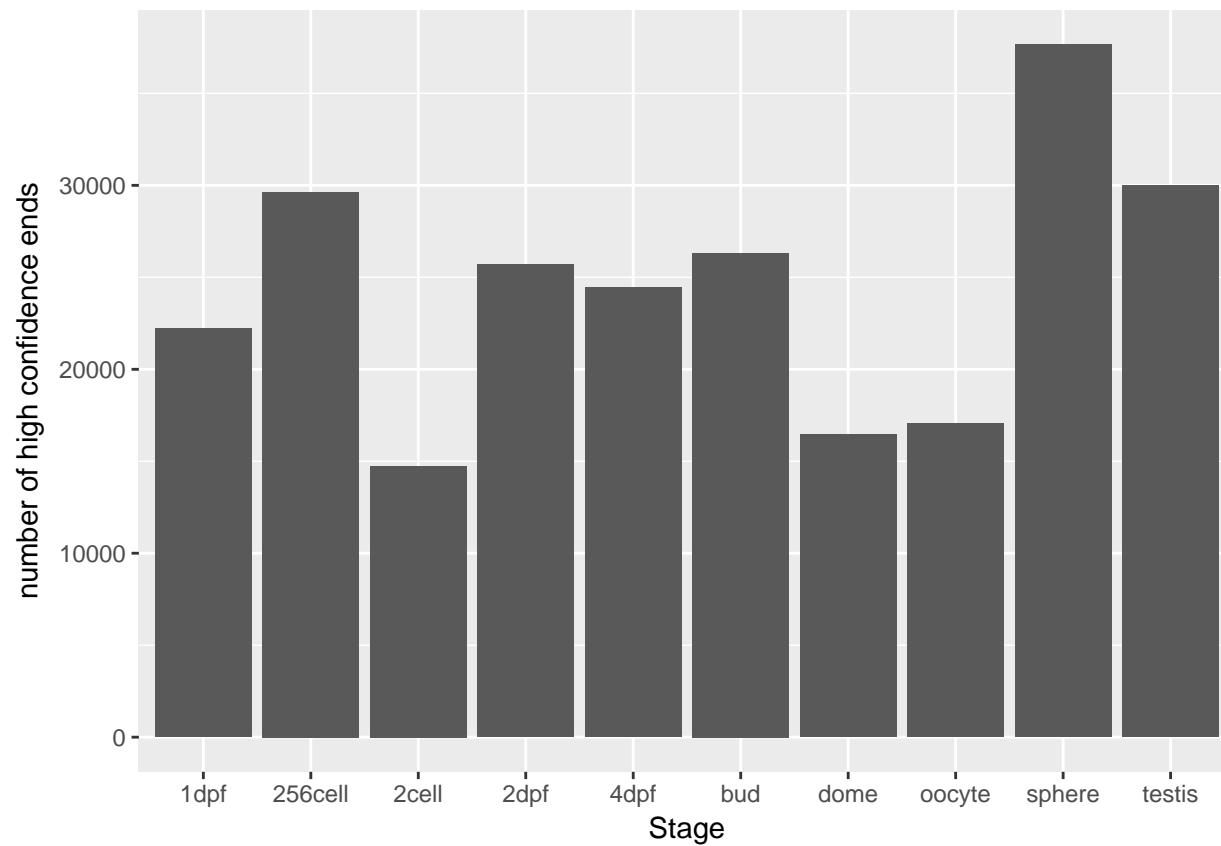# comparing counting windows from different stages
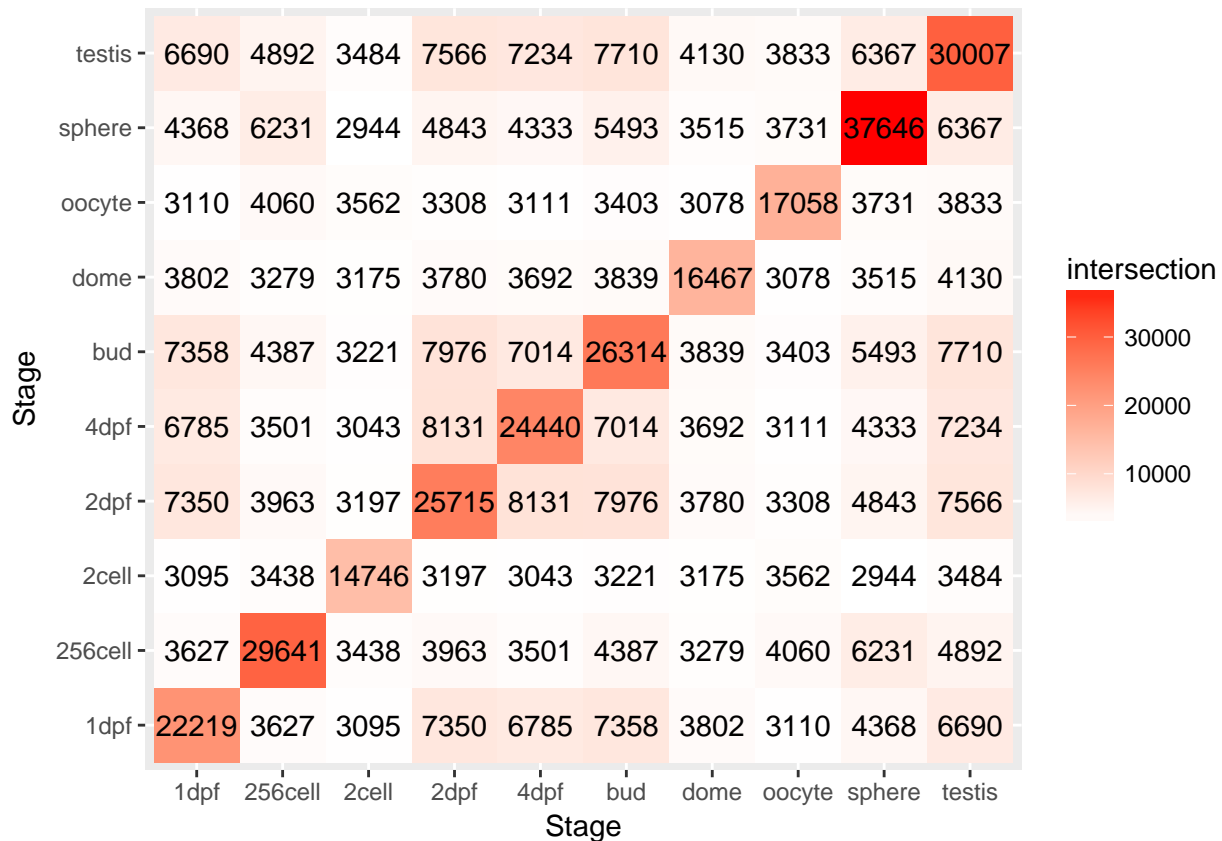
*Pooja Bhat*

*July 19, 2017*

## R Markdown



```
## pdf
##   2
```

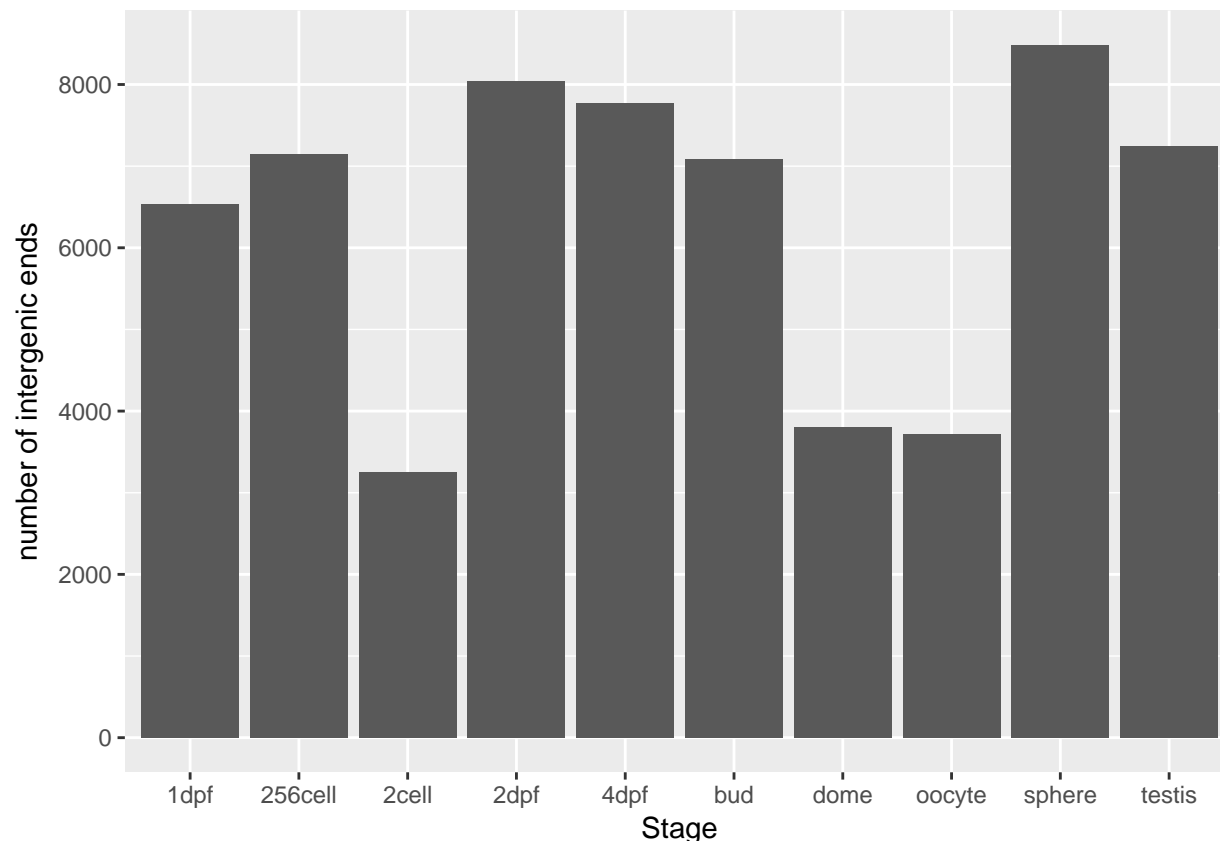| Stage | 1dpf | 256cell | 2cell | 2dpf | 4dpf | bud | dome | oocyte | sphere | testis |
|---|---|---|---|---|---|---|---|---|---|---|
| testis | 6690 | 4892 | 3484 | 7566 | 7234 | 7710 | 4130 | 3833 | 6367 | 30007 |
| sphere | 4368 | 6231 | 2944 | 4843 | 4333 | 5493 | 3515 | 3731 | 37646 | 6367 |
| oocyte | 3110 | 4060 | 3562 | 3308 | 3111 | 3403 | 3078 | 17058 | 3731 | 3833 |
| dome | 3802 | 3279 | 3175 | 3780 | 3692 | 3839 | 16467 | 3078 | 3515 | 4130 |
| bud | 7358 | 4387 | 3221 | 7976 | 7014 | 26314 | 3839 | 3403 | 5493 | 7710 |
| 4dpf | 6785 | 3501 | 3043 | 8131 | 24440 | 7014 | 3692 | 3111 | 4333 | 7234 |
| 2dpf | 7350 | 3963 | 3197 | 25715 | 8131 | 7976 | 3780 | 3308 | 4843 | 7566 |
| 2cell | 3095 | 3438 | 14746 | 3197 | 3043 | 3221 | 3175 | 3562 | 2944 | 3484 |
| 256cell | 3627 | 29641 | 3438 | 3963 | 3501 | 4387 | 3279 | 4060 | 6231 | 4892 |
| 1dpf | 22219 | 3627 | 3095 | 7350 | 6785 | 7358 | 3802 | 3110 | 4368 | 6690 |

intersection: 30000, 20000, 10000

```
## pdf
##   2
```

I wanted to enxt compare the high confidence intergenic ends that we identify in the different stages. This is still very exploratory and I want to see the deviation in the data we have.

```r
completePath_intergenicEnds = paste0(completePath,"onlyIntergenic_90percent_n100.bed")
intergenicEnds = lapply(completePath_intergenicEnds,function(x) read.delim(x,stringsAsFactors = F,header
names(intergenicEnds) = stages



library(reshape)
nuMberOfIntergenicEnds = melt(lapply(intergenicEnds,nrow))
colnames(nuMberOfIntergenicEnds) = c("nuMberOfIntergenicEnds","stage")
write.table(nuMberOfIntergenicEnds,"/Volumes/groups/ameres/Pooja/Projects/zebrafishAnnotation/zebrafish_

library(ggplot2)
q = ggplot(nuMberOfIntergenicEnds,aes(x=stage,y=nuMberOfIntergenicEnds)) + geom_bar(stat = "identity")
print(q)
```

looks like there is a relaitonship between the number of ends and the number of intergenic ends. This could just be a function of the number of polyA reads used in the samples or the sequencing depth. So the initial threshold to identify priming sites will probably have to be set differently for different samples.

```r
path_preprocessing = paste0(path_allStages,stages,"/output/polyAmapping_allTimepoints/logs/")
preProcessingFile=c()
for(i in 1:length(path_preprocessing)){
  preProcessingFile =  c(preProcessingFile,paste0(path_preprocessing[i],list.files(path_preprocessing[i]

}

preProcessingStats = lapply(preProcessingFile,function(x) read.table(x,stringsAsFactors = F))
names(preProcessingStats) = stages

preProcessingStats= lapply(preProcessingStats,function(x) x[1:5,1])

preProcessingStats_split = lapply(preProcessingStats,function(x) strsplit(x,":",T))
preProcessingStats_split = lapply(preProcessingStats_split,function(x) lapply(x,function(y) y[2]))

preProcessingStats_split_melt = melt(preProcessingStats_split)
preProcessingStats_split_melt$value = as.numeric(as.character(preProcessingStats_split_melt$value))
sampleNames = c("initialFile","adapterTrimmed","fivePrimeTrimming","polyAcontaining","finalFile")
preProcessingStats_split_melt$sample = sampleNames
#preProcessingStats_split_melt$stage = rep(stages,each = 5)
library(dplyr)
```
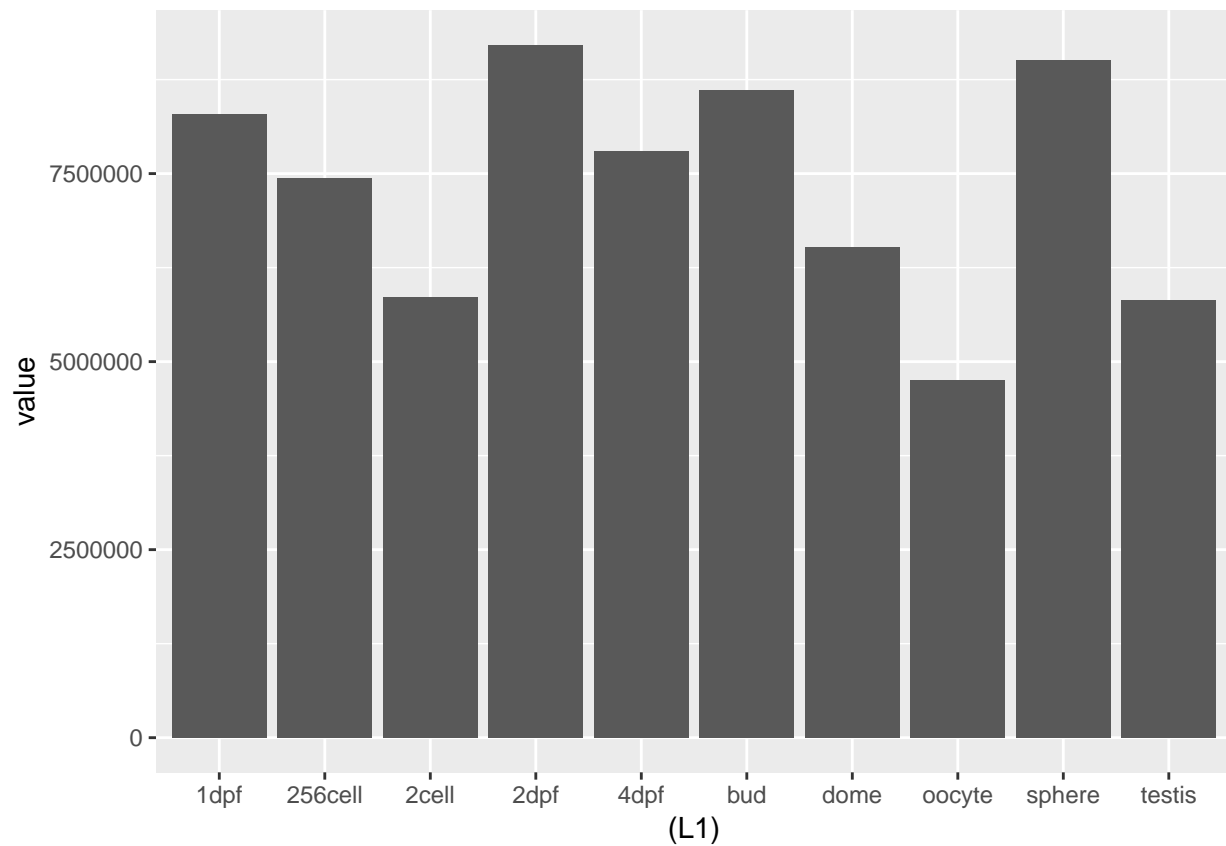
```
## Warning: package 'dplyr' was built under R version 3.4.1
```

```
##
## Attaching package: 'dplyr'

## The following object is masked from 'package:reshape':
##
##      rename

## The following objects are masked from 'package:stats':
##
##      filter, lag

## The following objects are masked from 'package:base':
##
##      intersect, setdiff, setequal, union
```

```r
filalFile = preProcessingStats_split_melt %>% filter(sample=="finalFile")
ggplot(filalFile,aes(x=(L1),y=value)) + geom_bar(stat="identity")
```



It looks like the number of ends we identify is based directly on the polyA readthrough. So we must either :

1. Normalize the cutoff to the number of polyA reads.
2. Consider even 1 read as indication of a priming site.