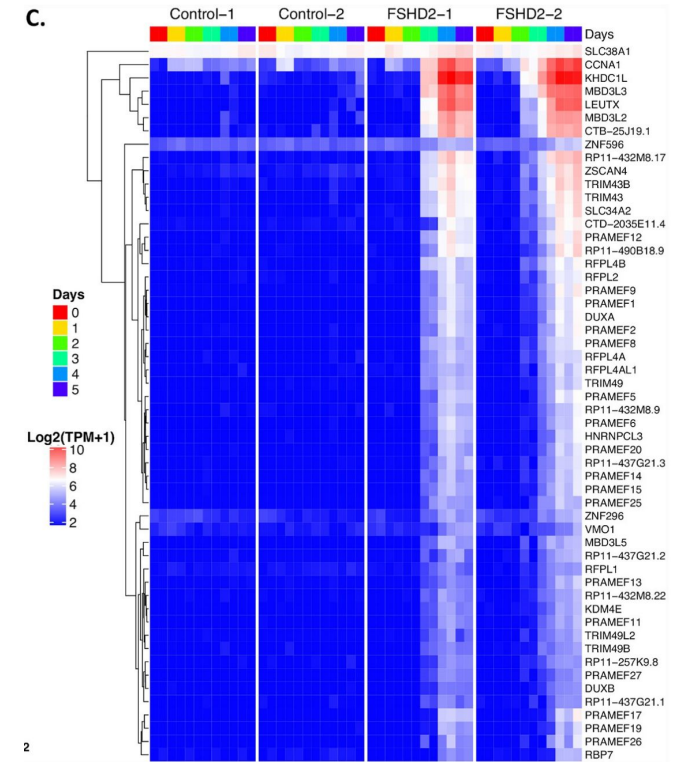


Intro to genomics visualizations

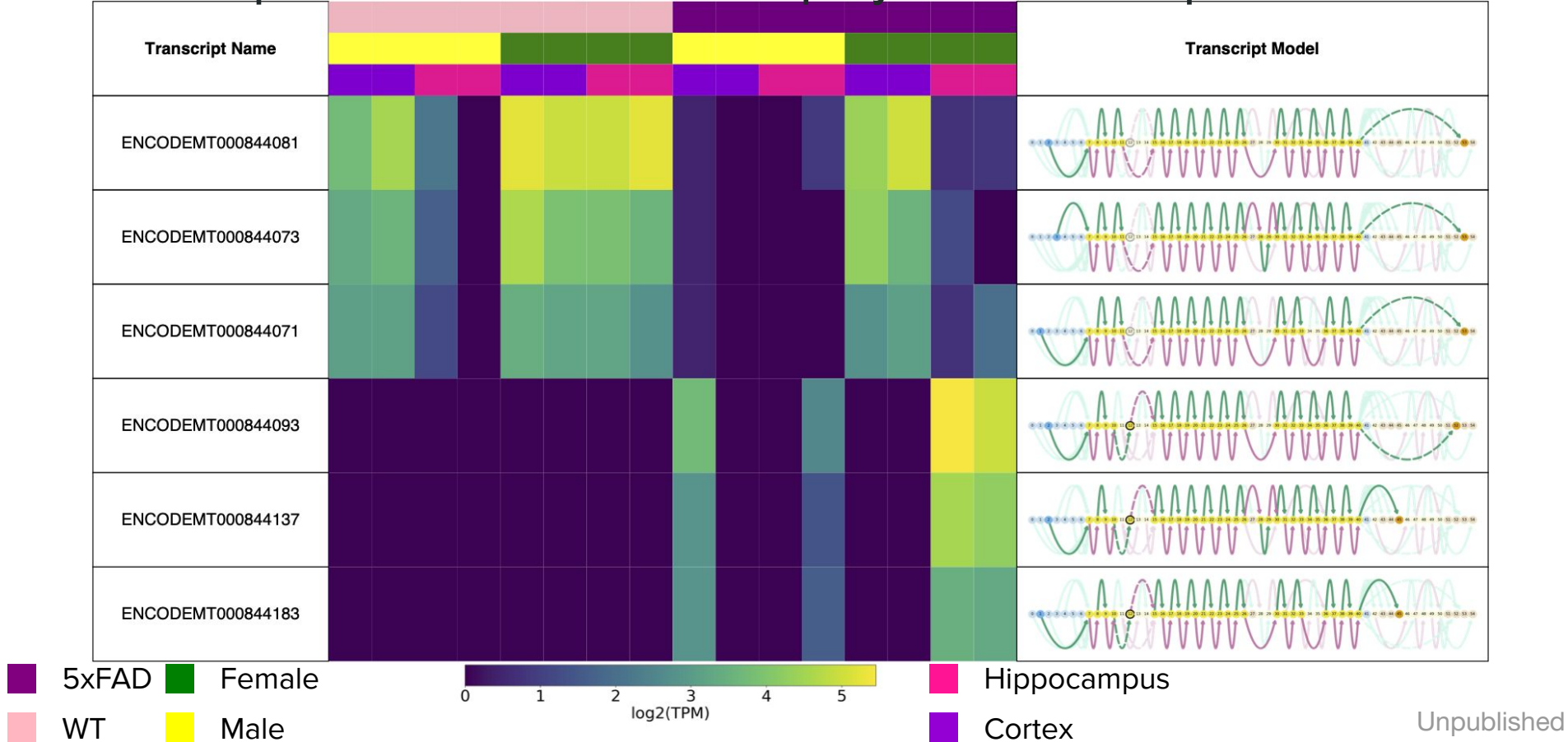
COSMOS 2021
7/13/21

Heatmaps are often used to display trends in expression or signal data

- Displays normalized expression or signal information (ie. TPM, $\log_2(\text{TPM})$, accessibility signal etc.)
- Often row or column-normalized (minimum and maximum color in each sample are used to normalize values in that row)
- Expression profiles often hierarchically clustered and displayed with dendrogram to order genes in an informative way

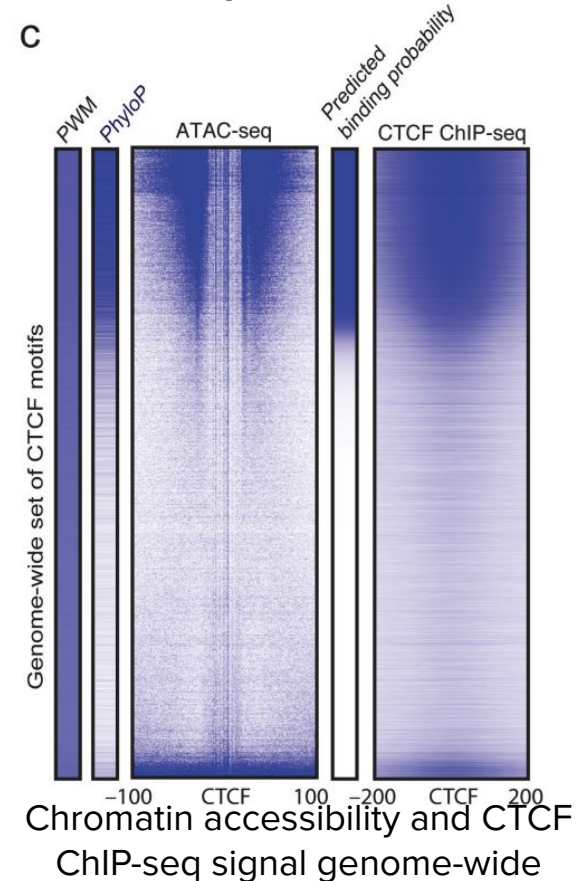


Heatmaps are often used to display trends in expression data



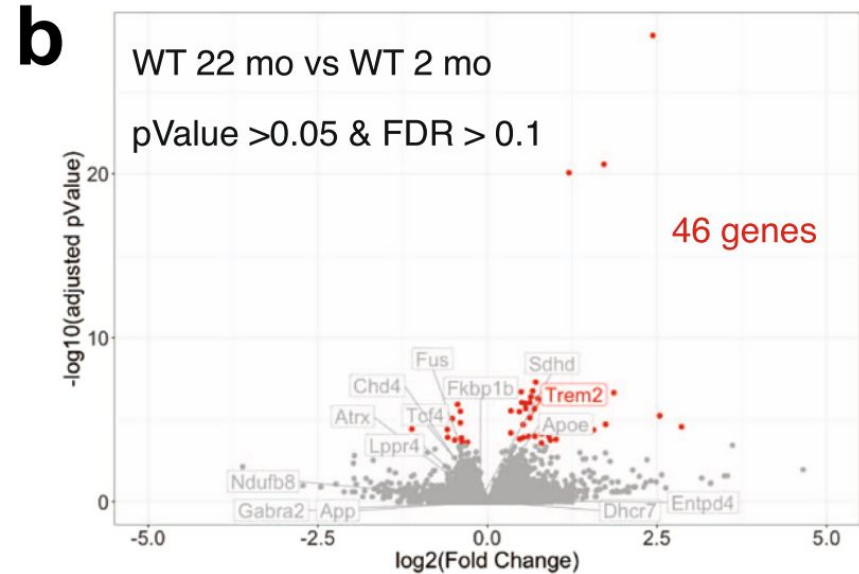
Heatmaps are often used to display trends in expression or signal data

- Displays normalized expression information (ie. TPM, $\log_2(\text{TPM})$, accessibility signal etc.)
- Often row or column-normalized (minimum and maximum color in each sample are used to normalize values in that row)
- Expression profiles often hierarchically clustered and displayed with dendrogram to order genes in an informative way



Volcano plots are used to visualize differences in signal or expression between 2 conditions

- Used to visualize significance and fold change of signal or expression
- Often used for visualizing the results of differential gene expression testing



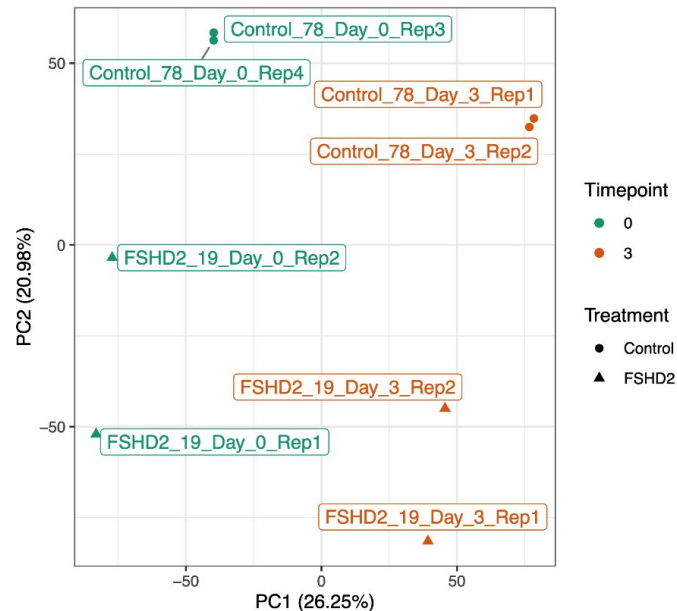
Volcano plot of differentially-expressed genes between 2 month and 22 month old mouse

Lower dimensional representations of expression or signal data can reveal high dimensional similarity between datasets

- High dimensional data can include expression data or signal data
- PCA - principal component analysis
- tSNE - t-distributed stochastic neighbor embedding
- UMAP - uniform manifold approximation and projection

Lower dimensional representations of expression or signal data can reveal high dimensional similarity between datasets

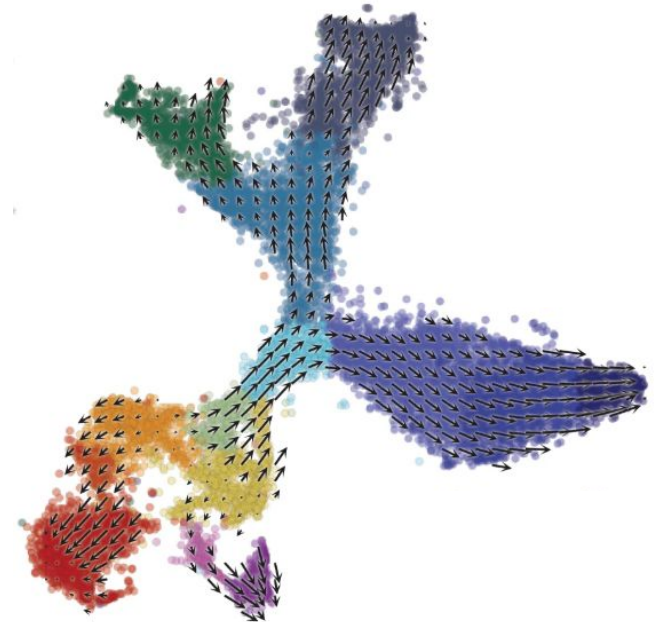
- High dimensional data can include expression data or signal data
- **PCA - principal component analysis**
- tSNE - t-distributed stochastic neighbor embedding
- UMAP - uniform manifold approximation and projection



PCA representation of multiple bulk RNA-seq gene expression datasets

Lower dimensional representations of expression or signal data can reveal high dimensional similarity between datasets

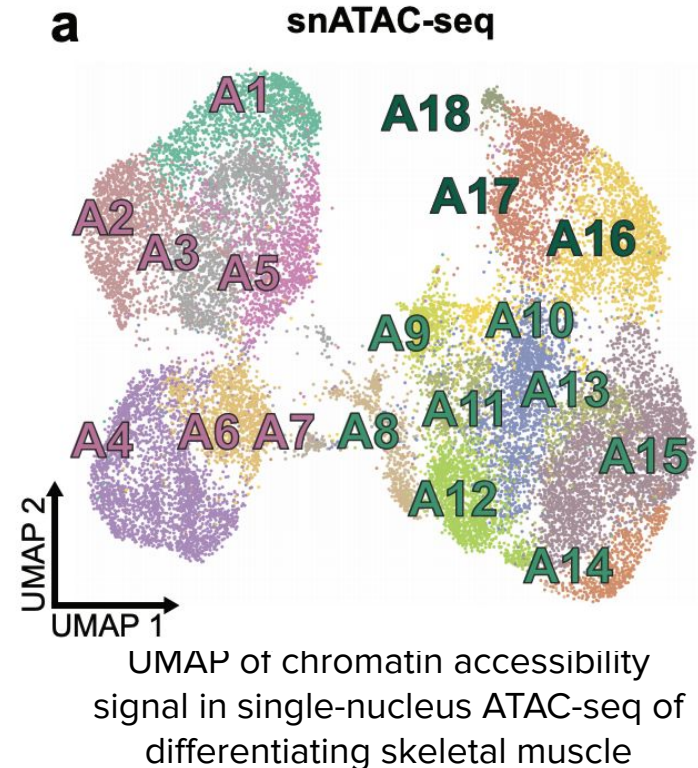
- High dimensional data can include expression data or signal data
- PCA - principal component analysis
- **tSNE - t-distributed stochastic neighbor embedding**
- UMAP - uniform manifold approximation and projection



tSNE representation gene expression profiles in single-cell RNA-seq in developing brain

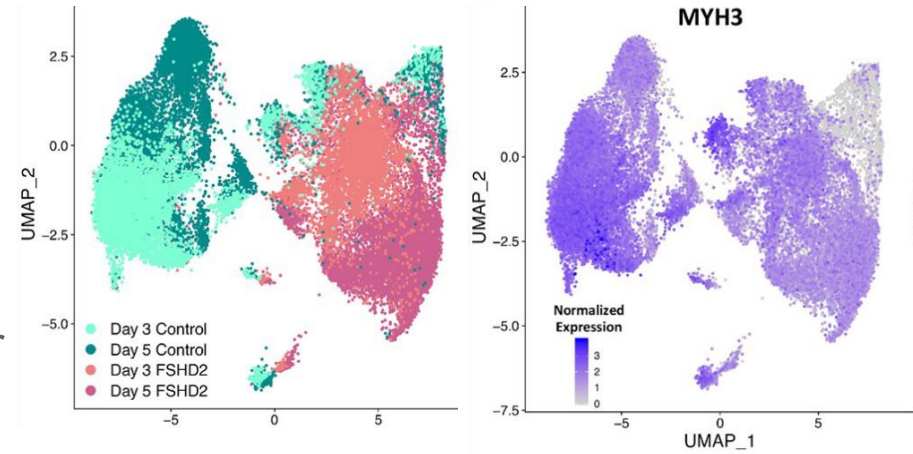
Lower dimensional representations of expression or signal data can reveal high dimensional similarity between datasets

- High dimensional data can include expression data or signal data
- PCA - principal component analysis
- tSNE - t-distributed stochastic neighbor embedding
- **UMAP - uniform manifold approximation and projection**



Lower dimensional representations of expression or signal data can reveal high dimensional similarity between datasets

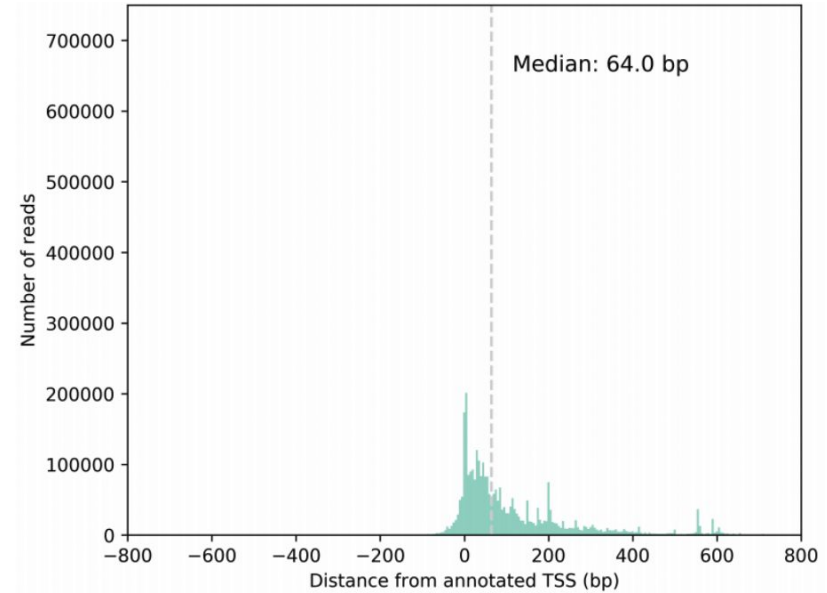
- High dimensional data can include expression data or signal data
- PCA - principal component analysis
- tSNE - t-distributed stochastic neighbor embedding
- **UMAP - uniform manifold approximation and projection**



UMAP of gene expression in FSHD vs. control single-nucleus RNA-seq colored by sample (left), expression of *MYH3* (right).

Histograms and density functions can be used to summarize distributions in data

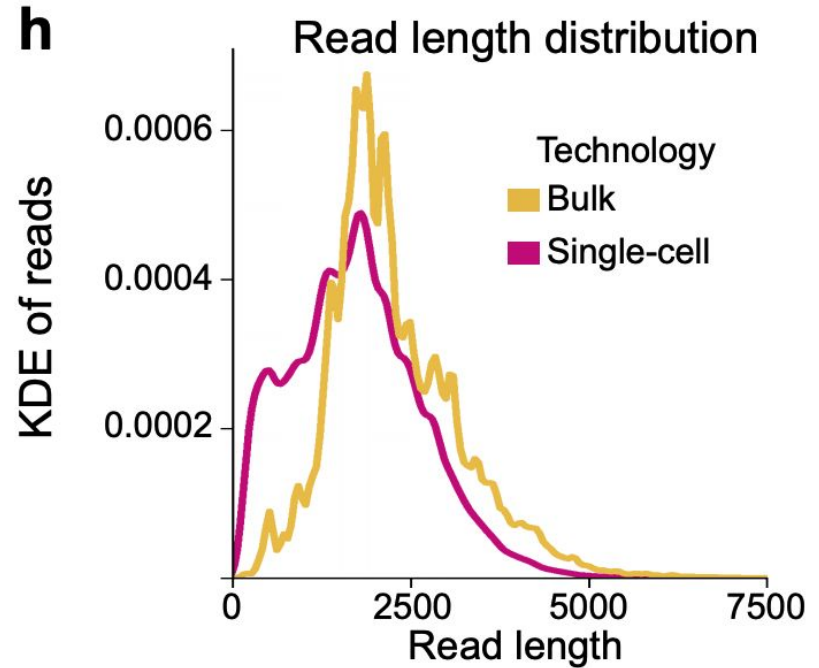
- Plots the density or raw number of occurrences per binned or continuous measurements



Histogram of distance from annotated TSS per read

Histograms and density functions can be used to summarize distributions in data

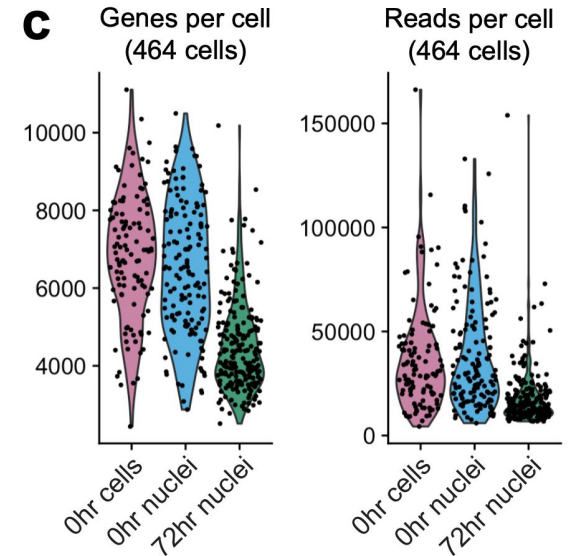
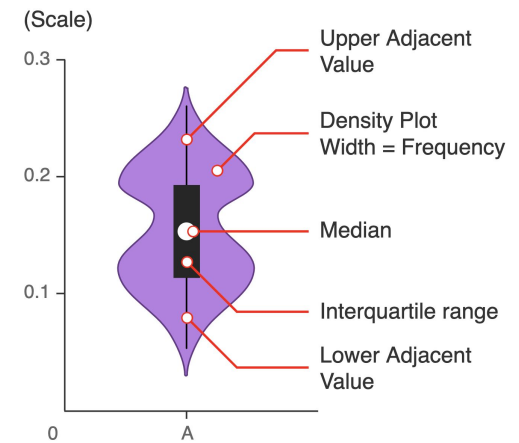
- Plots the density or raw number of occurrences per binned or continuous measurements



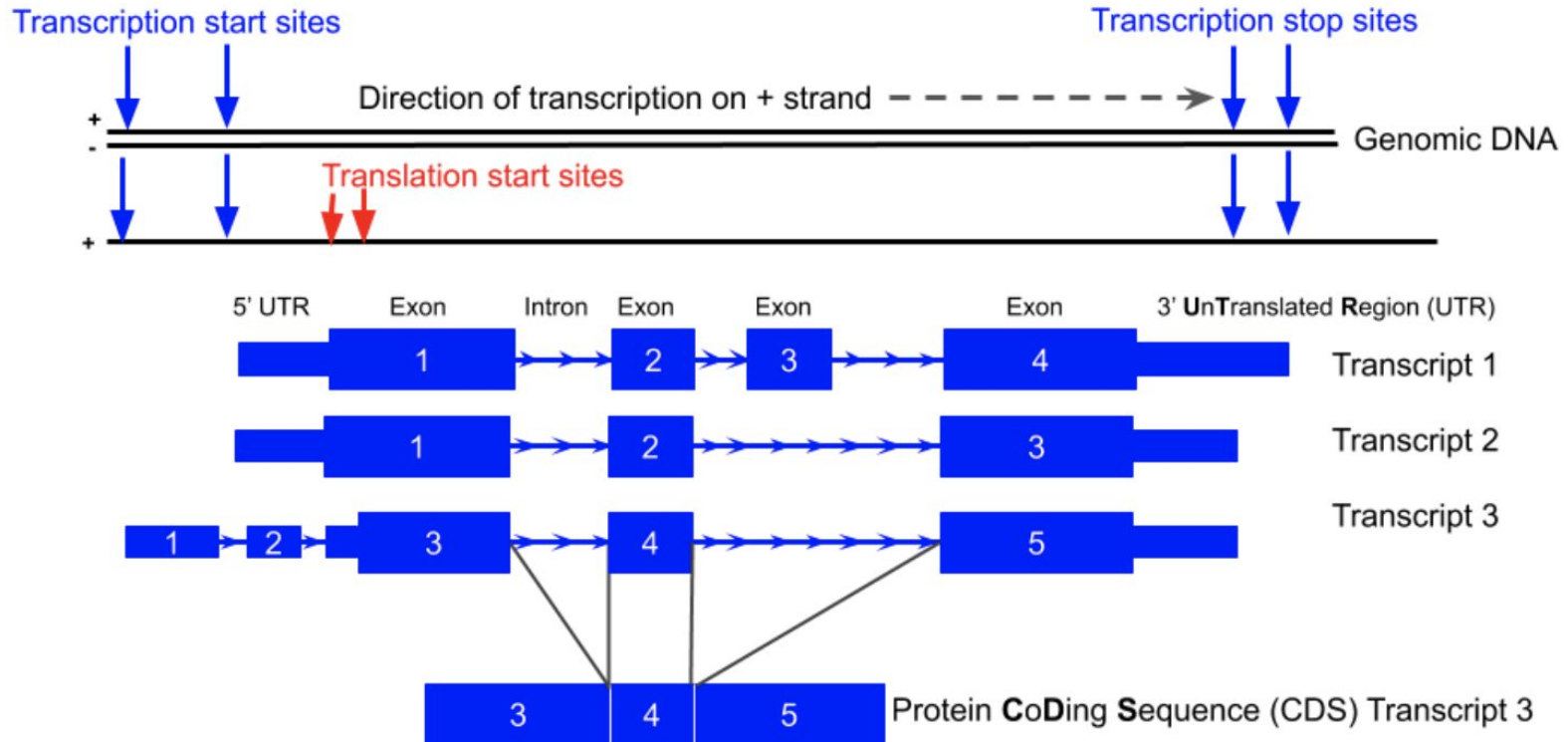
Kernel density estimate of distribution of read lengths between two experimental protocols

Violin plots are useful in single cell data analysis

- Similar to box plot but with kernel density plot on each side
- Show full distribution of data
- Visualize expression levels of a gene across all cells, or quality control metrics (genes detected per cell, reads detected per cell)

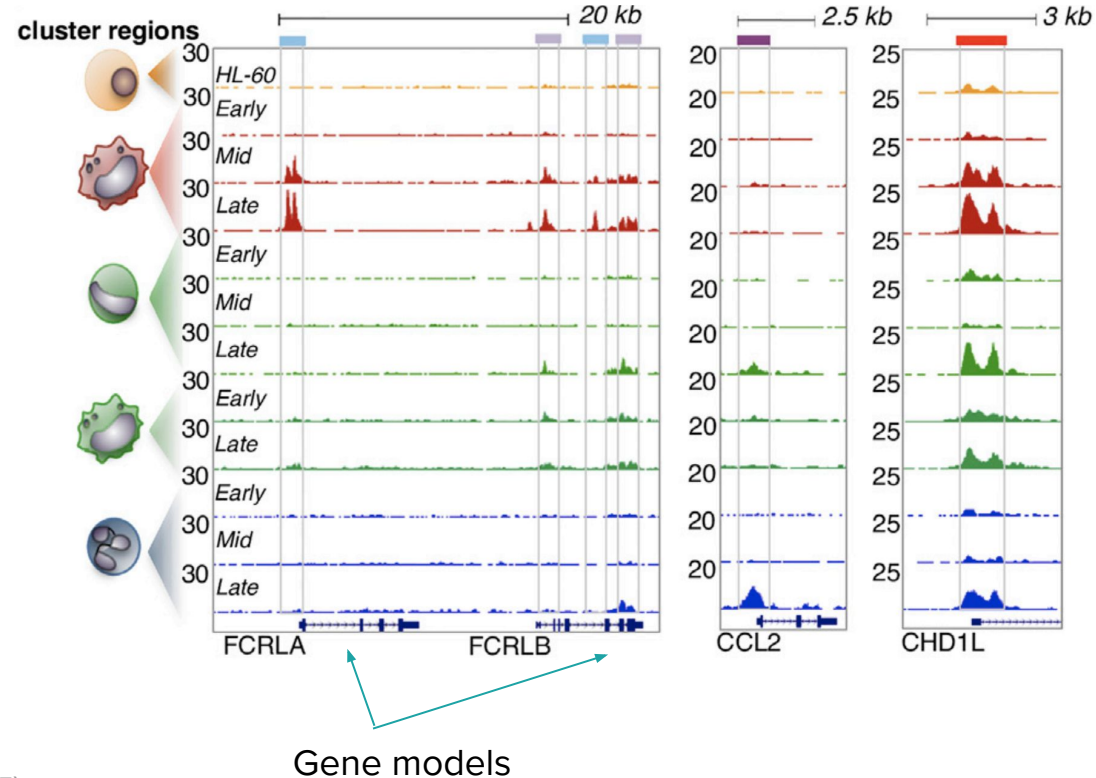


Traditional browser-style representation of transcript and gene models



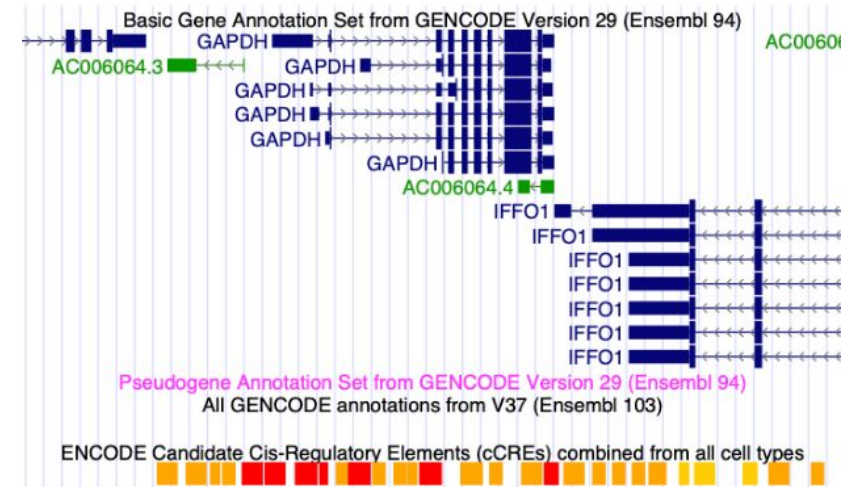
Signal such as chromatin accessibility can be visualized as “peaks” on the genome browser

- Gene expression measured in discrete counts but chromatin data involves “peak calling” of signal across genomic location
- Zoom into interesting regions (like promoters) to visualize signal



Discrete regions without signal (BED regions) can also be visualized on the genome browser

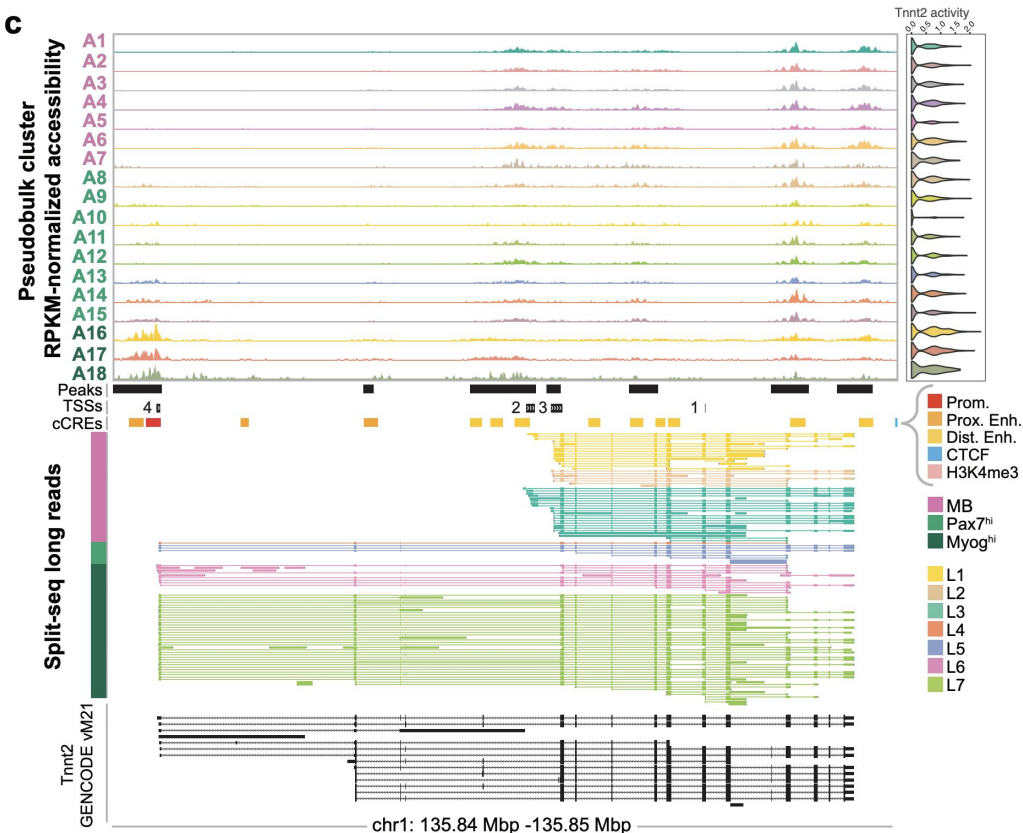
- Useful for regions with annotated function
- ENCODE cCREs (candidate cis-regulatory elements) - regions that are thought to have some regulatory impact on gene expression



hg38 ENCODE cCREs
at the *GAPDH* locus

The genome browser facilitates visualization of multiple types of data at the same location

- Open chromatin (scATAC) signal
- gene “activity” violin plots
- transcript models
- TSS BED regions
- ENCODE cCREs
- Integrate ATAC and RNA data modalities in one figure



Network diagrams show relationships between genes

- Visualize system-level gene regulation
- Genes are represented by nodes connected by edges
- Difficult to construct and often require experimental validation

g

