

Quantifying the effect of experimental perturbations in single-cell RNA-sequencing data using graph signal processing

Daniel B. Burkhardt^{1,†}, Jay S. Stanley III^{3,†}, Alexander Tong², Ana Luisa Perdigoto⁴,
Scott A. Gigante^{1,2,3}, Kevan C. Herold⁴, Guy Wolf^{6,7,‡}, Antonio J. Giraldez^{1,‡},
David van Dijk^{1,2,‡}, Smita Krishnaswamy^{1,2,‡*}

¹Department of Genetics; ²Department of Computer Science;

³Computational Biology & Bioinformatics Program;

⁴Department of Immunobiology; ⁵Department of Internal Medicine; Yale University, New Haven, CT, USA

⁶Department of Mathematics and Statistics, Université de Montréal, Montreal, QC, Canada

⁷Mila – Quebec AI Institute, Montreal, QC, Canada

*Corresponding Author E-mail: smitta.krishnaswamy@yale.edu

[†] These authors contributed equally. [‡] These authors contributed equally.

Abstract

Single-cell RNA-sequencing (scRNA-seq) is a powerful tool to quantify transcriptional states in thousands to millions of cells. It is increasingly common for scRNA-seq data to be collected in multiple conditions to measure the effect of an experimental perturbation. However, quantifying differences between scRNA-seq datasets remains an analytical challenge. Previous efforts at quantifying such differences focus on discrete regions of the transcriptional state space such as clusters of cells. Here, we describe a continuous measure of the effect of an experiment across the transcriptomic space with single cell resolution. First, we use the manifold assumption to model the cellular state space as a graph with cells as nodes and edges connecting cells with similar transcriptomic profiles. Next, we calculate an Enhanced Experimental Signal (EES) that estimates the likelihood of observing cells from each condition at every point in the manifold. We show that the EES has useful properties for analysis of single cell perturbation studies. We show that we can use the magnitude and frequency of the EES, using an algorithm we call vertex frequency clustering, to identify specific populations of cells that are or are not affected by an experimental treatment at the appropriate level of granularity. Using these selected populations we can derive gene signatures of affected populations of cells. We demonstrate both algorithms using a combination of biological and synthetic datasets. Implementations are provided in the MELD Python package, which is available at <https://github.com/KrishnaswamyLab/MELD>.

1 Introduction

As single-cell RNA-sequencing (scRNA-seq) has become more accessible, the design of single-cell experiments has become increasingly complex. Researchers regularly use scRNA-seq to quantify the effect of a drug, gene knockout, or other experimental perturbation on a biological system. However, quantifying the

differences between single-cell datasets collected from multiple experimental conditions remains an analytical challenge [1]. This task is hindered by the heterogeneity and noise in both the data and the effects of a given perturbation. More specifically, each single-cell dataset comprises several intrinsic structures of heterogeneous cells, and the effect of the experimental condition could be diffuse across all cells or isolated to particular areas of the cellular state space. Further, technical noise from scRNA-seq measurements, stochastic biological heterogeneity, and uneven exposure to a perturbation can frustrate attempts to understand differences between single-cell datasets.

To address this, we develop a signal over the cellular manifold that quantifies the conditional likelihood that each cell would appear in a given sample condition. Our approach relies on manifold model of single-cell data, which treats the transcriptomic state space as a continuous low-dimensional manifold, or set of manifolds, to characterize cellular heterogeneity and dynamic biological processes [2–8]. In the manifold model, the biologically valid combinations of gene expression are represented as a locally Euclidean topological space, such as a two-dimensional sheet embedded in three dimensions.

Our goal is to quantify the effect of an experimental perturbation on every single cell state observed in matched experimental and control scRNA-seq samples of the same biological system. We explicitly define and calculate an *enhanced experimental signal* (EES), which quantifies the effect of an experimental perturbation across the manifold as a change in the probability of observing each transcriptomic profile in the treatment condition relative to the control. We assume that the cell profiles observed in each experiment are sampled from an underlying multivariate probability density function over the transcriptomic state space that describes the likelihood of observing any cell state in a given condition. For example, it is more likely to observe neuronal cells in a sample of brain tissue than in a peripheral blood sample. Next, we assume that the effect of an experimental perturbation is to change this underlying probability density. For example, if you knock out a gene, some neuronal types or even transcriptional states of the same type may be more or less likely to be observed. The key observation here is that we expect to observe a continuous spectrum of changes in probability across the cellular manifold (**Figure 1**). Because the effect of an experiment is continuous, we seek to estimate this effect across all the observed regions of the manifold, namely at each single-cell profile sampled from each condition.

Although several methods exist for jointly analyzing multiple single-cell datasets [9–11], all previous works quantifying differences between datasets rely on an initial aggregation of the data prior to downstream analysis calculating differential abundance or differential gene expression between conditions. Most published analyses of multiple scRNA-seq samples follow the same basic steps [12–19]. First, datasets are merged applying either batch normalization [18, 19] or a simple concatenation of data matrices [12–17]. Next, clusters are identified by grouping either sets of cells or modules of genes. Finally, within each cluster, the cells from each condition are used to calculate statistical measures, such as fold-change between samples. However, reducing the experimental signal to the level of clusters sacrifices the power of single-cell data. In particular, we demonstrate cases in the following sections where subsets of a cluster are enriched and others subsets are depleted, but in the published analysis these nuances were missed because the analysis focused on fold-change in abundance of each cluster.

Instead of quantifying the effect of a perturbation on clusters, we focus on the level of single cells. First, we use the manifold assumption to create a simplified data model, a cell similarity graph where nodes are cells and edges connect cells with similar transcriptomic profiles [20]. We then apply tools from the emerging field of graph signal processing [21] to compute the EES as the likelihood of observing a given cell in the treatment condition relative to the control. An EES value of 0.5 indicates a cell is equally likely to originate from either condition while values close to 1 or 0 are almost only found in the experimental or control conditions respectively.

In the sections that follow, we show that the EES has useful information for the analysis of experimental conditions in scRNA-seq. First, the EES can be used as a measure of transcriptional response to a perturbation on a cell-by-cell basis to identify the cell states most and least affected by an experimental treatment. Second, we show that the frequency composition of the EES can be used as the basis for a clustering algorithm we call *vertex frequency clustering* (VFC), which identifies populations of cells that are transcriptionally similar and are similarly affected (either enriched, depleted, or unchanged) between conditions. In other words, the identification of the affected population is done using the EES at the level of granularity pertinent to the perturbation response, rather than at a predetermined granularity based on data geometry alone. Third, we obtain gene signatures of a perturbation by performing differential expression between vertex frequency clusters with varying EES distributions. We show that these signatures outperform signatures obtained by direct comparisons of two experimental conditions.

To demonstrate these advantages, we apply this analysis to a variety of biological datasets, including T-cell receptor stimulation [16], CRISPR mutagenesis in the developing zebrafish embryo [18], and a newly published dataset of interferon-gamma stimulation in human pancreatic islets. We also provide a set of quantitative comparisons for both algorithms using ground truth simulated scRNA-seq data. In each case, we demonstrate the ability of the EES to identify trends across experimental conditions and identify instances where use of the EES and vertex frequency clustering improves over published analytic techniques.

Implementations of the EES algorithm and vertex frequency clustering are provided in the Python package MELD, so named for its utility in joint analysis of single-cell datasets. MELD is open-source and available on GitHub at <https://github.com/KrishnaswamyLab/MELD>.

2 Results

2.1 Overview of the EES algorithm

We propose a novel framework for quantifying differences in cell states observed across single-cell experiments. Our work is inspired by recent successes in applying manifold learning to scRNA-seq analysis [20]. The manifold model is a useful approximation for the cellular transcriptomic space because not all combinations of gene expression are biologically valid. Instead, valid cellular states are intrinsically low-dimensional with smooth transitions between similar states. Single-cell data generally consists of several disconnected manifolds that each are locally continuous. The power of scRNA-seq as a measure of an experimental treatment is that it provides observations of cell state at thousands to millions of points along the manifold in each condition. In this context, our goal is to quantify the change in enrichment of cell states along the manifold as a result of the experimental treatment (**Figure 1**).

For an intuitive understanding, we first consider a simple experiment with one treatment condition and one control. We seek to calculate the conditional probability that each cell would be observed in either the experimental or control condition over a manifold approximated from all cells from both conditions. This conditional probability can be used as a measure of the effect of the experimental perturbation because it indicates for each cell how much more likely we are to observe that cell state in the treatment condition relative to the control condition (**Figure 1**). We refer to this ratio as the *Enhanced Experimental Signal* (EES). The steps of the EES algorithm are given in **Algorithm 1**.

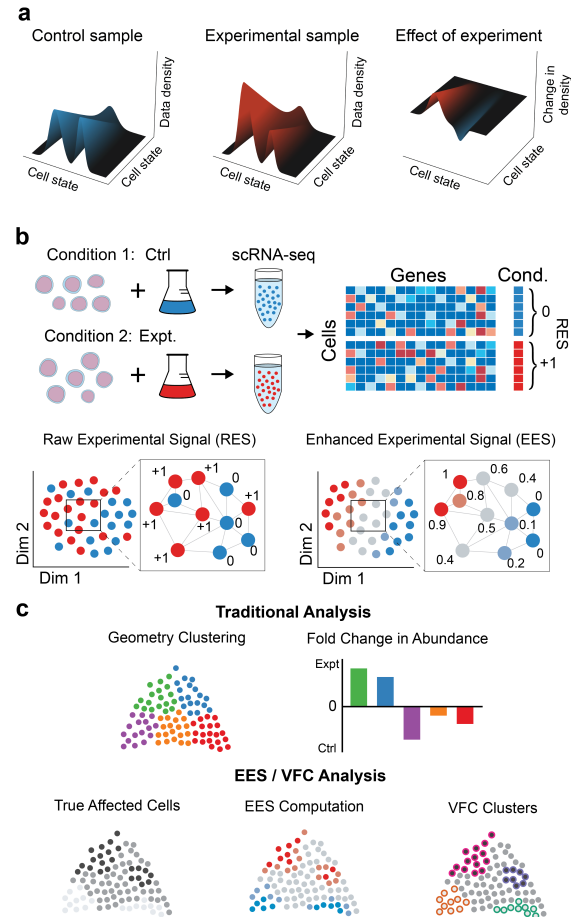


Figure 1: (a) To quantify the effect of an experiment, we model single cell experiments as samples from a probability density function (pdf) over the underlying transcriptomic cell state space manifold. The pdf for the control sample is the frequency with which cell states are observed in the control sample compared to the overall frequency of the cell state in both samples combined. In this context, the effect of an experimental perturbation is to alter this probability density and thus the data density in the experimental sample. Therefore, the effect of an experimental can be quantified as the change in the probability density in the experiment condition relative to the control. (b) The Enhanced Experimental Signal (EES) quantifies this effect by computing a kernel density estimate over the cell similarity graph on the Raw Experimental Signal (RES). The EES indicates the likelihood that a particular cell is from the experimental or control conditions. (c) In traditional analysis, the clusters are based solely on the data geometry and changes in abundance between conditions may not align with the true affected populations. Using the EES and VFC, we can identify the correct cluster resolution for downstream analysis.

As has been done previously, we first approximate the cellular manifold by constructing a simplified data geometry represented by an affinity graph between cells from both conditions [2–8]. In this graph, nodes are cells and the edges between nodes describe the transcriptional similarity between the cells. We then take a new approach to analyze the distribution of cells from each sample over the graph using graph signal processing [21]. A graph signal is any function that has a defined value for each node in a graph. For example, it is natural to represent gene expression values as signals on a graph. Here we use labels indicating the sample origin of each cell as a signal over the graph that we call the *Raw Experimental Signal*. The RES is a collection of one-hot indicator signals, with one signal per condition. Each signal has value 1 associated with each cell from the corresponding condition and value 0 elsewhere. In a simple two-sample experiment, the RES would comprise two one-hot signals, one for the control condition and one for the experimental condition. These one-hot signals are row-wise L1 normalized to normalize different numbers of cells sequenced in each sample.

To derive the *Enhanced Experimental Signal* (EES), we next calculate a kernel density estimate of the RES by applying a *low pass filter*, which can be thought of as averaging values of the RES across the edges on the graph, with higher weighting on edges connecting nearer neighbors. The output of this filter gives the EES (conditional probability), which is smooth on the graph, meaning neighboring cells will have similar probability estimates of being observed in a given condition. A visual representation of each step of the algorithm on real-world data can be found in **Figure S1**. We also describe a full pipeline for analysis of single cell datasets using MELD in **Section 4.3**.

Algorithm 1: The EES algorithm

Input: Dataset $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, $\mathbf{x}_i \in \mathbb{R}^m$; Condition labels \mathbf{y} s.t. y_i indicates the condition in which observation \mathbf{x}_i was sampled.

Output: Enhanced Experimental Signal $\tilde{\mathbf{Y}}_{norm} \in \mathbb{R}^{n \times d}$ where d is the number of unique conditions in \mathbf{y}

1. Build graph $\mathbf{G} = \{V, E\}$ by applying anisotropic or other kernel function on \mathbf{X} ;
 2. Instantiate One-Hot Indicator \mathbf{Y} , also referred to as the **RES**, with one column for each unique condition in \mathbf{y} ;
 3. Column-wise L1-normalize \mathbf{Y} to yield \mathbf{Y}_{norm} ;
 4. Apply EES filter over $(\mathbf{G}, \mathbf{Y}_{norm})$ to calculate $\tilde{\mathbf{Y}}$, the kernel density estimate of the data in each condition;
 5. Row-wise L1 normalize $\tilde{\mathbf{Y}}$ to yield $\tilde{\mathbf{Y}}_{norm}$ also referred to as the **EES**;
-

2.2 Graph construction approximates the underlying data manifold

The first step of the EES algorithm is to create a cell similarity graph in which neighboring cells (i.e., cells with small distances between them) are connected by edges. The goal of graph construction is to approximate the underlying manifold from which the data was sampled [22]. There are many ways to construct such a graph, and in general the algorithm presented here can work over any such construction. The default graph construction implemented in the MELD toolkit quantifies cell similarity (i.e., the edge weights of the graph) using the α -decay kernel proposed in [3], which can be interpreted as a smooth k -Nearest Neighbors (kNN) kernel. However, in cases where batch normalization between replicates is required, we first apply a variant of Mutual Nearest Neighbors (MNN) to merge the datasets [9]. The use of batch correction algorithms with the MELD toolkit is discussed more fully in **Supplementary Note 7.1**. Details of graph construction and implications of various choices are discussed in detail in **Section 4.1.1**.

2.3 Derivation of the EES Filter

Having constructed a graph from the combined datasets, the next step in the EES algorithm is to estimate the density of each sample label over the graph. A popular non-parametric approach to estimating the data density is using a kernel density estimate (KDE), which relies on an affinity kernel function. To estimate the density of labels over a graph, we turn to the heat kernel, which uses diffusion to provide local adaptivity in regions of varying data density [23] such as is observed in single cell data. Here, we extend this kernel as a low pass filter over a graph to estimate the density of sample labels. This is a natural extension of an affinity kernel function because when applied as an integral operator to a signal, it acts as a low-pass filter. For example, filters based on Gaussian kernels are often used to blur or smooth images. To begin, we take the Gaussian KDE, which is well known tool for density estimation in \mathbb{R}^d . The EES generalizes this form to smooth manifolds. The full construction of this generalization is described in detail in **Section 4.1.10**, and a high level overview is provided here.

A kernel density estimator $f(x, t)$ with bandwidth $t > 0$ and kernel function $K(x, y, t)$ is defined as

$$\hat{f}(x, t) = \frac{1}{N} \sum_{i=1}^N K(x, X_i, t), \quad x \in \mathcal{X} \quad (1)$$

where X is the observed data, x is some point in $\mathcal{X} := \mathbb{R}^d$ (i.e., \mathcal{X} is defined as \mathbb{R}^d), and \mathcal{X} is endowed with the Gaussian kernel defined as

$$K(x, y, t) = \frac{1}{(4\pi t)^{d/2}} e^{-\|x-y\|_2^2/4t} \quad (2)$$

Thus, **Equation 2** defines the Gaussian KDE in \mathbb{R}^d . However, this function relies on the Euclidean distance $\|x - y\|_2^2$, which is derived from the kernel space in \mathbb{R}^d . Since manifolds are only locally Euclidean, we cannot apply this KDE directly to a general manifold.

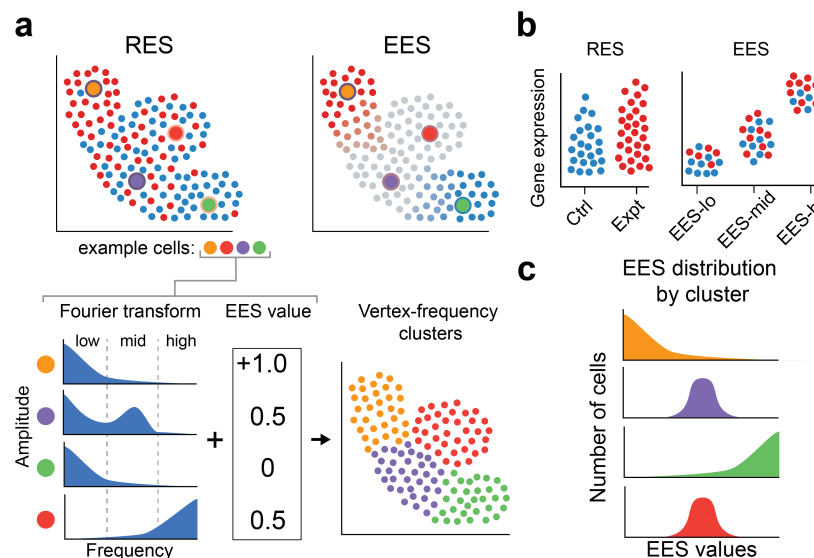


Figure 2: Vertex Frequency Analysis using the EES and RES (a) The Windowed Graph Fourier Transform of the RES and values of EES values at four example points shows distinct patterns between a transitional (blue) and unaffected (red) cell. This information is used in spectral clustering, resulting in Vertex Frequency Clustering. (b) Characterizing Vertex Frequency Clusters with the highest and lowest EES values elucidates gene expression changes associated with experimental perturbations. (c) Examining the distribution of EES scores in vertex-frequency clusters identifies cell populations most affected by a perturbation.

To generalize the Gaussian KDE to a manifold we need to define a kernel space (i.e., the range of a kernel operator) over a manifold. In \mathbb{R}^d the kernel space is often defined via infinite weighted sums of sines and cosines, also known as the Fourier series. However, this basis is not well defined for a Riemannian manifold, so we instead use the eigenbasis of the Laplace operator as our kernel basis. The derivation and implication of this extension is formally explored in **Section 4.1.10**. The key insight is that using this kernel space, the Gaussian KDE can be defined as a filter constructed from the eigenvectors and eigenvalues of the Laplace operator on a manifold. When this manifold is approximated using a graph, we define this KDE as a graph filter over the graph Laplacian given by the following equation:

$$\hat{f}(x, t) = e^{-t\mathcal{L}}x = \Psi h(\Lambda)\Psi^{-1}x \quad (3)$$

where t is the kernel bandwidth, \mathcal{L} is the graph Laplacian, x is the empirical density, Ψ and Λ are the eigenvectors and corresponding eigenvalues of \mathcal{L} , and $e^{-t\mathcal{L}}$ is the matrix exponential. This signal processing formulation can alternatively be formulated in an optimization with Tikhonov Regularization, which seeks to reconstruct the original signal while penalizing differences along edges of the graph. This connection is further explored in **Section 4.1.7**.

To achieve an efficient implementation of the filter in **Equation 3**, the MELD toolkit considers the spectral representation of the RES and uses a Chebyshev polynomial approximation [24] to efficiently compute the EES (see **Section 4.1.4**). The result is a highly scalable implementation. The EES can be calculated on a dataset of 50,000 cells in less than 8 minutes in a free Google Colaboratory notebook¹, with more than 7 minutes of that spent constructing a graph that can be reused for visualization [3] or imputation [4]. With the EES, it is now possible to identify the cells that are most and least affected by an experimental perturbation.

2.4 Vertex-frequency clustering identifies cell populations affected by a perturbation

A common goal for analysis of experimental scRNA-seq data is to identify subpopulations of cells that are responsive to the experimental treatment. Existing methods cluster cells by transcriptome alone and then attempt to quantify the degree to which these clusters are differentially represented in the two conditions. However, this is problematic because the granularity, or sizes, of these clusters may not correspond to the sizes of the cell populations that respond similarly to experimental treatment. Additionally, when partitioning data along a continuum, cluster boundaries are somewhat arbitrary and may not correspond to populations with distinct differences between conditions. Our goal is to identify clusters that are not only transcriptionally similar but also respond similarly to an experimental perturbation.

A naïve approach to identify such clusters would be to simply concatenate the EES to the gene expression data as an additional feature and cluster on these combined features. However, the magnitude of the EES does not give a complete picture of differences in response to a perturbation. For example, there are multiple ways for a cell to have an EES value of 0.5. In one case, it might be that there is a continuum of cells one end of which is enriched for the experimental condition, the other end for the control condition. In this case cells halfway through this continuum, or transitional cells will have an EES of 0.5 (we show an example of this in **Section 2.6**). Another scenario that would result in an EES value of 0.05 is even mixing of a population of cells between control and experimental conditions (with no transition), i.e., cells that are part of a non-responsive cell subtype that is unchanged between conditions (we show an example of this in **Section 2.8**). **Figure S2**. To differentiate between such response regimes we must consider not only the magnitude of the EES but also the frequency of the RES and how fast it changes over the manifold.

¹Freely available at colab.research.google.com, most instances provide a 4-core 2GHz CPU and 20GB of RAM.

Indeed in the transitional case the RES changes slowly or has *low frequency* over the manifold, and in the even-mixture case it changes frequently or has *high frequency* over the manifold.

As no contemporary method is suitable for resolving these cases, we developed an algorithm that integrates gene expression, the magnitude of EES, and the frequency response of the RES over the cellular manifold (**Figure S2d**). In particular, we cluster using local frequency profiles of the RES around each cell. This paradigm is motivated by the utility of analyzing cells based on different classes of heterogeneity. This method, which we call *vertex-frequency clustering* (VFC), is an adaptation of the signal-biased spectral clustering proposed by Shuman et al. [25]. The VFC algorithm provides a new feature basis for clustering based on the spectrogram [25] of the RES, which can be thought of as a histogram of frequency components of a graph signal. We observe that we can distinguish between non-responsive populations of cells with high frequency RES components and transitional populations with lower frequency RES components. The VFC feature basis combines this frequency information with the magnitude of the EES and the cell similarity graph to identify phenotypically similar populations of cells with uniform response to a perturbation. The algorithm is discussed in further detail in **Section 4.2**.

2.5 Quantitative validation of the EES and VFC algorithms

To validate the EES and VFC algorithms, we used a combination of simulated scRNA-seq data and synthetic experiments using previously published datasets. Because no previous benchmarks exist to quantify the ability of an algorithm to capture changes in density between samples, we needed to create a new framework for our comparisons. To create simulated scRNA-seq data, we used Splatter [26]. To ensure the algorithms worked on real scRNA-seq datasets, we also used two previously published datasets comprising Jurkat T cells [16] and cells from whole zebrafish embryos [18]. In each dataset, we created a ground truth probability distribution over all cells that determined the probability each cell would be observed in one of two simulated conditions. In each simulation, different populations of cells of varying sizes were depleted or enriched. Cells were then randomly split into two samples according to this ground truth probability and used as input to each algorithm. More detail on the comparison experiments is provided in **Section 4.7**.

We performed three sets of quantitative comparisons. First, we calculated the degree to which the EES algorithm captured the ground truth conditional probability distribution in each simulation. We found that the Pearson correlation between the EES and the simulated probabilities densities outperformed other graph smoothing algorithms by 10-52% on simulated data and 36-51% on real datasets (**Figure 3, Table 1**). We also determined that the EES is robust to the number of cells captured in the experiment with only a 10% decrease in performance when 65% of the cells in the T cell dataset were removed (**Figure S5**). We used results from these simulations to determine the optimal parameters for the EES algorithm (**Section 7.2**). Next, we quantified the accuracy of the VFC algorithm to identify clusters of cells that were enriched or depleted in each condition. When compared to six common clustering algorithms including Leiden [27] and CellHarmony [11], VFC was the top performing algorithm on every simulation on the T cell data and best performing on average on the zebrafish dataset with a 57% increase in average performance over Louvain, the next best algorithm (**Figures S6a-c & S7, Table 2**). Finally, we calculated how well VFC clusters could be used to calculate the gene signature of a perturbation. Gene signatures obtained using VFC and EES had compared to signatures obtained using direct comparison of two conditions—the current standard—and those obtained using other clustering algorithms (**Figure S6d**). These results confirm that EES and VFC outperform existing methods for analyzing multiple scRNA-seq datasets from different experimental conditions.

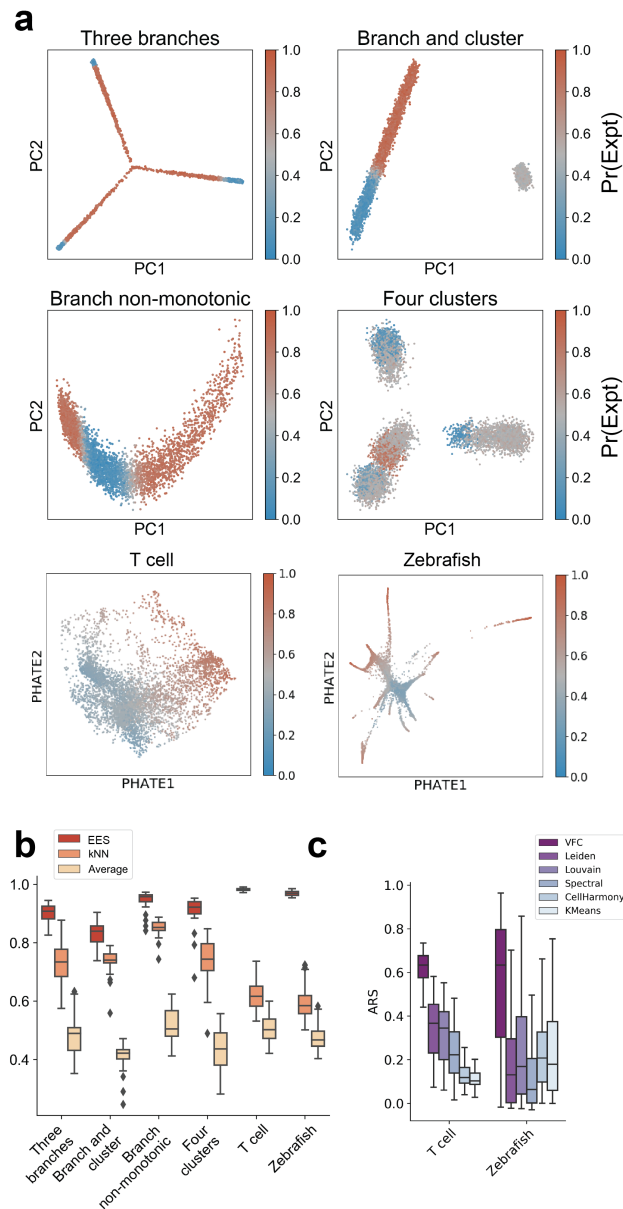


Figure 3: Quantitative comparison of the EES and VFC. **(a)** Single cell datasets were generated using Splatter [26] or taken from previously published experiments [16, 18]. Ground truth EES probabilities were randomly generated 20 times with varying noise and regions of enrichment for the simulated data and 100 random EES were generated for the real-world datasets. Each cell is colored the EES. **(b)** Comparison of the EES algorithm to kNN averaging of the RES and graph Averaging. **(c)** Comparison of VFC to popular clustering algorithms. Adjusted Rand Score (ARS) quantifies how accurately each method detects regions that were enriched, depleted, or unchanged in the experimental condition relative to the control. Higher values are better.

Dataset	EES	Graph Averaging	kNN Averaging
Branch and Cluster	0.82 (0.05)	0.41 (0.05)	0.73 (0.04)
Non-monotonic	0.94 (0.03)	0.52 (0.06)	0.85 (0.03)
Four clusters	0.91 (0.06)	0.44 (0.07)	0.76 (0.07)
Three Branches	0.90 (0.03)	0.48 (0.07)	0.73 (0.07)
T cells [16]	0.98 (0.01)	0.72 (0.06)	0.32 (0.04)
Zebrafish [18]	0.98 (0.01)	0.53 (0.07)	0.80 (0.07)

Table 1: Quantitative comparison of methods for label smoothing over a graph. 40 random seeds were used for each of 4 synthetic datasets. 100 random seeds were used to create random signals on the T cell and zebrafish datasets. Average Pearson Correlation with ground truth signal is displayed with standard deviation in parentheses. Top performing algorithm is bolded.

Dataset	VFC	Spectral	Louvain	Leiden	KMeans	CellHarmony
T cell [16]	0.62 (0.07)	0.23 (0.11)	0.31 (0.13)	0.34 (0.14)	0.11 (0.04)	0.13 (0.05)
Zebrafish [18]	0.53 (0.31)	0.13 (0.15)	0.23 (0.22)	0.19 (0.21)	0.23 (0.20)	0.22 (0.16)

Table 2: Quantitative comparison of clustering methods to identify the cell types affected by a simulated experimental perturbation using real world data.

2.6 The EES identifies a biologically relevant signature of T cell activation

To demonstrate the ability of the EES to identify a biologically relevant EES, we apply the algorithm to Jurkat T cells cultured for 10 days with and without anti-CD3/anti-CD28 antibodies as part of a Cas9 knock-out screen published by Datlinger et al. [16] (**Figure 4a**). The goal of this experiment was to characterize the transcriptional signature of T cell Receptor (TCR) activation and determine the impact of gene knockouts in the TCR pathway. First, we visualized cells using PHATE, a visualization and dimensionality reduction tool for single-cell RNA-seq data (**Figure 4b**) [3]. We observed a large degree of overlap in cell states between the stimulated and control conditions, as noted in the original study [16].

To determine a gene signature of the TCR activation, we considered anti-CD3/anti-CD28 stimulated cells with no CRISPR perturbation. First, we computed EES and VFC clusters on these samples. Then we derived a gene signature by performing differential expression analysis between VFC clusters with the highest and lowest EES values. We identified 2335 genes $q\text{-value} < 0.05$ as measured by a rank sum test with a Benjamini & Hochberg False Discovery Rate correction [28]. We then compared this signature to those obtained using the same methods from our simulation experiments. To determine the biological relevance of these signature genes, we performed gene set enrichment analysis on both gene sets using EnrichR [29]. Considering the GO terms highlighted by Datlinger et al. [16], we found that the MELD gene list has the highest combined score in all of the gene terms we examined (**Figure 4d**). These results show that the EES and VFC are capable of identifying a biologically relevant dimension of T cell activation at the resolution of single cells. Furthermore, the gene signature identified using the MELD toolkit outperformed standard differential expression analyses to identify the signature of a real-world experimental perturbation.

Finally, to quantitatively rank the impact of each Cas9 gene knockout on TCR activation we examined the distribution of EES values for all stimulated cells transfected with gRNAs targeting a given gene (**Figure S8**). We observed a large variation in the impact of each gene knockout consistent with the published results from Datlinger et al. [16]. Encouragingly, our results agree with the bulk RNA-seq validation experiment

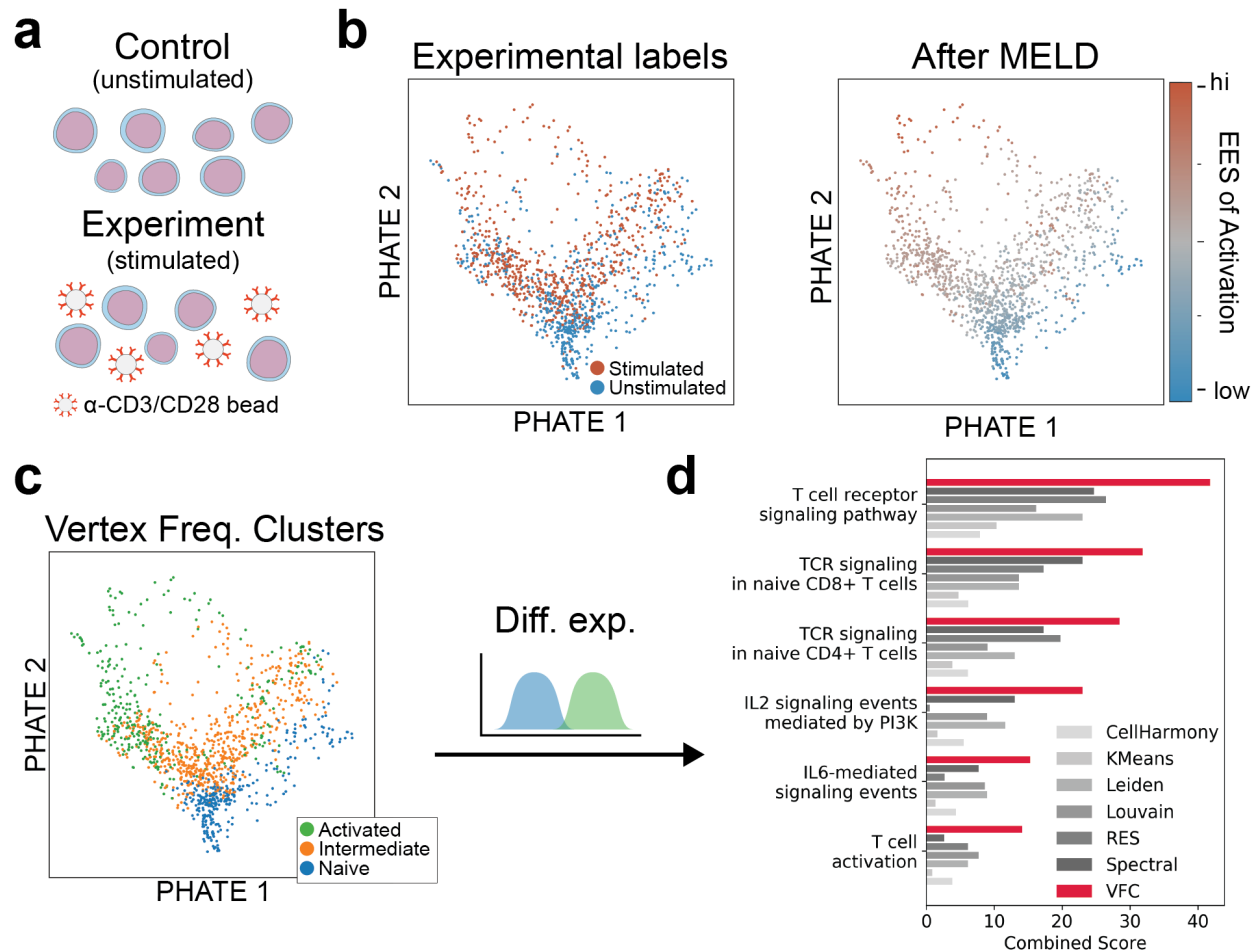


Figure 4: MELD recovers signature of TCR activation. (a) Jurkat T-cells were stimulated with α -CD3/CD28 coated beads for 10 days before collection for scRNA-seq. (b) Examining a PHATE plot, there is a large degree of overlap in cell state between experimental conditions. However, after MELD it is clear which cell states are prototypical of each experimental condition. (c) Vertex Frequency Clustering identifies an activated, a naive, and an intermediate population of cells. (d) Signature genes identified by comparing the activated to naive cells are enriched for annotations related to TCR activation using EnrichR analysis. Combined scores for the MELD gene signature are shown in red and scores for a gene signature obtained using the sample labels only are shown in grey.

of Datlinger et al. [16] showing strongest depletion of TCR response with knockout of kinases LCK and ZAP70 and adaptor protein LAT. We also find a slight increase in EES values (and therefore stimulation) in cells in which negative regulators of TCR activation are knocked out, including PTPN6, PTPN11, and EGR3. Together, these results show that the EES and VFC algorithms are suitable for characterizing a biological process such as TCR activation in the context of a complex Cas9 knockout screen.

2.7 VFC improves characterization of subpopulation response to *chd* loss-of-function

To demonstrate the utility of EES analysis applied to datasets composed of multiple cell types, we applied EES analysis to a recently published chordin loss-of-function experiment in zebrafish using CRISPR/Cas9 (Figure 5) [18]. In this system, loss of chordin function results in a ventralization phenotype characterized

by expansion of the ventral mesodermal tissues at the expense of the dorsally-derived neural tissues [30–32]. In the experiment published by Wagner et al. [18], zebrafish embryos were injected at the 1-cell stage with Cas9 and gRNAs targeting either chordin (*chd*), a BMP-antagonist required for developmental patterning, or tyrosinase (*tyr*), a control gene not expected to affect cell composition at these stages. Embryos were collected for scRNA-seq at 14–16 hours post-fertilization (hpf). We expect incomplete penetrance of the perturbation in this dataset because of the mosaic nature of Cas9 mutagenesis [33].

First, we calculate the EES between the chordin and tyrosinase conditions. Because the experiment was performed in triplicate with three paired *chd* and *tyr* samples, we first ran the EES algorithm for each replicate and then averaged the replicate-specific conditional likelihoods to calculate an EES per dataset (Figure S9). To characterize the effect of mutagenesis on various cell populations, we first examined the distribution of EES values across the 28 cell state clusters generated by Wagner et al. [18] (Figure 5b). We find that overall the most enriched clusters contain mesodermal cells and the most depleted clusters contain dorsally-derived neural cells matching the ventralization phenotype previously reported with *chd* loss-of-function [30–32]. However, we observe that several clusters have a wide range of EES values suggesting that there are cells in these clusters with differing responses to *chd*. Using VFC analysis we find that several of these clusters contain biologically distinct subpopulations of cells with divergent responses to *chd* knock out. Next, we examine three of these cases in depth and reveal previously unreported effects of *chd* loss-of-function within this dataset.

An advantage of using the EES and VFC is the ability to characterize the response to the perturbation at the proper resolution (Figure 2c). We infer that the resolution of the published clusters is too coarse because the distribution of EES values is very large for several of the clusters. For example the EES values within the Tailbud – Presomitic Mesoderm (TPM) range from 0.29–0.94 indicating some cells are strongly enriched while others are depleted. To disentangle these effects, we performed VFC subclustering for all clusters using the strategy proposed in Section 4.3. We found 12 of the 28 published clusters warranted further subclustering with VFC resulting in a total of 50 final cluster labels. To determine the biological relevance of the VFC clusters, we manually annotated each of the three largest clusters subdivided by VFC.

The Tailbud – Presomitic Mesoderm (TPM) cluster exhibits the largest range of EES values of all the clusters annotated by Wagner et al. [18]. In a PHATE visualization of the cluster, we observe many different branches of cell states, each with varying ranges of EES values (Figure 5c). Within the TPM cluster, we find four subclusters using VFC (Figure 5d). Using established markers [19], we identify these clusters as immature adaxial cells, mature adaxial cells, presomitic mesoderm cells, and hematopoietic cells (Figures 5c & S10). Examining the distribution of EES scores within each cell type, we conclude that the large range of EES values within the TPM cluster is due to largely non-overlapping distributions of scores within each of these subpopulations (Figure 5e). The immature and mature adaxial cells, which are embryonic muscle precursors, have low EES values indicating depletion of these cells in the *chd* condition which matches observed depletion of myotomal cells in chordin mutants [30]. Conversely, the presomitic mesoderm and hematopoietic mesoderm have high EES values, indicating that these cells are enriched in a chordin mutant. Indeed, expansion of the hematopoietic mesoderm has been observed in chordin morphants [34] and expansion of the presomitic mesoderm was observed in siblings of the *chd* embryos by Wagner et al. [18]. This heterogeneous effect was entirely missed by the fold-change analysis, since the averaging of all cells assigned to the TPM cluster caused the depletion of adaxial cells to be masked by the expansion of the presomitic and hematopoietic mesoderm.

Another advantage of vertex-frequency clustering is that we can now differentiate between a change in gene expression levels across conditions and a change in abundance of cells expressing a given gene between conditions. When we examined marker gene expression within each of the VFC subclusters, we

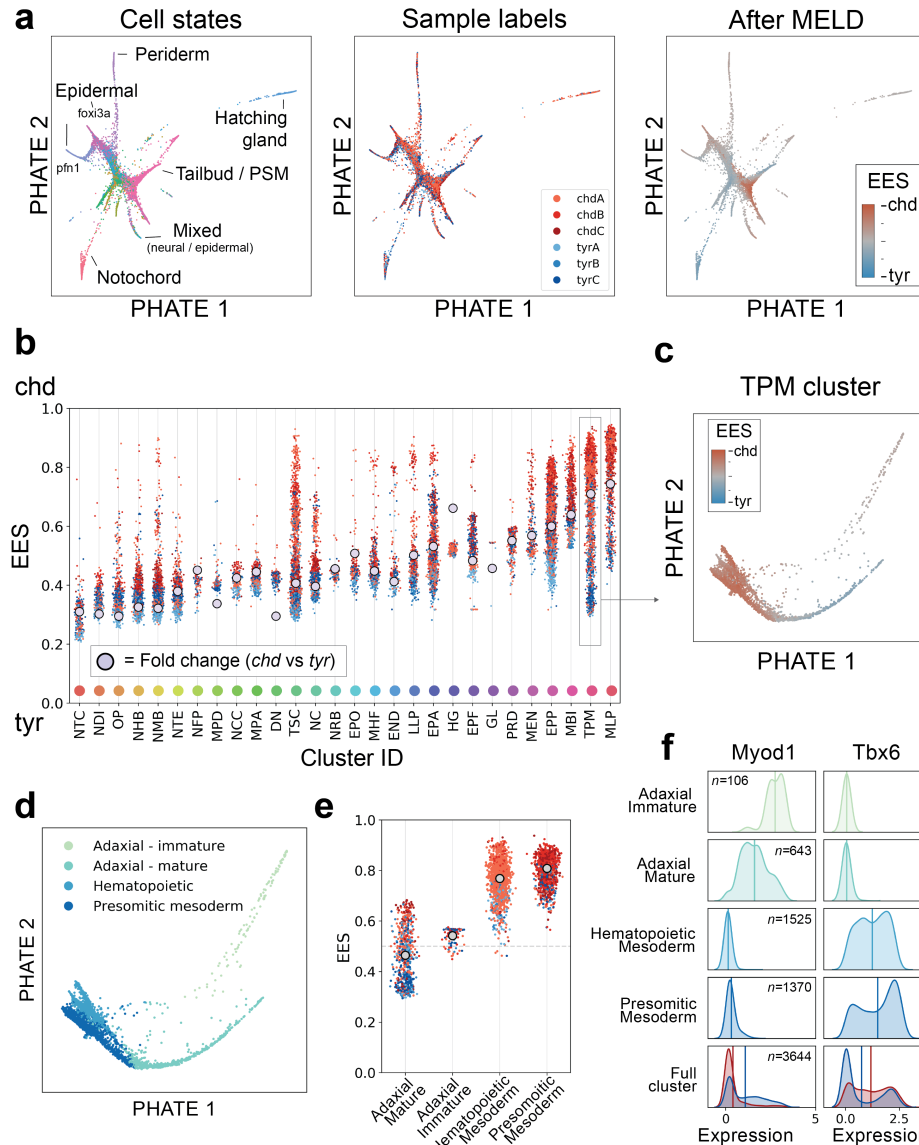


Figure 5: Characterizing chordin Cas9 mutagenesis with MELD. (a) PHATE shows a high degree of overlap of sample labels across cell types. Applying MELD to the mutagenesis vector reveals regions of cell states enriched in the *chd* or *tyr* conditions. (b) Using published cluster assignments², we show that the EES quantifies the effect of the experimental perturbation on each cell, providing more information than calculating fold-change in the number of cells between conditions in each cluster (grey dot), as was done in the published analysis. Color of each point corresponds to the sample labels in panel (a). Generally, average EES value aligns with the fold-change metric. However, we can identify clusters, such as the TPM or TSC, with large ranges of EES values indicating non-uniform response to the perturbation. (c) Visualizing the TPM cluster using PHATE, we observe several cell states with mostly non-overlapping EES values. (d) Vertex Frequency Clustering identifies four cell types in the TPM. (e) We see the range of EES values in the TPM cluster is due to subpopulations with divergent responses to the *chd* perturbation. (f) We observe that changes in gene expression between conditions is driven mostly by changes in abundance of subpopulations with the TPM cluster.

find different trends in expression in each cluster (**Figure 5f**). For example, *Myod1*, a marker of adaxial cells, is lowly expressed in the presomitic and hematopoietic mesoderm, but highly expressed in adaxial cells. Using a rank sum test, we find that *Myod1* is not differentially expressed between conditions within any of the VFC clusters despite there being differential expression using all cells in the TPM cluster (**Figure 5f**). We find a similar trend with *Tbx6*, a mesoderm marker that is not expressed in adaxial cells. We find *Tbx6* is differentially expressed between *chd* and *tyr* embryos within the whole cluster but not within the adaxial or presomitic mesoderm clusters. These results show that the observed change in expression of these genes in the published analysis was in fact due to changes in abundance of cell subpopulations that led to misleading differences in statistics calculated across multiple populations as a whole. Using the EES and VFC, we can identify more appropriate clusters.

We similarly analyzed the "Epidermal - pfn1 (EPP)" and "Tailbud - Spinal Cord (TSC)" clusters which had the 6rd and 3th largest standard deviation in EES values of all published clusters, respectively (**Figure S10**). We used VFC to break up the Epidermal - pfn1 cluster into two subclusters. Among the top differentially expressed genes between the resulting clusters we find *tbx2b*, *crabp2a*, and *pfn1*. *Crabp2a*, a marker of the neural plate border [19], is more lowly expressed in the cluster with higher EES values, suggesting that *chd* loss-of-function inhibits expression of *crabp2a*. This is consistent with previous studies showing a requirement of chordin for proper gene expression patterning within the neural plate [35, 36].

Within the Tailbud - Spinal Cord cluster we further identified three subpopulations of cells using VFC. Examining gene expression within the subclusters, we can see that the published cluster contains different populations of cells. One group expresses markers of the spinal cord (*neurog*, *elavl3*) and dorsal tissues (*olig3*, *pax6a/b*) with an average EES of 0.38, which is consistent with prior evidence that *chd* loss-of-function disrupts specification of the neuroectoderm and dorsal tissues such as the spinal cord [30]. Examining the two remaining subclusters, we see that these cells resemble cells found in both the TPM and Epidermal - Pfn1 clusters. One cluster exhibits high levels of *crabp2a* and EES values <0.5 similar to the neural plate border cells subpopulation within the Epidermal - Pfn1 cluster. Similarly, we find the remaining cluster expressed markers of the tailbud and presomitic mesoderm including *tbx6*, *sox2*, and *fgf8a*. Together, these results demonstrate the advantage of using the EES and vertex frequency clustering to quantify the effect of genetic loss-of-function perturbations in a complex system with many cell types.

2.8 Identifying the effect of IFN γ stimulation on pancreatic islet cells

Next to determine the ability of the EES and VFC algorithms to uncover new biology, we characterized a newly generated dataset of human pancreatic islet cells cultured for 24 hours with and without interferon-gamma (IFN γ), a system with significant clinical relevance to auto-immune diseases of the pancreas such as Type I Diabetes mellitus (T1D). The pathogenesis of T1D is generally understood to be caused by T cell mediated destruction of beta cells in the pancreatic islets [37] and previous reports suggest that islet-infiltrating T cells secrete IFN γ during the onset of T1D[38]. It has also been described that IFN γ -expressing T cells mediate rejection of pancreatic islet allografts [39]. Previous studies have characterized the effect of these cytokines on pancreatic beta cells using bulk RNA-sequencing[40], but no studies have addressed this system at single-cell resolution.

To better understand the effect of immune cytokines on islet cells, we cultured islet cells from three

²Abbreviations: MLP: Lateral plate, TPM: Tailbud - Presomitic mesoderm, HG: Hatching gland, MBI: Blood island, EPP: Epidermal - pfn1, MEN: Endothelial, PRD: Periderm, EPA: Epidermal anterior, EPO: Otic placode, LLP: Lateral line, EPF: Epidermal - foxi3a, GL: Germline, NRB: Rohon beard, NFP: Floorplate, MHF: Heart field, MPA: Pharyngeal arch, NCC: Neural crest - crestin, END: Endoderm, TSC: Tailbud - spinal cord, NC: Neural crest, NTE: Telencephalon, MPD: Pronephric duct, NHB: Hindbrain, NMB: Midbrain, NTC: Notochord, NDI: Diencephalon, DN: Neurons, OP: Optic

donors for 24 hours with and without IFN γ and collected cells for scRNA-seq. After filtering, we obtained 5,708 cells for further analysis. Examining the expression of marker genes for major cell types of the pancreas, we observed a noticeable batch effect associated with the donor ID, driven by the maximum expression of glucagon, insulin, and somatostatin in alpha, beta, and delta cells respectively (**Figure S11a**). To correct for this difference while preserving the relevant differences between donors, we applied the MNN kernel correction described in Section 4.1.1 to merge samples from each donor. Examining PHATE plots after batch correction, we observe three distinct populations of cells corresponding to alpha, beta, and delta cells (**Figure 6a**).

To quantify the effect of IFN γ treatment across these cell types, we calculated the EES of IFN γ stimulation using the same strategy to handle matched replicates as was done for the zebrafish data (**Figure 6a**). We then used established marker genes of islet cells [41] to identify three major populations of cells corresponding to alpha, beta, and delta cells (**Figures 6a-b & S11b**). We next applied vertex frequency clustering to each of the three endocrine cell types and identified a total of nine clusters. Interestingly, we found two clusters of beta cells with intermediate EES values. These clusters are cleanly separated on the PHATE plot of all islet cells (**Figure 6a**) and together the beta cells represent largest range of EES scores in the dataset. To further inspect these clusters, we consider a separate PHATE plot of the cells in the four beta cell clusters (**Figure 6e**). Examining the distribution of RES values in these intermediate cell types, we find that one cluster, which we label as *Non-responsive*, exhibits high frequency RES values indicative of a population of cells that does not respond to an experimental treatment. The *Responsive - Mid* cluster matches our characterization of a transitional population with a structured distribution of RES values. Supporting this characterization, we find a lack of upregulation in IFN γ -regulated genes such as STAT1 in the non-responsive cluster, similar to the cluster of beta cells with the lowest EES values (**Figure 6f**).

In order to understand the difference between the non-responsive beta cells and the responsive populations, we calculated differential expression of genes in the non-responsive clusters and all others as previously described [4]. The gene with the greatest difference in expression was insulin, the major hormone produced by beta cells, which is approximately 2.5-fold increased in the non-responsive cells (**Figure 6f**). This cluster of cells bears resemblance to a recently described “extreme” population of beta cells that exhibit elevated insulin mRNA levels and are found to be more abundant in diabetic mice [42, 43]. That these cells appear non-responsive to IFN γ stimulation and exhibit extreme expression of insulin suggests that the presence of extreme high insulin in a beta cell prior to IFN γ exposure may inhibit the IFN γ response pathway through an unknown mechanism.

We next characterized the gene expression signature of IFN γ treatment across all three endocrine cell types (**Figure 6c-d**). Using a rank sum test to identify genes that change the most between the clusters with highest and lowest EES values within each endocrine population, we identify 911 genes differentially expressed in all three cell types. This consensus signature includes activation of genes in the JAK-STAT pathway including STAT1 and IRF1 [44] and in the IFN-mediated antiviral response including MX1, OAS3, ISG20, and RSAD2 [45–47]. The activation of both of these pathways has been previously reported in beta cells in response to IFN γ [48, 49]. Furthermore, we observe a high degree of overlap in the IFN γ response between alpha and beta cells, but less so between delta cells and either alpha or beta cells. Examining the significantly differentially expressed genes, we find 2394 shared genes in the signatures of alpha and beta cells, but only 911 shared by alpha, beta, and delta cells. To confirm the validity of our gene signatures, we use EnrichR [29] to perform gene set enrichment analysis on the signature genes and find strong enrichment for terms associated with interferon signalling pathways (**Figure S11d**). From these results we conclude that although IFN γ leads to upregulation of the canonical signalling pathways in all three cell types, the response to stimulation in delta cells is subtly different to that of alpha or beta cells.

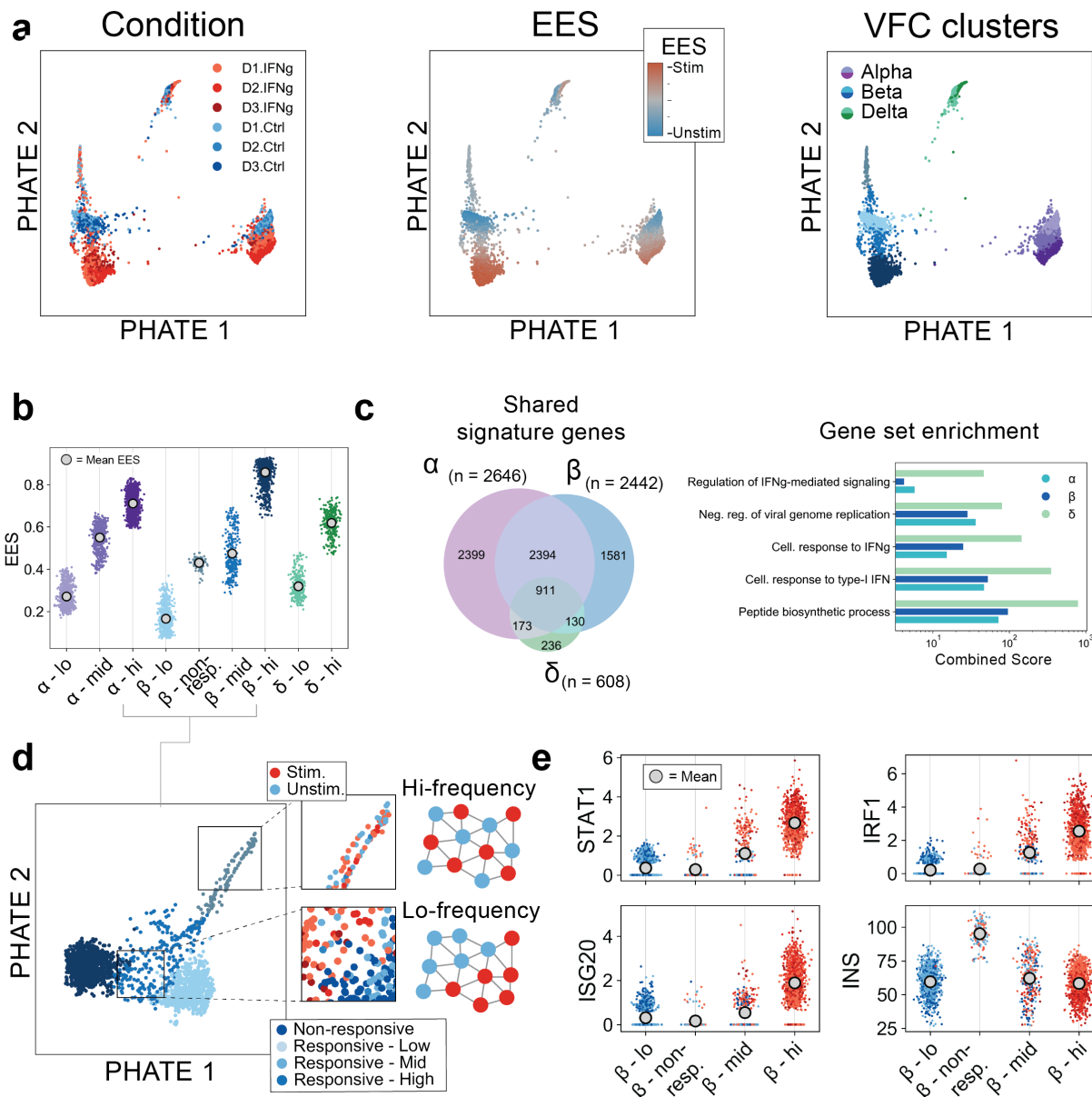


Figure 6: MELD characterizes the response to IFN γ in pancreatic islet cells. **(a)** PHATE visualization of pancreatic islet cells cultured for 24 hours with or without IFN γ . Vertex-frequency clustering identifies nine clusters corresponding to alpha, beta, and delta cells. **(b)** Examining the EES in each cluster, we observe that beta cells have a wider range of responses than alpha or delta cells. **(c)** We identify the signature of IFN γ stimulation by calculating differential expression between the VFC clusters with the highest and lowest EES values for each cell type. We find a high degree of overlap of the significantly differentially expressed genes between alpha and beta cells. **(d)** Results of gene set enrichment analysis for signature genes in each cell type. Beta cells have the strongest enrichment for IFN response pathway genes. **(e)** Examining the four beta cell clusters more closely, we observe two populations with intermediate EES values. These populations are differentiated by the structure of the RES in each cluster (outset). In the non-responsive cluster, the RES has very high frequency unlike the low frequency pattern in the transitional Responsive - mid cluster. **(f)** We find that the non-responsive cluster has low expression of IFN γ -regulated genes such as STAT1 despite containing roughly equal numbers of unstimulated and stimulated cells. This cluster is marked by approximately 40% higher expression of insulin.

Here, we applied EES analysis to a new dataset to identify the signature of IFN γ stimulation across alpha, beta, and delta cells. Furthermore, we used vertex frequency clustering to identify a population of beta cells with high insulin expression that appears unaffected by IFN γ stimulation. Together, these results demonstrate the utility of EES analysis to reveal novel biological insights in a clinically-relevant biological experiment.

2.9 Analysis of donor-specific composition

Although most of the analysis in this manuscript focuses on the two-sample condition, we show that it is possible to use the EES to quantify the differences between more than two conditions. In the islet dataset, we have samples of treatment and control scRNA-seq data from three different donors. To quantify the differences in cell profiles between samples, we first create a one-hot vector for each donor label and normalize across all three smoothed vectors. This produces a measure of how likely each transcriptional profile is to be observed in donor 1, 2, or 3. We then analyze each of these signals for each cluster examined in **Section 2.8 (Figure S12)**. We find that all of the alpha cell and delta cell clusters are depleted in donor 3 and the non-responsive beta cell cluster is enriched primarily in donor 1. Furthermore, the most highly activated alpha cell cluster is enriched in donor 2. As with the EES derived for the IFN γ response, it is also possible to identify donor-specific changes in gene expression, or clusters of cells differentially abundant between each donor. We propose that this strategy could be used to extend MELD analysis to experiments with multiple categorical experimental conditions, such as data collected from different tissues or stimulus conditions.

3 Discussion

When performing multiple scRNA-seq experiments in various experimental and control conditions, researchers often seek to characterize the cell types or sets of genes that change from one condition to another. However, quantifying these differences is challenging due to the subtlety of most biological effects relative to the biological and technical noise inherent to single-cell data. To overcome this hurdle, we designed the EES algorithm and vertex frequency clustering to quantify compositional differences between samples. The key innovation in the EES algorithm is quantifying the effect of a perturbation at the resolution of single cells using theory from manifold learning.

We have shown that our analysis framework improves over the current best-practice of clustering cells based on gene expression and calculating differential abundance and differential expression within clusters. Clustering prior to quantifying compositional differences can fail to identify the divergent responses of subpopulations of cells within a cluster. To identify clusters of cells with cohesive responses to a perturbation, we introduce a novel clustering algorithm, called Vertex-Frequency Clustering. Using the RES and EES, we derive clusters of cells as the correct cluster size to identify cells that are most enriched in either condition, cells transitioning between these states, and cells that are unaffected by an experimental perturbation. We show that gene signatures extracted using these clusters outperform those derived from direct comparison of two samples.

We demonstrated the application of EES and vertex frequency clustering analysis on single-cell datasets from three different biological systems and experimental designs. We provided a framework for handling paired experimental and control replicates and guidance on analysis of complex experimental designs with more than two conditions and in the context of a single-cell Cas9 knockout screen. In our analysis of the zebrafish dataset, we showed greatly improved resolution in our analysis facilitated by the use of the EES and VFC algorithms. In three cases, we show that the published clusters contained biologically relevant

subpopulations of cells with divergent responses the the experimental perturbation. We also described a previously unpublished dataset of pancreatic islet cells stimulated with IFN- γ and characterize a previously unreported subpopulation of β cells that appeared unresponsive to stimulation. We related this to emerging research describing a β cells subtype marked by high insulin mRNA expression and unique biological responses.

We anticipate MELD to have widespread use in many contexts since experimental labels can arise in many contexts. As we showed, if we have sets of single cell data from healthy individuals vs sick individuals, the EES could indicate cell types specific to disease. This framework could potentially be extended to patient level measurements where patients phenotypes as measured with clinical variables and laboratory values can be associated with enriched states in disease or treatment conditions. Indeed MELD has already seen use in several contexts [50–54]. To facilitate the application of these tools for future scRNA-seq analysis, we provide open-source Python implementations that inherit the Scikit-learn API in the MELD package on GitHub <https://github.com/KrishnaswamyLab/MELD>.

4 Methods

In this section, we will provide details about our computational methods for computing the EES, as well as extracting information from the EES by way of a method we call *vertex frequency clustering*. We will outline the mathematical foundations for each algorithm, explain how they relate to previous works in manifold learning and graph signal processing, and provide details of the implementations of each algorithm.

4.1 Computation of the EES

Computing the EES involves the following steps each of which we will describe in detail.

1. A cell similarity graph is built over the combined data from all samples where each node or vertex in the graph is a cell and edges in the graph connect cells with similar gene expression values.
2. The condition label for each cell is used to create the Raw Experimental Signal (RES).
3. The RES is then smoothed over the graph to calculate the EES using a graph filter called the EES filter.

4.1.1 Graph construction

The first step in the EES algorithm is to create a cell similarity graph. In single-cell RNA sequencing, each cell is measured as a vector of gene expression counts measured as unique molecules of mRNA. Following best practices for scRNA-seq analysis [1], we normalize these counts by the total number of Unique Molecular Indicators (UMIs) per cell to give relative abundance of each gene and apply a square-root transform. Next we compute the similarity all pairs of cells, by using their Euclidean distances as an input to a kernel function. More formally, we compute a similarity matrix W such that each entry W_{ij} encodes the similarity between cell gene expression vectors \mathbf{x}_i and \mathbf{x}_j from the dataset X .

In our implementation we use α -decaying kernel proposed by Moon et al. [3] because in practice it provides an effective graph construction for scRNA-seq analysis. However, in cases where batch, density, and technical artifacts confound graph construction, we also use a mutual nearest neighbor kernel as proposed by Haghverdi et al. [9].

The α -decaying kernel [3] is defined as

$$K_{k,\alpha}(x,y) = \frac{1}{2} \exp\left(-\left(\frac{\|x-y\|_2}{\varepsilon_k(x)}\right)^\alpha\right) + \frac{1}{2} \exp\left(-\left(\frac{\|x-y\|_2}{\varepsilon_k(y)}\right)^\alpha\right), \quad (4)$$

where x, y are data points, $\varepsilon_k(x), \varepsilon_k(y)$ are the distance from x, y to their k -th nearest neighbors, respectively, and α is a parameter that controls the decay rate (i.e., heaviness of the tails) of the kernel. This construction generalizes the popular Gaussian kernel, which is typically used in manifold learning, but also has some disadvantages alleviated by the α -decaying kernel, as explained in Moon et al. [3].

The similarity matrix effectively defines a weighted and fully connected graph between cells such that every two cells are connected and that the connection between cells x and y is given by $K(x, y)$. To allow for computational efficiency, we sparsify the graph by setting very small edge weights to 0.

While the kernel in **Equation 4** provides an effective way of capturing neighborhood structure in data, it is susceptible to batch effects. For example, when data is collected from multiple patients, subjects, or environments (generally referred to as “batches”), such batch effects can cause affinities within each batch are often much higher than between batches, thus artificially creating separation between them rather than follow the underlying biological state. To alleviate such effects, we adjust the kernel construction using an approach inspired by recent work from by Haghverdi et al. [9] on the Mutual Nearest Neighbors (MNN) kernel. We extend the standard MNN approach, which has previous been applied to the k -Nearest Neighbors kernel, to the α -decay kernel as follows. First, within each batch, the affinities are computed using **Equation 4**. Then, across batches, we compute slightly modified affinities as

$$K'_{k,\alpha}(x,y) = \min\left\{\exp\left(-\left(\frac{\|x-y\|_2}{\varepsilon'_k(x)}\right)^\alpha\right), \exp\left(-\left(\frac{\|x-y\|_2}{\varepsilon'_k(y)}\right)^\alpha\right)\right\},$$

where $\varepsilon'_k(x)$ are now computed via the k -th nearest neighbor of x in the batch containing y (and vice versa for $\varepsilon'_k(y)$). Next, a rescaling factor γ_{xy} is computed such that

$$\sum_{z \in \text{batch}(y)} \gamma_{xy} K'_{k,\alpha}(x, z) \leq \beta \sum_{z \in \text{batch}(x)} K_{k,\alpha}(x, z)$$

for every x and y , where $\beta > 0$ is a user configurable parameter. This factor gives rise to the rescaled kernel

$$K'_{k,\alpha,\beta}(x, y) = \begin{cases} K'_{k,\alpha}(x, y) & \text{if } \text{batch}(x) = \text{batch}(y) \\ \gamma_{xy} K'_{k,\alpha}(x, y) & \text{otherwise.} \end{cases}$$

Finally, the full symmetric kernel is then computed as

$$K'_{k,\alpha}(x, y) = K'_{k,\alpha}(y, x) = \min\{K'_{k,\alpha,\beta}(x, y), K'_{k,\alpha,\beta}(y, x)\},$$

and used to set the weight matrix for the constructed graph over the data. Note that this construction is a well-defined extension of (**Equation 4**), as it reduces back to that kernel when only a single batch exists in the data.

We also perform an anisotropic density normalization transformation so that the kernel reflects the underlying geometry normalized by density as in Coifman and Lafon [55]. The density normalized kernel $K_{k,\alpha}^q$ divides out by density, estimated by the sum of outgoing edge weights for each node is as follows,

$$K_{k,\alpha}^q = \frac{K'_{k,\alpha}(x, y)}{q(x)q(y)},$$

where

$$q(x) = \int_X K'_{k,\alpha} q(y) dy.$$

We use this density normalized kernel in all experiments. When the data is uniformly sampled from the manifold then the density around each point is constant then this normalization has no effect. When the density is non-uniformly sampled from the manifold this allows an estimation of the underlying geometry unbiased by density. This is especially important when performing density estimation from empirical distributions with different underlying densities. By normalizing by density, we allow for construction of the manifold geometry from multiple differently distributed samples and individual density estimation for each of these densities on the same support. This normalization is further discussed in **Section 4.1.10**.

4.1.2 Estimating density and conditional likelihood on a graph

Density estimation is difficult in high dimensions because the number of samples needed to accurately reconstruct density with bounded error is exponential in the number of dimensions. Since general high dimensional density estimation is an intrinsically difficult problem, additional assumptions must be made. A common assumption is that the data exists on a manifold of low intrinsic dimensionality in ambient space. Under this assumption a number of works on graphs have addressed density estimation limited to the support of the graph nodes [56–60]. Instead of estimating kernel density or histograms in D dimensions where D could be large, these methods rendered the data as a graph, and density is estimated each point on the graph (each data point) as some variant counting the number of points which lie within a radius of each point on the graph.

The EES algorithm also estimates density of a signal on a graph. We use a generalization of the standard heat kernel on the graph to estimate density (See **Section 4.1.7**). We draw analogs between the EES and Gaussian kernel density estimation on the manifold. Where the EES with a specific parameter set is equivalent to the Gaussian density estimate on the graph (See **Section 4.1.10**).

4.1.3 Graph Signal Processing

The EES algorithm leverages recent advances in graph signal processing (GSP) [21], which aim to extend traditional signal processing tools from the spatiotemporal domain to the graph domain. Such extensions include, for example, wavelet transforms [61], windowed Fourier transforms [25], and uncertainty principles [62]. All of these extensions rely heavily on the fundamental analogy between classical Fourier transform and graph Fourier transform (described in the next section) derived from eigenfunctions of the graph Laplacian, which is defined as

$$\mathcal{L} := D - W, \tag{5}$$

where D is the *degree* matrix, which is a diagonal matrix with $D_{ii} = d(i) = \sum_j W_{ij}$ containing the degrees of the vertices of the graph defined by W .

4.1.4 The Graph Fourier Transform

One of the fundamental tools in traditional signal processing is the Fourier transform, which extracts the frequency content of spatiotemporal signals [63]. Frequency information enables various insights into important characteristics of analyzed signals, such as pitch in audio signals or edges and textures in images. Common to all of these is the relation between frequency and notions of *smoothness*. Intuitively, a function is *smooth* if one is unlikely to encounter a dramatic change in value across neighboring points. A simple

way to imagine this is to look at the *zero-crossings* of a function. Consider, for example, sine waves $\sin ax$ of various frequencies $a = 2^k$, $k \in \mathbb{N}$. For $k = 0$, the wave crosses the x-axis (a zero-crossing) when $x = \pi$. When we double the frequency at $k = 1$, our wave is now twice as likely to cross the zero and is thus less smooth than $k = 0$. This simple zero-crossing intuition for smoothness is relatively powerful, as we will see shortly.

Next, we show that our notions of smoothness and frequency are readily applicable to data that is not regularly structured, such as single-cell data. The graph Laplacian \mathcal{L} can be considered as a graph analog of the Laplace (second derivative) operator ∇^2 from multivariate calculus. This relation can be verified by deriving the graph Laplacian from first principles.

For a graph \mathcal{G} on N vertices, its graph Laplacian \mathcal{L} and an arbitrary graph signal $\mathbf{f} \in \mathbb{R}^N$, we use **Equation 5** to write

$$\begin{aligned} (\mathcal{L} \mathbf{f})(i) &= ([D - W] \mathbf{f})(i) \\ &= d(i)\mathbf{f}(i) - \sum_j W_{ij}\mathbf{f}(j) \\ &= \sum_j W_{ij}(\mathbf{f}(i) - \mathbf{f}(j)). \end{aligned} \quad (6)$$

As the graph Laplacian is a weighted sum of differences of a function around a vertex, we may interpret it analogously to its continuous counterpart as the curvature of a graph signal. Another common interpretation made explicit by the derivation in **Equation 6** is that $(\mathcal{L}\mathbf{f})(i)$ measures the *local variation* of a function at vertex i .

Local variation naturally leads to the notion of *total variation*,

$$\mathbf{TV}(\mathbf{f}) = \sum_{i,j} W_{ij}(\mathbf{f}(i) - \mathbf{f}(j))^2,$$

which is effectively a sum of all local variations. $\mathbf{TV}(\mathbf{f})$ describes the global smoothness of the graph signal \mathbf{f} . In this setting, the more smooth a function is, the lower the value of the variation. This quantity is more fundamentally known as the *Laplacian quadratic form*,

$$\mathbf{f}^T \mathcal{L} \mathbf{f} = \sum_{i,j} W_{ij}(\mathbf{f}(i) - \mathbf{f}(j))^2. \quad (7)$$

Thus, the graph Laplacian can be used as an operator and in a quadratic form to measure the smoothness of a function defined over a graph. One effective tool for analyzing such operators is to examine their eigensystems. In our case, we consider the eigendecomposition $\mathcal{L} = \Psi \Lambda \Psi^{-1}$, with eigenvalues³ $\Lambda := \{0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N\}$ and corresponding eigenvectors $\Psi := \{\psi_i\}_{i=1}^N$. As the Laplacian is a square, symmetric matrix, the spectral theorem tells us that its eigenvectors in Ψ form an orthonormal basis for \mathbb{R}^N . Furthermore, the Courant-Fischer theorem establishes that the eigenvalues in Λ are local minima of $\mathbf{f}^T \mathcal{L} \mathbf{f}$ when $\mathbf{f}^T \mathbf{f} = 1$ and $\mathbf{f} \in U$ as $\dim(U) = i = 1, 2, \dots, N$. At each eigenvalue λ_i this function has $\mathbf{f} = \psi_i$. In summary, the eigenvectors of the graph Laplacian (1) are an orthonormal basis and (2) minimize the Laplacian quadratic form for a given dimension.

³Note that in this discussion we abuse notation by treating Λ as an ordered set of Laplacian eigenvalues and as the diagonal matrix with entries from the elements of this set. Similarly, Ψ is both the set of column eigenvectors $\{\psi_i\}_{i=1}^N$ as well as the $N \times N$ matrix $[\psi_1 \psi_2 \dots \psi_N]$ with eigenvector as a column.

Henceforth, we use the term *graph Fourier basis* interchangeably with graph Laplacian eigenvectors, as this basis can be thought of as an extension of the classical Fourier modes to irregular domains [21]. In particular, the ring graph eigenbasis is composed of sinusoidal eigenvectors, as they converge to discrete Fourier modes in one dimension. The graph Fourier basis thus allows one to define the *graph Fourier transform* (GFT) by direct analogy to the classical Fourier transform.

The GFT of a signal f is given by $\hat{f}(\lambda_\ell) = \sum_i f(i) \psi_\ell^T(i) = \langle \mathbf{f}, \psi_\ell \rangle$, which can also be written as the matrix-vector product

$$\hat{\mathbf{f}} = \Psi^T \mathbf{f}. \quad (8)$$

As this transformation is unitary, the inverse graph Fourier transform (IGFT) is $\mathbf{f} = \Psi \hat{\mathbf{f}}$. Although the graph setting presents a new set of challenges for signal processing, many classical signal processing notions such as filterbanks and wavelets have been extended to graphs using the GFT. We use the GFT to process, analyze, and cluster experimental signals from single-cell data using a novel graph filter construction and a new harmonic clustering method.

4.1.5 The EES Filter

In the EES algorithm, we seek to estimate the change in likelihood between two experimental labels along a manifold represented by a cell similarity graph. To estimate likelihood along the graph, we employ a novel graph filter construction, which we explain in the following sections. To begin, we review the notion of filtering with focus on graphs and demonstrate the filter in a low-pass setting. Next, we demonstrate the expanded version of the EES filter and provide an analysis of its parameters. Finally, we provide a simple solution to the EES filter that allows fast computation.

4.1.6 Filters on graphs

Filters can be thought of as devices that alter the spectrum of their input. Filters can be used as bases, as is the case with wavelets, and they can be used to directly manipulate signals by changing the frequency response of the filter. For example, many audio devices contain an equalizer that allows one to change the amplitude of bass and treble frequencies. Simple equalizers can be built simply by using a set of filters called a filterbank. In the EES algorithm, we use a tunable filter to amplify latent features on a single-cell graph.

Mathematically, graph filters work analogously to classical filters. Particularly, a filter takes in a signal and attenuates it according to a frequency response function. This function accepts frequencies and returns a response coefficient. This is then multiplied by the input Fourier coefficient at the corresponding frequency. The entire filter operation is thus a reweighting of the input Fourier coefficients. In low-pass filters, the function only preserves frequency components below a threshold. Conversely, high-pass filters work by removing frequencies below a threshold. Bandpass filters transfer frequency components that are within a certain range of a central frequency. The tunable filter in the EES algorithm is capable of producing any of these responses.

As graph harmonics are defined on the set Λ , it is common to define them as functions of the form $h : [0, \max(\Lambda)] \mapsto [0, 1]$. For example, a low pass filter with cutoff at λ_k would have $h(x) > 0$ for $x < \lambda_k$ and $h(x) = 0$ otherwise. By abuse of notation, we will refer to the diagonal matrix with the filter h applied to each Laplacian eigenvalue as $h(\Lambda)$, though h is not a set-valued or matrix-valued function. Filtering a signal \mathbf{f} is clearest in the spectral domain, where one simply takes the multiplication $\hat{\mathbf{f}}_{\text{filt}} = h(\Lambda) \hat{\mathbf{f}} = h(\Lambda) \Psi^T \mathbf{f}$.

Finally, it is worth using the above definitions to define a vertex-valued operator to perform filtering. As a graph filter is merely a reweighting of the graph Fourier basis, one can construct the *filter matrix*,

$$H = \Psi h(\Lambda) \Psi^T. \quad (9)$$

A manipulation using **Equation 8** will verify that $H\mathbf{f}$ is the WGFT of $\hat{\mathbf{f}}_{\text{filt}}$. This filter matrix will be used to solve the EES filter in approximate form for computational efficiency.

4.1.7 Laplacian Regularization

A simple assumption for recovering the EES signal from raw measurements is *smoothness*. In this model the latent signal is assumed to have a low amount of neighbor to neighbor variation. *Laplacian regularization* [64–72] is a simple technique that targets signal smoothness via the optimization

$$\mathbf{y} = \underset{\mathbf{z}}{\operatorname{argmin}} \underbrace{\|\mathbf{x} - \mathbf{z}\|_2^2}_a + \underbrace{\beta \mathbf{z}^T \mathcal{L} \mathbf{z}}_b. \quad (10)$$

Note that this optimization has two terms. The first term (a), called a *reconstruction penalty*, aims to keep the EES similar to the RES. The second term (b) ensures smoothness of the signal. Balancing these terms adjusts the amount of smoothness performed by the filter.

Laplacian regularization is a sub-problem of the EES filter that we will discuss for low-pass filtering. In the above, a reconstruction penalty (a) is considered alongside the Laplacian quadratic form (b), which is weighted by the parameter β . The Laplacian quadratic form may also be considered as the norm of the *graph gradient*, i.e.

$$\beta \mathbf{z}^T \mathcal{L} \mathbf{z} = \beta \|\nabla_G \mathbf{z}\|_2^2.$$

Thus one may view Laplacian regularization as a minimization of the edge-derivatives of a function while preserving a reconstruction. Because of this form, this technique has been cast as *Tikhonov regularization* [66, 73], which is a common regularization to enforce a low-pass filter to solve inverse problems in regression. In our results we demonstrate a EES filter that may be reduced to Laplacian regularization using a squared Laplacian.

In **Section 4.1.6** we introduced filters as functions defined over the Laplacian eigenvalues ($h(\Lambda)$) or as vertex operators in **Equation 9**. Minimizing optimization **Equation 10** reveals a similar form for Laplacian regularization. Although Laplacian regularization filter is presented as an optimization, it also has a closed form solution. We derive this solution here as it is a useful building block for understanding the EES. To begin,

$$\begin{aligned} \mathbf{y} &= \underset{\mathbf{z}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{z}\|_2^2 + \beta \mathbf{z}^T \mathcal{L} \mathbf{z} \\ &= \underset{\mathbf{z}}{\operatorname{argmin}} (\mathbf{x} - \mathbf{z})^T (\mathbf{x} - \mathbf{z}) + \beta \mathbf{z}^T \mathcal{L} \mathbf{z} \\ &= \underset{\mathbf{z}}{\operatorname{argmin}} \mathbf{x}^T \mathbf{x} + \mathbf{z}^T \mathbf{z} - 2\mathbf{x}^T \mathbf{z} + \beta \mathbf{z}^T \mathcal{L} \mathbf{z} \end{aligned}$$

Substituting $\mathbf{y} = \mathbf{z}$, we next differentiate with respect to \mathbf{y} and set this to 0,

$$\begin{aligned} 0 &= \nabla_{\mathbf{y}} (\mathbf{x}^T \mathbf{x} + \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{x} + \beta \mathbf{y}^T \mathcal{L} \mathbf{y}) \\ &= 2\mathbf{y} - 2\mathbf{x} + 2\beta \mathcal{L} \mathbf{y} \\ \mathbf{x} &= (\mathbf{I} + \beta \mathcal{L}) \mathbf{y}, \end{aligned}$$

so the global minima of (10) can be expressed in closed form as

$$\mathbf{y} = (\mathbf{I} + \beta \mathcal{L})^{-1} \mathbf{x}. \quad (11)$$

As the input \mathbf{x} is a graph signal in the vertex domain, the least squares solution (11) is a filter matrix $H_{\text{reg}} = (\mathbf{I} + \beta \mathcal{L})^{-1}$ as discussed in **Section 4.1.6**. The spectral properties of Laplacian regularization immediately follow as

$$\begin{aligned} H_{\text{reg}} &= (\mathbf{I} + \beta \mathcal{L})^{-1} \\ &= \Psi \frac{1}{1 + \beta \Lambda} \Psi^T. \end{aligned} \quad (12)$$

Thus Laplacian regularization is a graph filter with frequency response $h_{\text{reg}}(\lambda) = (1 + \beta \lambda)^{-1}$. **Figure S14b** shows that this function is a low-pass filter on the Laplacian eigenvalues with cutoff parameterized by β .

4.1.8 Tunable Filtering

Though simple low-pass filtering with Laplacian regularization is a powerful tool for many machine learning tasks, we sought to develop a filter that is flexible and capable of filtering the signal at any frequency. To accomplish these goals, we introduce the EES filter:

$$\begin{aligned} \mathbf{y} &= \underset{\mathbf{z}}{\text{argmin}} \|\mathbf{x} - \mathbf{z}\|_2^2 + \mathbf{z}^T \mathcal{L}_* \mathbf{z} \\ \text{where } \mathcal{L}_* &= \exp(\beta(\mathcal{L}/\lambda_{\max} - \alpha \mathbf{I})^\rho) - \mathbf{I} \end{aligned} \quad (13)$$

This filter expands upon Laplacian regularization by the addition of a new smoothness structure. Early and related work proposed the use of a power Laplacian smoothness matrix S in a similar manner as we apply here [66], but little work has since proven its utility. In our construction, α is referred to as modulation, β acts as a reconstruction penalty, and ρ is filter order. These parameters add a great deal of versatility to the EES filter, and we demonstrate their spectral and vertex effects in **Figure S14**, as well as provide mathematical analysis of the EES algorithm parameters in **Section 4.1.8**. Finally, in **Section 4.1.11** we discuss an implementation of the filter.

A similar derivation as **Section 4.1.7** reveals the filter matrix

$$H_{\text{EES}}(\mathcal{L}) = e^{-\beta(\mathcal{L}/\lambda_{\max} - \alpha \mathbf{I})^\rho}, \quad (14)$$

which has the frequency response

$$h_{\text{EES}}(\lambda) = e^{-\beta(\lambda/\lambda_{\max} - \alpha)^\rho}. \quad (15)$$

Thus, the value of the EES algorithm parameters in the vertex optimization (**Equation 13**) has a direct effect on the graph Fourier domain.

4.1.9 Parameter Analysis

β steepens the cutoff of the filter and shifts it more towards its central frequency (**Figure S14b**). In the case of $\alpha = 0$, this frequency is $\lambda_1 = 0$. This is done by scaling all frequencies by a factor of β . For stability reasons, we choose $\beta > 0$, as a negative choice of β yields a high frequency amplifier.

The parameters α and ρ change the filter from low pass to band pass or high pass. **Figure S14** highlights the effect on frequency response of the filters and showcases their vertex effects in simple examples. We begin our mathematical analysis with the effects of ρ .

ρ powers the Laplacian harmonics. This steepens the frequency response around the central frequency of the EES filter. Higher values of ρ lead to sharper tails (**Figure S14c, S14e**), limiting the frequency response outside of the target band, but with increased response within the band. Finally, ρ can be used to make a high pass filter by setting it to negative values (**Figure S14f**).

For the integer powers, a basic vertex interpretation of ρ is available. Each column of \mathcal{L}^k is k -hop localized, meaning that \mathcal{L}_{ij}^k is non-zero if and only if there exists a path length k between vertex i and vertex j (for a detailed discussion of this property, see Hammond et al. [61, section 5.2].) Thus, for $\rho \in \mathbb{N}$, the operator \mathcal{L}^ρ considers variation over a hop distance of ρ . This naturally leads to the spectral behavior we demonstrate in **Figure S14c**, as signals are required to be smooth over longer hop distances when $\alpha = 0$ and $\rho > 1$.

The parameter α removes values from the diagonal of \mathcal{L} . This results in a *modulation* of frequency response by translating the Laplacian harmonic that yields the minimal value for the problem (**Equation 13**). This allows one to change the central frequency, as α effectively modulates a band-pass filter. As graph frequencies are positive, we do not consider $\alpha < 0$. In the vertex domain, the effect of α is more nuanced. We study this parameter for $\alpha > 0$ by considering a modified Laplacian \mathcal{L}_* with $\rho = 1$.

To conclude, we propose a filter parameterized by reconstruction β (**Figure S14b**), order ρ (**Figure S14c, S14e**), and modulation α (**Figure S14d**). The parameters α and β are limited to be strictly greater than or equal to 0. When $\alpha = 0$, ρ may be any integer, and it adds more low frequencies to the frequency response as it becomes more positive. On the other hand, if ρ is negative and $\alpha = 0$, ρ controls a high pass filter. When $\alpha > 0$, the EES filter becomes a band-pass filter. In standard use cases we propose to use the parameters $\alpha = 0$, $\beta = 60$, and $\rho = 1$. Other parameter values are explored further in (**Figure S13**). We note that the results are relatively robust to parameter values around this default setting. All of our biological results were obtained using this parameter set, which gives a square-integrable low-pass filter. As these parameters have direct spectral effects, their implementation in an efficient graph filter is straightforward and presented in **Section 4.1.11**.

In contrast to previous works using Laplacian filters, our parameters allow analysis of signals that are combinations of several underlying changes occurring at various frequencies. For an intuitive example, consider that the frequency of various Google searches will vary from winter to summer (low-frequency variation), Saturday to Monday (medium-frequency variation), or morning to night (high-frequency variation). In the biological context such changes could manifest as differences in cell type abundance (low-frequency variation) and cell-cycle (medium-frequency variation) [74]. We illustrate such an example in **Figure S14a** by blindly separating a medium frequency signal from a low frequency contaminating signal over simulated data. Such a technique could be used to separate low- and medium-frequency components so that they can be analyzed independently. Each of the filter parameters is explained in more detail in **Section 4.1.8**.

4.1.10 Relation between the EES and Gaussian KDE through the Heat Kernel

Kernel density estimators (KDEs) are widely used as estimating density is one of the fundamental tasks in many data applications. The density estimate is normally done in ambient space, and there are many methods to do so with a variety of advantages and disadvantages depending on the application. We instead assume that the data is sampled from some low dimensional subspace of the ambient space, e.g. that the data lies along a manifold. The EES can be thought of as a Gaussian KDE over the discrete manifold formed by the data. This gives a density estimate at every sampled point for a number of distributions.

This density estimate, as the number of samples goes to infinity, should converge to the density estimate along a continuous manifold formed by the data. The case of data uniformly sampled on the manifold was explored in [22] proving convergence of the eigenvectors and eigenvalues of the discrete Laplacian to the eigenfunctions of the continuous manifold. Coifman and Maggioni [75] explored when the data is non-uniformly sampled from the manifold and provided a kernel which can normalize out this density which results in a Laplacian modeling the underlying manifold geometry, irrespective of data density. Building on these two works the EES allows us to estimate the manifold geometry using multiple samples with unknown distribution along it and estimate density and conditional density for each distribution on this shared manifold.

A general kernel density estimator (KDE) $f(x, t)$ with bandwidth $t > 0$ and kernel function $K(x, y, t)$ is defined as

$$\hat{f}(x, t) = \frac{1}{N} \sum_{i=1}^N K(x, X_i, t), \quad x \in \mathcal{X} \quad (16)$$

With $\mathcal{X} := \mathbb{R}^d$, and endowed with the Gaussian kernel,

$$K(x, y, t) = \frac{1}{(4\pi t)^{d/2}} e^{-\|x-y\|_2^2/4t}, \quad (17)$$

we have the Gaussian KDE in \mathbb{R}^d .

This kernel is of particular interest for its thermodynamic interpretation. Namely the Gaussian KDE is a space discretization of the unique solution to the heat diffusion partial differential equation (PDE) [23, 76]:

$$\frac{\partial}{\partial t} \hat{f}(x, t) = \frac{1}{2} \frac{\partial^2}{\partial x^2} \hat{f}(x, t), \quad x \in \mathcal{X}, t > 0, \quad (18)$$

with $\hat{f}(x, 0) = \frac{1}{N} \sum_{i=1}^n \delta_{X_i}$ where δ_x is the Dirac measure at x . This is sometimes called Green's function for the diffusion equation. Intuitively, $\hat{f}(x, t)$ can be thought of as measuring the heat after time t after placing units of heat on the data points at $t = 0$.

In fact the Gaussian kernel can be represented instead in terms of the eigenfunctions of the ambient space. With eigenfunctions ϕ and eigenvalues λ , the Gaussian kernel can be alternative expressed as:

$$K(x, y, t) = \sum_{n=0}^{\infty} e^{-t\lambda_n} \phi_n(x) \phi_n(y) \quad (19)$$

Of course for computational reasons we often prefer the closed form solution in (17). We now consider the case when \mathcal{X} instead consists of uniform samples from a Riemannian manifold \mathcal{M} embedded in \mathbb{R}^d such that $\mathcal{X} \subset \mathcal{M} \subset \mathbb{R}^d$. An analog to the Gaussian KDE in \mathbb{R}^d on a manifold is then the solution to the heat PDE restricted to the manifold, and again we can use the eigenfunction interpretation of the Green's function in (19), except replacing the eigenfunctions of the manifold.

The eigenfunctions of the manifold can be approximated through the eigenvectors of the discrete Laplacian. The solution of the heat equation on a graph is defined in terms of the discrete Laplacian $\mathcal{L} = \Psi \Lambda \Psi^{-1}$ as

$$\hat{K}_{\mathcal{L}}(x, y, t) = \delta_x e^{-t\mathcal{L}} \delta_y = \delta_x \Psi e^{-t\Lambda} \Psi^{-1} \delta_y \quad (20)$$

Where δ_x, δ_y are dirac functions at x and y respectively. This is equivalent to the EES when $\beta = t\lambda_{max}$, $\alpha = 0$, and $\phi = 1$.

When data \mathcal{X} is sampled uniformly from the manifold \mathcal{M} and the standard gaussian kernel is used to construct the graph, then Theorem 2.1 of Belkin and Niyogi [22] which proves the convergence of the eigenvalues of the discrete graph laplacian to the continuous laplacian implies (20) converges to the Gaussian KDE on the manifold.

However, real data is rarely uniformly sampled from a manifold. When the data is instead sampled from a smooth density $\mathcal{X} \sim q(x)$ over the manifold then the density must be normalized out to recover the geometry of the manifold. This problem was first tackled in Coifman and Lafon [55], by constructing an anisotropic kernel which divides out the density at every point. This correction allows us to estimate density over the underlying *geometry of the manifold* even in the case where data is not uniformly sampled. This allows us to use samples from multiple distributions, in our case distributions over cellular states, which allows a better estimate of underlying manifold utilizing all available data.

In practice, we combine two methods to construct a discrete Laplacian that reflects the underlying data geometry over which we estimate heat propagation and perform density estimation, as explained in **Section 4.1.1**.

4.1.11 Implementation

A naïve implementation of the EES algorithm would apply the matrix inversion presented in **Equation 14**. This approach is untenable for the large single-cell graphs that the EES algorithm is designed for, as H_{EES}^{-1} will have many elements, and, for high powers of ρ or non-sparse graphs, extremely dense. A second approach to solving **Equation 13** would diagonalize \mathcal{L} such that the filter function in **Equation 15** could be applied directly to the Fourier transform of input raw experimental signals. This approach has similar shortcomings as eigendecomposition is substantively similar to inversion. Finally, a speedier approach might be to use conjugate gradient or proximal methods. In practice, we found that these methods are not well-suited for EES filtering.

Instead of gradient methods, we use Chebyshev polynomial approximations of $h_{\text{EES}}(\lambda)$ to rapidly approximate and apply the EES filter. These approximations, proposed by Hammond et al. [61] and Shuman et al. [24], have gained traction in the graph signal processing community for their efficiency and simplicity. Briefly, a truncated and shifted Chebyshev polynomial approximation is fit to the frequency response of a graph filter. For analysis, the approximating polynomials are applied as polynomials of the Laplacian multiplied by the signal to be filtered. As Chebyshev polynomials are given by a recurrence relation, the approximation procedure reduces to a computationally efficient series of matrix-vector multiplications. For a more detailed treatment one may refer to Hammond et al. [61] where the polynomials are proposed for graph filters. For application of the EES filter to a small set of input RES, Chebyshev approximations offer the simplest and most efficient implementation of our proposed algorithm. For sufficiently large sets of RES, such as when considering hundreds of conditions, the computational cost of obtaining the Fourier basis directly may be less than repeated application of the approximation operator; in these cases, we diagonalize the Laplacian either approximately through randomized SVD or exactly using eigendecomposition, depending on user preference. Then, one simply constructs $H_{\text{EES}} = \Psi h_{\text{EES}}(\Lambda) \Psi^T$ to calculate the EES from the RES.

4.1.12 Summary of the EES algorithm

In summary, we have proposed a family of graph filters based on a generalization of Laplacian regularization framework to implement the computation of the EES. This optimization, which can be solved analytically, allows us to derive the EES, or conditional likelihood of the experimental label, as a smooth and

denoised signal, while also respecting multi-resolution changes in the likelihood landscape. As we show in **Section 4.7**, this formulation performs better at deriving the true conditional likelihood in quantitative comparisons than simpler label smoothing algorithms. Further, the EES algorithm it is efficient to compute.

The EES algorithm is implemented in Python 3 as part of the MELD package and is built atop the `scprep`, `graphtools`, and `pygsp` packages. We developed `scprep` efficiently process single-cell data, and `graphtools` was developed for construction and manipulation of graphs built on data. Fourier analysis and Chebyshev approximations are implemented using functions from the `pygsp` toolbox [77].

4.2 Vertex-frequency clustering

Next, we will describe the vertex frequency clustering algorithm for partitioning the cellular manifold into regions of similar response to experimental perturbation. For this purpose, we use a technique proposed in Shuman et al. [25] based on a graph generalization of the classical Short Time Fourier Transform (STFT). This generalization will allow us to simultaneously localize signals in both frequency and vertex domains. The output of this transform will be a spectrogram Q , where the value in each entry $Q_{i,j}$ indicates the degree to which the RES in the neighborhood around vertex i is composed of frequency j . We then concatenate the EES and perform k -means clustering. The resultant clusters will have similar transcriptomic profiles, similar EES values, and similar *frequency trends* of the RES. The frequency trends of the RES are important because they allow us to infer movements in the cellular state space that occur during experimental perturbation.

We derive vertex frequency clusters in the following steps:

1. We create the cell graph in the same way as is done to derive the EES in **Section 4.1.1**.
2. For each vertex in the graph (corresponding to a cell in the data), we create a series of localized windowed signals by masking the RES using a series of heat kernels centered at the vertex. Graph Fourier decomposition of these localized windows capture frequency of the RES at different scales around each vertex.
3. The graph Fourier representation of the localized windowed signals is thresholded using a *tanh* activation function to produce pseudo-binary signals.
4. These pseudo-binarized signals are summed across windows of various scales to produce a single $N \times N$ spectrogram Q . PCA is performed on the spectrogram for dimensionality reduction.
5. The EES is concatenated to the reduced spectrogram weighted by the $L2$ -norm of PC1 to produce \hat{Q} which captures both local RES frequency trends and changes in conditional density around each cell in both datasets.
6. k-Means is performed on the concatenated matrix to produce vertex-frequency clusters.

4.2.1 Analyzing frequency content of the RES

Before we go into further detail about the algorithm, it may be useful to provide some intuitive explanations for why the frequency content of the RES provides a useful basis for identifying clusters of cells affected by an experimental perturbation. Because the low frequency eigenvectors of the graph Laplacian identify smoothly varying axes of variance through a graph, we associate trends in the RES associated these low-frequency eigenvectors as biological transitions between cell states. This may correspond to the shift in T cells from naive to activated, for example. We note that at intermediate cell transcriptomic states between

the extreme states that are most enriched in either condition, we observe both low and middle frequency RES components, see the blue cell in the cartoon in **Figure 2a**. This is because locally, the RES varies from cell to cell, but on a large scale is varying from enriched in one condition to being enriched in the other. This is distinct from what we observe in our model when a group of cells are completely unaffected by an experimental perturbation. Here, we expect to find only high frequency variations in the RES and no underlying transition or low-frequency component. The goal of vertex frequency clustering is to distinguish between these four cases: enriched in the experiment, enriched in the control, intermediate transitional states, and unaffected populations of cells. We also want these clusters to have variable size so that even small groups of cells that may be differentially abundant are captured in our clusters.

4.2.2 Using the Windowed Graph Fourier Transform (WGFT) to identify local changes in RES frequency

While the graph Fourier transform is useful for exploring the frequency content of a signal, it is unable to identify how the frequency content of graph signals change locally over different regions of the graph. In vertex frequency clustering, we are interested in understanding how the frequency content of the RES changes in neighborhoods around each cell. In the time domain, the windowed Fourier transform identifies changing frequency composition of a signal over time by taking slices of the signal (e.g. a sliding window of 10 seconds) and applying a Fourier decomposition to each window independently (WFT) [63]. The result is a spectrogram Q , where the value in each cell $Q_{i,j}$ indicates the degree to which time-slice i is composed of frequency j . Recent works in GSP have generalized the constructions windowed Fourier transform to graph signals[25]. To extend the notion of a sliding window to the graph domain, Shuman et al. [25] write the operation of translation in terms of convolution as follows.

The *generalized translation operator* $T_i : \mathbb{R}^N \rightarrow \mathbb{R}^N$ of signal f to vertex $i \in \{1, 2, \dots, N\}$ is given by

$$(T_i f)(n) := \sqrt{N}(f * \delta_i)(n), \quad \delta_i(j) = \begin{cases} 1 & j = i \\ 0 & j \neq i \end{cases} \quad (21)$$

which convolves the signal f , in our case the RES, with a dirac at vertex i . Shuman et al. [25] demonstrate that this operator inherits various properties of its classical counterpart; however, the operator is not isometric and is affected by the graph that it is built on. Furthermore, for signals that are not tightly localized in the vertex domain and on graphs that are not directly related to Fourier harmonics (e.g., the circle graph), it is not clear what graph translation implies.

In addition to translation, a *generalized modulation operator* is defined by Shuman et al. [25] as $M_k : \mathbb{R}^N \rightarrow \mathbb{R}^N$ for frequencies $k \in \{0, 1, \dots, N-1\}$ as

$$(M_k f)(n) := \sqrt{N}f(n)U_k(n) \quad (22)$$

This formulation is analogous in construction to classical modulation, defined by pointwise multiplication with a pure harmonic – a Laplacian eigenvector in our case. Classical modulation translates signals in the Fourier domain; because of the discrete nature of the graph Fourier domain, this property is only weakly shared between the two operators. Instead, the generalized modulation M_k translates the *DC component* of f , $\hat{f}(0)$, to λ_k , i.e. $(\widehat{M_k f})(\lambda_k) = \hat{f}(0)$. Furthermore, for any function f whose frequency content is localized around λ_0 , $(M_k f)$ is localized in frequency around λ_k . Shuman et al. [25] details this construction and provides bounds on spectral localization and other properties.

With these two operators, a graph windowed Fourier atom is constructed[25] for any window function $g \in \mathbb{R}^N$

$$g_{i,k}(n) := (M_k T_i g)(n) = N U_k(n) \sum_{\ell=0}^{N-1} \hat{g}(\lambda_\ell) U_\ell^*(i) U_\ell(n). \quad (23)$$

We can then build a spectrogram $Q = (q_{ik}) \in \mathbb{R}^{N \times N}$ by taking the inner product of each $g_{i,k} \forall i \in \{1, 2, \dots, N\} \wedge \forall k \in \{0, 1, \dots, N-1\}$ with the target signal f

$$q_{ik} = S f(i, k) := \langle f, g_{i,k} \rangle. \quad (24)$$

As with the classical windowed Fourier transform, one could interpret this as segmenting the signal by windows and then taking the Fourier transform of each segment

$$q_i = \langle (T_i g \odot f), U \rangle \quad (25)$$

where \odot is the element-wise product.

4.2.3 Using heat kernels of increasing scales to produce the WGFT of the RES

To generate the spectrogram for clustering, we first need a suitable window function. We use the normalized heat kernel as proposed by Shuman et al. [25]

$$\hat{g}(\lambda) = C e^{-t\lambda}, \quad (26)$$

$$C = \|g\|_2^{-1}. \quad (27)$$

By translating this kernel, element-wise multiplying it with our target signal f and taking the Fourier transform of the result, we obtain a windowed graph Fourier transform of f that is localized based on the *diffusion distance* [25, 62] from each vertex to every other vertex in the graph.

For an input RES \mathbf{f} , signal-biased spectral clustering as proposed by Shuman et al. [25] proceeds as follows:

1. Generate the window matrix P_t , which contains as its columns translated and normalized heat kernels at the scale t
2. Column-wise multiply $F_t = P \odot \mathbf{f}$; the i -th column of F_t is an entry-wise product of the i -th window and \mathbf{f} .
3. Take the Fourier Transform of each column of F_t . This matrix, \hat{C}_t is the normalized WGFT matrix.

This produces a single WGFT for the scale t . At this stage, Shuman et al. [25] proposed to saturate the elements of \hat{C}_t using the activation function $\tanh(|\hat{C}_t|)$ (where $|\cdot|$ is an element-wise absolute value). Then, k-means is performed on this saturated output to yield clusters. This operation has connections to spectral clustering as the features that k-means is run on are coefficients of graph harmonics.

We build upon this approach to add robustness, sensitivity to sign changes, and scalability. Particularly, vertex-frequency clustering builds a set of activated spectrograms at different window scales. These scales are given by simulated heat diffusion over the graph by adjusting the time-scale t in **Equation 26**. Then, the entire set is combined through summation.

4.2.4 Combining the EES and WGFT of the RES

As discussed in **Section 2.4**, it is useful to consider the sign of the EES in addition to the frequency content of the RES. This is because if we consider two populations of cells, one of which is highly enriched in the experimental condition and another that is enriched in the control, we expect to find similar frequency content of the RES. Namely, both should have very low-frequency content, as indicated in the cartoon in **Figure 2a**. However, we expect these two populations to have very different EES values. To allow us to distinguish between these populations, we also include the EES in the matrix used for clustering.

We concatenate the EES as an additional column to the multi-resolution spectrogram Q . However, we want to be able to tune the clustering with respect to how much the EES affects the result compared to the frequency information in Q . Therefore, inspired by spectral clustering as proposed by [78], we first perform PCA on Q to get $k + 1$ principle components and then normalize the EES by the L_2 -norm of the first principle component. We then add the EES as an additional column to the PCA-reduced Q to produce the matrix \hat{Q} . The weight of the EES can be modulated by a user-adjustable parameter w , but for all experiments in this paper, we leave $w = 1$. Finally, \hat{Q} is used as input for k -means clustering.

The multiscale approach we have proposed has a number of benefits. Foremost, it removes the complexity of picking a window-size. Second, using the actual input signal as a feature allows the clustering to consider both frequency and sign information in the raw experimental signal. For scalability, we leverage the fact that P_t is effectively a diffusion operator and thus can be built efficiently by treating it as a Markov matrix and normalizing the graph adjacency by the degree.

4.2.5 Summary of the vertex frequency clustering algorithm

To identify clusters of cells that are transcriptionally similar and also affected by an experimental perturbation in the same way, we introduced an algorithm called vertex frequency clustering. Our approach builds on previous work by Shuman et al. [25] analyzing the local frequency content of the RES (raw experimental signal) as defined over the vertices of a graph. Here, we introduce two novel adaptations of the algorithm. First, we take a multiresolution approach to quantifying frequency trends in the neighborhoods around each node. By considering windowed signals that are large (i.e. contain many neighboring points) and small (i.e. very proximal on the graph), we can identify clusters both large and small that are similarly affected by an experimental perturbation. Our second contribution is the inclusion of the EES in our basis for clustering. This allows VFC to take into account the degree of enrichment of each group of cells between condition.

4.3 A pipeline for analyzing single cell data using MELD

Using the EES algorithm and VFC, it is now possible to propose a novel framework for analyzing single cell perturbation experiments. The goal of this framework is to identify populations of cells that are the most affected by an experimental perturbation and to characterize a gene signature of that perturbation. A schematic of the proposed pipeline is shown in **Figure S4**.

Prior to using the algorithms in MELD, we recommended first following established best practices for analysis of single cell data including exploratory analysis using visualization, preliminary clustering, and cluster annotation via differential expression analysis [1]. These steps ensure that the dataset is of high quality and comprises the cell types expected from the experimental setup. Following exploratory characterization, we propose the following analysis:

1. Calculate the EES for the experimental and control condition

2. Determine which exploratory clusters require subclustering with VFC by examining the EES distribution within each cluster, a visualization of the cluster, and the results of VFC with varying numbers of clusters
3. Create new cluster assignments using VFC
4. Annotate each cluster following best practices [1]
5. Characterize enrichment of cell populations using EES and gene signatures

The basic steps to calculate the EES for each condition is described in **Section 2.1**. In the case of multiple replicates, we recommend calculating the EES for each sample over a graph of all cells from all samples so long as there is sufficient overlap between samples. This overlap can be assessed using the k-nearest neighbor batch effect test described in Büttner et al. [79]. We then normalize the EES for matched experimental and control samples of the same replicate. A single EES for the experimental condition across replicates can be calculated by averaging the treatment condition EES across replicates. Variation in this signal across replicates can be used as a measure of consistency for the measured perturbation across cell types. The result of this step is an estimate of the probability that each cell would be observed in the experimental condition relative to the control.

Having calculated the EES, we next recommend determining which cell populations identified during exploratory analysis require further subclustering with VFC to identify cell types enriched or depleted in the experimental condition. Determining optimal cluster resolution for single cell analysis will vary across experiments depending on the biological system being studied and the goals of each individual researcher. Instead of providing a single measure to determine the number of clusters, we outline a general strategy as a guide for users of MELD.

To determine the number of VFC clusters, we suggest taking into consideration transcriptional variation within each coarse-grained cluster and the effect of the perturbation. First, using a dimensionality reduction tool such as PHATE, examine a two or three dimensional scatter plot of the cluster colored by the EES for each cell. Here, the goal is to identify either regions that have very different EES values or regions of data density separated by low-density regions suggesting the present of multiple subclusters to target with VFC. We also suggest examining the distribution the EES values within each cluster to determine if the cells in the cluster exhibit a restricted range of responses to the EES or large variation that would require subclustering. Finally, we recommend running VFC with various numbers of clusters (2-5 is often sufficient) and inspecting the output on a PHATE plot and/or with a swarm plot. In ambiguous cases, it may be helpful to perform differential expression analysis and gene set enrichment to determine whether or not each cluster is biologically relevant to the experimental question under consideration [1, 80]. Importantly, not all clusters need subclustering, and we emphasize the ideal cluster resolution will vary based on the goals of each analyst.

To determine the gene signature of the perturbation, we recommend quantifying the differences in expression between VFC clusters. For experiments with only a single cell type and 3-4 VFC clusters, it is often sufficient to perform differential expression analysis between the cluster most enriched in the experimental condition and the cluster most depleted in the experimental condition. An example of this analysis is provided in **Section 2.6**. For experiments with several cell types, we recommend calculating the gene signature between the enriched and depleted VFC clusters within each exploratory cluster. To obtain a consensus gene signature, a research may take the intersection of the gene signatures within exploratory cluster. An example of this analysis is provided in **Section 2.8**.

We note that the strategy for identifying gene signatures outlined in the previous paragraph differs from the current framework employed in recent papers (**Figure S3**). Instead of comparing expression between cells from the experimental condition and the control, we compare clusters of cells identified with VFC. The rationale for the framework presented here is that if VFC clusters are transcriptionally homogeneous and exhibit a uniform response to the perturbation, we expect differences in gene expression between conditions *within* each cluster to represent biological and technical noise. However, characterizing transcriptional differences *between* cells of different clusters regardless of condition of origin will yield a description of the cell states that vary between experimental conditions. We confirm that the gene signatures obtained in this manner are more accurate than between-sample comparisons in **Section 4.7**.

4.4 Processing and analysis of the T-cell datasets

Gene expression counts matrices prepared by Datlinger et al. [16] were accessed from the NCBI GEO database accession GSE92872. 3,143 stimulated and 2,597 unstimulated T-cells were processed in a pipeline derived from the published supplementary software. First, artificial genes corresponding to gRNAs were removed from the counts matrix. Genes observed in fewer than five cells were removed. Cell with a library size higher than 35,000 UMI / cell were removed. To filter dead or dying cells, expression of all mitochondrial genes was z-scored and cells with average z-score expression greater than 1 were removed. As in the published analysis, all mitochondrial and ribosomal genes were excluded. Filtered cells and genes were library size normalized and square-root transformed. To build a cell-state graph, 100 PCA dimensions were calculated and edge weights between cells were calculated using an alpha-decay kernel as implemented in the Graphtools library (www.github.com/KrishnaswamyLab/graphtools) using default parameters. To infer the EES, MELD was run on the cell state graph using the stimulated / unstimulated labels as input with the smoothing parameter $\beta = 60$. To identify a signature, the top and bottom VFC clusters by EES value were used for differential expression using a rank test as implemented in diffxpy [28] and a q-value cutoff of 0.05. GO term enrichment was performed using EnrichR using the gseapy Python package (<https://pypi.org/project/gseapy/>).

4.5 Processing and analysis of the zebrafish dataset

Gene expression counts matrices prepared by Wagner et al. [18] (the chordin dataset) were downloaded from NCBI GEO (GSE112294). 16079 cells from *chd* embryos injected with gRNAs targeting chordin and 10782 cells from *tyr* embryos injected with gRNAs targeting tyrosinase were accessed. Lowly expressed genes detected in fewer than 5 cells were removed. Cells with library sizes larger than 15,000 UMI / cell were removed. Counts were library-size normalized and square root transformed. Cluster labels included with the counts matrices were used for cell type identification.

During preliminary analysis, a group of 24 cells were identified originating exclusively from the *chd* embryos. Despite an average library size in the bottom 12% of cells, these cells exhibited 546-fold, 246-fold, and 1210-fold increased expression of Sh3Tc1, LOC101882117, and LOC101885394 respectively relative to other cells. To the best of our knowledge, the function of these genes in development is not described. These cells were annotated by Wagner et al. [18] as belonging to 7 cell types including the Tailbud – Spinal Cord and Neural – Midbrain. These cells were excluded from further analysis.

To generate a cell state graph, 100 PCA dimensions were calculated from the square root transformed filtered gene expression matrix of both datasets. Edge weights between cells on the graph were calculated using an alpha-decay kernel with parameters $knn=20$, $decay=40$. MAGIC was used to impute gene expression values using default parameters. MELD was run using the *tyr* or *chd* labels as input. The EES was

calculated for each of the 6 samples independently and normalized per replicate to generate 3 EESs. The average EES for the experimental condition was calculated and used for downstream analysis. To identify subpopulations within the published clusters, we manually examined a PHATE embedding of each sub-cluster, the distribution of EES values in each cluster, and the results of VFC subclustering with varying numbers of clusters. The decision to apply VFC was done on a per-cluster basis with the goal of identifying cell subpopulations with transcriptional similarity (as assessed by visualization) and uniform response to perturbation (as assessed by EES values). Cell types were annotated using sets of marker genes curated by Farrell et al. [19]. Changes in gene expression between VFC clusters was assessed using a rank sum test as implemented by `diffxpy`.

4.6 Generation, processing and analysis of the pancreatic islet datasets

Single-cell RNA-sequencing was performed on human islet cells from three different islet donors in the presence and absence of IFN γ . The islets were received on three different days. Cells were cultured for 24 hours with 25ng/mL IFN γ (R&D Systems) in CMRL 1066 medium (Gibco) and subsequently dissociated into single cells with 0.05% Trypsin EDTA (Gibco). Cells were then stained with FluoZin-3 (Invitrogen) and TMRE (Life Technologies) and sorted using a FACS Aria II (BD). The three samples were pooled for the sequencing. Cells were immediately processed using the 10X Genomics Chromium 3' Single-Cell RNA-sequencing kit at the Yale Center for Genome Analysis. The raw sequencing data was processed using the 10X Genomics Cell Ranger Pipeline. Raw data will be made available prior to publication.

Data from all three donors was concatenated into a single matrix for analysis. First, cells not expressing insulin, somatostatin, or glucagon were excluded from analysis using donor-specific thresholds. The data was square root transformed and reduced to 100 PCA dimensions. Next, we applied an MNN kernel to create a graph across all three donors with parameters `knn=5`, `decay=30`. This graph was then used for PHATE. The EES was calculated using MELD with default parameters. To identify coarse-grained cell types, we used previously published markers of islet cells [41]. We then used VFC to identify subpopulations of stimulated and unstimulated islet cells. To identify signature genes of IFN γ stimulation, we calculated differential expression between the clusters with the highest and lowest EES values within each cell type using a rank sum test as implemented in `diffxpy`. A consensus signature was then obtained by taking the intersection genes with `q-values < 0.05`. Gene set enrichment was then calculated using `gseapy`.

4.7 Quantitative comparisons

To generate single-cell data for the quantitative comparisons, we used Splatter. Datasets were all generated using the "Paths" mode so that a latent dimension in the data could be used to create the ground truth likelihood that each cell would be observed in the "experimental" condition relative to the "control". We focused on four data geometries: a tree with three branches, a branch and cluster with either end of the branch enriched or depleted and the cluster unaffected, a single branch with a middle section either enriched or depleted, and four clusters with random segments enriched or depleted. To create clusters, a multi-branched tree was created, and all but the tips of the branches were removed. The ground truth experimental signal was created using custom Python scripts taking the "Steps" latent variable from Splatter and randomly selecting a proportion of each branch or cluster between 10% and 80% of the data was enriched or depleted by 25%. These regions were divided into thirds to create a smooth transition between the unaffected regions and the differentially abundant regions. This likelihood ratio was then centered so that, on average, half the cells would be assigned to each condition. The centered ground truth signal was used to parameterize a

Bernoulli random variable and assign each cell to the experimental or control conditions. The data and RES were used as input to the respective algorithms.

To quantify the accuracy of the EES to approximate the ground truth likelihood ratio, we compared the EES, kNN-smoothed signal, or graph averaged signal to the ground truth likelihood of observing each cell in either of the two conditions. We used the Pearson's R statistic to calculate the degree to which these estimates approximate the likelihood ratio. Each of the four data geometries was tested 30 times with different random seeds for scRNA-seq simulation and RES generation.

We also performed EES comparisons using the T cell and zebrafish datasets described above. The preprocessed data was used to generate a three-dimensional PHATE embedding that was z-score normalized. We then used a combination of PHATE dimensions to create a ground truth probability each cell would be observed in the experimental or control condition. Cells were then assigned to either condition based on this probability as described above. We ran the same comparisons as on the simulated data with 100 random seeds per dataset.

To quantify the accuracy of VFC to detect the regions of the dataset that were enriched, depleted, or unaffected between conditions, we calculated the Adjusted Rand Score between the ground truth regions with enriched, depleted, or unchanged likelihood ratios between conditions. VFC was compared to k-Means, Spectral Clustering, Louvain, Leiden, and CellHarmony. As Leiden and Louvain do not provide a method to control the number of clusters, we implemented a binary search to identify a resolution parameter that provides the target number of clusters. Although Cell Harmony relies on an initial Louvain clustering, the tool does not implement Louvain with a tuneable resolution. It is also not possible to provide an initial clustering to CellHarmony, so we resorted to cutting Louvain at the level closest to our target number of clusters. Finally, because CellHarmony does not reconcile the disparate cluster assignments in the reference and query datasets, and because not all cells in the query dataset may be aligned to the reference we needed to generate manually new cluster labels for cells in the query dataset so that the method could be compared to other clustering tools.

To characterize the ability of MELD analysis using the EES and VFC to characterize gene signatures of a perturbation dataset, we returned to the T cell dataset. We again used the same setup to create synthetically 3 regions with different sampling probabilities in the dataset using PHATE clusters as above. Because one of these clusters has no differential abundance between conditions, we calculated the ground truth gene expression signature between the enriched and depleted clusters only using diffxpy [28]. To calculate the gene signature for each clustering method, we performed differential expression between the most enriched cluster in the experimental condition and the most depleted cluster in the experimental condition (or highest and lowest EES for MELD). We also considered directly performing two-sample comparison using the sample labels. To quantify the performance of each method, we used the area under the receiving operator characteristic (AUCROC) to compare the q-values produced using each method to the ground truth q-values. This process was repeated over 100 random seeds. The AUCROC curves and performance of each method relative to VFC is displayed in **Figure S6d,e**.

5 Data availability

Gene expression counts matrices prepared by Datlinger et al. [16] were accessed from the NCBI GEO database accession GSE92872. Gene expression counts matrices prepared by Wagner et al. [18] were downloaded from NCBI GEO accession GSE112294. The new pancreatic islets datasets will be made available on NCBI GEO prior to publication and this section will be revised to include the accession number.

6 Code availability

Code for the EES and VFC algorithms implemented in Python is available as part of the MELD package on GitHub <https://github.com/KrishnaswamyLab/MELD> and on the Python Package Index (PyPI). The GitHub repository also contains tutorials, code to reproduce the analysis of the zebrafish dataset, and code associated with several of the quantitative comparisons.

References

- [1] Malte D Luecken and Fabian J Theis. Current best practices in single-cell RNA-seq analysis: A tutorial. *Molecular Systems Biology*, 15(6):e8746, June 2019. ISSN 1744-4292. doi: 10.15252/msb.20188746.
- [2] Caleb Weinreb, Samuel Wolock, Allon M. Klein, and Bonnie Berger. SPRING: A kinetic interface for visualizing high dimensional single-cell expression data. *Bioinformatics*, 34(7):1246–1248, April 2018. ISSN 1367-4803. doi: 10.1093/bioinformatics/btx792.
- [3] Kevin R. Moon, David van Dijk, Zheng Wang, Scott Gigante, Daniel B. Burkhardt, William S. Chen, Kristina Yim, Antonia van den Elzen, Matthew J. Hirn, Ronald R. Coifman, Natalia B. Ivanova, Guy Wolf, and Smita Krishnaswamy. Visualizing structure and transitions in high-dimensional biological data. *Nature Biotechnology*, 37(12):1482–1492, December 2019. ISSN 1546-1696. doi: 10.1038/s41587-019-0336-3.
- [4] David van Dijk, Roshan Sharma, Juozas Nainys, Kristina Yim, Pooja Kathail, Ambrose J. Carr, Cassandra Burdziak, Kevin R. Moon, Christine L. Chaffer, Diwakar Pattabiraman, Brian Bieri, Linas Mazutis, Guy Wolf, Smita Krishnaswamy, and Dana Pe’er. Recovering Gene Interactions from Single-Cell Data Using Data Diffusion. *Cell*, 174(3):716–729.e27, July 2018. ISSN 0092-8674. doi: 10.1016/j.cell.2018.05.061.
- [5] Karthik Shekhar, Sylvain W. Lapan, Irene E. Whitney, Nicholas M. Tran, Evan Z. Macosko, Monika Kowalczyk, Xian Adiconis, Joshua Z. Levin, James Nemesh, Melissa Goldman, Steven A. McCarroll, Constance L. Cepko, Aviv Regev, and Joshua R. Sanes. Comprehensive Classification of Retinal Bipolar Neurons by Single-Cell Transcriptomics. *Cell*, 166(5):1308–1323.e30, August 2016. ISSN 0092-8674, 1097-4172. doi: 10.1016/j.cell.2016.07.054.
- [6] Jacob H. Levine, Erin F. Simonds, Sean C. Bendall, Kara L. Davis, El-ad D. Amir, Michelle D. Tadmor, Oren Litvin, Harris G. Fienberg, Astraea Jager, Eli R. Zunder, Rachel Finck, Amanda L. Gedman, Ina Radtke, James R. Downing, Dana Pe’er, and Garry P. Nolan. Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis. *Cell*, 162(1):184–197, July 2015. ISSN 0092-8674. doi: 10.1016/j.cell.2015.05.047.
- [7] Chen Xu and Zhengchang Su. Identification of cell types from single-cell transcriptomes using a novel clustering method. *Bioinformatics (Oxford, England)*, 31(12):1974–1980, June 2015. ISSN 1367-4811. doi: 10.1093/bioinformatics/btv088.
- [8] Etienne Becht, Leland McInnes, John Healy, Charles-Antoine Dutertre, Immanuel W. H. Kwok, Lai Guan Ng, Florent Ginhoux, and Evan W. Newell. Dimensionality reduction for visualizing single-cell data using UMAP. *Nature Biotechnology*, 37(1):38–44, January 2019. ISSN 1546-1696. doi: 10.1038/nbt.4314.

- [9] Laleh Haghverdi, Aaron T. L. Lun, Michael D. Morgan, and John C. Marioni. Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nature Biotechnology*, April 2018. ISSN 1546-1696. doi: 10.1038/nbt.4091.
- [10] Andrew Butler, Paul Hoffman, Peter Smibert, Efthymia Papalexi, and Rahul Satija. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nature Biotechnology*, 36(5):411–420, May 2018. ISSN 1546-1696. doi: 10.1038/nbt.4096.
- [11] Erica A. K. DePasquale, Daniel Schnell, Phillip Dexheimer, Kyle Ferchen, Stuart Hay, Kashish Chetal, Íñigo Valiente-Alandí, Burns C. Blaxall, H. Leighton Grimes, and Nathan Salomonis. cellHarmony: Cell-level matching and holistic comparison of single-cell transcriptomes. *Nucleic Acids Research*, 47(21):e138–e138, December 2019. ISSN 0305-1048. doi: 10.1093/nar/gkz789.
- [12] Anoop P. Patel, Itay Tirosh, John J. Trombetta, Alex K. Shalek, Shawn M. Gillespie, Hiroaki Wakimoto, Daniel P. Cahill, Brian V. Nahed, William T. Curry, Robert L. Martuza, David N. Louis, Orit Rozenblatt-Rosen, Mario L. Suvà, Aviv Regev, and Bradley E. Bernstein. Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science*, 344(6190):1396–1401, June 2014. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.1254257.
- [13] Greg Finak, Andrew McDavid, Masanao Yajima, Jingyuan Deng, Vivian Gersuk, Alex K. Shalek, Chloe K. Slichter, Hannah W. Miller, M. Juliana McElrath, Martin Prlic, Peter S. Linsley, and Raphael Gottardo. MAST: A flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biology*, 16, 2015. ISSN 1474-7596. doi: 10.1186/s13059-015-0844-5.
- [14] Itay Tirosh, Benjamin Izar, Sanjay M. Prakadan, Marc H. Wadsworth, Daniel Treacy, John J. Trombetta, Asaf Rotem, Christopher Rodman, Christine Lian, George Murphy, Mohammad Fallahi-Sichani, Ken Dutton-Regester, Jia-Ren Lin, Ofir Cohen, Parin Shah, Diana Lu, Alex S. Genshaft, Travis K. Hughes, Carly G. K. Ziegler, Samuel W. Kazer, Aleth Gaillard, Kellie E. Kolb, Alexandra-Chloé Villani, Cory M. Johannessen, Aleksandr Y. Andreev, Eliezer M. Van Allen, Monica Bertagnolli, Peter K. Sorger, Ryan J. Sullivan, Keith T. Flaherty, Dennie T. Frederick, Judit Jané-Valbuena, Charles H. Yoon, Orit Rozenblatt-Rosen, Alex K. Shalek, Aviv Regev, and Levi A. Garraway. Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq. *Science*, 352(6282):189–196, April 2016. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.aad0501.
- [15] Diego Adhemar Jaitin, Assaf Weiner, Ido Yofe, David Lara-Astiaso, Hadas Keren-Shaul, Eyal David, Tomer Meir Salame, Amos Tanay, Alexander van Oudenaarden, and Ido Amit. Dissecting Immune Circuits by Linking CRISPR-Pooled Screens with Single-Cell RNA-Seq. *Cell*, 167(7):1883–1896.e15, December 2016. ISSN 0092-8674. doi: 10.1016/j.cell.2016.11.039.
- [16] Paul Datlinger, André F. Rendeiro, Christian Schmidl, Thomas Krausgruber, Peter Traxler, Johanna Klughammer, Linda C. Schuster, Amelie Kuchler, Donat Alpar, and Christoph Bock. Pooled CRISPR screening with single-cell transcriptome readout. *Nature Methods*, January 2017. ISSN 1548-7091. doi: 10.1038/nmeth.4177.
- [17] Xin Gao, Deqing Hu, Madelaine Gogol, and Hua Li. ClusterMap: Comparing analyses across multiple Single Cell RNA-Seq profiles. *bioRxiv*, page 331330, June 2018. doi: 10.1101/331330.

- [18] Daniel E. Wagner, Caleb Weinreb, Zach M. Collins, James A. Briggs, Sean G. Megason, and Allon M. Klein. Single-cell mapping of gene expression landscapes and lineage in the zebrafish embryo. *Science*, page eaar4362, April 2018. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.aar4362.
- [19] Jeffrey A. Farrell, Yiqun Wang, Samantha J. Riesenfeld, Karthik Shekhar, Aviv Regev, and Alexander F. Schier. Single-cell reconstruction of developmental trajectories during zebrafish embryogenesis. *Science*, 360(6392):eaar3131, June 2018. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.aar3131.
- [20] Kevin R. Moon, Jay S. Stanley, Daniel Burkhardt, David van Dijk, Guy Wolf, and Smita Krishnaswamy. Manifold learning-based methods for analyzing single-cell RNA-sequencing data. *Current Opinion in Systems Biology*, 7:36–46, February 2018. ISSN 2452-3100. doi: 10.1016/j.coisb.2017.12.008.
- [21] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst. The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains. *IEEE Signal Processing Magazine*, 30(3):83–98, May 2013. ISSN 1053-5888. doi: 10.1109/MSP.2012.2235192.
- [22] Mikhail Belkin and Partha Niyogi. Convergence of Laplacian Eigenmaps. In B. Schölkopf, J. C. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, pages 129–136, 2006.
- [23] Z. I. Botev, J. F. Grotowski, and D. P. Kroese. Kernel density estimation via diffusion. *Annals of Statistics*, 38(5):2916–2957, October 2010. ISSN 0090-5364, 2168-8966. doi: 10.1214/10-AOS799.
- [24] David I Shuman, Pierre Vandergheynst, and Pascal Frossard. Chebyshev polynomial approximation for distributed signal processing. In *Distributed Computing in Sensor Systems and Workshops (DCOSS), 2011 International Conference on*, pages 1–8. IEEE, 2011.
- [25] David I Shuman, Benjamin Ricaud, and Pierre Vandergheynst. Vertex-frequency analysis on graphs. *Applied and Computational Harmonic Analysis*, 40(2):260–291, March 2016. ISSN 1063-5203. doi: 10.1016/j.acha.2015.02.005.
- [26] Luke Zappia, Belinda Phipson, and Alicia Oshlack. Splatter: Simulation of single-cell RNA sequencing data. *Genome Biology*, 18(1):174, September 2017. ISSN 1474-760X. doi: 10.1186/s13059-017-1305-0.
- [27] V. A. Traag, L. Waltman, and N. J. van Eck. From Louvain to Leiden: Guaranteeing well-connected communities. *Scientific Reports*, 9(1):5233, March 2019. ISSN 2045-2322. doi: 10.1038/s41598-019-41695-z.
- [28] David Fischer. Theislab/diffxpy. Theis Lab, June 2020.
- [29] Maxim V. Kuleshov, Matthew R. Jones, Andrew D. Rouillard, Nicolas F. Fernandez, Qiaonan Duan, Zichen Wang, Simon Koplev, Sherry L. Jenkins, Kathleen M. Jagodnik, Alexander Lachmann, Michael G. McDermott, Caroline D. Monteiro, Gregory W. Gundersen, and Avi Ma’ayan. Enrichr: A comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Research*, 44 (Web Server issue):W90–W97, July 2016. ISSN 0305-1048. doi: 10.1093/nar/gkw377.

- [30] M. Hammerschmidt, F. Pelegri, M. C. Mullins, D. A. Kane, F. J. van Eeden, M. Granato, M. Brand, M. Furutani-Seiki, P. Haffter, C. P. Heisenberg, Y. J. Jiang, R. N. Kelsh, J. Odenthal, R. M. Warga, and C. Nusslein-Volhard. Dino and mercedes, two genes regulating dorsal development in the zebrafish embryo. *Development*, 123(1):95–102, December 1996. ISSN 0950-1991, 1477-9129.
- [31] Stefan Schulte-Merker, Kevin J. Lee, Andrew P. McMahon, and Matthias Hammerschmidt. The zebrafish organizer requires *chordino*. *Nature*, 387(6636):862–863, June 1997. ISSN 1476-4687. doi: 10.1038/43092.
- [32] Shannon Fisher and Marnie E. Halpern. Patterning the zebrafish axial skeleton requires early *chordin* function. *Nature Genetics*, 23(4):442–446, December 1999. ISSN 1546-1718. doi: 10.1038/70557.
- [33] Shuo-Ting Yen, Min Zhang, Jian Min Deng, Shireen J. Usman, Chad N. Smith, Jan Parker-Thornburg, Paul G. Swinton, James F. Martin, and Richard R. Behringer. Somatic mosaicism and allele complexity induced by CRISPR/Cas9 RNA injections in mouse zygotes. *Developmental Biology*, 393(1):3–9, September 2014. ISSN 0012-1606. doi: 10.1016/j.ydbio.2014.06.017.
- [34] Anskar Y. H. Leung, Eric M. Mendenhall, Tommy T. F. Kwan, Raymond Liang, Craig Eckfeldt, Eleanor Chen, Matthias Hammerschmidt, Suzanne Grindley, Stephen C. Ekker, and Catherine M. Verfaillie. Characterization of expanded intermediate cell mass in zebrafish chordin morphant embryos. *Developmental Biology*, 277(1):235–254, January 2005. ISSN 0012-1606. doi: 10.1016/j.ydbio.2004.09.032.
- [35] Ben Steventon, Claudio Araya, Claudia Linker, Sei Kuriyama, and Roberto Mayor. Differential requirements of BMP and Wnt signalling during gastrulation and neurulation define two steps in neural crest induction. *Development*, 136(5):771–779, March 2009. ISSN 0950-1991, 1477-9129. doi: 10.1242/dev.029017.
- [36] Carolin Schille and Alexandra Schambony. Signaling pathways and tissue interactions in neural plate border formation. *Neurogenesis*, 4(1), February 2017. ISSN 2326-2133. doi: 10.1080/23262133.2017.1292783.
- [37] Anastasia Katsarou, Soffia Gudbjörnsdottir, Araz Rawshani, Dana Dabelea, Ezio Bonifacio, Barbara J. Anderson, Laura M. Jacobsen, Desmond A. Schatz, and Åke Lernmark. Type 1 diabetes mellitus. *Nature Reviews Disease Primers*, 3:17016, March 2017. ISSN 2056-676X. doi: 10.1038/nrdp.2017.16.
- [38] V. Ablamunits, D. Elias, T. Reshef, and I. R. Cohen. Islet T cells secreting IFN- γ in NOD mouse diabetes: Arrest by p277 peptide treatment. *Journal of Autoimmunity*, 11(1):73–81, February 1998. ISSN 0896-8411. doi: 10.1006/jaut.1997.0177.
- [39] Andrew S. Diamond and Ronald G. Gill. An Essential Contribution by IFN- γ to CD8+ T Cell-Mediated Rejection of Pancreatic Islet Allografts. *The Journal of Immunology*, 165(1):247–255, July 2000. ISSN 0022-1767, 1550-6606. doi: 10.4049/jimmunol.165.1.247.
- [40] Miguel Lopes, Burak Kutlu, Michela Miani, Claus H. Bang-Berthelsen, Joachim Størling, Flemming Pociot, Nathan Goodman, Lee Hood, Nils Welsh, Gianluca Bontempi, and Decio L. Eizirik. Temporal profiling of cytokine-induced genes in pancreatic β -cells by meta-analysis and network inference. *Genomics*, 103(4):264–275, April 2014. ISSN 0888-7543. doi: 10.1016/j.ygeno.2013.12.007.

- [41] Mauro J. Muraro, Gitanjali Dharmadhikari, Dominic Grün, Nathalie Groen, Tim Dielen, Erik Jansen, Leon van Gurp, Marten A. Engelse, Francoise Carlotti, Eelco J.P. de Koning, and Alexander van Oudenaarden. A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Systems*, 3(4):385–394.e3, October 2016. ISSN 2405-4712. doi: 10.1016/j.cels.2016.09.002.
- [42] Yurong Xin, Giselle Dominguez Gutierrez, Haruka Okamoto, Jinrang Kim, Ann-Hwee Lee, Christina Adler, Min Ni, George D. Yancopoulos, Andrew J. Murphy, and Jesper Gromada. Pseudotime Ordering of Single Human β -Cells Reveals States of Insulin Production and Unfolded Protein Response. *Diabetes*, 67(9):1783–1794, September 2018. ISSN 0012-1797, 1939-327X. doi: 10.2337/db18-0365.
- [43] Lydia Farack, Matan Golan, Adi Egozi, Nili Dezorella, Keren Bahar Halpern, Shani Ben-Moshe, Immacolata Garzilli, Beáta Tóth, Lior Roitman, Valery Krizhanovsky, and Shalev Itzkovitz. Transcriptional Heterogeneity of Beta Cells in the Intact Pancreas. *Developmental Cell*, 48(1):115–125.e4, January 2019. ISSN 1534-5807. doi: 10.1016/j.devcel.2018.11.001.
- [44] Chilakamarti V Ramana, M. Pilar Gil, Robert D Schreiber, and George R Stark. Stat1-dependent and -independent pathways in IFN- γ -dependent signaling. *Trends in Immunology*, 23(2):96–101, February 2002. ISSN 1471-4906. doi: 10.1016/S1471-4906(01)02118-4.
- [45] Anthony J. Sadler and Bryan R. G. Williams. Interferon-inducible antiviral effectors. *Nature reviews. Immunology*, 8(7):559–568, July 2008. ISSN 1474-1733. doi: 10.1038/nri2314.
- [46] Katherine A. Fitzgerald. The Interferon Inducible Gene: Viperin. *Journal of Interferon & Cytokine Research*, 31(1):131–135, January 2011. ISSN 1079-9907. doi: 10.1089/jir.2010.0127.
- [47] Zhiwei Zheng, Lin Wang, and Jihong Pan. Interferon-stimulated gene 20-kDa protein (ISG20) in infection and disease: Review and outlook. *Intractable & Rare Diseases Research*, 6(1):35–40, February 2017. ISSN 2186-3644. doi: 10.5582/irdr.2017.01004.
- [48] Monica Hultcrantz, Michael H. Hühn, Monika Wolf, Annika Olsson, Stella Jacobson, Bryan R. Williams, Olle Korsgren, and Malin Flodström-Tullberg. Interferons induce an antiviral state in human pancreatic islet cells. *Virology*, 367(1):92–101, October 2007. ISSN 0042-6822. doi: 10.1016/j.virol.2007.05.010.
- [49] Andrew F. Stewart, Mehboob A. Hussain, Adolfo García-Ocaña, Rupangi C. Vasavada, Anil Bhushan, Ernesto Bernal-Mizrachi, and Rohit N. Kulkarni. Human β -Cell Proliferation and Intracellular Signaling: Part 3. *Diabetes*, 64(6):1872–1885, June 2015. ISSN 0012-1797. doi: 10.2337/db14-1843.
- [50] Xinyue Chen, Daniel B. Burkhardt, Amaleah A. Hartman, Xiao Hu, Anna E. Eastman, Chao Sun, Xujun Wang, Mei Zhong, Smita Krishnaswamy, and Shangqin Guo. MLL-AF9 initiates transformation from fast-proliferating myeloid progenitors. *Nature Communications*, 10(1), December 2019. doi: 10.1038/s41467-019-13666-5.
- [51] Emily V. Dutrow, Deena Emera, Kristina Yim, Severin Uebbing, Acadia A. Kocher, Martina Krenzer, Timothy Nottoli, Daniel B. Burkhardt, Smita Krishnaswamy, Angeliki Louvi, and James P. Noonan. The Human Accelerated Region HACNS1 modifies developmental gene expression in humanized mice. *bioRxiv*, page 2019.12.11.873075, December 2019. doi: 10.1101/2019.12.11.873075.

- [52] Katherine E. Savell, Jennifer J. Tuscher, Morgan E. Zipperly, Corey G. Duke, Robert A. Phillips, Allison J. Bauman, Saakshi Thukral, Faraz A. Sultan, Nicholas A. Goska, Lara Ianov, and Jeremy J. Day. A dopamine-induced gene expression signature regulates neuronal function and cocaine response. *Science Advances*, 6(26):eaba4221, June 2020. ISSN 2375-2548. doi: 10.1126/sciadv.aba4221.
- [53] Katherine Minjee Chung, Jaffarguriqbal Singh, Lauren Lawres, Kimberly Judith Dorans, Cathy Garcia, Daniel B. Burkhardt, Rebecca Robbins, Arjun Bhutkar, Rebecca Cardone, Xiaojian Zhao, Ana Babic, Sara A. Vayrynen, Andressa Dias Costa, Jonathan A. Nowak, Daniel T. Chang, Richard F. Dunne, Aram F. Hezel, Albert C. Koong, Joshua J. Wilhelm, Melena D. Bellin, Vibe Nylander, Anna L. Gloyn, Mark I. McCarthy, Richard G. Kibbey, Smita Krishnaswamy, Brian M. Wolpin, Tyler Jacks, Charles S. Fuchs, and Mandar Deepak Muzumdar. Endocrine-Exocrine Signaling Drives Obesity-Associated Pancreatic Ductal Adenocarcinoma. *Cell*, 181(4):832–847.e18, May 2020. ISSN 0092-8674. doi: 10.1016/j.cell.2020.03.062.
- [54] Neal G. Ravindra, Mia Madel Alfajaro, Victor Gasque, Victoria Habet, Jin Wei, Renata B. Filler, Nicholas C. Huston, Han Wan, Klara Szigeti-Buck, Bao Wang, Guilin Wang, Ruth R. Montgomery, Stephanie C. Eisenbarth, Adam Williams, Anna Marie Pyle, Akiko Iwasaki, Tamas L. Horvath, Ellen F. Foxman, Richard W. Pierce, David van Dijk, and Craig B. Wilen. Single-cell longitudinal analysis of SARS-CoV-2 infection in human airway epithelium. *bioRxiv*, July 2020. doi: 10.1101/2020.05.06.081695.
- [55] Ronald R. Coifman and Stéphane Lafon. Diffusion maps. *Applied and Computational Harmonic Analysis*, 21(1):5–30, July 2006. ISSN 1063-5203. doi: 10.1016/j.acha.2006.04.006.
- [56] Y. P Mack and M Rosenblatt. Multivariate k-nearest neighbor density estimates. *Journal of Multivariate Analysis*, 9(1):1–15, March 1979. ISSN 0047-259X. doi: 10.1016/0047-259X(79)90065-4.
- [57] Gérard Biau, Frédéric Chazal, David Cohen-Steiner, Luc Devroye, and Carlos Rodríguez. A weighted k-nearest neighbor density estimate for geometric inference. *Electronic Journal of Statistics*, 5:204–237, 2011. ISSN 1935-7524. doi: 10.1214/11-EJS606.
- [58] Yi-Hung Kung, Pei-Sheng Lin, and Cheng-Hsiung Kao. An optimal k-nearest neighbor for density estimation. *Statistics & Probability Letters*, 82(10):1786–1791, October 2012. ISSN 0167-7152. doi: 10.1016/j.spl.2012.05.017.
- [59] Ulrike Von Luxburg and Morteza Alamgir. Density estimation from unweighted k-nearest neighbor graphs: A roadmap. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 225–233. Curran Associates, Inc., 2013.
- [60] Bernard W. Silverman. *Density Estimation for Statistics and Data Analysis*. Routledge, February 2018. ISBN 978-1-315-14091-9. doi: 10.1201/9781315140919.
- [61] David K Hammond, Pierre Vandergheynst, and Rémi Gribonval. Wavelets on graphs via spectral graph theory. *Applied and Computational Harmonic Analysis*, 30(2):129–150, 2011.
- [62] Nathanael Perraudin, Benjamin Ricaud, David Shuman, and Pierre Vandergheynst. Global and local uncertainty principles for signals on graphs. *arXiv preprint arXiv:1603.03030*, 2016.

- 1257 [63] Stephane Mallat. *A Wavelet Tour of Signal Processing: The Sparse Way*. Academic Press, December
1258 2008. ISBN 978-0-08-092202-7.
- 1259 [64] Dengyong Zhou and Bernhard Schölkopf. A regularization framework for learning from graph data.
1260 In *ICML workshop on statistical relational learning and Its connections to other fields*, volume 15,
1261 pages 67–8, 2004.
- 1262 [65] Jihun Ham, Daniel D Lee, and Lawrence K Saul. Semisupervised alignment of manifolds. In *AISTATS*,
1263 pages 120–127, 2005.
- 1264 [66] Mikhail Belkin, Irina Matveeva, and Partha Niyogi. Regularization and semi-supervised learning
1265 on large graphs. In *International Conference on Computational Learning Theory*, pages 624–638.
1266 Springer, 2004.
- 1267 [67] Rie K. Ando and Tong Zhang. Learning on Graph with Laplacian Regularization. In B. Schölkopf,
1268 J. C. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, pages
1269 25–32. MIT Press, 2007.
- 1270 [68] Kilian Q Weinberger, Fei Sha, Qihui Zhu, and Lawrence K. Saul. Graph Laplacian Regularization
1271 for Large-Scale Semidefinite Programming. In B. Schölkopf, J. C. Platt, and T. Hoffman, editors,
1272 *Advances in Neural Information Processing Systems 19*, pages 1489–1496. MIT Press, 2007.
- 1273 [69] X. He, M. Ji, C. Zhang, and H. Bao. A Variance Minimization Criterion to Feature Selection Using
1274 Laplacian Regularization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(10):
1275 2013–2025, October 2011. ISSN 0162-8828. doi: 10.1109/TPAMI.2011.44.
- 1276 [70] X. Liu, D. Zhai, D. Zhao, G. Zhai, and W. Gao. Progressive Image Denoising Through Hybrid Graph
1277 Laplacian Regularization: A Unified Framework. *IEEE Transactions on Image Processing*, 23(4):
1278 1491–1503, April 2014. ISSN 1057-7149. doi: 10.1109/TIP.2014.2303638.
- 1279 [71] J. Pang, G. Cheung, A. Ortega, and O. C. Au. Optimal graph laplacian regularization for natural
1280 image denoising. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing*
1281 *(ICASSP)*, pages 2294–2298, April 2015. doi: 10.1109/ICASSP.2015.7178380.
- 1282 [72] Jiahao Pang and Gene Cheung. Graph Laplacian Regularization for Image Denoising: Analysis in the
1283 Continuous Domain. *IEEE Transactions on Image Processing*, 26(4):1770–1785, April 2017. ISSN
1284 1057-7149, 1941-0042. doi: 10.1109/TIP.2017.2651400.
- 1285 [73] Nathanaël Perraudin, Johan Paratte, David Shuman, Lionel Martin, Vassilis Kalofolias, Pierre Van-
1286 dergheynst, and David K. Hammond. GSPBOX: A toolbox for signal processing on graphs. *ArXiv*
1287 *e-prints*, August 2014.
- 1288 [74] Martin Barron and Jun Li. Identifying and removing the cell-cycle effect from single-cell rna-
1289 sequencing data. *Scientific reports*, 6:33892, 2016.
- 1290 [75] Ronald R Coifman and Mauro Maggioni. Diffusion wavelets. *Applied and Computational Harmonic*
1291 *Analysis*, 21(1):53–94, 2006.
- 1292 [76] Probal Chaudhuri and J. S. Marron. Scale Space View of Curve Estimation. *The Annals of Statistics*,
1293 28(2):408–428, 2000. ISSN 0090-5364.

- [77] Nathanaël Perraudin, Nicki Holighaus, Peter L Søndergaard, and Peter Balazs. Designing Gabor windows using convex optimization. *arXiv preprint arXiv:1401.6033*, 2014.
- [78] Andrew Y Ng, Michael I Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems*, pages 849–856, 2002.
- [79] Maren Büttner, Zhichao Miao, F. Alexander Wolf, Sarah A. Teichmann, and Fabian J. Theis. A test metric for assessing single-cell RNA-seq batch correction. *Nature Methods*, 16(1):43–49, January 2019. ISSN 1548-7105. doi: 10.1038/s41592-018-0254-1.
- [80] Robert A. Amezcua, Aaron T. L. Lun, Etienne Becht, Vince J. Carey, Lindsay N. Carpp, Ludwig Geistlinger, Federico Marini, Kevin Rue-Albrecht, Davide Risso, Charlotte Soneson, Levi Waldron, Hervé Pagès, Mike L. Smith, Wolfgang Huber, Martin Morgan, Raphael Gottardo, and Stephanie C. Hicks. Orchestrating single-cell analysis with Bioconductor. *Nature Methods*, 17(2):137–145, February 2020. ISSN 1548-7105. doi: 10.1038/s41592-019-0654-x.
- [81] M. D. Luecken, M. Büttner, K. Chaichoompu, A. Danese, M. Interlandi, M. F. Mueller, D. C. Strobl, L. Zappia, M. Dugas, M. Colomé-Tatché, and F. J. Theis. Benchmarking atlas-level data integration in single-cell genomics. *bioRxiv*, page 2020.05.22.111161, May 2020. doi: 10.1101/2020.05.22.111161.

7 Supplementary Notes

7.1 Applying MELD analysis to single cell datasets with a batch effect

When jointly analyzing single cell datasets collected in different samples, difficulty may arise due to systematic changes in gene expression profiles between biologically equivalent cells [79]. These changes may be technical in nature (e.g. differences in the reverse transcription efficiency during library preparation) or biological (e.g. changes in sample preparation cause unexpected changes in biological state of otherwise equivalent cells). Regardless of the cause, the unifying feature of batch effects is that they confound the analysis a given research wants to perform. As such, it is unsurprising that dozens of batch normalization tools have been developed for single cell data [81]. However, it is important to emphasize that what constitutes a batch effect is dependent on the biological question in which a researcher is interested. Some analysts might be uninterested in variation caused by a change in cell media composition between samples, but other researchers might want to study these differences. Batch normalization tools have no way to know what variation is biologically relevant to the specific hypotheses of a given experiment and thus risk removing meaningful experimental signal when “correcting” measured values. This is problematic for analysis using MELD, because the goal of the toolkit is to quantify the differences that exist between samples without regard for the specific interests of given hypothesis. As such, we do not recommend using batch correction along the experimental axis (i.e. between experimental and control conditions) before running MELD. However, recognizing that in some cases batch correction is essential, we describe several considerations for performing MELD analysis on batch-corrected data.

For the EES algorithm to accurately estimate conditional probability of each sample, we assume that the graph learned from single cell data approximates the underlying cell state manifold. In **Section 20** we describe the use of an anisotropic kernel that normalizes for varying sampling density across cell states. However, some batch correction methods, such as mutual nearest neighbors [9], rely on the construction of a graph with artificially inflated weights between nodes from different samples. This graph no longer

models the cell states an experiment measured, but rather enforces similarities between cells based on the heuristic of the chosen normalization model. We provide no theoretical guarantees that a graph learned from batch corrected data will accurately model the underlying probability densities of each condition.

In practice when analyzing islet cells collected from multiple donors, that applying batch correction methods across the donor label improves our ability to capture a signal of IFN γ stimulation (**Section 2.8**). It is important to note that in this case, batch correction applied to a label that is orthogonal to the experimental axis. We have not examined the accuracy of the EES algorithm when batch correction is applied between experimental and control samples, although it is our expectation that this will likely remove biological signal. We recommend any user considering applying batch correction methods prior to running MELD analysis follow these steps:

1. To determine if a batch effect exists, confirm that cells from one sample are not finding appropriate neighbors in another following the strategy outlined by Büttner et al. [79].
2. To characterize the effect, identify which genes change the most between the samples
3. Confirm that the genes that are different are not relevant to the biological question under investigation
4. Apply batch correction
5. Confirm that relevant biological differences are still present using MELD analysis
6. If the biological differences are not present, repeat from step 1 with less batch correction. If you hit your personal recursion limit, consider that you don't actually want to do batch correction
7. If biological differences are present, then confirm that previous batch effect has been corrected and proceed to downstream analysis

7.2 Parameter search for the EES algorithm

To determine the optimal set of parameters for the EES algorithm, we performed a parameter search using splatter-generated datasets. For each of the four dataset structures, we generated 10 datasets with different random seeds and 10 random ground-truth EES per dataset for a total of 400 datasets per combination of parameters. A coarse-grained grid search revealed that setting $\alpha = 0$ and $\rho = 1$ performed best regardless of the β parameter. This is expected because with these settings, the MELD filter is the standard heat kernel. A fine-grained search over parameters for β showed that optimal values were between 50-75 (**Figure S13**). We chose a value of 60 as the default in the MELD toolkit and this was used for all experiments. We would like to note that the optimal β parameter will vary with dataset structure and the number of cells. **Figure S13b** shows how the optimal β values varies as a function of the number of cells generated using splatter while keeping the underlying geometry the same.

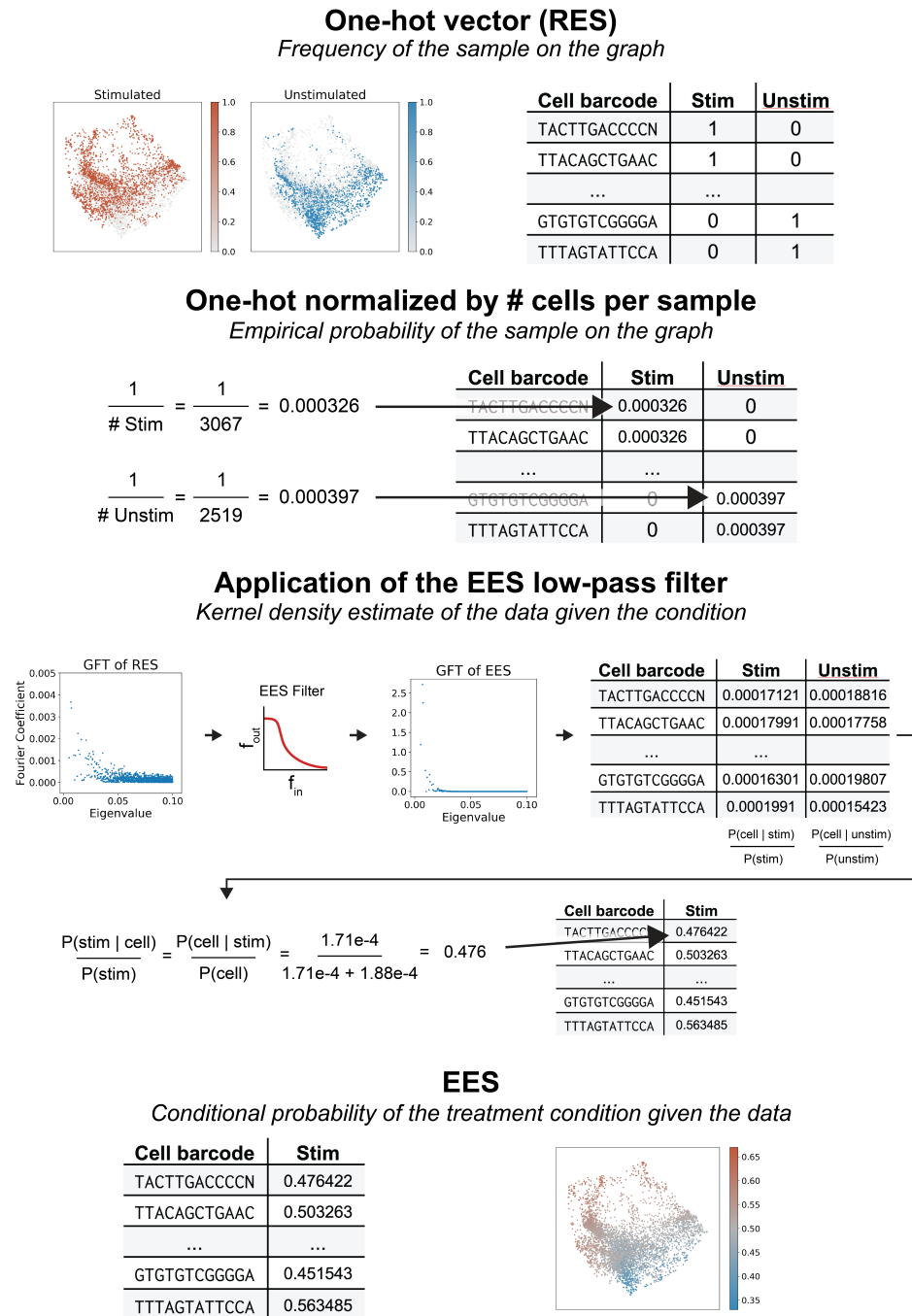


Figure S1: A step-by-step visual representation of the EES algorithm using data from Datlinger et al. [16]. The sample labels are used to create a one-hot indicator vector for each condition. These one-hot vectors are then column-wise L1-normalized such that the sum of each vector is 1. This gives each sample equal weight over the manifold despite a potential uneven number of cells in each condition. Next, the EES filter is used to calculate a kernel density estimate for each condition. These density estimates are then row-wise L1-normalized to yield the conditional probability that each cell would be observed in each condition. The conditional probability of the experimental condition relative to the control is used as the EES for two-condition experiments.

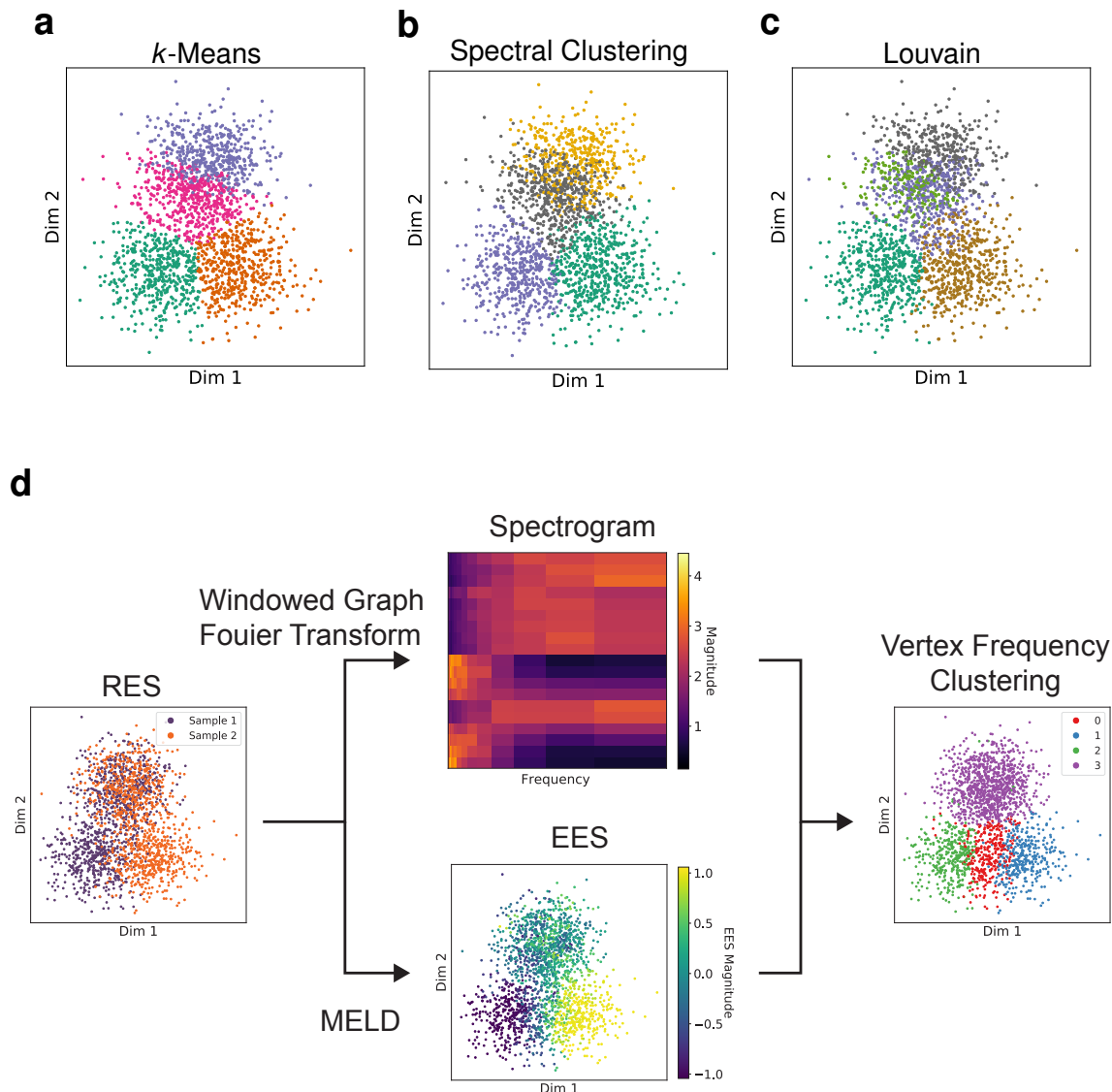


Figure S2: Vertex-Frequency clustering with MELD. A Gaussian mixture model was used to generate $N = 1000$ points in a mixture of three Gaussian distributions. This experiment is representative of a two-cell type experiment (split by Dim 2) in which one sample changes (bottom clusters) along Dim 1 due to the experiment while the other remains mixed (top clusters). **(a)** *k*-Means clustering separates the left and right experimental groups but splits the upper group erroneously. **(b)** Spectral clustering replicates the performance of *k*-Means in this example. **(c)** Louvain modularity clustering splits the mixture into five groups, with the same lower separations as before but with three groups in the upper cell type. **(d)** Vertex-Frequency clustering recovers a new cluster type. Briefly, the RES (left) is used for (1) a windowed graph Fourier Transform to obtain vertex-frequency information (above, logarithmically downsampled for clarity) and (2) MELD, which generates a continuous profile of the simulated experimental effect. These measures are concatenated together and clustered with *k*-Means. The clusters (right) separate the two cell types (purple and green/red/blue), and finds a separate grouping of cells that are in transition from green to blue, shown in red. One may see that in the spectrogram the green and blue groups are found on relatively low frequency patterns (bottom half of spectrogram, mostly black bands), whereas the medium frequency transition is well separated (middle of bottom bands). The well-mixed, nonresponsive population is entirely high frequency (top half).

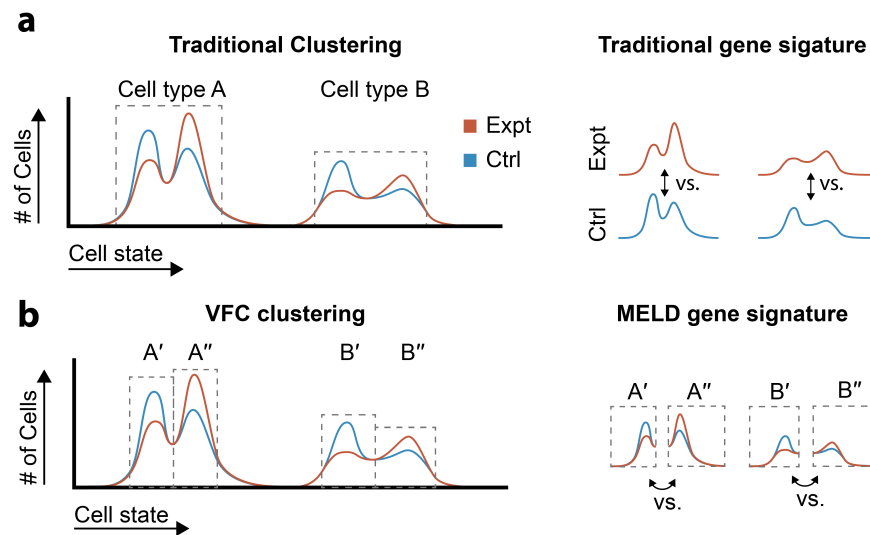


Figure S3: Identifying gene signatures using the EES and VFC. **(a)** In traditional gene signature analysis, clusters are identified based on data geometry and may not capture subpopulations of cells with varying response to a perturbation. In this framework, gene signatures are calculated by comparing cells from the experimental and control condition within each cluster. **(b)** To identify gene signatures of a perturbation with the MELD toolkit, we propose first partitioning cell populations with divergent responses to an experimental perturbation prior to differential expression analysis. We then assume that the differences within each VFC cluster is noise. Differential expression can either be calculated between subclusters identified by VFC (as shown) or by comparing each VFC cluster to the rest of the dataset independently.

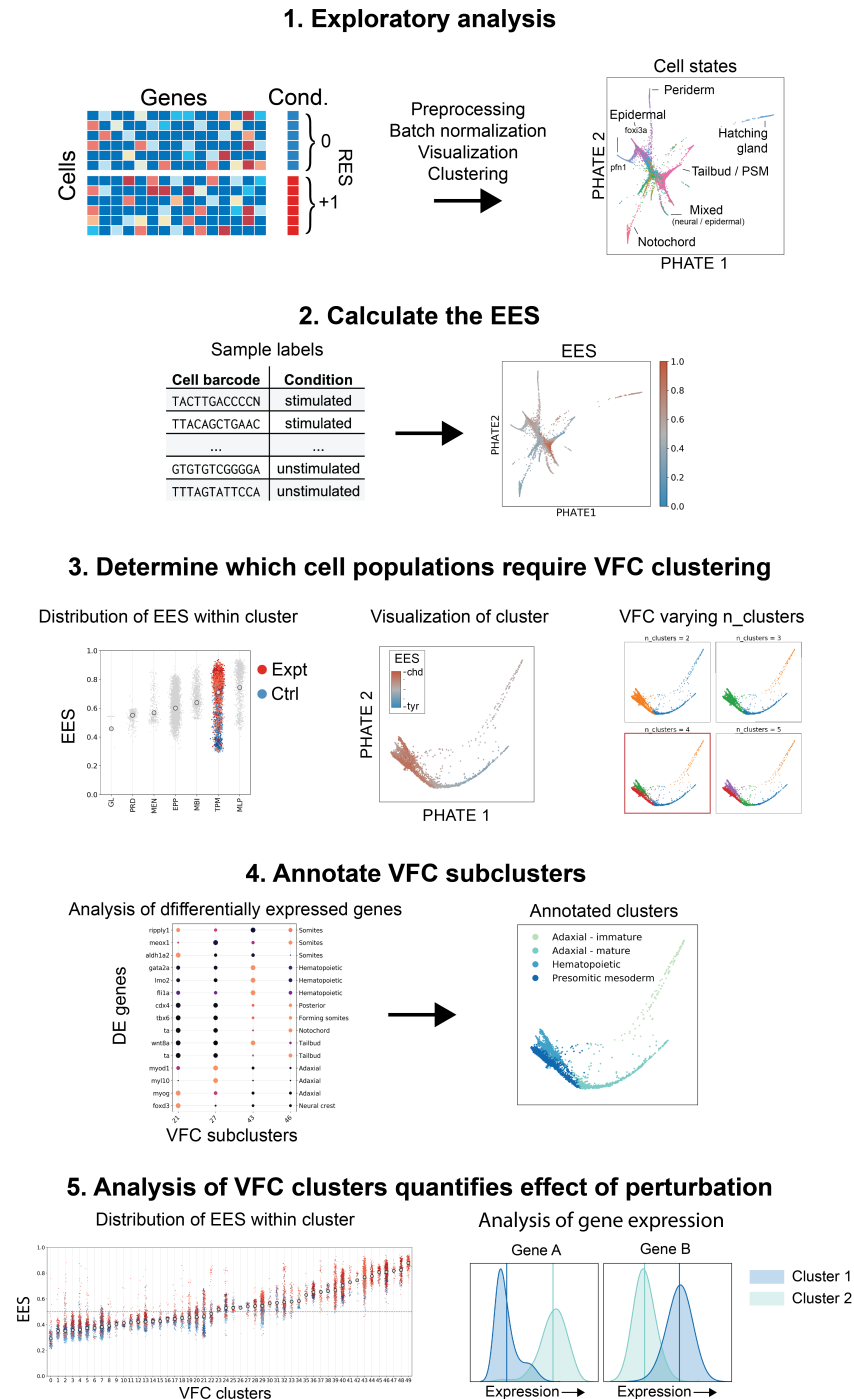


Figure S4: Overview of a pipeline for single cell analysis using MELD. (1.) Initial exploratory analysis of the dataset should follow established best practices to identify coarse-grained cell populations [1, 80]. (2.) Calculating the EES provides a measure for each cell describing the probability that cell would be observed in the experimental condition relative to the control. (3.) To identify populations most affected by a perturbation, we consider several sources of information regarding biological heterogeneity and the effect of the perturbation within each exploratory cluster. We then apply VFC at the determined cluster resolution. (4.) To assess the biological relevance of each VFC cluster, standard methods for cluster annotation can be applied. (5.) To characterize the gene signature of the perturbation, we compare expression differences between VFC clusters with varying EES distributions.

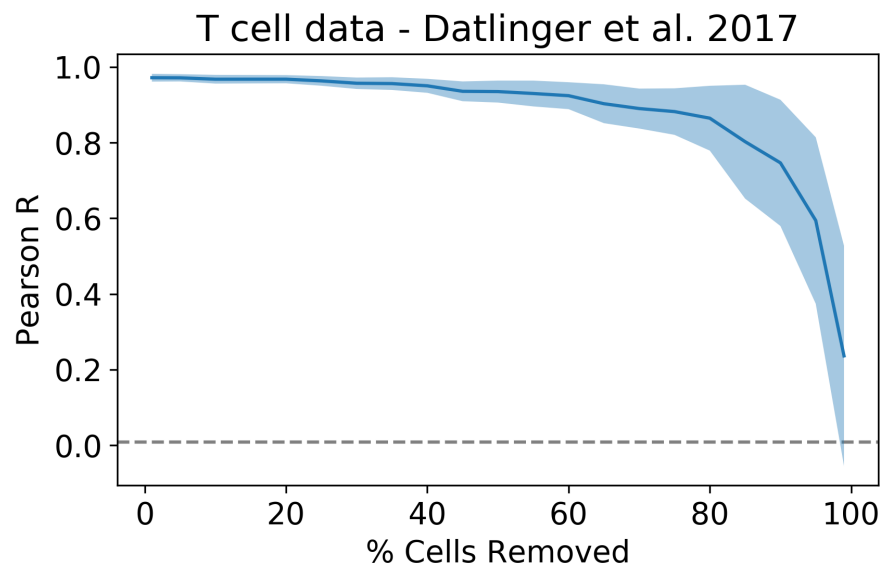


Figure S5: Result of down-sampling on accurately recovering simulated EES values. Using the procedure described in **Section 4.7**, we generated 100 random ground truth EES values and then removed between 1-99% of the cells in the dataset before running the EES algorithm normally. The average Pearson's R is shown as a function of the number of cells removed prior to running the EES algorithm. The shaded area demarks ± 1 standard deviation. We observe an average correlation > 0.9 for all experiments with at least 35% of the data present, or 1956 out of 5591 cells.

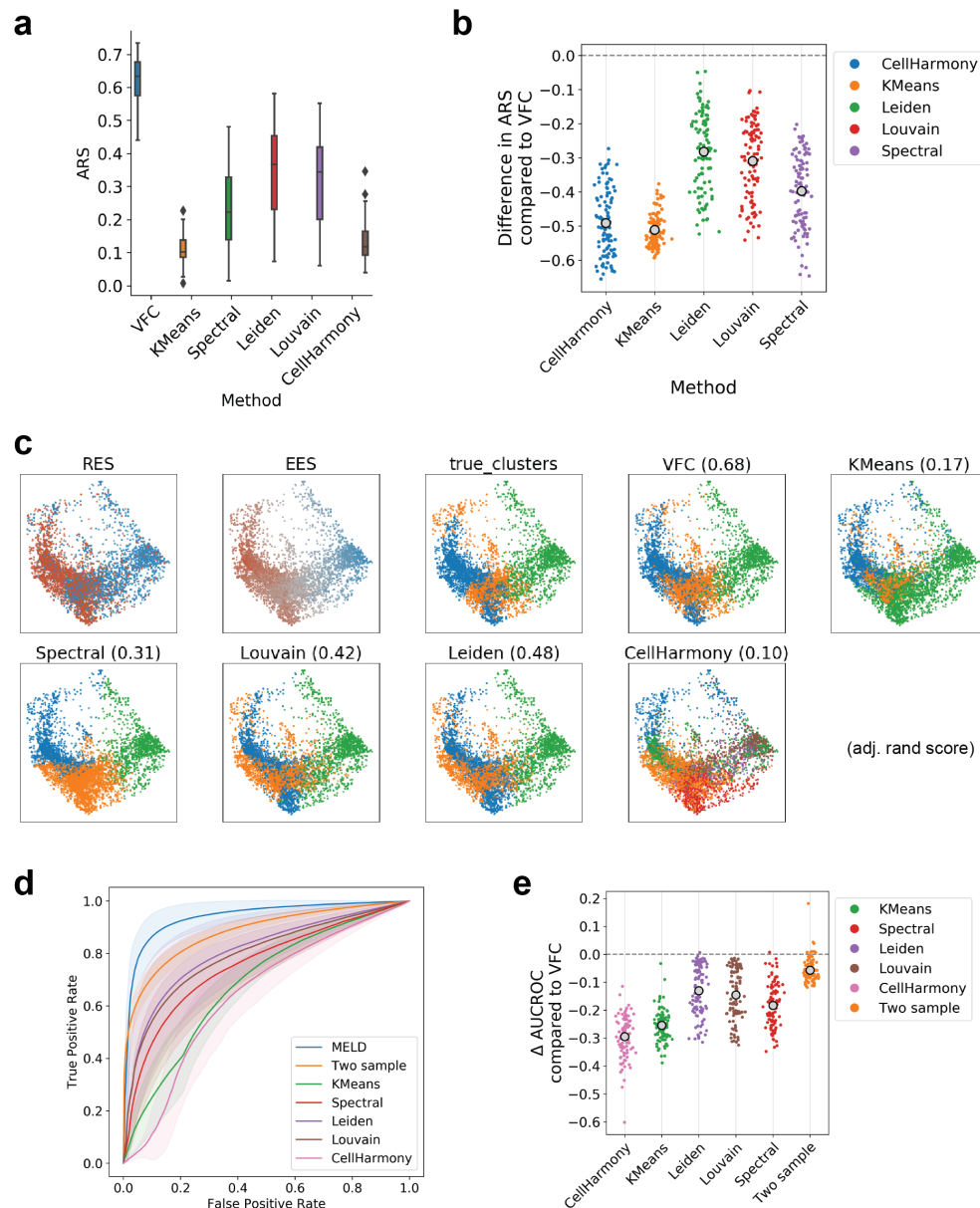


Figure S6: VFC accurately identifies cell populations affected by a perturbation in T cell data from Datlinger et al. [16]. **(a)** To create ground truth clusters, we artificially enriched and depleted various cell populations in either the experimental or control condition. Here we show the Adjusted Rand Score (ARS) over 100 simulations for 6 methods. For ARS, values close to 1 indicate perfect correspondence with ground truth, and values close to 0 indicate random labelling. VFC is the top performing method. **(b)** Because each simulation produced varying ARS scores for each method due to random seeds, we also consider the difference on performance between each method and VFC on each simulation. In none of 100 random seeds did any method outperform VFC. **(c)** The sample labels, EES, and clustering results for one randomly selected simulation. **(d)** Receiver operating characteristic (ROC) curves for the gene expression signatures described in Section 4.7. The Area Under the Curve of the ROC (AUCROC) indicates the overall performance of each strategy for identifying a gene signature. MELD is the top performing approach followed by direct comparison of the two samples. **(e)** As above, we consider the difference in AUCROC over each of 100 simulations between MELD and each method. In only 4 simulations does another method outperform MELD by more than 0.01.

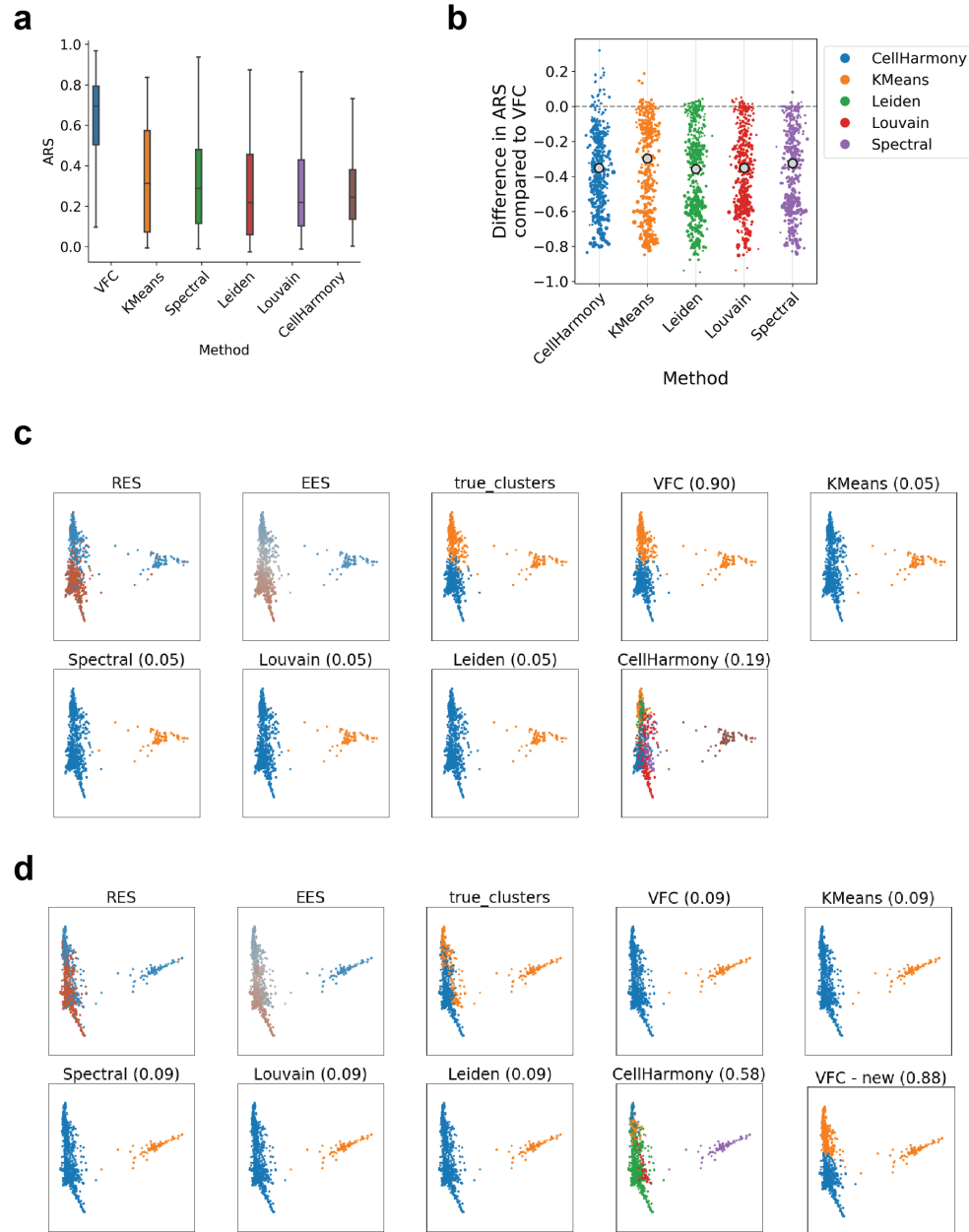


Figure S7: Quantitative comparison of clustering algorithms using zebrafish data from Wagner et al. [18]. **(a)** To create ground truth clusters, we artificially enriched and depleted various cell populations in either the experimental or control condition. Here we show the Adjusted Rand Score (ARS) over 100 simulations for 6 methods. VFC is the top performing method on average. **(b)** Difference on performance between each method and VFC on each simulation. **(c)** The sample labels, EES, and clustering results for the simulation in which VFC performed best relative to other methods. **(d)** The sample labels, EES, and clustering results for the simulation in which VFC performed best relative to other methods. We found that by adjusting the weighting of the EES from 1 (default) to 2, VFC becomes the top performing algorithm on this case ('VFC - new').

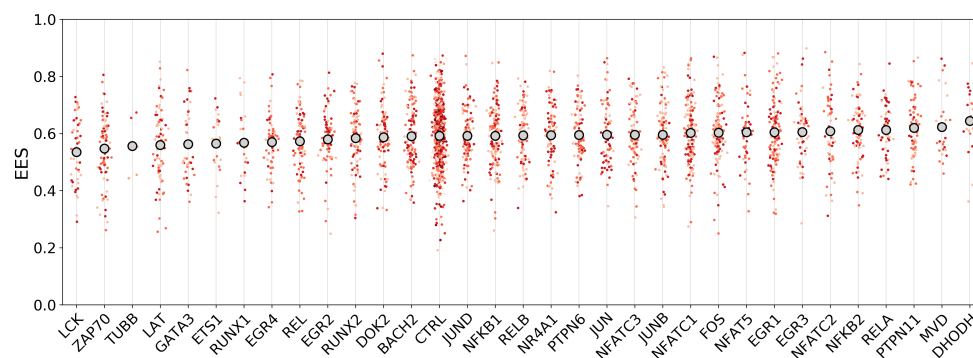


Figure S8: Quantitative analysis of Cas9 perturbations in T cells [16] using the EES. Each plot shows the distribution of EES values for all stimulated cells transfected with gRNAs targeting a specific gene. The shade of each cell indicates the different gRNAs targeting the same gene. To determine the impact of the gRNA on the TCR activation pathway, we rank each gene by the average EES value. We observed a large variation in the impact of each gene knockout consistent with the published results from Datlinger et al. [16]. Encouragingly, our results agree with their bulk RNA-seq validation experiment showing greatest depletion of TCR response with knockout of kinases LCK and ZAP70 and adaptor protein LAT. We also find a slight increase in EES values (and therefore stimulation) in cells in which negative regulators of TCR activation are knocked out, including PTPN6, PTPN11, and EGR3.

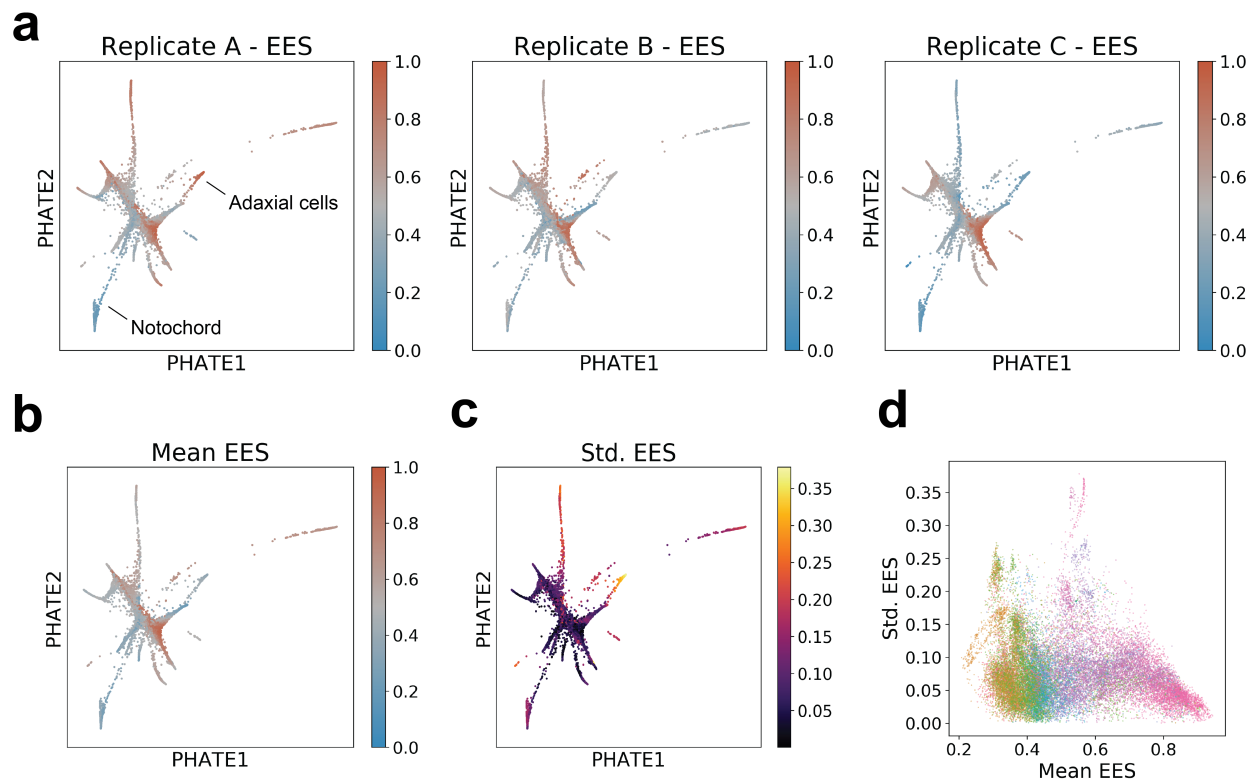


Figure S9: Analysis of replicates within the zebrafish data generated by Wagner et al. [18]. **(a)** Because the EES is calculated by independently filtering a one-hot indicator vector for each condition, to calculate the EES for each replicate, we simply row-normalize the smoothed vectors for the two signals indicating matched experimental / control pairs. For example, the Replicate A - EES is calculated by normalizing the "chdA" and "tyrA" filtered indicator vectors. We notice comparing replicates that the EES for a given cell population may vary. For example, the Adaxial cell population is enriched in the Chd condition in Replicate A, but depleted in Replicate C. Similarly, cells in the Notochord population are depleted in the Chd condition in Replicates A and C, but show minimal change in abundance in Replicate B. **(b)** The average EES across all replicates is shown for each cell on a PHATE embedding. **(c)** The standard deviation of the EES across all replicates is shown for each cell on a PHATE embedding. Regions that have higher values exhibit greater variation in their response to the experimental perturbation. We should trust the average EES values for these cells less than for cells with little variation in EES values. **(d)** A biaxial scatter plot showing the relationship between mean EES and standard deviation in the EES for each cell. Color indicates the cluster labels from **Figure 5a**. We observe that for cells with the highest EES values, the standard deviation is smaller than for cells with EES values close to 0.5, creating a slight negative Pearson correlation of -0.18.

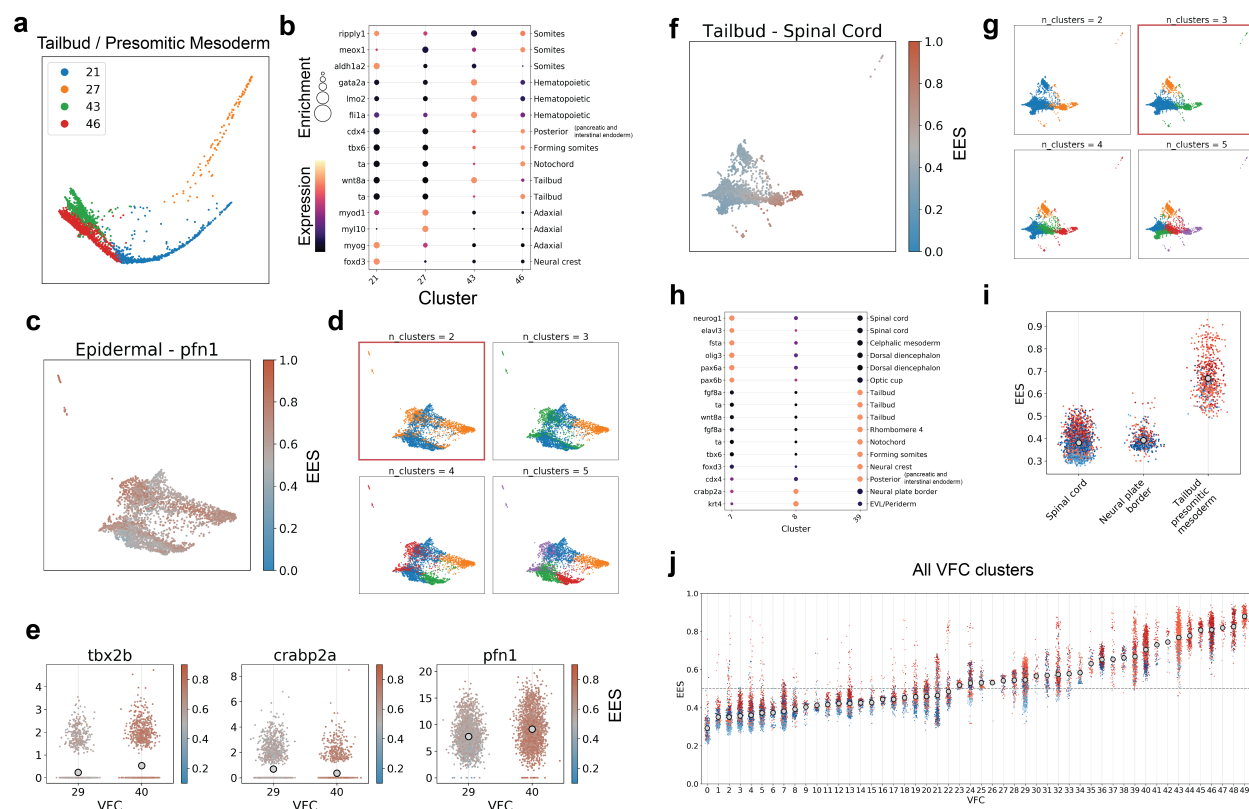


Figure S10: Characterization of vertex-frequency clusters in the zebrafish dataset. **(a)** Raw vertex-frequency cluster assignments on a PHATE visualization of the Tailbud - Presomitic Mesoderm cluster. **(b)** Normalized expression of previously identified marker genes of possible subtypes of the Tailbud - Presomitic Mesoderm [19]. The color of the dot for each gene in each cluster indicates the expression level and the size of the dot corresponds to the normalized Wasserstein distance between expression within cluster to all other clusters. **(c)** Distribution of EES values within the "Epidermal - pfn1" cluster identified by Wagner et al. [18] shown on a PHATE plot. **(d)** Four different values of "n_clusters" that was used to create different VFC clusters with the "Epidermal - pfn1" cluster. We selected n_clusters = 2 because this identified a population of cells with similar EES values and localization on the PHATE embedding. **(e)** Expression of three significantly differentially expressed genes between the two VFC subpopulations detected in the "Epidermal - pfn1" population. Tbx2b and Crabp2a were identified as markers of the epidermis and neural plate border respectively by Farrell et al. [19]. Because we observed differential expression of these two markers between the VFC subclusters suggests the "Epidermal - pfn1" cells identified by Wagner et al. [18] actually comprises cells originating from two distinct cell populations. **(f)** Distribution of EES values within the "Tailbud - Spinal Cord" cluster identified by Wagner et al. [18] shown on a PHATE plot. **(g)** Four different values of n_clusters that was used to create different VFC clusters within the "Tailbud - Spinal Cord" cluster. We selected n_clusters = 3 because this identified populations of cells with similar EES values and localization on the PHATE embedding. **(h)** Same plot as in **(b)** for the subclusters of the "Tailbud - Spinal Cord". **(i)** Distribution of EES values within each VFC subcluster show that the three subclusters are biologically distinct with differing responses to the experimental perturbation. **(j)** Repeating the VFC subclustering process for all cells, we identified a total of 50 clusters within the zebrafish dataset generated by Wagner et al. [18]. Compared to the plot in **Figure 5b**, we observed a more restricted distribution of EES values within each cluster suggesting these labels represent populations of cells that are more homogeneous with respect to the experimental perturbation.

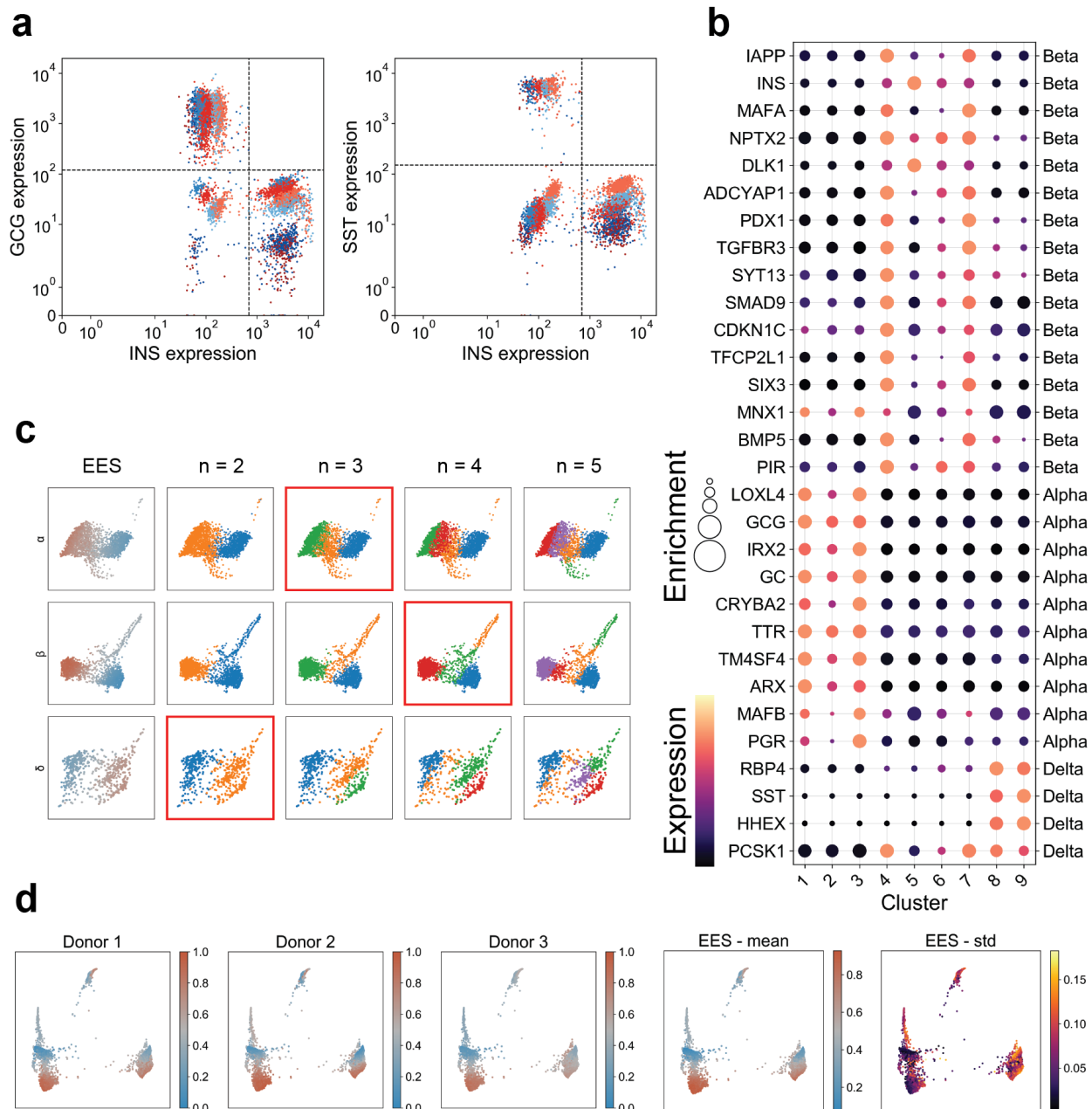


Figure S11: Analysis of pancreatic islet cells from three donors. **(a)** Library-size normalized expression of insulin (INS), glucagon (GCG), and somatostatin (SST) shows donor-specific batch effect across islet cells. **(b)** Normalized expression of previously identified marker genes of alpha, beta, and delta cells[41] in each cluster. The color of the dot for each gene in each cluster indicates the expression level after MAGIC and the size of the dot corresponds to the normalized Wasserstein distance between expression within cluster to all other clusters. **(c)** Results of VFC using varying numbers of clusters for each of the three cell types. The red box denotes the selected level of clustering for each cell type. **(d)** The EES is calculated independently for each donor and then averaged to obtain the EES used in the main analysis. We also calculate the standard deviation of the EES for each cell.

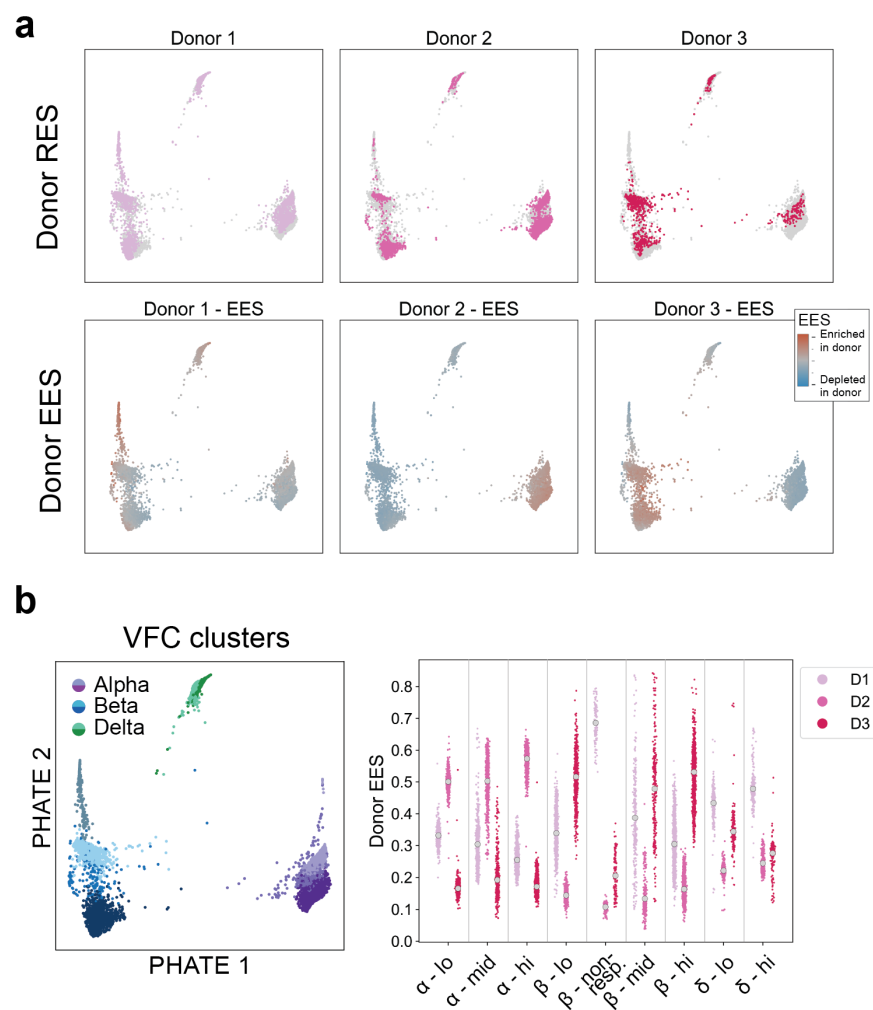


Figure S12: Analysis of islet cell profiles across donors. **(a)** The RES and EES associated with each donor from which islet cells were obtained. **(b)** Comparison of the EES values within each vertex frequency cluster identifies changes in enrichment for each cluster in various donors. For example, the β - non-responsive cluster is strongly enriched in donor 1.

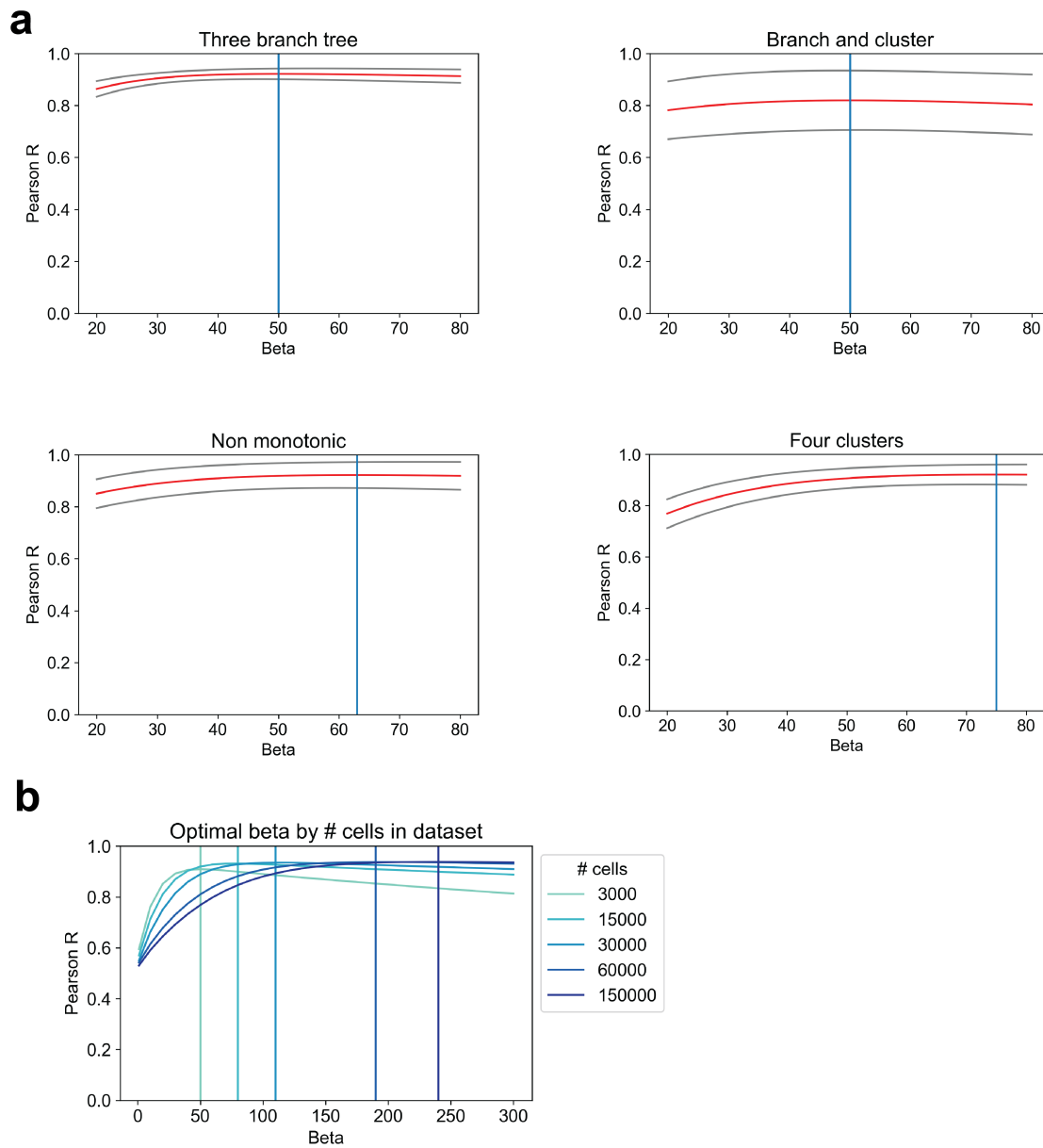


Figure S13: Selecting parameters for MELD. **(a)** Results of a parameter search over the β parameter using the four datasets described in Section 4.7. The red line shows the average performance over 10 different datasets of each geometry with one standard deviation marked by the grey lines. We observe reasonably consistent performance of the EES algorithm across all datasets using a β value between 50-75. We chose a value of 60 as the default in the MELD package and used this setting for all experiments. **(b)** We observe that the optimal β parameter for a dataset varies with the number of cells in the dataset. We suggest increasing the default beta parameter for datasets larger than 30,000 cells.

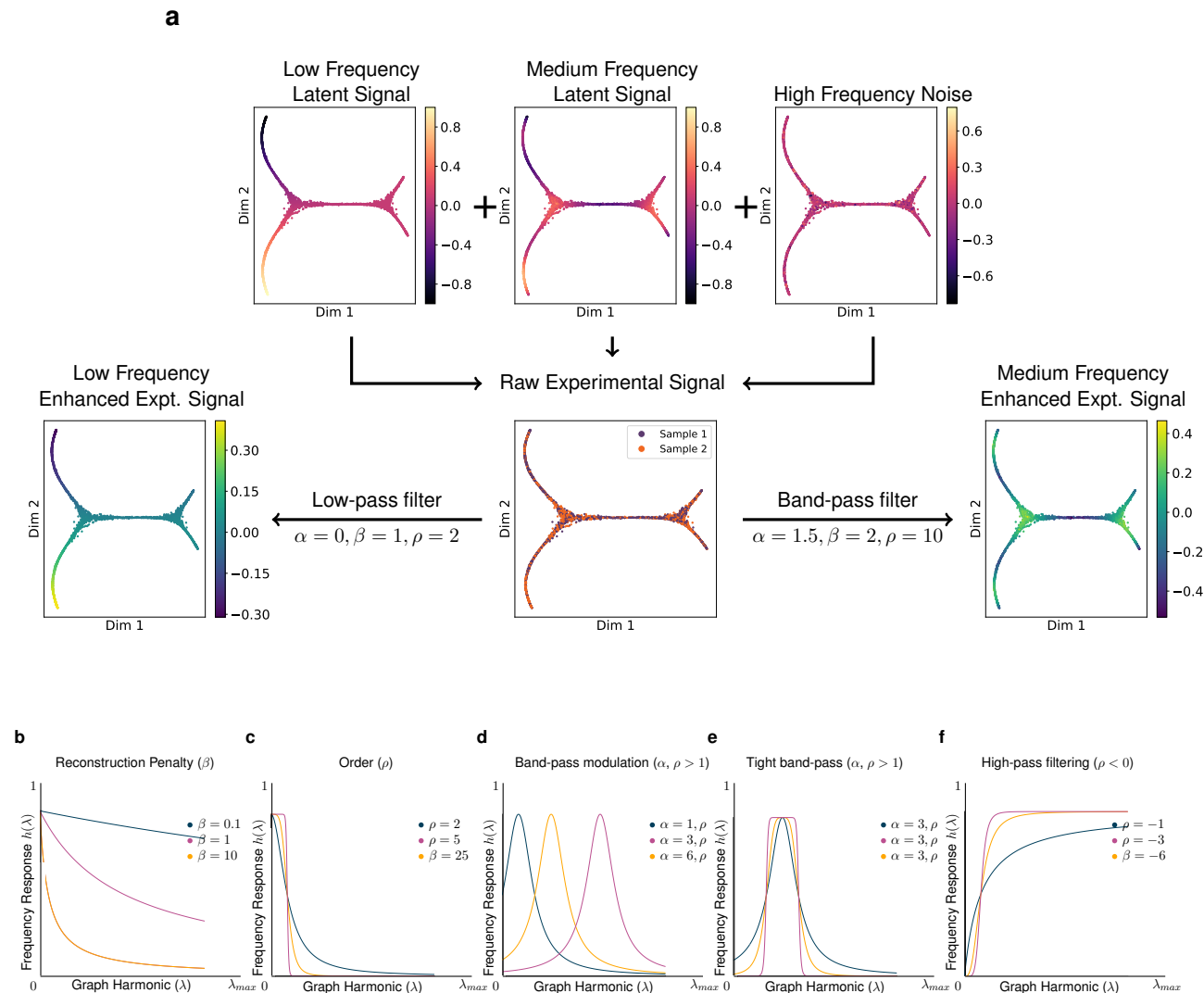


Figure S14: Source Separation and Parameter Analysis with the MELD filter. **(a)** A raw experimental signal (center) is obtained that is a binarized observation of a low frequency latent signal (top left), a medium frequency latent signal (top middle), and high frequency noise (top right). Analysis of the RES alone is intractable as it is corrupted by noise and experimental binarization. MELD low-pass filters (bottom left) to separate a longitudinal trajectory and band-pass filters (bottom right) to yield the periodic signature of the medium frequency latent signal. Parameters used for this analysis are supplied beneath the corresponding arrows. **(b)** Reconstruction penalty β controls a low-pass filter. For this demonstration, $\alpha = 0, \rho = 1$. This filter is equivalent to Laplacian regularization. **(c)** Order ρ controls the filter squareness. This parameter is used in the low-pass filter of **(a)**. For this demonstration, $\beta = 1, \alpha = 0$. **(d)** Band-pass modulation via α . When ρ is even valued, α modulates the central frequency of a band-pass filter. This parameter is used in **(a)** to separate a medium-frequency source from a low-frequency source. **(e)** α and ρ combine to make square band-pass filters. For **(d)** and **(e)**, $\beta = 1$. **(f)** Negative values of ρ yield a high-pass filter. For **(b-f)**, Laplacian harmonics for a general normalized Laplacian are plotted on the x-axis. The frequency response of the filter given by the colored parameters is on the y-axis.