# BarcodeBender for Mapping Nuclei in Slide-Tags
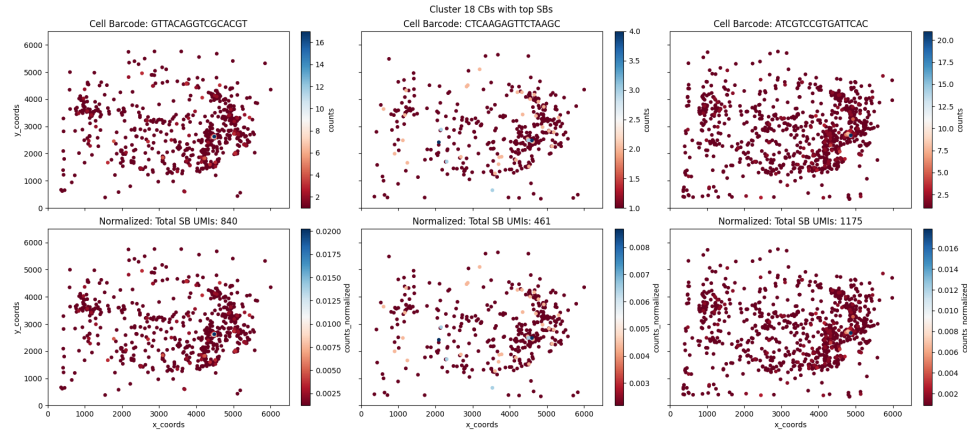
Jake Shapiro

November 2023

## 1 Barcode Bender V0 - Inference on Cluster 18 Nuclei

As a first run, I want to develop a model to infer spatial positions of nuclei that can be confidently mapped, which are of the same cell type, and with as few parameters as possible. For this, I will run inference on droplets from the deeply-sequenced gel-2 Slide-Tags run that were able to be mapped by DBSCAN, which were assigned to the same cell type cluster, and only use the top 10k SBs by total UMIs in those droplets. The DBSCAN mapping of these nuclei is shown below:
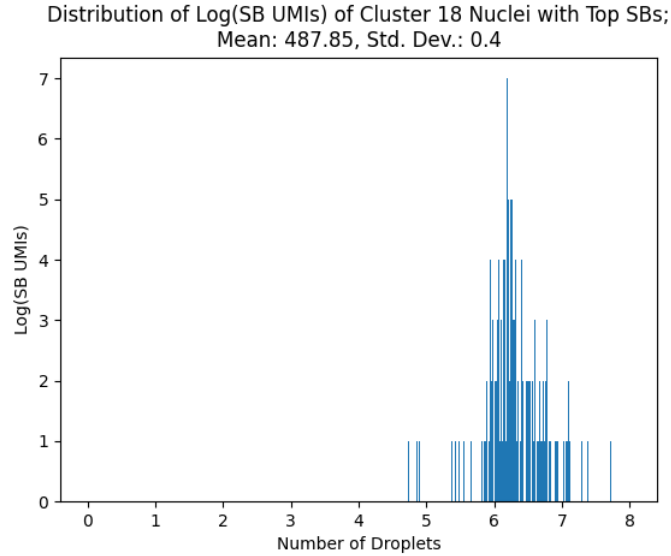


This nuclei-SB set provides several advantages. First, these are all droplets that were called by CellRanger, since the DBSCAN pipeline is run only on CellRanger-called droplets. Second, it is likely that if DBSCAN was able to map these nuclei, then there is a strong spatial signal that BarcodeBender can pick up on and (hopefully) use to improve upon DBSCAN. Third, by having the
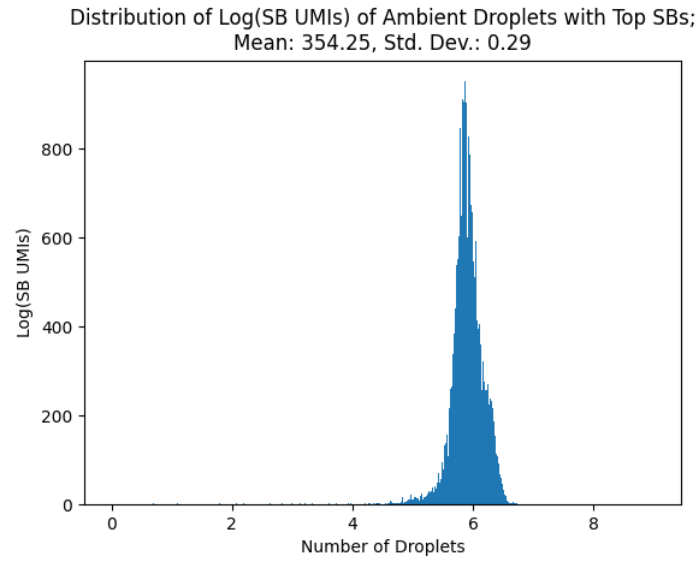
nuclei be of the same cell type, we remove cell type as a confounding variable. There also seems to be a defined boundary for the mapping of these nuclei. Finally, by only retaining the top 10k SBs, we reduce the number of parameters and the complexity of the model. A sampling of these CBs with raw, normalized, and background-removed plots is shown below and there seems to be a decently-strong spatial signal in each (darker blue spots):



Specifically, this dataset contains 261 nuclei and 9,497 SBs. The log-normalized distribution of SB UMIs in these nuclei is shown below:



For ambient droplets, I will use droplets of rank 40,000 - 100,000 by GEX UMIs, whose log-normalized SB UMI distribution is shown below:

Distribution of Log(SB UMIs) of Ambient Droplets with Top SBs;
Mean: 354.25, Std. Dev.: 0.29

These droplets almost certainly do not contain nuclei, based on GEX UMIs, but still have lots of SB UMIs. Thus, they likely faithfully capture the profile of ambient SBs in solution and can be used to estimate a prior for the ambient SB distribution.

The model I will use is as follows:

$$c_{ji}^{\mathrm{obs}} = c_{ji}^{\mathrm{nuc}} + c_{ji}^{\mathrm{noise}}$$

$$c_{ji}^{\mathrm{nuc}} \sim \mathrm{Poisson}[\epsilon_j^{\mathrm{capt}}\epsilon_j^{\mathrm{perm}}k_{ji}]$$

$$\epsilon_j^{\mathrm{capt}} \sim \mathrm{Beta}(\epsilon_\alpha^{\mathrm{capt}}, \epsilon_\beta^{\mathrm{capt}})$$

$$\epsilon_j^{\mathrm{perm}} \sim \mathrm{Beta}(\epsilon_\alpha^{\mathrm{perm}}, \epsilon_\beta^{\mathrm{perm}})$$

$$k_{ji} = \rho_i k(|\mathbf{r}_i - \mathbf{R}_j|; \sigma_i)$$

$$\rho_i \sim \mathrm{LogNormal}(\rho_{loc}, \rho_{scale})$$

$$\sigma_i \sim \mathrm{LogNormal}(\sigma_{loc}, \sigma_{scale})$$

$$c_{ji}^{\mathrm{noise}} \sim \mathrm{Poisson}[\epsilon_j^{capt}d_j^{\mathrm{drop}}\chi_i^{\mathrm{amb}}]$$

$$d_j^{\mathrm{drop}} \sim \mathrm{LogNormal}(d_\mu^{\mathrm{drop}}, d_\sigma^{\mathrm{drop}})$$

where:

- $c_{ji}^{\mathrm{obs/nuc/noise}}$ are the counts of SB $i$ in droplet $j$ due to observation, capture by the nucleus, or from noise, respectively. The observed counts due to nucleus capture and noise are Poisson-distributed because each SB has a small probability of being capture by the nucleus or droplet.

- $\epsilon_j^{\mathrm{capt}} \in [0,1]$ is the mRNA capture efficiency during RT of droplet j. It is distributed according to a Beta distribution with parameters $\epsilon_\alpha^{\mathrm{capt}}, \epsilon_\beta^{\mathrm{capt}}$.

- $\epsilon_j^{\mathrm{perm}} \in [0,1]$ is the permeability of nucleus j. It is distributed according to a Beta distribution with parameters $\epsilon_\alpha^{\mathrm{perm}}, \epsilon_\beta^{\mathrm{perm}}$.

- $k_{ji}$ is the number of SBs from bead $i$ at the location of nucleus $j$. It is sampled from the distribution $\rho_i k(|\mathbf{r}_i - \mathbf{R}_j|; \sigma_i)$ where $\rho_i$ is the total number of SBs shed by bead $i$, $\mathbf{r}_i$ is the location of bead $i$ and is known from *in-situ* sequencing, $\sigma_i$ is the diffusion radius of SB $i$, and $\mathbf{R}_j$ is the location of nucleus $j$.

- $d_j^{\mathrm{drop}} > 0$ is the "size" of droplet $j$ (the larger, the more ambient it scoops). It is distributed according to a log-normal distribution with parameters $d_\mu^{\mathrm{drop}}$ and $d_\sigma^{\mathrm{drop}}$.

- $\chi_i^{\mathrm{ambient}} \in [0,1]$ is the relative abundance of SB $i$ in ambient solution. It will be initialized from the normalized distribution of SBs in CBs of rank

40,000 - 100,000 by total GEX UMIs (almost certainly no nuclei, still in the ambient SB regime).

The learned parameters are:

- $\epsilon_\alpha^{\mathrm{capt}}, \epsilon_\beta^{\mathrm{capt}} > 0$

- $\epsilon_\alpha^{\mathrm{perm}}, \epsilon_\beta^{\mathrm{perm}} > 0$

- $\mathbf{R}_j \in \mathbb{R}^2$

- $\rho_{loc}, \rho_{scale}, \sigma_{loc}, \sigma_{scale} > 0$

- $d_\mu^{\mathrm{drop}}$ and $d_\sigma^{\mathrm{drop}} > 0$

- $\chi^{\mathrm{ambient}}, \chi_i^{\mathrm{ambient}} \in [0,1], \sum_i \chi_i^{\mathrm{ambient}} = 1$

All parameters except $\mathbf{r}_i$ (again, known from *in-situ* sequencing will be learned using variational inference. $\chi^{\mathrm{ambient}}$ may also be treated as fixed because we know what the background looks like from the droplets in the ambient regime. A graphical representation of the generative model is below (hand-drawn, sorry):