# Introduction

The RNASeq have been applied to study human cancer, giving the branch for sub typing the cancer subtypes, through the molecular subtyping tree from the analysis, we can gain the information for the cancer progression and diagnosis, in advance, aim to the subtype of the gene expression clustering, we can design the optimized cancer treatment for specific subtype, and lead to more successful way to treat. it also begun to discover the molecular heterogeneity of cancer, and to be more integrative to analyze cancer transcriptome data combine with other data sources will give more information for digging into the possible way to dissect the epigenetic event for design the treatment. In this project, it will discuss in analytical approaches, including compare the subtype analyze from the RNASeq data compare with microarray subtypes, transcriptional network analysis.

# Materials and Methods

Obtained the raw RNASeq data from 1208 breast cancer specimens, then convert the bam file into fasta file for analysis. Splice junction detection and alignment by STAR2, then perform post-alignment QC by RNA-SeQC. The result is the aligned bam file, it need to perform quality filtering by Samtools, and count the fragment by HT-Seq after aligned to the GRCH38 ("human Genome Reference Consortium version 38) with STAR. Then we will generate the "Manifest" file for input to GDC Data transfer tool, download Manifest(cut –f1 gdc_manifest_20161118_234400_htseq_only.txt > uuids_htseq_only.txt), then use grep to capture htseq row include header. Execute GDC Data Transfer Tool command with Manifest file as argument chmod +x <executable>gdc-client download -m gdc_manifest_20161118_234400_htseq_only.txt and download sample meta-data by cut –f1 gdc_manifest_20161118_234400_htseq_only.txt > uuids_htseq_only.txt # and remove header, then perl reformat_ids.pl uuids_htseq_only.txt > Payload.txt. And we need to reformat data to a JSON query request by command curl --request POST --header "Content-Type: application/json" --data @Payload.txt 'https://gdc-api.nci.nih.gov/files' > File_metadata.11.23.16.txt. For convert the gene identifier from ensembl_gene_id to hgnc_symbol and subtype the pam50 gene set, we need to call the biomart package. The command is execute as follows:

library(biomaRt)

ensembl <- useEnsembl(biomart="ensembl", dataset="hsapiens_gene_ensembl")
head(listAttributes(ensembl))

crossref.df <- getBM(attributes=c('ensembl_gene_id','hgnc_symbol'),mart=ensembl)
write.table(crossref.pam50,file='ensembl_hgnc.pam50.txt',row.names=F,sep="\t",quote=F)

Then we will link those pam50 only ensembl_gene_id to our data set. The subtype define by Microarray is available for 51 of the samples in the set of 119. We will compare the subtype of microarray with the subtype tree by our RNASeq data. 11 normal part solid tissue from TCGA database is downloaded for relating the clusters in our analysis to the status of the samples as normal and their predicted subtypes. After process the raw data, we need to identified the previous reported cancer subtype: normal, luminal A, luminal B, basal or HER2+ by the normalized heat map drawing. Table the subtype to visualize the correlation of the subtypes and the clustering. To process the htseq files into DESeq dataset, use function DESeqDataSetFromHTSeqCount, and conduct log transformation by rlogTransformation for desirable cluster expression data, then use assay() to get the matrix from DESeq dataset. Subset the data for only pam50 dataset by convert row names from Ensemble ids to hence symbols by the commands as follow:

```
library(data.table)

d <- setDT(d,keep.rownames = T)[]

a<-merge(d,pam50,by.x='rn',by.y='ensembl_gene_id')
```

Downloading library gplots for making the heat map by the command as follows:
Create metric x from a for heat map, then make the heat map as follows:

```
library(gplots)

heatmap.2(x, col=redgreen(100),Rowv = TRUE, Colv = T, scale="none",
trace="none",margin=c(10,6))
```
(Fig 2)

```
heatmap.2(x-median(x),col=redgreen(100),Rowv = TRUE, Colv = T, scale="none",
trace="none",margin=c(10,6))
```
(Fig 3)

We also need to make the histogram of  the raw counts per library which is an important guide to check how well are those samples can be expected to separate breast cancer subtypes.  The total number of reads mapped to all genes in the GRCH38 reference are sued and show in the figure 1.

The last step is to download the sigclust2 to see the clustering of those cancer subtypes quality by the p-value shows at the node of each separation. The command is perform as follows:

```
library(sigclust2)

shc.obj <- shc(t(x), metric="euclidean", linkage="ward.D2")
```

plot(shc.obj)

the figure is corrected by FWER, so the p-values showed on the nodes are below 0.05 cutoff(Fig 4).


## Result and Discussion

To observe the cluster from normalized heat map (Fig 3) and compare with the matching microarray subtypes, it shows LumA is most close to Normal solid tumor, basal is most distal from all other types, which also includes HER2 in the same cluster. After LumA, LumB is the most related subtype to LumA which also include one HER2 and some LumA subtypes nested in the main LumB cluster.  This result may lead to potentially different cancer subtype according to the the section of LamA, LamB and HER2 are quite disperse. In the HER2 subtype, this group got highest expression for ERBB2, which are important oncogene in tumorigenesis pathway. For the ER pathway, ERS1 got high expression in most branches, but in HER2 and BASAL subtypes, it expressed neutral. In progesterone receptor gene(PGR), HER2 and basal subtypes got lower expression, while some LumA and LumB subtypes have hight expressions. HER2 subtype p value for branching got 2.43e-36, which is a good quality for sub typing. Also all other subtypes are in a good quality control(p-value) for branching as a subtype.

The result branching to basal and HER2 from the microarray sub typing are not very differnt, the other subtypes could use the most highly expression gene in pam50 gene set to catalog a more  organized cancer subtype for correlated oncogene and coordinated to design the treatment for clinical trend.
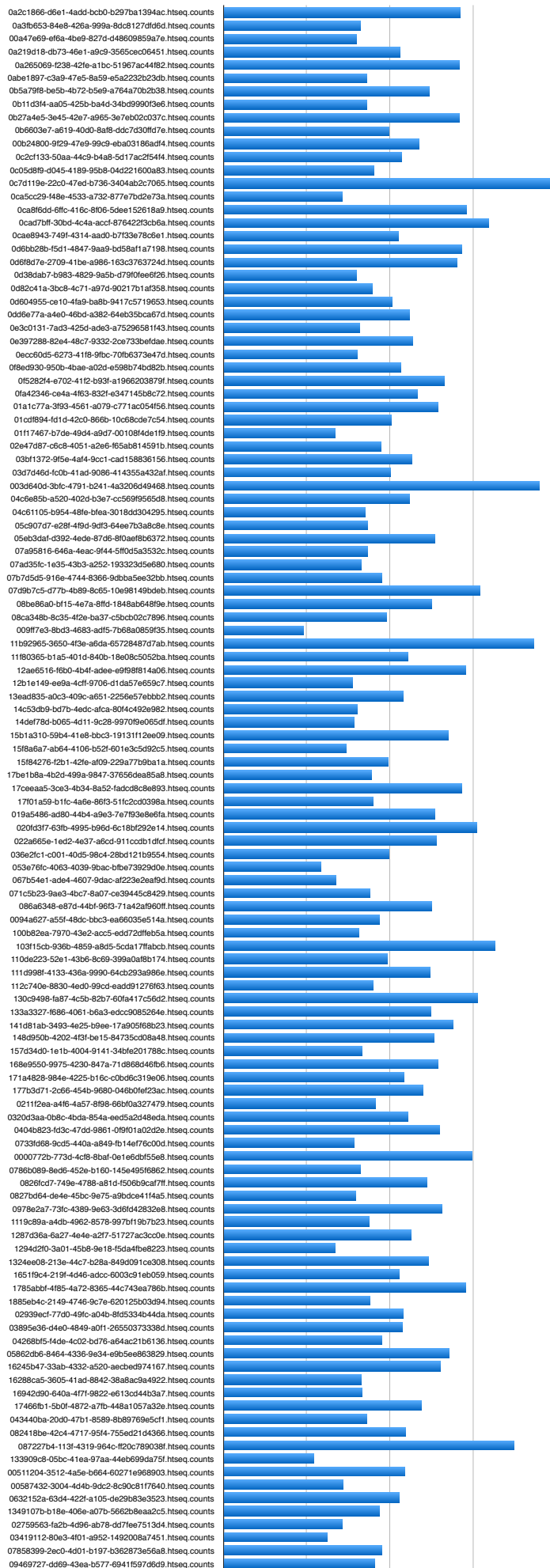
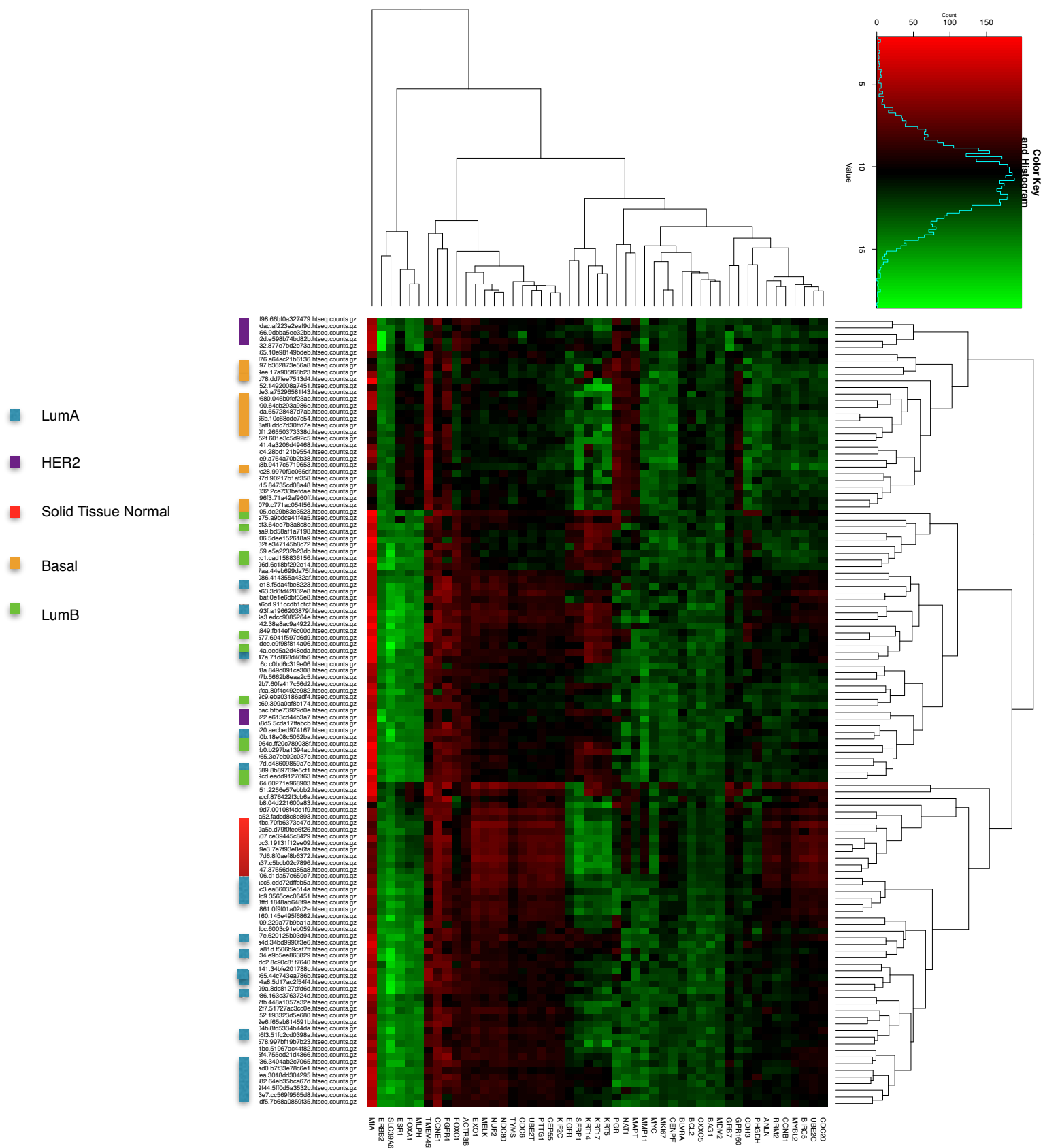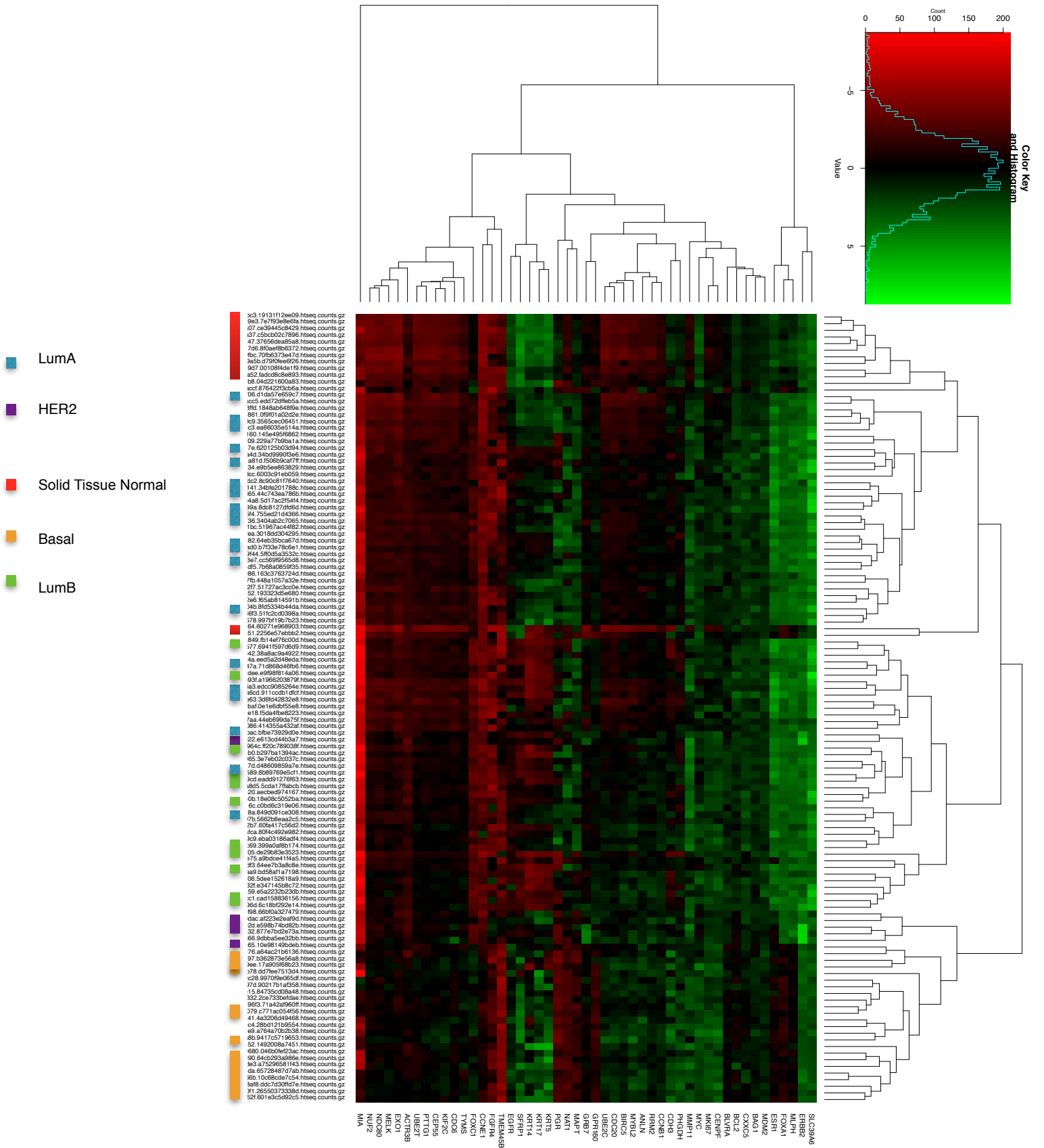Fig1. The histogram shows the sum of reads mapped to all genes in

Fig 2. Heatmap embedded with sample and PAM50 gene set dendrograms on axes using by log transformed counts.

Fig 3. Heatmap in rlog transformed counts after been scaled to the median expression value in the PAM50 gene set labeled with tumor subtype.
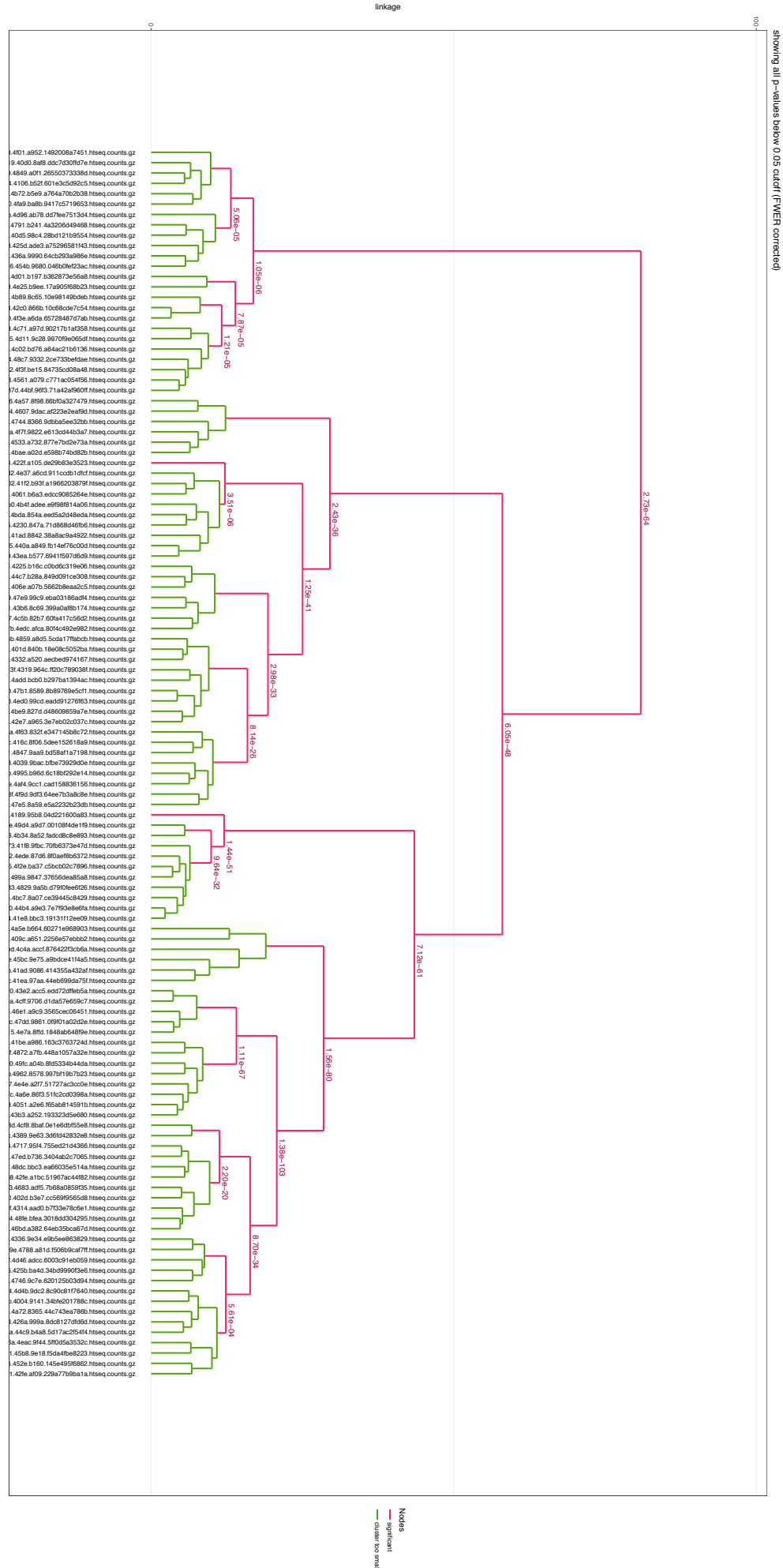
Fig 4. Sample dendrogram showing the significant clusters in analysis by sigclust2 with p-value on each node.