are genetically
252    programmed and may, consequently, respond to treatment. In the final case study, we observed
253    the opposite scenario, in which treatment by gilteritinib induced mutationally similar cells to
254    disperse into different myeloid compartments, highlighting the challenge of targeting these
255    malignant cells for eradication. Our results thus demonstrate that genotype or immunophenotype
256    alone is insufficient to predict the evolution of proteogenomic states in AML.
257
258    DAb-seq employs targeted primers to amplify specific genomic regions and panels of antibodies.
259    While both readouts enable massive multiplexing of queried targets, practical and economic
260    constraints necessitate *a priori* knowledge of which loci and epitopes to profile. As such, the
261    strength of DAb-seq is not unbiased feature discovery, as with scRNA-seq, but rather sensitive
262    and precision analysis of actionable information. Furthermore, as with all targeted methods of
263    DNA genotyping, DAb-seq cannot exclude the possibility that disease-relevant mutations occur
264    beyond the sequenced loci or in immunophenotypic markers not included in the panels. In the case
265    of pediatric AML, it is therefore impossible for us to conclude if the $FLT3^{wt}/KRAS^{wt}$ blast
266    population is driven by epigenetic changes or unmapped genomic aberrations. Nevertheless, the
267    sensitivity of DAb-seq, and its low genotyping drop-out, allows identification of co-occurring
268    mutations, including heterozygous mutations that are notoriously difficult for RNA-based
269    approaches. Moreover, DAb-seq firmly places genomic mutations in understood phenotypic
270    contexts, which is vital for understanding how they program the disease and, ultimately, treatments
271    select for them.
272
273    In the era of personalized medicine, treatment decisions are increasingly based on DNA mutation
274    status, such as targeted EGFR inhibitors or protein expression like HER2 or PD-L1 status. To fully
275    leverage the capabilities of modern profiling techniques, however, information across all available
276    domains must be integrated to optimize the therapeutic strategy for a given patient. Indeed, our
277    findings underscore the importance of utilizing both genotype and immunophenotype to fully

278 characterize disease and assess efficacy of treatment. For example, CAR T-cell therapy derives
279 specificity from protein expression, yet would fail to elicit a complete response if pathogenic
280 genotypes were distributed across multiple phenotypic clusters. Such a scenario would require
281 joint single-cell profiling as in DAb-seq to unravel. As multiomic single-cell technologies like
282 DAb-seq become available, it will be feasible to use comprehensive precision analysis to
283 deconvolute the subtlety of each patient's cancer and thereby select the best treatment regimen.
284
285
286 **Methods**
287
288 *Conjugation of antibodies to oligonucleotide barcodes*
289 Monoclonal antibodies were conjugated to azide-modified oligonucleotides using a copper-free
290 click chemistry reaction as described previously[39]. Monoclonal antibodies were resuspended to
291 100 μg in 100 μL PBS. See Supplementary Table 3 for a complete list of antibodies and
292 oligonucleotide barcode sequences. Antibodies were incubated with DBCO-PEG5-NHS Ester
293 linker (Click Chemistry Tools, cat. no. A102P) at a 4:1 molar ratio linker:antibody for 2 h at room
294 temperature. Following incubation, the antibody-linker solution was washed once in a 50 kDa
295 cellulose spin filter (Millipore Sigma, cat. no. UFC505024). DNA oligonucleotides with a 5' azide
296 modification (Integrated DNA Technologies) were reconstituted in water and added to the washed
297 antibodies at a 2.5:1 molar ratio oligonucleotide:antibody. Following a 16 h incubation, the
298 conjugated antibodies were washed three times in a 50 kDa filter to remove unreacted
299 oligonucleotides. All antibody conjugates were run on a Bioanalyzer Protein 230 electrophoresis
300 chip (Agilent Technologies, cat. no. 5067-1517) to verify successful conjugation.
301
302 *Cell culture and PBMC processing for control experiments*
303 The following three cell lines were used in the initial control experiment: Raji (ATCC, CCL-86),
304 Jurkat (ATCC, TIB-152), K562 (ATCC, CCL-243). Cells were cultured under the supplier's
305 recommended conditions. PBMCs from a single healthy donor were sourced commercially
306 (iXCells Biotechnologies, cat. no. 10HU-003) and stored at -80°C until use. Prior to staining, the
307 cultured cell lines and PBMCs were washed once in PBS with 5% fetal bovine serum (FBS)
308 (Thermo Fisher, cat. no. 10082147). For the control experiment, the three cell lines were combined
309 at an equal ratio.
310
311 *Collection of patient samples*
312 Patients included in this study were treated at the University of California, San Francisco (UCSF),
313 and peripheral blood or bone marrow was stored in the UCSF tumor bank. Samples were processed
314 immediately after collection to isolate mononuclear cells. Sample collection was in accordance
315 with the Declaration of Helsinki under institutional review board-approved tissue banking
316 protocols. Written informed consent was obtained from all patients.
317
318 *Thawing patient samples*
319 A protocol was optimized to maximize recovery of viable cells from patient samples. Cryovials
320 containing patient tissue (peripheral blood or bone marrow aspirate) were warmed by hand and
321 carefully transferred dropwise to a 50 mL tube containing 40 mL of cold DMEM media (Thermo
322 Fisher, cat. no. 11995040) with 20% FBS and 2 mM EDTA. The tube was centrifuged at 700 rpm
323 at 4°C for 7 min with no brake. The supernatant was discarded, and the cells were resuspended in

324    10 mL of warmed RPMI-1640 media (Thermo Fisher, cat. no. A1049101) with 10% FBS. The
325    solution was strained through a 70 µm cell strainer (Corning, cat. no. 431751) to remove any large
326    cell aggregates and the tube was centrifuged a second time at 700 rpm at 4°C for 5 min with low
327    brake. The supernatant was discarded, and the cells were resuspended in PBS with 5% FBS for
328    staining.
329
330    *Cell staining using oligonucleotide-conjugated antibodies*
331    For each sample, 2 million cells were added to a 5 mL DNA LoBind tube (Eppendorf, cat. no.
332    0030108310), centrifuged at 400 x g for 4 min, and resuspended in 180 µL PBS with 5% FBS.
333    Cells were blocked for 10 min on ice following addition of 10 µL Fc blocking solution (BioLegend,
334    cat. no. 422301), 4 µL of a 1% dextran sulfate solution (Research Products International, cat. no.
335    D20020), and 4 µL of 10 mg/mL salmon sperm DNA (Invitrogen, cat. no. 15632011). Cells were
336    stained for 30 min on ice with 0.5 µg of each conjugated antibody. After incubation, five rounds
337    of washing were performed to remove excess antibody. For each wash, 5 mL PBS with 5% FBS
338    was added to the tube and centrifuged at 400 x g for 4 min. Stained cells were resuspended in
339    Mission Bio cell buffer at a final concentration of 3 M/mL prior to microfluidic encapsulation.
340
341    *Microfluidic single-cell DNA genotyping and antibody capture*
342    A commercial single-cell DNA genotyping platform (Mission Bio, Tapestri) was used to perform
343    microfluidic encapsulation, lysis, and barcoding according to the manufacturer's protocol for the
344    acute myeloid leukemia V1 panel. Where noted, modifications were made to enable co-capture of
345    oligonucleotide-labeled antibodies. Stained cells were loaded into a microfluidic cartridge and co-
346    encapsulated into droplets with a lysis buffer containing protease and mild detergent. Droplets
347    were incubated in a thermal cycler for 1 h at 50°C to digest all cellular proteins, followed by 10
348    min at 80°C to heat-inactivate the protease.  To enable antibody capture during the barcoding stage,
349    the antibody tags were designed with 3' complementarity to one of the *RUNX1* gene forward
350    primers and the corresponding reverse primer was omitted from the reverse primer pool.
351    Supplementary Table 2 lists the sequences of the forward and reverse primers in the DNA panel.
352    Lysed cells in droplets were transferred to the barcoding module of the microfluidic cartridge in
353    addition to polymerase mix, the modified reverse primer pool, barcoded hydrogel beads, and oil
354    for droplet generation. The droplets were placed under a UV lamp (Analytik Jena, Blak-Ray
355    XX15L) for 8 min to cleave the single-stranded PCR primers containing unique cell barcodes from
356    the hydrogel beads. To amplify DNA targets and capture antibody tags, droplets were thermal
357    cycled using the following program: 95°C for 10 m; 20 cycles of (95°C for 30 s, 72°C for 10 s,
358    61°C for 4 min, 72°C for 30 s); 72°C for 2 min; 4°C hold.
359
360    *Single-cell DNA amplicon and antibody tag sequencing library preparation*
361    Recovery and cleanup of single-cell libraries proceeded according to the Mission Bio V1 protocol
362    with additional modifications for antibody library preparation. The 8 PCR tubes containing
363    barcoded droplets were pooled as pairs and treated with Mission Bio Extraction Agent. Water was
364    added to each tube and the aqueous fraction transferred to a new 1.5 mL DNA LoBind tube.
365    Ampure XP beads (Beckman Coulter, cat. no. A63881) were added at a 0.75X volume ratio
366    beads:PCR product for size selection. The supernatant from the size selection step, containing
367    library fragments shorter than ~200 bp, was retained and used for antibody library preparation,
368    while the remaining beads with bound DNA panel library fragments were washed twice with 80%
369    EtOH    and    eluted    in    30    µL    water.    A    biotinylated    capture    oligonucleotide

370    (/5Biosg/GGCTTGTTGTGATTCGACGA/3C6/, Integrated DNA Technologies) complementary
371    to the 5' end of the antibody tags was added to the retained supernatant to a final concentration of
372    0.6 μM. The supernatant-probe solution was heated to 95°C for 5 min to denature the PCR product,
373    then snap-cooled on ice for probe hybridization. 10 μL of streptavidin beads (Thermo Fisher, cat.
374    no. 65001) were washed according to the manufacturer's protocol and added to each tube of PCR
375    product. Following a 15 min incubation at room temperature, the beads were isolated by magnetic
376    separation, washed two times in PBS, and resuspended in 30 μL water. PCR was performed on the
377    purified DNA panel and antibody tags to produce sequencing libraries. For each tube of purified
378    DNA panel, 50 μL reactions were prepared containing 4 ng of barcoded product in 15 μL water,
379    25 μL Mission Bio Library Mix, and 5 μL each of custom P5 and Nextera P7 primers (N7XX),
380    both at 4 μM stock concentration. The reactions were thermal cycled using the following program:
381    95°C for 3 min; 10 cycles of (98°C for 20 s, 62°C for 20 s, 72°C for 45 s); 72°C for 2 min; 4°C
382    hold. For each tube of purified antibody tags, identical reactions were prepared, instead using 15
383    μL bead-bound template, 5 μL antibody tag-specific P7 primer at 4 μM, and 20 cycles of
384    amplification. See Supplementary Table 4 for a complete listing of custom library preparation
385    primers. Following amplification, both the DNA panel and antibody tag libraries were cleaned
386    with 0.7X Ampure XP beads and eluted in 12 μL water.
387

388    *Next-generation sequencing*
389    All DNA panel and antibody tag libraries were run on a Bioanalyzer High Sensitivity DNA
390    electrophoresis chip (Agilent Technologies, cat. no. 5067-4626) to verify complete removal of
391    primer-dimer products. Libraries were quantified by fluorometer (Qubit 3.0, Invitrogen) and
392    sequenced on Illumina next-generation sequencing platforms with a 20% spike-in of PhiX control
393    DNA (Illumina, cat. no. FC-110-3001). All sequencing runs used a dual-index configuration and
394    a custom Read 1 primer (5' GCCTGTCCGCGGAAGCAGTGGTATCAACGCAGAGTAG 3',
395    Integrated DNA Technologies). The 3-cell control sample was sequenced on an Illumina MiSeq
396    using a v2 300-cycle kit in 2x150 bp paired-end mode (Illumina, cat. no. MS-102-2002). For the
397    patient samples, DNA panel and antibody tag libraries were sequenced separately to maximize
398    cost-effectiveness. DNA panels were sequenced with an Illumina NovaSeq 6000 SP 300-cycle Kit
399    (Illumina, cat. no. 20027465) in 2 x 150 bp paired-end mode. Antibody tag libraries were
400    sequenced with an Illumina NextSeq 550 75-cycle High Output Kit (Illumina, cat. no. 20024906)
401    in paired-end mode, using 38 cycles for Read 1 and 39 cycles for Read 2.
402

403    *Bioinformatic pipeline for single-cell DNA genotyping and antibody tag counting*
404    Sequencing data was processed using a custom pipeline available on GitHub (see Code
405    Availability). For all reads, combinatorial cell barcodes were parsed from Read 1 using cutadapt
406    (v2.4) and matched to a barcode whitelist. Barcode sequences within a Hamming distance of 1
407    from a whitelist barcode were corrected.
408

409    For the DNA genotyping libraries, reads with valid barcodes were trimmed with cutadapt to
410    remove 5' and 3' adapter sequences and demultiplexed into single-cell FASTQ files using the
411    script "demuxbyname" from the BBMap package (v.38.57). Valid cell barcodes were selected
412    using the inflection point of the cell rank plot in addition to the requirement that 60% of DNA
413    intervals were covered by a minimum of 8 reads. FASTQ files for valid cells were aligned to the
414    hg19 build of the human genome reference using bowtie2 (v2.3.4.1). The single-cell alignments
415    in BAM format were filtered (properly mapped, mapping quality > 2, primary alignment), sorted,

416    and indexed with samtools (v1.8). GVCF files were produced for all cells using HaplotypeCaller
417    from the GATK suite (v.4.1.3.0). Joint genotyping was performed on all genomic intervals in
418    parallel (excluding primer regions) using GATK GenotypeGVCFs. For longitudinal patient
419    samples, cells from all timepoints were joint genotyped as a multi-sample cohort. Genotyped
420    intervals from all cells were combined into a single variant call format (VCF) file, and multiallelic
421    records were split and left-aligned using bcftools (v1.9). Variants were annotated with ClinVar
422    metadata (v.20190805) and SnpEff functional impact predictions (v4.3t). Variant records for all
423    cells were exported to HDF5 format using a condensed representation of the genotyping calls (0:
424    wildtype; 1: heterozygous alternate; 2: homozygous alternate; 3: no call).
425
426    The antibody tag libraries were processed identically for cell barcode demultiplexing. For reads
427    with valid cell barcodes, 8 bp antibody barcodes and 10 bp unique molecular identifiers (UMIs)
428    were extracted from Read 2 using cutadapt with the requirement that all UMI bases had a minimum
429    quality score of 20. Antibody barcode sequences within a Hamming distance of 1 from known
430    antibody barcodes were corrected. UMI sequences were grouped by cell and antibody and counted
431    using the UMI-tools package (v.0.5.3, "adjacency" method). UMI counts of antibodies for each
432    cell barcode were exported in tabular format for further analysis.
433
434    *Cell and genotype filtering*
435    Cell barcodes were additionally filtered according to antibody counts. Valid barcode groups were
436    required to have a minimum of 100 antibody UMIs by the adjacency counting method and a
437    maximum IgG1 count no greater than five times the median IgG1 count of the associated DAb-
438    seq experiment. For each valid cell barcode, all variants were filtered according to the quality and
439    sequence depth reported by GATK. Genotyping calls were required to have a minimum quality of
440    30 and total depth of 10; variant entries below these thresholds were marked as "no call" and
441    excluded from analyses.
442
443    *Antibody-based embedding and clustering*
444    To correct for technical effects in the raw antibody counts and batch variability between
445    experiments from the same patient but different time points, a linear regression over all cells from
446    the same patient was performed. Specifically, to all entries $c_{ij}$ of the UMI corrected antibody count
447    matrix $\mathbf{c}$, where $i$ is the cell index and $j$ the antibody index, one pseudocount was added and the
448    matrix was log-transformed. A matrix of quality metrics $\mathbf{q}$ with cells as rows and four columns
449    (total antibody reads, total antibody counts after UMI correction, IgG1 count and total amplicon
450    reads) was log-transformed, column-wise normalized, and mean-centered. A singular value
451    decomposition was performed on the transformed matrix $\mathbf{q}$ and the left-singular vectors retained
452    as design matrix. Each column vector $\mathbf{c}_j$ was then regressed with either the first three, two, or one
453    left-singular vectors, for patient samples, PMBC or cell lines respectively as regressors. The vector
454    of residuals $\mathbf{u}_j$ is then the corrected antibody signal of antibody $j$ (Extended Data Figure 1).
455
456    A UMAP embedding in two dimensions of the corrected antibody signal was done in Python 2.7
457    using the umap-learn[30] (v0.3.10) and scanpy[40] (v.1.4.4.post1) packages, with the minimum
458    distance parameter set to 0.1 for the pediatric patient and 0.2 for all other samples and default
459    parameters otherwise. To construct the underlying nearest neighbor graph from the corrected
460    antibody count matrix, 15 or 16 nearest neighbors based on the first 16 to all principal components

461 were used. The scanpy implementation of the Leiden algorithm[35] with resolution set to 0.1 for the
462 three cell line experiment and 1 otherwise was used to assign cells to phenotypic compartments.
463 For the gradient analysis of the pediatric Patient with AML (Figure 4), only cells belonging to
464 Leiden communities with blast phenotype were retained and the singular value decomposition of
465 the remaining rows of **u** was calculated. Cells were then ordered by their value of the second left-
466 singular vector. Antibody counts and genotype fractions along the gradient were averaged with a
467 moving window of 200 cells. Similarly, the average position of the cells in the two-dimensional
468 UMAP embedding was estimated by smoothing x and y coordinates with a moving window of the
469 same length. A 3rd-order spline was placed through the smoothed cell position to indicate the
470 orientation of the gradient in the UMAP embedding.
471
472 **Code Availability**
473
496
497 **Author Contributions**
498
499 B.D. and C.L.D. performed the experiments, sequenced the samples, analyzed the data, and wrote
500 the initial draft of the manuscript. C.A.C.P. assisted with patient sample processing. H.N.V. and
501 A.A. revised the manuscript. C.S. treated the patients and obtained samples. All authors read,
502 reviewed, and approved the manuscript.
503
504 **References**
505
506 1.      Arber, D. A. *et al.* The 2016 revision to the World Health Organization classification of

507        myeloid neoplasms and acute leukemia. *Blood* **127**, 2391–2405 (2016).

508    2.    Marusyk, A., Almendro, V. & Polyak, K. Intra-tumour heterogeneity: a looking glass for
509        cancer? *Nat. Publ. Gr.* **12**, 323–334 (2012).

510    3.    Suvà, M. L. & Tirosh, I. Single-Cell RNA Sequencing in Cancer: Lessons Learned and
511        Emerging Challenges. *Mol. Cell* **75**, 7–12 (2019).

512    4.    Landau, D. A., Carter, S. L., Getz, G. & Wu, C. J. Clonal evolution in hematological
513        malignancies and therapeutic implications. *Leukemia* **28**, 34–43 (2014).

514    5.    Patel, J. P. *et al.* Prognostic Relevance of Integrated Genetic Profiling in Acute Myeloid
515        Leukemia. *N. Engl. J. Med.* **366**, 1079–1089 (2012).

516    6.    Buckley, S. A. & Walter, R. B. Antigen-specific immunotherapies for acute myeloid
517        leukemia. *Hematology* **2015**, 584–595 (2015).

518    7.    García-Dabrio, M. C. *et al.* Complex Measurements May Be Required to Establish the
519        Prognostic Impact of Immunophenotypic Markers in AML. *Am. J. Clin. Pathol.* **144**, 484–
520        492 (2015).

521    8.    Papaemmanuil, E. *et al.* Genomic Classification and Prognosis in Acute Myeloid
522        Leukemia. *N. Engl. J. Med.* **374**, 2209–2221 (2016).

523    9.    Paguirigan, A. L. *et al.* Single-cell genotyping demonstrates complex clonal diversity in
524        acute myeloid leukemia. *Sci. Transl. Med.* **7**, 281re2 (2015).

525    10.   Wang, L. *et al.* Integrated single-cell genetic and transcriptional analysis suggests novel
526        drivers of chronic lymphocytic leukemia. *Genome Res.* **27**, 1300–1311 (2017).

527    11.   Pellegrino, M. *et al.* High-throughput single-cell DNA sequencing of acute myeloid
528        leukemia tumors with droplet microfluidics. *Genome Res.* **28**, 1345–1352 (2018).

529    12.   Smith, C. C. *et al.* Heterogeneous resistance to quizartinib in acute myeloid leukemia
530        revealed by single-cell analysis. *Blood* **130**, 48–58 (2017).

531    13.   De Zen, L. *et al.* Quantitative multiparametric immunophenotyping in acute lymphoblastic
532        leukemia: correlation with specific genotype. I. ETV6/AML1 ALLs identification.
533        *Leukemia* **14**, 1225–1231 (2000).

534    14.   van Galen, P. *et al.* Single-Cell RNA-Seq Reveals AML Hierarchies Relevant to Disease
535        Progression and Immunity. *Cell* **0**, 1–17 (2019).

536    15.   Giustacchini, A. *et al.* Single-cell transcriptomics uncovers distinct molecular signatures
537        of stem cells in chronic myeloid leukemia. *Nat. Med.* **23**, 692–702 (2017).

538    16.   Nam, A. S. *et al.* Somatic mutations and cell identity linked by Genotyping of
539        Transcriptomes. *Nature* **571**, 355–360 (2019).

540    17.   Shalek, A. K. *et al.* Single-cell transcriptomics reveals bimodality in expression and
541        splicing in immune cells. *Nature* **498**, 236–240 (2013).

542    18.   Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of Individual
543        Cells Using Nanoliter Droplets. *Cell* **161**, 1202–1214 (2015).

544    19.   Klein, A. M. *et al.* Droplet Barcoding for Single-Cell Transcriptomics Applied to
545        Embryonic Stem Cells. *Cell* **161**, 1187–1201 (2015).

546    20.   Zeisel, A. *et al.* Brain structure. Cell types in the mouse cortex and hippocampus revealed
547        by single-cell RNA-seq. *Science* **347**, 1138–42 (2015).

548    21.   Schiebinger, G. *et al.* Optimal-Transport Analysis of Single-Cell Gene Expression
549        Identifies Developmental Trajectories in Reprogramming. *Cell* **176**, 928-943.e22 (2019).

550    22.   Deng, Q., Ramsköld, D., Reinius, B. & Sandberg, R. Single-cell RNA-seq reveals
551        dynamic, random monoallelic gene expression in mammalian cells. *Science (80-. ).* **343**,
552        193–196 (2014).

23. Mansour, M. R. *et al.* Oncogene regulation. An oncogenic super-enhancer formed through somatic mutation of a noncoding intergenic element. *Science* **346**, 1373–7 (2014).

24. Rheinbay, E. *et al.* Analyses of non-coding somatic drivers in 2,658 cancer whole genomes. *Nature* **578**, 102–111 (2020).

25. Shahi, P., Kim, S. C., Haliburton, J. R., Gartner, Z. J. & Abate, A. R. Abseq: Ultrahigh-throughput single cell protein profiling with droplet microfluidic barcoding. *Sci. Rep.* **7**, 1–12 (2017).

26. Schuurhuis, G. J. *et al.* Minimal/measurable residual disease in AML: a consensus document from the European LeukemiaNet MRD Working Party. *Blood* **131**, 1275–1291 (2018).

27. Wood, B. L. Flow Cytometric Monitoring of Residual Disease in Acute Leukemia. in 123–136 (Humana Press, Totowa, NJ, 2013). doi:10.1007/978-1-62703-357-2_8

28. Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res.* **47**, D941–D947 (2019).

29. Appelbaum, F. R. & Bernstein, I. D. Gemtuzumab ozogamicin for acute myeloid leukemia. *Blood* **130**, 2373–2376 (2017).

30. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. (2018).

31. Campana, D. & Behm, F. G. Immunophenotyping of leukemia. *J. Immunol. Methods* **243**, 59–75 (2000).

32. Sexauer, A. *et al.* Terminal myeloid differentiation in vivo is induced by FLT3 inhibition in FLT3/ITDAML. *Blood* **120**, 4205–4214 (2012).

33. McMahon, C. M. *et al.* Gilteritinib induces differentiation in relapsed and refractory FLT3-mutated acute myeloid leukemia. *Blood Adv.* **3**, 1581–1585 (2019).

34. Yun, H. D. *et al.* Erythroid differentiation of myeloblast induced by gilteritinib in relapsed FLT3-ITD–positive acute myeloid leukemia. *Blood Advances* **3**, 3709–3712 (2019).

35. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing well-connected communities. *Sci. Rep.* **9**, 5233 (2019).

36. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. (2008). doi:10.1088/1742-5468/2008/10/P10008

37. Buscarlet, M. *et al.* DNMT3A and TET2 dominate clonal hematopoiesis and demonstrate benign phenotypes and different genetic predispositions. *Blood* **130**, 753–762 (2017).

38. Jaiswal, S. *et al.* Clonal Hematopoiesis and Risk of Atherosclerotic Cardiovascular Disease. *N. Engl. J. Med.* **377**, 111–121 (2017).

39. Gong, H. *et al.* Simple Method To Prepare Oligonucleotide-Conjugated Antibodies and Its Application in Multiplex Protein Detection in Single Cells. *Bioconjug. Chem.* **27**, 217–225 (2016).

40. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* **19**, 15 (2018).
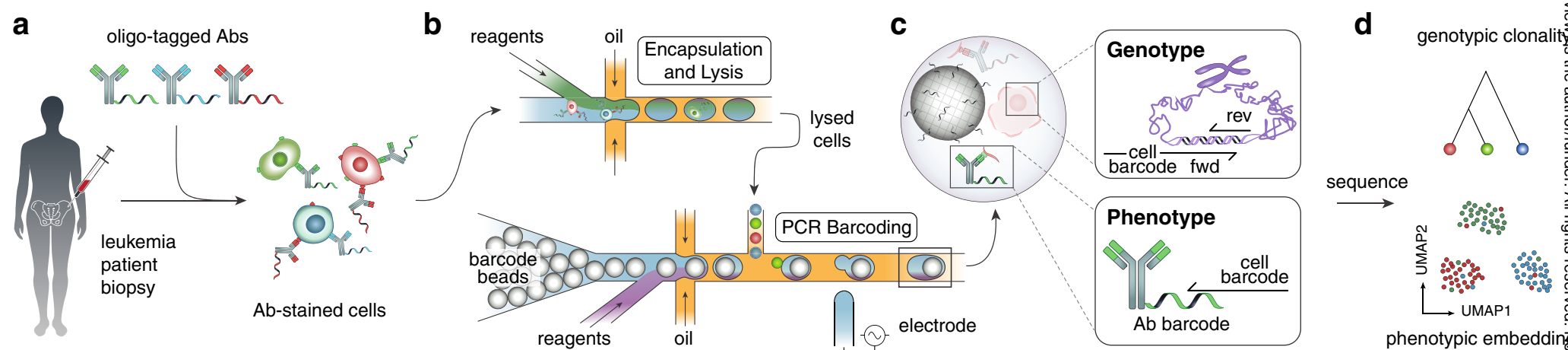
**Figure 1: The DAb-seq workflow.**
**a**, Bone marrow aspirates of patients with AML contain healthy and malignant cells that exhibit diverse genotypes and immunophenotypes. These cells are stained with antibodies labeled with DNA tags. **b**, Stained cells are paired and encapsulated with a barcode bead on a Mission Bio Tapestri instrument. **c**, In each droplet, a PCR labels antibody tags and genomic DNA targets simultaneously with a unique cell index. **d**, Sequencing the barcoded amplicons and antibody tags yields coupled single-cell immunophenotype and genotype data for thousands of cells.
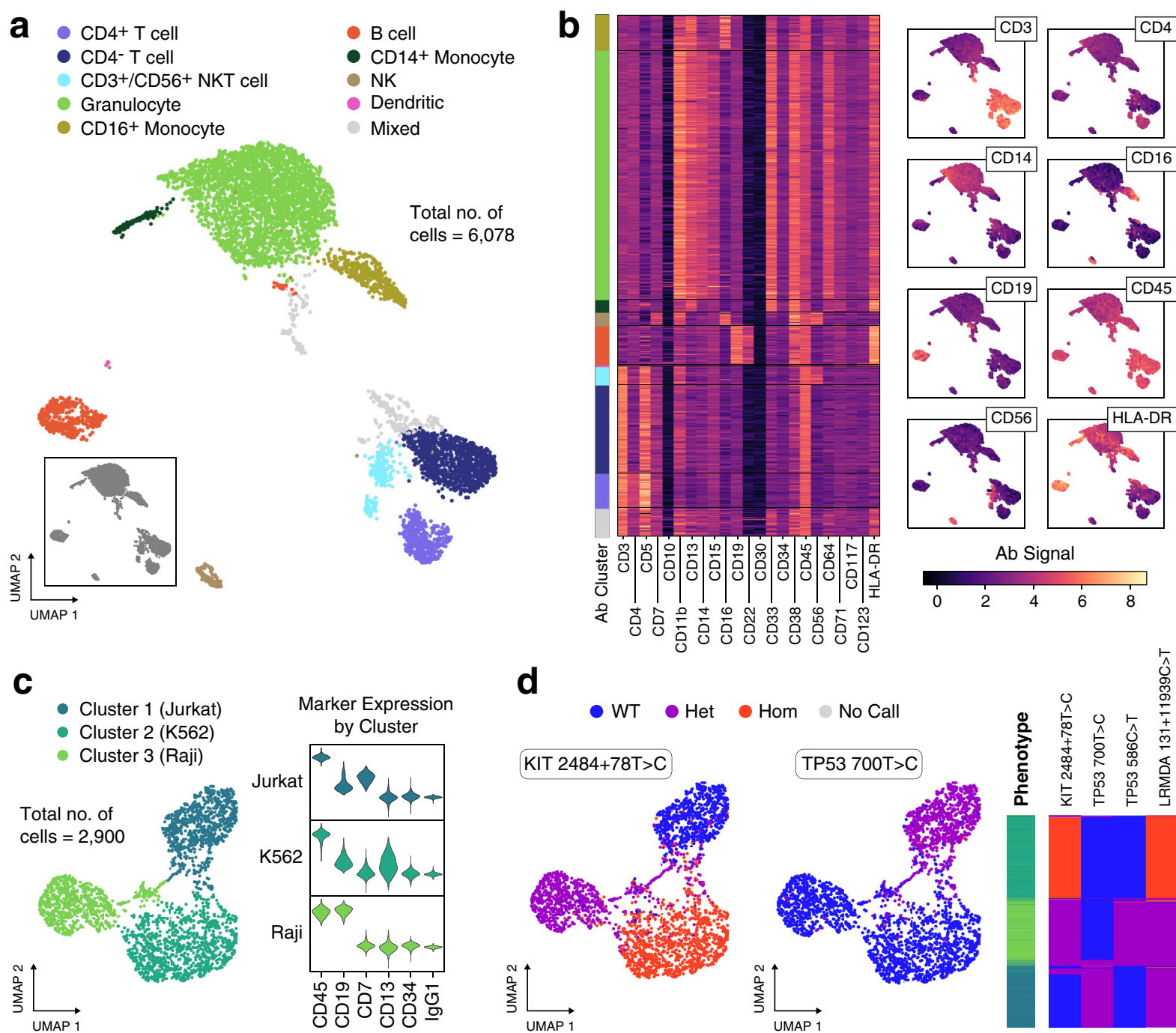
**Figure 2: DAb-seq enables simultaneous discrimination of single cells by their immunophenotype and genotype.** **a**, DAb-seq workflow performed on PBMCs from a healthy donor using a panel of 23 antibodies. Leiden clustering and two-dimensional UMAP embedding of the antibody tag data reveals expected blood compartments. Compartments are annotated based on detected marker expression. **b**, Heatmap of the corrected log-transformed antibody counts for each cell and antibody. Cells are ordered based on Leiden clusters. Overlay of corrected log-transformed antibody counts with the UMAP embedding highlights compartment-specific expression. **c**, Correspondence of antibody signal with genomic polymorphisms in DAb-seq experiments tested on a mixture of three cell lines and a panel of six antibodies. Cells cluster by antibody signal as shown in the UMAP embedding. **d**, Detected single nucleotide polymorphisms in these cells map to the phenotypic cell clusters as shown in the UMAP embedding and a heatmap, where rows correspond to single cells. The first column of the heatmap indicates assigned phenotype cluster, and the remaining columns indicate the genotyping call at the labeled loci.
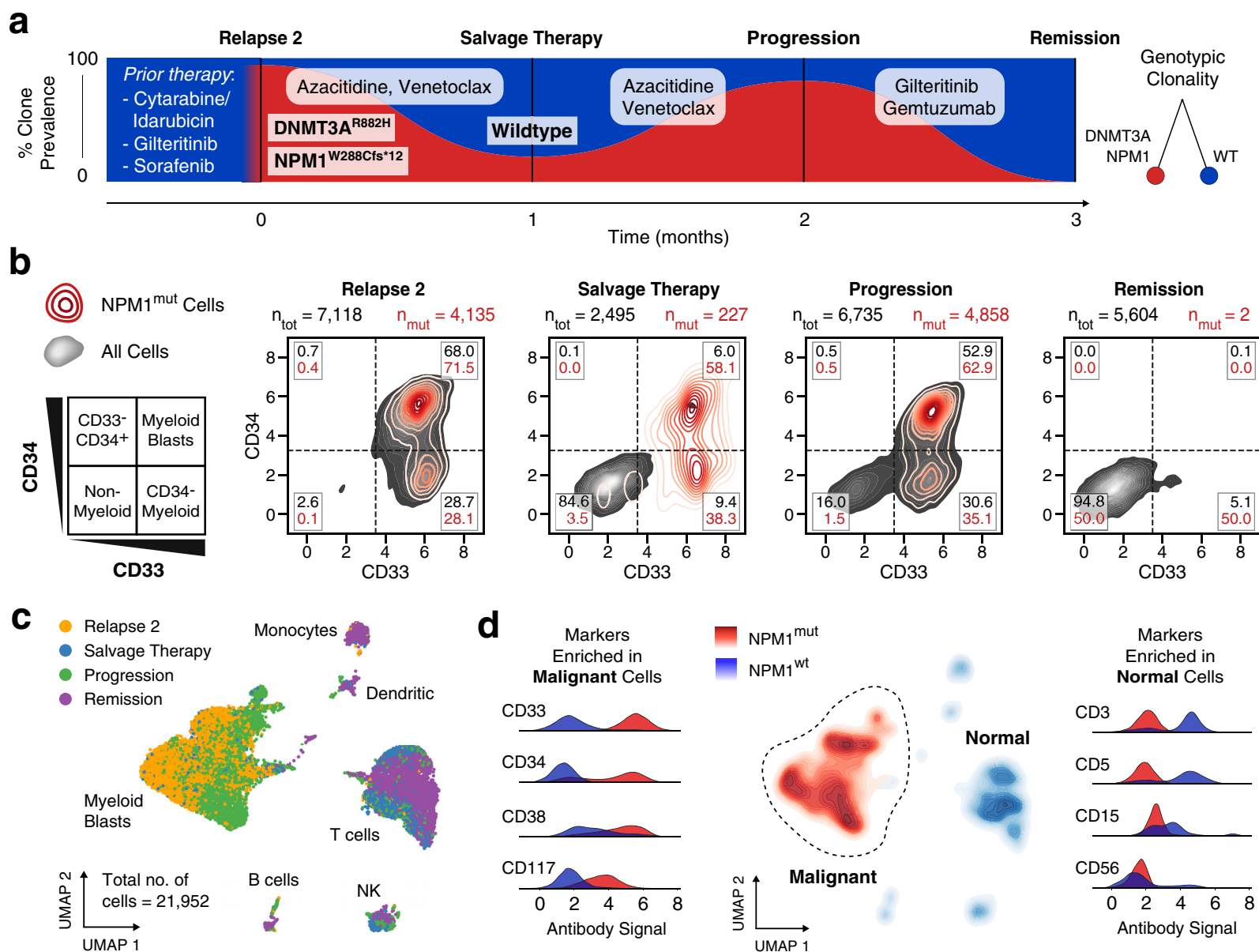
**Figure 3: AML blasts exhibit a stable genotype and phenotype through treatment.**
**a**, DAb-seq performed on four bone marrow aspirates of a patient with AML during disease progression as indicated in the fishplot (black lines). The patient received multiple rounds of chemotherapy prior to the experiment (Supplementary Table 1). The fraction of blast cells with NPM1 W288Cfs*12 (NPM1mut) mutation for each sampled time point detected by DAb-seq are shown in red. **b**, Scatter plots with kernel densities show CD33 and CD34 signal for all cells (grey) and NPM1mut cells (red) for each of the sampled time points. The percentage of normal and mutant cells within each gate are listed. Virtually gating cells highlights a persisting CD33+ blast population which is eradicated with gemtuzumab, a CD33-targeted therapy. **c**, UMAP embedding based on the log-transformed and corrected antibody counts from all cells labeled by timepoint indicates that the high-dimensional immunophenotype of the blasts is stable over the sampled timepoints. **d**, The genotype of each cell at the NPM1 locus is plotted as a kernel density estimate using the UMAP coordinates from c. Antibody signals enriched among malignant and normal populations are plotted as kernel densities using all cells and labeled by genotype.
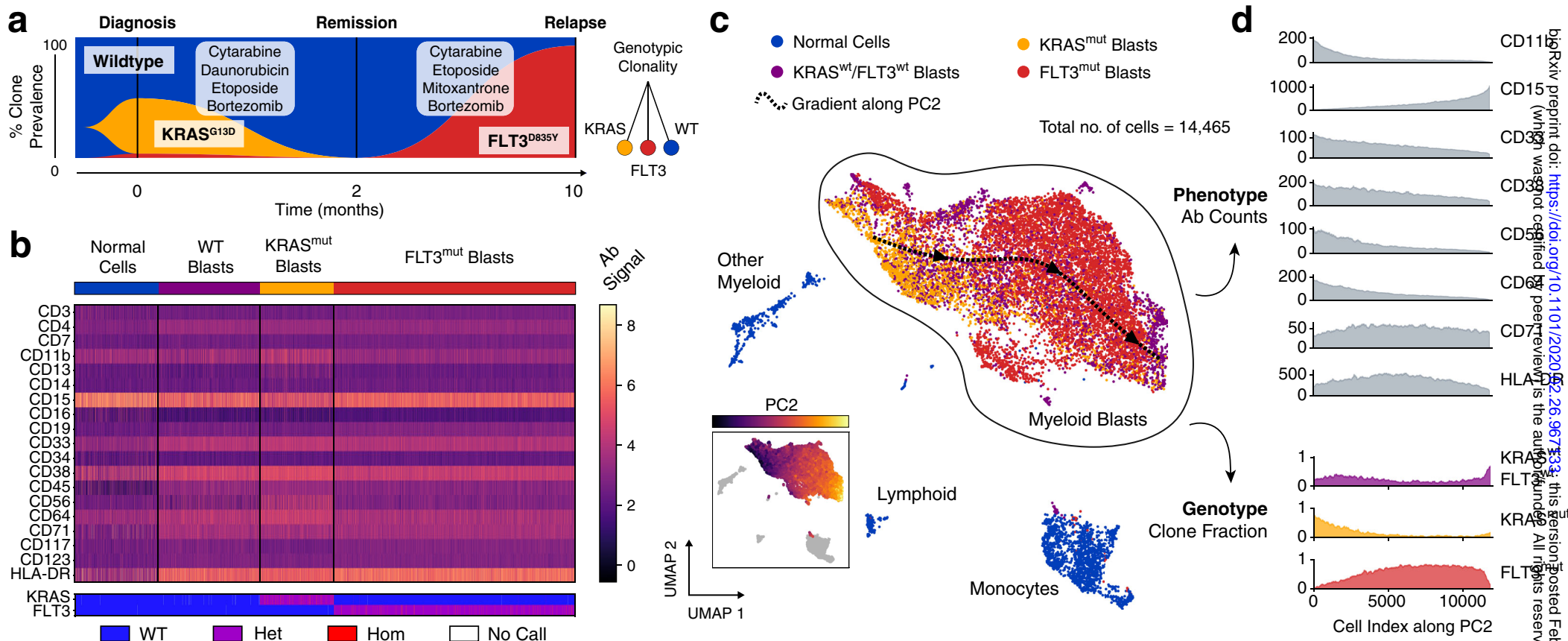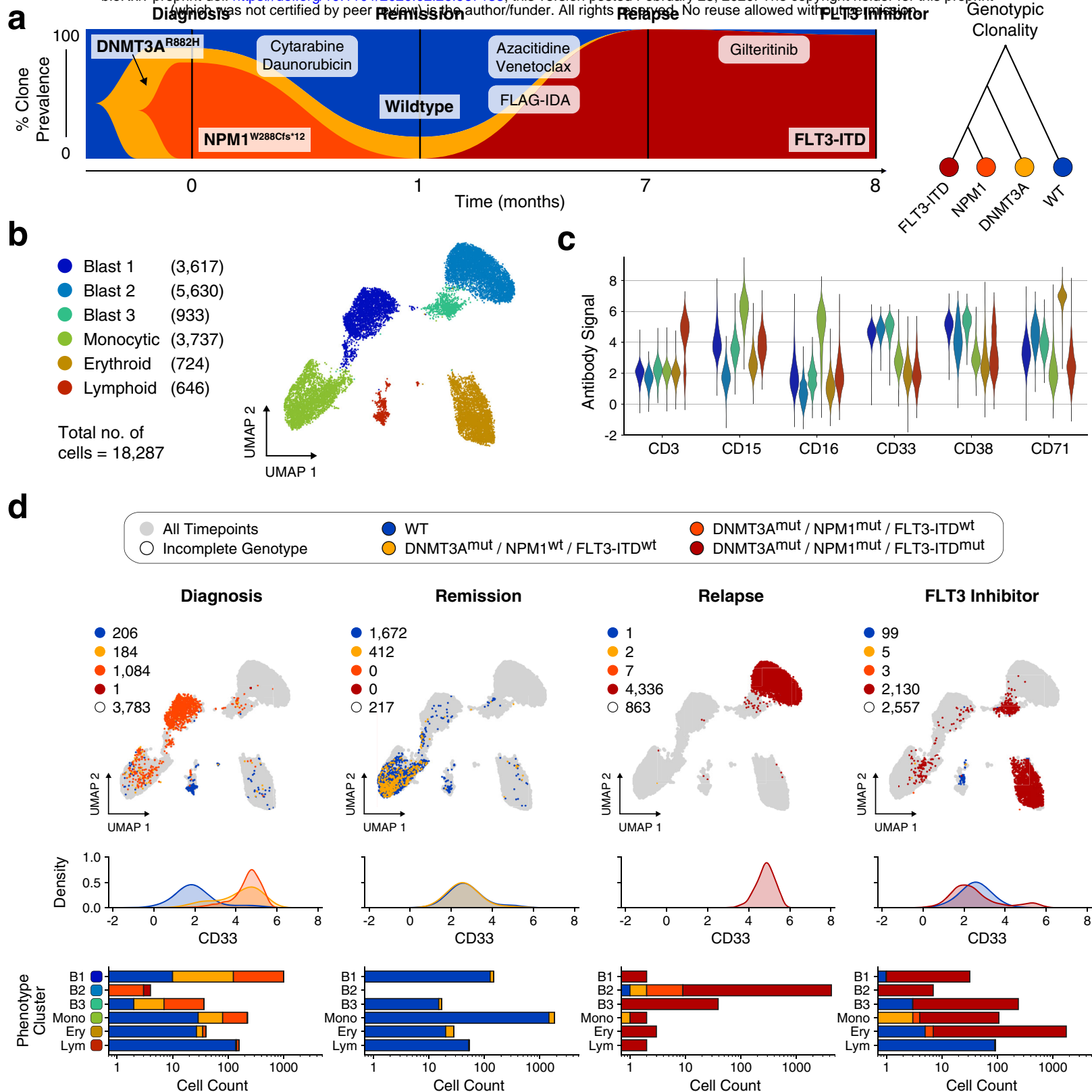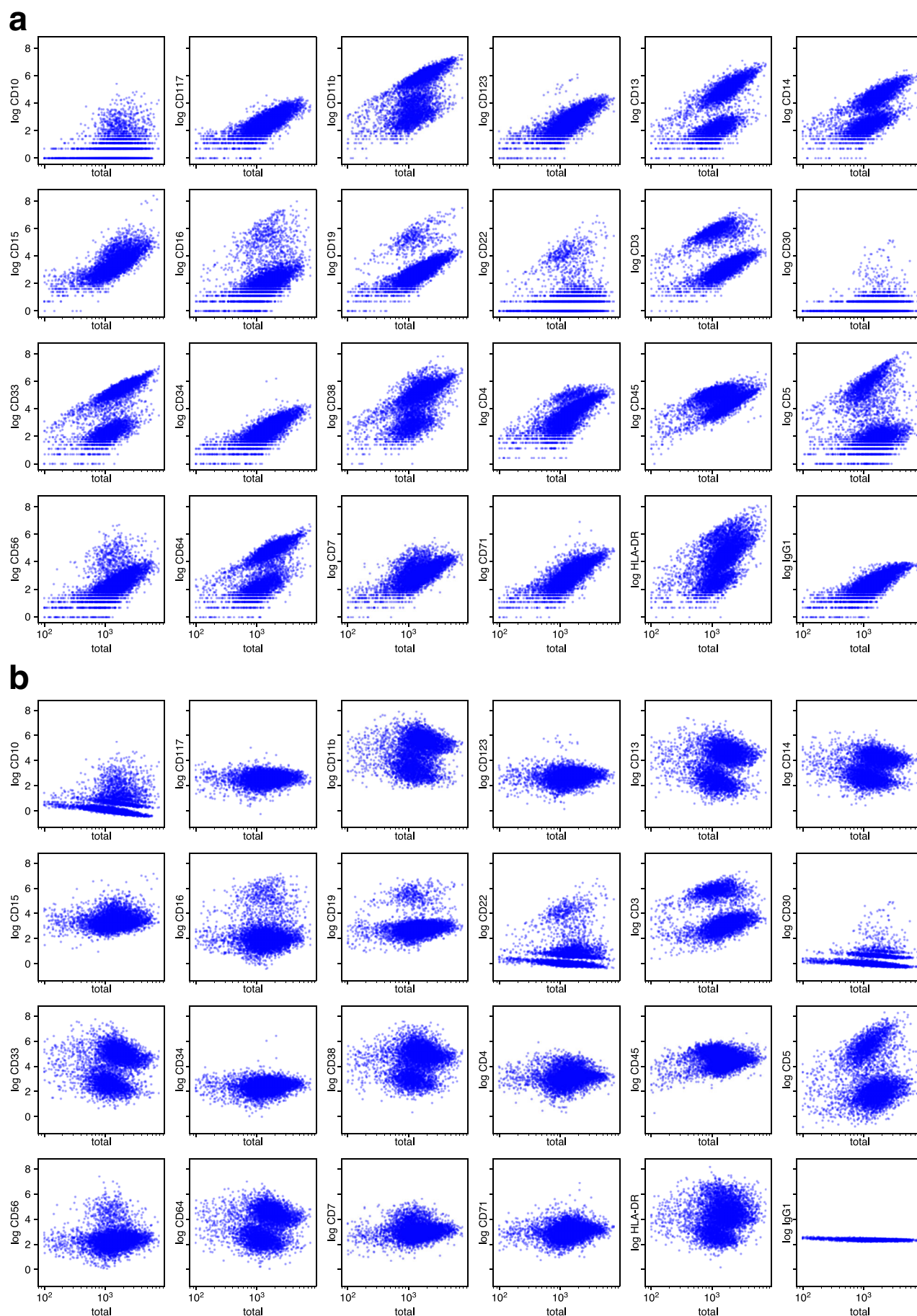
**Figure 4: Distinct genetic subclones form an overlapping immunophenotypic continuum in a case of pediatric AML.**
**a**, Three timepoints sampled with DAb-seq during treatment comprise a mixture of independent clones (KRAS G13D heterozygous blasts, yellow; FLT3 D835Y blasts, red). The wildtype compartment contains additional cells with a blast-like immunophenotype lacking detectable mutations. **b**, Heatmap of log-transformed corrected antibody counts and genotyping calls for the KRAS and FLT3 loci for each cell across all timepoints. The heatmap is grouped by genotype. Cells with wildtype genotype but blast-like immunophenotype are labeled separately. **c**, UMAP embedding of all cells from all time points based on log-transformed corrected antibody counts. Color indicates mutation status as in a. The blast compartment is overlaid with a spline approximating the gradient of the 2nd principal component of the antibody count matrix (shown in inlet figure) and indicates a gradual change in immunophenotype. **d**, Moving average expression of antibodies and fraction of mutated cells sorted by the 2nd principal component of the antibody count matrix. The overlapping phenotypic continuum between the genetically distinct blast clones is apparent.

**Figure 5: Decoupling of blast phenotype and genotype in response to FLT3 inhibitor therapy.**
**a**, Fishplot showing observed fraction of cells with distinct genetic mutations for each sampled time point. The co-occurrence of the three mutations in the single-cell data is consistent with a linear model of mutation accumulation. **b**, UMAP embedding of all cells based on measured antibody signal. The cells segregate into six distinct phenotypic clusters with multiple blast compartments. **c**, Average expression of each cell cluster for a selection of markers. **d**, Top row: Same UMAP embedding as in b given as grey outline. For each sampled time point, observed cells are plotted and colored according to the detected genotype. Blasts distribute among multiple phenotypic compartments in the final time point following FLT3 inhibitor treatment. Middle row: Kernel density plot of the CD33 antibody signal resolved by time point and genotype. Cells from genotypic compartments with less than 10 cells per time point are not plotted. Bottom row: Bar chart depicting genotypic composition of each phenotypic cluster in b resolved by time point.

**Extended Data Figure 1: Antibody count bias correction by linear regression.**
**a**, Raw UMI counts for each antibody and cell are plotted versus total antibody count from the same cell. A clear correlation between the two is visible. A similar slope is visible for the isotype control (bottom row, rightmost column), suggesting technical bias. **b**, Same plots as in a after correcting for global droplet performance by linear regression (see Methods). Correlation with total antibody counts is reduced.