

ANALYSIS OF GENES EXPRESSED IN SACCHAROMYCES CEREVISIAE IN AEROBIC AND ANAEROBIC CONDITIONS

SUBHAM DAS - 18399

Indian Institute of Science Education and Research, Bhopal

Abstract- RNA-seq of Baker's Yeast was analyzed using bioinformatics tools developed over the years to study differential gene expression under two (or more) different conditions. In this study, we analyze the gene expression under aerobic and anaerobic conditions. We carefully study two differentially expressed genes, **THI72** and **YPS1**, which change in their expression levels going from one condition to another. We delve into the details of the procedure and observe how the files are being analyzed, and then we visually study the output of our two genes of interest to understand their functioning. We discuss the output and future implications of our observations.

PURPOSE: To analyze the differentially expressed genes in baker's yeast (*saccharomyces cerevisiae*) and study their trend of expression by identifying two significant genes of interest of which one is upregulated and the other downregulated in the fermentation process.

I. INTRODUCTION

Here we describe a transcriptomic analysis of fermentation of *Saccharomyces cerevisiae*. The dataset was obtained from an experiment where fermentation was carried out in carefully controlled environmental conditions in a bioreactor to reduce transcriptomic responses that would be due to factors other than the oxygen levels. Further, we analyze two genes of interest - **THI72** and **YPS1**, the former being up-regulated and the latter being down-regulated from aerobic to anaerobic condition. Using the protocol of **Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks** by Trapnell et al., we compare gene and transcript expression under the two given conditions. Finally, CummeRbund - a tool for visualizing RNA-seq analysis results, Integrative Genomics Viewer (IGV) - to observe the coverage of a gene (whether up or down-regulated) in the two conditions, and Gene Ontology Profiling - enrichment analysis to see how the differentially expressed genes contribute in different biological processes, molecular functions, cellular compartments, and pathways they are involved in.

II. METHODS USED TO ANALYZE RNA-seq AND PIPELINE OF PROCEDURE

RNA-seq experiments must be analyzed with robust, efficient and statistically principled algorithms. In view of this the following tools have been utilized to study our dataset:

1. Tophat2
2. Bowtie2
3. Cufflinks (cuffmerge, cuffdiff)
4. CummeRbund library in R
5. Integrative Genome browser (IGV)
6. Gene Ontology Profiling (g:profiler)

The RNA-seq dataset was obtained by NextSeq 500 paired end sequencing. Further details of the dataset include:

Platform - Illumina
Model - NextSeq 500

Layout - Paired
Source – Transcriptomic

The dataset (FASTQ files) was obtained from ENA browser, extracted and analysis performed.

The following processes/commands were run and relevant output files obtained:

1. TopHat - Reads for each condition are mapped to the reference genome. After running TopHat, the resulting alignment files are provided to Cufflinks.

OUTPUT: acceptedhits.bam (BAM files)

2. Cufflinks - Used to generate a transcriptome assembly for each condition.

OUTPUT: Transcripts (GTF files)

3. Cuffmerge - Assemblies are then merged together using the Cuffmerge utility, which is included with the Cufflinks package. This merged assembly provides a uniform basis for calculating gene and transcript expression in each condition.

OUTPUT: merged.gtf (GTF file)

4. Cuffdiff - The reads and the merged assembly are fed to Cuffdiff, which calculates expression levels and tests the statistical significance of observed changes. Cuffdiff also performs an additional layer of differential analysis. By grouping transcripts into biologically meaningful groups (such as transcripts that share the same transcription start site (TSS)), Cuffdiff identifies genes that are differentially regulated at the transcriptional or post-transcriptional level.

OUTPUT: diff_out folder with FPKM_TRACKING, COUNT_TRACKING, READ_TRACKING, and DIFF files. These are a set of tabular files which can be indexed and visualized in Cummebund.

5. CummeRbund - Provides functions for creating commonly used expression plots such as volcano, scatter and box plots.

The entire process pipeline can be summarized by the flowchart:

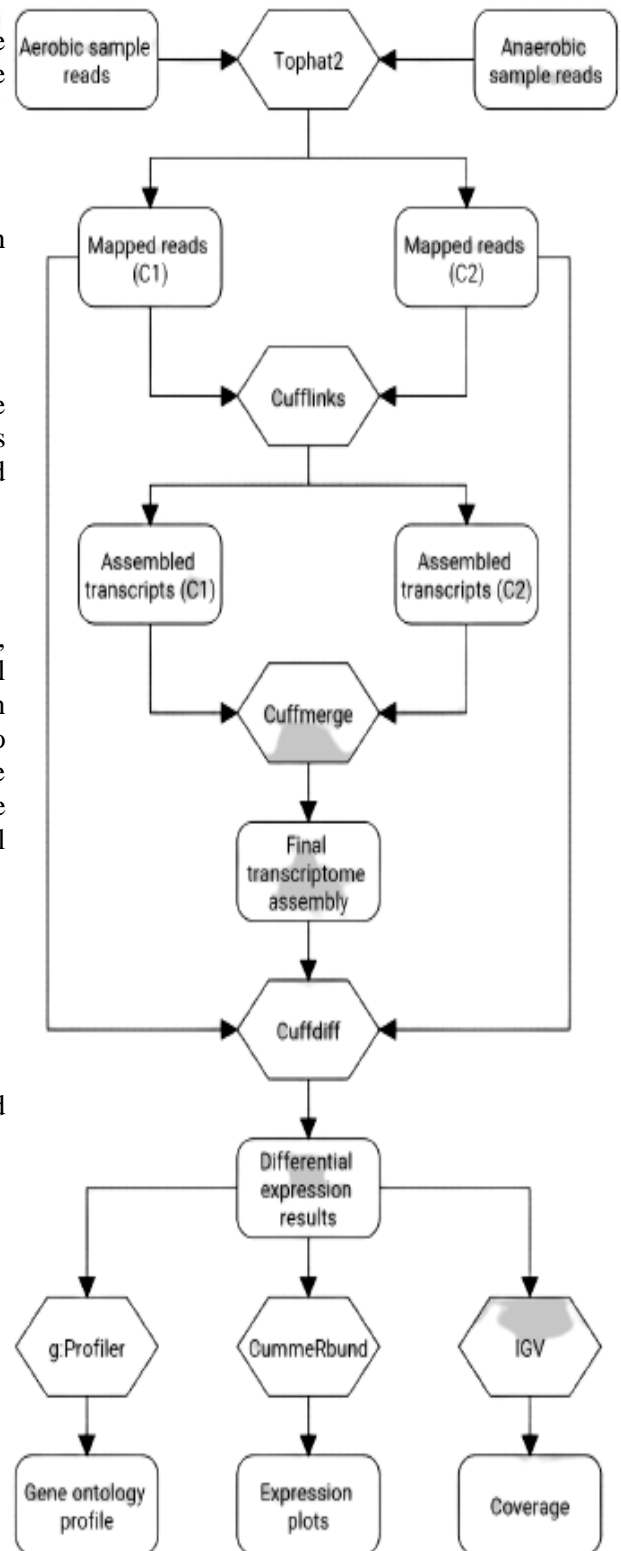


Fig.1 | An overview of the protocol.

III. RESULTS

The analysis provided a list of 6516 differentially expressed genes of which 1093 were significantly expressed ($P < 0.05$) with $\text{abs}(\log_2\text{fold_change}) > 1$.

```
> cuff_data ~  
Cuffset instance with:  
  2 samples  
 6516 genes  
 8057 isoforms  
 7741 TSS  
 6643 CDS  
 6516 promoters  
 7741 splicing  
 6050 relCDS  
  
> |
```

Fig.2 | Total expressed genes

Other significant data output of coding sequence (CDS), regulated genes (relCDS), isoform, promoter, splicing, and transcription start site (TSS) is as follows:

```
> library(cummeRbund)  
> cuff_data <- readCufflinks('diff_out')  
> gene_diff_data <- diffData(genes(cuff_data))  
> sig_gene_data <- subset(gene_diff_data, (significant == 'yes'))  
> nrow(sig_gene_data)  
[1] 3893  
> isoform_diff_data <- diffData(isoforms(cuff_data), 'c1', 'c2')  
> sig_isoform_data <- subset(isoform_diff_data, (significant == 'yes'))  
> nrow(sig_isoform_data)  
[1] 4333  
> tss_diff_data <- diffData(TSS(cuff_data), 'c1', 'c2')  
> sig_tss_data <- subset(tss_diff_data, (significant == 'yes'))  
> nrow(sig_tss_data)  
[1] 4304  
> cds_diff_data <- diffData(CDS(cuff_data), 'c1', 'c2')  
> sig_cds_data <- subset(cds_diff_data, (significant == 'yes'))  
> nrow(sig_cds_data)  
[1] 4007  
> promoter_diff_data <- distvalues(promoters(cuff_data))  
> sig_promoter_data <- subset(promoter_diff_data, (significant == 'yes'))  
> nrow(sig_promoter_data)  
[1] 232  
> splicing_diff_data <- distvalues(splicing(cuff_data))  
> sig_splicing_data <- subset(splicing_diff_data, (significant == 'yes'))  
> nrow(sig_splicing_data)  
[1] 31  
> relCDS_diff_data <- distvalues(relCDS(cuff_data))  
> sig_relCDS_data <- subset(relCDS_diff_data, (significant == 'yes'))  
> nrow(sig_relCDS_data)  
[1] 132  
> sig_gene_data <- subset(gene_diff_data, (significant == 'yes'),  $\text{abs}(\log_2\text{fold\_change}) > 1$ )  
> nrow(sig_gene_data)  
[1] 3893  
> final_sig_gene_data <- subset(sig_gene_data,  $\text{abs}(\log_2\text{fold\_change}) > 1$ )  
> nrow(final_sig_gene_data)  
[1] 1093  
> |
```

Fig.3 | Significantly expressed genes

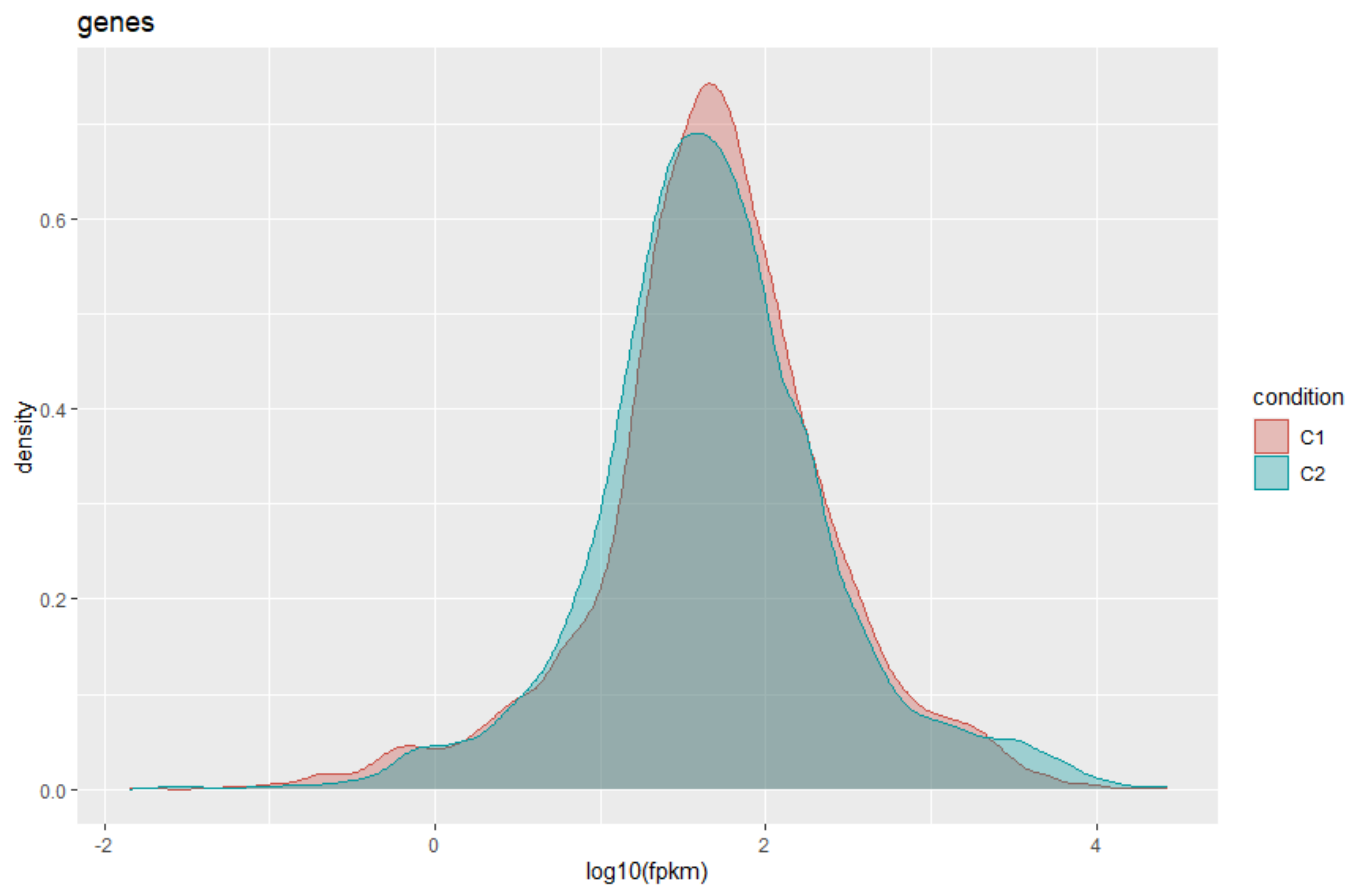


Fig.4 | CummeRbund plot of the expression level distribution for all genes in experimental conditions C1 (Aerobic condition) and C2 (Anaerobic condition). FPKM, fragments per kilobase of transcript per million fragments mapped.

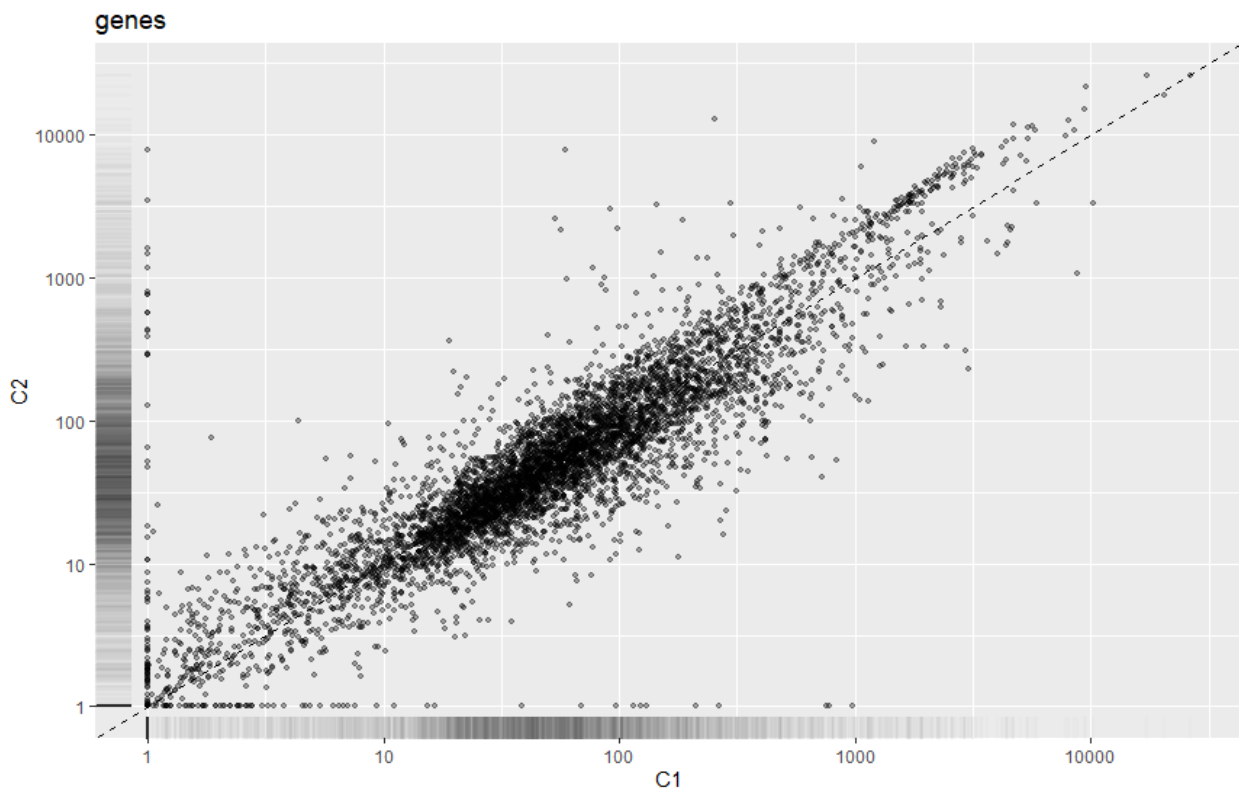


Fig.5 | CummeRbund scatter plots highlight general similarities and specific outliers of genes between aerobic and anaerobic conditions.

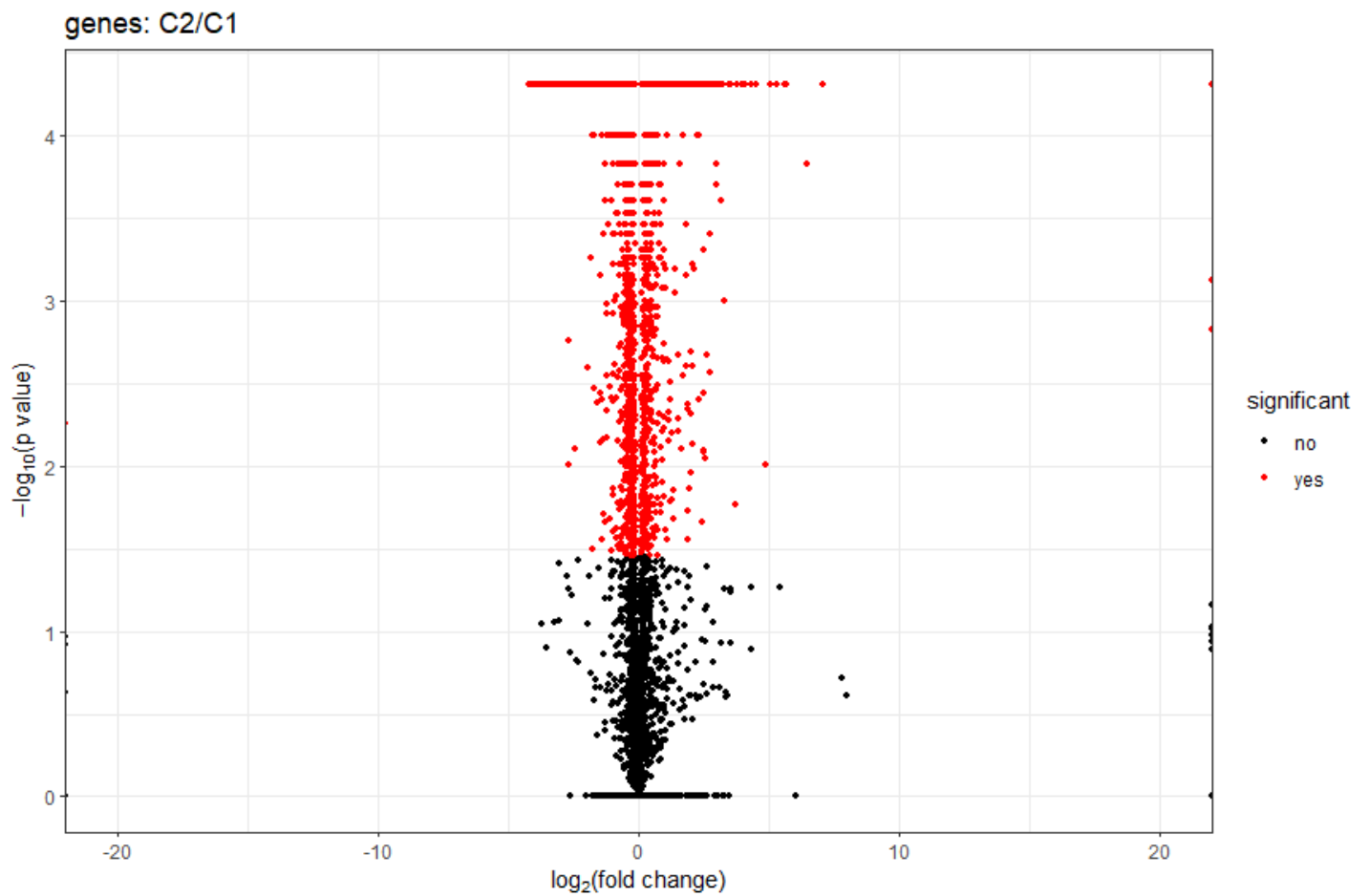


Fig.6 | CummeRbund volcano plot reveal genes that differ significantly between the pairs of conditions.

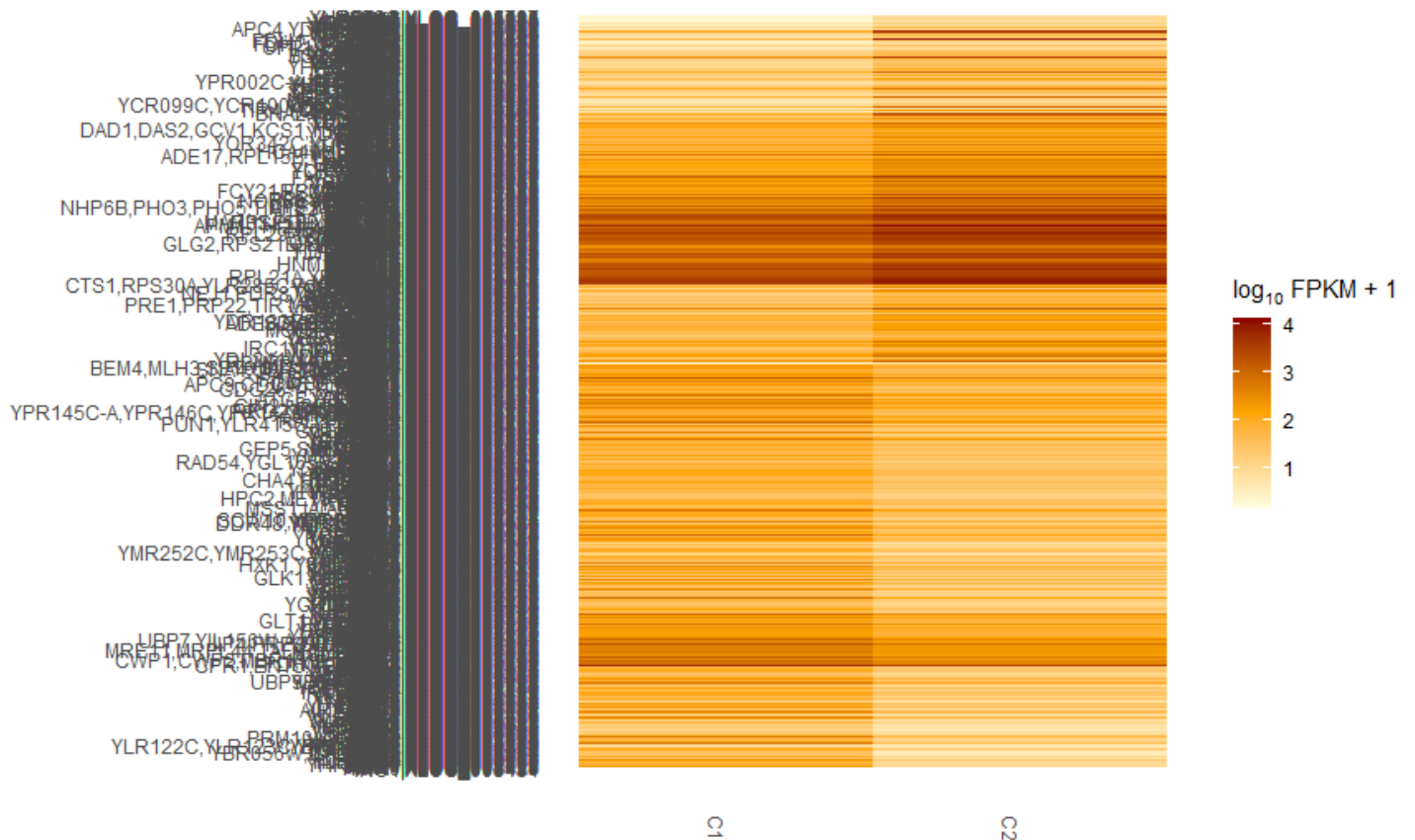


Fig.7 | Heat map of differentially expressed genes

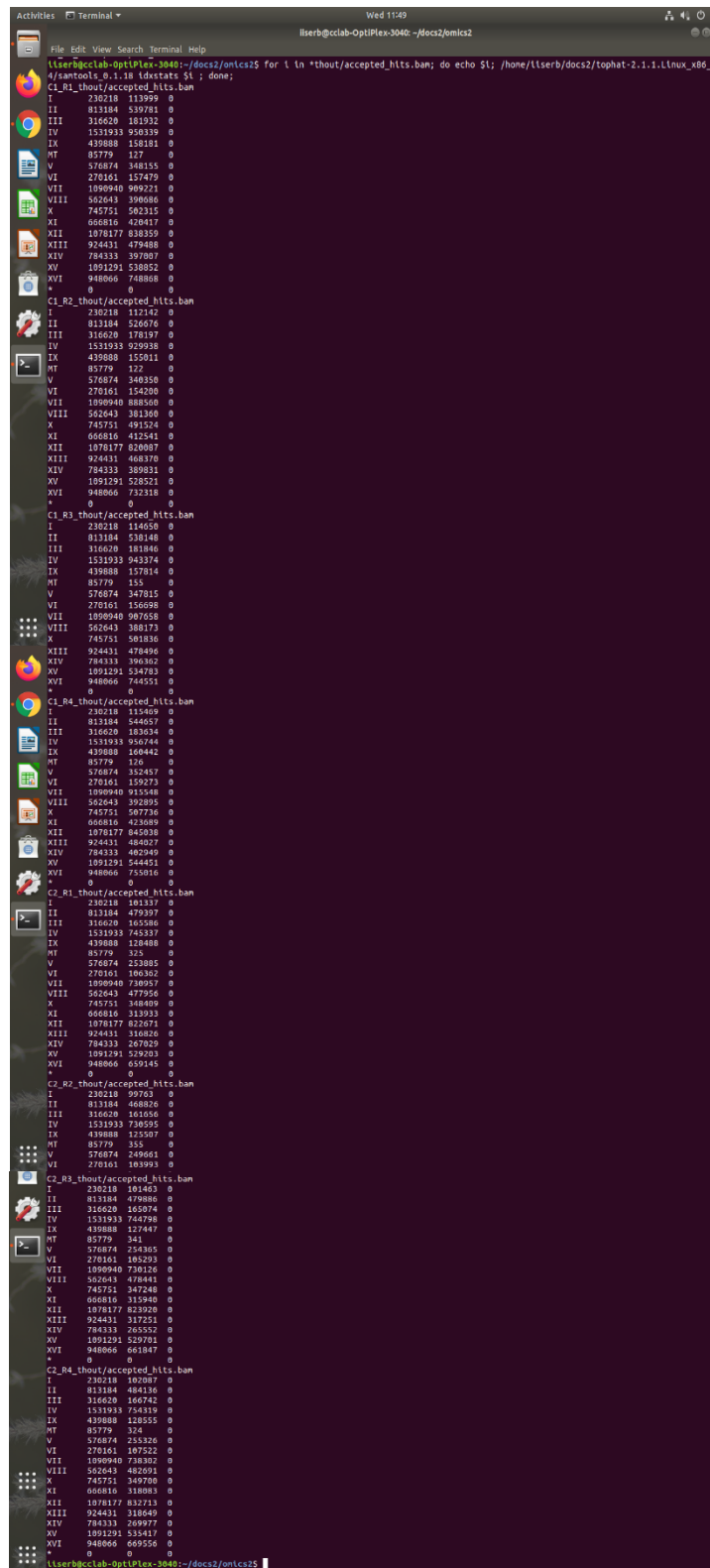


Fig.8 | Report of the number of fragments that map to each chromosome.

After the global gene differentiation visualization was performed, we used the output file: **gene_exp.diff** to sort the log2(fold_change) column, which is the effect size estimate. This value indicates how much the gene or transcript's expression seems to have changed between the two conditions. A high negatives value indicates a highly down-regulated gene, while a positive value indicates an up-regulated gene. Using this data two genes of interest were selected **THI72** and **YPS1**, and further analyzed.

The expression bar plot, isoform bar, line plot, and read coverage from IGV are as follows:

(1) THI72

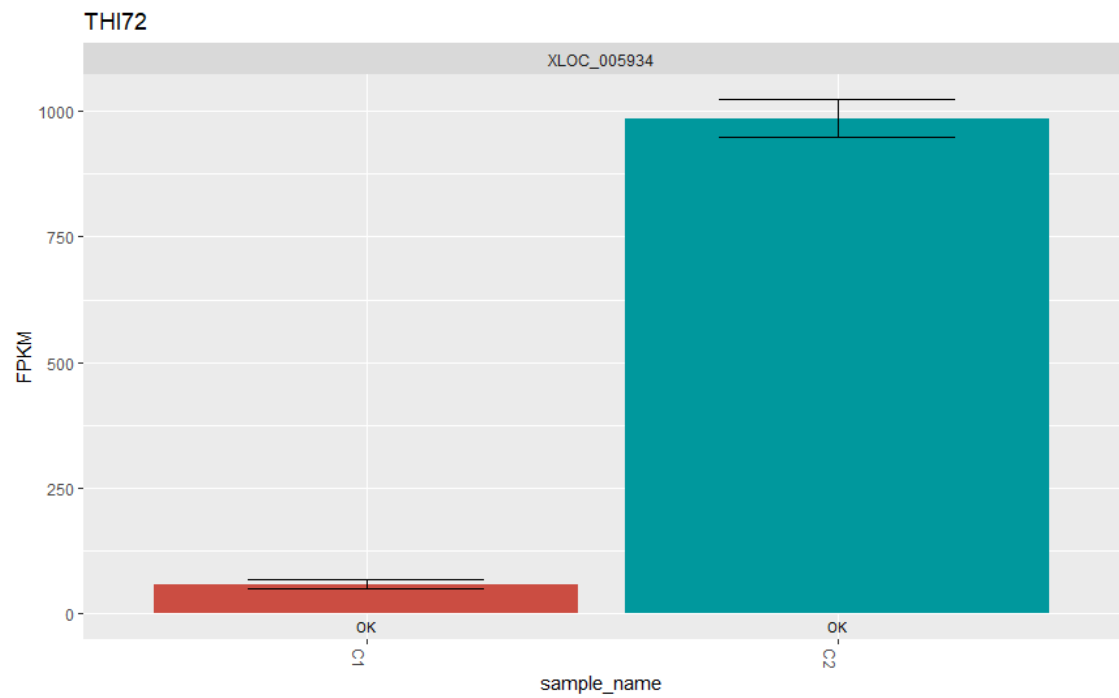


Fig.9 | Differential analysis results. Expression plot shows clear differences in the expression of THI72 across conditions C1 (aerobic) and C2 (anaerobic), measured in FPKM

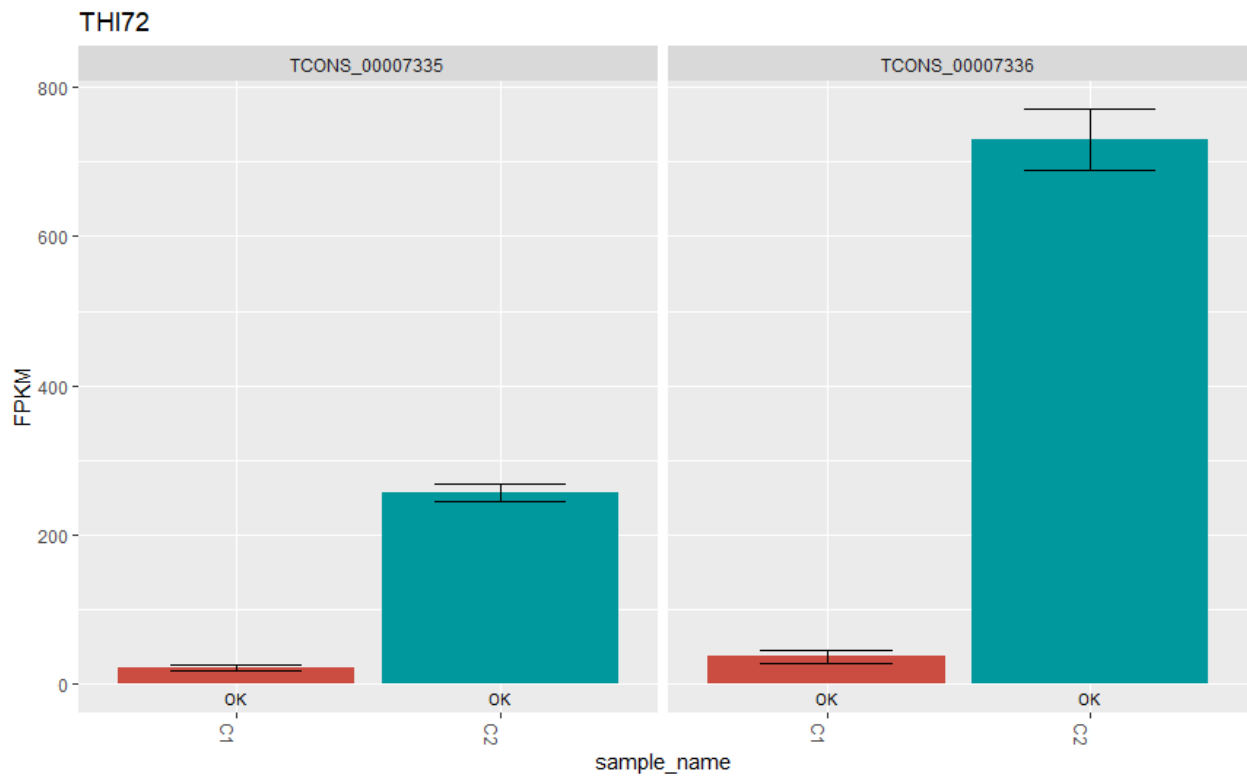


Fig.10 | Changes in THI72 expression are attributable to a large increase in the expression of one of two alternative isoforms.

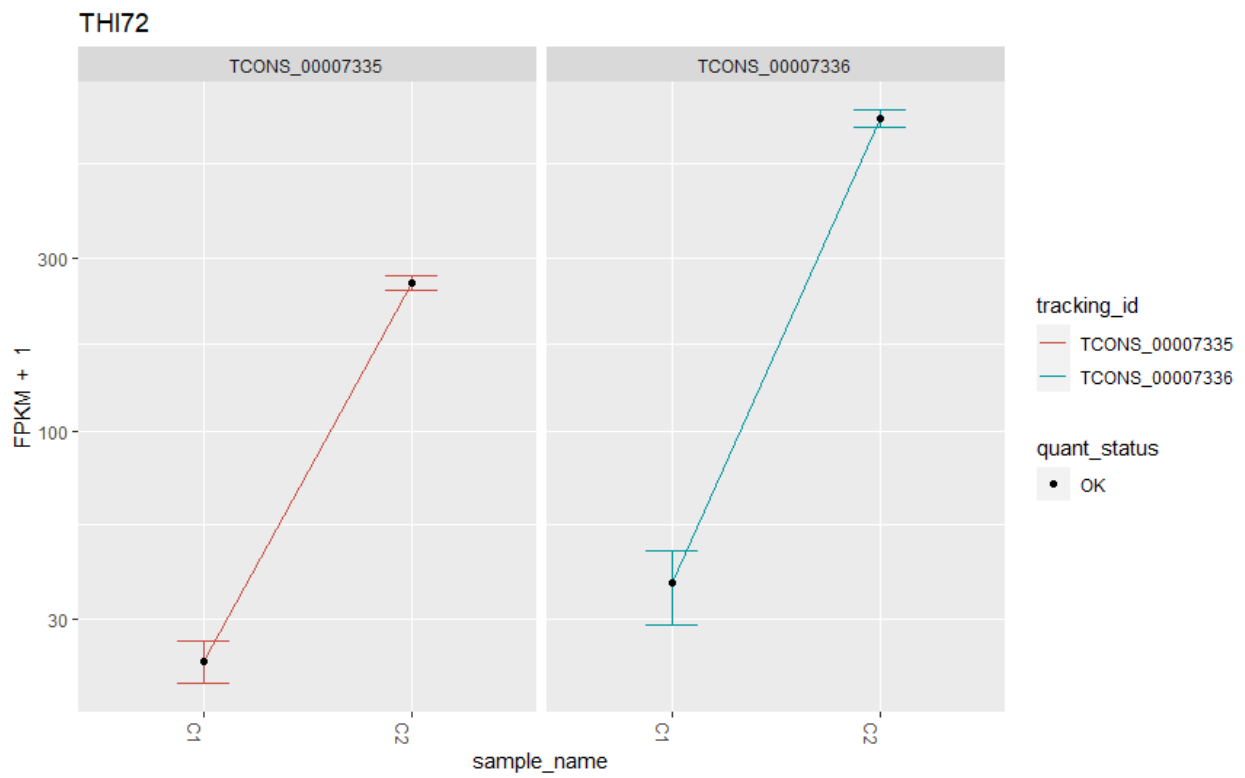


Fig.11 | Figure shows linear plot of expression of THI72 isoforms. The gene is up-regulated.

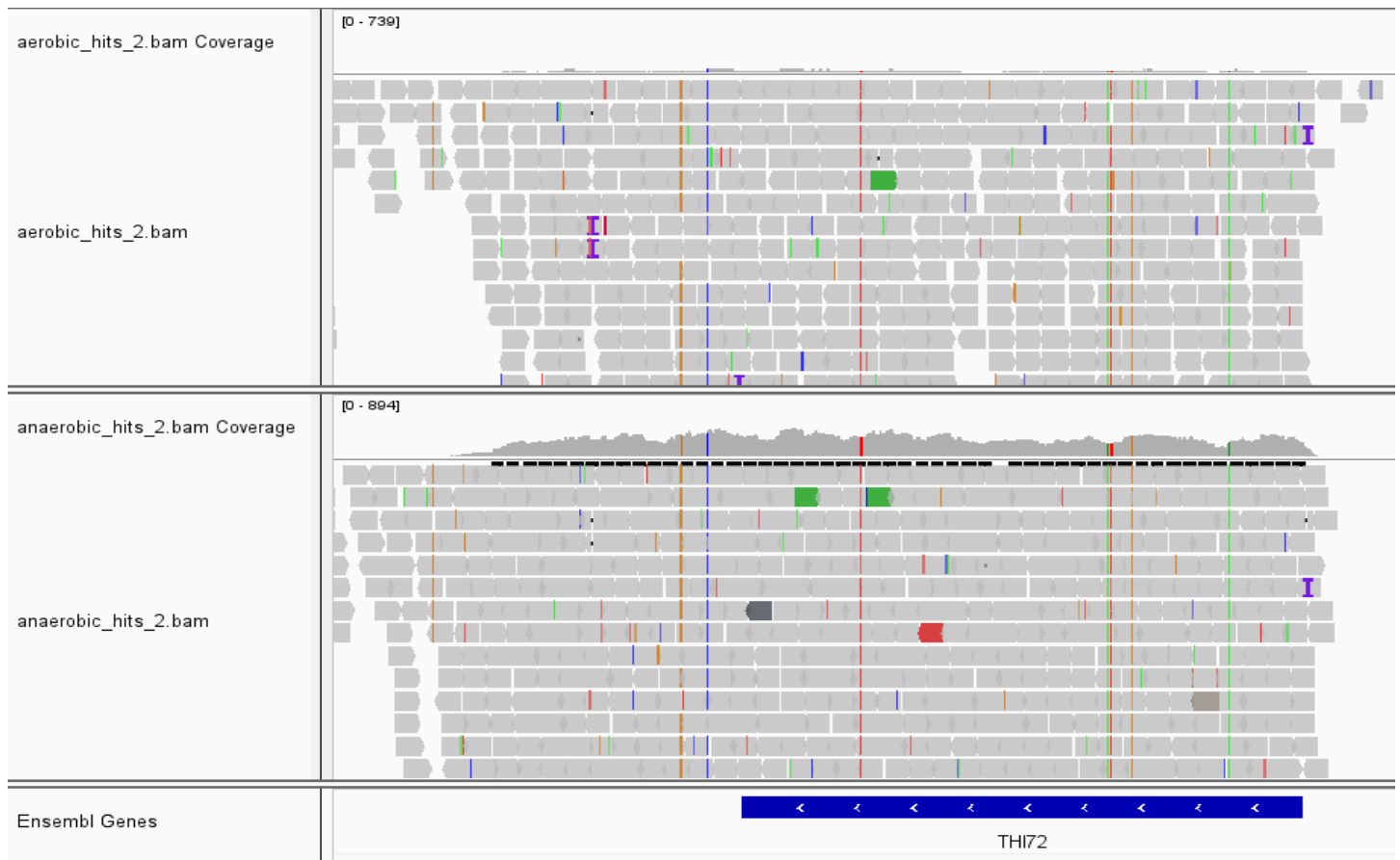


Fig.12 |The read coverage, viewed through the genome browsing application IGV42, shows an increase in sequencing reads originating from the gene in condition C2.

(2) YPS1

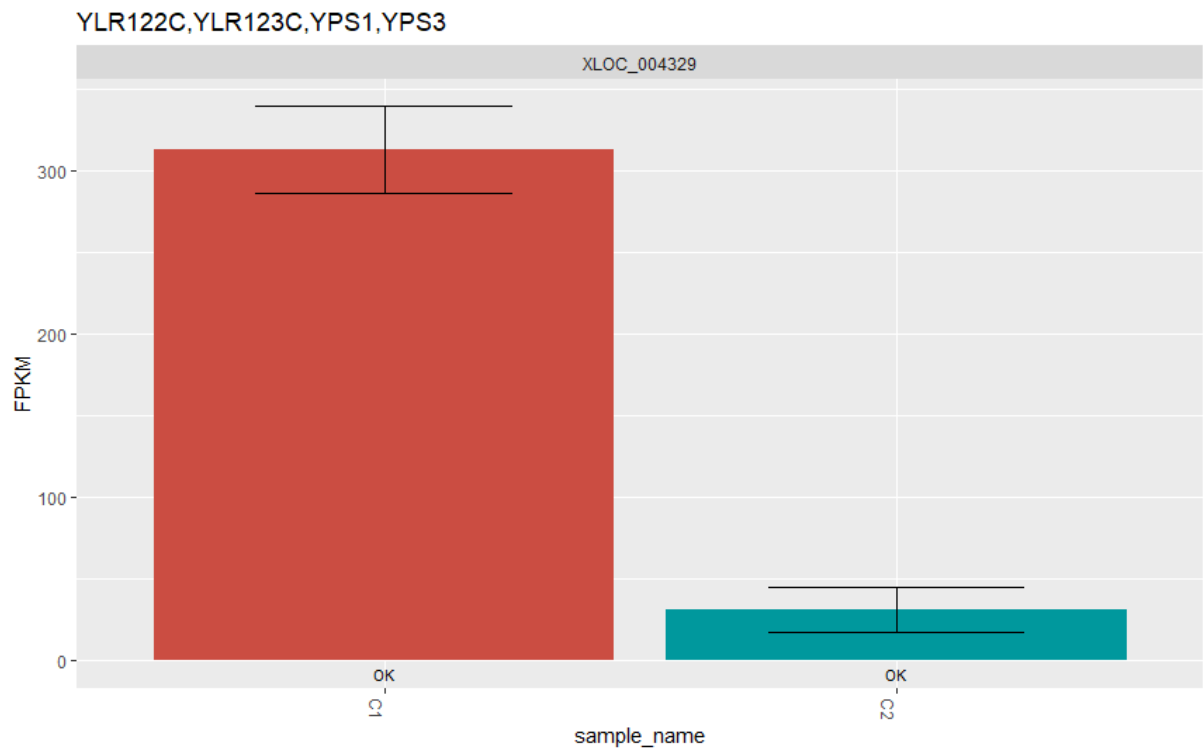


Fig.13 | Differential analysis results. Expression plot shows clear differences in the expression of YPS1 across conditions C1 (aerobic) and C2 (anaerobic), measured in FPKM

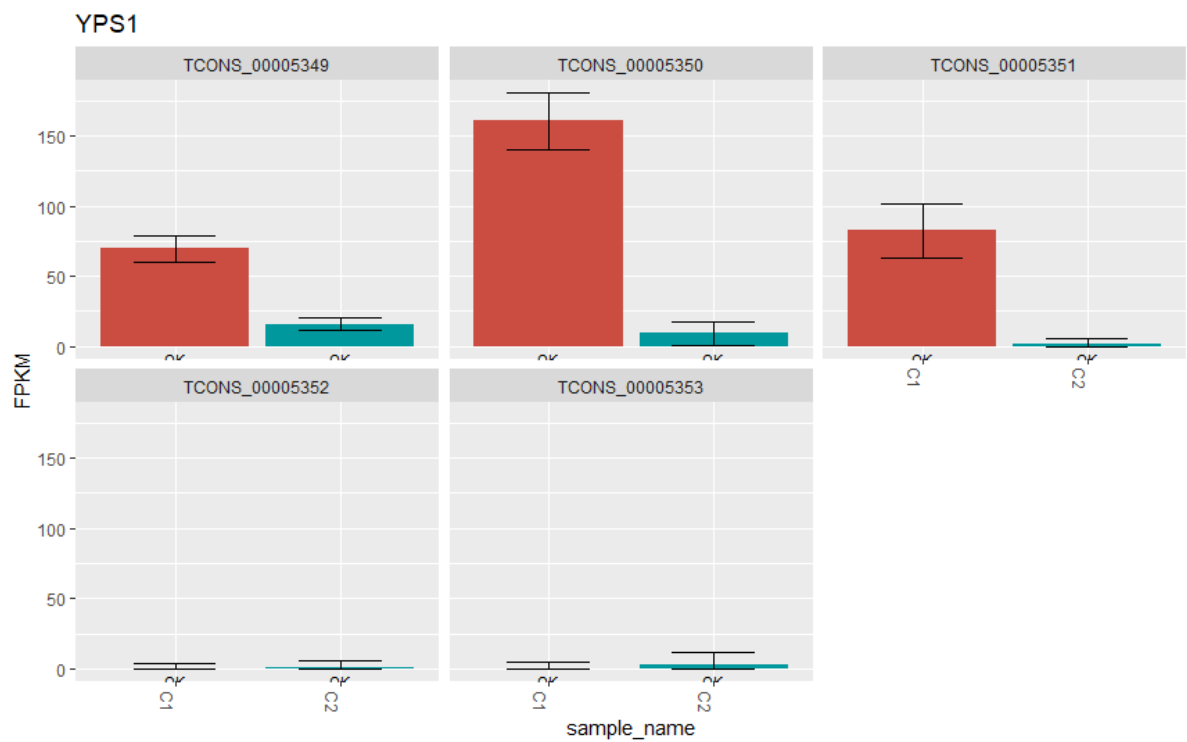


Fig.14 | Changes in YPS1 expression are attributable to a large increase in the expression of one of five alternative isoforms, of which two are insignificant.

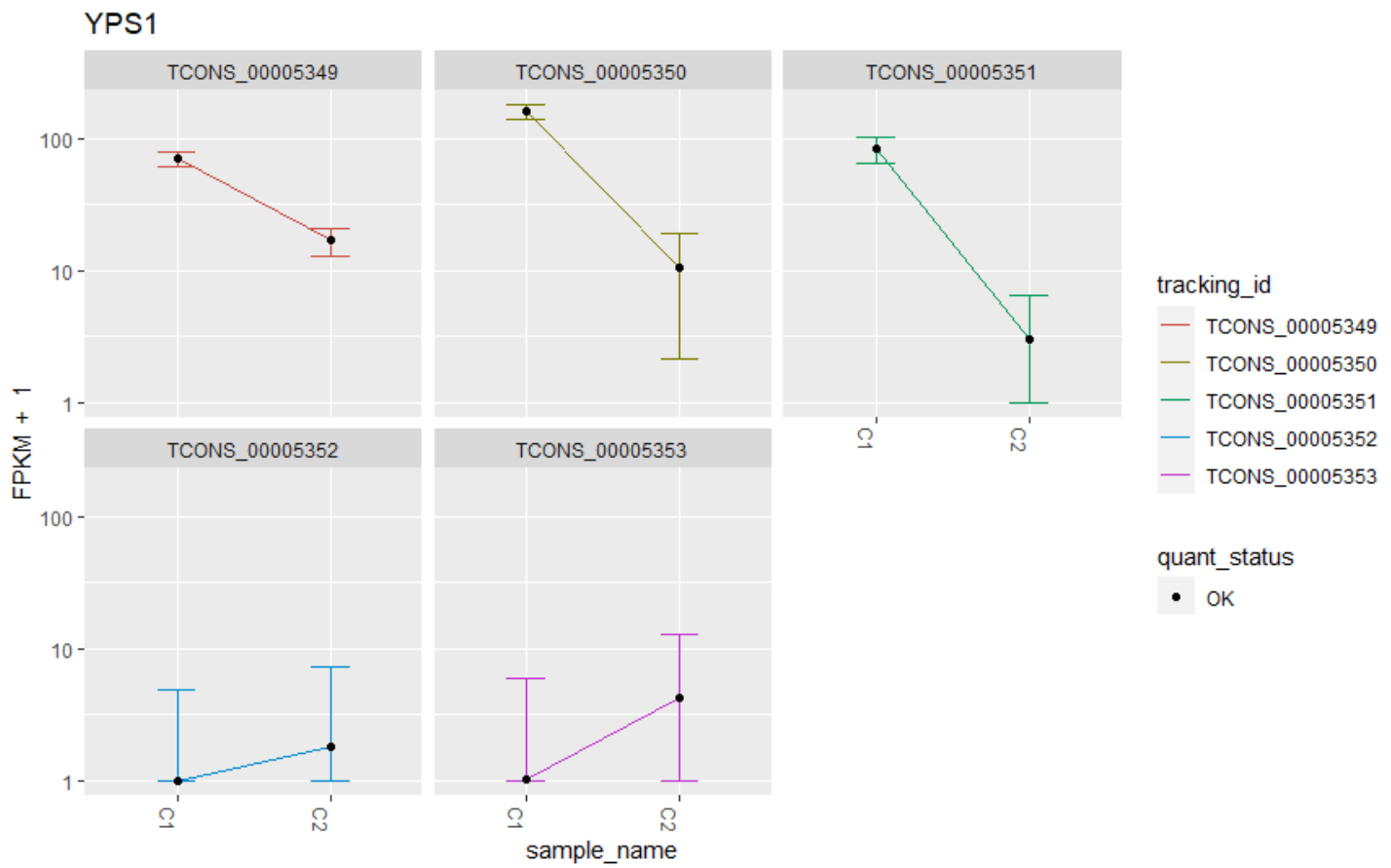


Fig.15 | Figure shows linear plot of expression of YPS1 isoform. The gene is down-regulated. Bottom two plots can be ignored as they are insignificant.

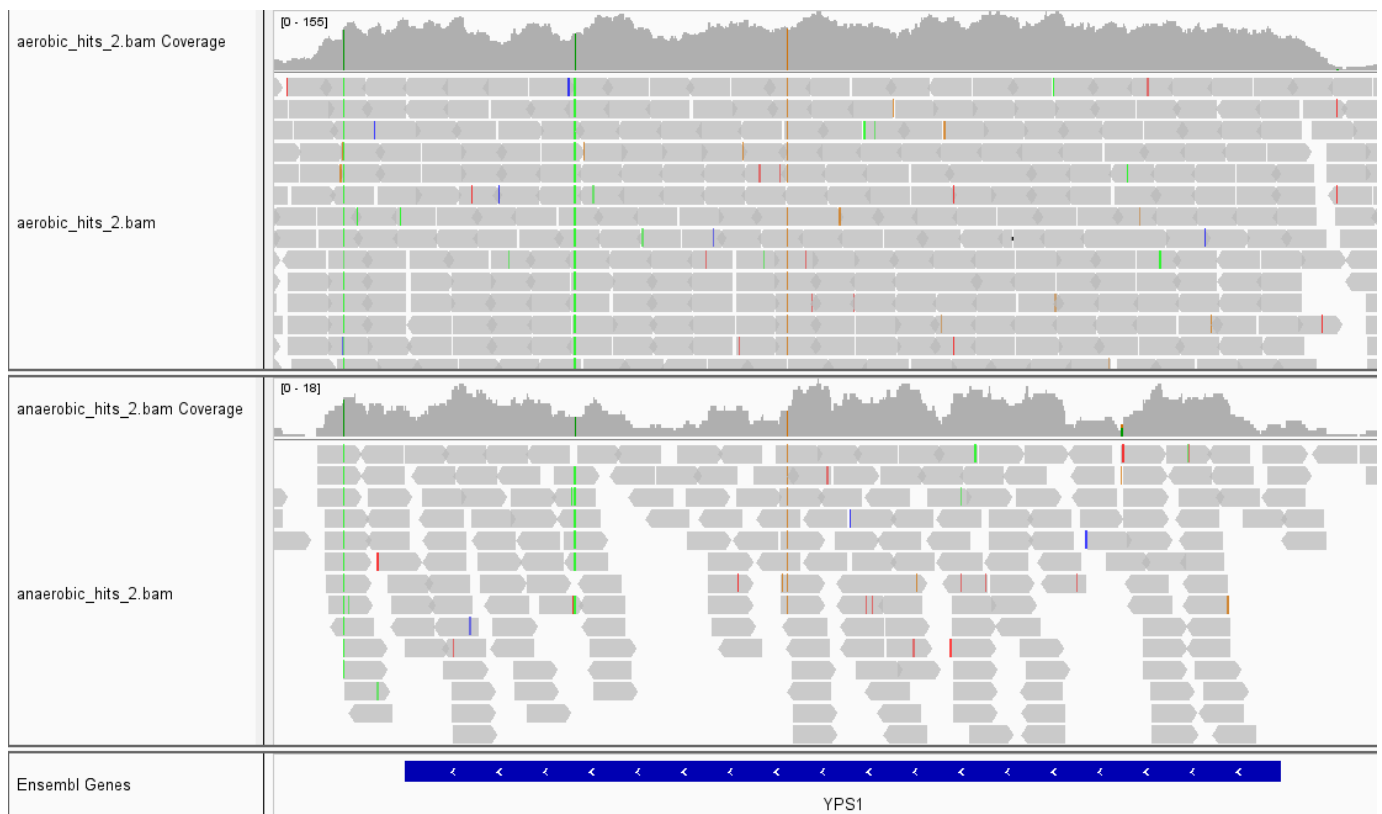
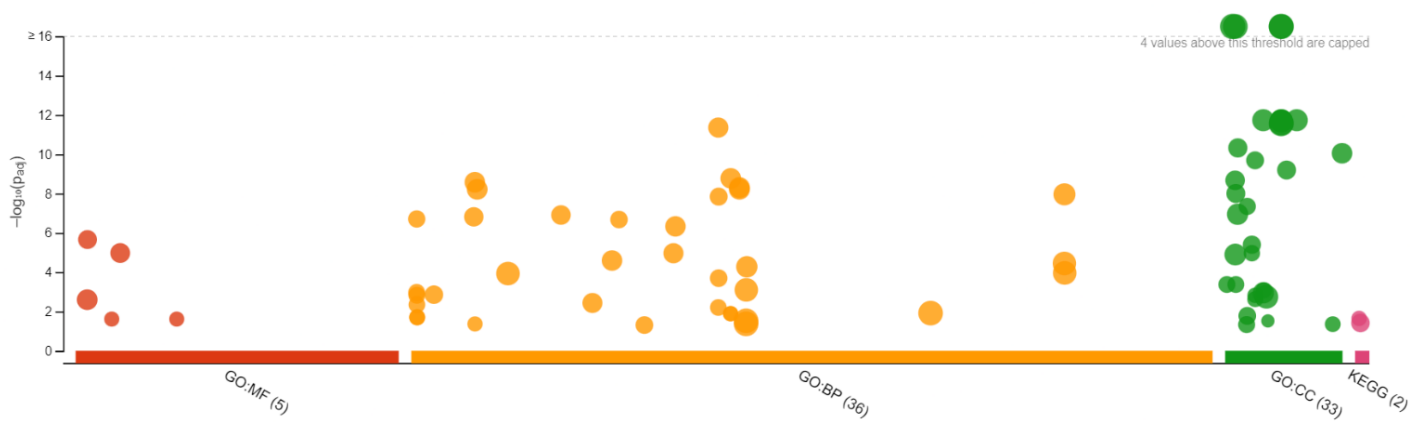


Fig.16 | The read coverage, viewed through the genome browsing application IGV42, shows a decrease in sequencing reads originating from the gene in condition C2.

Finally, a list was created of only the names of the significantly expressed genes and uploaded to **g:Profiler** to check the gene ontology expression (large scale effects on cellular pathways and functions). GO profiler performs functional enrichment analysis, also known as over-representation analysis (ORA) or gene set enrichment analysis, on the input gene list. It maps genes to known functional information sources and detects statistically significantly enriched terms.

We obtained data on:

- (i) GO molecular function
- (ii) GO cellular component
- (iii) GO biological process
- (iv) KEGG pathway



version	e104_eg51_p15_3922dba	g:Profiler
date	10/28/2021, 11:34:30 PM	
organism	scerevisiae	

Fig.17 | Gene ontology profile plot

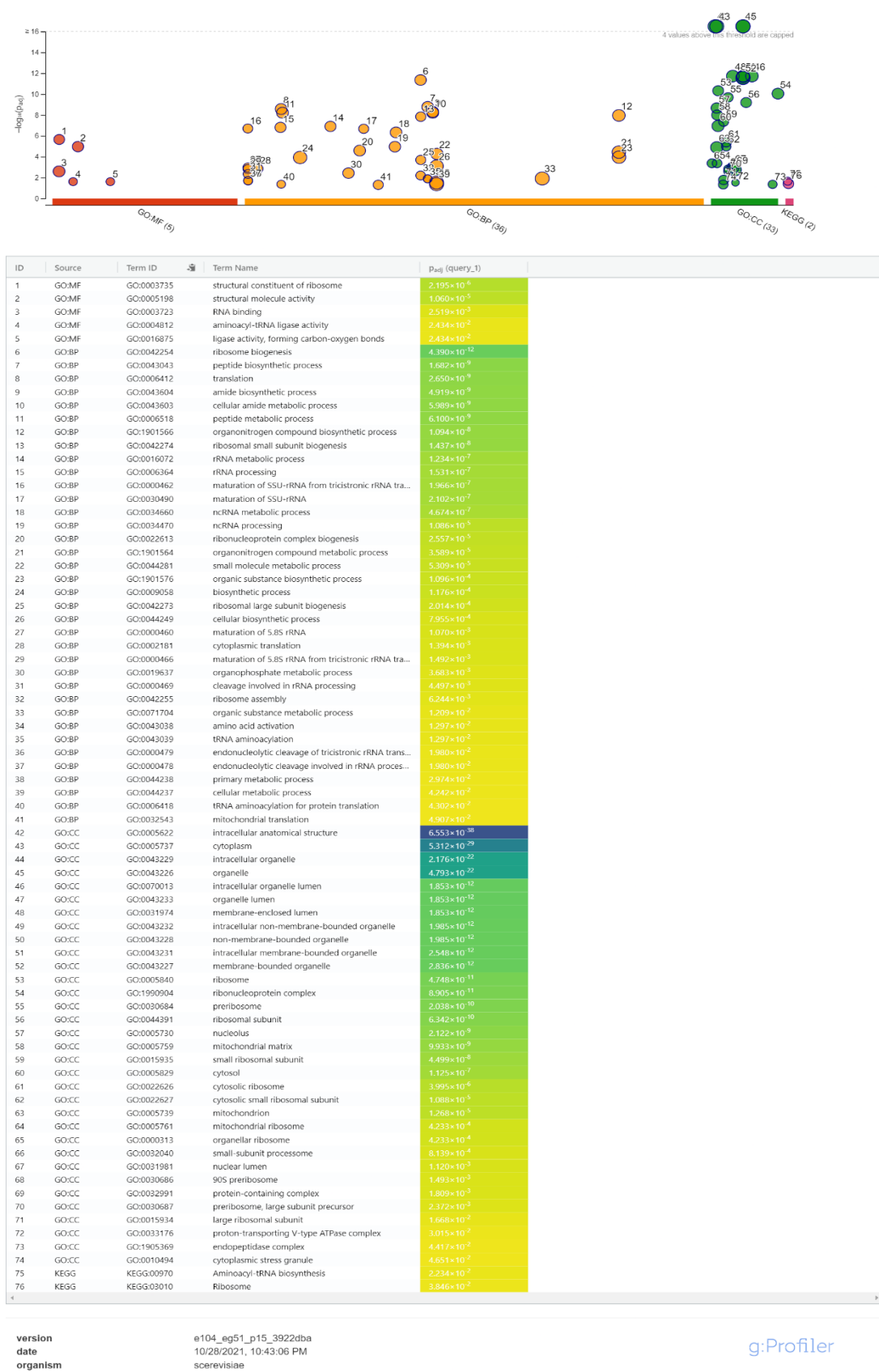


Fig.18 | Differentially expressed genes involved in enriched processes in *S. cerevisiae* in response to anaerobic conditions compared to aerobic conditions

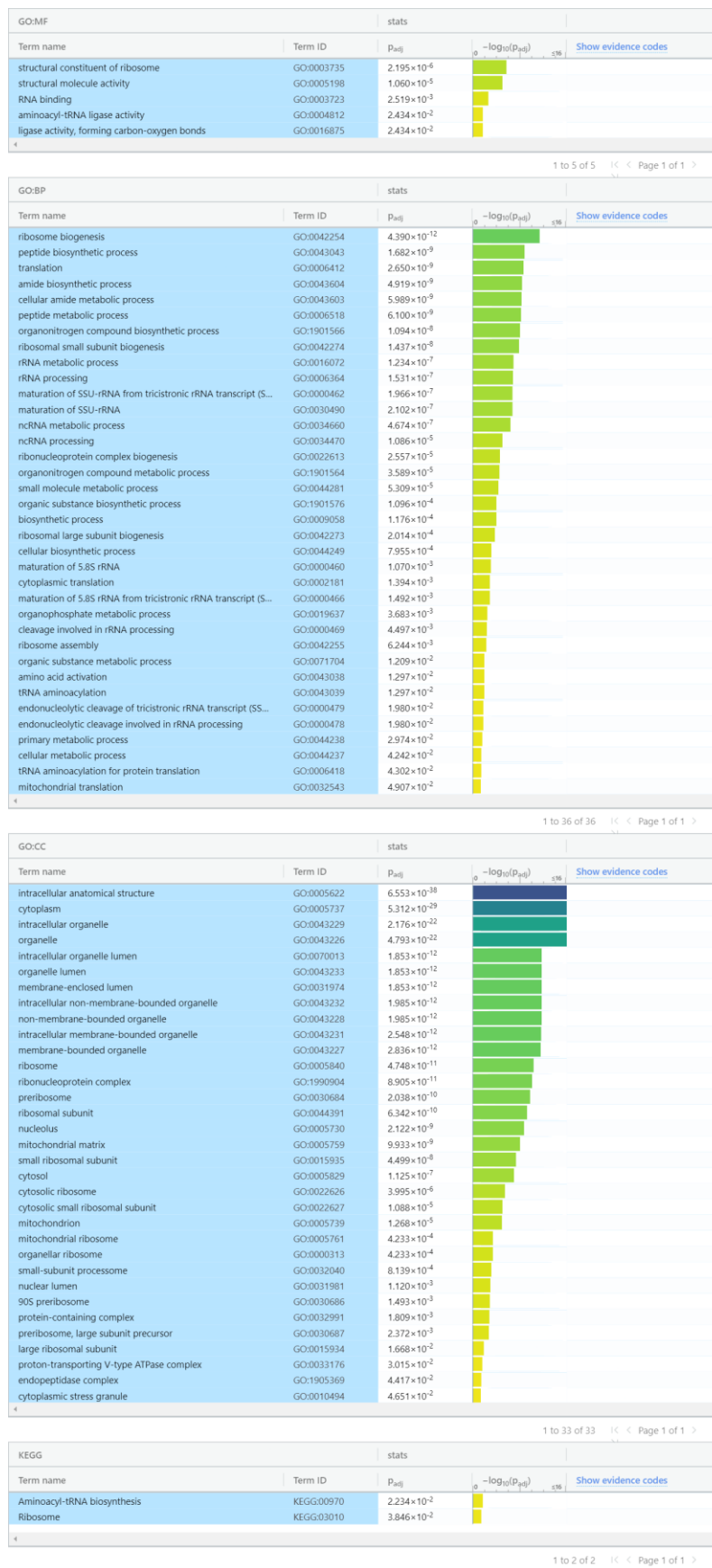


Fig.19 | Enriched GO terms (Biological Processes (BP), Molecular Function (MF), and Cellular component (CC)) in *Saccharomyces cerevisiae* under both the conditions.

IV. DISCUSSION

Low-affinity thiamine transporter, **THI72**, is responsible for the intake of thiamine. Its primary function is the uptake of closely related compounds, and that thiamine transport is a secondary activity of these proteins. The analysis reported a 17 fold ($\log_2_fold_change = 4.06446$) increase in the expression of **THI72**.

Aspartic protease, a hyperglycosylated member of the yapsin family of proteases, **YPS1**, stays attached to the plasma membrane via a glycosylphosphatidylinositol (GPI) anchor. It is involved in nutrient limitation-induced cleavage of the extracellular inhibitory domain of signaling mucin Msb2p, resulting in activation of the filamentous growth MAPK pathway. It is also involved with other yapsins in the cell wall integrity response; role in KEX2-independent processing of the alpha-factor precursor. Studies have shown that yeast cells grow equally well in nutrient-rich environment under aerobic and anaerobic conditions, Taherzadeh et al. Filamentous growth restricts cell multiplication of budding yeast cells under nutrient-poor conditions. During fermentation, yeast cells have a surplus of nutrients that can be utilized by fermenting to ethanol by partial oxidation, and therefore cells prefer the budding form of reproduction over filamentous growth. The analysis reported a 10 fold ($\log_2_fold_change = -3.31345$) decrease in the expression of **YPS1**.

Of the two genes selected, **THI72** is up-regulated when conditions change from aerobic to anaerobic, while **YPS1** is down-regulated. From the gene ontology profiling data, we can predict how the processes that use these genes would be affected in the duration of the experiment. All processes utilizing the **THI72** gene would be expressed highly, and the changes in the cell would occur accordingly, while those processes in which **YPS1** is involved would be reduced. The extent to which the process would be affected depends on how significantly the gene is involved in that particular process.

Output from Cuffdiff tells us whether a gene's observed fold change is significant. In the case of **YPS1**, it has significantly differentially expressed isoforms. However, as a gene's expression level is the sum of the expression levels of its isoforms, and some **YPS1** isoforms are increased while others are decreased, the fold change in overall gene expression is modest. While for **THI72**, which has only two isoforms, that to both having increased expression levels, the fold change in overall gene expression is relatively high.

The experiment reveals markedly differentially expressed genes and transcripts between the two conditions. The protocol does not result in more spurious genes and transcripts than expected (the default false discovery rate for Cuffdiff is 5%). However, poorly replicated conditions, inadequate depth or quality of sequencing, and errors in the underlying annotation used to quantify genes and transcripts can all lead to artifacts during differential analysis.

S. cerevisiae is better able to deal with the fermentation environment possibly due to its efficient competitive uptake of sterols, copper and iron, accompanied by cell wall remodeling to accommodate additional mannoproteins and PAU proteins (information inferred from gene ontology expression). These strategies allow the yeast to regulate membrane fluidity and cell wall porosity, and withstand an anaerobic, high ethanol environment.

V. FUTURE DIRECTIONS AND IMPLICATIONS

The protocol followed to obtain differential gene expression for baker's yeast under the given conditions (aerobic and anaerobic) can be extrapolated to another species and any required condition. For example, we can follow this protocol to study human RNA-seq in a diseased and healthy state and compare the expression of different genes. This data would be helpful in developing a cure or modifying pathways to alter the function.

The strong cell wall-related responses in *S. cerevisiae* suggest the importance of this organelle in the cellular response to other species. In particular, the data support that the regulation of adhesion properties may play a central role in modulating the physical and ecological interactions between species

This analysis revealed genes possibly responsible for adapting cells better to anaerobic conditions. Thus, important cellular pathways and key players in fermentation can be studied and industrially utilized for optimal wine production.

REFERENCES

- [1] Trapnell, C., Roberts, A., Goff, L. et al. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc* 7, 562–578 (2012). <https://doi.org/10.1038/nprot.2012.016>
- [2] Shekhawat, K., Patterton, H., Bauer, F.F. et al. RNA-seq based transcriptional analysis of *Saccharomyces cerevisiae* and *Lachancea thermotolerans* in mixed-culture fermentations under anaerobic conditions. *BMC Genomics* 20, 145 (2019). <https://doi.org/10.1186/s12864-019-5511-x>
- [3] Taherzadeh MJ, Karimi K. Pretreatment of lignocellulosic wastes to improve ethanol and biogas production: a review. *Int J Mol Sci.* 2008 Sep;9(9):1621-51. doi: 10.3390/ijms9091621. Epub 2008 Sep 1. PMID: 19325822; PMCID: PMC2635757.