

S1 Section Implementation of Molecular cross-validation [4]

The most straightforward method of denoising data is to use a principal component analysis (PCA); a subset of principal components (PCs) which explain a desired amount of variance are identified, and the expression data is projected onto this subset. This removes the effect of low-variance, presumably noise-driven components. The choice of how many PCs to retain greatly influences the structure of the data and the amount of remaining biological variation. DEWÄKSS utilizes PCA only to constructing the kNN-G; after the graph has been constructed denoising proceeds using the full expression matrix.

We explored selecting a reasonable number of PCs using the method Molecular cross-validation (MCV), presented by Batson et al. [4] (for use in projection algorithms). MCV takes in count data and partitions the individual counts randomly into two partitions, $(\mathbf{X}_J, \mathbf{X}_{J^c}) \in \mathbb{R}^{N, M}$. A normalization pipeline then normalizes each data partition individually to yield $(\bar{\mathbf{X}}_J, \bar{\mathbf{X}}_{J^c})$ and carries out PCA on $\bar{\mathbf{X}}_J$.

For different numbers of PCs, we project and then reconstruct $\bar{\mathbf{X}}_J$ and then calculate the Mean Squared Error (MSE) between the reconstruction and $\bar{\mathbf{X}}_{J^c}$. This is done as follows: let $\bar{\mathbf{X}}_J$ be the normalized expression matrix and \mathbf{V} be the orthogonal matrix with columns as principal vectors. Then the mapping $\mathbf{Z} = \bar{\mathbf{X}}_{J^c} \mathbf{V}$ gives the projected features \mathbf{Z} that we can map back to $\bar{\mathbf{X}}_J = \mathbf{Z} \mathbf{V}^T$. Using all columns of \mathbf{V} will map $\bar{\mathbf{X}}_J \rightarrow \bar{\mathbf{X}}_J$; however, using a subset \mathbf{V}_K of the PCs we get

$$\check{\mathbf{X}} = \bar{\mathbf{X}}_J \mathbf{V}_{1:k} \mathbf{V}_{1:k}^T \tag{S19}$$

which we can use to calculate the MSE:

$$\text{MSE} = \left\| \check{\mathbf{X}} - \bar{\mathbf{X}}_{J^c} \right\|. \tag{S20}$$

We can calculate the MSE for each number of PCs $k \in [1, K]$ and find k that corresponds to the minimum MSE:

$$k^* = \arg \min_k \left\| \bar{\mathbf{X}}_J \mathbf{V}_{1:k} \mathbf{V}_{1:k}^T - \bar{\mathbf{X}}_{J^c} \right\| \tag{S21}$$

We applied this algorithm using *TruncatedSVD* [32] to perform PCA and reconstruct the input matrix using the operation in equation 21.