# Supplementary Information

## Table of Contents

## Supplementary Figures

## Supplementary Figure 1.

**a**



**b**



**c**



30

31

32 **Supplementary Fig. 1: Validation of Alleloscope genotyping results for the P6335**
33 **colorectal cancer sample with linked-reads sequencing data.** (a) Segmentation of the
34 pooled scDNA-seq data using the HMM algorithm. (b) $\hat{\theta}_i$ values recapitulate CNV carriers
35 that are detected using only coverage for four chromosomal regions on the P6335 tumor
36 sample. Different colors represent different genotype clusters. Phasing accuracy for each
37 region is shown in the title. (c) $\hat{\theta}_i$ calculation using known SNP phases from paired linked-
38 reads sequencing data. Genotyping accuracy is labeled in the plots. The colors follow the
39 clustering results from b. The color scheme is the same as that in Fig. 2 and
40 Supplementary Fig. 5.

**Supplementary Figure 2.**



41
42 **Supplementary Fig. 2: Heatmaps of allele-specific genotypes and haplotype-**
43 **specific genotypes from CHISEL for the P5931 sample.**

**Supplementary Figure 3.**



44

45  **Supplementary Fig. 3: Phasing accuracy for the CNA regions in the P6198 sample**
46  **by comparing to the matched linked-reads sequencing data.** LOH: segments with
47  any LOH events. Amp: segments with amplifications that lead to allelic imbalance. Ctrl:
48  control segments without allelic imbalance.

**Supplementary Figure 4.**

**a**



**b**



49

**Supplementary Fig. 4: Segmentation plot and genotype heatmap for the P6198 sample.** The color scheme is the same as that in Fig. 2 and Supplementary Fig. 5.

5

**Supplementary Figure 5.**



**Supplementary Fig. 5: Segmentation plot and genotype heatmap for the P6335 sample.** The color scheme is the same as that in Fig. 2 and Supplementary Fig. 5.

# Supplementary Figure 6.



**Supplementary Fig. 6: Segmentation plot and genotype heatmap for the BC10X sample.** (a) Genome segmentation using HMM on the pooled coverage signals across the cells. (b) Genotype profiles of five example regions. The coloring scheme is same as that in part (c). (c) Hierarchical clustering of single-cell ASCN genotypes reveals complex subclone structure. Genotypes of the five regions in three example cells from the three major subclones are shown in the left. Different colors represent different genotypes. In the color panel, M and m represent the "Major haplotype" and "minor haplotype" respectively.

**Supplementary Figure 7.**

64

**Supplementary Fig. 7: Segmentation plot and genotype heatmap for the P5846 sample.** The color scheme is the same as that in Fig. 2 and Supplementary Fig. 5.
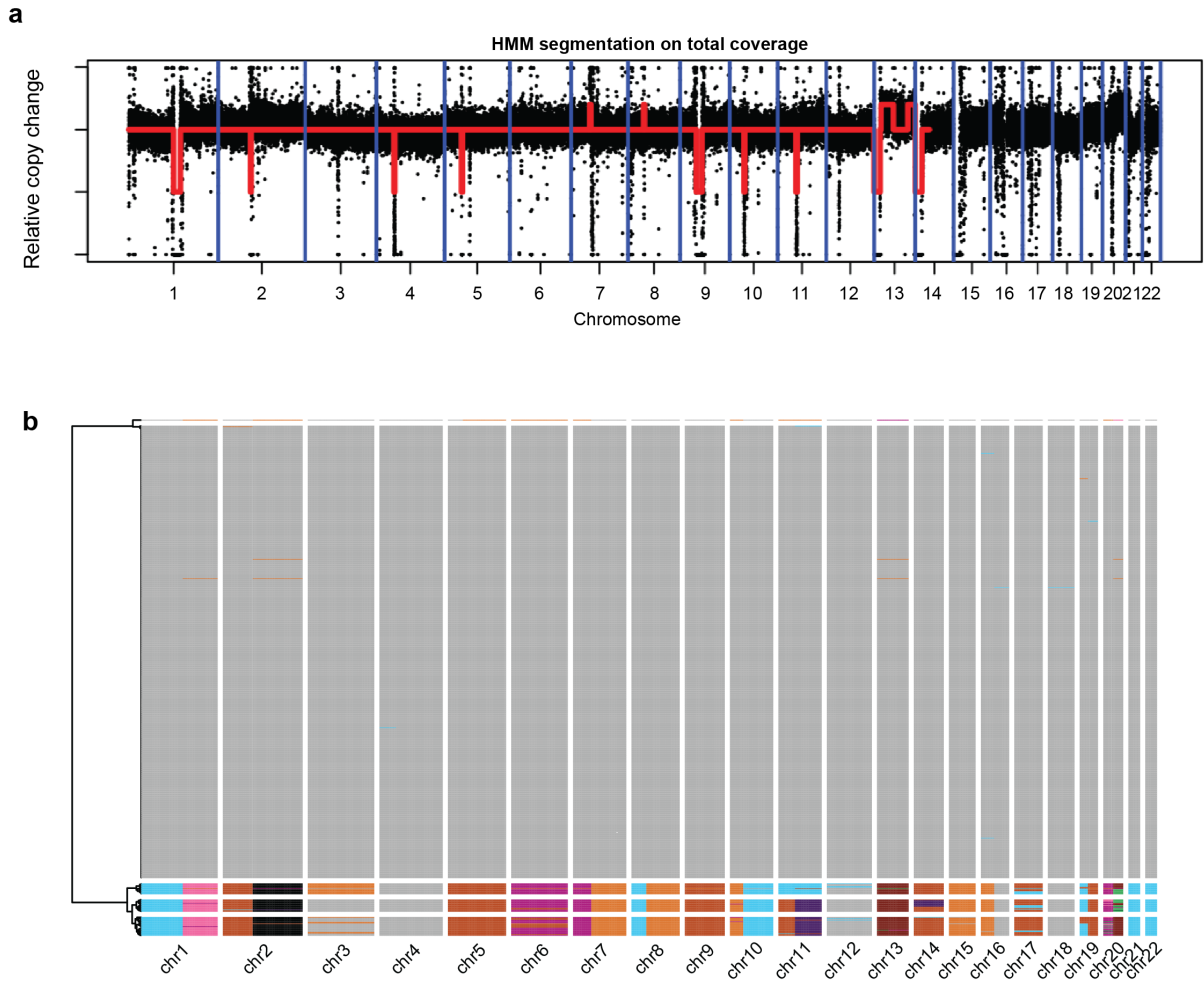
**Supplementary Figure 8.**

a



b



**Supplementary Fig. 8: Segmentation plot and genotype heatmap for the P5847 sample.** The color scheme is the same as that in Fig. 2 and Supplementary Fig. 5.
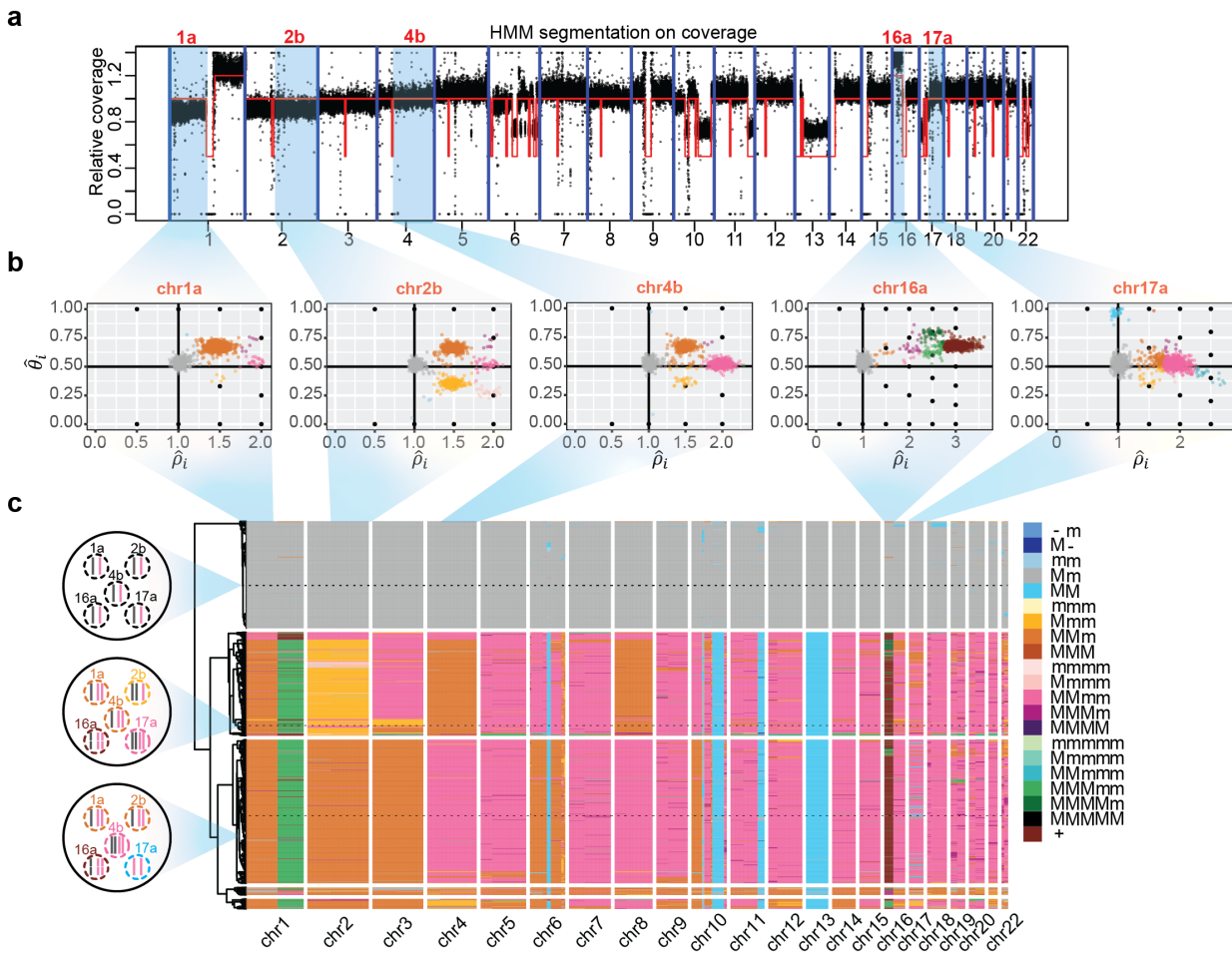
**Supplementary Figure 9.**

a



b

**Supplementary Fig. 9: Segmentation plot and genotype heatmap for the P5915 sample.** The color scheme is the same as that in Fig. 2 and Supplementary Fig. 5.

**Supplementary Figure 10.**



**Supplementary Fig. 10: Segmentation plot and genotype heatmap for the P6461 sample.** The color scheme is the same as that in Fig. 2 and Supplementary Fig. 5.
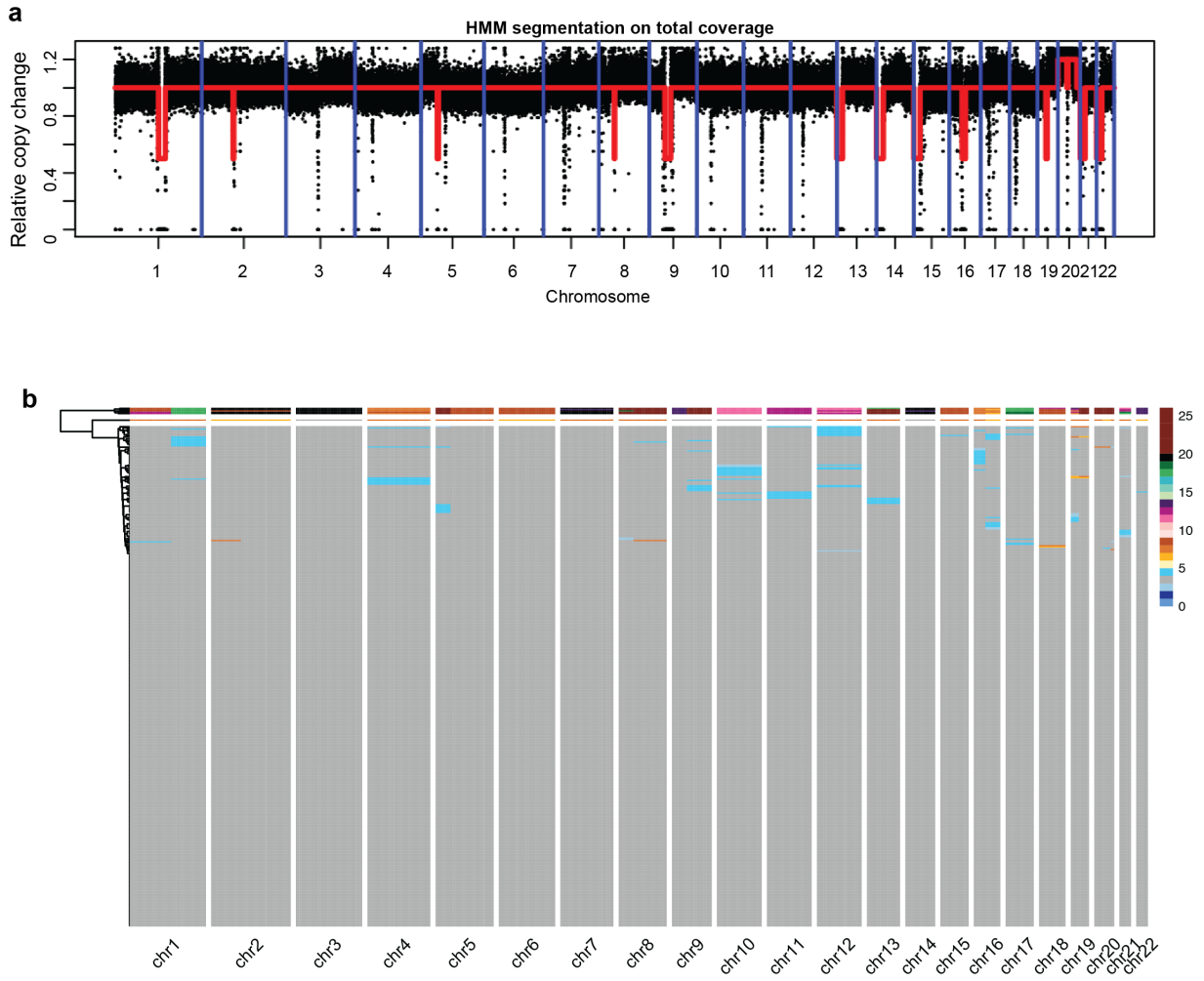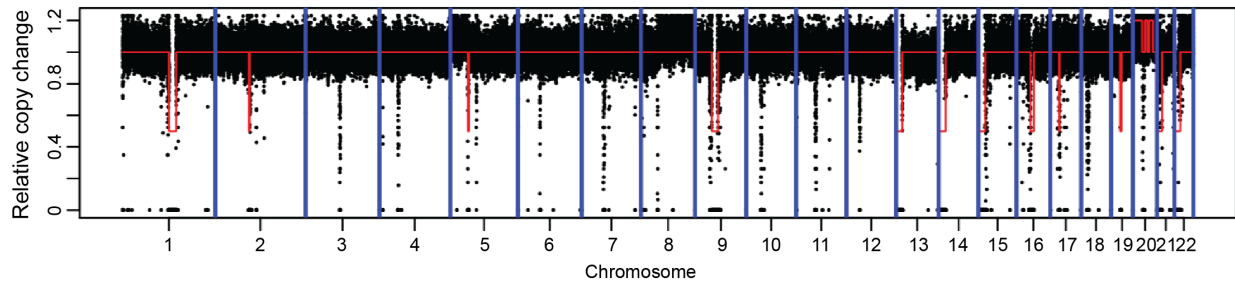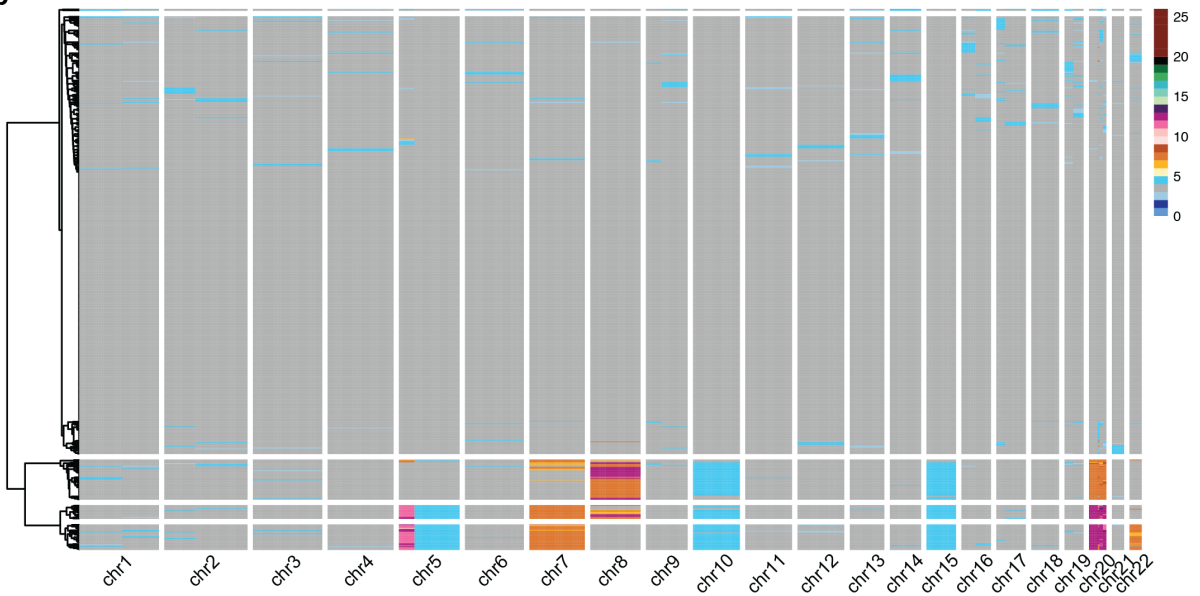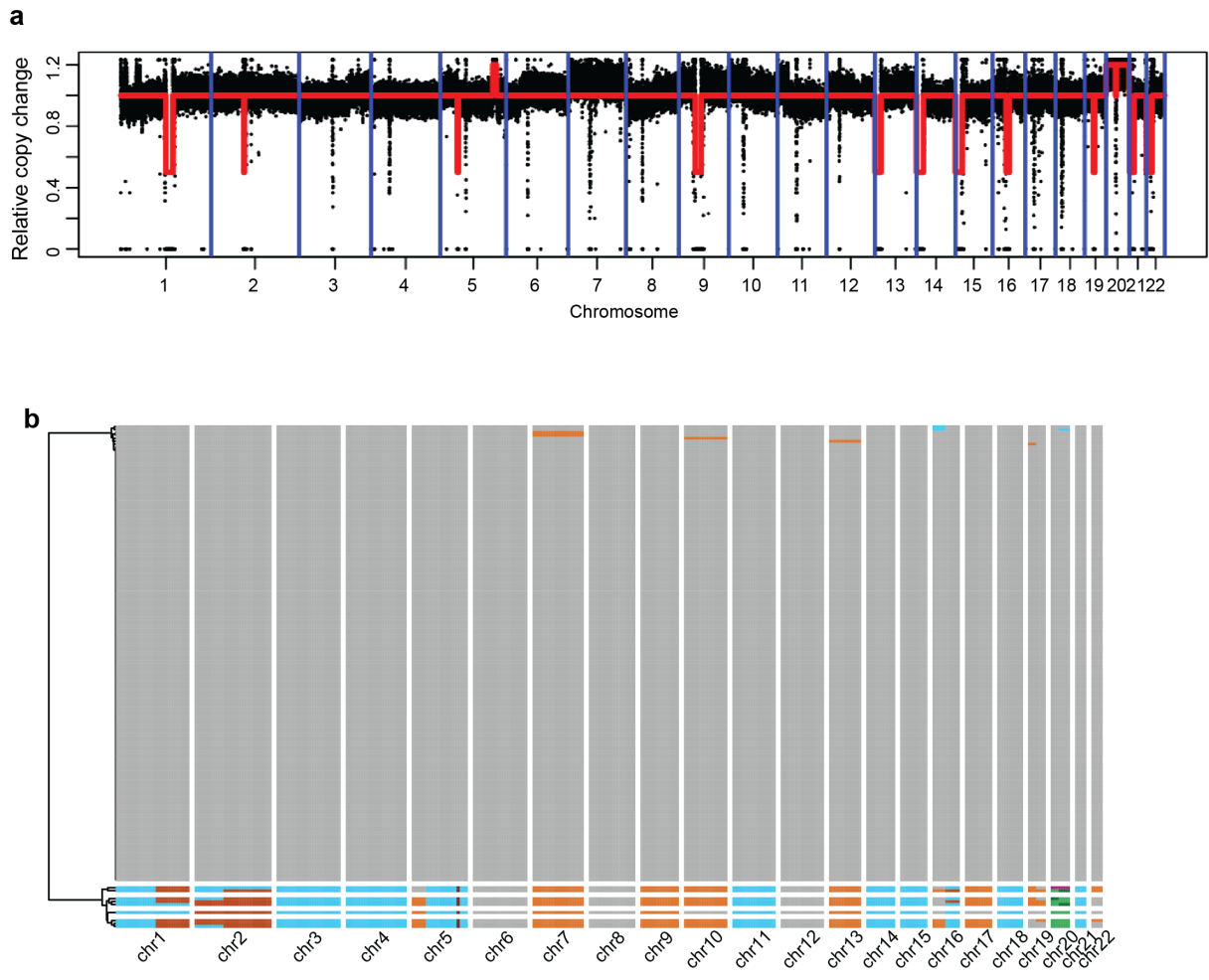
**Supplementary Figure 11.**



**Supplementary Fig. 11: Single cell genotyping of CNV events by Alleloscope for scATAC-seq data of a basal cell carcinoma sample (SU006[1]).** (a) Genotype profiles of six example regions. The regions were taken from the segmentation of whole exome sequencing (WES) data. Each dot represents a cell-specific $(\hat{\rho}_i, \hat{\theta}_i)$ pair. Cells are colored by annotation derived from peak signals[1]. Two tumor cell clusters, identified using ATAC peaks, are labeled by red and blue; fibroblasts (Fibro) are labeled by grey. Density contours of the three cell subpopulations are also shown. (b) Hierarchical clustering of cells in scATAC-seq by $\hat{\theta}_i$ reveals that the two tumor subpopulations are differentiated by peak signals that don't correlate with broad copy number events.

12

**Supplementary Figure 12.**



86

**Supplementary Fig. 12: Confidence scores for the genotype assignment of each cell in each region for the SNU601 scDNA-seq dataset.**

**Supplementary Figure 13.**

**Supplementary Fig. 13: Distribution of the posterior confidence scores of subclone assignment for the 2,753 cells from SNU601 scATAC-seq.**

**Supplementary Figure 14.**



92

**Supplementary Fig. 14: Power for the detection of 1 copy deletion and 1 copy amplification for data of varying coverage (per base), heterozygous SNP count, and number of cells.** The heterozygous SNP count reflects the size of the region: larger regions contain more heterozygous loci. Cells were clustered based on the minimum distance of $\hat{\theta}_i$ to the canonical values. Top: phasing accuracy, defined as the proportion of SNPs with $\hat{I}_j$ correctly estimated; bottom: cell CNV state accuracy, defined as the proportion of cells that are correctly assigned to carrier state. Amp: amplification. Del: deletion. Line types represent different proportions (0.5%, 0.1% and 0.05%) of carrier cells. The number of SNPs, coverage, number of cells and purity were set as 10,000, 0.03, 1000, and 0.5 if not specified.

103  **<u>Supplementary Methods</u>**

104  <u>Simulations and Power Analysis</u>

105  For a simulated region, let n be the number of cells, m be the number of heterozygous

106  SNPs, $\theta$ be the major haplotype proportion, and $\mu_i$ be the total coverage of cell i sampled

107  from the cells on chr7 in the P5931 tumor sample. For cell i, we simulated total coverage

108  of SNP j ($\mu_{ij}$) using a Poisson distribution

109  $$\mu_{ij} \sim Poisson(\mu_i),$$

110  where $i = 1 \sim n$. Parallelly, phases of SNP j ($I_j$) were simulated under a Bernoulli

111  distribution

112  $$I_j \sim Bernoulli(0.5),$$

113  where $I_j$ indicates whether a reference allele is on the major haplotype for SNP j, and

114  j=1~m. Using $\mu_{ij}$ and $I_j$, simulated read counts of reference alleles of SNP j in cell i ($A_{ij}$)

115  were simulated under a Binomial distribution

116  $$A_{ij} \sim Binomial(\mu_{ij}, p_{ij}),$$

117  where $p_{ij}$ is the proportion of the reference allele at loci j in cell i with the values shown in

118  the following table

| $p_{ij}$ | cell i with CNA | cell i without CNA |
|---|---|---|
| $I_j = 1$ | $\theta$ | 0.5 |
| $I_j = 0$ | $1 - \theta$ | 0.5 |

119

120    Then simulated read counts of alternative alleles of SNP j in cell i ($B_{ij}$) were retrieved by

121    $$B_{ij} = \mu_{ij} - A_{ij}$$

122    In the first simulation used to illustrate distribution of the estimates from Alleloscope, we

123    fixed the cell number n to be 1,000, the SNP number m to be 10,000 which are typical in

124    real datasets. $\theta$ was set to be 1 and 0.66 for cells carrying deletion and one-copy

125    amplification respectively with the purity equal to 0.5. On the simulated $A_{ij}$ and $B_{ij}$

126    matrices Alleloscope estimated phases for each SNP and CAN states for each cell.

127    Distribution of the estimated values versus the true values are visualized using boxplots.

128    To know the effects of SNP numbers, cell coverage, cell numbers, and purity, power

129    analysis was performed for one-copy deletion and one-copy amplification scenarios.

130    We assessed the accuracy for phasing and cell-level CNA state estimation under the

131    following scenarios: SNP numbers from 1,000 to 50,000, mean coverage from 0.01 to

132    0,05 for each cell, cell number from 500 to 2500. For different scenarios, we assessed

133    the effect of three purity: 0.5, 0.1, and 0.01, reflecting from larger subclones to rare

134    subclones. All parameters remained the same as those in the previous paragraph except

135    for the parameters that were assessed. Phasing accuracy was calculated by comparing

136    true $I_j$'s and estimated $\hat{I}_j$'s in the region. If $\hat{I}_j \geq 0.5$, the values were considered as 1.

137    Otherwise, the values were considered 0. On the other hand, the accuracy of cell CNA

138    state estimation was the clustering accuracy using the estimated $\hat{\theta}_i$ values. Cells with $\hat{\theta}_i$

139    values smaller than the midpoints between true $\theta$ of normal cells ($\theta_o = 0.5$) and true $\theta$ of

140    carriers ($\theta_{del} = 1$ ; $\theta_{amp} = 0.66$ ) were considered as normal cells; otherwise, cells were

141  considered as carriers. The clustering accuracy was calculated by comparing the clusters

142  to the true cell states.

143  **<u>Reference</u>**

144  1.   Satpathy, A.T. et al. Massively parallel single-cell chromatin landscapes of human

145       immune cell development and intratumoral T cell exhaustion. *Nat Biotechnol* **37**,

146       925-936 (2019).

147