

# Rotation in Murre lab

# Building Riboprofiling Pipeline

Pin-Chung (Tony) Cheng

# Workflow

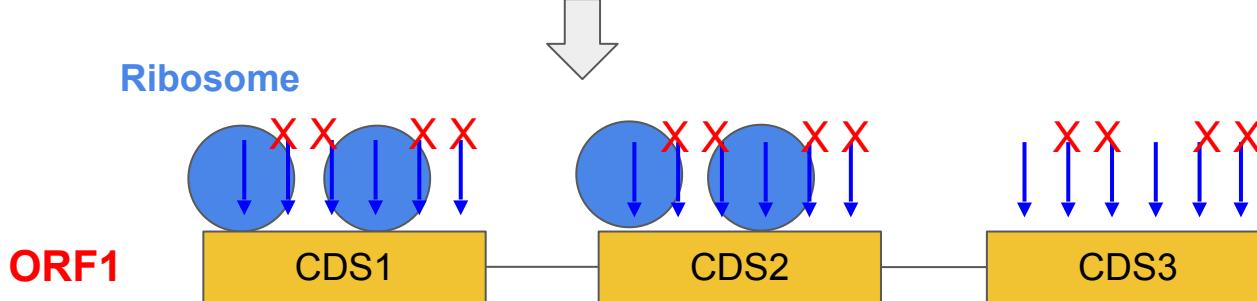
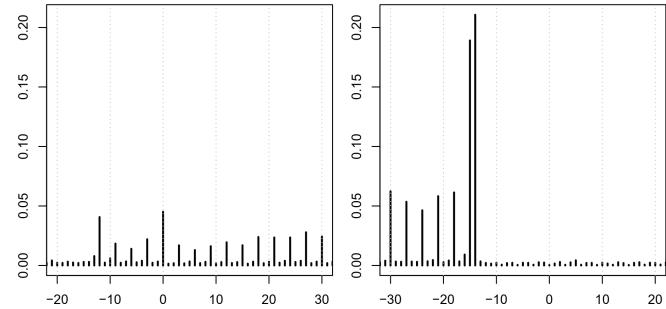


offsetCorrect.sam -> filtered reads files  
repre\*.gtf -> predicted ORF file

- Demultiplex
- Read alignment
- Filter rRNA reads
- Filter read length

...

...



# Pre-filter reads on terminal

```
41 # Step 0: Run on terminal
42
43 # Create read bed file from sam file
44 sh_name = paste0(script.basename, "/Codes/sam_to_bed_strand_v2.sh")
45 read_basename = sub(".sam", "", basename(opt$file))
46 pre_bed_name = paste0(script.basename, "/Temp/", read_basename, "_pre_bedtools_intersect.tsv")
47
48 system(paste(sh_name, opt$file, pre_bed_name))
49
50 # Create gtf bed file and filter exon for pre bedtools intersect in linux
51 repre_gtf_bed = paste0(script.basename, "/Temp/", read_basename, "_repre_gtf2bed.bed")
52
53 # need a lot of \ to escape quotation mark
54 # uniq again, becasue orfID can have many same exon coordinate > falsely increase read count
55 system(paste("grep \\'exon\\s\\\\'\\'", opt$gtf, "| awk \\'{print $1 \\"\\t\\\" $4 \\"\\t\\\" $5 \\"\\t\\\" $7}\\' | uniq >", repre_gtf_bed))
56
57
58 # Run bedtools intersect, -s: match strand, -wa: if A intersect B, keep A
59 post_bed_name = paste0(script.basename, "/Temp/", read_basename, "_post_bedtools_intersect.tsv")
60
61 system(paste("bedtools intersect -a", pre_bed_name, "-b", repre_gtf_bed, ">", post_bed_name))
62
63
```

```

66 # Step 1: Read in "repre.valid.ORF.*gtf" and prepare for intersect
67
68 gtf = read.table(opt$gtf, sep="\t", stringsAsFactors=FALSE)
69
70 # Remove duplicated rows, if input is from multiple gtf file concatenated from multiple samples
71 gtf = distinct(gtf)
72
73 # Select the columns we need and filter only "CDS"
74 gtf_selected = gtf %>%
75   select(V1, V3, V4, V5, V8, V9, V7) %>%
76   filter(V3 == "CDS")
77
78
79 # extract orf_id
80 gtf_selected$V9 = str_match(gtf_selected$V9, "\\s(ENST[0-9]+.*:chr.+);\\stranscript_id")[,2]
81
82 colnames(gtf_selected) = c("chr", "CDS", "str", "end", "frame", "orf_id", "strand_col")
83
84 # remove CDS column
85 gtf_bed = gtf_selected[, -2]
86
87 # change to appropriate column data type
88 gtf_bed$frame = as.numeric(gtf_bed$frame)
89 gtf_bed$orf_id = as.character(gtf_bed$orf_id)
90
91 gtf_bed = arrange(gtf_bed, orf_id)
92
93
94
95 # -----
96 # Step 1.2: Create a merged ORF table for annotation and a dataframe for intersect
97
98 # find identical elements in a vector, return a list of positions(vector), need to sort the vector first
99 find_identical_element = function(x){
100
101   listv = list()
102
103   # first index
104   first_idx_v = which(!duplicated(x))
105   last_idx_v = c((first_idx_v - 1)[-1], length(x))

```

Read in gtf file,  
manipulate dataframe

```

149 # Step 1.3: Create Annotation table for merged_ORF to geneID
150
151 # Create mORF_table by adding column "geneID" "tx_id" "merged_orfID" > for annotation
152 tx_id_vec = str_extract(gtfo$orf_id, "ENST[0-9]+")
153
154 idx = match(tx_id_vec, anno$target_id)
155
156 mORF_table = gtf_bed %>%
157   mutate(tx_id = tx_id_vec, ens_gene = anno$sens_gene[idx], ext_gene = anno$ext_gene[idx],
158         merged_orfID = as.character(mORF_list[gtfo$orf_id]))
159
160
161 # Use distinct() to merge same merged_ORF (CDSs string)
162 mORF_table = mORF_table %>%
163   select(tx_id, ens_gene, ext_gene, merged_orfID) %>% # Note: merged_ORF can match many tx_id, here just show one
164   distinct(merged_orfID, .keep_all = TRUE) %>%
165   arrange(ext_gene)
166
167
168 # create a vector of number to add to geneID, eg. A2M.1, A2M.2, A2M.3, AAAS.1 ...
169 gene_count = mORF_table %>% count(ext_gene)
170
171 gene_nvec = c()
172 for (i in 1:nrow(gene_count)){
173
174   # create a vector for each geneID, starting from 1
175   v = 1:gene_count$n[i]
176
177   gene_nvec = append(gene_nvec, v)
178 }
179
180
181 # Add number to geneID
182 mORF_table = mORF_table %>%
183   mutate(ens_gene = paste0(ens_gene, ".", gene_nvec), ext_gene = paste0(ext_gene, ".", gene_nvec))
184
185 # Replace output NA.$number with NA
186 NA_idx = str_extract(mORF_table$ext_gene, "NA.*")
187 mORF_table$ens_gene[NA_idx] = NA
188
189

```

Different ORFs have same set of  
CDSs > create merge ORF

```

270 # -----
271 # Step 3: Intersect reads and mergeORF
272 # query (reads) is the first argument, subejct(ORF) the second
273
274 list = bedtools_intersect(reads, mORF_df, ignore_strand=FALSE)
275
276 overlap = list$overlap
277
278
279 # Use dplyr package, do all steps by piping:
280   # Add read_position, cds_start, cds_frame by index;
281   # Add remainder column by calculating (read_pos - cds_start - cds frame)/3 and output remainder
282   # Add tx_id column by index
283
284 result = overlap %>%
285   mutate(read_pos = reads[overlap$bed1_idx, 2], cds_start = mORF_df[overlap$bed2_idx, 2],
286         cds_frame = mORF_df[overlap$bed2_idx, 4], remainder = (read_pos - cds_start - cds_frame) %% 3,
287         orfID = mORF_df[overlap$bed2_idx, 5])
288
289 # Select result that are in frame
290 result_inframe = result %>%
291   filter(remainder == 0)
292
293
294 # Number of in-frame reads
295 paste0(" in_frame reads: ", nrow(result_inframe))
296 paste0(" total reads: ", nrow(result))
297
298
299 result_f1 = result %>%
300   filter(remainder==0) %>%
301   count(orfID)
302
303 colnames(result_f1) = c("merged_orfID", "f1Num")
304
305 result_f2 = result %>%
306   filter(remainder==1) %>%
307   count(orfID)
308
309 colnames(result_f2) = c("merged_orfID", "f2Num")
310

```

Intersect,  
 filtering read by frame f1, f2, f3,  
 count

# Run Rscript in one step

```
# Run steps on server terminal
```

```
# Copy package “Ribo_Tony” to your directory:
```

```
cp -r /home/tonycheng/Tools/Ribo_Tony PATH/TO/YOUR_DIRECTORY
```

```
# Create environment:
```

```
conda env create -f tony_r_env.yml
```

```
# Activate environment:
```

```
conda activate tony_r
```

# Run Rscript in one step

```
# Input: offset corrected sam file and re.pred.valid.gtf file (can be a concatenated gtf from multiple samples)
```

```
# Output: a table with Gene ID and read counts on codon frame 1, 2, 3 position
```

```
Rscript PATH/TO/Ribo_Tony/Intersect_mergeORF.R -f ${PATH/TO/*offsetCorrect.sam} -g ${PATH/TO/repre*.gtf} -o ${PATH/TO/OUTPUT_FILENAME.tsv}
```

Example code:

```
Rscript /home/tonycheng/Tools/Ribo_Tony/Intersect_mergeORF.R -f /home/tonycheng/Data/Fleur_data/sam/GSC2_t2.sam -g /home/tonycheng/Data/Fleur_data/gtf/GSC_total_uniq.repre.valid.ORF.gtf -o /home/tonycheng/Data/Fleur_data/GSC2_t2_mergeORF.tsv
```

# Prepare a concatenated gtf file

```
# Concatenate, sort, and remove duplicate row
```

```
cat ${gtf file1} ${gtf file2} ${gtf file3} | sort | uniq > ${output_filename}
```

```
# Create a directory with all the sample gtf files
```

```
mkdir ${folder_name}
```

Example code:

```
cat /PATH/TO/DIRECTORY/*.gtf | sort | uniq > GSC_total_uniq.repre.valid.ORF.gtf
```

\* Read counts on frame 1

Gene ID

Merged ORF: set of CDSs

f1Num	f2Num	f3Num	ens_gene	ext_gene	merged_orfID
23718	1739	2894	ENSG00000026025.2	VIM.2	17229423:17229985:0 17230650:17230710:1
:17236379:2 17237230:17237268:0 chr10;+					
21846	1520	2609	ENSG00000026025.11	VIM.11	17229423:17229985:0 17230650:17230710:1
13003	778	1421	ENSG00000026025.10	VIM.10	17229423:17229985:0 17230650:17230785:1
9816	791	1373	ENSG00000026025.8	VIM.8	17230633:17230710:0 17233587:17233682:0
9556	791	1338	ENSG00000026025.9	VIM.9	17230663:17230710:0 17233587:17233682:0
8355	2446	6811	ENSG00000222293.1	RNU2-36P.1	86423020:86423034:0 chr9;-
6464	485	929	ENSG00000026025.4	VIM.4	17230633:17230710:0 17233587:17233682:0
6204	485	894	ENSG00000026025.5	VIM.5	17230663:17230710:0 17233587:17233682:0
4947	411	704	ENSG00000026025.12	VIM.12	17233605:17233682:0 17233770:17233931:0
4896	6	774	ENSG00000062716.1	VMP1.1	59841265:59841285:0 chr17;+
4765	325	568	ENSG00000026508.7	CD44.7	35139304:35139370:0 35176575:35176740:2
:35201195:2 35201671:35201787:2 35204512:35204640:2 35206112:35206243:2 35208105:35208206:2 352					
29330:1 chr11;+					
4764	325	567	ENSG00000026508.3	CD44.3	35139304:35139370:0 35176575:35176740:2
:35201787:2 35204512:35204640:2 35206112:35206243:2 35208105:35208206:2 35209965:35210054:2 352					
4743	325	566	ENSG00000026508.8	CD44.8	35139304:35139370:0 35176575:35176740:2
:35211449:2 35214852:35214914:2 35219316:35219387:2 35221654:35221732:2 35229129:35229330:1 chr					
4741	325	566	ENSG00000026508.9	CD44.9	35139304:35139370:0 35176575:35176740:2
:35219387:2 35221654:35221732:2 35229129:35229330:1 chr11;+					
4723	325	566	ENSG00000026508.1	CD44.1	35139304:35139370:0 35176575:35176740:2
:35221732:2 35229129:35229330:1 chr11;+					
4376	342	523	ENSG00000135046.2	ANXA1.2	73158536:73158601:0 73158695:73158803:0
:73165209:0 73166097:73166192:2 73167497:73167555:2 73169032:73169154:0 73170051:73170104:0 chr					
4075	273	452	ENSG00000026508.17	CD44.17	35139304:35139370:0 35176575:35176740:2
4060	521	263	ENSG00000150991.7	UBC.7	124911717:124912180:2 124913321:1249137
3890	347	440	ENSG00000134531.2	EMP1.2	13211511:13211588:0 13213479:13213575:0
3789	284	606	ENSG00000026025.13	VIM.13	17230663:17230710:0 17233587:17233682:0
3715	222	468	ENSG00000111640.3	GAPDH.3	6534833:6534861:0 6536494:6536593:1 653
2;+					
3673	300	445	ENSG00000135046.1	ANXA1.1	73158794:73158803:0 73159329:73159423:2
:73166192:2 73167497:73167555:2 73169032:73169154:0 73170051:73170104:0 chr9;+					
3620	693	485	ENSG00000242265.15	PEC10.13	94656587:94656591:0 94663334:94

# Differential Expression Analysis

# Fluer's Data

# Two patient sample: GSC2, GSC20

# Three conditions: t=0, t=1, t=2

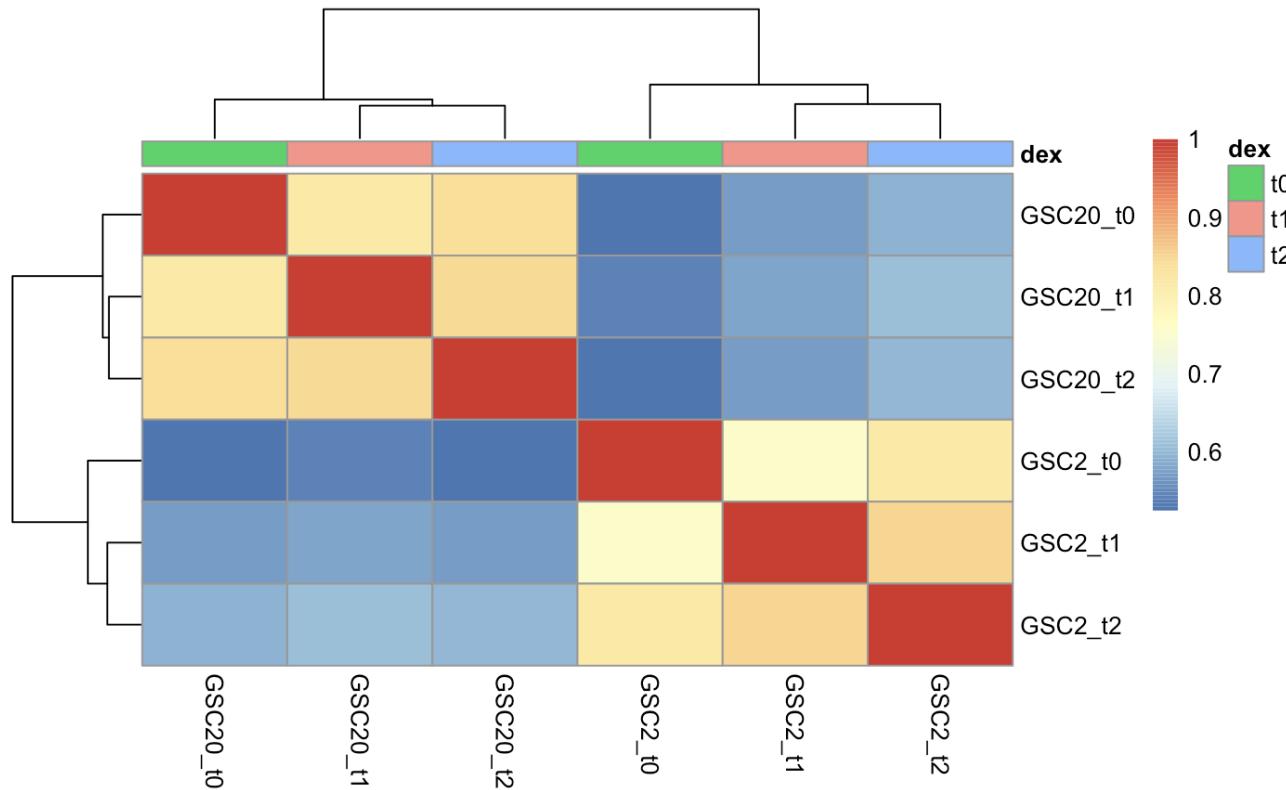
Gene <chr>	GSC2_t0 <int>	GSC2_t1 <int>	GSC2_t2 <int>	GSC20_t0 <int>	GSC20_t1 <int>	GSC20_t2 <int>
VIM.2	3611	5900	23718	2071	3315	4103
VIM.11	3418	5561	21846	2014	3133	4006
VIM.10	2192	4006	13003	1471	2242	2986
RNU2-57P.1	1466	1570	2486	NA	NA	NA
RNU2-26P.1	1415	1620	2349	NA	0	NA
VIM.8	1270	1775	9816	592	898	1165
VIM.9	1245	1745	9556	557	883	1120
RNU2-36P.1	883	770	8355	1	32	2
VIM.4	772	1129	6464	400	610	813
VIM.5	747	1099	6204	365	595	768

1-10 of 38,707 rows

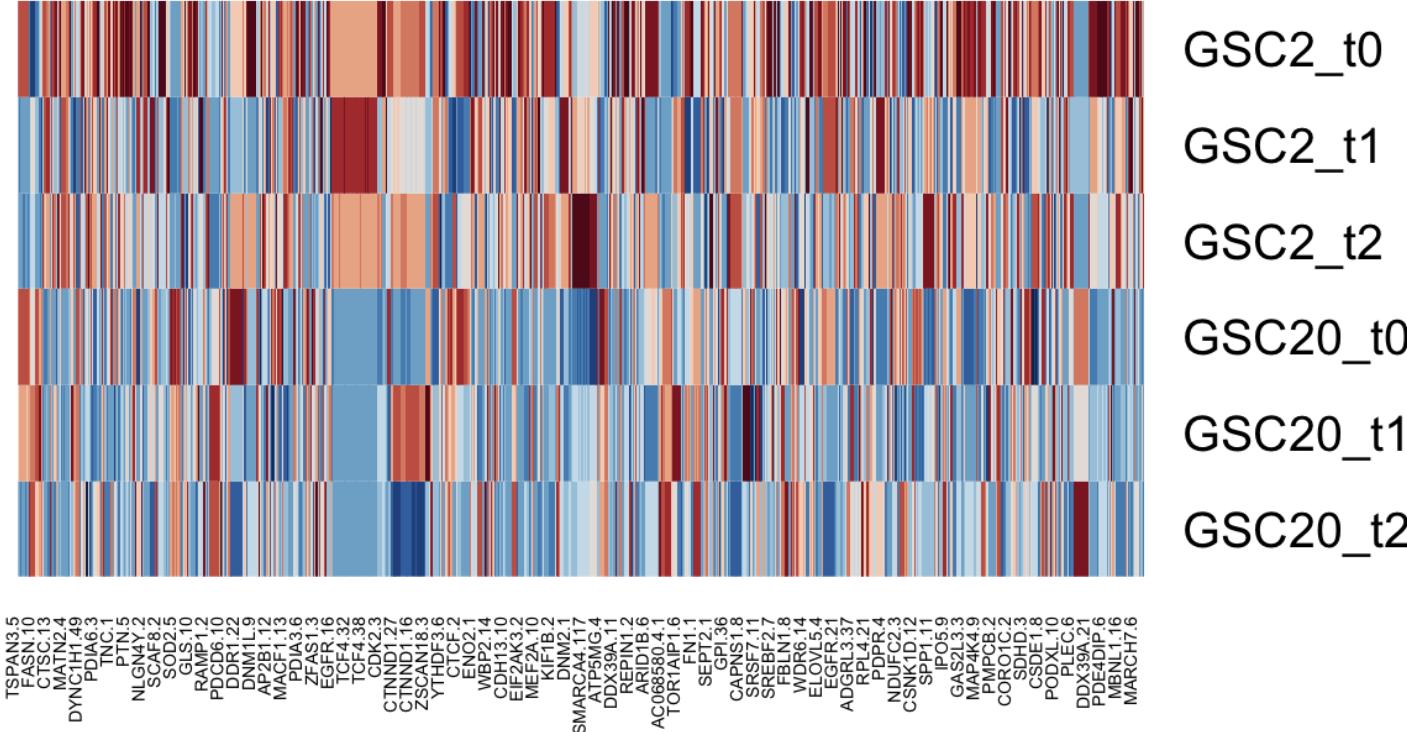
Previous 1 2 3 4 5 6 ... 100 Next

# Heatmap

Samples are more correlated between patients rather than conditions



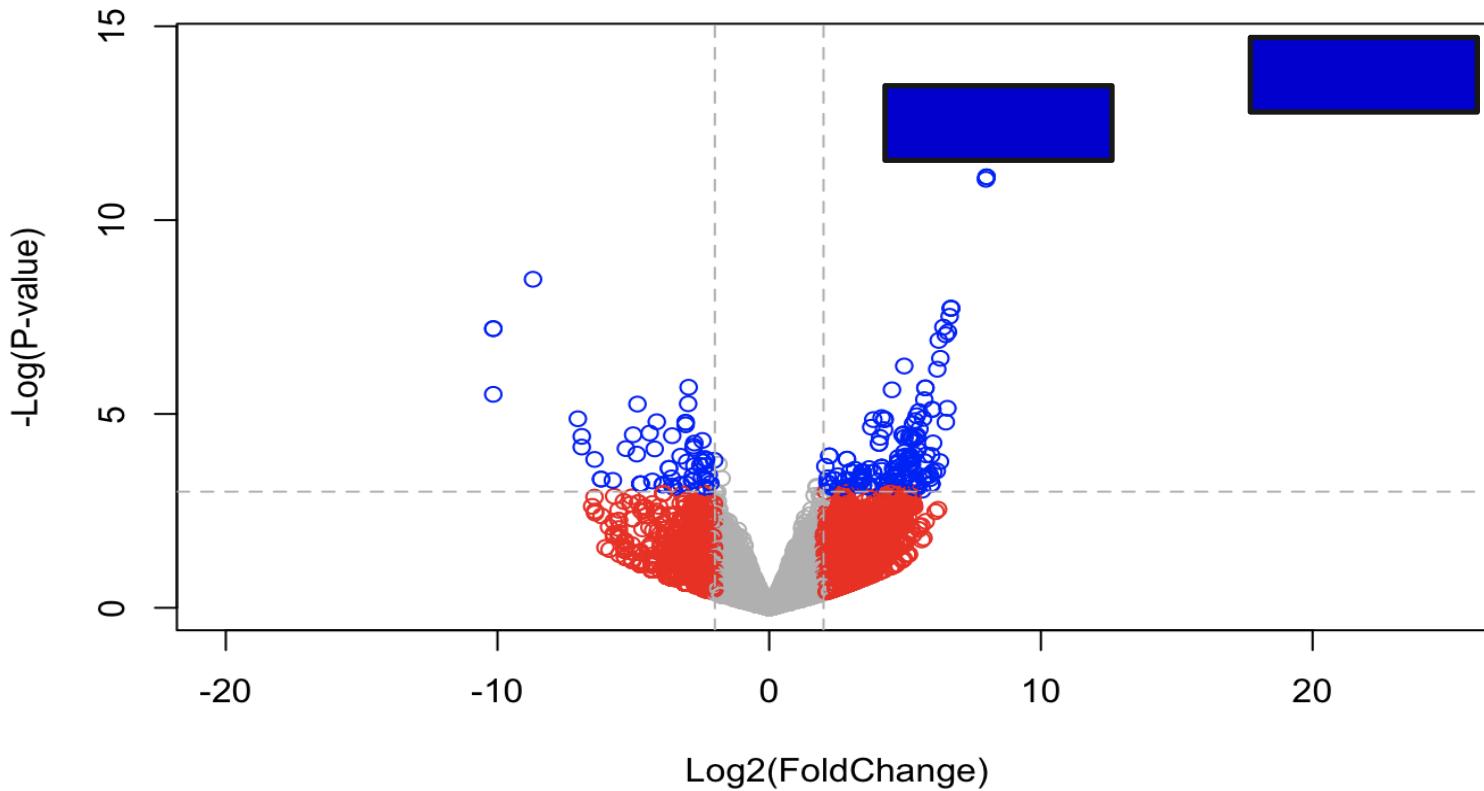
# Heatmap



# t2 vs t0: order by log2FoldChange

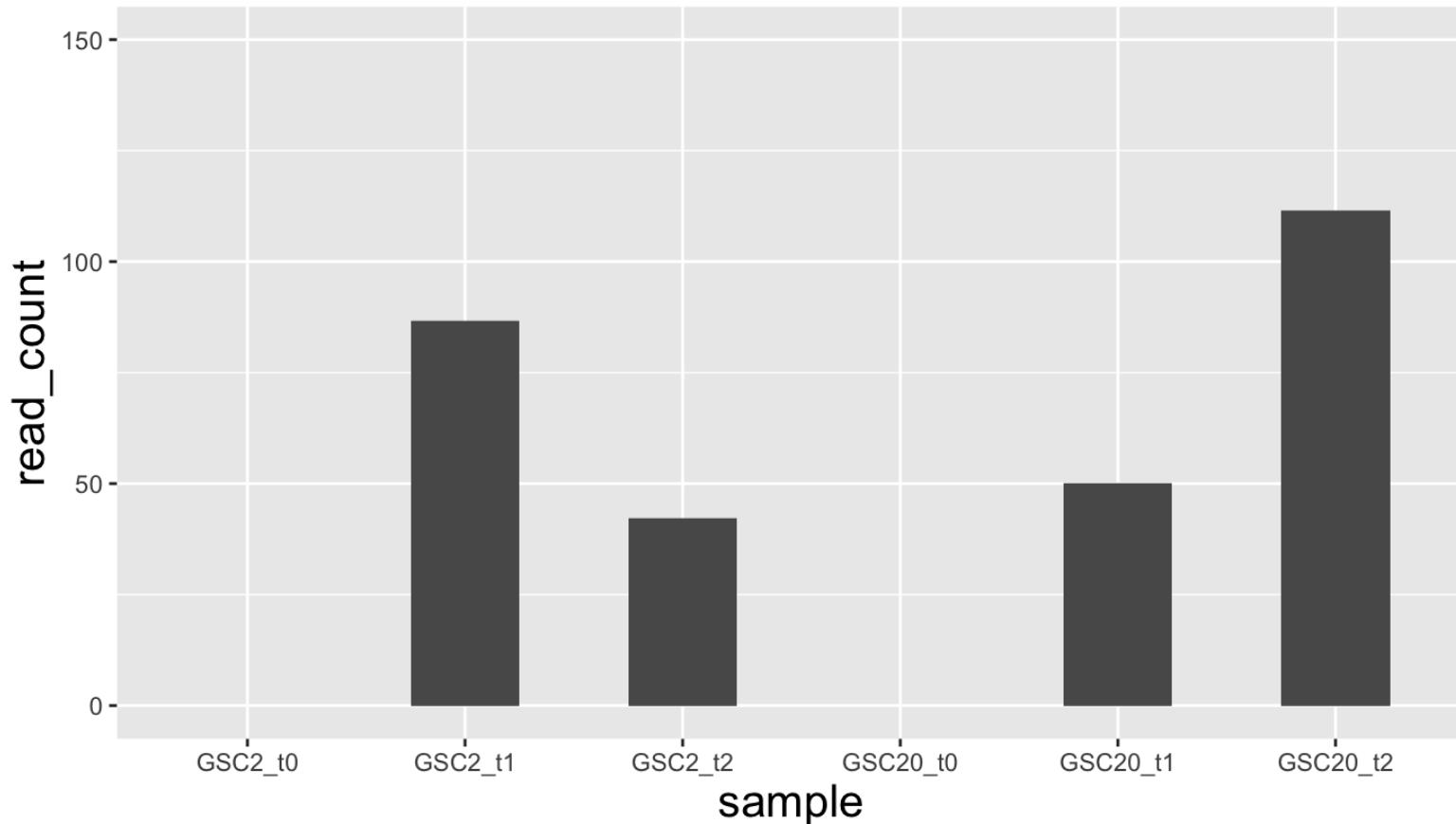
	baseMean <dbl>	log2FoldChange <dbl>	IfcSE <dbl>	stat <dbl>	pvalue <dbl>	padj <dbl>
	275.19264	24.027155	4.784769	5.021591	5.124517e-07	0.01978679
	48.43936	7.999540	1.847291	4.330417	1.488272e-05	0.15295725
	48.43936	7.999540	1.847291	4.330417	1.488272e-05	0.15295725
	46.13576	7.977077	1.848003	4.316593	1.584557e-05	0.15295725
	14.37565	6.691227	1.904044	3.514218	4.410506e-04	0.99999629
	14.37565	6.691227	1.904044	3.514218	4.410506e-04	0.99999629
	14.37565	6.691227	1.904044	3.514218	4.410506e-04	0.99999629
	15.65169	6.640647	1.919939	3.458780	5.426281e-04	0.99999629
	16.89030	6.578488	1.965008	3.347818	8.145057e-04	0.99999629
	10.20756	6.559565	2.378785	2.757527	5.824039e-03	0.99999629

# Volcano Plot





## Read Count



```
rep("Thank you!", 100)
```

