

Project 2: Transcriptional Profile of Mammalian Cardiac Regeneration with mRNA-seq

Divya Venkatraman, Garima Lohani, Marlene Tejeda, Xudong Han
(Programmer) (Analyst) (Data Curator) (Biologist)

Group 1

TA : Kritika Karri

INTRODUCTION:

The adult mammalian heart has a limited capacity for self-renewal following injury [1]. Shortly after birth, mammalian cardiac myocytes enter cellular senescence. Unlike adult mammalian hearts, neonatal mice have the capacity to regenerate their hearts in response to injury; however, this capability is lost a week after birth [2]. In the regenerating neonatal mouse heart, cardiac myocytes demonstrate loss of sarcomere structures, and a significant amount of these cells re-enter the cell cycle. This is indicated by phosphorylated histone 3 (pH3) and up-regulation of aurora B kinase which are key regulators of mitosis [1]. Thus, understanding the mechanisms by which myocytes naturally re-enter the cell cycle during regeneration is fundamental to understanding the molecular processes that prevent regeneration in the adult heart. Here, we analyzed the raw data of short mRNA sequencing from newborn mice (P0), compared the gene expression level with the data from adult (8–10 weeks of age) mice, and we identified significantly differentially expressed genes. We analyzed the functional annotation clustering of these genes through DAVID 6.7 (Database for Annotations, Visualization and Integrated Discovery)[9] and built a heatmap using these genes among all samples in different development stages. We found that our results were very similar compared to what was reported in the paper, which indicates that these genes might be a potential signature for development or even regeneration of mice cardiac myocytes.

DATA:

CD1 neonatal mice (Charles Rivers Laboratories, MA) were sacrificed by decapitation at P0, P4, and P7 and isoflurane overdose at 8-10 weeks of age to obtain cardiac myocytes for analysis [1]. Cardiac myocytes were obtained from the ventricle of the heart and the apex of the heart from both neonatal mice and adult mice. For whole heart ventricle isolation of both neonatal and adult mice, the whole heart was removed washed with cold PBS and stored with liquid nitrogen. Following, heart atria were separated from the heart and ventricles were processed for RNA sequencing. Two ventricles were used for each replicate for RNA sequencing. The apex of the heart was also removed for sequencing. Researchers also utilized sham operated mice as a control. These mice had their chests opened, hearts exposed, but not resected. Apical resection was collected at 24 hours and 7 days post-surgery using the Neonatal Heart Dissociation Kit [1]. Furthermore, the neonatal cardiac myocytes were dissociated from whole mouse hearts (P0 and P4) and then purified using the Neonatal Cardiomyocyte isolation kit. Each biological replicate contained hearts from 5 to 10 pups. Per timepoint and experimental treatment, three biological replicates were generated for RNAseq. RNA-Seq libraries were generated from cardiac myocytes isolated using the Miltenyi purification system and prepared using the Illumina True Seq sample preparation kit. Forty nucleotide-long pair-end sequencing was performed by using Illumina HiSeq 2000.

There were various samples that the researchers analyzed over the course of this project. However, we only focused on P0 and adult samples which amounted to 8 samples. These samples were found in GEO **Series GSE64403**. Data was assessed to be high quality per the base sequence quality and basic statistics in FastQC Report.

METHODS:

Data Acquisition

Module load sratoolkit was used to extract the SRA(Sequence Read Archive) format to FASTQ. Commands such as fastq--dump and split were used [10]. Since the data was paired end there were two FASTQ files produced. To ensure the format was as expected, head command was used to ensure the header of both files were the same.

Aligning the reads against the reference genome and obtaining quality control metrics

We aligned the reads against the mm9 mouse reference genome using TopHat[3]. To run TopHat we created a batch job using a qsub file and submitted it to the cluster. The tophat arguments used were (-r 200 -G /project/bf591/project_2/reference/annot/mm9.gtf --segment-length=20 --segment-mismatches=1 --no-novel-juncs -o P0_1_tophat -p 16). -r gives the expected inner distance between pairs. -G takes in a set of gene model annotations in the format of a gtf file. TopHat first extracts the transcript sequences and uses Bowtie to align reads to this virtual transcriptome first. Only the reads that do not fully map to the transcriptome are then mapped on the genome. -p gives the number of processor threads used to align the genes [3]. The TopHat command produced a BAM file with the aligned reads.

We used three utilities from the RseQC package [5] , geneBody_coverage, inner_distance, bam_stat to obtain some quality control(QC) metrics for the alignment. These utilities input sorted and indexed bam files. We used samtools [8] to index the BAM file.

Quantifying gene expression and finding differentially expressed genes using Cufflinks

After we performed the alignment using TopHat, we used cufflinks [6] to count how reads map to genomic regions defined by an annotation that was given in a gtf file. In order to run cufflinks, we created another batch job using a qsub file. The cufflinks arguments used were (--compatible-hits-norm -G /project/bf528/project_2/reference/annot/mm9.gtf -b /project/bf528/project_2/reference/mm9.fa -u -o P0_1_cufflinks -p 16). Like in TopHat, -G gives the gene model annotation file and -p is for the number of threads used to run cufflinks. -b instructs cufflinks to run its bias detection algorithm which improves accuracy of the transcript abundance estimates [7]. --compatible-hits-norm tells cufflinks to use only those fragments compatible with some reference transcript towards the number of mapped hits used in the Fragments Per Kilobase of transcript per Million mapped reads (FPKM) denominator. The cufflinks command gave us a file called genes.fpkm_tracking which contains the quantified alignments in FPKM for all genes.

We then used cuffdiff [3] to find significant changes in transcript expression. It produced a file called gene_exp.diff which contains gene-level differential expression. This file was used in the further analysis.

Identifying differentially expressed genes associated with myocyte differentiation:

The identification of differentially expressed genes using cuffdiff produced a file as mentioned above. This file contained differentially expressed gene statistics for P0 vs Adult. The file was read into R studio. Then, the data was sorted in ascending order according to q-value. The list of the top ten differentially expressed groups was assessed. The significant genes were subset into different data frames. Next, we plotted histogram for log2_foldchange for all the genes and the significant genes. The differences were observed and analysed.

Further, the subset of significant genes was separated into upregulated and downregulated genes using log2_foldchange column. Then, the upregulated files and downregulated files were saved into text files.

In the following steps, the upregulated genes and downregulated genes were analysed in DAVID (Database for Annotations, Visualization and Integrated Discovery) 6.8 Function Annotation Clustering. This tool helped us obtain enriched gene sets which were organized into clusters. The Mus Musculus was selected as the species. For the Gene Ontology group, only GOTERM_BP_FAT, GOTERM_MF_FAT, and GOTERM_CC_FAT groups were selected.

Comparing our data with those from the reference paper

We selected the highlighted genes in Fig.1D in the reference paper [1] and compared the expression level among these genes between P0 and Adult samples from our data. We took the average value of FPKM between P0_1 and P0_2, Adult_1 and Adult_2 respectively, and built three line charts to show the expression difference among these genes.

We also compared our gene functional annotation clustering data with those from the reference paper. We selected the clusters with enrichment scores greater than $-\log_{10}(0.05)$ and reported the overlapping GO terms which are recorded in supplementary online table I from the reference paper. We made an augmented table based on our DAVID table by adding one more column called Overlap.

Finally, we built a heat map using the P0_1 FPKM data we created, combined with the data from the other 7 samples given by the instructor. Each file may contain duplicated genes in more than one row, so we merged these repeated genes into one row by adding their FPKM value together. We selected the top 200 differentially expressed genes between P0 and Ad according to their P-value. All the code is available on our group github repository or BU SCC.

RESULTS:

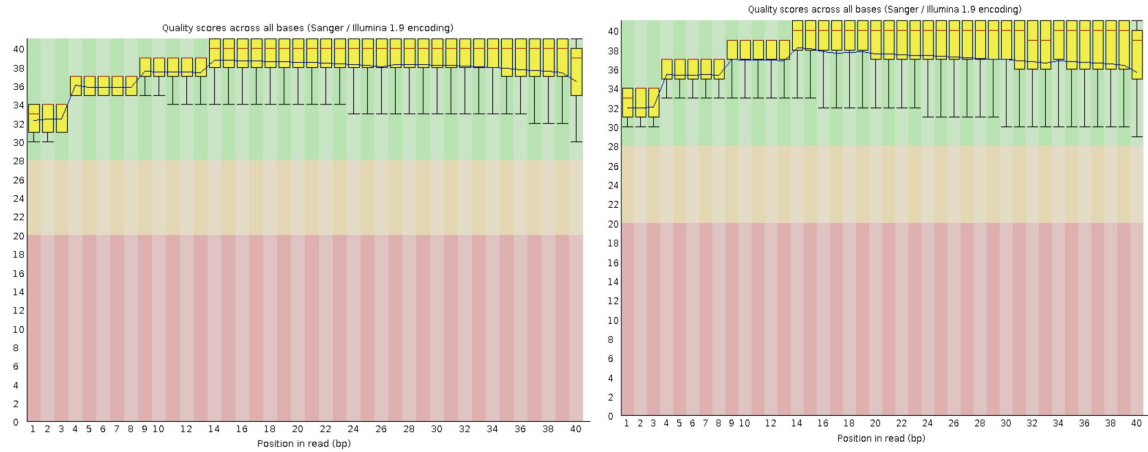
Data Quality

a.

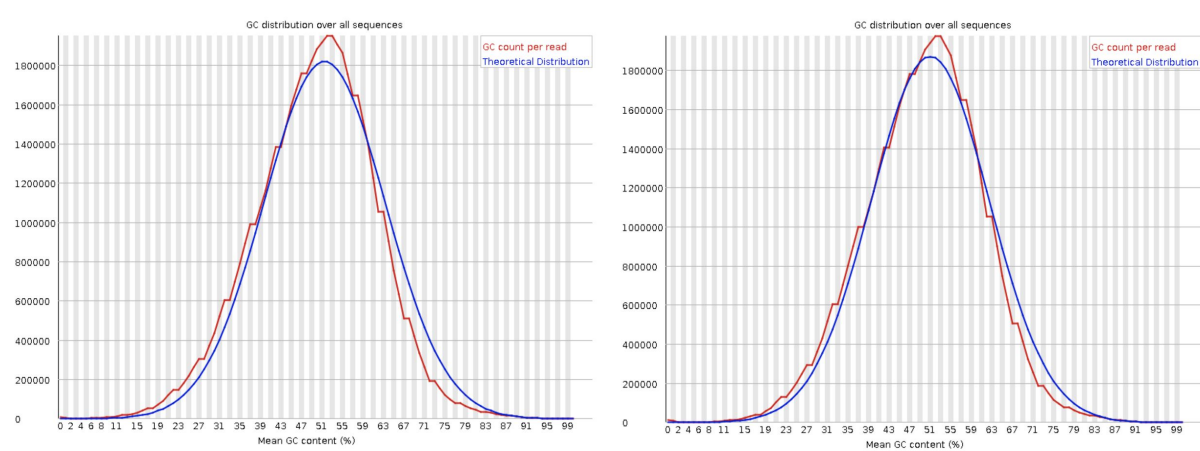
Measure	Value
Filename	PO_1.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	21577562
Sequences flagged as poor quality	0
Sequence length	40
%GC	49

Measure	Value
Filename	PO_2.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	21577562
Sequences flagged as poor quality	0
Sequence length	40
%GC	49

b.



c.



d.

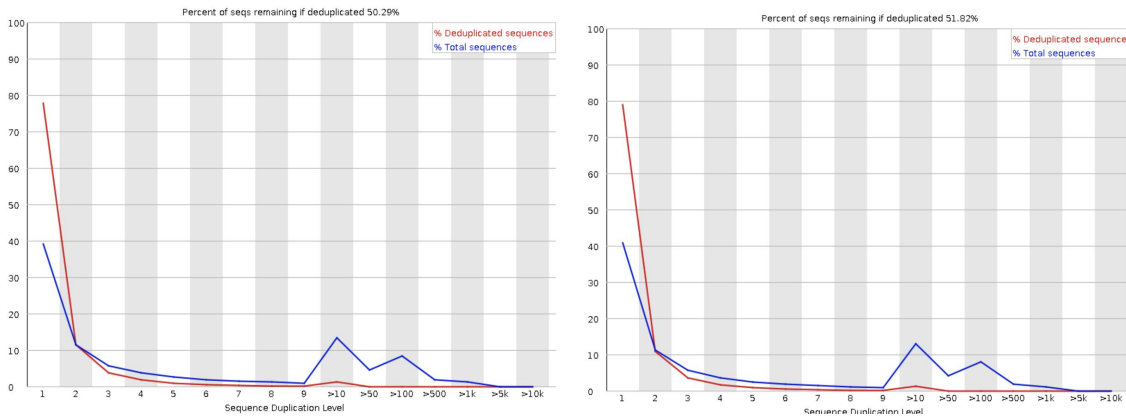


Figure 1. Data Quality Assessment. PO_1 to the left and right PO_2. a. Basic Statistics. b. Quality score per base pair c. GC distribution d. Percent of sequences remaining if duplicated

In figure 1.a there were no sequences flagged as poor quality. The y-axis in figure 1.b shows the quality scores. The higher the quality score then the better the quality of base pairs. The background of the graph divides the y-axis into very good quality scores (green), scores of reasonable quality (orange), and reads of poor quality (red). Quality score remained above 30 per base sequence. Indicating good sequence quality. In figure 1.c it displays the number of reads vs. percentage of bases G and C per read. An unusually-shaped distribution indicates a contaminated library. In the data we can observe there is no shift in the distribution. In figure 1.d we know that normal high-throughput libraries should contain diverse sequences. Finding that a single sequence is very over-represented means that it is highly biologically significant or indicates that the library is contaminated. We can see that there was some duplication sequence. However, there were no overrepresented sequences.

Alignment Summary and QC Metrics

Once reads were aligned using TopHat we used the command samtools flagstat [8] to check the results of the alignment. The following was reported by the command.

```
49706999 + 0 in total (QC-passed reads + QC-failed reads)
8317665 + 0 secondary
0 + 0 supplementary
0 + 0 duplicates
49706999 + 0 mapped (100.00% : N/A)
41389334 + 0 paired in sequencing
20878784 + 0 read1
20510550 + 0 read2
29422646 + 0 properly paired (71.09% : N/A)
39936472 + 0 with itself and mate mapped
```

1452862 + 0 singletons (3.51% : N/A)
1387382 + 0 with mate mapped to a different chr
704916 + 0 with mate mapped to a different chr (mapQ>=5)

These results showed that there were no QC-failed reads and 71.09% of total reads were properly paired.

The RseQC geneBody_coverage.py [5] utility calculated the RNA-seq reads coverage over the gene body and produced a coverage graph. The coverage graph can tell us if there is any 5' or 3' bias. The graph in Figure 2 shows that there is a 3' bias in our data.

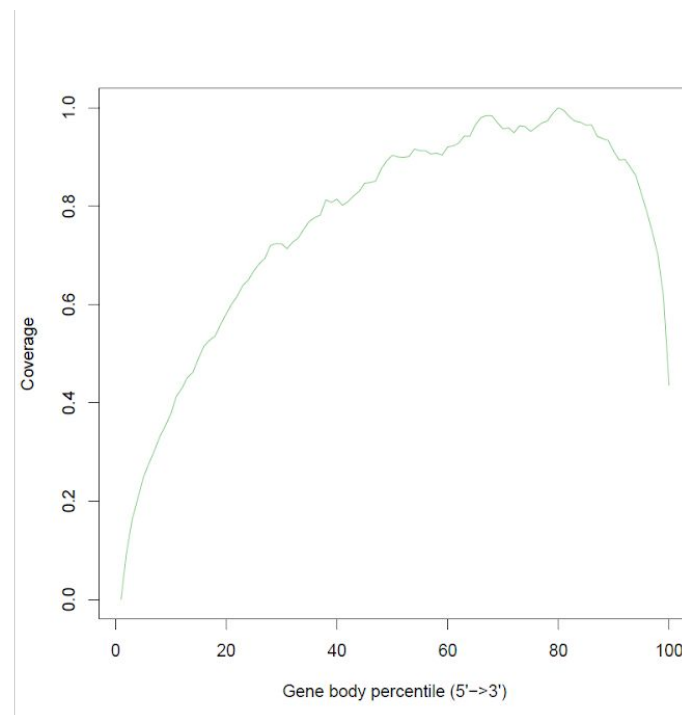


Figure 2 : The graph shows the coverage of reads across the gene body .

The inner_distance.py utility in the RseQC package [5] calculates the insert size which is the mRNA length between two fragments or read pairs. The utility produced a density graph which showed the density of each insert size. Figure 3 shows that insert sizes between 50 – 100 have the highest density.

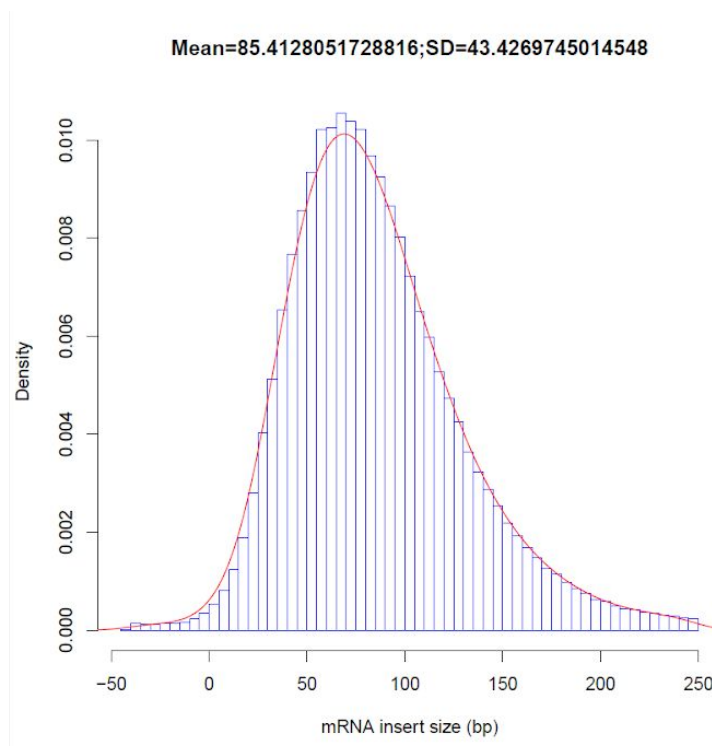


Figure 3 : Histogram showing the density of each insert size

The bam_stat utility summarises the mapping statistics of a BAM file. Table 1 contains the results of bam_stat. These results are similar to the samtools flagstat results. It shows the number of reads that have been mapped and properly paired. It also shows the number of reads that are spliced and non-spliced.

	No. of Reads
Total records:	49706999
QC failed:	0
Optical/PCR duplicate:	0
Non primary hits	8317665
Unmapped reads:	0
mapq < mapq_cut (non-unique):	2899954
mapq >= mapq_cut (unique):	38489380
Read-1:	19409941
Read-2:	19079439
Reads map to '+':	19236824

Reads map to '-':	19252556
Non-splice reads:	33099839
Splice reads:	5389541
Reads mapped in proper pairs:	27972916
Proper-paired reads map to different chrom:	4

Table 1: Results of bam_stat utility showing number of reads for each mapping metric.

Gene Expression Quantification with cufflinks

We used cufflinks to obtain the FPKM values for the genes. The gene Mir5105 has an FPKM of 2303120 which caused the bar plot in Figure 4 to be skewed.

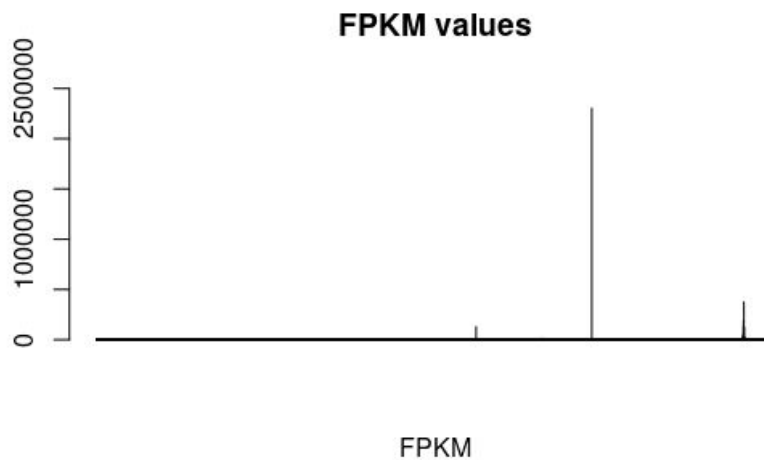


Figure 4: Bar plot of FPKM value of the genes

Differentially Expressed Genes associated with Myocyte Differentiation

The FPKM stands for fragments per kilobase million. These normalized read counts for the sequencing depth, which means sequencing runs with more depth will have more reads mapping to each gene, and length of the gene, which means that longer genes will have more reads mapping to them. In paired end sequencing, both ends can map, giving two reads per fragment or sometimes only one end of the paired end has a quality read and maps. FPKM tracks fragments so that one with two reads is not counted twice. [3]

The $\log_2\text{fold_change}$ is the (base2) log of (FKPM_Adult/FPKM_P0). It is a measurement of changes in gene expression between Ad v P0. Therefore, the $\log_2\text{fold_change} > 0$ imply upregulated genes and $\log_2\text{fold_change} < 0$ imply downregulated genes. And the $\log_2\text{fold_change}$ of zero means no change in

expression. The top ten differentially expressed genes were all upregulated genes and had the same p value and q value (Table 2). Out of 36329 differentially expressed (DE) genes, 5193 were significantly differentially expressed genes. The number of upregulated genes were 2760 whereas there were 2433 downregulated genes (Figure 5A).

Top Ten Differentially Expressed Genes and Statistics						
	gene	FPKM_P0	FPKM Adult	log2fold_change	p-value	q-value
1	Rblcc1	12.19	31.94	1.38	5e-05	0.000318974
2	Pcmd1	13.36	30.17	1.17	5e-05	0.000318974
3	Adhfe1	13.54	27.03	0.99	5e-05	0.000318974
4	Tmem70	36.59	85.04	1.21	5e-05	0.000318974
5	Gsta3	0.41	7.11	4.10	5e-05	0.000318974
6	Lmbrd1	6.70	13.31	0.99	5e-05	0.000318974
7	<u>Dst</u>	18.94	54.22	1.51	5e-05	0.000318974
8	Plekha2	26.63	72.03	1.43	5e-05	0.000318974
9	Mrpl30	55.01	130.53	1.24	5e-05	0.000318974
10	Tmem182	46.02	108.74	1.24	5e-05	0.000318974

Table 2: The top ten differentially expressed gene associated with myocyte differentiation according to lowest q values, where FPKM_P0 is the fragment per kilobase million for Postnatal Day 0, FPKM_Adult is the fragment per kilobase million for Adult sample, and log2fold_change is the $\log_2(\text{FPKM_Adult}/\text{FPKM_P0})$.

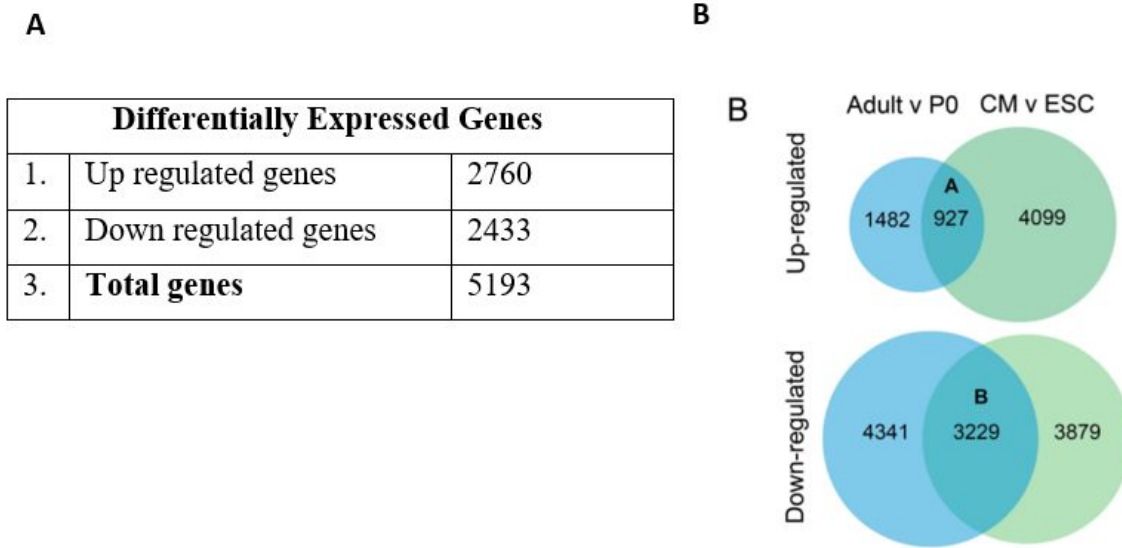


Figure 5: The number of differentially expressed genes for Adult v P0. A) Number of Differentially Expressed Genes at significance $p < 0.01$ in our study, B) Number of genes reported for Adult Vs P0 in Figure 1B from O'meara et al.

A large peak at $\log_2\text{fold_change}$ of zero could be observed in the histogram of $\log_2\text{fold_change}$ for all differentially expressed genes (Figure 6A). However, the histogram of $\log_2\text{fold_change}$ for significantly differentially expressed genes did not have the large peak at $\log_2\text{fold_change}$ of zero (Figure 6B). The large peak at $\log_2\text{fold_change}$ of zero represents the insignificant genes. So, these insignificant differentially expressed genes, which had large frequency, were removed in the second histogram.

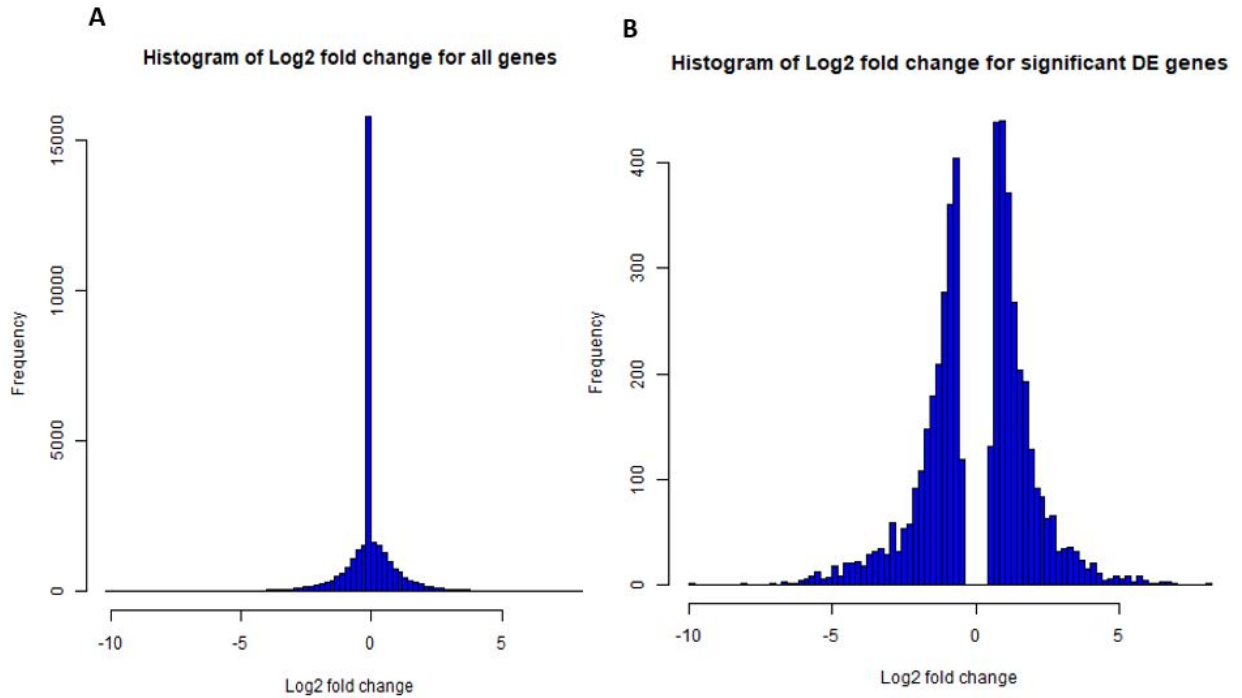


Figure 6: Histogram of log₂_fold change, A) All differentially expressed genes, B) Significant differentially expressed genes

DAVID Results:

The following table summarizes the top cluster from DAVID 6.8 Function Annotation Clustering (Table 3). The top GO terms (gene ontology) for upregulated genes were mitochondrion, respiratory chain, mitochondrial protein complex, organic acid metabolic process, extracellular organelle and sarcomere. And the top clusters for downregulated genes were cell cycle, chromosome, RNA metabolic processes, chromosome organization, single-organism organelle organization, and DNA repair. The highest enrichment score for upregulated genes was 55.01. And the highest enrichment score for downregulated genes was 21.91.

	Upregulated genes		Downregulated genes	
Cluster	GO term	Enrichment scores	GO term	Enrichment scores
1	mitochondrion	55.01	cell cycle	21.91
2	respiratory chains	24.43	chromosome	21.00
3	mitochondrial protein complex	23.38	RNA metabolic process	20.56
4	organic acid metabolic process	23.23	chromosome organization	17.52
5	extracellular organelle	14.78	single-organism organelle organization	16.06
6	sarcomere	11.17	DNA repair	15.69

Table 3: The result of DAVID 6.8 Function Annotation Clustering for significant upregulated genes and downregulated genes, where GO stands for Gene Ontology

C Common Up and Downregulated Gene Enrichment Terms

A Up-regulated		B Down-regulated	
Enrichment term	Score	Enrichment term	Score
Mitochondria	14.35	Non-membrane bound organelle	88.91
Sarcomere	8.50	Nuclear Lumen	88.91
Sarcoplasm	6.03	RNA processing	59.78
Respiration/Metabolism	4.98	Cell Cycle	59.78
Glycolysis	4.39	DNA repair	59.78

Figure 7: Gene enrichment terms reported in Figure 1C from O'meara et al.

The data we obtained partially overlapped with the reference paper

After we get FPKM data for P0_1 and gene expression data between P0 and Adult, we firstly focus on the highlighted genes in Fig.1D from the reference. These genes are divided into three groups (Fig.8), including genes associated with sarcomere, mitochondria and cell cycle GO terms. Nearly all the genes have been detected in our data except for Mpc1 in the mitochondria group and Bora in the cell cycle group. As we can see, all the genes in the sarcomere and mitochondria groups are up-regulated (Fig.8A-B) and all the genes in the cell cycle group are down-regulated (Fig.8C). The trends of changes in these genes perfectly match the trends in the reference paper, which indicates that genes involved in sarcomere and mitochondria are very likely to be up-regulated during the differentiation of rat cardiac myocytes, while the genes involved in the cell cycle are very likely to be down-regulated.

We also compared the top functional annotation clusters (enrichment score larger than $-\log_{10}(0.05)$) for all up-regulated genes and down-regulated genes to those from the reference paper. We generated an augmented table with one more column based on our DAVID table, which records whether the GO term is overlapped. We listed a tiny part of our tables below (Table4), those GO terms are all from the annotation cluster with the highest enrichment score for up-regulated genes between P0 and Ad. The complete tables can be seen in the TableS1 and TableS2 in the SCC. In up-regulated genes, there are 545 overlapped GO terms in the total 1863 GO terms. In down-regulated genes, there are 594 overlapped GO terms in the total 1698 GO terms. Over 70% of GO terms in the top clusters with the highest enrichment score of both files are labeled as overlapped.

Finally, we built a heatmap for the whole samples from P0_1 to Ad_2 according to their expression level based on the top 200 differentially expressed genes between P0 and Adult (Fig 9). As we can see, all the samples in different development stages (columns) are paired together in the heatmap, which is very similar to the one from the reference paper, and the hierarchical clustering of these 8 samples is consistent with the course of neonatal mice development (from P0 to Ad). From this heatmap, most of the genes we selected are persistently up-regulated or down-regulated during the mice development, which indicates that those genes may be the representative signature to reflect the development or regeneration of mice cardiac myocytes.

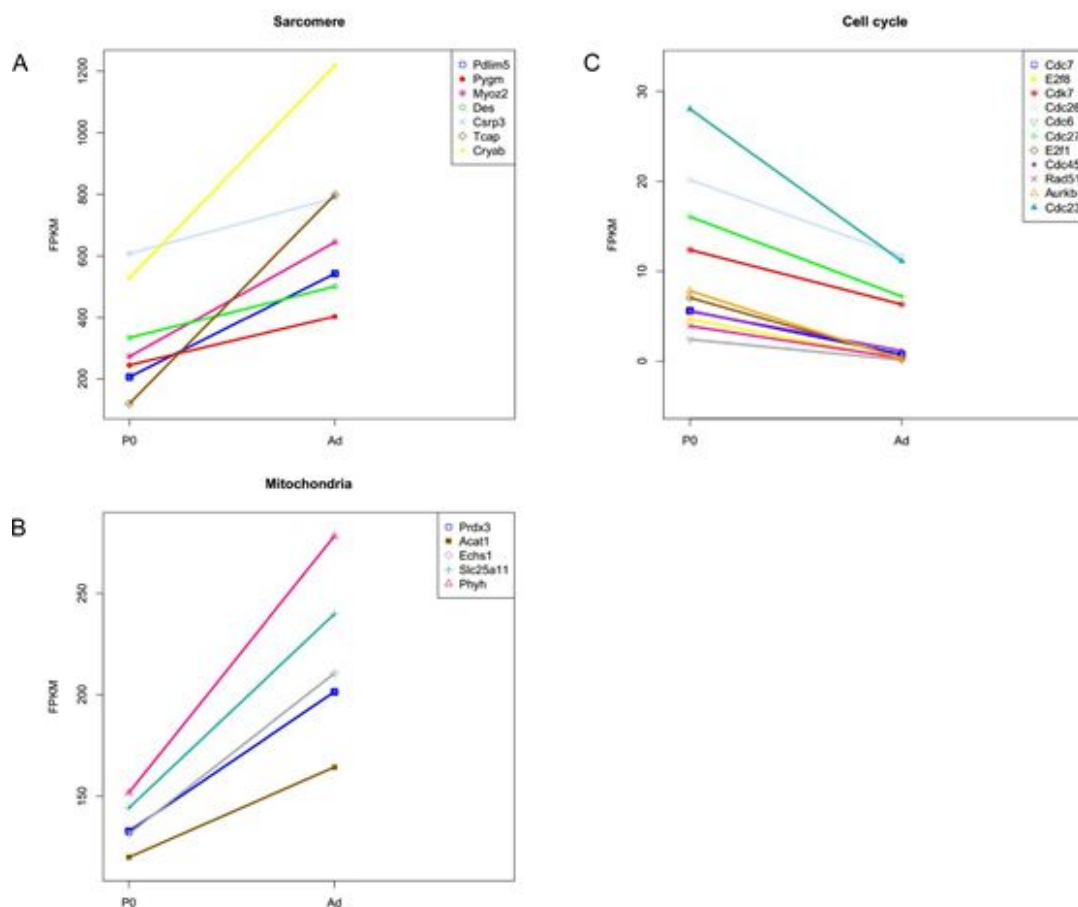


Figure 8. The expression changes among highlighted genes. The vertical axis is the FPKM of these genes. In sarcomere and mitochondria groups (A-B), all the genes are up-regulated, while all the genes in the cell cycle group (C) are down-regulated.

Enrichment Score: 55.01787753192606				
Term	Count	%	PValue	Overlap
GO:0005739~mitochondrion	569	21.74245	6.89E-94	Yes
GO:0044429~mitochondrial part	338	12.91555	2.20E-77	Yes
GO:0005743~mitochondrial inner membrane	193	7.374857	1.33E-57	Yes
GO:0005740~mitochondrial envelope	249	9.514712	1.94E-54	Yes
GO:0031966~mitochondrial membrane	237	9.056171	9.38E-54	Yes
GO:0019866~organelle inner membrane	200	7.642339	1.30E-52	Yes
GO:0031975~envelope	321	12.26595	1.90E-43	Yes
GO:0098798~mitochondrial protein complex	99	3.782958	4.23E-43	No
GO:0031967~organelle envelope	319	12.18953	5.28E-43	Yes
GO:0044455~mitochondrial membrane part	106	4.050439	3.27E-38	No

Table 4: The augmented table for the cluster with the highest enrichment score for the up-regulated genes. The last column reports whether the GO terms are overlapped with the reference paper. The enrichment score and P-value are all calculated via DAVID online tool.

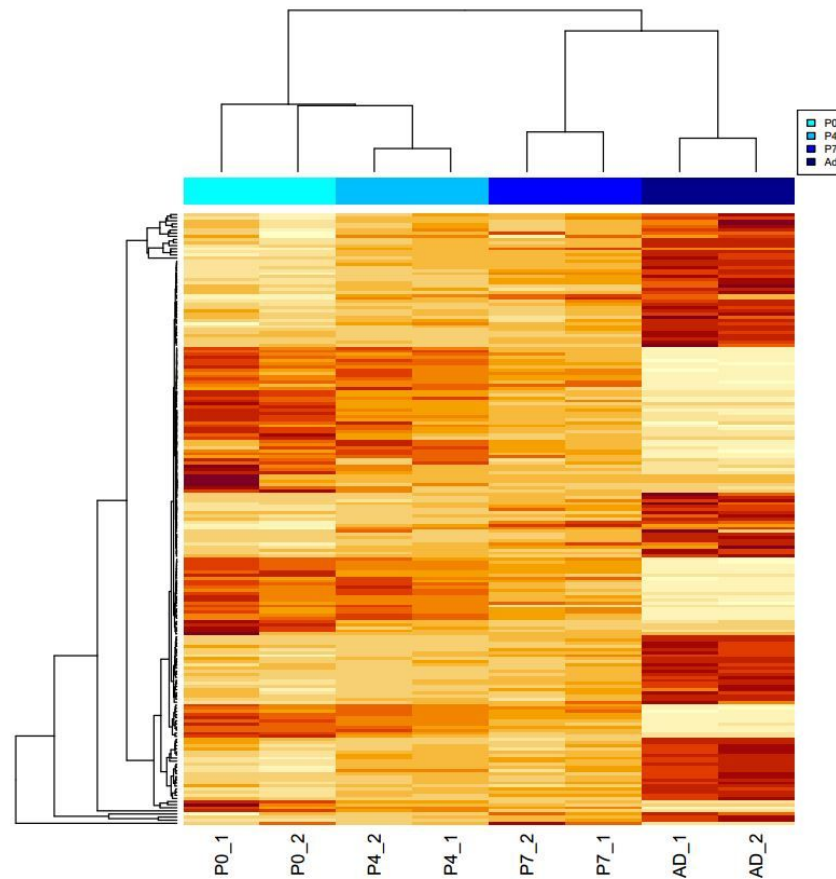


Figure 9. Hierarchical clustering of all expressed genes over different development stages. Each row represents a gene that is significantly differentially expressed between P0 and Adult. Each column represents a sample from different development stages of rat cardiac myocytes. Darker colour means the higher expression level.

DISCUSSIONS:

Samples collected from neonatal mice (P0) and adult mice(8-10 weeks old) to obtain cardiac myocytes for our analysis. Overall, the data quality was good for both P0_1 and PO_2 which is what we would expect. We observe that there are no unmapped reads in our alignment and approximately 71% of total reads were properly paired. The FPKM values observed for the data are very similar to the reference paper with the exception of the Mir5105 gene which shows a very high FPKM value that was not reported by the reference paper. To perform alignment and gene expression quantification, we submitted qsub scripts which ran on the cluster and output was analysed the next day. The TopHat script took time to run, hence was done

The number of differentially expressed genes found in our study do not agree with the reference paper O'meara et al (Figure 5A and 5B). In the paper O'meara et al, there were 2409 upregulated genes in comparison to 2760 upregulated genes found in our study . There is a difference of more than three hundred upregulated genes. Furthermore, in our study we found 2433 downregulated genes (Figure 5A) in

comparison to 7570 downregulated genes reported in the reference paper. So, in the case of downregulated genes, there was a difference of more than five thousand genes. There is a large difference in the number of upregulated genes as well as down regulated genes from reference paper.

The GO terms were found to overlap with the paper O'meara et al (Table 3 and Figure 7). The mitochondria in the reference paper can be related to mitochondrion and mitochondrial protein complexes. Also, it should be noted that respiratory chains would be related to the GO term of Respiration/Metabolism in reference paper. Sacromere was also found as common upregulated genes. The common downregulated genes were cell cycle and DNA repair. Although the terms 'Chromosome' and 'Chromosome organization' are not the same term within the reference paper, these could be related to 'Nuclear Lumen' in O'meara et al..

The Enrichment scores found in our study were not the same as reported in paper O'meara et al (Table 3 and Figure 7). These DAVID annotation and enrichment results are not exactly the same because the DAVID version used in the reference paper is different from the DAVID we use in this study. The DAVID version used in the reference paper is not known. Also the gene set used in the paper O'meara et al is different from our study. Further, the reference paper involved both in vivo and in vitro cardiac myocyte differentiation whereas our study was limited to in vivo maturation condition.

In our comparison analyses, we found our data can almost duplicate the results in the reference paper. For the highlighted genes in the reference's Fig.1D, the regulation for the expression of these genes between P0 and Ad are all in the same direction with the reference, except for two genes with 0 FPKM values. For the functional annotation part, we found GO terms in higher enrichment score are more likely to overlap with the online table i of the reference. And in our heatmap, we selected top 200 differentially expressed genes and we identified these genes can reflect the entire course of development for mice cardiac myocytes, according to the persistent changes of the genes' expression and the hierarchical clustering order of all 8 samples, which indicates that these genes may be the signature for the development or even regeneration of mice cardiac myocytes.

CONCLUSION:

We aimed to use the data from short mRNA sequencing from newborn mice, compared the gene expression level with the data from adult mice, and identified significantly differentially expressed genes. In analyzing a gene subset, the total number of up and down regulated differentially expressed genes did not identically match the original O'meara et al. The overall observation was that the GO Enrichment Terms were related to each other as well as the paper. In conclusion, although the exact gene counts and top genes were not exactly the same, the overall biological functions were able to be observed between this study's groups and the original O'meara et al. paper.

REFERENCES:

- [1] O'meara, C. C., Wamstad, J. A., Gladstone, R. A., Fomovsky, G. M., Butty, V. L., Shrikumar, A., ... & Lee, R. T. (2015). Transcriptional reversion of cardiac myocyte fate during mammalian cardiac regeneration. *Circulation research*, 116(5), 804-815.
- [3] Trapnell, C., Williams, B., Pertea, G. et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* 28, 511–515 (2010).
- [4] Cole Trapnell, Lior Pachter, Steven L. Salzberg, TopHat: discovering splice junctions with RNA-Seq, *Bioinformatics*, Volume 25, Issue 9, 1 May 2009, Pages 1105–1111
- [5] Wang, L., Wang, S., & Li, W. (2012). RSeQC: quality control of RNA-seq experiments. *Bioinformatics* (Oxford, England), 28(16), 2184–2185.
- [6] Trapnell, C., Hendrickson, D., Sauvageau, M. et al. Differential analysis of gene regulation at transcript resolution with RNA-seq. *Nat Biotechnol* 31, 46–53 (2013).
- [7] Roberts, A., Trapnell, C., Donaghey, J. et al. Improving RNA-Seq expression estimates by correcting for fragment bias. *Genome Biol* 12, R22 (2011).
- [8] Li H.*, Handsaker B.*, Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin R. and 1000 Genome Project Data Processing Subgroup (2009) The Sequence alignment/map (SAM) format and SAMtools. *Bioinformatics*, 25, 2078-9
- [9] Huang DW, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene lists using DAVID Bioinformatics Resources. *Nature Protoc.* 2009;4(1):44-57. [[PubMed](#)]
- [10] “Toolkit Documentation : Software : Sequence Read Archive : NCBI/NLM/NIH.” n.d. Accessed March 18, 2020. https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=toolkit_doc.