

Cell Systems, Volume 7

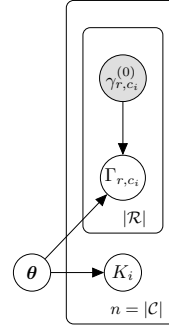
Supplemental Information

Statistical Binning for Barcoded Reads

Improves Downstream Analyses

Atiya Shajii, Ibrahim Numanagić, Christopher Whelan, and Bonnie Berger

Variable	Description
\mathcal{C}	set of all clouds in connected component
\mathcal{R}	set of all reads mapping to some cloud in \mathcal{C}
$\theta = (\theta_{c_1}, \dots, \theta_{c_n})$	vector of cloud weights
K_i	number of reads generated by cloud c_i
Γ_{r,c_i}	event that read r truly originates from cloud c_i
γ_{r,c_i}	$\Pr(\Gamma_{r,c_i} \theta)$
$\gamma_{r,c_i}^{(0)}$	$\Pr(\Gamma_{r,c_i})$ (prior based on edit distance, mate, etc.)



$$\theta = (\theta_{c_1}, \dots, \theta_{c_n}) \sim \text{Dir}(\mathbf{1})$$

$$K_i | \theta \sim \text{Cloud}(\theta_{c_i})$$

$$\Gamma_{r,c_i} \sim \text{Ber}(\gamma_{r,c_i}^{(0)})$$

$$\Gamma_{r,c_i} | \theta \sim \text{Ber}(\gamma_{r,c_i})$$

Figure S1: Graphical representation of EMA’s latent variable model involved in barcoded read alignment, as described in STAR Methods. Related to Figure 1. θ denotes the vector of cloud weights; K_i denotes the number of reads generated by cloud $c_i \in \mathcal{C}$; Γ_{r,c_i} denotes whether read $r \in \mathcal{R}$ maps to cloud c_i , and $\gamma_{r,c_i}^{(0)}$ is a prior on this event based on barcode-oblivious information like edit distance, mate alignment, etc.

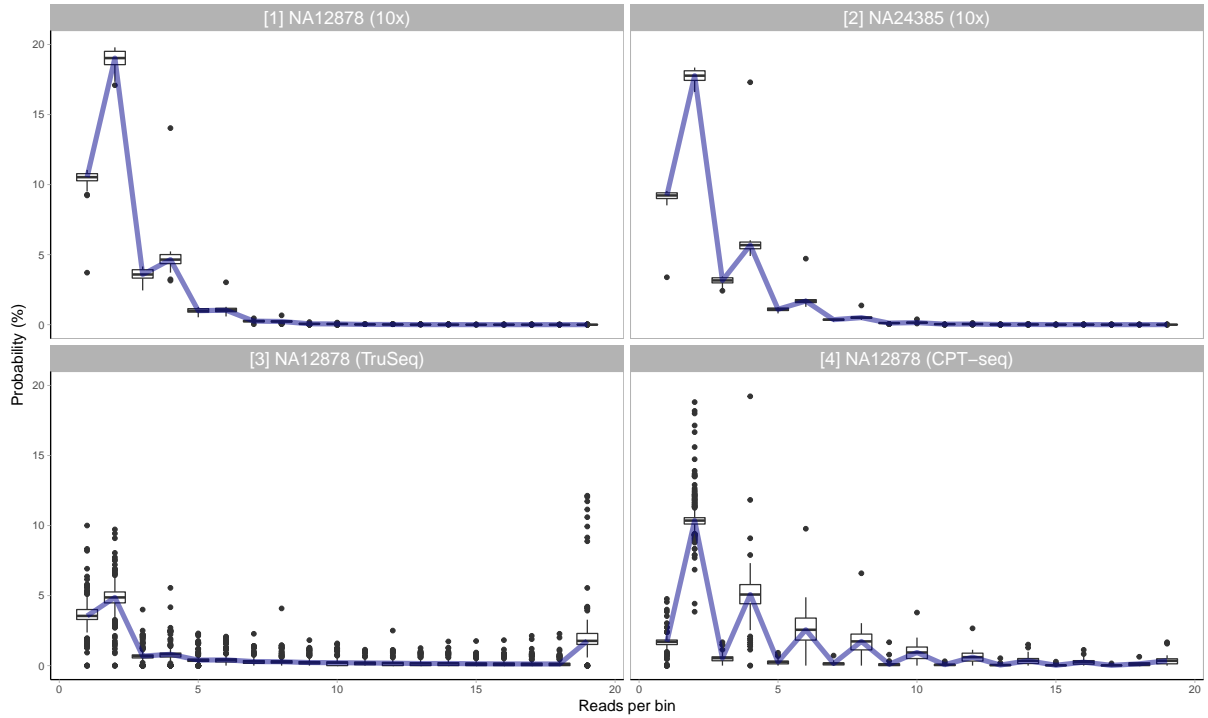


Figure S2: Distribution of the number of reads in a 1kb window within a cloud. Related to Figure 1. The first row shows the distribution of two 10x data samples, while the bottom row shows TruSeq SLR’s and CPT-seq’s distributions. We only consider the clouds in which no reads have multiple alignments within the cloud. The box plots correspond to different bin offsets within the cloud.

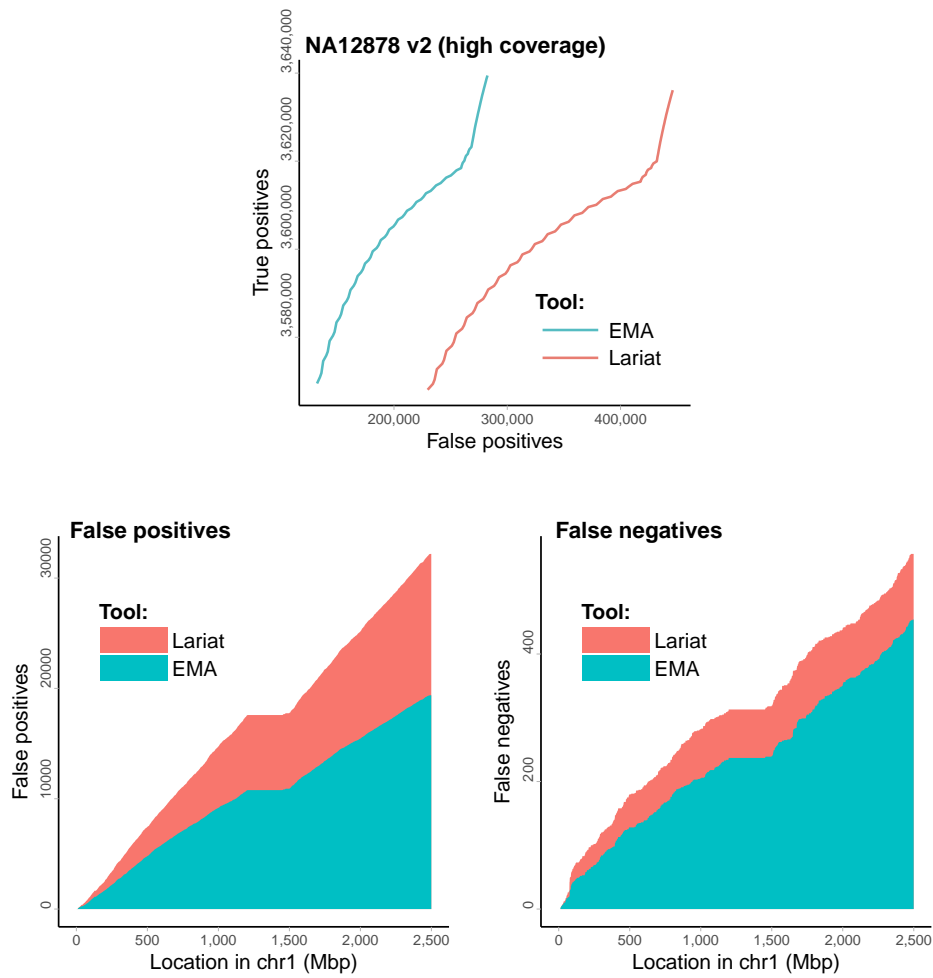


Figure S3: Genotyping accuracy for EMA as compared to Lariat on a high-coverage NA12878 10x dataset. Related to Figure 2. The top plot shows true positives as a function of false positives, and the bottom two plots are cumulative histograms of true and false positives throughout chromosome 1. We note that EMA's improvement is even more substantial with higher coverage.