

# comparing counting windows from different stages

*Pooja Bhat*

*July 19, 2017*

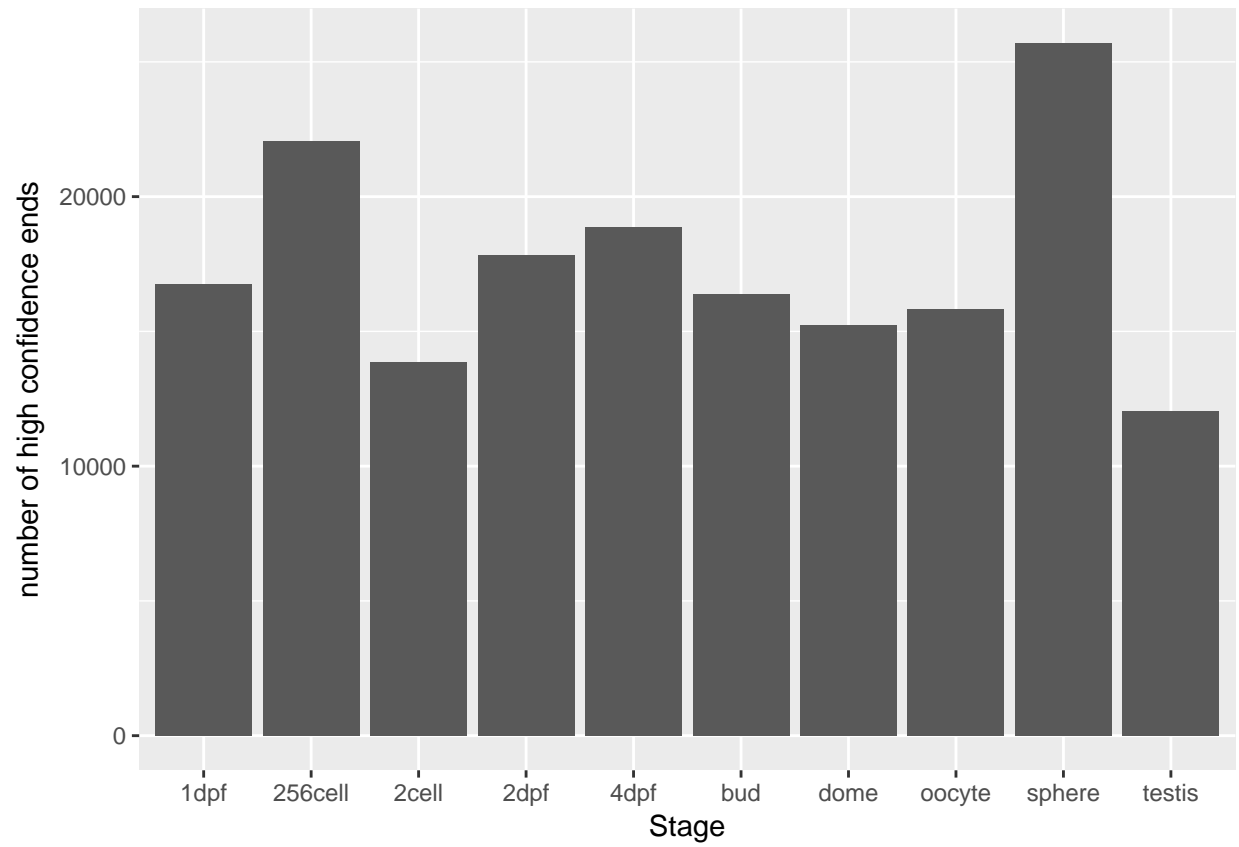
## R Markdown

```
## Warning: package 'dplyr' was built under R version 3.4.1
##
## Attaching package: 'dplyr'
##
## The following objects are masked from 'package:stats':
##
##   filter, lag
##
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
## -----
## You have loaded plyr after dplyr - this is likely to cause problems.
## If you need functions from both plyr and dplyr, please load plyr first, then dplyr:
## library(plyr); library(dplyr)
## -----
##
## Attaching package: 'plyr'
##
## The following objects are masked from 'package:dplyr':
##
##   arrange, count, desc, failwith, id, mutate, rename, summarise,
##   summarize
## Loading required package: stats4
## Loading required package: BiocGenerics
## Loading required package: parallel
##
## Attaching package: 'BiocGenerics'
##
## The following objects are masked from 'package:parallel':
##
##   clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,
##   clusterExport, clusterMap, parApply, parCapply, parLapply,
##   parLapplyLB, parRapply, parSapply, parSapplyLB
## The following objects are masked from 'package:dplyr':
##
##   combine, intersect, setdiff, union
## The following objects are masked from 'package:stats':
##
##   IQR, mad, sd, var, xtabs
```

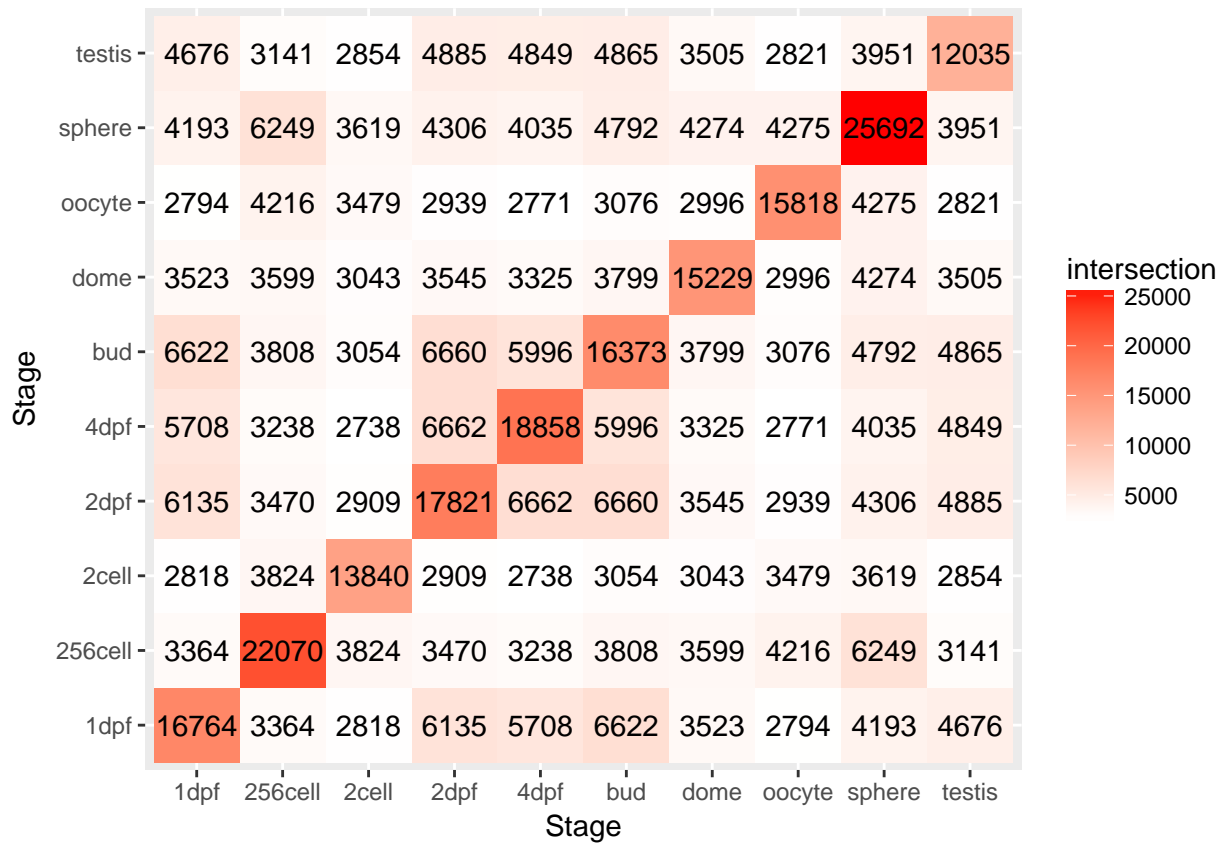
```

## The following objects are masked from 'package:base':
##
##   anyDuplicated, append, as.data.frame, cbind, colMeans,
##   colnames, colSums, do.call, duplicated, eval, evalq, Filter,
##   Find, get, grep, grepl, intersect, is.unsorted, lapply,
##   lengths, Map, mapply, match, mget, order, paste, pmax,
##   pmax.int, pmin, pmin.int, Position, rank, rbind, Reduce,
##   rowMeans, rownames, rowSums, sapply, setdiff, sort, table,
##   tapply, union, unique, unsplit, which, which.max, which.min
## Loading required package: S4Vectors
##
## Attaching package: 'S4Vectors'
## The following object is masked from 'package:plyr':
##
##   rename
## The following objects are masked from 'package:dplyr':
##
##   first, rename
## The following object is masked from 'package:base':
##
##   expand.grid
## Loading required package: IRanges
##
## Attaching package: 'IRanges'
## The following object is masked from 'package:plyr':
##
##   desc
## The following objects are masked from 'package:dplyr':
##
##   collapse, desc, slice
## Loading required package: GenomeInfoDb
##
## Attaching package: 'reshape'
## The following objects are masked from 'package:S4Vectors':
##
##   expand, rename
## The following objects are masked from 'package:plyr':
##
##   rename, round_any
## The following object is masked from 'package:dplyr':
##
##   rename

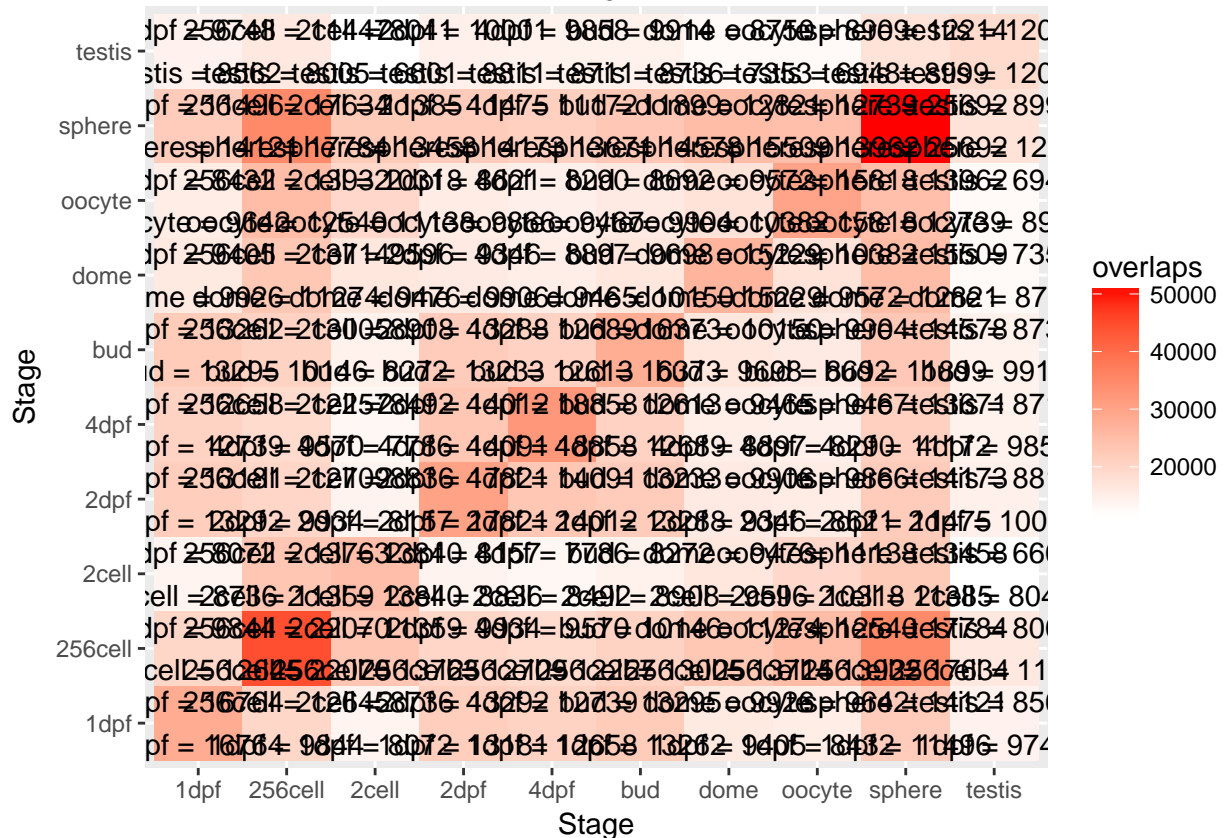
```



## pdf  
## 2



## pdf  
## 2



## pdf

```
## 2
```

```
completePath_countingWindows = paste0(completePath,"allAnnotations.bed")

allCountingWindows = lapply(completePath_countingWindows,function(x) read.delim(x,stringsAsFactors = F,
names(allCountingWindows) = stages
allCombinations$queryOverlaps_CW = NA
allCombinations$subjectOverlaps_CW = NA
allCombinations$querySubject_CW = NA
for(i in 1:nrow(allCombinations)){

  sample1 = allCountingWindows[which(stages == allCombinations[i,1])]
  sample2 = allCountingWindows[which(stages == allCombinations[i,2])]
  granges_1 = makeGRangesFromDataFrame(df = sample1,keep.extra.columns = T,seqnames.field = "V1",start.

  granges_2 = makeGRangesFromDataFrame(df = sample2,keep.extra.columns = T,seqnames.field = "V1",start.

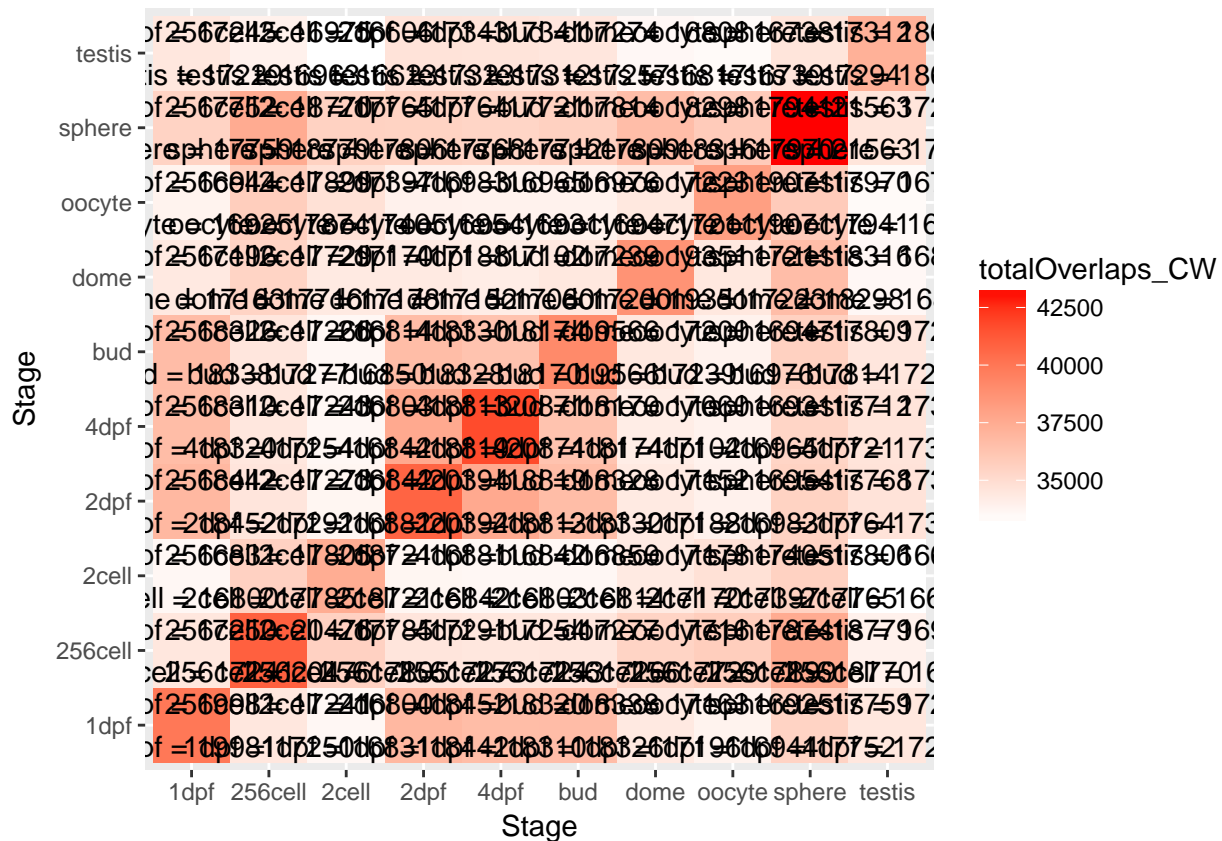
  sample1_overlap = sample1[[1]][queryHits(findOverlaps(granges_1,granges_2)),]
  sample1_overlap= sample1_overlap[!duplicated(sample1_overlap),]
  sample2_overlap = sample2[[1]][subjectHits(findOverlaps(granges_1,granges_2)),]
  sample2_overlap= sample2_overlap[!duplicated(sample2_overlap),]
  allCombinations$queryOverlaps_CW[i] = nrow(sample1_overlap)
  allCombinations$subjectOverlaps_CW[i] = nrow(sample2_overlap)
  allCombinations$querySubject_CW[i] = paste(paste(names(sample1),"=",nrow(sample1_overlap)),paste(n

}

allCombinations$totalOverlaps_CW = allCombinations$queryOverlaps_CW + allCombinations$subjectOverlaps_CW

overlapCountingWindows = ggplot(allCombinations, aes(Var1, Var2)) + geom_tile(aes(fill = totalOverlaps_CW))

print(overlapCountingWindows)
```



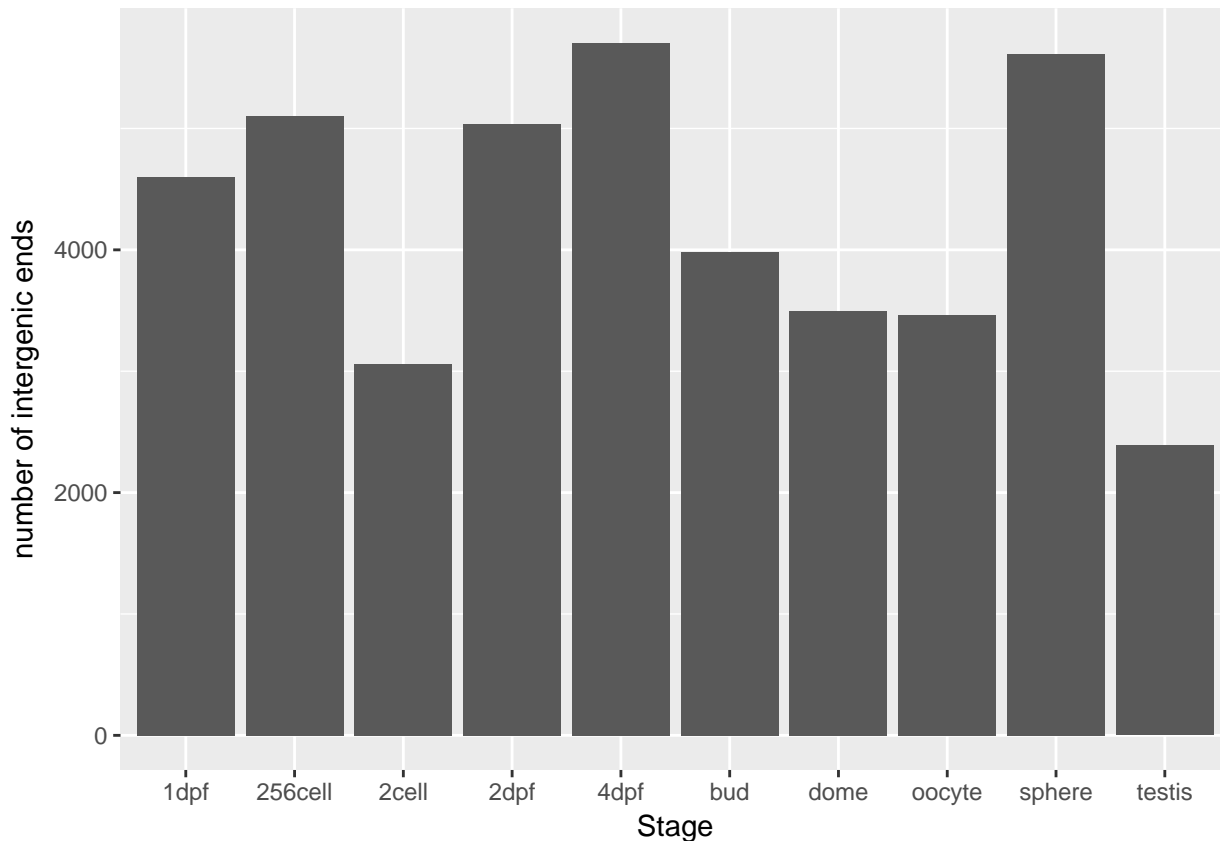
So there is a high degree of overlap in the final counting windows.

I wanted to next 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.

```
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)
numOfIntergenicEnds = melt(lapply(intergenicEnds,nrow))
colnames(numOfIntergenicEnds) = c("numOfIntergenicEnds","stage")
write.table(numOfIntergenicEnds,"/Volumes/groups/amerres/Pooja/Projects/zebrafishAnnotation/zebrafish.

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



looks like there is a relationship 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.

```
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],

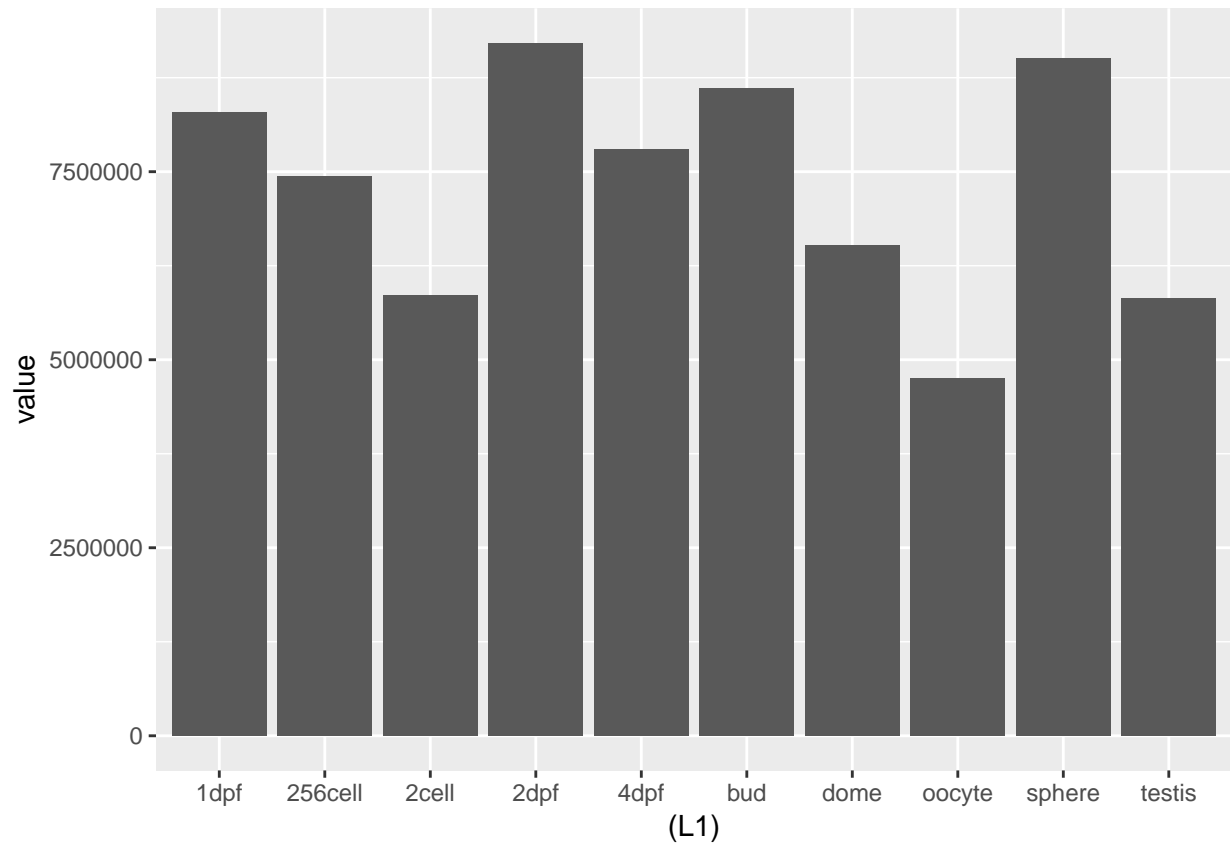
})

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

preProcessingStats_split = lapply(preProcessingStats,function(x) strsplit(x$V1,":",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)
finalFile = preProcessingStats_split_melt %>% filter(sample=="finalFile")
ggplot(finalFile,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.