

# SA12: Hydra URD Granular Mucous & Zymogen Gland Cells

*Jeff Farrell*  
*October 1, 2018*

```
# Load required libraries
library(URD)

## Loading required package: ggplot2
## Warning: package 'ggplot2' was built under R version 3.4.4
## Loading required package: Matrix
## Warning: package 'Matrix' was built under R version 3.4.2

# Set main
opts_chunk$set(root.dir = "~/Dropbox/HydraURDSubmission/")
main.path <- "~/Dropbox/HydraURDSubmission/"

# Build output directory
dir.create(paste0(main.path, "URD/GranularZymogen/"), recursive = T)
```

## Load subsetting URD data for this lineage

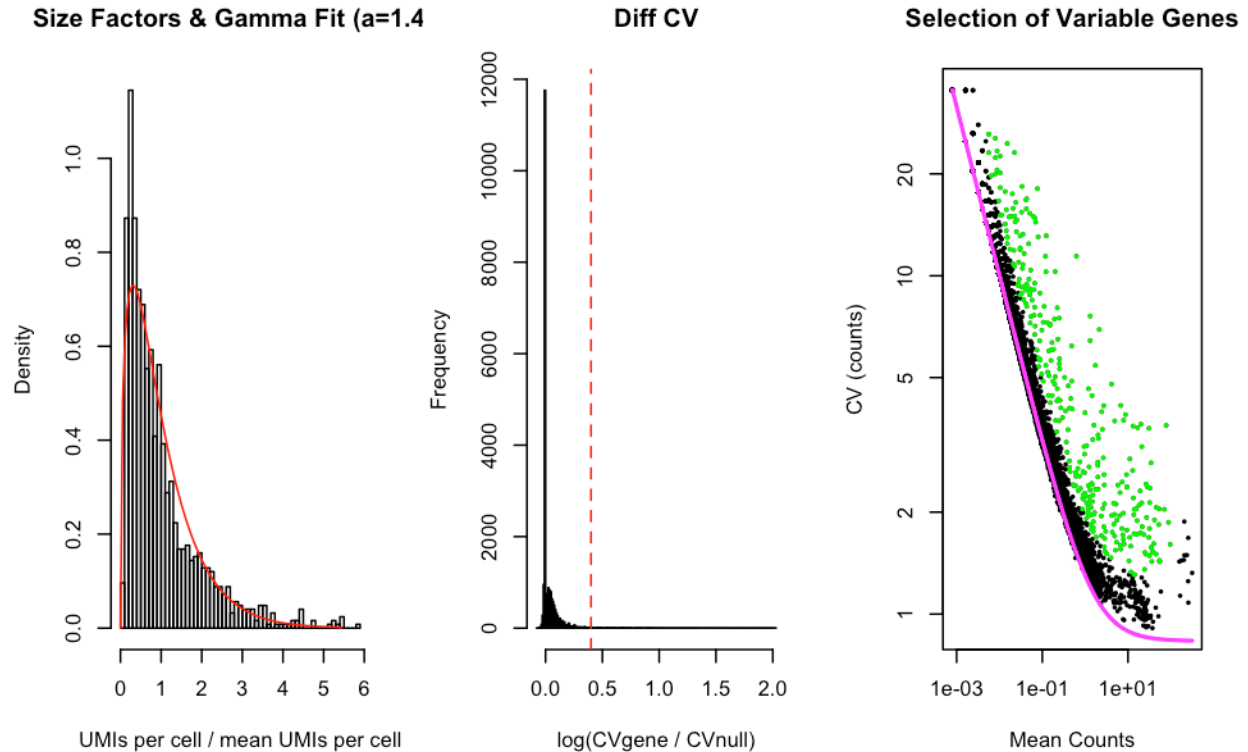
```
hydra.zymogen <- readRDS(file = paste0(main.path, "objects/Hydra_URD_Input_Zymogen.rds"))
```

## Variable genes

Since we've subsetting, it's best to calculate a list of variable genes specific to this lineage.

```
# Determine variable genes (315 genes, 288 were variable across whole IC lineage)
hydra.zymogen <- findVariableGenes(object = hydra.zymogen, set.object.var.genes = T,
  diffCV.cutoff = 0.4, do.plot = T)

## Warning in xy.coords(x, y, xlabel, ylabel, log): 11693 x values <= 0
## omitted from logarithmic plot
```



## Graph Clustering

Generated several more fine-grained clusterings in order to have better options for choosing root and tip cells.

```
# Set random seed so clusters get the same names every time
set.seed(19)
# Graph clustering with various parameters
hydra.zymogen <- graphClustering(hydra.zymogen, num.nn = c(20, 30, 40, 50), do.jaccard = T,
  method = "Infomap")
hydra.zymogen <- graphClustering(hydra.zymogen, num.nn = c(5, 10, 15, 20, 30), do.jaccard = T,
  method = "Louvain")
hydra.zymogen <- graphClustering(hydra.zymogen, num.nn = c(6, 7, 8), do.jaccard = T,
  method = "Louvain")
# Record the clustering
pdf(paste0(main.path, "URD/GranularZymogen/Clusters-GranularZymogen-Infomap20.pdf"))
plotDim(hydra.zymogen, "Infomap-20", label.clusters = T, legend = F)
dev.off()

## quartz_off_screen
## 2
```

## Diffusion map

### Calculate the transition probabilities

Calculate a diffusion map specific to the granular mucous trajectory. In these very specific datasets (i.e. those cropped to a single lineage), using more nearest neighbors than the square root of the data seems to work better, probably because there are fewer cell states to consider. (Here we use 75 NNs, versus prediction of 35.) We used a global sigma that is slightly smaller than the 'autodetected' one (here, 6).

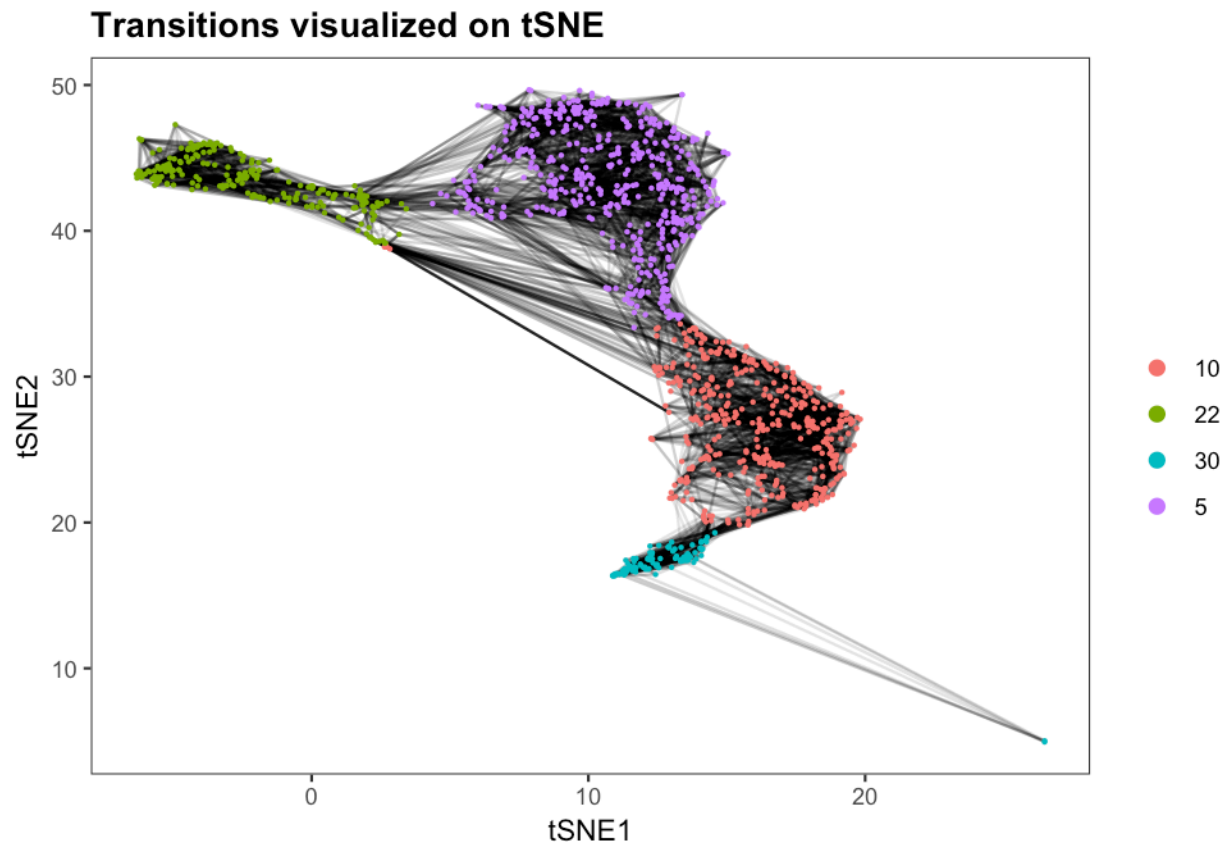
```
# Calculate diffusion maps without doublets or terminal nematocytes, using the
# more restricted variable gene list that lacks CCA genes
hydra.zymogen <- calcDM(hydra.zymogen, knn = 75, sigma.use = 6)

## [1] "Using provided global sigma 6"
```

## Evaluation of diffusion map

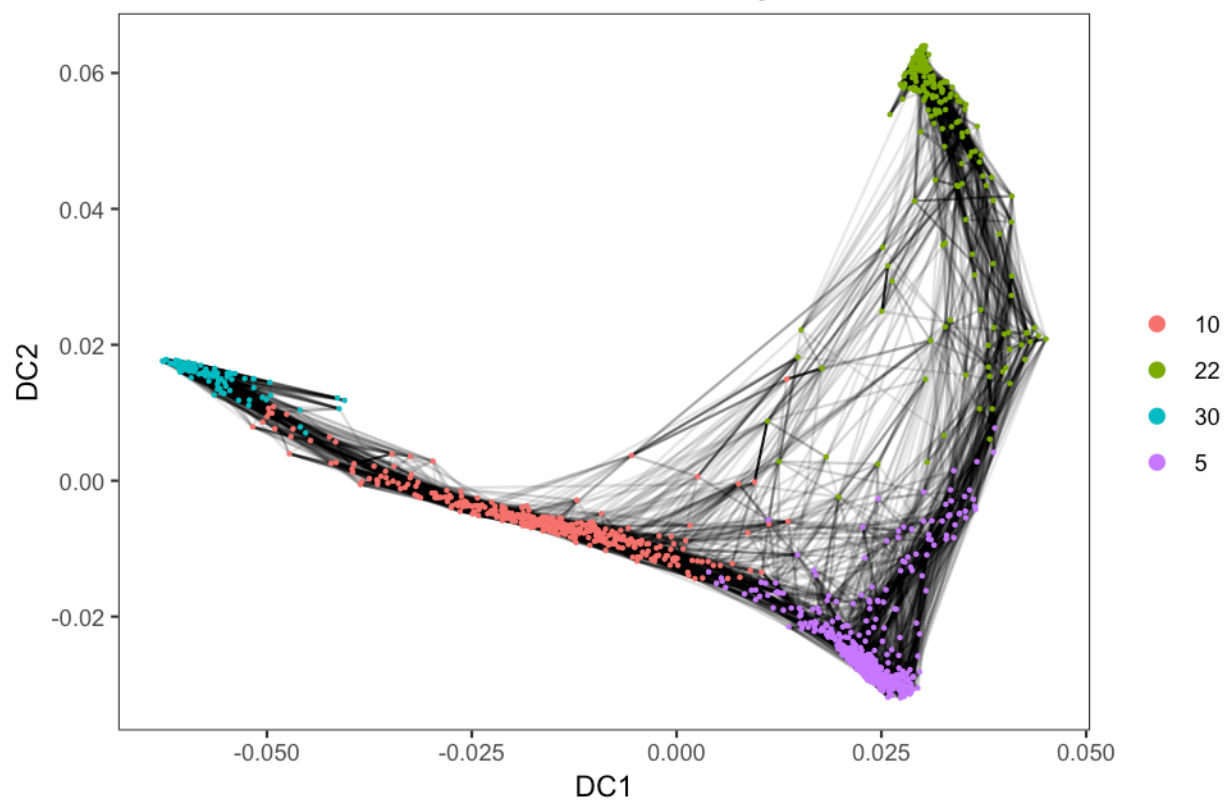
The diffusion map looks good and represents the ends of the differentiation process strongly as tips. In DC1 and DC2, it is evident that there are multiple pathways through the data – both a continuous head to foot differentiation that we aim to capture here, and also separate paths to multiple parts of the body column that are captured as a branching structure in the interstitial tree, and can be seen in DCs 3 and 4.

```
# Make transition plots
plotDim(hydra.zymogen, "res.1.5", transitions.plot = 5000, plot.title = "Transitions visualized on tSNE")
```



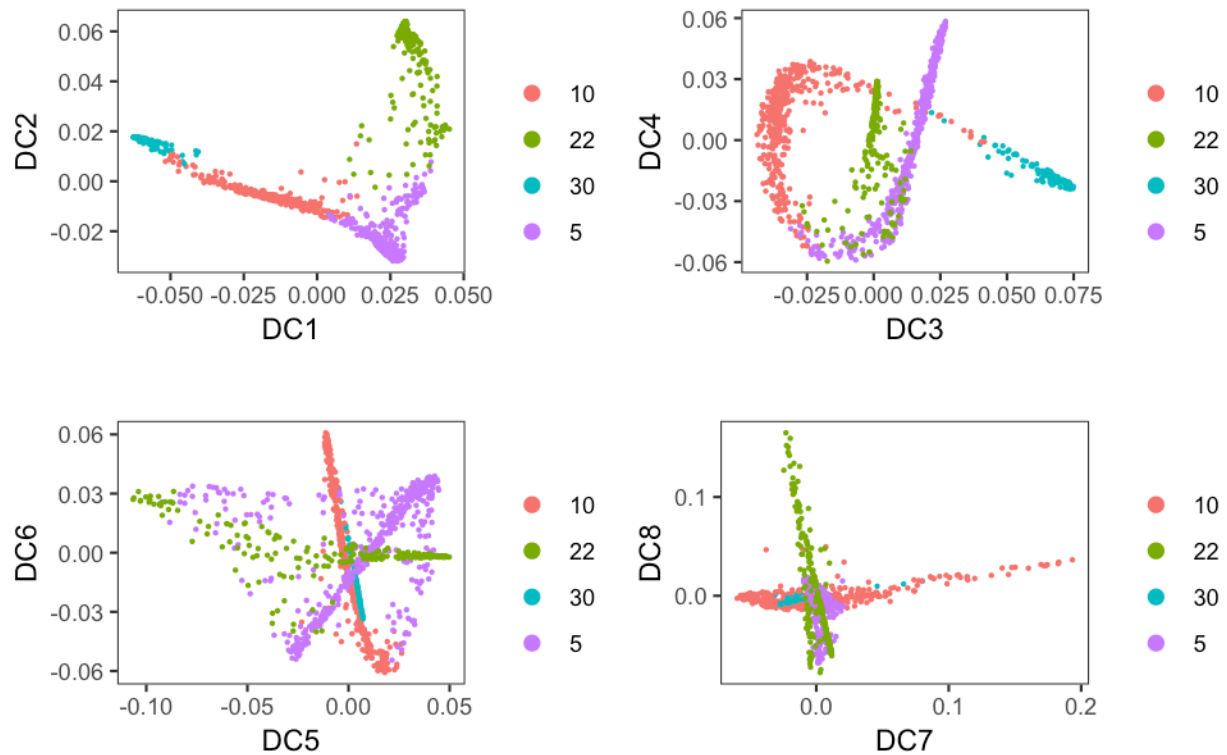
```
plotDim(hydra.zymogen, "res.1.5", transitions.plot = 10000, reduction.use = "dm",
plot.title = "Transitions visualized on diffusion map")
```

## Transitions visualized on diffusion map



```
# Make array plots  
plotDimArray(hydra.zymogen, label = "res.1.5", reduction.use = "dm", dims.to.plot = 1:8,  
  plot.title = "", outer.title = "Pairs of Diffusion Components")
```

## Pairs of Diffusion Components



## Calculate pseudotime

Here, we treat the process as a continuous head to food set of transitions. We calculate pseudotime from each end (starting with the head as root, then starting with the foot as root) and average the two pseudotimes to produce a single ordering across the entire granular mucous gland cell and zymogen gland cell lineage.

Since the simulation process is stochastic, we load our previously calculated results for consistency, but the following commands were run previously:

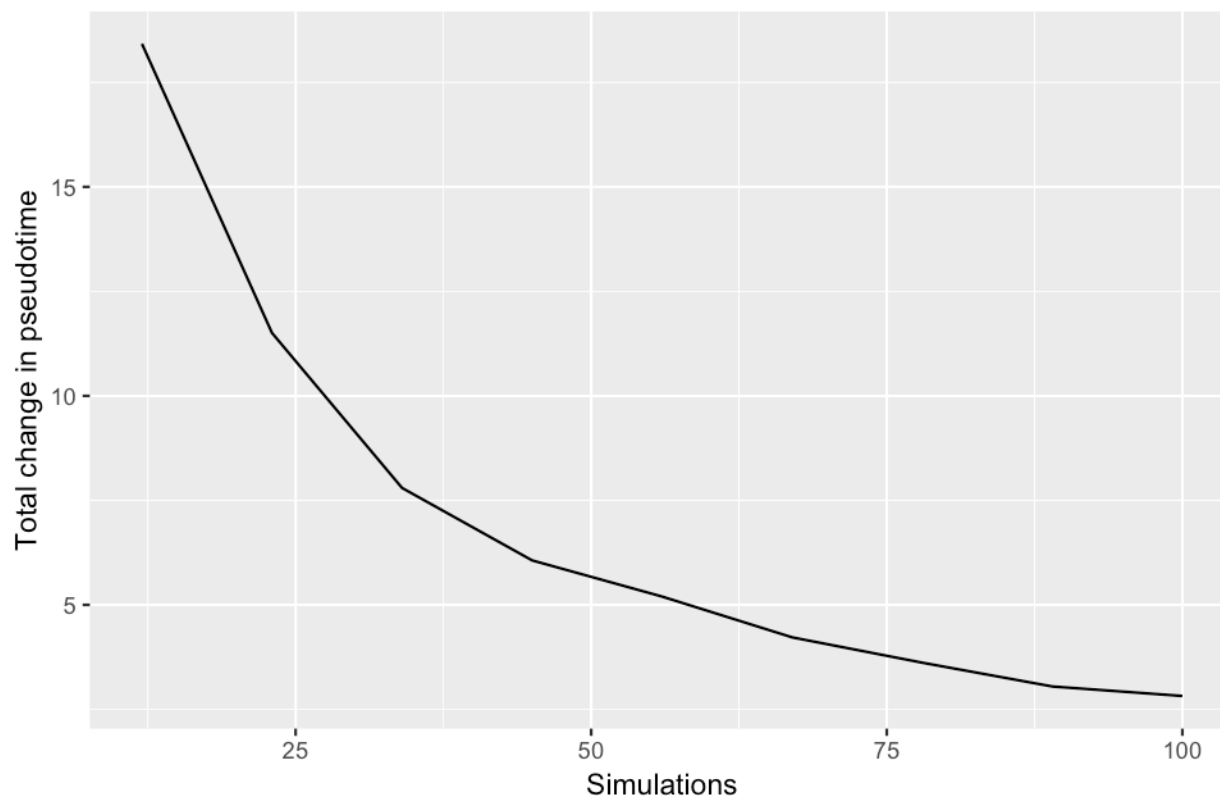
```
# Perform the 'flood' simulations to determine cells' pseudotime
zymogen.flood.13 <- floodPseudotime(hydra.zymogen, root.cells = cellsInCluster(hydra.zymogen,
  "Infomap-20", "13"), n = 100, minimum.cells.flooded = 2, verbose = T)
zymogen.flood.14 <- floodPseudotime(hydra.zymogen, root.cells = cellsInCluster(hydra.zymogen,
  "Infomap-20", "14"), n = 100, minimum.cells.flooded = 2, verbose = T)
```

We then process the random simulations to convert them to a 0-1 range pseudotime and verify that enough simulations have been performed.

```
# Process the floods to derive pseudotime
hydra.zymogen <- floodPseudotimeProcess(hydra.zymogen, zymogen.flood.13, floods.name = "pseudotime.13")
hydra.zymogen <- floodPseudotimeProcess(hydra.zymogen, zymogen.flood.14, floods.name = "pseudotime.14")
```

```
# Check that enough pseudotime simulations were run -- is change in pseudotime
# reaching an asymptote?
pseudotimePlotStabilityOverall(hydra.zymogen)
```

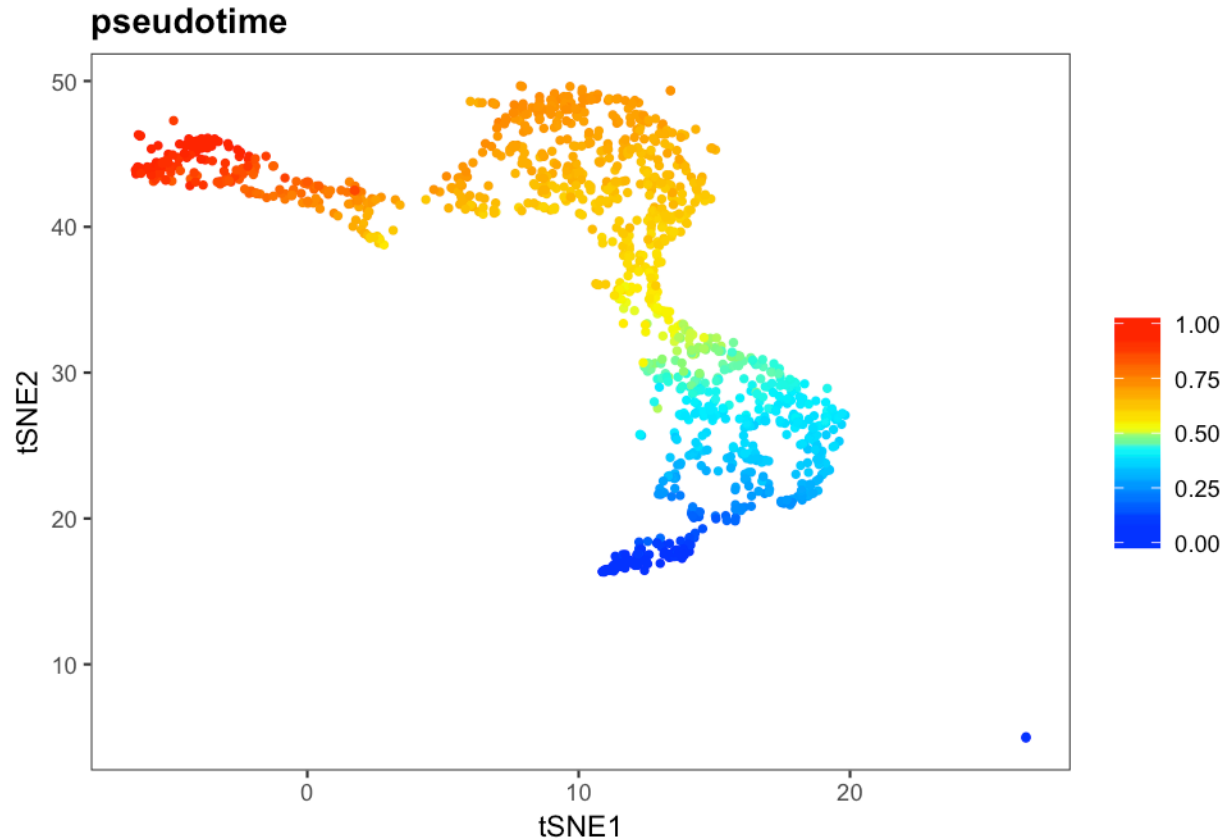
## Overall Pseudotime Stability



```
# Combine the two into a single pseudotime so that you don't get a bunch of cells
# squished up at 0.
hydra.zymogen@pseudotime$pseudotime.13.norm <- hydra.zymogen@pseudotime$pseudotime.13/max(hydra.zymogen@pseudotime$pseudotime.13)
hydra.zymogen@pseudotime$pseudotime.14.norm <- 1 - (hydra.zymogen@pseudotime$pseudotime.14/max(hydra.zymogen@pseudotime$pseudotime.14))
hydra.zymogen@pseudotime$pseudotime <- rowMeans(hydra.zymogen@pseudotime[, c("pseudotime.13.norm",
  "pseudotime.14.norm")])
```

Pseudotime was calculated from the head to the foot. It has good temporal ordering across this trajectory.

```
# Inspect pseudotime on the tSNE plot
plotDim(hydra.zymogen, "pseudotime")
```



## Find spatially varying genes

Since in this trajectory, “pseudotime” is really a proxy for spatial location (head to foot), we can find the spatially varying genes by finding those that vary in pseudotime. We group cells 5 at a time and calculate a spline curve that fits the mean expression (vs. pseudotime) of each group of 5 cells. We then consider those genes spatially varying that are: (1) expressed (present in at least 1% of cells, mean expression spline curve > 0.5 at some point), (2) vary significantly (their spline curve varies at least 33%), (3) are well fit (noise is usually poorly fit, here we threshold on the sum of squared residuals).

## Calculate varying genes

```
# Consider all genes expressed in 1% of cells
frac.exp <- rowSums(hydra.zymogen@logupx.data > 0)/ncol(hydra.zymogen@logupx.data)
expressed.genes <- names(frac.exp)[which(frac.exp > 0.01)]

# Calculate spline fit
expressed.spline.5cell <- geneSmoothFit(hydra.zymogen, method = "spline", pseudotime = "pseudotime",
  cells = colnames(hydra.zymogen@logupx.data), genes = expressed.genes, moving.window = 1,
  cells.per.window = 5, spar = 0.875)

## [1] "2018-11-01 15:14:56: Calculating moving window expression."
## [1] "2018-11-01 15:15:38: Generating un-scaled fits."
## [1] "2018-11-01 15:16:06: Generating scaled fits."
## [1] "2018-11-01 15:16:33: Reducing mean expression data to same dimensions as spline fits."

# Which genes have a spline curve that crosses 0.5?
max.mean.spline.5cell <- apply(expressed.spline.5cell$mean.smooth, 1, max)
genes.wellexp.5cell <- names(which(max.mean.spline.5cell > 0.5))

# Look at how good the spline fits are
spline.fit.5cell <- apply(expressed.spline.5cell$scaled.smooth - expressed.spline.5cell$scaled.expression.red,
```

```

1, function(i) sum(i^2))
spline.fit.norm.5cell <- spline.fit.5cell/ncol(expressed.spline.5cell$scaled.expression.red)
genes.wellfit.5cell <- names(which(spline.fit.norm.5cell <= 0.045))

# Which genes change in their scaled log2 mean expression value by at least 33%?
change.scale.5cell <- apply(expressed.spline.5cell$scaled.smooth, 1, function(x) diff(range(x)))
genes.scale.5cell <- names(which(change.scale.5cell >= 0.33))

# Take the intersection of those genes and use them in the heatmap & analysis
varying.genes.5cell <- intersect(intersect(genes.wellexp.5cell, genes.scale.5cell),
genes.wellfit.5cell)

```

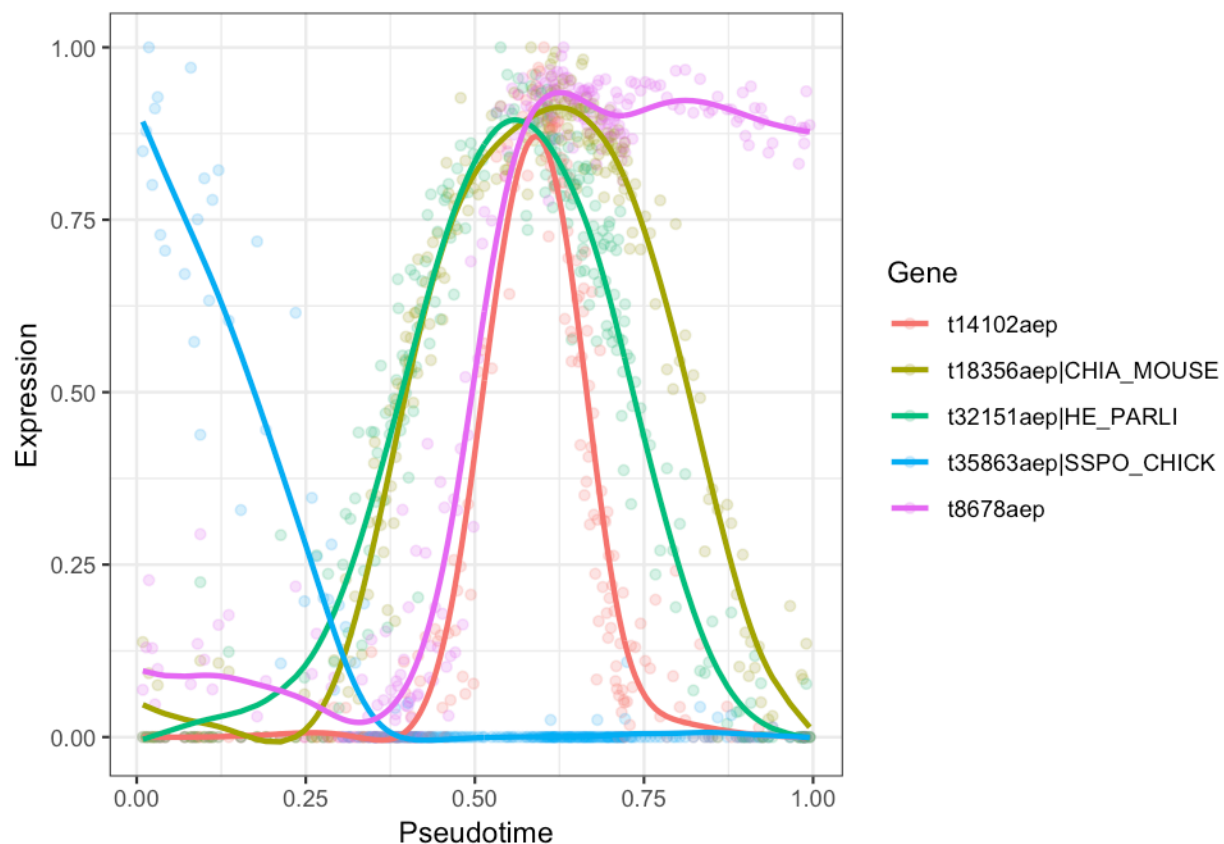
## Spline plots

We can then plot the expression of genes and their fit splines.

```

genes.spline.plot <- c("t14102aep", "t8678aep", "t35863aep|SSPO_CHICK", "t18356aep|CHIA_MOUSE",
"t32151aep|HE_PARLI")
plotSmoothFit(expressed.spline.5cell, genes = genes.spline.plot, scaled = T)

```



## Heatmaps

Finally, we can make heatmaps of the expression of all of the genes that we found to vary spatially. These we save as PDF output because they do not perform well inside of R Markdown.

```

# Heatmap Basics
cols <- (scales::gradient_n_pal(RColorBrewer::brewer.pal(9, "YlOrRd")))(seq(0, 1,
length.out = 50))

# Get the scaled data, z-score it, and set values <0 to 0, the hierarchical
# cluster. This emphasizes clustering on regions of peak expression for a gene.
se <- expressed.spline.5cell$scaled.expression[varying.genes.5cell, ]
se.sd <- apply(se, 1, sd)

```



```

se.mean <- apply(se, 1, mean)
se.z <- sweep(sweep(se, 1, se.mean, "-"), 1, se.sd, "/")
se.z[se.z < 0] <- 0
h.sez <- as.dendrogram(hclust(dist(se.z)))

# Find a 'peak pseudotime' for each gene by taking the weighted average of each
# window's pseudotime, weighted by the expression z-score (>0) within it
pt.wm <- apply(se.z, 1, function(w) {
  weighted.mean(x = as.numeric(colnames(se.z)), w = w)
})

# Make some heatmaps
row.font.size = 0.12
pdf(paste0(main.path, "URD/GranularZymogen/GranularZymogen-Varying.pdf"), width = 17,
    height = 22)
gplots::heatmap.2(as.matrix(se), Rowv = reorder(h.sez, max(pt.wm) - pt.wm, agglo.FUN = median),
  Colv = F, dendrogram = "none", col = cols, trace = "none", density.info = "none",
  key = F, cexCol = 0.8, cexRow = row.font.size, margins = c(5, 8), lwid = c(0.3,
    4), lhei = c(0.4, 4), labCol = NA)
title("Granular Mucous (Head to Foot)", line = -1, adj = 0.48, cex.main = 4)
dev.off()

## quartz_off_screen
##                2

```

## Save results

```

# Save the results
saveRDS(expressed.spline.5cell, file = paste0(main.path, "URD/GranularZymogen/Splines-GranularZymogen.rds"))
write(varying.genes.5cell, file = paste0(main.path, "URD/GranularZymogen/Genes-GranularZymogen-Varying.txt"))
saveRDS(hydra.zymogen, file = paste0(main.path, "URD/GranularZymogen/Hydra_URD_GranularZymogen.rds"))

```