

# Functional analysis of structural variants in single cells using Strand-seq

Received: 29 October 2021

Accepted: 7 October 2022

Published online: 24 November 2022

 Check for updates

Hyobin Jeong <sup>1,18,20</sup>, Karen Grimes <sup>1,2,20</sup>, Kerstin K. Rauwolf<sup>3</sup>, Peter-Martin Bruch <sup>4,5,6</sup>, Tobias Rausch <sup>1,5</sup>, Patrick Hasenfeld <sup>1</sup>, Eva Benito<sup>1</sup>, Tobias Roider <sup>1,4,5</sup>, Radhakrishnan Sabarinathan<sup>7</sup>, David Porubsky <sup>8,9,19</sup>, Sophie A. Herbst <sup>4,5</sup>, Büşra Erarslan-Uysal<sup>5,10</sup>, Johann-Christoph Jann <sup>11</sup>, Tobias Marschall <sup>12</sup>, Daniel Nowak <sup>11</sup>, Jean-Pierre Bourquin<sup>3</sup>, Andreas E. Kulozik <sup>5,10</sup>, Sascha Dietrich<sup>4,5,6,13</sup>, Beat Bornhauser <sup>3</sup>, Ashley D. Sanders <sup>1,14,15,16,21</sup> ✉ & Jan O. Korbel <sup>1,5,17,21</sup> ✉

Somatic structural variants (SVs) are widespread in cancer, but their impact on disease evolution is understudied due to a lack of methods to directly characterize their functional consequences. We present a computational method, scNOVA, which uses Strand-seq to perform haplotype-aware integration of SV discovery and molecular phenotyping in single cells by using nucleosome occupancy to infer gene expression as a readout. Application to leukemias and cell lines identifies local effects of copy-balanced rearrangements on gene deregulation, and consequences of SVs on aberrant signaling pathways in subclones. We discovered distinct SV subclones with dysregulated Wnt signaling in a chronic lymphocytic leukemia patient. We further uncovered the consequences of subclonal chromothripsis in T cell acute lymphoblastic leukemia, which revealed c-Myb activation, enrichment of a primitive cell state and informed successful targeting of the subclone in cell culture, using a Notch inhibitor. By directly linking SVs to their functional effects, scNOVA enables systematic single-cell multiomic studies of structural variation in heterogeneous cell populations.

The mutational landscapes of numerous cancers were recently cataloged<sup>1,2</sup>, revealing that somatic SVs represent around 55% of driver mutations<sup>2,3</sup>. Somatic mutational processes generate a broad spectrum of SVs from simple (for example, deletions and inversions) to complex classes (for example, chromothripsis)<sup>4–8</sup>, and these SVs are important drivers of malignancy, metastasis and relapse<sup>9–12</sup>. However, with the exception of focal deletions and amplifications, somatic SVs have proven difficult to characterize functionally in cancer genomic surveys<sup>1–3,13</sup>. Studies integrating transcriptome and whole genome sequencing (WGS) data have inferred SV functional outcomes<sup>13–16</sup>, but these typically require large cohorts and do not account for intratumor heterogeneity (ITH)<sup>3</sup>. Instead, SV effects can be measured directly by reading both genotype and molecular phenotype in the same cell,

using single-cell multiomics<sup>17–21</sup>. Several such methods have been developed<sup>17–20</sup>, but these do not presently account for small (<10 Mb) somatic copy number alterations (SCNAs), balanced SVs and complex rearrangement events like chromothripsis<sup>4,5,7,22</sup>, which has limited efforts to functionally characterize the most common class of driver mutations in cancer.

To address this, we developed scNOVA (single-cell nucleosome occupancy and genetic variation analysis)—a method enabling functional characterization of the full spectrum of somatic SV classes. scNOVA uses Strand-seq<sup>23</sup> in two ways: (1) it uses the DNA fragmentation pattern resulting from micrococcal nuclease (MNase) digestion<sup>23</sup> to directly measure nucleosome occupancy (NO) and indirectly infer patterns of gene activity, and (2) it couples this ‘molecular phenotype’

with SVs discovered by single-cell tri-channel processing (scTRIP, which jointly models read-orientation, read depth and haplotype-phase<sup>24</sup>) in the same cell. MNase digests the linker DNA between nucleosomes, leaving nucleosome-protected DNA intact, to enable genome-wide inference of NO by measuring sequence read counts<sup>25–28</sup>. Previous work has shown that active enhancers and transcribed genes exhibit reduced NO<sup>25–30</sup>. However, the relationships between NO and SV landscapes in cancer remain unexplored. scNOVA addresses this by integrating SVs and NO along the genome of a cell, to functionally characterize SVs in heterogeneous samples.

## Results

### NO classifies cell types and predicts gene activity changes

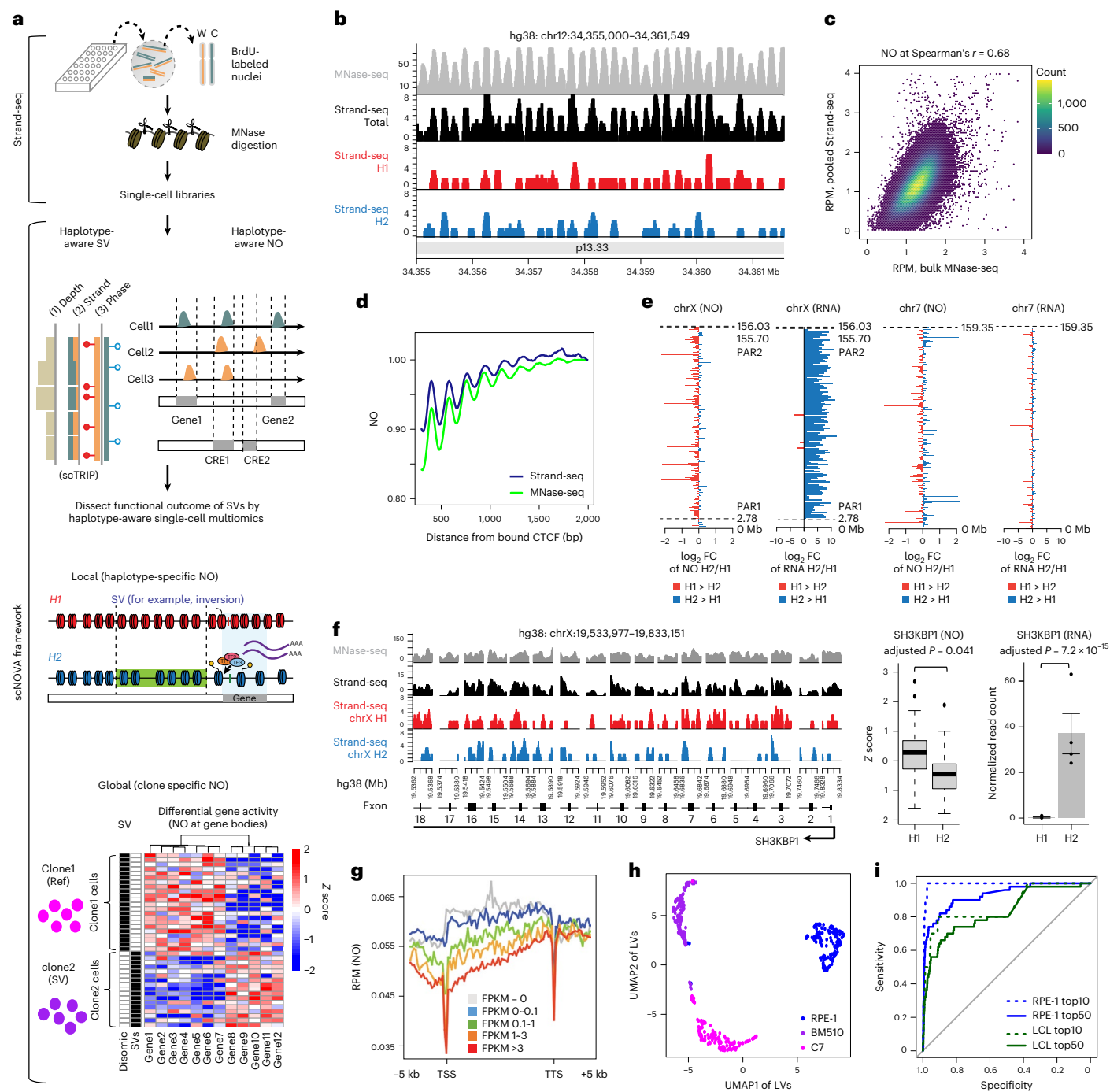
**Strand-seq data reveals NO.** We hypothesized that NO patterns derived from MNase fragmentation during Strand-seq library preparation could represent a readout to allow functional characterization of SVs (Fig. 1a and Extended Data Fig. 1). To test this, we evaluated whether Strand-seq data revealed nucleosome positioning through comparison with bulk MNase-seq data. We used the NA12878 lymphoblastoid cell line (LCL), which has both datatypes available, and pooled 95 Strand-seq libraries (sequenced to a median of 540,379 mapped nonduplicate reads per single cell<sup>31</sup>; Supplementary Table 1), into a ‘pseudobulk’ track, allowing direct comparison with the corresponding MNase-seq dataset (sequenced to 19-fold genomic coverage<sup>32</sup>). We measured NO along the genome (Methods) and found Strand-seq and MNase-seq were highly concordant in terms of uniformity of coverage and inferred nucleosome positions at DNase I hypersensitive sites (Spearman’s  $r = 0.68$ ) (Fig. 1b,c). Nucleosome positioning near the binding site of CTCF<sup>26,28</sup> (a key chromatin organizer) closely matched between both assays (Fig. 1d and Supplementary Fig. 1), and estimated nucleosome repeat lengths<sup>28</sup> were highly concordant (Supplementary Fig. 1). In addition, both assays measured NO in all 15 chromatin states identified by the Roadmap Epigenome Consortium<sup>33</sup>. Among these chromatin states, Strand-seq and MNase-seq revealed the highest NO signals on average for the polycomb-repressed state and the bivalent enhancer state, whereas the lowest average NO signals were consistently seen for the active transcription start site (TSS) state (Extended Data Fig. 2). This indicates that Strand-seq enables direct measurement of NO to reveal a ‘molecular readout’. We thus developed the scNOVA framework, which harnesses Strand-seq to measure NO genome-wide and couples this with SVs discovered in the same sequenced cell (Fig. 1a).

As Strand-seq resolves its measurements by haplotype<sup>31</sup>, we considered that haplotype-specific differences in NO (haplotype-specific NO) resulting from random monoallelic expression, germline SNPs and local effects of SVs could be harnessed for scNOVA. To assess the utility of haplotype-resolved NO, we phased 24,652,658 of 49,205,197 (50.1%) of the NA12878 Strand-seq read fragments, and pooled these reads to generate pseudobulk NO tracks for each chromosomal haplotype (denoted ‘H1’ and ‘H2’, respectively; Fig. 1b). Using the female-derived NA12878 cell line, we compared haplotype-specific NO to haplotype-resolved gene expression measurements from bulk RNA-seq data<sup>34</sup> (Methods). We identified a significant increase of NO in gene bodies mapping to H1 compared with H2 across the X chromosome (adjusted  $P = 0.0012$ ; Wilcoxon rank sum test), suggesting that H1 represents the inactive X chromosome. These data were consistent with haplotype-resolved gene expression measurements at loci subject to X-inactivation<sup>35</sup>, whereas genes escaping X-inactivation did not exhibit haplotype-specific NO (Fig. 1e,f and Supplementary Fig. 3). We also investigated whether Strand-seq data is informative of haplotype-specific NO at cis-regulatory elements (CREs), and identified a 1.4-fold enrichment for allele-specific CRE binding on the X chromosome ( $P = 0.015$ ; hypergeometric test; based on 718 CREs with haplotype-specific NO genome-wide; 10% false discovery rate (FDR)) (Supplementary Fig. 2). Moreover, CREs with haplotype-specific NO

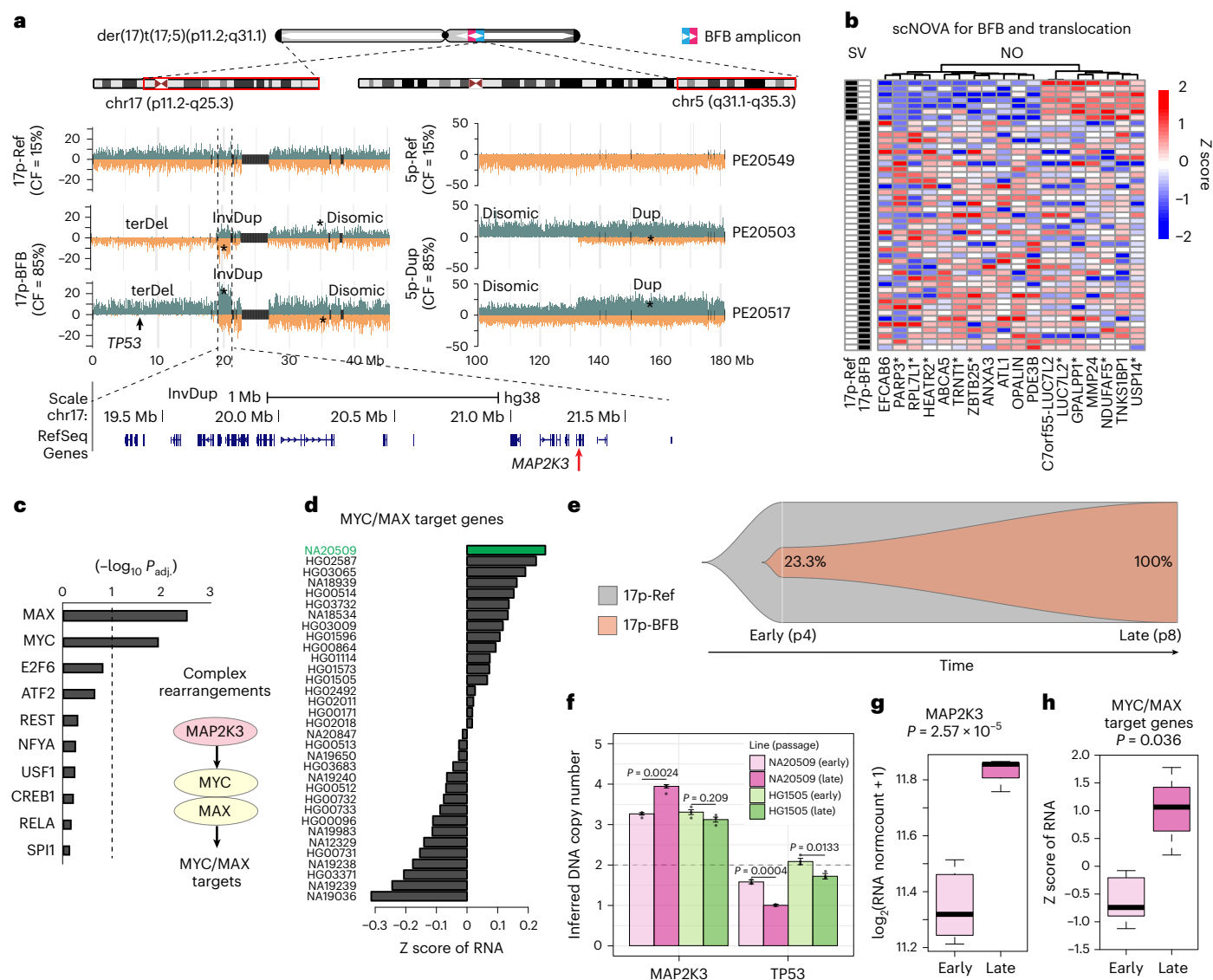
were significantly over-represented near genes showing allele-specific expression in the genome ( $P < 0.0018$ , hypergeometric test; Supplementary Fig. 2). These data suggest that haplotype-specific NO, a signal directly obtained from Strand-seq datasets, reflects biological gene regulation patterns in the genome.

**Cell-typing.** Since NO within gene bodies reflects gene activity in MNase-seq data<sup>28</sup>, we hypothesized that Strand-seq based NO patterns could be used to infer gene expression. To investigate this, we tested whether NO globally reflects cellular gene expression patterns in the retinal pigment epithelium-1 (RPE-1) cell line, for which we previously generated both Strand-seq and RNA-seq data<sup>24</sup>. To profile NO globally, we pooled 33 million read fragments (including phased and nonphased reads) from 79 Strand-seq libraries into pseudobulk NO tracks. We identified an inverse correlation between NO at gene bodies and gene expression ( $P < 2.2 \times 10^{-16}$ ; Spearman’s  $r$  of up to  $-0.24$ ; Fig. 1g and Supplementary Fig. 4), where highly expressed genes showed significantly lower NO within their gene bodies (and vice versa). We next explored the utility of NO for cell-type inference (‘cell-typing’), based on the activity of lineage-specific genes, by implementing a multivariate dimensionality reduction framework. We performed in silico mixing of Strand-seq libraries from different LCLs and RPE cell lines, and built a classifier that separates distinct cell types by partial least squares discriminant analysis (PLS-DA). We used a training set of 179 mixed libraries, and initially considered 19,629 features, which reflect ENSEMBL<sup>36</sup> genes with sufficient read coverage (Methods). After feature selection, 1,738 features were retained. We then used a nonoverlapping set of 123 cells to assess performance, all of which scNOVA classified accurately (area under the curve (AUC) = 1; Extended Data Fig. 3). Our framework also discriminated between cells from three related RPE cell lines derived from the same donor, which exhibit distinct SV landscapes<sup>24,37</sup> (AUC = 0.96; Fig. 1h) indicating that scNOVA enables accurate cell-typing.

**Gene activity changes between cell populations.** Having established that scNOVA can use the expression of lineage-specific genes for cell-typing, we evaluated if it could predict gene expression differences between defined cell populations, such as subclones bearing distinct SVs. We devised a module that integrates deep convolutional neural networks and negative binomial generalized linear models (Supplementary Figs. 5 and 6), to measure differential gene activity between two defined cell populations. To benchmark this module, we mixed Strand-seq libraries from different cell lines in silico, creating ‘pseudoclones’, and evaluated the predicted changes in gene activity between defined pseudoclones (each composed of cells from one cell line) by analyzing NO at gene bodies (Supplementary Fig. 7 and Extended Data Fig. 4). We first compared RPE-1 to the HG01573 LCL line, and defined the ground truth of expression using RNA-seq. We found that the differential gene activity score of scNOVA (Methods) was highly predictive of the ten most differentially expressed genes, where analyses of pseudoclones comprising 156 RPE-1 and 46 HG01573 libraries revealed an AUC of 0.93 (we observed a similar performance when analyzing the 50 most differentially expressed genes; Fig. 1i). Gene activity changes inferred included well-known markers of epithelial (for example, *EGFR*, *VCAM*) and lymphoid (for example, *CD74*, *CD100*) cell types (Supplementary Table 2). The scNOVA predictions were informative also when we simulated minor subclones present with clonal frequency (CF) = 20%, CF = 5% and CF = 1.3%, resulting in AUCs of 0.92, 0.79 and 0.68, respectively (Extended Data Fig. 4). We obtained similar results when applying scNOVA to pseudoclones derived from different (genetically related) RPE cell lines (Supplementary Fig. 7). These benchmarking exercises suggest that scNOVA can accurately infer gene activity changes between defined cell populations, suggesting that this framework can be used to functionally characterize subclonal SVs.



**Fig. 1** Haplotype-aware single-cell multiomics to functionally characterize SVs. **a**, Leveraging Strand-seq, scNOVA performs SV discovery and then, using phased NO tracks, identifies functional effects of SVs locally (via evaluation of haplotype-specific NO) and globally (clone-specific NO). Orange, Strand-seq reads mapped to the Watson (W) strand; green, reads mapped to the Crick (C) strand. **b**, Strand-seq-based NO tracks in NA12878 reveal nucleosome positions well-concordant with bulk MNase-seq, depicted for a chromosome 12 locus with relatively regular nucleosome positioning<sup>92</sup>. Red, NO tracks mapping to haplotype 1 (H1); blue, H2; black, combining phased and unphased reads; gray, MNase-seq. The y-axis depicts the mean read counts at each bp in 10 bp bins. **c**, Correlated NO at consensus DNase I hypersensitive sites<sup>33</sup> for NA12878. **d**, Averaged nucleosome patterns at CTCF binding sites<sup>34</sup> in NA12878, using pseudobulk Strand-seq and MNase-seq. **e**, FCs of haplotype-resolved NO in gene bodies plotted for chromosome X and chromosome 7 (a representative autosome) in NA12878. FCs of haplotype-resolved RNA expression measurements are shown to the right. **f**, Pseudobulk haplotype-phased NO track of exons of the representative chromosome X gene *SH3KBP1* based on Strand-seq. Boxplots comparing H1 and H2 use two-sided Wilcoxon rank sum tests followed by Benjamini–Hochberg multiple testing (FDR) correction (boxplots defined by minima = 25th percentile – 1.5 × interquartile range (IQR), maxima = 75th percentile + 1.5 × IQR, center = median and bounds of box = 25th and 75th percentile;  $n = 47$  single cells). Bar charts show haplotype-specific RNA expression of *SH3KBP1* (two-sided likelihood ratio test followed by FDR correction;  $n = 4$  biological replicates; data are presented as mean values  $\pm$  s.e.m.). **g**, Inverse correlation of NO at gene bodies and gene expression. NO is based on pseudobulk Strand-seq libraries from RPE-1. Gene bodies were scaled to the same length. **h**, Cell-typing based on NO at gene bodies (AUC = 0.96). Cell line codes: Blue, RPE-1; Purple, BM510; Magenta, C7; LV, latent variable. **i**, Receiver operating characteristics for inferring altered gene activity by analyzing NO at gene bodies, using pseudobulk Strand-seq libraries from in silico cell mixing.



**Fig. 2 | Linking subclonal SVs to their functional consequences in LCLs.**

**a**, Complex SVs in NA20509, with BFB-mediated rearrangements (17p) and a terminal dispersed duplication (5q) present with CF = 85%, shown for representative single cells. Ref. cells lacking complex SVs; InvDup, inverted duplication; terDel, terminal deletion. Reads are mapped to the W (orange) or C (green) strand. Gray, single cell IDs. **b**, Heatmap of 18 genes with altered gene activity amongst subclones, based on scNOVA ('17p-BFB', SV subclone; '17p-Ref', 17p not rearranged). Asterisks denote TF targets of c-Myc and Max. **c**, Gene set over-representation analysis for TF target genes showing significant enrichment of c-Myc and Max targets in the 17p-BFB subclone.  $P_{adj}$ , adjusted P value. Right, Model for c-Myc/Max target activation in NA20509 based on scNOVA, combined with previous knowledge. **d**, Mean RNA-seq expression Z scores of c-Myc/Max target genes across 33 LCLs. **e**, Fishplot showing CF changes over long-term

culture from 23.3% (7 of 33 cells; p4) to 100% (30 of 30 cells; p8). **f**, qPCR verifies clonal expansion of the BFB clone in p8 compared with p4 ( $P$  value based on FDR-corrected two-sided unpaired  $t$ -tests;  $n = 3$ ). HG1505, control cell line with a somatically stable *MAP2K3* locus. Note that for both NA20509 and HG1505 the germline copy number of the *MAP2K3* locus was consistently estimated to be three. Data are presented as mean values  $\pm$  s.e.m. **g**, RNA-seq shows significant increase of *MAP2K3* at p8 versus p4 (FDR-corrected two-sided Wald test, based on DESeq2;  $n = 5$  and three biological replicates for p4 and p8, respectively). **h**, Mean RNA expression Z scores of c-Myc/Max target genes in NA20509 (differences between p4 and p8 were evaluated using a two-sided Wilcoxon rank sum test;  $n = 5$  and three biological replicates for p4 and p8, respectively). Boxplot was defined by minima = 25th percentile - 1.5  $\times$  IQR, maxima = 75th percentile + 1.5  $\times$  IQR, center = median and bounds of box = 25th and 75th percentile (**g-h**).

### Functional outcomes of SVs in cell lines

To test this, we set out to investigate the functional outcomes of somatic SV landscapes in a panel of LCL samples<sup>38</sup> ( $N = 25$ ) from the 1000 Genomes Project<sup>39</sup> (1KGP). Single-cell SV discovery in 1,372 Strand-seq libraries generated for this panel (Supplementary Table 1) discovered 205 somatic SVs, with 24 of 25 (96%) LCLs showing at least one SV subclone—a sevenfold increase compared to a previous report<sup>40</sup> (Supplementary Table 3 and Supplementary Data). Of all the cell lines, 13 (52%) contained an SV subclone above 10% CF. This included the widely used NA12878 cell line<sup>34,39</sup>, in which we discovered a subclonal

500 kb deletion at19q13.12 (CF = 21%) that was mutually exclusive with two 22q11.2 deletions seen at CFs of 21% and 57%, respectively (Supplementary Figs. 9 and 10). The 22q11.2 SVs mapped to the well-known site of IGL recombination occurring during normal B cell development<sup>41</sup>. We hence focused on the 19q13.12 event, which resulted in the loss of a copy of *ZNF382*—a tumor suppressor and repressor of c-Myc<sup>42</sup>. Application of scNOVA measured significantly increased activity of *ERCC6*—a target gene of the c-Myc/Max transcription factor (TF) dimer<sup>43</sup>—and decreased activity of *PIEZO2* and *TRAPPC9*, in cells harboring this deletion (10% FDR; Supplementary Table 2).

To validate these findings, we reanalyzed Fluidigm and Smart-seq single-cell RNA-seq (scRNA-seq) datasets generated for NA12878 (refs. 44,45). We employed several established tools for SCNAs discovery from scRNA-seq data<sup>46–48</sup> (Supplementary Table 4), all of which failed to discover any of the SV subclones seen in this cell line (Supplementary Table 4). Yet, upon directly inputting the respective SV breakpoint coordinates into the CONICsmat tool<sup>46</sup>, we succeeded in identifying the 19q13.12 deletion (denoted ‘19q-Del’) through ‘targeted SCNA recalling’. We next pursued differential gene expression analyses by scRNA-seq, comparing 19q-Del cells to unaffected (‘19q-Ref’) cells, and verified overexpression of *ERCC6* in 19q-Del cells (10% FDR; Supplementary Fig. 10). For *PIEZO2* and *TRAPPC9*, the scRNA-seq-based expression trends were consistent with scNOVA (Supplementary Fig. 10), but did not reach the FDR threshold. A search for the over-represented TF targets amongst the differentially active genes identified c-Myc and Max as the most over-represented TFs in 19q-Del cells (10% FDR; Supplementary Fig. 10). These results indicate that scNOVA can functionally characterize SVs inaccessible to scRNA-seq-based SCNA discovery.

We next focused on NA20509, the LCL with the most abundant SV subclone (85% CF). Somatic SVs in NA20509 arose primarily through the breakage-fusion-bridge-cycle (BFB) process<sup>24,49</sup> involving a 49 Mb terminal duplication on 5q, and a 2.5 Mb inverted duplication on 17p with an adjacent terminal deletion (terDel) (Fig. 2a). The 5q and 17p segments became fused into a derivative chromosome of around 115 Mb (Supplementary Fig. 13), which probably stabilized the BFB. We searched for global gene activity changes in this ‘17p-BFB’ subclone compared with the nonrearranged cells (‘17p-Ref’) and identified 18 dysregulated genes (Fig. 2b). Testing for gene set over-representation<sup>50</sup> (Methods) revealed an enrichment of the target genes of c-Myc/Max heterodimers (10% FDR; Fig. 2c), that is, the same TFs we observed in the 19q-Del subclone in NA12878. Consistent with this, we identified somatic copy-number gain of *MAP2K3*, which encodes a gene activating c-Myc/Max<sup>51</sup>, resulting from the BFB (Fig. 2a).

We performed several orthogonal analyses to validate these findings. First, we verified all somatic SVs using deep WGS data generated for the IKGp sample panel<sup>52</sup> (Supplementary Fig. 13). Second, we analyzed RNA-seq data<sup>38</sup> for this LCL panel, which revealed that NA20509 exhibits the highest *MAP2K3* expression and the highest c-Myc/Max target expression (Supplementary Fig. 14 and Fig. 2d). Third, we followed the 17p-BFB subclone in culture, by subjecting early (p4) and late passage (p8) cells to Strand-seq, which revealed outgrowth of the 17p-BFB subclone (CF = 23% at p4, CF = 100% at p8;  $P < 0.00001$ , Fisher’s exact test; Fig. 2e), suggesting these cells have a proliferative advantage. Quantitative real-time PCR experiments verified this clonal outgrowth pattern (Fig. 2f).

Since the functional impact of SVs on clonal expansion is unexplored in LCLs, we more deeply characterized the molecular phenotypes of 17p-BFB cells by pursuing RNA-seq in p4 and p8 cultures. We observed increased *MAP2K3* expression (1.39-fold, 10% FDR) at p8,

consistent with *MAP2K3* dysregulation as a result of copy-number gain in the 17p-BFB subclone (Fig. 2g and Supplementary Note). Pathway-level analysis showed deregulation of c-Myc/Max target genes following clonal expansion ( $P = 0.036$ ; Wilcoxon rank sum test; Fig. 2h and Supplementary Fig. 14). Collectively, these data link the outgrowth of SV subclones to the deregulation of c-Myc/Max targets, which could represent a common driver of clonal expansion in LCLs.

### Local effects of copy-balanced driver SVs in leukemia

To deconvolute the effects of driver SVs in patients, we applied scNOVA to analyze the local consequences of balanced SVs, which are widespread in leukemia<sup>3,53</sup>. We analyzed primary cells from a patient with acute myeloid leukemia (AML) (32-year-old male; patient-ID = AML\_1) bearing a balanced t(8;21) translocation that results in *RUNX1-RUNX1T1* gene fusion<sup>54</sup>. We sorted CD34<sup>+</sup> cells from AML\_1 (Supplementary Fig. 15), and sequenced 42 Strand-seq libraries. SV discovery revealed a 46,XY,t(8;21)(q22;q22) karyotype (Fig. 3a, Supplementary Fig. 16 and Supplementary Table 3) consistent with clinical diagnosis. We fine-mapped the translocation breakpoint to intron 1 of *RUNX1T1* and intron 5 of *RUNX1* (Supplementary Fig. 17), and subsequently identified haplotype-specific NO at 11 genes, genome-wide (10% FDR; Supplementary Table 2). This included *RUNX1T1*, which showed reduced NO on the derivative (H2) haplotype (Fig. 3b), consistent with increased gene activity mediated as a local effect of the translocation<sup>55</sup>. The remaining genes did not reside near a detected somatic SV, suggesting other factors (such as germline SNPs; Supplementary Fig. 17) may have affected their NO.

To systematically investigate potential local effects, we used a sliding window (Methods) to measure NO on both sides of the translocation breakpoint. We observed decreased NO, suggesting increased chromatin accessibility, from the breakpoint junction up to the respective nearest topologically associating domain (TAD) boundaries (Fig. 3c). This signal was most pronounced in an enhancer-rich region around 0.8 to 1.1 Mb upstream of *RUNX1* originating from chromosome 21 ( $P < 0.003$ ; likelihood ratio test, adjusted using permutations; Fig. 3c), found to physically interact with the *RUNX1* promoter in CD34<sup>+</sup> cells<sup>56</sup>. Within this segment, we identified two CREs with significantly reduced NO (10% FDR; Exact test) (Fig. 3d and Supplementary Table 5), which may foster *RUNX1-RUNX1T1* expression. Chromosome-wide analysis showed haplotype-specific NO patterns were restricted to the fused TAD (Fig. 3e,f), in line with these patterns resulting from the translocation.

We also revisited Strand-seq datasets with previously reported copy-neutral SVs, including the BM510 cell line in which copy-neutral interchromosomal SVs resulted in *TP53-NTRK3* gene fusion<sup>24</sup>. In agreement with the oncogenic role of *TP53-NTRK3* (ref. 24), scNOVA identified *NTRK3* upregulation as the only significant local effect (10% FDR), consistent with allele-specific *TP53-NTRK3* expression measured on the rearranged haplotype (Extended Data Fig. 5). Second, we revisited

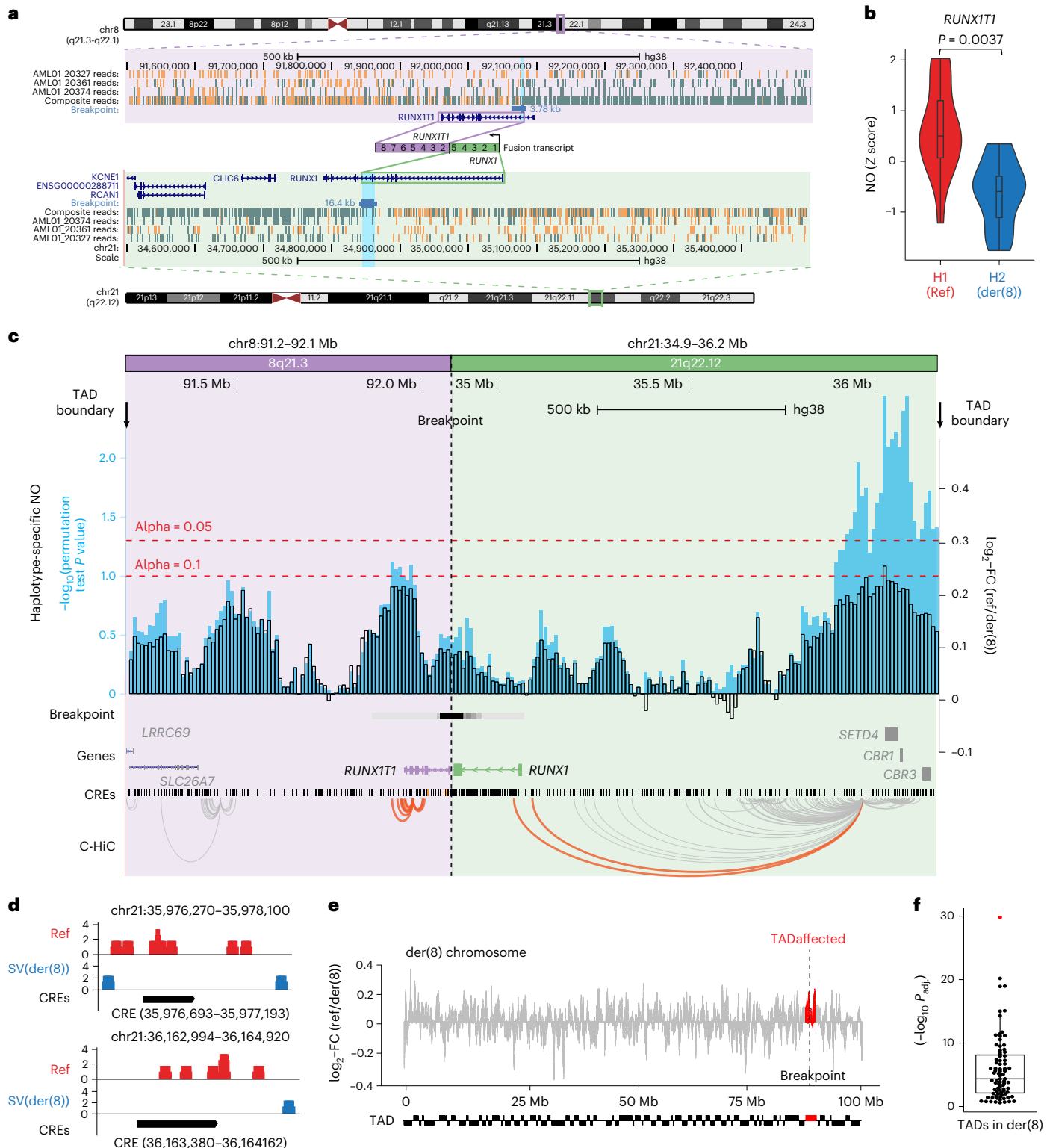
### Fig. 3 | Haplotype-specific NO analysis shows local effects of a copy-neutral driver SV in AML.

**a**, Balanced t(8;21) translocation in AML\_1, discovered based on strand cosegregation ( $P$  value = 0.00003 for translocation discovery using strand cosegregation<sup>24</sup>, FDR-adjusted Fisher’s exact test; Supplementary Fig. 16). The SV breakpoint was fine-mapped to the region highlighted in light blue. Composite reads shown were taken from all informative cells in which reads could be phased (WC or CW configuration; Methods). **b**, A violin plot demonstrates haplotype-specific NO at the *RUNX1T1* gene body (10% FDR; two-sided Wilcoxon rank sum test followed by Benjamini–Hochberg multiple correction;  $n = 17$  single cells; boxplot was defined by minima = 25th percentile – 1.5 × IQR, maxima = 75th percentile + 1.5 × IQR, center = median and bounds of box = 25th and 75th percentile), consistent with aberrant activity of the locus on der(8). **c**, Haplotype-specific NO around the SV breakpoint. FCs of haplotype-specific NO, measured between the *RUNX1-RUNX1T1* containing derivative chromosome (der(8)) and corresponding regions on the

unaffected homolog (Ref), are shown in black, and  $-\log_{10}(P$  values) in light blue. Enhancer-target gene physical interactions based on chromatin conformation capture<sup>36,93</sup> are depicted in orange (interactions involving *RUNX1* and *RUNX1T1*) and gray (involving other loci). **d**, Significant CREs located within the distal peak region, demonstrating haplotype-specific absence of NO on der(8) at 10% FDR, suggesting increased CRE accessibility on der(8). Within the segment around 0.8 to 1.1 Mb upstream of *RUNX1*, which showed pronounced haplotype-specific NO, we tested 69 CREs for haplotype-specific NO, which identified two significant CREs. **e**, Haplotype-specific NO measured between der(8) and corresponding regions of the unaffected homolog. Red, regions corresponding to the fused TAD. **f**, A beeswarm plot shows that the fused TAD (red) is an outlier in terms of haplotype-specific NO on der(8) ( $P$  values based on Kolmogorov–Smirnov tests;  $n = 83$  TADs in der(8); boxplot was defined by minima = 25th percentile – 1.5 × IQR, maxima = 75th percentile + 1.5 × IQR, center = median and bounds of box = 25th and 75th percentile).

a 2.6 Mb inversion mapping to 14q32 in a T-cell acute lymphoblastic leukemia (T-ALL) patient-derived xenograft (T-ALL\_P1)<sup>24</sup>. scNOVA discovered downregulation of *BCL11B*, a known haploinsufficient T-ALL tumor suppressor<sup>57</sup>, as a significant local effect of this balanced inversion, supporting allele-specific silencing of *BCL11B* on the rearranged haplotype as measured by RNA-seq<sup>24</sup> (Extended Data Fig. 6). These data collectively show that scNOVA allows linking balanced SVs to their local functional consequences—a functionality not provided by any previous single-cell multiomic method<sup>20</sup>.

**Dissecting functional effects of heterogeneous somatic SVs**  
We next set out to functionally dissect a leukemia sample with unknown genetic drivers, by characterizing B-cells from a 61-year-old patient with chronic lymphocytic leukemia (CLL) (CLL\_24)<sup>58</sup>. Analysis of 86 Strand-seq libraries revealed an unprecedented level of somatic SVs, with 11 different karyotypes represented by 13 SVs occurring in subclones with CFs of 1–5% (Supplementary Table 3). This vastly exceeds intrapatient diversity estimates for CLLs from the Pan-Cancer Analysis of Whole Genomes (PCAWG), where maximally three subclones



were reported<sup>59</sup>, highlighting how Strand-seq provides access to SVs escaping discovery by WGS<sup>3,24</sup>. Chromosome 10q showed especially pronounced subclonal heterogeneity; we identified seven partially overlapping deletions ranging from 2 to 31 Mb in size, and residing proximal to the fragile site *FRA10B*<sup>60</sup> (Fig. 4a and Supplementary Fig. 18). These SVs clustered into a 1.4 Mb ‘minimal segment’ at 10q24.32, arising independently from both haplotypes (Fig. 4b). While previous studies reported somatic 10q24.32 deletions in 1–4% of CLLs<sup>61–63</sup>, molecular analysis of this recurrent somatic SV has so far been lacking.

We first compared all cells bearing a 10q24.32 deletion (‘10q-Del’,  $N = 11$ ) to cells lacking such SV (‘10q-Ref’,  $N = 75$ ), hence disregarding the fine-scale subclonal structure of CLL\_24, and predicted 115 dysregulated genes (Fig. 4c and Supplementary Table 2). Next, we performed molecular phenotype analysis using MsigDB<sup>64</sup> (Methods), which revealed that 10q-Del cells exhibit increased activity in several leukemia-relevant signaling pathways, including Wnt, c-Met (a pathway promoted by Wnt signaling<sup>65</sup>), B cell receptor (BCR) signaling, phosphatidylinositol (3,4,5)-trisphosphate (PIP3) signaling and the CREB pathway (10% FDR; Fig. 4d). RNA-seq data available for 178 CLLs<sup>62</sup> and stratified by 10q24.32 status, revealed upregulation of Wnt and c-Met signaling—but not of BCR, PIP3 and CREB signaling—in CLLs exhibiting 10q24.32 deletions (10% FDR; CLLs with 10q-Del:  $N = 4$ ; 10q-Ref:  $N = 174$ ; Fig. 4e and Supplementary Fig. 24). These data therefore suggest a link between 10q24.32 deletion and the promotion of Wnt signaling.

We further tested whether the different 10q-Del events seen in CLL\_24 subclones have led to distinct functional outcomes, focusing on three subclones represented by at least two cells: ‘SCa’, showing one interstitial deletion directly at the minimal segment; ‘SCb’, harboring a terDel, with the breakpoint located at the minimal segment boundary and ‘SCc’, containing two interstitial deletions, at the minimal segment and at 10q23.31 (Fig. 4b and Supplementary Table 3). Molecular phenotype analysis of each subclone identified 109, 206 and 266 differentially active genes, respectively (Supplementary Table 2), with the most pronounced levels of Wnt upregulation in SCb and SCc (Fig. 4f). SCb showed the highest activation of c-Met, BCR and PIP3 signaling, whereas CREB signaling was highest in SCc (Supplementary Fig. 21). This suggests that deletion location and length at 10q24.32 affect their molecular consequences, and furthermore illustrates the ability of scNOVA to predict molecular differences in subclones represented by as few as two cells.

To more deeply characterize the CLL\_24 subclones, we generated CITE-seq (cellular indexing of transcriptomes and epitopes by single-cell sequencing) data, which couples scRNA-seq with protein surface marker measurements<sup>66</sup>. Again, we attempted SCNA discovery in the scRNA-seq data, which failed to detect any SCNAs, or subclones, in CLL\_24 (Supplementary Table 4). However, targeted SCNA recalling<sup>46</sup> identified 82 CITE-seq cells harboring the greater than 31 Mb 10q-terDel of SCb (‘10q-terDel’), whereas the deletions in SCa (2.2 Mb) and SCc (2.1 Mb and 1.9 Mb, respectively) escaped detection (Extended Data

Fig. 7 and Supplementary Notes). Having recovered the SCb subclone in the CITE-seq data, we performed single-cell gene set enrichment analysis<sup>67</sup> (Methods), which verified that all pathways inferred by scNOVA (Wnt, c-Met, BCR, PIP3 and CREB) are upregulated in 10q-terDel cells (Fig. 4d,g). A gene regulatory network analysis<sup>68</sup> comparing 10q-terDel with 10q-Ref cells identified 43 differentially active TFs (FDR 10%; Fig. 4h) and a functional enrichment analysis<sup>69</sup> showed over-representation of Wnt signaling, BCR signaling and the PD-1 checkpoint pathway (Supplementary Table 16 and Fig. 4h); the PD-1 checkpoint pathway has been linked to immune resistance and transformation of CLL to aggressive lymphoma<sup>70,71</sup>. Since somatic lesions mediating PD-1 expression in CLL have remained elusive, we used the CITE-seq data to analyze PD-1 protein expression, which demonstrated upregulation of PD-1 in 10q-terDel-containing cells as the only significant hit at the protein level (Fig. 4i). Notably, *NFATC1*, a TF predicted to be differentially active by both scNOVA and CITE-seq, regulates Wnt<sup>72</sup>, PIP3 (refs. <sup>73,74</sup>), CREB<sup>75</sup> and BCR signaling<sup>76</sup> as well as PD-1 expression<sup>77</sup>, and thus may contribute to global pathway dysregulation in CLL\_24. Our analysis reveals subtle pathway activities of somatic deletions present at low CF (Fig. 4f,j), and collectively implicates 10q24.32 deletions in dysregulated Wnt signaling—a crucial pathway for CLL pathogenesis<sup>78</sup>.

### Functional characterization of subclonal chromothripsis

While chromothripsis is a widespread mutational process in cancer<sup>3,4,22</sup>, this process is not ascertained by previous single-cell multiomic methods, and its molecular outcomes remain largely elusive<sup>3,79</sup>. We previously discovered a subclonal chromothripsis event<sup>24</sup> in T-ALL\_P1 that affects most of 6q (denoted ‘6q-CT’; CF = 30%) (Fig. 5a and Supplementary Table 3); however, the consequences of this complex rearrangement were uncharacterized. Using scNOVA, we identified 12 genes with differential NO between 6q-CT and 6q-Ref cells (denoted the ‘CT gene signature’; 10% FDR; Fig. 5a,b and Supplementary Table 2). A closer analysis showed 27 TF genes overlapping the chromothriptic region (Fig. 5a). Gene set over-representation testing using the target genes of these TFs revealed that c-Myb, product of the *MYB* oncogene, was significantly enriched among the genes included in the CT gene signature (10% FDR; adjusted  $P = 0.00015$ ; Fig. 5b,c and Supplementary Table 6). The *MYB* gene is located within a region that was duplicated (and inverted) as a result of 6q-CT, suggesting a potential dosage effect (Fig. 5a). Corroborating these predictions, we performed RNA-seq in a panel of 13 T-ALLs, amongst which T-ALL\_P1 showed the highest expression of c-Myb targets (Fig. 5d and Supplementary Table 7). We also verified that *MYB* is allele-specifically expressed from the SV-affected haplotype ( $P = 0.0317$ ; likelihood ratio test; Supplementary Fig. 30), which together nominates *MYB* as a candidate driver gene dysregulated as a consequence of 6q-CT.

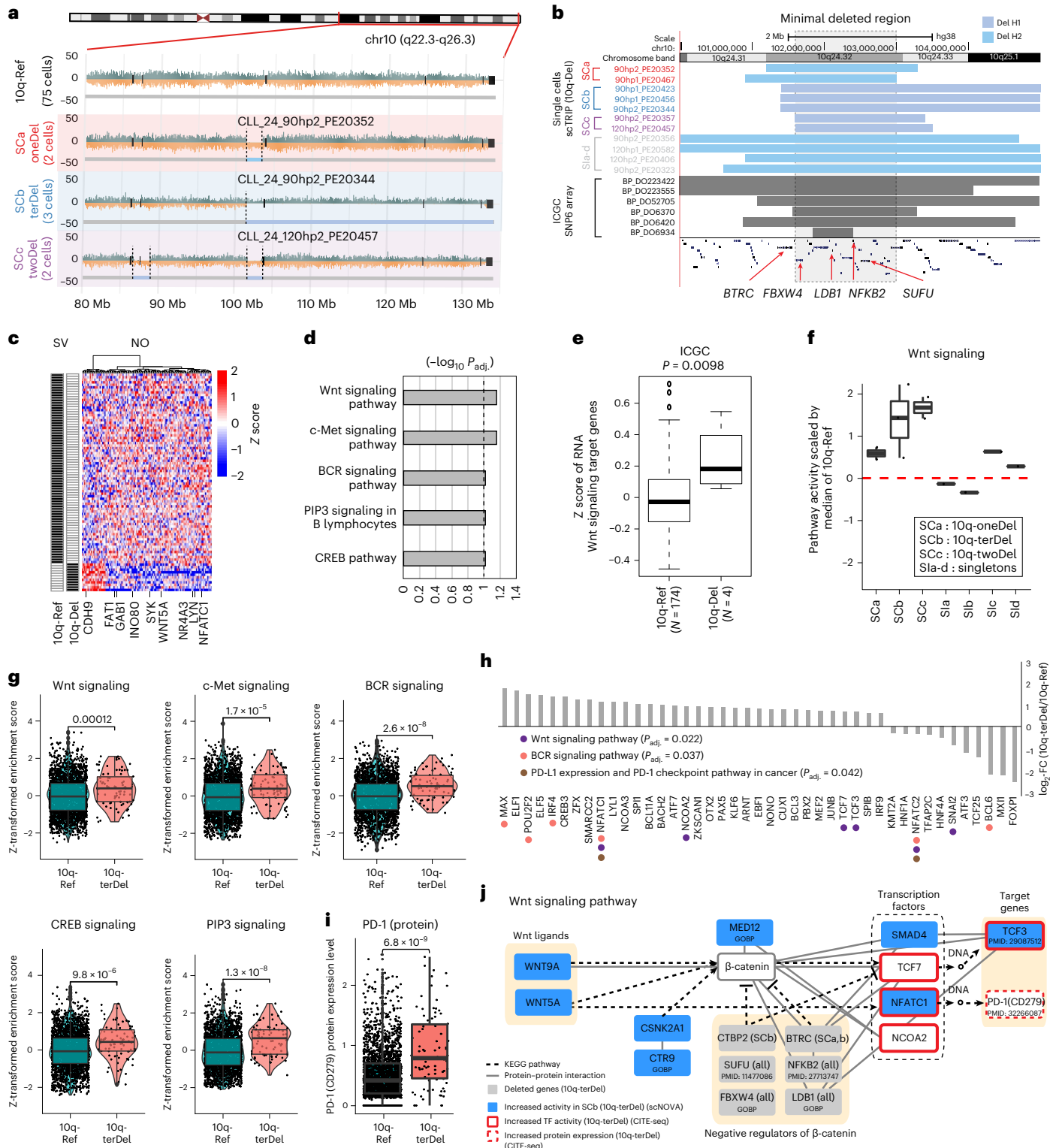
To more deeply characterize this sample, we generated scRNA-seq data for T-ALL\_P1 (5,504 cells; Fig. 6a). Since scRNA-seq-based SCNAs discovery<sup>46–48</sup> missed the 6q-CT event (Supplementary Table 4), we

**Fig. 4 | Deconvoluting consequences of subclonal SV heterogeneity in a CLL primary sample.** **a**, Single-cell SV discovery in CLL\_24. All cells exhibiting deletions (10q-Del) shown in Supplementary Fig. 18. 10q-Ref, cells bearing a not rearranged 10q. **b**, Minimal deleted region (chr10:101615000-103028000; hg38), displaying recurrent deletions in a separate cohort of CLLs<sup>62</sup>. **c**, Heatmap of genes with altered activity in 10q-Del based on scNOVA (alternative mode; 10% FDR). Genes from all significant pathways reported in **d** are highlighted. **d**, Pathway modules with differential activity, in cells exhibiting 10q-Del (10% FDR). **e**, Minimal deleted region-bearing CLL samples from the International Cancer Consortium (ICCG) demonstrate overexpression of Wnt signaling genes compared with 10q-Ref ( $P = 0.0098$ ; two-sided likelihood ratio test;  $n = 174$  and  $n = 4$  independent CLL samples for 10q-Ref and 10q-Del, respectively). **f**, Pathway activities ( $(-1) \times Z$  score of NO) derived from jointly modeled NO at the gene bodies of Wnt signaling pathway genes for each SV-bearing CLL\_24 cell. Sla-Sld correspond to single cells exhibiting a deletion at 10q24 not shared by any other cell.  $n = 2, 3, 2$  and 1 cells

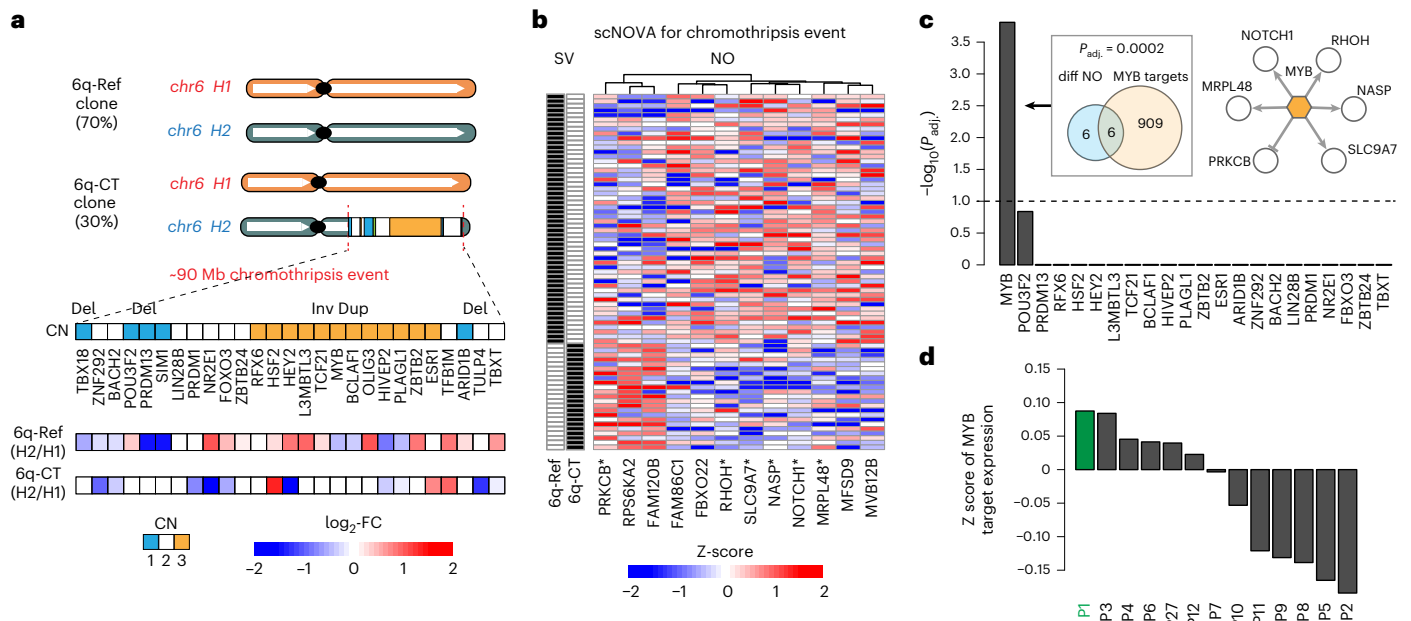
are depicted in the plot for SCa, SCb, SCc and Sla-Sld, respectively. **g**, Single-cell gene set enrichment scores for five leukemia-related pathways from CITE-seq. Enrichment scores for 10q-terDel ( $n = 82$ ) and 10q-Ref ( $n = 2,381$ ) cells were compared using two-sided  $t$ -tests. **h**, Chart depicting 43 differentially active TFs between 10q-terDel and 10q-Ref cells based on DoRothEA<sup>68</sup>. Genes involved in the pathways over-represented by these TFs are annotated using colored dots. **i**, Differentially expressed surface protein CD279 (PD-1) in 10q-terDel ( $n = 82$ ) compared with 10q-Ref ( $n = 2,381$ ) cells based on a two-sided Wilcoxon rank sum test. **j**, Wnt pathway diagram showing the altered genes or TFs in SCb (10q-terDel) identified by scNOVA (blue nodes) and CITE-seq (red borders). Gray, known (see PubMedIDs) and computationally predicted regulators (based on Gene Ontology Biological Process (GOBP)) of Wnt signaling that are deleted in SCb. Throughout the figure, boxplots were defined by minima = 25th percentile – 1.5 × IQR, maxima = 75th percentile + 1.5 × IQR, center = median and bounds of box = 25th and 75th percentile.

again performed targeted SCNA recalling (Supplementary Notes) generating confident calls for 838 (around 15%) cells in the scRNA-seq dataset (the remaining 4,666 cells lacked a confident assignment; 'NA'). Out of these 838 cells, 729 were predicted to harbor the 6q-CT event, and 109 were called 6q-Ref. Unsupervised clustering<sup>80</sup> of the scRNA-seq data stratified by 6q status (Methods) revealed that 6q-CT cells (as predicted through targeted recalling) were enriched in two expression clusters (clusters 3 and 7;  $P = 3.43 \times 10^{-5}$  and  $1.15 \times 10^{-3}$ ; FDR-adjusted Fisher's exact test; Fig. 6d and Supplementary Fig. 34),

in line with a distinctive expression profile. To corroborate this, we applied UCell<sup>81</sup> to assign cells into '6q-CT' or '6q-Ref' based on the CT gene signature, which confirmed enrichment of 6q-CT in clusters 3 and 7 (Fig. 6c,d;  $P = 3.39 \times 10^{-38}$  and  $P = 2.15 \times 10^{-4}$ ; FDR-adjusted Fisher's exact test). Trajectory analysis<sup>82</sup> showed the 6q-CT cells (as defined by UCell) were enriched for DNearly (double-negative early;  $P = 2.78 \times 10^{-13}$ ), DNQ (double-negative quiescent;  $P = 1.27 \times 10^{-5}$ ) and DPP (double-positive proliferating;  $P = 1.88 \times 10^{-7}$ ) T cells (FDR-corrected Fisher's exact tests; Fig. 6b and Supplementary







**Fig. 5 | scNOVA identifies functional effects of a subclonal chromothripsis event.** **a**, The 27 TF genes located in a segment that underwent chromothripsis<sup>24</sup> on 6q in T-ALL\_P1. Haplotype-specific NO measurements, which scNOVA generated for CREs assigned to the nearest genes, are depicted below. FC of normalized haplotype-specific NO is shown for each subclone. 6q-CT, subclone bearing chromothripsis on 6q; 6q-Ref, subclone bearing a not rearranged chromosome 6. **b**, Heatmap of 12 genes with differential activity between subclones in T-ALL\_P1, based on scNOVA (denoted CT gene signature). Asterisks

denote TF targets highlighted in **c**, TF target over-representation analyses for CT gene signature, revealing c-Myb as the only significant hit. Venn diagram depicting enrichment of c-Myb targets ( $P$  value based on an FDR-adjusted hypergeometric test). Upper right, network with c-Myb and its target genes based on scNOVA, combined with previous knowledge. **d**, Mean  $Z$  scores of c-Myb target gene expression measured by bulk RNA-seq in a panel of 13 T-ALL-derived samples. T-ALL\_P1 (P1) exhibited the overall highest expression of c-Myb targets.

Fig. 35), and depleted of mature CD4<sup>+</sup> T cells ( $P = 1.45 \times 10^{-11}$ , Supplementary Fig. 35). This suggests a potential differentiation block at the progenitor stage as a result of 6q-CT and, more generally, that 6q-CT cells bear a distinctive molecular phenotype as a result of the chromothriptic rearrangements.

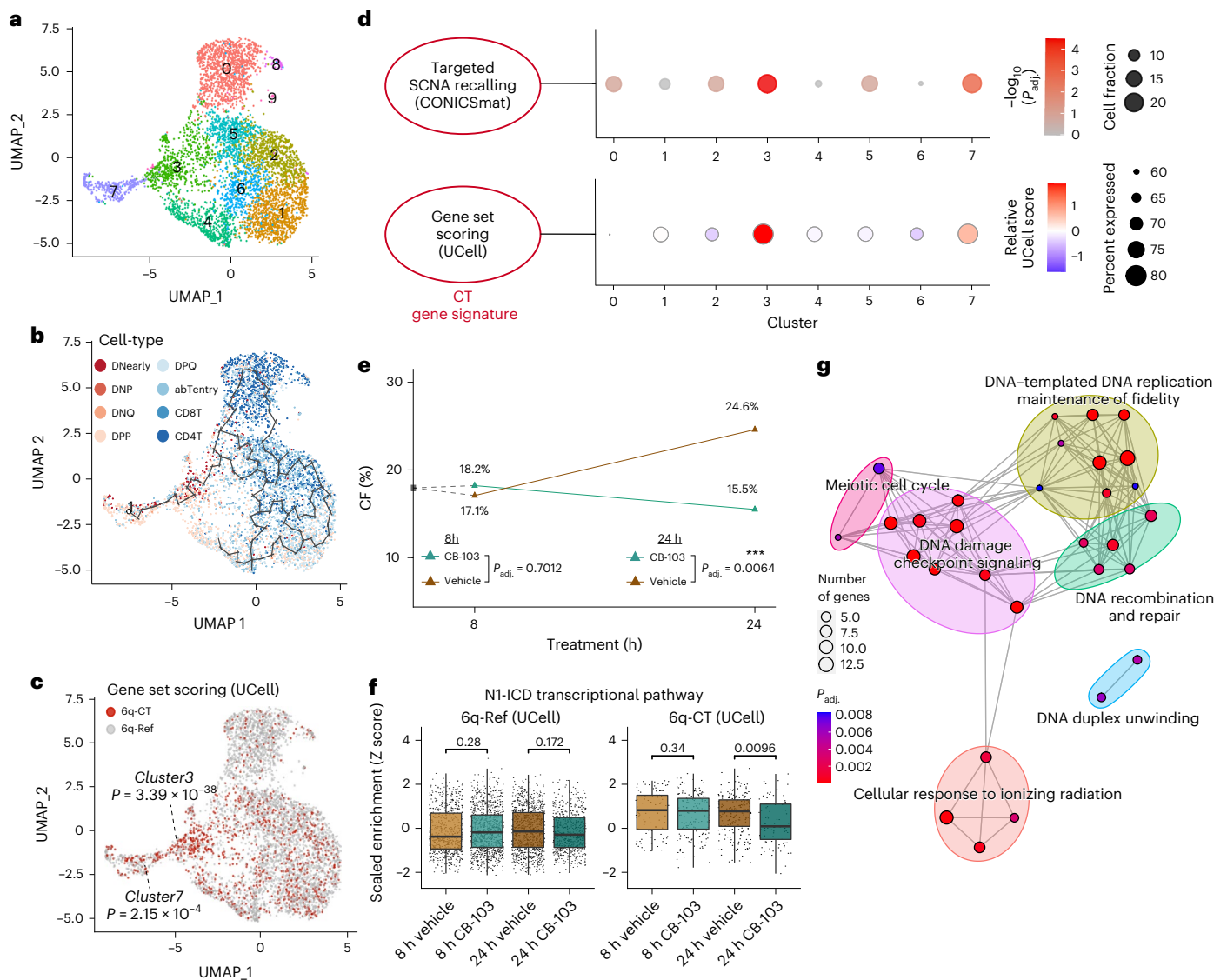
Having identified c-Myb pathway activation as a consequence of 6q-CT in TALL\_P1, we hypothesized this molecular phenotype could guide drug targeting in cell culture. We selected *NOTCH1* as a suitable candidate for targeting this subclone because this c-Myb target was (1) inferred by scNOVA to be highly upregulated in 6q-CT cells (Fig. 5b) and (2) is targetable by different compounds and strategies<sup>83</sup>. We treated T-ALL\_P1 cell cultures with the CB-103 pan-NOTCH small-molecule inhibitor (targeting the *Notch1* intracellular domain (NI-ICD)<sup>84,85</sup>) or a vehicle control for 8 h and 24 h (Methods). Using scRNA-seq (3,663 single cells) to analyze drug response patterns, we inferred 6q-CT and 6q-Ref cells at each timepoint by transferring the cell annotation labels from the untreated (reference) sample with Seurat<sup>80</sup> (Fig. 6c and Supplementary Fig. 37). After 24 h in culture, vehicle-treated T-ALL\_P1 cells showed a 45% relative increase in the 6q-CT subclone compared to 8 h (CF of 17.1% to 24.6%;  $P = 0.0180$ ; FDR-adjusted Fisher's exact test), indicating that 6q-CT cells expanded clonally. By contrast, upon CB-103 treatment, the CF of the 6q-CT subclone was reduced at 24 h (to CF = 15.5%;  $P = 0.0064$ ; Fig. 6e and Supplementary Fig. 38), indicating that 6q-CT cells were preferentially lost with NI-ICD inhibition. Additionally, we observed specific depletion of the REACTOME NI-ICD gene set only in 6q-CT cells after 24 h of CD-103 treatment, consistent with specific subclone targeting ( $P = 0.0096$ ; FDR-adjusted Wilcoxon rank sum test; Fig. 6f and Supplementary Fig. 39). These results highlight the potential of scNOVA to functionally characterize highly complex classes of DNA rearrangement (that is, chromothripsis events), and to clinically target subclones bearing complex cancer driver SVs.

## Discussion

The functional characterization of SVs is of critical importance for precision oncology<sup>1–3</sup>. Our method characterizes a wide spectrum of SV classes<sup>24</sup>, and couples these with NO analysis to link somatic SVs to local or global gene activity changes. Accounting for balanced SVs, scNOVA allows the investigation of copy-number stable (that is, euploid) malignancies previously inaccessible to single-cell multiomics<sup>3,20</sup> (Supplementary Table 12). Strand-seq derived SCNA calls were far better resolved compared to scRNA-seq based calls (Supplementary Table 4), suggesting a more limited utility of scRNA-seq data for discovering SCNA drivers in cancer, with the exception of malignancies displaying extremely high levels of chromosomal instability with particularly large-scale SCNAs<sup>3,86</sup>.

We uncovered unprecedented karyotypic diversity in a CLL sample, comprising distinct deletions at 10q24.32, which we link to leukemia-related signaling pathways, particularly Wnt signaling. Read depth based profiling of SCNAs is prone to underreport such subclonal structural diversity<sup>3</sup>. Enrichment of cases bearing 10q24.32 deletions amongst relapsed/refractory and high-risk CLL<sup>87</sup> suggests a potential role of Wnt pathway dysregulation mediated through 10q24.32 in disease progression. Whether the *FRA10B* fragile site is involved in the formation of these deletions remains to be seen and requires larger cohorts. Interestingly, CLL\_24 exhibits a SNP (rs118137427; 3.7% allele frequency in Europeans) within *FRA10B* associated with the acquisition of 10q-terDel in normal blood<sup>88</sup>. Based on the PCAWG resource comprising 94 CLLs<sup>2</sup>, rs118137427 is seen in 2 out of 4 (50%) CLLs with 10q24.32 deletions, but in only 6 of 90 (6.7%) CLLs with 10q-Ref ( $P = 0.035$ ; Fisher's exact test), suggesting a possible link between SNPs at *FRA10B* and ITH in leukemia that warrants future investigation.

Our framework readily functionally characterizes complex rearrangements previously inaccessible to single-cell multiomics<sup>3</sup>. Complex somatic SVs are prevalent in cancer and linked with aggressive



**Fig. 6 | Targeting the chromothriptic subclone in cell culture.** **a**, Uniform manifold approximation and projection (UMAP) of scRNA-seq data showing ten unsupervised clusters in T-ALL\_P1. **b**, Overlay of gene set-derived cell-type annotation and inferred lineage trajectory onto this UMAP. **c**, Single cell types whose expression profiles matched the CT gene signature (gene set UCell score > (median score + s.d.)) are assigned to ‘6q-CT’ and shown in red; the remaining cells did not meet the threshold for the CT gene signature (assigned ‘6q-Ref’ status). *P* values depict enrichment of 6q-CT cells in clusters 3 and 7. **d**, Significant enrichment of 6q-CT cells in clusters 3 and 7 based on scRNA-seq. Upper panel, dot plot showing the significance of over-representation of 6q-CT calls in scRNA clusters based on targeted SCNA recalling (*P* values based on FDR-adjusted Fisher’s exact tests). Lower panel, gene set-level expression summary for the CT gene signature, derived using UCell<sup>81</sup> with the directionality of expression changes taken into account. **e**, CF of 6q-CT cells after treatment with Notch inhibitor CB-103 (green) and vehicle control (brown) along a time course 8 h before and 24 h after treatment. CF was estimated by transferring

gene set based CT annotations obtained from the scRNA-seq of T-ALL\_P1 before treatment to the scRNA-seq of T-ALL\_P1 after treatment. Changed CF (%) at 24 h compared with 8 h is shown in the plot on top of the 24 h datapoints. For each timepoint, the difference of CF under vehicle and CB-103 was evaluated by Fisher’s exact test (results are based on pairwise comparisons). **f**, Scaled enrichment scores obtained by single-cell gene set enrichment analysis for the ‘N1-ICD transcriptional pathway’ gene set. Scores across treatment conditions (vehicle versus CB-103) were compared using two-sided Wilcoxon rank sum tests. (Boxplot was defined by minima = 25th percentile – 1.5 × IQR, maxima = 75th percentile + 1.5 × IQR, center = median and bounds of box = 25th and 75th percentile; *n* = 665, 978, 915 and 556 cells for 6q-Ref from 8 h Vehicle, 8 h CB-103, 24 h Vehicle and 24 h CB-103; *n* = 91, 157, 213 and 88 cells for 6q-CT for each condition, respectively). **g**, Network representation of GOBPs enriched by differentially expressed genes in 6q-CT compared with 6q-Ref cells under CB-103 treatment (24 h), subtracting any genes not specific to the drug treatment.

tumor phenotypes<sup>2,3,22</sup> underlining significant potential of scNOVA for the comprehensive functional characterization of cancer cells. Since scNOVA does not require coupling distinct experimental modalities in each individual cell, it overcomes important methodological challenges<sup>20</sup>, including data sparseness and higher costs from generating data for more than one modality<sup>20,89</sup>. Additionally, the coverage achieved by Strand-seq enables the analysis of haplotype-specific NO

along the entire genome (Supplementary Fig. 41), providing advantages over classical allele-specific analyses that are restricted to regionally phased SNPs<sup>15</sup>.

Nonetheless, important challenges remain, and the full spectrum of mutations arising in an individual cell is likely to remain inaccessible to a single method in the foreseeable future. Strand-seq does not capture SVs less than 200 kb that more rarely acts as cancer drivers<sup>2</sup>.

Additionally, while scNOVA infers differentially active genes, it does not span the same dynamic expression range as scRNA-seq (Supplementary Table 12). This suggests that pairing scNOVA with targeted SCNA calling by scRNA-seq can provide added value by allowing variants outside the detection range of other methods to be characterized. Finally, Strand-seq requires dividing cells for BrdU labeling<sup>23</sup> (Fig. 1a), and is therefore not applicable for nondividing cells or fixed samples. However, it can be used for dividing cells in organoids, primary fresh frozen progenitor cells, cells in regenerating tissues and cancer samples amenable to culture. Our study used cell lines for benchmarking followed by proof-of-principle application in patient samples. Generalization of these analyses to larger cohorts will allow systematic investigation of the roles subclonal SVs play in leukaemogenesis.

We foresee a wide variety of potential future applications. Our framework offers potential for studies on the determinants and consequences of chromosomal instability in cancer, and may promote research into the interplay of genetic and nongenetic cancer determinants<sup>20</sup>. It likewise could be used to advance surveys of precancerous lesions<sup>3,90</sup>. Additionally, scNOVA may offer value in precision oncology by exposing subclonal driver alterations along with their targetable functional outcomes, to target cancer subclones in patients. Furthermore, SVs can accidentally arise in key model cell lines, as we demonstrate for widely used LCLs, and the features of scNOVA are ideally suited to functionally characterize unwanted heterogeneity in such samples. Unwanted somatic SVs also arise as a by-product of CRISPR-Cas9 genome editing, which generates micronuclei and chromosome bridges in human primary cells, structures that initiate the formation of chromothripsis<sup>91</sup>. scNOVA could promote the safety of therapeutically relevant genome editing in the future, by enabling the simultaneous detection and functional characterization of such potentially pathogenic editing outcomes.

In summary, scNOVA moves directly from SV landscapes to their functional consequences in heterogeneous cell populations. By making a broad spectrum of somatic SVs accessible for functional characterization genome-wide, this single-cell multiomic framework serves as a foundation for deciphering the impact of somatic rearrangement processes in cancer.

## Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41587-022-01551-4>.

## References

- Priestley, P. et al. Pan-cancer whole-genome analyses of metastatic solid tumours. *Nature* **575**, 210–216 (2019).
- ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analysis of whole genomes. *Nature* **578**, 82–93 (2020).
- Cosenza, M. R., Rodriguez-Martin, B. & Korbil, J. O. Structural variation in cancer: role, prevalence, and mechanisms. *Annu. Rev. Genomics Hum. Genet.* **23**, 123–152 (2022).
- Stephens, P. J. et al. Massive genomic rearrangement acquired in a single catastrophic event during cancer development. *Cell* **144**, 27–40 (2011).
- Rausch, T. et al. Genome sequencing of pediatric medulloblastoma links catastrophic DNA rearrangements with TP53 mutations. *Cell* **148**, 59–71 (2012).
- Baca, S. C. et al. Punctuated evolution of prostate cancer genomes. *Cell* **153**, 666–677 (2013).
- Umbreit, N. T. et al. Mechanisms generating cancer genome complexity from a single cell division error. *Science* **368**, eaba0712 (2020).
- Hadi, K. et al. Distinct classes of complex structural variation uncovered across thousands of cancer genome graphs. *Cell* **183**, e32 (2020).
- Minussi, D. C. et al. Breast tumours maintain a reservoir of subclonal diversity during expansion. *Nature* **592**, 302–308 (2021).
- Watkins, T. B. K. et al. Pervasive chromosomal instability and karyotype order in tumour evolution. *Nature* **587**, 126–132 (2020).
- Viswanathan, S. R. et al. Structural alterations driving castration-resistant prostate cancer revealed by linked-read genome sequencing. *Cell* **174**, e19 (2018).
- McPherson, A. et al. Divergent modes of clonal spread and intraperitoneal mixing in high-grade serous ovarian cancer. *Nat. Genet.* **48**, 758–767 (2016).
- Weischenfeldt, J. et al. Pan-cancer analysis of somatic copy-number alterations implicates IRS4 and IGF2 in enhancer hijacking. *Nat. Genet.* **49**, 65–74 (2016).
- Liu, Y. et al. Discovery of regulatory noncoding variants in individual cancer genomes by using cis-X. *Nat. Genet.* **52**, 811–818 (2020).
- PCAWG Transcriptome Core Group. et al. Genomic basis for RNA alterations in cancer. *Nature* **578**, 129–136 (2020).
- Northcott, P. A. et al. Enhancer hijacking activates GF11 family oncogenes in medulloblastoma. *Nature* **511**, 428–434 (2014).
- Dey, S. S., Kester, L., Spanjaard, B., Bienko, M. & van Oudenaarden, A. Integrated genome and transcriptome sequencing of the same cell. *Nat. Biotechnol.* **33**, 285–289 (2015).
- Macaulay, I. C. et al. G&T-seq: parallel sequencing of single-cell genomes and transcriptomes. *Nat. Methods* **12**, 519–522 (2015).
- Yin, Y. et al. High-throughput single-cell sequencing with linear amplification. *Mol. Cell* **76**, e10 (2019).
- Nam, A. S., Chaligne, R. & Landau, D. A. Integrating genetic and non-genetic determinants of cancer evolution by single-cell multi-omics. *Nat. Rev. Genet.* **22**, 3–18 (2020).
- Nam, A. S. et al. Somatic mutations and cell identity linked by genotyping of transcriptomes. *Nature* **571**, 355–360 (2019).
- Cortés-Ciriano, I. et al. Comprehensive analysis of chromothripsis in 2,658 human cancers using whole-genome sequencing. *Nat. Genet.* **52**, 331–341 (2020).
- Sanders, A. D., Falconer, E., Hills, M., Spierings, D. C. J. & Lansdorp, P. M. Single-cell template strand sequencing by Strand-seq enables the characterization of individual homologs. *Nat. Protoc.* **12**, 1151–1176 (2017).
- Sanders, A. D. et al. Single-cell analysis of structural variations and complex rearrangements with tri-channel processing. *Nat. Biotechnol.* **38**, 343–354 (2020).
- Schones, D. E. et al. Dynamic regulation of nucleosome positioning in the human genome. *Cell* **132**, 887–898 (2008).
- Lai, B. et al. Principles of nucleosome organization revealed by single-cell micrococcal nuclease sequencing. *Nature* **562**, 281–285 (2018).
- Struhl, K. & Segal, E. Determinants of nucleosome positioning. *Nat. Struct. Mol. Biol.* **20**, 267–273 (2013).
- Teif, V. B. et al. Genome-wide nucleosome positioning during embryonic stem cell development. *Nat. Struct. Mol. Biol.* **19**, 1185–1192 (2012).
- Lam, F. H., Steger, D. J. & O’Shea, E. K. Chromatin decouples promoter threshold from dynamic range. *Nature* **453**, 246–250 (2008).
- Shivaswamy, S. et al. Dynamic remodeling of individual nucleosomes across a eukaryotic genome in response to transcriptional perturbation. *PLoS Biol.* **6**, e65 (2008).
- Porubský, D. et al. Direct chromosome-length haplotyping by single-cell sequencing. *Genome Res.* **26**, 1565–1574 (2016).

32. Kundaje, A. et al. Ubiquitous heterogeneity and asymmetry of the chromatin environment at regulatory elements. *Genome Res.* **22**, 1735–1747 (2012).
33. Roadmap Epigenomics Consortium et al. Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330 (2015).
34. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
35. Loda, A., Collombet, S. & Heard, E. Gene regulation in time and space during X-chromosome inactivation. *Nat. Rev. Mol. Cell Biol.* **23**, 231–249 (2022).
36. Yates, A. D. et al. Ensembl 2020. *Nucleic Acids Res.* **48**, D682–D688 (2020).
37. Mardin, B. R. et al. A cell-based model system links chromothripsis with hyperploidy. *Mol. Syst. Biol.* **11**, 828 (2015).
38. Ebert, P. et al. Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* **372**, eabf7117 (2021).
39. 1000 Genomes Project Consortium et al. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
40. Shirley, M. D. et al. Chromosomal variation in lymphoblastoid cell lines. *Hum. Mutat.* **33**, 1075–1086 (2012).
41. Mraz, M. et al. The origin of deletion 22q11 in chronic lymphocytic leukemia is related to the rearrangement of immunoglobulin lambda light chain locus. *Leuk. Res.* **37**, 802–808 (2013).
42. Dang, S. et al. Dynamic expression of ZNF382 and its tumor-suppressor role in hepatitis B virus-related hepatocellular carcinogenesis. *Oncogene* **38**, 4804–4819 (2019).
43. Li, Z. et al. A global transcriptional regulatory role for c-Myc in Burkitt's lymphoma cells. *Proc. Natl Acad. Sci. USA* **100**, 8164–8169 (2003).
44. Marinov, G. K. et al. From single-cell to cell-pool transcriptomes: stochasticity in gene expression and RNA splicing. *Genome Res.* **24**, 496–510 (2014).
45. Li, H. et al. Reference component analysis of single-cell transcriptomes elucidates cellular heterogeneity in human colorectal tumors. *Nat. Genet.* **49**, 708–718 (2017).
46. Müller, S., Cho, A., Liu, S. J., Lim, D. A. & Diaz, A. CONICS integrates scRNA-seq with DNA sequencing to map gene expression to tumor sub-clones. *Bioinformatics* **34**, 3217–3219 (2018).
47. Fan, J. et al. Linking transcriptional and genetic tumor heterogeneity through allele analysis of single-cell RNA-seq data. *Genome Res.* **28**, 1217–1227 (2018).
48. Patel, A. P. et al. Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science* **344**, 1396–1401 (2014).
49. McClintock, B. The stability of broken ends of chromosomes in *Zea mays*. *Genetics* **26**, 234–282 (1941).
50. Kuleshov, M. V. et al. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* **44**, W90–W97 (2016).
51. Yang, X. et al. Discovery of the first chemical tools to regulate MKK3-mediated MYC activation in cancer. *Bioorg. Med. Chem.* **45**, 116324 (2021).
52. Byrská-Bishop, M. et al. High coverage whole genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios. *Cell* **185**, 3426–3440.e19 (2022).
53. Zhang, Y. & Rowley, J. D. Chromatin structural elements and chromosomal translocations in leukemia. *DNA Repair* **5**, 1282–1297 (2006).
54. Erickson, P. et al. Identification of breakpoints in t(8;21) acute myelogenous leukemia and isolation of a fusion transcript, AML1/ETO, with similarity to *Drosophila* segmentation gene, runt. *Blood* **80**, 1825–1831 (1992).
55. Xiao, Z. et al. Molecular characterization of genomic AML1-ETO fusions in childhood leukemia. *Leukemia* **15**, 1906–1913 (2001).
56. Mifsud, B. et al. Mapping long-range promoter contacts in human cells with high-resolution capture Hi-C. *Nat. Genet.* **47**, 598–606 (2015).
57. Gutierrez, A. et al. The BCL11B tumor suppressor is mutated across the major molecular subtypes of T-cell acute lymphoblastic leukemia. *Blood* **118**, 4169–4173 (2011).
58. Döhner, H. et al. Genomic aberrations and survival in chronic lymphocytic leukemia. *N. Engl. J. Med.* **343**, 1910–1916 (2000).
59. D'Ente, S. C. et al. Characterizing genetic intra-tumor heterogeneity across 2,658 human cancer genomes. *Cell* **184**, 2239–2254.e39 (2021).
60. Hewett, D. R. et al. FRA10B structure reveals common elements in repeat expansion and chromosomal fragile site genesis. *Mol. Cell* **1**, 773–781 (1998).
61. Edelmann, J. et al. High-resolution genomic profiling of chronic lymphocytic leukemia reveals new recurrent genomic alterations. *Blood* **120**, 4783–4794 (2012).
62. Puente, X. S. et al. Non-coding recurrent mutations in chronic lymphocytic leukaemia. *Nature* **526**, 519–524 (2015).
63. Malek, S. N. The biology and clinical significance of acquired genomic copy number aberrations and recurrent gene mutations in chronic lymphocytic leukemia. *Oncogene* **32**, 2805–2817 (2013).
64. Liberzon, A. et al. Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27**, 1739–1740 (2011).
65. Boon, E. M. J., van der Neut, R., van de Wetering, M., Clevers, H. & Pals, S. T. Wnt signaling regulates expression of the receptor tyrosine kinase met in colorectal cancer. *Cancer Res.* **62**, 5126–5128 (2002).
66. Stoekius, M. et al. Simultaneous epitope and transcriptome measurement in single cells. *Nat. Methods* **14**, 865–868 (2017).
67. Borchering, N. et al. Mapping the immune environment in clear cell renal carcinoma by single-cell genomics. *Commun. Biol.* **4**, 122 (2021).
68. Garcia-Alonso, L., Holland, C. H., Ibrahim, M. M., Turei, D. & Saez-Rodriguez, J. Benchmark and integration of resources for the estimation of human transcription factor activities. *Genome Res.* **29**, 1363–1375 (2019).
69. Kamburov, A. & Herwig, R. ConsensusPathDB 2022: molecular interactions update as a resource for network biology. *Nucleic Acids Res.* **50**, D587–D595 (2022).
70. Böttcher, M. et al. Control of PD-L1 expression in CLL-cells by stromal triggering of the Notch-c-Myc-EZH2 oncogenic signaling axis. *J. Immunother. Cancer* **9**, e001889 (2021).
71. Wang, Y. et al. Distinct immune signatures in chronic lymphocytic leukemia and Richter syndrome. *Blood Cancer J.* **11**, 86 (2021).
72. Fromigué, O., Haÿ, E., Barbara, A. & Marie, P. J. Essential role of nuclear factor of activated T cells (NFAT)-mediated Wnt signaling in osteoblast differentiation induced by strontium ranelate. *J. Biol. Chem.* **285**, 25251–25258 (2010).
73. Moon, J. B. et al. Akt induces osteoclast differentiation through regulating the GSK3 $\beta$ /NFATc1 signaling cascade. *J. Immunol.* **188**, 163–169 (2012).
74. Nurieva, R. I. et al. A costimulation-initiated signaling pathway regulates NFATc1 transcription in T lymphocytes. *J. Immunol.* **179**, 1096–1103 (2007).
75. Park, H.-J., Baek, K., Baek, J.-H. & Kim, H.-R. The cooperation of CREB and NFAT is required for PTHrP-induced RANKL expression in mouse osteoblastic cells. *J. Cell. Physiol.* **230**, 667–679 (2015).
76. Li, L. et al. B-cell receptor-mediated NFATc1 activation induces IL-10/STAT3/PD-L1 signaling in diffuse large B-cell lymphoma. *Blood* **132**, 1805–1817 (2018).

77. Oestreich, K. J., Yoon, H., Ahmed, R. & Boss, J. M. NFATc1 regulates PD-1 expression upon T cell activation. *J. Immunol.* **181**, 4832–4839 (2008).
78. Staal, F. J. T., Famili, F., Garcia Perez, L. & Pike-Overzet, K. Aberrant Wnt signaling in leukemia. *Cancers (Basel)* **8**, 78 (2016).
79. Korbel, J. O. & Campbell, P. J. Criteria for inference of chromothripsis in cancer genomes. *Cell* **152**, 1226–1236 (2013).
80. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.* **36**, 411–420 (2018).
81. Andreatta, M. & Carmona, S. J. UCell: robust and scalable single-cell gene signature scoring. *Comput. Struct. Biotechnol. J.* **19**, 3796–3798 (2021).
82. Street, K. et al. Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics. *BMC Genomics* **19**, 477 (2018).
83. Majumder, S. et al. Targeting notch in oncology: the path forward. *Nat. Rev. Drug Discov.* **20**, 125–144 (2021).
84. Study of CB-103 in adult patients with advanced or metastatic solid tumours and haematological malignancies. <https://clinicaltrials.gov/ct2/show/NCT03422679> (2017).
85. Lehal, R. et al. Pharmacological disruption of the Notch transcription factor complex. *Proc. Natl Acad. Sci. USA* **117**, 16292–16301 (2020).
86. Drews, R. M. et al. A pan-cancer compendium of chromosomal instability. *Nature* **606**, 976–983 (2022).
87. Edelmann, J. et al. Genomic alterations in high-risk chronic lymphocytic leukemia frequently affect cell cycle key regulators and NOTCH1-regulated transcription. *Haematologica* **105**, 1379–1390 (2020).
88. Loh, P.-R. et al. Insights into clonal haematopoiesis from 8,342 mosaic chromosomal alterations. *Nature* **559**, 350–355 (2018).
89. Zhu, C., Preissl, S. & Ren, B. Single-cell multimodal omics: the power of many. *Nat. Methods* **17**, 11–14 (2020).
90. Forsberg, L. A., Gisselsson, D. & Dumanski, J. P. Mosaicism in health and disease - clones picking up speed. *Nat. Rev. Genet.* **18**, 128–142 (2017).
91. Leibowitz, M. L. et al. Chromothripsis as an on-target consequence of CRISPR-Cas9 genome editing. *Nat. Genet.* **53**, 895–905 (2021).
92. Gaffney, D. J. et al. Controls of nucleosome positioning in the human genome. *PLoS Genet.* **8**, e1003036 (2012).
93. Fishilevich, S. et al. GeneHancer: genome-wide integration of enhancers and target genes in GeneCards. *Database (Oxford)* **2017**, bax028 (2017).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022

<sup>1</sup>Genome Biology Unit, European Molecular Biology Laboratory (EMBL), Heidelberg, Germany. <sup>2</sup>Faculty of Biosciences, EMBL and Heidelberg University, Heidelberg, Germany. <sup>3</sup>Division of Pediatric Oncology, University Children's Hospital, Zürich, Switzerland. <sup>4</sup>Department of Hematology, Oncology and Rheumatology, Heidelberg University Hospital, Heidelberg, Germany. <sup>5</sup>Molecular Medicine Partnership Unit, European Molecular Biology Laboratory, University of Heidelberg, Heidelberg, Germany. <sup>6</sup>Department of Hematology and Oncology, University Hospital Düsseldorf, Düsseldorf, Germany. <sup>7</sup>National Centre for Biological Sciences, Tata Institute of Fundamental Research, Bangalore, India. <sup>8</sup>Center for Bioinformatics, Saarland University, Saarbrücken, Germany. <sup>9</sup>Max Planck Institute for Informatics, Saarbrücken, Germany. <sup>10</sup>Department of Pediatric Oncology, Hematology, and Immunology, University of Heidelberg and Hopp Children's Cancer Center, Heidelberg, Germany. <sup>11</sup>Department of Hematology and Oncology, Medical Faculty Mannheim of the Heidelberg University, Heidelberg, Germany. <sup>12</sup>Institute for Medical Biometry and Bioinformatics, Medical Faculty, Heinrich Heine University, Düsseldorf, Germany. <sup>13</sup>Department of Translational Medical Oncology, National Center for Tumor Diseases (NCT) Heidelberg and German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>14</sup>Berlin Institute for Medical Systems Biology, Max Delbrück Center for Molecular Medicine in the Helmholtz Association (MDC), Berlin, Germany. <sup>15</sup>Berlin Institute of Health (BIH), Berlin, Germany. <sup>16</sup>Charité-Universitätsmedizin, Berlin, Germany. <sup>17</sup>Bridging Research Division on Mechanisms of Genomic Variation and Data Science, German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>18</sup>Present address: Hanyang Institute of Bioscience and Biotechnology, Hanyang University, Seoul, Republic of Korea. <sup>19</sup>Present address: Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA. <sup>20</sup>These authors contributed equally: Hyobin Jeong, Karen Grimes <sup>21</sup>These authors jointly supervised this work: Ashley D. Sanders, Jan O. Korbel. ✉ e-mail: [ashley.sanders@mdc-berlin.de](mailto:ashley.sanders@mdc-berlin.de); [jan.korbel@embl.org](mailto:jan.korbel@embl.org)

## Methods

### Strand-seq library preparation

NA20509 Strand-seq libraries were prepared as previously described<sup>94</sup>. Strand-seq libraries of primary leukemia samples were generated as follows: peripheral blood mononuclear cells of a previously untreated female CLL patient (routine diagnostics: *IGHV* unmutated, no *TP53* mutation, no detected alteration in 6q21, 8q24, 11q22.3, 12q13, 13q14 and 17p13) were isolated after obtaining informed consent. Cells were isolated and cultured using previously established protocols<sup>95</sup>. CLL cells were cultured at  $1 \times 10^6$  cells  $\text{ml}^{-1}$  in Roswell Park Memorial Institute (RPMI) medium (Gibco by Life technologies), supplemented with 10% human serum (PAN BIOTECH), 1% Pen/Strep (GIBCO by Life Technologies) and 1% Glutamine (GIBCO by Life Technologies). Cells were stimulated with  $1 \mu\text{g ml}^{-1}$  Resiquimod (Enzo) and  $50 \text{ ng ml}^{-1}$  IL-2 (Sigma). BrdU ( $40 \mu\text{M}$ ; Sigma) was incorporated for 90 h and 120 h, respectively, to perform nontemplate strand labeling. Single nuclei from each time-point were sorted into 96-well plates using a BD FACSMelody cell sorter, followed by Strand-seq library preparation (described below). In the case of the AML sample, frozen primary mononuclear cells from a bone marrow aspirate were thawed and stained with CD34-APC (clone 581; Biolegend), CD38-PeCy7 (clone HB7; eBioscience), CD45Ra-FITC (clone HII100; eBioscience), CD90-PE (clone 5E10; eBioscience) and LIVE/DEAD Fixable Near-IR Dead Cell Stain (ThermoFisher). Single, viable, CD34<sup>+</sup> cells (Supplementary Fig. 15) were sorted using a BD FACS Aria Fusion Cell Sorter into ice-cold serum-free expansion medium (SFEM) supplemented with  $100 \text{ ng ml}^{-1}$  SCF and Flt3 (Stem Cell Technologies),  $20 \text{ ng ml}^{-1}$  IL-3, IL-6, G-CSF and TPO (Stem Cell Technologies). Cells were plated in Corning Costar Ultra-Low Attachment 96-well flat-bottom plates (Sigma) at  $1 \times 10^5$  cells  $\text{ml}^{-1}$  in warm medium as above. At 24 h after culture,  $40 \mu\text{M}$  BrdU was added. Nuclei were isolated after 43 h total culture time, and BrdU-incorporating nuclei sorted into 96-well plates followed by Strand-seq library preparation. All Strand-seq libraries were automatically prepared using a Biomek FXP liquid handling robotic system, as described previously<sup>23,96</sup>. Libraries were sequenced on an Illumina NextSeq 500 sequencing platform (MID-mode, 75 base pair (bp) paired-end sequencing protocol).

### Strand-seq data preprocessing

Reads from Strand-seq (fastq) libraries were aligned to the hg38 assembly using BWA<sup>97</sup>, as previously described<sup>24</sup>. Sequence reads with low quality (MAPQ < 10), supplementary reads and duplicated reads were removed. Single-cell library selection was performed as described previously<sup>24</sup>. The single-cell footprints of different SV classes were discovered using the principle of scTRIP of Strand-seq data using the MosaicCatcher computational pipeline with default settings<sup>24</sup>.

### Coupling NO measurements and SV discovery in the same cell with scNOVA

We developed scNOVA as a computational framework for coupling discovered somatic SVs with analyses of NO profiles in the same cell. The scNOVA workflow covers a set of different operations from single-cell SV discovery (using the previously described scTRIP method<sup>24</sup>) to NO profiling at CREs, and gene as well as pathway dysregulation inference based on NO at gene bodies, and can be used in a haplotype-aware or -unaware manner (Extended Data Fig. 1). To maximize reusability, interoperability and reproducibility we combined all scNOVA modules into a coherent workflow using snakemake. Alternatively, these modules can be executed individually.

**Data analysis and operational definition utilized for NO.** We operationally defined NO closely following definitions from a previous study<sup>28</sup>. NO maps were calculated by counting how many reads from the Strand-seq libraries (which typically comprise mono-nucleosomal fragments around 140–180 bp in size; see Supplementary Table 1 and Supplementary Fig. 1) covered a given bp based on aligning reads to the

GRCh38 (hg38) genome assembly with BWA<sup>97</sup>. Genomic regions with unusual (such as artificially high) coverage were considered artifacts, and were automatically excluded ('blacklisted') by our Strand-seq analysis workflow as previously described<sup>24</sup>. No further peak calling or smoothing was conducted, and no assumptions on the length of the nucleosomal DNA were made to derive NO maps, as nucleosome boundaries were determined on both sides of the nucleosome by paired-end sequencing<sup>28</sup>. For the calculation of NO around bound CTCF binding sites (downloaded from ENCODE<sup>34</sup>), the averaged profile was scaled<sup>28</sup> to yield an NO equal to 1 at position  $-2,000$  bp from the center of the bound CTCF site.

**Cell type classification.** We generated feature sets from the NO at the body of genes (defined as the region from the TSS to the transcription termination site, which includes exons and introns) at the single-cell level. When several sequencing batches from the same samples were available, we applied batch correction to the NO count matrix using ComBat-seq<sup>98</sup>. NO in gene body regions was normalized by segmental copy number status, and by library size to obtain reads per million, which we transformed into  $\log_2$  scale. This feature set was used for the unsupervised dimension reduction plot (Extended Data Fig. 3) and for training of a supervised classification model based on PLS-DA<sup>99</sup>.

**Haplotype-phasing of single-cell NO tracks.** As previously described, Strand-seq directly resolves its underlying sequence reads onto haplotypes ranging from telomere to telomere<sup>31</sup> (chromosome-length haplotyping). scNOVA phases NO profiles onto a chromosomal homolog using the StrandPhaseR algorithm<sup>31</sup>, which is employed wherever the template strand segregation pattern of a chromosome enables unambiguous haplotype-phasing, that is, for Watson/Crick (WC) or Crick/Watson (CW) template state configurations in Strand-seq libraries<sup>31,96</sup>. Haplotype-specific analyses pursued by scNOVA employ phased reads (normalized by locus copy number), whereas the inference of gene activity changes uses both phased reads (from chromosomes with a WC or CW configuration) and unphased reads (from chromosomes with a CC or WW configuration<sup>31,96</sup>).

**Inference of haplotype-specific NO and identification of local effects of SVs.** To dissect local effects of SVs, the scNOVA framework performs a genome-wide haplotype-specific NO analysis at gene bodies in pseudobulk, which yields a haplotype-specific NO matrix. Using this matrix, scNOVA then scans up to  $\pm 1$  Mb around each somatic SV breakpoint to infer local effects of these breakpoints on haplotype-specific gene activity, using FDR-adjusted Wilcoxon rank sum tests. Once a local effect on gene activity is identified, scNOVA additionally provides the option to locally scan for CREs exhibiting haplotype-specific NO. To do so, user-provided CRE positions from the cell type of interest are used by scNOVA to calculate haplotype-specific NO at CREs, and the Exact test (10% FDR) is used for significance testing.

**Inference of genome-wide changes in gene activity.** This haplotype-unaware module of scNOVA considers all reads—whether phased or not—to infer gene activity alterations via analysis of differential patterns of NO along gene bodies. scNOVA obtains gene loci from ENSEMBL (GRCh38.81), converted into bed format (Genebody\_hg38.81.bed). Strand-seq reads falling within the start and end position of genes (Genebody\_hg38.81.bed) were identified with the Deeptool multiBamSummary function<sup>100</sup>, using the following parameters: [multiBamSummary BED-file -BED Genebody\_hg38.81.bed -bamfiles Input.bam -extendReads -outRawCounts output.tab -out output.npz]. The scNOVA gene dysregulation inference module contains two steps: Step 1 filters out genes unlikely to be expressed ('not expressed', NEs), whereas Step 2 infers dysregulated (that is, differentially expressed) genes between subclones using a generalized linear model.

using Strand-seq read count data, with these read counts then being normalized using the median-of-ratios method from DESeq2 (ref. <sup>103</sup>). For each member in the biological pathway gene sets from MSigDB<sup>64</sup>, scNOVA then computes mean normalized NO values, in each single-cell, as a proxy for pathway-level NO. Lowly variable genes (s.d. <80%) are removed. Pathway-level NO is compared between cells with and without SVs using linear mixed model fitting followed by likelihood ratio testing, and controlling the FDR at 10%. For linear mixed model fitting, SV status is defined as a fixed effect and different Strand-seq library batches are defined as random effects, by scNOVA.

### Quantitative real-time PCR

NA20509 was ordered from Coriell and taken into culture at passage four. The late passage was grown until passage eight in a time span of 8 weeks. HG01505 was taken into culture at passage five and was grown until passage nine within a total time span of 6 weeks. DNA, RNA and protein were isolated with the NucleoSpin TriPrep Mini kit (740966.50) according to the manufacturer's protocol. qPCR was performed on genomic DNA. PCR primers for *MAP2K3* and *TP53* were obtained from Sigma. qPCR was performed using BD SYBR Green PCR Master Mix (4309155) with a final primer concentration of 300 nM each and 10 ng input gDNA. A *GAPDH* control region was used as a normalizer. The primer sequences for DNA qPCR are provided in Supplementary Table 17.

### Drug treatment with CB-103

Primary human T-ALL cells were recovered from cryopreserved bone marrow aspirates of patients enrolled in the ALL-BFM 2009 study. Patient-derived xenografts were generated as previously described by intrafemoral injection of 1 million viable primary ALL cells in NSG mice<sup>104</sup>. Patient-derived xenografts (T-ALL\_P1)<sup>24</sup> cells were frozen until processing. Human hTERT immortalized primary bone marrow mesenchymal stroma cells (MSC; provided by D. Campana) were cultured in Roswell Park Memorial Institute (RPMI) 1640 medium supplemented with 10% heat-inactivated fetal bovine serum, L-glutamine (2 mM), penicillin/streptomycin (100 IU ml<sup>-1</sup>) and hydrocortisone (1 μM). MSCs were seeded in 24-well plates at a concentration of 500,000 cells per well in 1 ml Aim V medium. After 24 h, T-ALL cells were added at a concentration of 1.5 million cells per well in 1 ml Aim V. CB-103 (MedChemExpress, HY-135145) or DMSO (vehicle) as control was added after an additional 24 h at a concentration of 10 μM. After 8 h and 24 h, cells were trypsinized, collected and frozen in 90% fetal bovine serum/10% DMSO.

### Single-cell RNA sequencing and data processing

For scRNA-seq library preparation, cryopreserved cells were thawed rapidly at 37 °C and resuspended in 10 ml warm RPMI medium with 100 μg ml<sup>-1</sup> Dnase I. Cells were centrifuged for 5 mins at 300g, and resuspended in ice-cold PBS with 2% fetal bovine serum and 5 mM EDTA. Cells were stained on ice with anti-murine-CD45-PE (mCD45) (clone 30-F11; BioLegend; 1:20) in the dark for 30 mins. 1:100 4,6-diamidino-2-phenylindole (DAPI) was added and incubated in the dark for 5 mins before sorting. Triple negative cells (4,6-diamidino-2-phenylindole-mCD45-GFP<sup>-</sup>) were sorted (Supplementary Fig. 32) using a BD FACSAria fusion cell sorter into ice-cold 0.03% bovine serum albumin (BSA) in PBS. All isolated cells were used immediately for scRNA-seq libraries, which were generated as per the standard 10x Genomics Chromium 3' (v.3.1 Chemistry) protocol. Completed libraries were sequenced on a NextSeq5000 sequencer (HIGH-mode, 75 bp paired-end).

Sequenced transcripts were aligned to both human and mouse genomes (GRCh38 and mm10) and quantified into count matrices using cell ranger mkfastq and count workflows (10x Genomics, v.3.1.0, default parameters). The R package Seurat<sup>80</sup> (v.4.0.3) was used for quality control of single cells and unsupervised clustering of the data. Briefly, human cells were separated from multiplets/mouse contamination based on greater than 97 % of their reads aligning to GRCh38.

Further filtering for high quality cells accepted only those with more than 200 but less than 20,000 total RNA counts, and a percentage of mitochondrial reads less than 10% for the untreated data, and less than 40% for the drug-treated samples. Finally, remaining mouse transcripts were removed before further analysis.

In the untreated data, normalization, scaling and regression of mitochondrial read percentage was carried out using the scTransform package<sup>105</sup>. Dimensionality reduction and differential expression analysis of identified clusters was performed as standard using Seurat. Trajectory analysis was performed using Monocle3 (ref. <sup>106</sup>). In the drug treatment data, individual Seurat objects that had been quality controlled as above were normalized by scTransform<sup>105,107</sup> and then integrated to correct for batch effects and allow for comparative analysis. To re-annotate clusters from the untreated data in the drug treatment data, the TransferData() function from Seurat<sup>80</sup> was used to project labels from our reference (that is, untreated data) onto the integrated drug treatment data. Single-cell gene set enrichment analysis was performed using the R package 'escape'<sup>67</sup>.

### Cellular indexing of transcriptomes and epitopes by single-cell sequencing

A peripheral blood-derived sample (CLL\_24) was recovered from cryopreservation as previously described<sup>108</sup> to reach viability above 90%. Then, 5 × 10<sup>5</sup> viable cells were stained by a premixed cocktail of oligonucleotide-conjugated antibodies (Supplementary Table 14) and incubated at 4 °C for 30 min. We provided dilution used for each antibody in Supplementary Table 14. Cells were washed three times with ice-cold washing buffer. After completion, bead-cell suspensions, synthesis of complementary DNA and single-cell gene expression and antibody-derived tag (ADT) libraries were performed using a Chromium single cell v.3.1.3' kit (10x Genomics) according to the manufacturer's instructions. Then, 3' gene expression and ADT libraries were pooled in a ratio of 3:1 aiming for 40,000 reads (gene expression) and 15,000 reads per cell (ADT), respectively. Sequencing was performed on a NextSeq 500 (Illumina). After sequencing, the cell ranger wrapper function (10x Genomics, v.6.1.1) cellranger mkfastq was used to demultiplex and to align raw base-call files to the human reference genome (hg38). The obtained FASTQ files were counted by the cellranger count command. If not otherwise indicated default settings were used. Single-cell gene set enrichment analysis was performed using the R package 'escape'<sup>67</sup>.

### Single-cell gene signature scoring using UCell

The activity of the scNOVA-identified gene set from T-ALL\_P1 in scRNA-seq data was profiled using the UCell package<sup>81</sup>. Briefly, signature genes considered were those with either increased (implying decreased expression) or decreased (implying increased expression) nucleosome occupancy (see Fig. 5b), or genes encoding TFs whose targets showed differential nucleosome occupancy (see Fig. 5c). The following gene set was used for T-ALL\_P1: 'PRKCB-', 'RPS6KA2-', 'FAM120B-', 'FAM86C1+', 'FBXO22+', 'RHOH+', 'SLC9A7+', 'NASP+', 'NOTCH1+', 'MRPL48+', 'MFSD9+', 'MVB12B+', 'MYB+' (with '+' for upregulated, and '-' for downregulated). The score per single cell for the entire directional gene set was calculated using the AddModuleScore\_UCell() function. Cells were considered to be 'active' for the signature genes if their respective UCell score was greater than or equal to the median UCell score of the entire dataset, plus the s.d. Similarly, for T-cell cell-type labeling, marker gene sets for T-cell subsets were obtained from Park et al.<sup>109</sup> and single cells were scored for their activity in each gene set. Cells were labeled by their best-fit cell type, that is the cell-type whose gene set gave the highest UCell score.

### Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

Sequencing data from this study can be retrieved from the European Genome-phenome Archive (EGA) and the European Nucleotide Archive (ENA). LCL data are available under the following accessions: Strand-seq ([PRJEB39750](https://www.ebi.ac.uk/ena/browser/view/PRJEB39750), [PRJEB55038](https://www.ebi.ac.uk/ena/browser/view/PRJEB55038)); RNA-seq ([ERP123231](https://www.ebi.ac.uk/ena/browser/view/ERP123231)); WGS ([PRJEB37677](https://www.ebi.ac.uk/ena/browser/view/PRJEB37677)). C11 cell line data are available under the accession [PRJEB55012](https://www.ebi.ac.uk/ena/browser/view/PRJEB55012). Leukemia patient data and human primary cells derived data were deposited in the European Genome-phenome Archive (EGA) under the following accession numbers: skin fibroblast ([EGAS00001006498](https://www.ebi.ac.uk/ena/browser/view/EGAS00001006498)); cord blood ([EGAS00001006567](https://www.ebi.ac.uk/ena/browser/view/EGAS00001006567)). T-ALL Strand-seq and scRNA-seq ([EGAS00001003365](https://www.ebi.ac.uk/ena/browser/view/EGAS00001003365)), CLL Strand-seq ([EGAS00001004925](https://www.ebi.ac.uk/ena/browser/view/EGAS00001004925)), AML Strand-seq ([EGAS00001004903](https://www.ebi.ac.uk/ena/browser/view/EGAS00001004903)), T-ALL bulk RNA-seq ([EGAS00001003248](https://www.ebi.ac.uk/ena/browser/view/EGAS00001003248)), CLL bulk RNA-seq ([EGAS00001005746](https://www.ebi.ac.uk/ena/browser/view/EGAS00001005746)), CLL CITE-seq ([EGAS00001004925](https://www.ebi.ac.uk/ena/browser/view/EGAS00001004925)). Access to human patient data is governed by the EGA Data Access Committee.

## Code availability

The computational code of our analytical framework scNOVA is available open source at <https://github.com/jeongdo801/scNOVA>, with no restrictions on reuse. Other software used: Mosaicatcher (<https://github.com/friendsofstrandseq/mosaicatcher-pipeline>), StrandPhaseR (<https://github.com/daewoooo/StrandPhaseR>), InferCNV (<https://github.com/broadinstitute/inferCNV/>), HoneyBADGER (<https://jef.works/HoneyBADGER/>), CONICSmat (<https://github.com/diazlab/CONICS>), NucTools (<https://homeveg.github.io/nuc-tools>), Delly2 (<https://github.com/dellytools/delly>), BWA (v.0.7.15), STAR (v.2.7.9a), SAMtools (v.1.3.1), biobambam2 (v.2.0.76), deepTools (v.2.5.1), perl (v.5.16.3), Python (v.3.7.4), cuDNN (v.7.6.4.38), CUDA (v.10.1.243), TensorFlow (v.1.15.0), scikit-learn (v.0.21.3), matplotlib (v.3.1.1), R v.4.0.0, DESeq2, FlowJo and BD FACSDiva.

## References

94. Porubsky, D. et al. Recurrent inversion polymorphisms in humans associate with genetic instability and genomic disorders. *Cell* **185**, 1986–2005.e26 (2022).
95. Dietrich, S. et al. Drug-perturbation-based stratification of blood cancer. *J. Clin. Invest.* **128**, 427–445 (2018).
96. Falconer, E. et al. DNA template strand sequencing of single-cells maps genomic rearrangements at high resolution. *Nat. Methods* **9**, 1107–1112 (2012).
97. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
98. Zhang, Y., Parmigiani, G. & Johnson, W. E. ComBat-seq: batch effect adjustment for RNA-seq count data. *NAR Genom. Bioinform.* **2**, lqaa078 (2020).
99. Boulesteix, A.-L. & Strimmer, K. Partial least squares: a versatile tool for the analysis of high-dimensional genomic data. *Brief. Bioinform.* **8**, 32–44 (2007).
100. Ramirez, F. et al. deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.* **44**, W160–W165 (2016).
101. Ulz, P. et al. Inferring expressed genes by whole-genome sequencing of plasma DNA. *Nat. Genet.* **48**, 1273–1278 (2016).
102. Eraslan, G., Avsec, Ž., Gagneur, J. & Theis, F. J. Deep learning: new computational modelling techniques for genomics. *Nat. Rev. Genet.* **20**, 389–403 (2019).
103. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
104. Schmitz, M. et al. Xenografts of highly resistant leukemia recapitulate the clonal composition of the leukemogenic compartment. *Blood* **118**, 1854–1864 (2011).
105. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol.* **20**, 296 (2019).

106. Cao, J. et al. The single-cell transcriptional landscape of mammalian organogenesis. *Nature* **566**, 496–502 (2019).
107. Stuart, T. et al. Comprehensive integration of single-cell data. *Cell* **177**, e21 (2019).
108. Roeder, T. et al. Dissecting intratumour heterogeneity of nodal B-cell lymphomas at the transcriptional, genetic and drug-response levels. *Nat. Cell Biol.* **22**, 896–906 (2020).
109. Park, J.-E. et al. A cell atlas of human thymic development defines T cell repertoire formation. *Science* **367**, eaay3224 (2020).
110. Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat. Methods* **14**, 975–978 (2017).
111. Weirauch, M. T. et al. Determination and inference of eukaryotic transcription factor sequence specificity. *Cell* **158**, 1431–1443 (2014).
112. Nagel, S. et al. Activation of TLX3 and NKX2-5 in t(5;14)(q35;q32) T-cell acute lymphoblastic leukemia by remote 3'-BCL11B enhancers and coregulation by PU.1 and HMGA1. *Cancer Res.* **67**, 1461–1471 (2007).
113. Xaus, J. et al. The expression of MHC class II genes in macrophages is cell cycle dependent. *J. Immunol.* **165**, 6364–6371 (2000).

## Acknowledgements

We thank A. Krebs, J. Zaugg, K. Rippe and I. Cortés-Ciriano for providing thoughtful feedback on the development of scNOVA. We also thank M. Paulsen (Flow Cytometry Core Facility) for assistance in cell sorting, B. Raeder for assisting in Strand-seq library preparation, and the EMBL Genomics Core Facility for assisting in single-cell automation (J. Zimmermann and V. Benes) and scRNA-seq library preparation (L. Villacorta). Finally, we thank W. Höps for assistance with single-cell analysis, as well as M. Happich and P. Richter-Pechanska for assistance with RNA-seq analysis. Principal funding came from the European Research Council (ERC Consolidator grant no. 773026, to J.O.K.). Funding also came from the an ERC Starting Grant (grant number 336045) to J.O.K., the National Institutes of Health (grant no. 2U24HG007497-05) to J.O.K. and T.M., the Baden-Württemberg Stiftung (for supporting the projects 'Epigenetics in T-ALL' and 'SV\_Surveillance') to J.O.K. and A.E.K., a Volkswagen Foundation grant (VW - 95826) to J.O.K. and the German Federal Ministry of Education and Research (grant no. 031A537B; de.NBI project) to J.O.K. H.J. and A.D.S. acknowledge fellowships through the Alexander von Humboldt Foundation. We thank the Human Genome Structural Variation Consortium for providing early access to deep bulk RNA-seq data from several LCLs (generated using funds provided by NHGRI Grant 2U24HG007497-05). D.N. is an endowed Professor of the German José-Carreras-Foundation (DJCLSH03/01). J.C.J. was funded by a Gerok position of the 'Deutsche Forschungsgemeinschaft' (DFG) (NO 817/5-2, FOR2033, NICHEM). K.K.R. received postdoctoral funding from the Deutsche Krebshilfe (Mildred-Scheel-Fellowship).

## Author contributions

H.J., K.G., A.D.S. and J.O.K. designed the study (including conceptualization of haplotype-specific NO analysis, cell-type classification and altered gene activity using Strand-seq data). H.J., K.G., A.D.S. and J.O.K. developed the scNOVA computational method. H.J., K.G., A.D.S. and J.O.K. performed single-cell SV discovery. A.D.S. and P.H. performed LCL Strand-seq experiments; K.G., P.-M.B. and S.D. performed CLL Strand-seq experiments; K.G., J.-C.J. and D.N. performed AML Strand-seq experiments and K.G.,



K.K.R. and P.H. performed T-ALL Strand-seq experiments. T.R. carried out WGS-based SV discovery and verification. H.J., D.P. and T.M. performed haplotype-phasing. LCL scRNA-seq analysis was carried out by H.J.; CLL scRNA-seq analysis by K.G., H.J. and T.R. and T-ALL scRNA-seq analysis by K.G., H.J. and K.K.R. K.K.R., K.G. and B.B. performed drug treatment experiments. K.G., P.H. and E.B. analyzed LCL clonal expansion. LCL RNA-seq analysis was carried out by H.J.; CLL RNA-seq analysis by H.J., S.H., P.-M.B. and S.D. and T-ALL RNA-seq analysis by H.J., B.E.-U. and A.E.K. R.S. and J.O.K. performed PCAWG SV driver spectrum analysis. The manuscript was written by H.J., K.G., A.D.S. and J.O.K., with additional contributions from all authors.

### Competing interests

The following authors have previously disclosed a patent application (no. EP19169090) that is relevant to this manuscript: A.D.S., J.O.K., T.M. and D.P. The remaining authors declare no competing interests.

### Additional information

**Extended data** is available for this paper at <https://doi.org/10.1038/s41587-022-01551-4>.

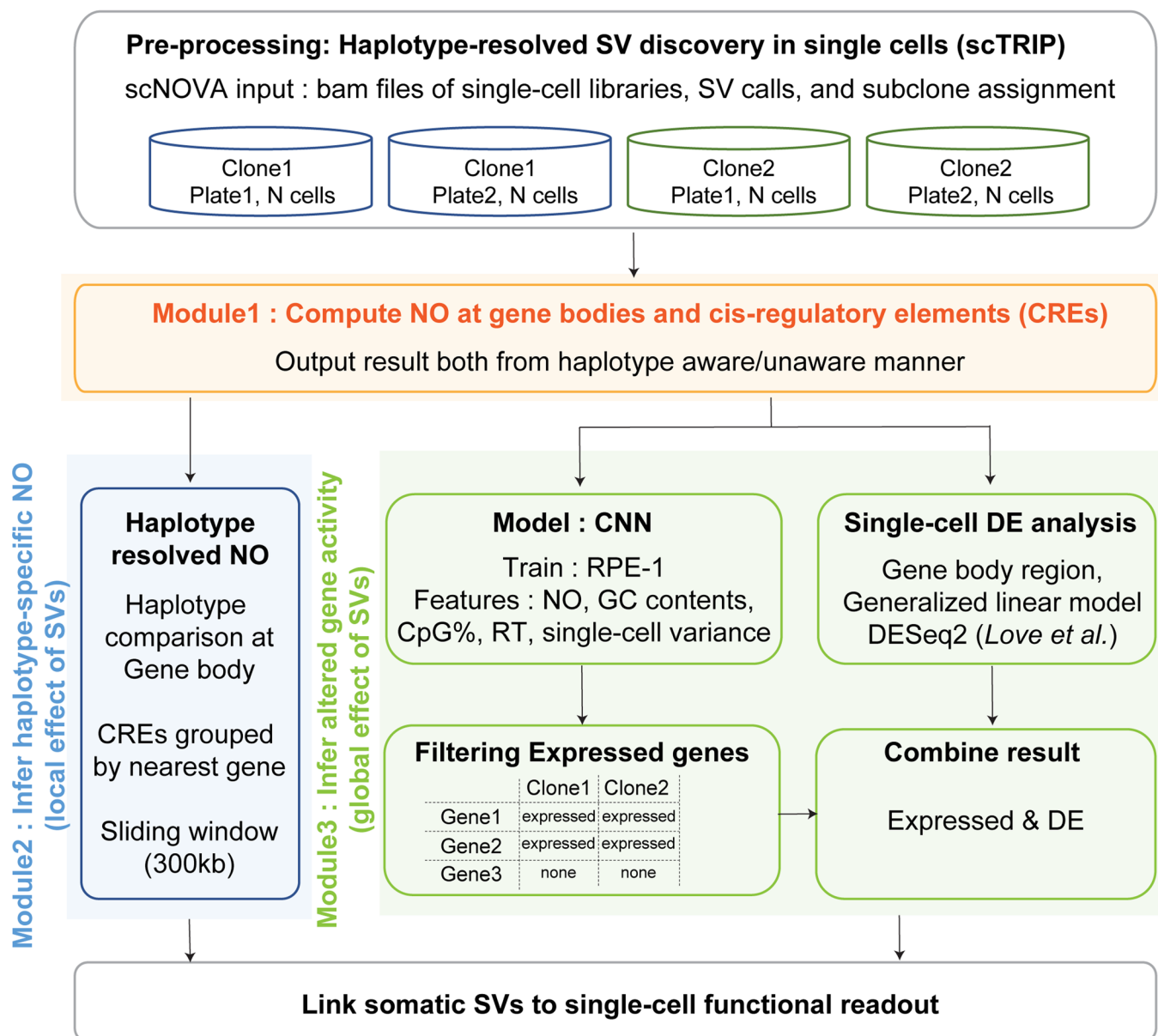
**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41587-022-01551-4>.

**Correspondence and requests for materials** should be addressed to Ashley D. Sanders or Jan O. Korb.

**Peer review information** *Nature Biotechnology* thanks Jonas Demeulemeester, Elisa Oricchio, Peter Park and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

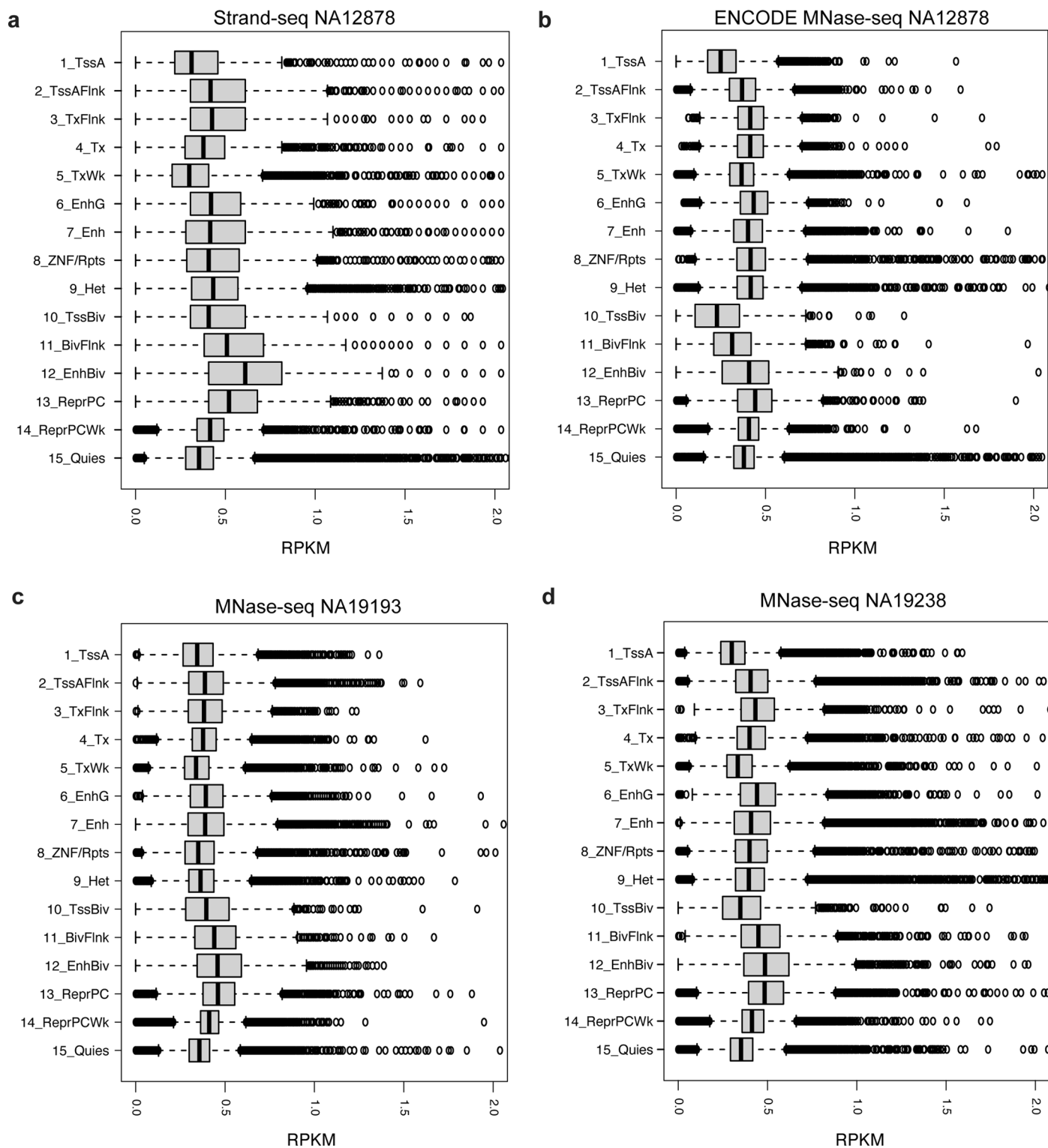
**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

## Computational pipeline of scNOVA



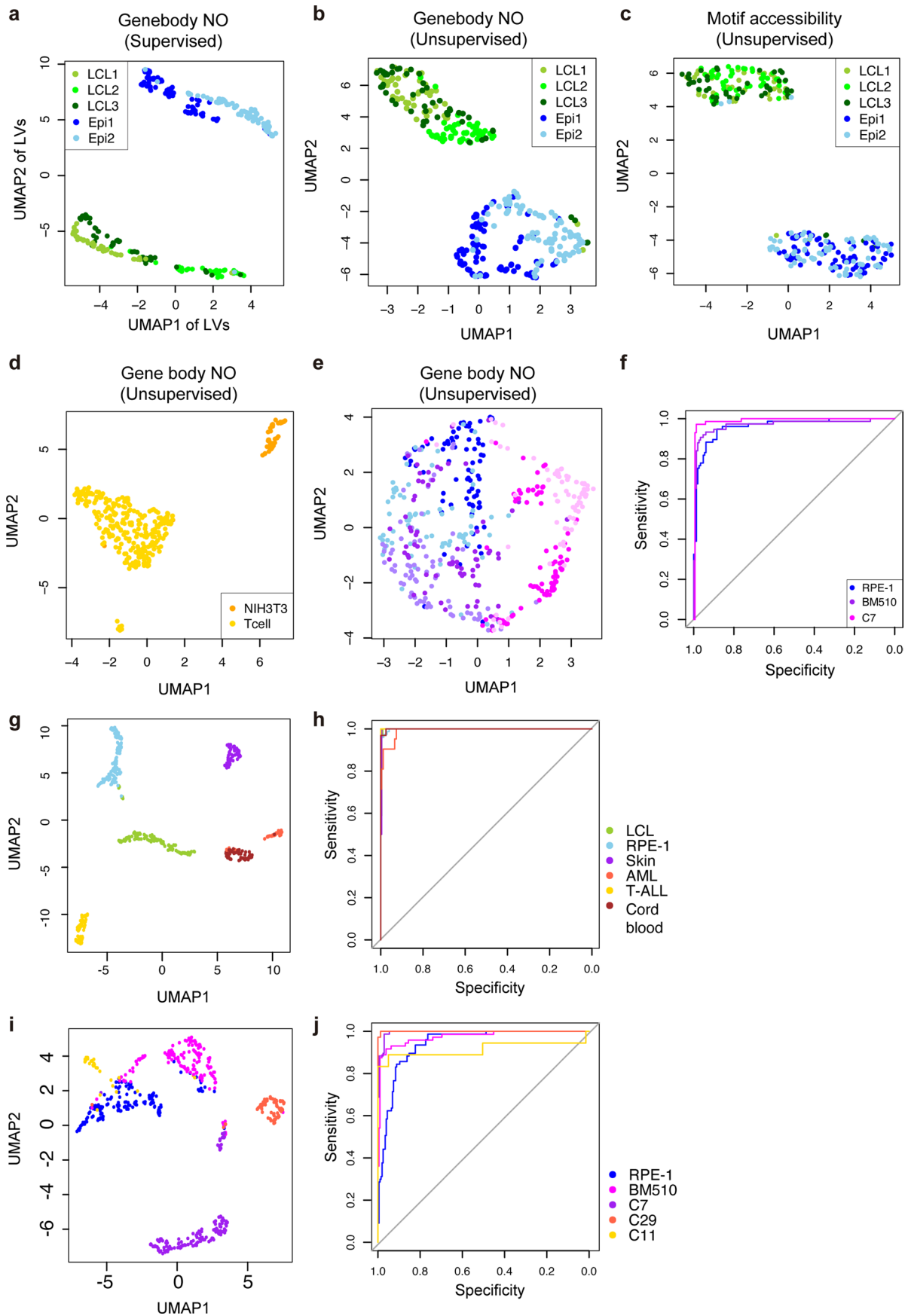
**Extended Data Fig. 1 | Overview of components of the scNOVA computational workflow.** scNOVA employs single cell tri-channel processing (scTRIP) as realized in the MosaiCatcher pipeline to perform haplotype-aware somatic SV discovery<sup>24</sup>. Modules of scNOVA enable single-cell multomics of these somatic SVs, including inference of haplotype-specific NO to investigate local (*cis*) effect of SVs, and inference of altered gene/pathway activity to investigate global (*trans*) effect of SVs detectable between genetically distinct subclones. To infer alterations in gene activity, scNOVA integrates deep convolutional neural

network (CNN) based machine learning, and negative binomial generalized linear models. The framework dissects intra-sample genetic heterogeneity at single-cell resolution, measures the local haplotype-specific impact of somatic SVs, can be used to explore global gene dysregulation in SV-containing cells, can discriminate between genetically-distinct subclones, and can uncover shared functional consequences of heterogeneous SVs affecting the same chromosomal interval.



**Extended Data Fig. 2 | Read depth of Strand-seq and MNase-seq data stratified into 15 chromatin states defined by Roadmap epigenome consortium<sup>33</sup>.** 15 chromatin states based on the NA12878 cell line were utilized in this genome-wide analysis. Plots generated represent Strand-seq data from NA12878 ( $n = 95$  cells) (a), and publicly available MNase-seq from NA12878, NA19193, and NA19238 ( $n = 1$  sample each) (b-d). The bulk MNase-seq experiment of NA12878 was pursued using single-end SOLID sequencing reads, and that of NA19193 and NA19238 was done using paired-end Illumina reads. The X-axis in the box plot indicates reads per kilobase per million (RPKM) measured for each genomic segment annotated by one of the 15 chromatin states. Abbreviations for chromatin states<sup>33</sup> are: TssA - Active TSS, TssAFlnk - Flanking Active TSS, TxFlnk - Transcription at gene 5' and 3', Tx - Strong transcription, TxWk - Weak transcription, EnhG - Genic enhancers, Enh - Enhancers, ZNF/Rpts - ZNF genes & repeats, Het - Heterochromatin, TssBiv - Bivalent/Poised TSS, BivFlnk - Flanking Bivalent TSS/Enh, EnhBiv - Bivalent Enhancer, ReprPC - Repressed PolyComb, ReprPCWk - Weak Repressed PolyComb, Quies - Quiescent/Low. Boxplots were defined by minima = 25th percentile - 1.5X interquartile range (IQR), maxima = 75th percentile + 1.5X IQR, center = median, and bounds of box = 25th and 75th percentile. Both Strand-seq and MNase-seq assays measured NO in all fifteen chromatin states. Among these chromatin states, Strand-seq and MNase-seq revealed the highest NO signals on average for the polycomb repressed state and the bivalent enhancer state; whereas the lowest average NO signals were consistently seen for the active transcription start site (TSS) state.

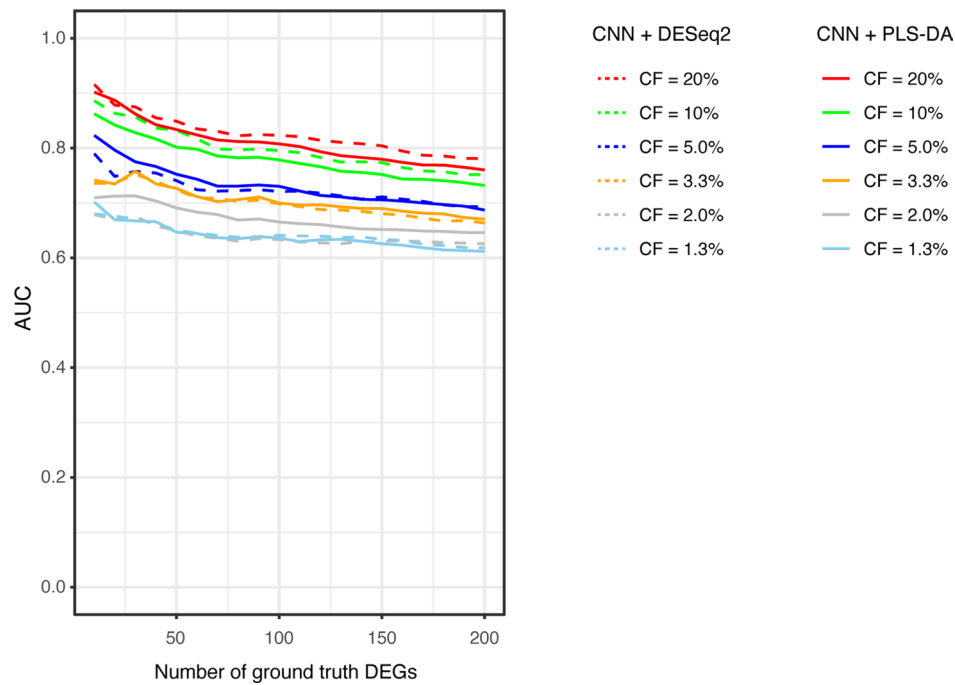
transcription, EnhG - Genic enhancers, Enh - Enhancers, ZNF/Rpts - ZNF genes & repeats, Het - Heterochromatin, TssBiv - Bivalent/Poised TSS, BivFlnk - Flanking Bivalent TSS/Enh, EnhBiv - Bivalent Enhancer, ReprPC - Repressed PolyComb, ReprPCWk - Weak Repressed PolyComb, Quies - Quiescent/Low. Boxplots were defined by minima = 25th percentile - 1.5X interquartile range (IQR), maxima = 75th percentile + 1.5X IQR, center = median, and bounds of box = 25th and 75th percentile. Both Strand-seq and MNase-seq assays measured NO in all fifteen chromatin states. Among these chromatin states, Strand-seq and MNase-seq revealed the highest NO signals on average for the polycomb repressed state and the bivalent enhancer state; whereas the lowest average NO signals were consistently seen for the active transcription start site (TSS) state.



Extended Data Fig. 3 | See next page for caption.

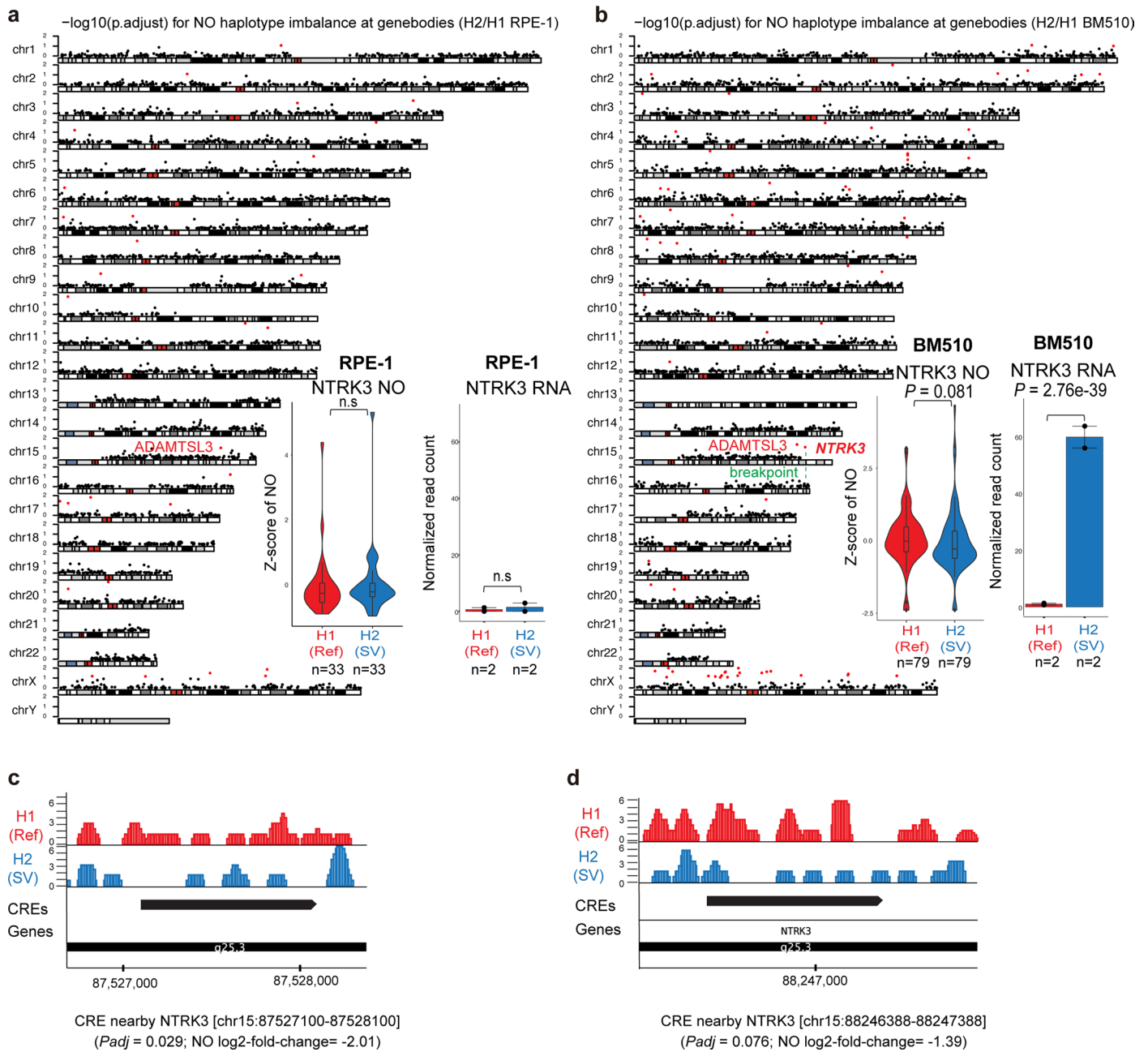
**Extended Data Fig. 3 | Utility of NO for cell-typing.** (a) Cell-typing based on NO at gene bodies (AUC = 1). Epi1: RPE-1 replicate 1 (79 cells); Epi2: replicate 2 (77 cells); LCL1: HGO1573 (46 cells); LCL2: HG02018 (50 cells), LCL3: NA19036 (50 cells); LV: latent variable. (b) UMAP visualization of Strand-seq libraries based on NO at gene-bodies (normalized by segmental ploidy status<sup>24</sup>). (c) We also explored dimensionality reduction of Strand-seq libraries based on DNA motif accessibility. Using the chromVAR package<sup>110</sup>, single-cell NO profiles for 2 kb DNase I hypersensitive sites (DHSs) were transformed into a deviation Z-score, which measures how likely a certain motif accessibility would occur when randomly sampling sets of peaks with similar GC content and read depth. For each single-cell, the deviation Z-score was calculated for 870 human TF motifs from the cisBP database<sup>111</sup>. These dimensionality reduction plots suggest that batch effect within the same cell type (three individuals in LCL, and two batches in

RPE-1 sequenced separately) is minimal, and far less than the cell-type dependent variability. (d) UMAP using scMNase-seq<sup>26</sup>, including 45 NIH3T3 cells and 272 murine naive T cells, based on NO at the gene-bodies. (e) UMAP of RPE-1 (the originally commercially available cell line) and its transformed derived<sup>37</sup> cell lines (BM510 and C7). Two biological replicates were sequenced for each cell line. (f) Receiver operating characteristic (ROC) using the PLS-DA based classifier. AUC for classifying each cell line was 0.9614, 0.9694, and 0.9892 for RPE-1, BM510, and C7 respectively. (g-h) Cell-typing for LCL, RPE-1, skin fibroblast, AML, T-ALL, and umbilical cord blood cells (g), and ROC curve depicting classification performance (overall AUC = 0.998) (h). (i-j) Cell-typing in five RPE-1 derived cell lines<sup>37</sup> (RPE-1, BM510, C7, C29, and C11) (i), and ROC curve depicting classification performance (overall AUC = 0.9648) (j).



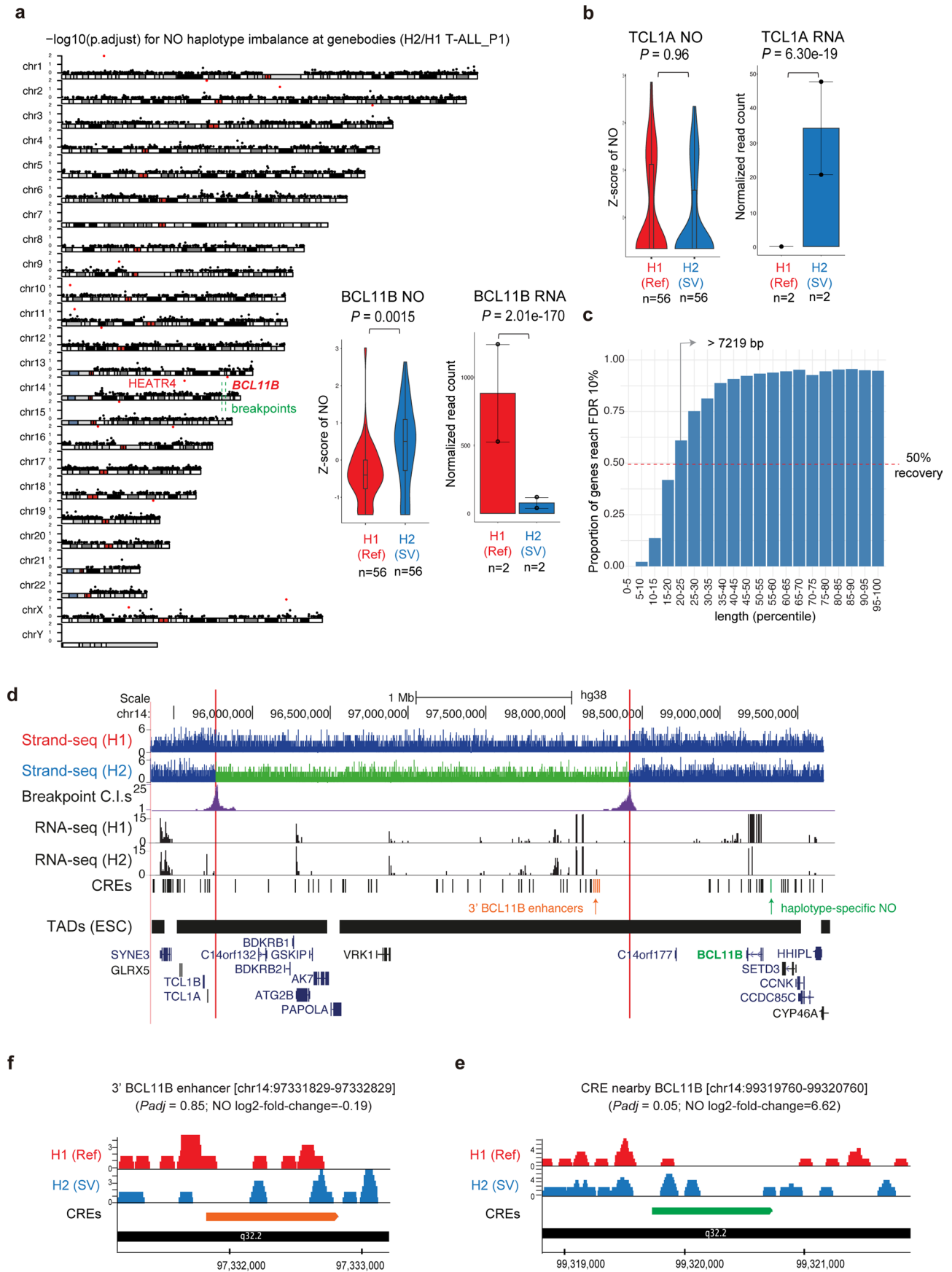
**Extended Data Fig. 4 | *In silico* downsampling experiments.** We performed *in silico* cell mixing of RPE-1 and HG01573 cells to simulate application of scNOVA to different cell fractions (CFs). In this analysis six different CF ranges were considered (20, 10, 5, 3.3, 2, and 1.3). For each *in silico* cell mixing experiment, a total of 150 single cells were randomly subsampled for the major pseudo-clone (containing RPE-1 cells) and the minor pseudo-clone (HG01573 cells), by controlling the minor pseudo-clone CF at 20, 10, 5, 3.3, 2, and 1.3%, respectively.

AUC, area under the curve. DEGs, differentially expressed genes. For each CF, we performed random subsampling of single-cell libraries 10 times, and depicted the respective mean AUC in the plot. Two different analysis modes - default (dashed lines, CNN with negative binomial generalized linear model), and alternative (solid lines, CNN with PLS-DA) are depicted. When the CF is larger than 10%, the default mode performs better, whereas for CFs smaller than 10%, the alternative mode outperforms the default mode.



**Extended Data Fig. 5 | Haplotype-specific NO analysis in RPE-1 and BM510.** (a) Haplotype-specific NO analysis of NO at gene bodies genome-wide in RPE-1 (a) and BM510 (b). For each chromosomal karyogram, the y-axis indicates the significance of haplotype-specific NO for each gene ( $-\log_{10} p.adjust$ ). All the significant genes were indicated in red dots (FDR 10%; two-sided wilcoxon rank sum test followed by Benjamini Hochberg multiple correction; derived from  $n = 33$  cells and  $n = 79$  cells for RPE-1 and BM510, respectively; Boxplots were defined by minima = 25th percentile - 1.5X interquartile range (IQR), maxima = 75th percentile + 1.5X IQR, center = median, and bounds of box = 25th and 75th percentile.). *NTRK3* (identified in BM510) is the only significant gene adjacent

to an SV breakpoint. Haplotype-resolved RNA expression at the *NTRK3* locus is depicted using bar graphs in the right panel (two-sided likelihood ratio test followed by Benjamini Hochberg multiple correction;  $n = 2$  biological replicates; Data are presented as mean values  $\pm$  SEM). (c-d) Haplotype-specific NO analysis at CREs. Browser track depicts the haplotype-resolved NO of the not rearranged (Ref) homolog in red, and the SV homolog in blue. scNOVA identified two CREs with significant haplotype-specific NO, including an intergenic CRE spanning chr15:87527100-87528100 ( $p.adjust = 0.029$ , log<sub>2</sub>-fold change = -2.01) (c) and an intronic CRE at chr15:88246388-88247388 ( $p.adjust = 0.076$ , log<sub>2</sub>-fold change = -1.39) (d).

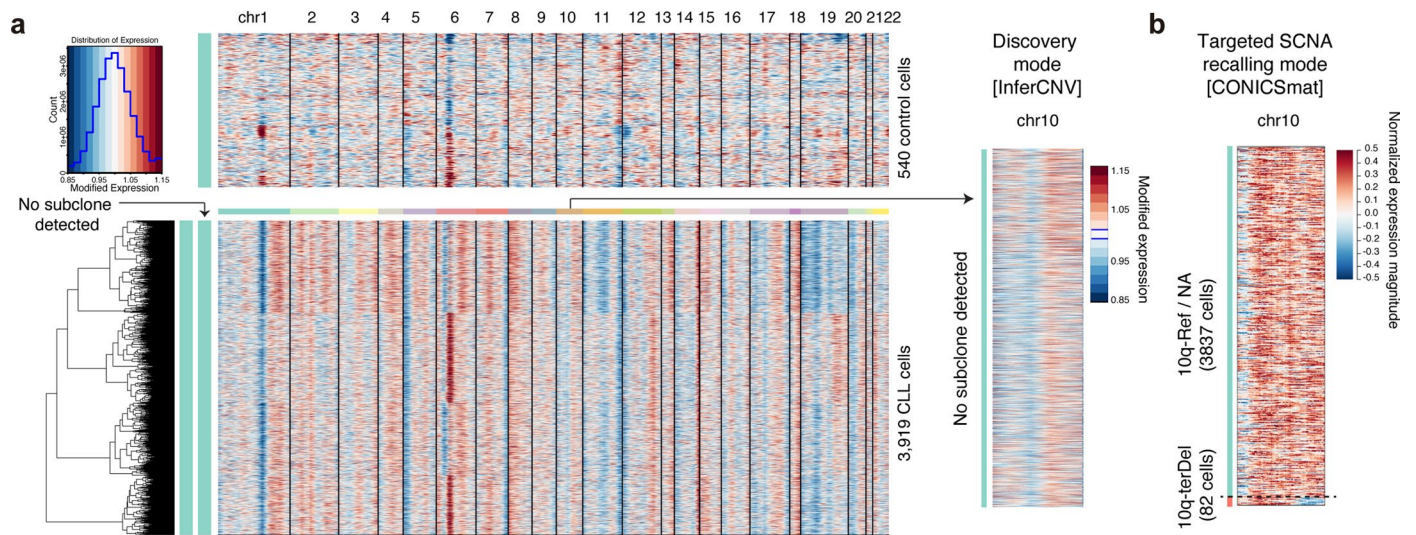


Extended Data Fig. 6 | See next page for caption.



**Extended Data Fig. 6 | Haplotype-specific NO analysis in T-ALL\_P1.** (a) For each chromosomal karyogram, the y-axis indicates the significance of haplotype-specific NO at each gene ( $-\log_{10} p_{\text{adjust}}$ ). Genes with haplotype-specific NO are indicated using red dots (FDR 10%). An inset figure depicts haplotype-specific NO (two-sided wilcoxon rank sum test and Benjamini Hochberg multiple correction;  $n = 56$  cells) and RNA expression at the *BCL11B* gene locus (two-sided likelihood ratio test and Benjamini Hochberg multiple correction;  $n = 2$  biological replicates), which has a nearby somatic SV (within 1 Megabase) and represents the (only) predicted local SV effect. (b) We did not measure haplotype-specific NO for *TCL1A* (two-sided wilcoxon rank sum test and Benjamini Hochberg multiple correction;  $n = 56$  cells), a small gene with 4229 bp in size, in spite of its haplotype-specific gene expression<sup>24</sup> (two-sided likelihood ratio test and Benjamini Hochberg multiple correction;  $n = 2$  biological replicates). Boxplots were defined by minima=25th percentile-1.5X interquartile range (IQR), maxima=75th percentile+1.5X IQR, center=median, and bounds of box=25th and

75th percentile. For bargraphs, data are presented as mean values  $\pm$  SEM (a-b). (c) Simulation analysis revealed a minimum gene length (7219 bp) needed to robustly detect haplotype-specific NO at gene bodies, a gene length met by 80% of genes in the genome (Supplementary Notes). (d) Inversion breakpoints and rearranged TADs. Known 3' *BCL11B* enhancers<sup>112</sup> are depicted in orange. In the not rearranged haplotype, they are located proximal to *BCL11B*, but in the inverted haplotype these enhancers they are located far away from *BCL11B*, and proximal to *TCL1A* in the different TAD boundary. (e) Application of scNOVA identified an intergenic CRE near the *BCL11B* with haplotype-specific NO. The browser track depicts the haplotype-resolved NO of the not rearranged (Ref) homolog in red and the SV homolog in blue. (f) The known 3' *BCL11B* enhancer does not show significant haplotype-specific NO, but the inversion physically relocates these enhancers to the far distance from the *BCL11B*. A representative CRE is shown amongst four CREs overlapping with known 3' *BCL11B* enhancers.



**Extended Data Fig. 7 | Inference of SCNAs using CITE-seq data from the CLL\_24 sample.** (a) InferCNV<sup>48</sup> analysis of 3,919 high quality CLL cells, and 540 control cells (cells sequenced by CITE-seq not originating from the B-cell lineage; see Supplementary Fig. 25), profiled by CITE-seq. This analysis did not discover any subclones in CLL\_24. (Note that the high variability observed on the 6p-arm, not only seen in CLL cells but also in control cells, likely arose from the

presence of *MHC* genes in this locus, whose expression is cell cycle dependent<sup>113</sup>.) (b) CONICSmat based targeted SCNA recalling of the 10q-terDel (previously discovered in SCb; see Fig. 4b) using the high-resolution breakpoints derived from Strand-seq. Use of these SV breakpoints allowed CONICSmat to confidently call the 10q-terDel in 82 single cells from the CITE-seq data.

## Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

- |     |           |
|-----|-----------|
| n/a | Confirmed |
|-----|-----------|
- The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement
  - A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
  - The statistical test(s) used AND whether they are one- or two-sided  
*Only common tests should be described solely by name; describe more complex techniques in the Methods section.*
  - A description of all covariates tested
  - A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
  - A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
  - For null hypothesis testing, the test statistic (e.g.  $F$ ,  $t$ ,  $r$ ) with confidence intervals, effect sizes, degrees of freedom and  $P$  value noted  
*Give  $P$  values as exact values whenever suitable.*
  - For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
  - For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
  - Estimates of effect sizes (e.g. Cohen's  $d$ , Pearson's  $r$ ), indicating how they were calculated

*Our web collection on [statistics for biologists](#) contains articles on many of the points above.*

### Software and code

Policy information about [availability of computer code](#)

Data collection	Our study uses sequencing data output by illumina sequencers and hence no special software was used for collecting it.
Data analysis	<p>The computational code of our analytical framework is hosted on GitHub (see <a href="https://github.com/jeongdo801/scNOVA">https://github.com/jeongdo801/scNOVA</a>). All code is available freely for academic research.</p> <p>Other software used:  Mosaiccatcher (<a href="https://github.com/friendsofstrandseq/mosaiccatcher-pipeline">https://github.com/friendsofstrandseq/mosaiccatcher-pipeline</a>), StrandPhaseR (<a href="https://github.com/daewoooo/StrandPhaseR">https://github.com/daewoooo/StrandPhaseR</a>), InferCNV (<a href="https://github.com/broadinstitute/inferCNV/">https://github.com/broadinstitute/inferCNV/</a>), HoneyBADGER (<a href="https://jef.works/HoneyBADGER/">https://jef.works/HoneyBADGER/</a>), CONICSmats (<a href="https://github.com/diazlab/CONICS">https://github.com/diazlab/CONICS</a>), NucTools (<a href="https://homeveg.github.io/nuctools">https://homeveg.github.io/nuctools</a>), Delly2 (<a href="https://github.com/dellytools/delly">https://github.com/dellytools/delly</a>), BWA (v0.7.15), STAR (v2.7.9a), SAMtools (v1.3.1), biobambam2 (v2.0.76), deeptools (v2.5.1), perl (v5.16.3), Python (v3.7.4), cuDNN (v7.6.4.38), CUDA (v10.1.243), TensorFlow (v1.15.0), scikit-learn (v0.21.3), matplotlib (v3.1.1), R version 4.0.0, DESeq2, FlowJo, BD FACSDiva™</p>

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.

## Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

Sequencing data from this study can be retrieved from the European Genome-phenome Archive (EGA), and the European Nucleotide Archive (ENA) [accessions: LCL data are available under the following accessions: Strand-seq (PRJEB39750, PRJEB55038); RNA-seq (ERP123231); WGS (PRJEB37677). C11 cell line data are available under the accession PRJEB55012. Leukemia patient data and human primary cells derived data were deposited in the European Genome-phenome Archive (EGA), under the following accession numbers: skin fibroblast (EGAS00001006498); cord blood (EGAS00001006567). T-ALL Strand-seq and scRNA-seq (EGAS00001003365), CLL Strand-seq (EGAS00001004925), AML Strand-seq (EGAS00001004903), T-ALL bulk RNA-seq (EGAS00001003248), CLL bulk RNA-seq (EGAS00001005746), CLL CITE-seq (EGAS00001004925).] Access to human patient data is governed by the EGA Data Access Committee.

## Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

- Life sciences       Behavioural & social sciences       Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://nature.com/documents/nr-reporting-summary-flat.pdf)

## Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	No sample-size calculation was performed, since this study focuses on establishing a novel methodology rather than on performing statistical tests between groups of samples.
Data exclusions	We excluded low quality single-cell libraries that showed very low, uneven coverage, or an excess of 'background reads' yielding noisy single cell data prior to analysis. Cells with incomplete BrdU incorporation or cells undergoing more than one DNA synthesis phase under BrdU exposure are largely excluded during cell sorting and thus get only rarely sequenced during Strand-seq experiments.
Replication	To ensure reproducibility of our computational findings we have organized our main workflow using Snakemake, a widely used workflow manager, and we provide the workflow description (Snakefile) along with a Bioconda environment that facilitates easy installation of all dependencies (with well-defined versions). We have repeated the analyses of our datasets and can confirm consistent and reproducible results from these workflows. To ensure reproducibility of our experimental findings, we generated replicates wherever possible, which confirmed reproducibility of the result.
Randomization	Does not apply (there are no experimental groups in our study)
Blinding	Does not apply. (this study focuses on intra-sample comparison rather than performing statistical tests between groups of samples).

## Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

### Materials & experimental systems

n/a	Involved in the study
<input type="checkbox"/>	<input checked="" type="checkbox"/> Antibodies
<input type="checkbox"/>	<input checked="" type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input type="checkbox"/>	<input checked="" type="checkbox"/> Human research participants
<input type="checkbox"/>	<input checked="" type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

### Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input type="checkbox"/>	<input checked="" type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

## Antibodies

Antibodies used

FACS (clone, manufacturer, catalogue number, lot number): APC mouse anti-human CD34 (clone 581; Biolegend; #343509; Lot:

## Antibodies used

B260867), PeCy7 mouse anti-human CD38 (clone HB7; eBioscience; #15538396; Lot: 1974952), FITC mouse anti-human CD45Ra (clone HI100; eBioscience; #15526406; Lot: 4329359), PE mouse anti-human CD90 (clone 5E10; eBioscience; # 15526836; Lot:1982684), PE anti-murine CD45 (clone 30-F11; Biolegend; #103106; Lot: B361031). CITE-seq (clone, manufacturer, catalogue number, barcode number): CD10 (HI10a; Biolegend; 312231; 0062), CD103 (Ber-ACT8 ; Biolegend; 350231; 0145), CD11b (ICRF44; Biolegend; 301353; 0161), CD11c (S-HCL-3; Biolegend; 371519; 0053), CD127 (A019D5; Biolegend; 351352; 0390), CD134 (Ber-ACT35; Biolegend; 350033; 0158), CD137 (4B4-1; Biolegend; 309835; 0355), CD150 (A12 (7D4); Biolegend; 306313; 0870), CD152 (BN13; Biolegend; 369619; 0151), CD16 (3G8; Biolegend; 302061; 0083), CD161 (HP-3G10; Biolegend; 339945; 0149), CD183 (G025H7; Biolegend; 353745; 0140), CD184 (12G5; Biolegend; 306531; 0366), CD185 (J252D4; Biolegend; 356937; 0144), CD19 (HIB19; Biolegend; 302259; 0050), CD194 (L291H4; Biolegend; 359423; 0071), CD195 (J418F1; Biolegend; 359135; 0141), CD197 (G043H7; Biolegend; 353247; 0148), CD2 (TS1/8; Biolegend; 309229; 0367), CD20 (2H7; Biolegend; 302359; 0100), CD21 (Bu32; Biolegend; 354915; 0181), CD22 (S-HCL-1; Biolegend; 363514; 0393), CD223 (11C3C65; Biolegend; 369333; 0152), CD23 (EBVCS-5; Biolegend; 338523; 0897), CD24 (ML5; Biolegend; 311137; 0180), CD244 (C1.7; Biolegend; 329527; 0189), CD25 (BC96; Biolegend; 302643; 0085), CD27 (O323; Biolegend; 302847; 0154), CD273 (24F.10C12; Biolegend; 329619; 0008), CD274 (29E.2A3; Biolegend; 329743; 0007), CD278 (C398.4A; Biolegend; 313555; 0171), CD279 (EH12.2H7; Biolegend; 329955; 0088), CD28 (CD28.2; Biolegend; 302955; 0386), CD29 (TS2/16; Biolegend; 303027; 0369), CD3 (UCHT1; Biolegend; 300475; 0034), CD31 (W5M59; Biolegend; 303137; 0124), CD32 (FUN-2; Biolegend; 303223; 0142), CD357 (108-17; Biolegend; 371225; 0360), CD366 (F38-2E2; Biolegend; 345047; 0169), CD38 (HIT2; Biolegend; 303541; 0389), CD39 (A1; Biolegend; 328233; 0176), CD4 (RPA-T4; Biolegend; 300563; 0072), CD43 (CD43-10G7; Biolegend; 343209; 0357), CD44 (IM7; Biolegend; 103045; 0073), CD45 (HI30; Biolegend; 304064; 0391), CD45RA (HI100; Biolegend; 304157; 0063), CD45RO (UCHL1; Biolegend; 304255; 0087), CD47 (CC2C6; Biolegend; 323129; 0026), CD48 (BJ40; Biolegend; 336709; 0029), CD5 (UCHT2; Biolegend; 300635; 0138), CD56 (QA17A16; Biolegend; 392421; 0084), CD57 (QA17A04; Biolegend; 393319; 0168), CD62L (DREG-56; Biolegend; 304847; 0147), CD69 (FN50; Biolegend; 310947; 0146), CD7 (CD7-6B7; Biolegend; 343123; 0066), CD70 (113-16; Biolegend; 355117; 0027), CD73 (AD2; Biolegend; 344029; 0577), CD79b (CB3-1; Biolegend; 341415; 0187), CD86 (IT2.2; Biolegend; 305443; 0006), CD8a (RPA-T8; Biolegend; 301067; 0080), CD95 (DX2; Biolegend; 305649; 0156), Kappa (MHK-49 ; Biolegend; 316531; 0894), KLRG1 (SA231A2; Biolegend; 367721; 0153), Lambda (MHL-38 ; Biolegend; 316627; 0898), TIGIT (A15153G; Biolegend; 372725; 0089), Isotype Ctrl (MOPC-21; Biolegend; 400199; 0090), Isotype Ctrl (HTK888; Biolegend; 400973; 0241), Isotype Ctrl (MPC-11; Biolegend; 400373; 0092), Isotype Ctrl (RTK4530; Biolegend; 400673; 0095), Isotype Ctrl (MOPC-173; Biolegend; 400285; 0091)

## Validation

All antibodies were validated for the specific application by the manufacturer and validation data is available on the manufacturer's website.

## FACS

CD34 CD34 <https://www.biolegend.com/fr-ch/products/apc-anti-human-cd34-antibody-6090> DOI: 10.1538/expanim.49.97  
 CD38 CD38 <https://www.thermofisher.com/antibody/product/CD38-Antibody-clone-HB7-Monoclonal/25-0388-42> DOI: 10.1016/j.stem.2021.02.001  
 CD45Ra PTPRC <https://www.thermofisher.com/antibody/product/CD45RA-Antibody-clone-HI100-Monoclonal/14-0458-82> DOI: 10.1080/2162402X.2017.1371399  
 CD90 THY1 <https://www.thermofisher.com/antibody/product/CD90-Thy-1-Antibody-clone-eBio5E10-5E10-Monoclonal/12-0909-42> DOI: 10.1186/s41232-017-0049-2  
 CD45 PTPRC <https://www.biolegend.com/fr-fr/products/pe-anti-mouse-cd45-antibody-100> DOI: 10.4049/jimmunol.176.11.6532  
 CITE-seq  
 CD10 MME <https://www.biolegend.com/en-us/search-results/totalseq-a0062-anti-human-cd10-antibody-15949?GroupID=BLG5905> doi: 10.1084/jem.181.6.2271  
 CD103 ITGAE <https://www.biolegend.com/nl-be/products/totalseq-a0145-anti-human-cd103-integrin-%CE%B1e-antibody-16194> doi: 10.1538/expanim.49.97  
 CD11b ITGAM <https://www.biolegend.com/en-us/products/totalseq-a0161-anti-human-cd11b-antibody-15927> doi: 10.4049/jimmunol.1302846  
 CD11c ITGAX <https://www.biolegend.com/de-de/products/totalseq-a0053-anti-human-cd11c-antibody-16231> doi: 10.1016/j.jri.2011.01.014  
 CD127 IL7R <https://www.biolegend.com/en-us/search-results/totalseq-a0390-anti-human-cd127-il-7alpha-antibody-15943?GroupID=BLG9274> doi: 10.1038/nbt.3973  
 CD134 TNFRSF4 <https://www.biolegend.com/en-us/products/totalseq-a0158-anti-human-cd134-ox40-antibody-16437> doi: 10.2215/CJN.06460612  
 CD137 TNFRSF9 <https://www.biolegend.com/en-us/products/totalseq-trade-a0355-anti-human-cd137-4-1bb-antibody-16737> doi: 10.4049/jimmunol.165.5.2903  
 CD150 SLAMF1 <https://www.biolegend.com/fr-ch/search-results/totalseq-a0870-anti-human-cd150-slam-antibody-18039> doi: 10.1016/j.cell.2021.12.018  
 CD152 CTLA4 <https://www.biolegend.com/de-at/products/totalseq-a0151-anti-human-cd152-ctla-4-antibody-15707> doi: 10.1084/jem.176.6.1595  
 CD16 FCGR3A <https://www.biolegend.com/de-de/products/totalseq-a0083-anti-human-cd16-antibody-15765> doi: 10.1189/jlb.0408244  
 CD161 KLRB1 <https://www.biolegend.com/nl-be/products/totalseq-a0149-anti-human-cd161-antibody-16156> doi: 10.1084/jem.188.5.867  
 CD183 CXCR3 <https://www.biolegend.com/nl-nl/search-results/totalseq-a0140-anti-human-cd183-cxcr3-antibody-16163> DOI: 10.1016/j.cell.2021.12.018  
 CD184 CXCR4 <https://www.biolegend.com/en-gb/products/totalseq-a0366-anti-human-cd184-cxcr4-antibody-17277> DOI: 10.1074/jbc.M610931200  
 CD185 CXCR5 <https://www.biolegend.com/ja-jp/products/totalseq-a0144-anti-human-cd185-cxcr5-antibody-16330> DOI: 10.7554/eLife.63632  
 CD19 CD19 <https://www.biolegend.com/nl-nl/products/totalseq-a0050-anti-human-cd19-antibody-15777> DOI: 10.3324/haematol.2009.013151  
 CD194 CCR4 <https://www.biolegend.com/de-de/products/totalseq-a0071-anti-human-cd194-ccr4-antibody-16170> doi: 10.1016/j.xpro.2021.100900  
 CD195 CCR5 <https://www.biolegend.com/en-us/search-results/totalseq-a0141-anti-human-cd195-ccr5-antibody-16161> DOI: 10.1016/j.cell.2021.12.018  
 CD197 CCR7 <https://www.biolegend.com/en-gb/search-results/totalseq-a0148-anti-human-cd197-ccr7-antibody-16352?GroupID=BLG9613> DOI: 10.1016/j.cell.2019.05.031

CD2 CD2 <https://www.biolegend.com/de-de/clone-search/totalseq-a0367-anti-human-cd2-antibody-16714> DOI:<https://doi.org/10.1074/jbc.271.10.5369>

CD20 MS4A1 <https://www.biolegend.com/ja-jp/products/totalseq-a0100-anti-human-cd20-antibody-16173>  
DOI: 10.1203/01.PDR.0000130480.51066.FB

CD21 CR2 <https://www.biolegend.com/fr-ch/products/totalseq-a0181-anti-human-cd21-antibody-16203>  
GroupID=ImportedGROUP1 DOI: 10.1002/eji.1830260714

CD22 CD22 <https://www.biolegend.com/en-us/products/totalseq-a0393-anti-human-cd22-antibody-15936> DOI: 10.1016/j.coi.2005.03.005

CD223 LAG3 <https://www.biolegend.com/fr-ch/products/totalseq-a0152-anti-human-cd223-lag-3-antibody-16157>  
GroupID=BLG14890 DOI: 10.1016/j.cell.2021.12.018

CD23 FCER2 <https://www.biolegend.com/en-us/search-results/totalseq-a0897-anti-human-cd23-antibody-18091> DOI: 10.1128/JVI.46.3.800-807.1983

CD24 CD24 <https://www.biolegend.com/en-us/products/totalseq-a0180-anti-human-cd24-antibody-16331> PMID: 14581365

CD244 CD244 <https://www.biolegend.com/en-us/search-results/totalseq-a0189-anti-human-cd244-2b4-antibody-16196>  
GroupID=BLG8490 DOI: 10.1182/blood-2011-02-339135

CD25 IL2RA <https://www.biolegend.com/fr-lu/products/totalseq-a0085-anti-human-cd25-antibody-15770> DOI: 10.1016/j.cell.2019.05.031

CD27 CD27 <https://www.biolegend.com/ja-jp/products/totalseq-a0154-anti-human-cd27-antibody-16174> DOI: 10.1016/j.cell.2019.05.031

CD273 PDCD1LG2 <https://www.biolegend.com/de-de/products/totalseq-a0008-anti-human-cd273-b7-dc-pd-l2-antibody-15932>  
DOI: 10.4049/jimmunol.170.3.1257

CD274 PDCD1LG1 <https://www.biolegend.com/de-at/products/totalseq-a0007-anti-human-cd274-b7-h1-pd-l1-antibody-16195>  
GroupID=BLG5404 DOI: 10.4049/jimmunol.170.3.1257

CD278 ICOS <https://www.biolegend.com/en-gb/products/totalseq-a0171-anti-human-mouse-rat-cd278-icos-antibody-17152>  
DOI: 10.4049/jimmunol.171.2.783

CD279 PDCD1 <https://www.biolegend.com/de-at/products/totalseq-a0088-anti-human-cd279-pd-1-antibody-15772> DOI: 10.4049/jimmunol.181.10.6707

CD28 CD28 <https://www.biolegend.com/de-de/products/totalseq-a0386-anti-human-cd28-antibody-16787> DOI: 10.1016/j.febslet.2006.11.044

CD29 ITGB1 <https://www.biolegend.com/en-us/products/totalseq-a0369-anti-human-cd29-antibody-16664>?GroupID=BLG10310  
DOI: 10.1182/blood-2004-07-2598

CD3 CD3E <https://www.biolegend.com/de-at/products/totalseq-a0034-anti-human-cd3-antibody-15707> DOI: 10.4049/jimmunol.180.11.7431

CD31 PECAM1 <https://www.biolegend.com/en-gb/products/totalseq-a0124-anti-human-cd31-antibody-16332> DOI: 10.1182/blood-2006-10-047092

CD32 FCGR2A <https://www.biolegend.com/en-us/products/totalseq-a0142-anti-human-cd32-antibody-16168> DOI: 10.1182/blood-2010-11-316158

CD357 TNFRSF18 <https://www.biolegend.com/fr-lu/products/totalseq-a0360-anti-human-cd357-gitr-antibody-17349>  
GroupID=BLG15183 DOI: 10.1016/j.cell.2021.12.018

CD366 HAVCR2 <https://www.biolegend.com/fr-fr/products/totalseq-a0169-anti-human-cd366-tim-3-antibody-17350> DOI: 10.1182/blood-2008-02-142596

CD38 CD38 <https://www.biolegend.com/fr-fr/products/totalseq-a0389-anti-human-cd38-antibody-16899> DOI: 10.1538/expanim.49.97

CD39 ENTPD1 <https://www.biolegend.com/it-it/products/totalseq-a0176-anti-human-cd39-antibody-16204> DOI: 10.7554/eLife.63632

CD4 CD4 <https://www.biolegend.com/fr-ch/products/totalseq-a0072-anti-human-cd4-antibody-15762> DOI: 10.7554/eLife.63632

CD43 SPN <https://www.biolegend.com/fr-ch/products/totalseq-a0357-anti-human-cd43-antibody-17546> PMID: 7507092

CD44 CD44 <https://www.biolegend.com/nl-nl/products/totalseq-a0073-anti-mouse-human-cd44-antibody-15923>  
DOI: 10.1186/1479-5876-7-89

CD45 PTPRC <https://www.biolegend.com/en-gb/search-results/totalseq-a0391-anti-human-cd45-antibody-15934>  
GroupID=GROUP658 DOI: 10.1038/emboj.2012.192

CD45RA PTPRC <https://www.biolegend.com/it-it/products/totalseq-a0063-anti-human-cd45ra-antibody-15775> DOI: 10.4049/jimmunol.0901967

CD45RO PTPRC <https://www.biolegend.com/nl-nl/products/totalseq-a0087-anti-human-cd45ro-antibody-15771> DOI: 10.4049/jimmunol.180.11.7431

CD47 CD47 <https://www.biolegend.com/de-at/products/totalseq-a0026-anti-human-cd47-antibody-15957> PMID: 10572074

CD48 CD48 <https://www.biolegend.com/en-gb/products/totalseq-a0029-anti-human-cd48-antibody-15942> DOI: 10.1189/jlb.0611308

CD5 CD5 <https://www.biolegend.com/en-gb/clone-search/totalseq-a0138-anti-human-cd5-16333>?GroupID=BLG5902 DOI: 10.1073/pnas.1001515107

CD56 NCAM1 <https://www.biolegend.com/en-us/products/totalseq-a0084-anti-human-cd56-recombinant-antibody-15766>  
GroupID=GROUP28 DOI: 10.1186/s40364-020-00253-w

CD57 B3GAT1 <https://www.biolegend.com/fr-fr/products/totalseq-a0168-anti-human-cd57-recombinant-antibody-17680>  
DOI: 10.1016/j.xpro.2021.100900

CD62L SELL <https://www.biolegend.com/de-at/products/totalseq-a0147-anti-human-cd62l-antibody-16334> DOI: 10.4049/jimmunol.181.9.6563

CD69 CD69 <https://www.biolegend.com/de-at/products/totalseq-a0146-anti-human-cd69-antibody-16200> DOI: 10.1093/toxsci/kfp224

CD7 CD7 <https://www.biolegend.com/fr-fr/products/totalseq-a0066-anti-human-cd7-antibody-15944> DOI: 10.1038/nbt.3973

CD70 CD70 <https://www.biolegend.com/it-it/products/totalseq-a0027-anti-human-cd70-antibody-16184> DOI: 10.1182/blood-2009-08-239145

CD73 NT5E <https://www.biolegend.com/en-us/search-results/totalseq-a0577-anti-human-cd73-ecto-5-nucleotidase-antibody-16773>  
DOI: 10.1016/j.joen.2011.05.022

CD79b CD79B <https://www.biolegend.com/nl-nl/products/totalseq-a0187-anti-human-cd79b-ig%CE%B2-antibody-16433>  
DOI: 10.1016/j.cell.2019.05.031

CD86 CD86 <https://www.biolegend.com/en-gb/products/totalseq-a0006-anti-human-cd86-antibody-15937>?GroupID=BLG11941

DOI: 10.7554/eLife.63632  
 CD8a CD8A <https://www.biolegend.com/fr-ch/products/totalseq-a0080-anti-human-cd8a-antibody-15763>  
 DOI: 10.1186/1479-5876-7-89  
 CD95 FAS <https://www.biolegend.com/de-de/search-results/totalseq-a0156-anti-human-cd95-fas-antibody-16363> DOI: 10.4049/jimmunol.0903133  
 Kappa IGKC <https://www.biolegend.com/it-it/products/totalseq-a0894-anti-human-ig-light-chain-kappa-antibody-17854>  
 DOI: 10.1016/j.bbmt.2013.06.007  
 KLRG1 KLRG1 <https://www.biolegend.com/en-us/search-results/totalseq-a0153-anti-human-klrg1-afa-antibody-16530?GroupID=GROUP28> DOI: 10.7554/eLife.63632  
 Lambda IGLC2 <https://www.biolegend.com/it-it/search-results/totalseq-a0898-anti-human-ig-light-chain-lambda-antibody-18163>  
 DOI: 10.1016/j.cell.2021.12.018  
 TIGIT TIGIT <https://www.biolegend.com/nl-be/products/totalseq-a0089-anti-human-tigit-antibody-15773> DOI: 10.1038/s41388-018-0288-y  
 Isotype Ctrl Mouse IgG1, κ <https://www.biolegend.com/ja-jp/products/totalseq-a0090-mouse-igg1-kappa-isotype-control-15774>  
 DOI: 10.1182/blood-2011-02-339135  
 Isotype Ctrl Armenian Hamster IgG <https://www.biolegend.com/de-de/products/totalseq-a0241-armenian-hamster-igg-isotype-ctrl-17278?Clone=HTK888> DOI: 10.1189/jlb.1107802  
 Isotype Ctrl Mouse IgG2b, κ <https://www.biolegend.com/it-it/products/totalseq-a0092-mouse-igg2b-kappa-isotype-control-15778>  
 DOI: 10.4049/jimmunol.180.12.7989  
 Isotype Ctrl Rat IgG2b, κ <https://www.biolegend.com/it-it/products/totalseq-a0095-rat-igg2b-kappa-isotype-ctrl-16228>  
 DOI: 10.4049/jimmunol.181.1.104  
 Isotype Ctrl Mouse IgG2a, κ <https://www.biolegend.com/nl-be/products/totalseq-a0091-mouse-igg2a-kappa-isotype-control-15779>  
 DOI: 10.1016/j.immuni.2020.08.004

## Eukaryotic cell lines

### Policy information about cell lines

Cell line source(s)	hTERT RPE-1 cells were purchased from ATCC (CRL-4000) and checked for mycoplasma contamination. C11 cells were derived in-house (from hTERT RPE-1 cells) as described previously (PMID: 32268084). GM20509, and HG01505 cell lines were purchased from Coriell and taken into culture at passage 4 (early) and passage 8 (late).
Authentication	Authentication was performed by confirming the presence of known somatic DNA rearrangements in these cell lines (e.g. the previously-reported unbalanced translocation in the case of RPE-1). Additional authentication was done at the level of SNPs shared between the cell lines, which are derived from the same anonymous donor.
Mycoplasma contamination	The cell lines tested negative for mycoplasma
Commonly misidentified lines (See ICLAC register)	No commonly misidentified lines used in this study.

## Human research participants

### Policy information about studies involving human research participants

Population characteristics	We included one AML sample donor (AML_1), one CLL sample donor (CLL24), and one T-ALL sample donor (TALL P1) in this study. Patient AML_1 was 32 years of age, male, and the sample was obtained as a diagnostic bone marrow from the first aspiration of an AML with a t(8;21) translocation, arising after cytostatic therapy for testicular cancer. CLL_24 patient was 61 years old age, female, and the sample was obtained from the peripheral blood mononuclear cells. TALL patient P1 was 12 years of age, female, diagnosed with acute lymphoblastic leukemia (ALL) and the relapse sample was obtained for analysis.
Recruitment	Information about enrollment for P1 (AIEOP-ALL BFM 2009) can be found here: <a href="https://www.kinderkrebsinfo.de/fachinformationen/studieenportal/abgeschlossene_studieen_register/aieop_bfmm_all_2009/indeex_ger.html">https://www.kinderkrebsinfo.de/fachinformationen/studieenportal/abgeschlossene_studieen_register/aieop_bfmm_all_2009/indeex_ger.html</a> or: <a href="https://clinicaltrials.gov/t2/show/NCT01117441">https://clinicaltrials.gov/t2/show/NCT01117441</a> There was no sample selection bias which may affect to the results.
Ethics oversight	Samples used in this analysis have received approval from the relevant institutional review boards and ethics committees (University of Kiel). Written informed consent had been obtained from all the patients and the experiments conformed to the principles set out in the WMA Declaration of Helsinki and the Department of Health and Human Services Belmont Report. The protocols used in this study received approval from the relevant institutional review boards and local ethics committees. Written informed consent was obtained from patients, and all experiments were consistent with current bioethical policies. T-ALL experiments were approved by the ethics commission of the Kanton Zurich (approval number 2014-0383).

Note that full information on the approval of the study protocol must also be provided in the manuscript.

## Clinical data

### Policy information about clinical studies

All manuscripts should comply with the ICMJE [guidelines for publication of clinical research](#) and a completed [CONSORT checklist](#) must be included with all submissions.

Clinical trial registration	P1: <a href="https://clinicaltrials.gov/ct2/show/NCT01117441">https://clinicaltrials.gov/ct2/show/NCT01117441</a>
Study protocol	The patients were treated according to the respective protocols AIEOP-ALL 2009 (P1) (details of which can be found above, in

Study protocol	Recruitment section)
Data collection	This study is prospective, controlled, randomized and multi-centered. More than 70 clinics in Germany, Austria and Switzerland are participating in the study. NCT01117441 (P1) Clinical Trial title: International Collaborative Treatment Protocol For Children And Adolescents With Acute Lymphoblastic Leukemia; Allocation: Randomized; Enrollment: 4750 participants; Study Start Date: June 2010; Estimated Completion Date : December 2021
Outcomes	Since the study is still ongoing (in the case of NCT01117441 recruitment is not yet completed) and the follow-up time is too short, outcomes are not yet available and not applicable.

## Flow Cytometry

### Plots

Confirm that:

- The axis labels state the marker and fluorochrome used (e.g. CD4-FITC).
- The axis scales are clearly visible. Include numbers along axes only for bottom left plot of group (a 'group' is an analysis of identical markers).
- All plots are contour plots with outliers or pseudocolor plots.
- A numerical value for number of cells or percentage (with statistics) is provided.

### Methodology

Sample preparation	Primary human T-ALL cells were recovered from cryopreserved bone marrow aspirates of patients enrolled in the ALL-BFM 2009 study. Patient-derived xenografts (PDX) were generated by intrafemoral injection of 1 Million viable primary ALL cells in NSG mice. PDX-derived (P1) cells were frozen until processing. Primary human AML cells were isolated via bone marrow aspiration from the donor AML_1 during diagnosis, and frozen until processing. All samples were then processed as follows: cryopreserved cells were thawed rapidly at 37 °C and resuspended in 10 ml warm Roswell Park Memorial Institute (RPMI) medium with 100 µg/ml Dnase I. Cells were centrifuged for 5 mins at 300 g, and resuspended in ice-cold phosphate buffered saline (PBS) with 2% foetal bovine serum (FBS) and 5mM EDTA. Samples were then stained on ice in the dark for 30 mins as follows: T-ALL P1 was stained with anti-murine-CD45-PE (mCD45) (clone 30-F11; BioLegend; 1:20), and DAPI (1:200) was added for 5 mins prior to sorting; AML_1 was stained with CD34-APC (clone 581; Biolegend), CD38-PeCy7 (clone HB7; eBioscience), CD45Ra-FITC (clone HI100; eBioscience), CD90-PE (clone 5E10; eBioscience), and LIVE/DEAD™ Fixable Near-IR Dead Cell Stain (Thermofisher). After staining, cells were washed once in 4 ml ice-cold PBS with 2% FBS and 5 mM EDTA and centrifuged at 300 g for 5 mins. Cells were resuspended in ice-cold PBS with 2% FBS and 5 mM EDTA for sorting.
Instrument	BD FACSAria™ Fusion Cell Sorter
Software	FlowJo, BD FACSDiva™
Cell population abundance	Due to limited sample material, post-sort purities were not re-assessed using flow cytometry. Instead, this was done by gating and quantification of populations using Flowjo (Supplemental Figures S21 and S27). In the case of AML_1: 83.5% of the total events were included after gating out debris in the FSC-A vs SSC-A plot; 95.4% of these events were within the Single Cells gate (based on SSC-W vs SSC-A); 51.5% of these Single Cells were gated as Viable Cells (based on Fixable Live/Dead Near-IR viability stain vs FSC-A); and the final sorting population of CD34+ cells represented 46.4 % of the Viable Cells (based on CD34-APC vs CD38-PeCy7). In TALL P1: 86.7% of the total events were retained after debris exclusion (based on FSC-A vs SSC-A); 96% of these events were gated as Single Cells (based on SSC-W vs SSC-A); 20.8% of the Single Cells were gated as Viable Cells (based on DAPI viability stain vs FSC-A); and of these Viable Cells, 72.8% were gated as human, T-ALL cells (based on low murine CD45-PE expression).
Gating strategy	For TALL P1: FSC-A vs SSC-A was the starting gate, wherein debris was excluded. Next, Single Cells were gated based on the exclusion of outliers in SSC-W vs SSC-A. Viable Cells were then gated within this population based on a low staining for DAPI (DAPI-viability vs FSC-A). Finally, human T-ALL cells were discriminated from murine immune cells based on a low expression of murine CD45 (murine CD45 - PE vs GFP). The full gating strategy is depicted in Supplemental Figure S27. For AML_1: The first gate excluded any cellular debris based on FSC-A vs SSC-A. These cells were then sub-gated to identify only Single Cells, based on removal of outliers from the SCC-W vs SSC-A plot. Viable Cells were gated within the Single Cells based on a low intracellular staining for the viability stain Fixable LIVE/DEAD near-IR (Fixable LIVE/DEAD near-IR Viability vs FSC-A). Finally, the ultimate sorting population of CD34+ AML cells was gated based on a high expression of CD34 (CD34-APC vs CD38 PeCy7). The full gating strategy is depicted in Supplemental Figure S21.
<input checked="" type="checkbox"/> Tick this box to confirm that a figure exemplifying the gating strategy is provided in the Supplementary Information.	