

# **Perform differential expression analysis on fibrotic and non-fibrotic patients under 4 different treatments on HPC**

## **Author**

Sam (Cheng-Hsiang) Lu

Email: Cheng-Hsiang.Lu@cshs.org

## **Mentor**

David Casero

Email: David.Casero@cshs.org

## **Background**

### **Inflammatory bowel diseases(IBD)**

Inflammatory bowel diseases (IBD) are a group of chronic conditions that cause inflammation and damage to the digestive tract. The two main types of IBD are Crohn's Disease and Ulcerative Colitis.

Crohn's disease can affect any part of the digestive tract, like small or large intestine, and can cause symptoms such as abdominal pain, diarrhea, weight loss, and fatigue. It can also cause complications such as fistulas (abnormal connections between different parts of the intestine) and strictures (narrowing of the intestine).

Ulcerative colitis, on the other hand, affects only the colon and rectum and causes symptoms such as bloody diarrhea, abdominal pain, and a frequent need to pass stools. It can also lead to complications such as inflammation of the skin, eyes, and joints.

Both Crohn's disease and ulcerative colitis are chronic conditions, meaning they can last for a lifetime and require ongoing treatment to manage symptoms and prevent complications.

### **Induced Pluripotent Stem Cells(iPSCs)**

Induced pluripotent stem cells (iPSCs) are a type of stem cell that are generated in the laboratory by reprogramming adult cells, such as skin or blood cells, to a pluripotent state. A pluripotent state means that the cells have the potential to develop into any type of cell in the body, just like embryonic stem cells.

iPSCs offer several advantages as they can be generated from the patient's own cells, avoiding issues with immune rejection, ethical concerns and the need for embryos.

iPSCs can be used to study the underlying causes of diseases, test new drugs and therapies, and potentially generate replacement tissues or organs for transplantation.

## Aim

In this project, we will be analyzing RNA-seq data from 19 samples, comprising of 10 samples with fibrotic complications and 9 non-fibrotic samples. Each sample has undergone two runs and 4 different treatments(untreated, TGFb, TNFa, and TGF-b+TNF-a), resulting in a total of 151 samples(1 library failed). We used induced pluripotent stem cells (iPSC) to differentiate into myofibroblasts and stimulated the system with different signals to observe its development. The objective is to investigate the effect of four different treatments: untreated, TGF-B, TNF-A, and TGF-B+TNF-A on the development of the system. In the end, we will perform differential expression analysis to identify the genes that are differentially expressed in fibrotic and non-fibrotic samples under 4 treatments.

## Pipelines

### In HPC

**Convert 151 Fastq to Fasta files** First, put all fastq.gz files into one folder and list all fastq files' name in fastqfiles.txt.

```
ls *q.gz > fastqfiles.txt
```

Cut redundant suffix “\_R1\_trimmed” and list all fastq files' name in libraryname.txt and prefix.txt.

```
ls *q.gz | cut -f 1 -d '.' | sed 's/_R1_trimmed//g' >libraryname.txt
```

```
ls *q.gz | cut -f 1 -d '.' | sed 's/_R1_trimmed//g' > prefix.txt
```

Form a table with 3 columns: fastqfiles.txt libraryname.txt prefix.txt.

```
paste fastqfiles.txt libraryname.txt prefix.txt > tofastatable.txt
```

Create small-sized fasta-formatted files. To submit this job to the cluster on HPC, you need to read the file, library, and prefix. Once you have done that, run the script “generatefastaFromFastaqz” which combine the script “DCfastaqTo-fastaLibraryId.pl”. This results in small-sized fasta-formatted files contain only one header and one sequence per read. You can find all scripts in the “scripts” folder.

```
cat tofastatable.txt | awk {print}' | while read file library prefix ; do qsub -cwd -o $PWD
done
```

This is what each fasta-formatted file would look like:

```
>417iP34untreatedM_S96_1
TNGGCTGGCTAAAGAAGTGAGTATGACCCAGAGGCCAGAGAGGGCAGGGAGAGAATGCCTGGCCACTT
>417iP34untreatedM_S96_2
CNCTTCTTTCGGCGTAGCTCATCAACCTCATATGGTGTGAGCCTTGCTACTTCCGGATGTTCCACATAA
>417iP34untreatedM_S96_3
CNTCCGCCTTTCTGTTCTGTTTTGCCGTCCCCTCCCGCCTGCTGCTGGCCCTGCCTCGCCCCGGGCGGC
>417iP34untreatedM_S96_4
ANCCATGTGAAAAACAGAACTGTAGGGGTTTTTTGTTGTTTGTGTTTGTGAGATGGGGTCTTGCTCTGT
>417iP34untreatedM_S96_5
ANTCCGATGTTCTGCAATTTCTGTGCCCTAGGTTGAACCTCTTTCAGCATTGCACTAGCATCTTCATC
>417iP34untreatedM_S96_6
CNCGACATTAGAAGGTTTTCTGTGGATGGATCGGGCACCGTCTTCTCATATTCCTTTTGGAGACCC
>417iP34untreatedM_S96_7
CNCCTTGCGCACCTCCACCATGACCAGCCCTTCTTCCACCAAGCCCAGCCCCACATCGCCCTTGAATC
>417iP34untreatedM_S96_8
TNAGGTCTCACTGTCCTGGCCTGACCTTCAGCTCTCCAACACTGGGCCCCGCCGGTCTCCGGGAGCCA
>417iP34untreatedM_S96_9
CTGGGGTTATGAGTTTATAGTTGGGAACCTTCTTACAGAGTTTATCATAGGTAGCTTTGTCAAACAAGA
>417iP34untreatedM_S96_10
CCAGTGAGAAGACAGCTTTGCAGTCACACTGGAGATCAGAGTTCAGGCTGCAGCATGTCACCAACGCC
```

**Generate auxiliary files and directories for each sample** Put a list of names of all fasta files in the directory and save them in a text file named “fastafiles.txt”.

```
ls *fasta.gz > fastafiles.txt
```

Cut the redundant suffix “.fasta.gz” from the names of all fasta files and generate a new list of file names with the suffix removed in a text file named “targetdirectories.GTF.txt”.

```
cat fastafiles.txt |sed 's/\.fasta\.gz//g' > targetdirectories.GTF.txt
```

Create a separate directory for each sample listed in “fastafiles.txt”.

```
cat fastafiles.txt |while read line ; do mkdir ${line/.fasta/.gz/GTFpass1/} ; done
```

**Form the submission script called “sendmyof”** Add a shebang line at the beginning of your script file named “sendmyof” to indicate the interpreter that should be used to execute the script.

```
echo '#!/bin/bash/' > sendmyof
```

The command below runs the “generatesendscriptSingleGTFParam” script with several input parameters to map the RNA-seq data with STAR. The input parameters include the list of target directories containing the input data (“targetdirectories.GTF.txt”), the subdirectory name (“GTFpass1”), a parameter file containing settings for STAR alignment (“Parameters.txt”), a prefix for output files (“myof”), the path to the STAR index directory (“/home/luc/RNASEQ\_MASTER/Hsapiens/GRC38/INDEXES/GRCh38.primary.33.basicselected.STAR2.7”), the path to the input data directory (“/home/luc/iPSC/MYOFIBROBLAST/”), the amount of free memory to use (“mem\_free=32G”), and the number of threads to use (“8”). In the end, it will generate a “processLaneSingleGTFParam” file and run the STAR package in each sample’s folder.

```
./generatesendscriptSingleGTFParam targetdirectories.GTF.txt GTFpass1 Parameters.txt myof /
```

Change sendmyof into executable mode and run sendmyof.

```
chmod a+x sendmyof
. sendmyof
```

It will take less than one day to run through 151 samples and generate each sample a folder which contain every output from STAR.

**Create a table summarizing the mapping statistics for each sample**

Change directory into one sample file which ends with “GTFpass1”. Extract the first column from the mapping statistics file and store it in “temp2.txt”.

```
grep "|" 008iP22TGFbM_S71GTFpass1/008iP22TGFbM_S71GTFpass1Log.final.out | cut -f 1 -d "|" |
```

The first column from the mapping statistics file.

```

        Started job on
        Started mapping on
        Finished on
Mapping speed, Million of reads per hour

        Number of input reads
        Average input read length
        UNIQUE READS:
        Uniquely mapped reads number
        Uniquely mapped reads %
        Average mapped length
        Number of splices: Total
Number of splices: Annotated (sjdb)
        Number of splices: GT/AG
        Number of splices: GC/AG
        Number of splices: AT/AC
        Number of splices: Non-canonical
        Mismatch rate per base, %
        Deletion rate per base
        Deletion average length
        Insertion rate per base
        Insertion average length
        MULTI-MAPPING READS:
        Number of reads mapped to multiple loci
        % of reads mapped to multiple loci
        Number of reads mapped to too many loci
        % of reads mapped to too many loci
        UNMAPPED READS:
Number of reads unmapped: too many mismatches
        % of reads unmapped: too many mismatches
        Number of reads unmapped: too short
        % of reads unmapped: too short
        Number of reads unmapped: other
        % of reads unmapped: other
        CHIMERIC READS:
        Number of chimeric reads
        % of chimeric reads

```

Create an empty temporary file for storing intermediate results.

```
rm tempprev.txt
touch tempprev.txt
```

Extract the total mapped reads from each subsequent mapping statistics file and combine with previous results.

```
ls *pass1/*final.out | while read line ; do
grep "|" $line | cut -f 2 > temp.txt
paste tempprev.txt temp.txt > tempnew.txt
mv tempnew.txt tempprev.txt
done
```

Remove the first column and write the final results to a file called “mappingstats-Firstpass.txt”.

```
cut -f 2- tempprev.txt | awk 'NR>3 {print}' > tempnew.txt
mv tempnew.txt tempprev.txt
paste temp2.txt tempprev.txt > mappingstatsFirstpass.txt
```

The mappingstatsFirstpass.txt would look like this:

Mapping speed, Million of reads per hour	452.92	344.45	460.17	487.96	331.34	441.99	560.78	487.93	532.33	351.69	381.99	447.89	463.27
Number of input reads (4604608)	40855906	43588067	47982561	39392754	33886137	72122988	53266037						
Average input read length	75	75	75	75	75	75	75	75	75	75	75	75	75
Uniquely mapped reads number	42647348	36668193	39982295	43570276	36214287	31289999	66061689	48198947					
Uniquely mapped reads %	92.62%	89.75%	91.73%	90.80%	91.93%	92.34%	91.60%	90.47%	93.77%	90.29%	90.29%	92.04%	91.95%
Average mapped length	75.21	75.21	75.20	75.20	75.19	75.21	75.20	75.20	75.20	75.24	75.24	75.23	75.22
Number of splices: Total	14107761	11743321	12025357	12627415	11851597	10196588	20128263	13942120					
Number of splices: Annotated (sjdb)	14029840	11673663	11945482	12536747	11787471	10148891	19995363						
Number of splices: GT/AG	14015583	11657575	11928403	12513568	11777228	10127320	19961004	13827884					
Number of splices: GC/AG	78050	74878	84368	87757	62434	59661	145119	96183	129381	88755	60185	78428	94829
Number of splices: AT/AC	7048	6261	6785	8291	6179	5252	12249	9152	15682	8499	6008	7801	9790
Number of splices: Non-canonical	7080	5487	5801	7799	5756	4355	9871	8901	10933	5153	3764	6609	7788
Mismatch rate per base, %	0.27%	0.28%	0.29%	0.28%	0.28%	0.31%	0.28%	0.29%	0.29%	0.28%	0.27%	0.28%	0.30%
Deletion rate per base	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%
Deletion average length	1.47	1.47	1.47	1.48	1.48	1.47	1.46	1.47	1.48	1.50	1.48	1.48	1.47
Insertion rate per base	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%
Insertion average length	1.61	1.52	1.43	1.43	1.63	1.60	1.47	1.43	1.62	1.42	1.38	1.44	1.59
Number of reads mapped to multiple loci	2390878	3287402	2498107	3093230	2029231	1702575	4120922	3568843	3090452	2530456	1853908	3009727	2483263
% of reads mapped to multiple loci	5.21%	8.05%	5.73%	6.45%	5.15%	5.02%	5.71%	6.78%	4.88%	5.64%	6.11%	6.45%	5.28%
Number of reads mapped to too many loci	222693	287795	314715	343251	343822	264218	498200	397171	262921	547542	254116	316223	407743
% of reads mapped to too many loci	0.48%	0.51%	0.72%	0.72%	0.67%	0.78%	0.69%	0.75%	0.32%	1.20%	0.68%	0.79%	0.50%
Number of reads unmapped: too many mismatches	1269	1809	1322	1446	1105	1001	2452	1565	1799	1824	1173	1452	1228
% of reads unmapped: too many mismatches	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
Number of reads unmapped: too short	759644	674719	771365	950053	788663	614921	1410179	1081555	902701	1286763	815091	960438	1181777
% of reads unmapped: too short	1.65%	1.65%	1.77%	1.98%	2.00%	1.81%	1.96%	2.03%	1.09%	2.69%	2.07%	2.36%	2.55%
Number of reads unmapped: other	16971	16710	26173	24305	15726	13423	25546	25946	17237	24268	22820	24661	12953
% of reads unmapped: other	0.04%	0.04%	0.05%	0.05%	0.04%	0.04%	0.05%	0.05%	0.05%	0.08%	0.05%	0.05%	0.05%
Number of chimeric reads	0	0	0	0	0	0	0	0	0	0	0	0	0
% of chimeric reads	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%

**Counts** Generate a directory called “COUNTS” and copy all gene count files to this folder and then clean all file names.

```
mkdir COUNTS
cp *pass1/*PerGene* COUNTS/
```

```
ls *tab|while read line ; do mv $line ${line/GTFpass1ReadsPerGene.out/} ; done
```

For each count file, extract and create the five count tables.

```
ls *.tab | while read line ; do
echo $line
cat $line | awk 'NR==3{print}' | cut -f 2- > ${line/tab/nofeature\.tab}
cat $line | awk 'NR==4{print}' | cut -f 2- > ${line/tab/ambiguous\.tab}
cat $line | awk 'NR>4{print}' | cut -f 2 > ${line/tab/nostrand\.tab}
cat $line | awk 'NR>4{print}' | cut -f 3 > ${line/tab/sense\.tab}
cat $line | awk 'NR>4{print}' | cut -f 4 > ${line/tab/antisense\.tab}
done
```

Make a Geneid list from one of the count tables as “countsannot\_GRCh38.primary.Selected.Geneid.txt”.

```
ls 008iP22TGFBM_S71.tab | head -1 | while read line; do
cut -f 1 $line | awk 'NR>4{print}' > countsannot_GRCh38.primary.Selected.Geneid.txt
done
```

Create a file listing the names of all samples as “RBarretTNFATGFBsamples.txt”.

```
ls *.sense.tab | sed 's/\.sense.tab//g' | tr -s " " "\n" | sed 's/_1//g' > RBarretTNFATGFBsamples.txt
```

Make count tables for sense, anti-sense, nostrand, ambiguous, and nofeature reads.

```
****
# combine all sense counts into RBarretTNFATGFB_sense.ALL.cnt
paste *.sense.tab > RBarretTNFATGFB_sense.ALL.cnt

# combine all antisense counts into RBarretTNFATGFB_antisense.ALL.cnt
paste *.antisense.tab > RBarretTNFATGFB_antisense.ALL.cnt

# combine all nostrand counts into RBarretTNFATGFB_nostrand.ALL.cnt
paste *.nostrand.tab > RBarretTNFATGFB_nostrand.ALL.cnt

# combine all ambiguous counts into RBarretTNFATGFB_ambiguous.cnt
cat *ambiguous.tab > RBarretTNFATGFB_ambiguous.cnt

# combine all nofeature counts into RBarretTNFATGFB_nofeature.cnt
cat *nofeature.tab > RBarretTNFATGFB_nofeature.cnt
```

Next, I am going to use “RBarretTNFATGFB\_antisense.ALL.cnt” file for the further analysis.

## In MATLAB

**Transfer data and import annotation** Transfer the counts, annotation, and mappability data to your local laptop.

```
RBarretTNFATGFBcnt = textread('RBarretTNFATGFB_antisense.ALL.cnt','');
RBarretsamplesTNFATGFB = textread('RBarretTNFATGFBsamples.txt','%s');
RBarretsampleskeysTNFATGFB = textread('samplekeys_Sam.txt','%s');
```

```
% calculate the sum of the counts in RBarretTNFATGFBcnt, divides the result by 1000000, and
RBarretTNFATGFBmeta_seqdepth=round(sum(RBarretTNFATGFBcnt)/1000000);
```

You can find these annotation in the “mappability and R code” folder.

```
Gencode_33_Selected_MappSS=textread('mappability and R code/gencode.v33.Selected.ReadsPerGene.txt','s');
Gencode_33_Selected_MappUS=textread('mappability and R code/gencode.v33.Selected.ReadsPerGene.txt','s');
Gencode_33_Selected_Geneid=textread('mappability and R code/gencode.v33.annotation.Selected.Geneid.txt','s');
Gencode_33_Selected_Biotype=textread('mappability and R code/gencode.v33.annotation.Selected.Biotype.txt','s');
Gencode_33_Selected_Genename=textread('mappability and R code/gencode.v33.annotation.Selected.Genename.txt','s');
```

**Compile counts** First, initialize a new variable called RBarretTNFATGFBTPM with the same count data as RBarretTNFATGFBcnt. Then, iterates over each gene in the count data matrix. For each gene, the corresponding row in RBarretTNFATGFBTPM is updated by dividing the count data by the read counts from the “Gencode\_33\_Selected\_MappSS”, multiplying by 1000, and storing the result in RBarretTNFATGFBTPM.

Finally, iterates over each sample in the TPM data matrix. For each sample, the corresponding column in RBarretTNFATGFBTPM is updated by dividing the values in the column by the sum of the values in the column, multiplying by 1,000,000, and storing the result in RBarretTNFATGFBTPM. This step **normalizes the TPM values** across samples and scales the resulting values to TPM.

```
RBarretTNFATGFBTPM = RBarretTNFATGFBcnt;
for i=1:size(RBarretTNFATGFBcnt,1)
% divid the gene count matrix RBarretTNFATGFBcnt by Gencode_33_Selected_MappSS matrix, which
RBarretTNFATGFBTPM(i,:) = RBarretTNFATGFBcnt(i,:)/Gencode_33_Selected_MappSS(i)*1000;
end
% set any NaN or Inf values resulting from the normalization process to 0
RBarretTNFATGFBTPM(isnan(RBarretTNFATGFBTPM)) = 0;
RBarretTNFATGFBTPM(isinf(RBarretTNFATGFBTPM)) = 0;
for i=1:size(RBarretTNFATGFBTPM,2)
% scale the TPM values so that the sum of expression values across each sample of the matrix
RBarretTNFATGFBTPM(:,i) = RBarretTNFATGFBTPM(:,i)/sum(RBarretTNFATGFBTPM(:,i))*1000000;
end
```

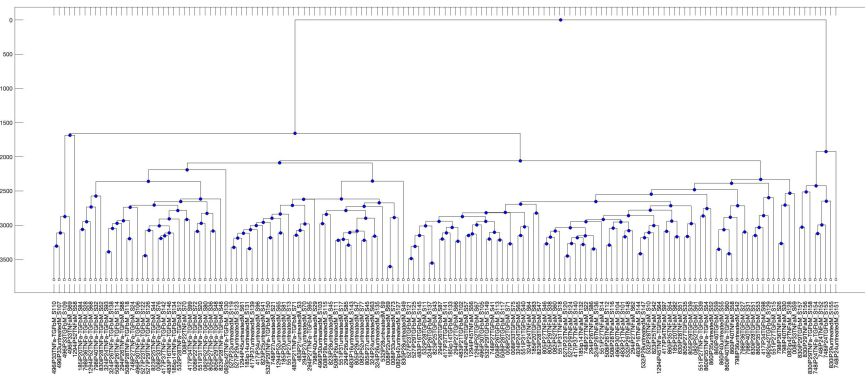


**Make the first dendrogram** Make a dendrogram to visualize the relationships among samples in the RBarretTNFATGFB dataset based on their gene expression profiles.

1. I generates a random selection of 1,000 genes from the TPM data matrix.
2. Calculates the pairwise distances between the selected genes.
3. Creates a for loop that iterates 9,999 times. For each iteration, a new random selection of 1,000 genes is generated, and the pairwise distances between these genes are added to the previous 'thisdist' calculation.
4. Converts the one-dimensional distance vector 'thisdist' into a distance matrix 'thisdistmat' using the 'squareform' function.
5. Generates a hierarchical clustering tree based on the distance matrix 'thisdistmat'.

Overall, I perform a clustering analysis on a subset of genes in the RBarretTNFATGFB dataset to visualize the relationships among samples based on their gene expression profiles.

```
thisrand = unique(randi([1 size(RBarretTNFATGFBTPM,1)],1,1000));
thisdist = pdist(RBarretTNFATGFBTPM(thisrand,:))';
for i=1:9999
    thisrand = unique(randi([1 size(RBarretTNFATGFBTPM,1)],1,1000));
    thisdist = thisdist+pdist(RBarretTNFATGFBTPM(thisrand,:))';
end
thisdistmat = squareform(thisdist/10000);
thistree = seqlinkage(thisdistmat,'average', RBarretsamplesTNFATGFB)
plot(thistree, 'ORIENTATION', 'top')
```



**All biotypes counts percents** This part of codes is performing all biotypes counts percents across multiple samples.

```
% use the unique function and stored the allbiotypes variable.
allbiotypes = unique(Gencode_33_Selected_Biotype);

% create A cell array allbiotypeslength to store the lengths of each biotype name.
allbiotypeslength = cell(length(allbiotypes),1);

% two new matrices, allbiotypescounts and allbiotypescountspercents, are initialized with zeros.
allbiotypescounts = zeros(length(allbiotypes),size(RBarretTNFATGFBCnt,2));
allbiotypescountspercents = zeros(length(allbiotypes),size(RBarretTNFATGFBCnt,2));

for i=1:length(allbiotypes)
% finds all the indices of Gencode_33_Selected_Biotype that match the current biotype. Then
temp = strmatch(allbiotypes{i}, Gencode_33_Selected_Biotype);
% Stored the length of the temp vector represents the number of genes with the current biotype.
allbiotypeslength{i} = length(temp);
if length(temp)>1
% Sum the expression values for all genes with the current biotype across all samples. The result is stored in allbiotypescounts.
allbiotypescounts(i,:) = sum(RBarretTNFATGFBCnt(strmatch(allbiotypes{i}, Gencode_33_Selected_Biotype)),2);
allbiotypescountspercents(i,:) = allbiotypescounts(i,:)/sum(RBarretTNFATGFBCnt)*100;
end
end

dlmwrite('allbiotypescountspercents.txt', allbiotypescountspercents,'delimiter','\t')
writetable(cell2table(allbiotypes),'allbiotypes.txt','WriteVariableNames',0)
```

Using Excel, create a spreadsheet using “allbiotypes.txt” and “allbiotypescountspercents.txt”, and calculate the minimum, maximum, and average values for each biotype. You can access my completed spreadsheet [here](#). Notably, protein\_coding genes exhibit an average of 98.65% among the various biotypes, consistent with my expectations.

Biotypes	min	max	average
IG_C_gene	0.000000000	0.000446580	0.000013340
IG_D_gene	0.000000000	0.000000000	0.000000000
IG_J_gene	0.000000000	0.000013813	0.000000235
IG_V_gene	0.000000000	0.000073058	0.000010604
Mt_rRNA	0.228350000	1.165600000	0.487934371
Mt_tRNA	0.000141280	0.002049800	0.000492376
TEC	0.011838000	0.029033000	0.018976556
TR_C_gene	0.000000000	0.002667100	0.000287637
TR_D_gene	0.000000000	0.000011542	0.000000580
TR_J_gene	0.000000000	0.000076445	0.000009219
TR_V_gene	0.000000000	0.000061956	0.000011904
lncRNA	0.537580000	1.612500000	0.825408079
miRNA	0.002352100	0.014811000	0.004202784
misc_RNA	0.000641580	0.002443400	0.001048973
<b>protein_coding</b>	<b>97.675000000</b>	<b>99.144000000</b>	<b>98.658132450</b>
rRNA	0.000000000	0.000028828	0.000005066
ribozyme	0.000000000	0.000007357	0.000000293
sRNA	0.000000000	0.000004680	0.000000135
scRNA	0.000000000	0.000000000	0.000000000
scaRNA	0.000002616	0.000324330	0.000036518
snRNA	0.000262070	0.001464200	0.000543351
snoRNA	0.001727500	0.004335800	0.002883131
vaultRNA	0.000000000	0.000000000	0.000000000

**Protein coding genes** The “allbiotypes.txt” file contains multiple biotypes. For the next step, I will only retain the “protein\_coding” biotype.

```
allbiotypes=unique(Gencode_33_Selected_Biotype);
% finds the 15th unique value, which is protein_coding, of Gencode_33_Selected_Biotype in the
proteinencodingindx = strmatch(allbiotypes{15}, Gencode_33_Selected_Biotype);
biotypeindx = proteinencodingindx;

% creates an array additionalgenes contains the indices of genes that have certain prefixes
additionalgenes = [strmatch('MT-',Gencode_33_Selected_Genename) ; strmatch('H1',Gencode_33_Selected_Genename)];
```

```

% creates an array nonadditionalgenes with the same length as the Gencode_33_Selected_Genename
nonadditionalgenes = 1:length(Gencode_33_Selected_Genename);
% remove the indices of genes in additionalgenes from the nonadditionalgenes array.
nonadditionalgenes(additionalgenes) = [];

% mappableindx contains the indices of elements in the Gencode_33_Selected_MappSS array that
mappableindx = find(Gencode_33_Selected_MappSS>50);

% a new variable finalIndexGeneric which is the intersection of three other variables: biotype,
finalIndexGeneric = intersect(biotypeindx,intersect(nonadditionalgenes,mappableindx));
% find the indices of rows in RBarretTNFATGFBcnt that have a sum greater than 150.
countindx = find(sum(RBarretTNFATGFBcnt')>150);

% update finalIndexGeneric to be the intersection of finalIndexGeneric and countindx.
finalIndexGeneric=intersect(finalIndexGeneric,countindx);

% create a new variable RBarretTNFATGFBcnt_GMask which is a subset of RBarretTNFATGFBcnt corresponding to
RBarretTNFATGFBcnt_GMask = RBarretTNFATGFBcnt(finalIndexGeneric,:);

Gencode_33_Selected_Geneid_GMask = Gencode_33_Selected_Geneid(finalIndexGeneric);
Gencode_33_Selected_Genename_GMask = Gencode_33_Selected_Genename(finalIndexGeneric);
Gencode_33_Selected_MappSS_GMask = Gencode_33_Selected_MappSS(finalIndexGeneric);
Gencode_33_Selected_MappUS_GMask = Gencode_33_Selected_MappUS(finalIndexGeneric);

% normalize the expression data like we do previously
RBarretTNFATGFBExpression_GMask = RBarretTNFATGFBcnt_GMask;
for i=1:size(RBarretTNFATGFBExpression_GMask,2)
    RBarretTNFATGFBExpression_GMask(:,i) = RBarretTNFATGFBcnt_GMask(:,i)/sum(RBarretTNFATGFBcnt_GMask(:,i));
end
for i=1:size(RBarretTNFATGFBExpression_GMask,1)
    RBarretTNFATGFBExpression_GMask(i,:) = RBarretTNFATGFBExpression_GMask(i,:)/Gencode_33_Selected_MappSS_GMask(i,:);
end
RBarretTNFATGFBExpression_GMask(isnan(RBarretTNFATGFBExpression_GMask)) = 0;
RBarretTNFATGFBExpression_GMask(isinf(RBarretTNFATGFBExpression_GMask)) = 0;

% RBarretTNFATGFBcpm_GMask contains the expression data normalized only by CPM, using the same method as above
RBarretTNFATGFBcpm_GMask = zeros(size(RBarretTNFATGFBcnt_GMask));
for i=1:size(RBarretTNFATGFBcnt_GMask,2)
    RBarretTNFATGFBcpm_GMask(:,i) = RBarretTNFATGFBcnt_GMask(:,i)/sum(RBarretTNFATGFBcnt_GMask(:,i));
end

% RBarretTNFATGFBtpm_GMask contains the expression data normalized only by TPM.
RBarretTNFATGFBtpm_GMask = RBarretTNFATGFBcnt_GMask;
for i=1:size(RBarretTNFATGFBcnt_GMask,1)
    RBarretTNFATGFBtpm_GMask(i,:) = RBarretTNFATGFBcnt_GMask(i,:)/Gencode_33_Selected_MappSS_GMask(i,:);
end

```

```

end
RBarretTNFATGFBTPM_GMask(isnan(RBarretTNFATGFBTPM_GMask)) = 0;
RBarretTNFATGFBTPM_GMask(isinf(RBarretTNFATGFBTPM_GMask)) = 0;
for i=1:size(RBarretTNFATGFBTPM_GMask,2)
RBarretTNFATGFBTPM_GMask(:,i) = RBarretTNFATGFBTPM_GMask(:,i)/sum(RBarretTNFATGFBTPM_GMask(:,i));
end

```

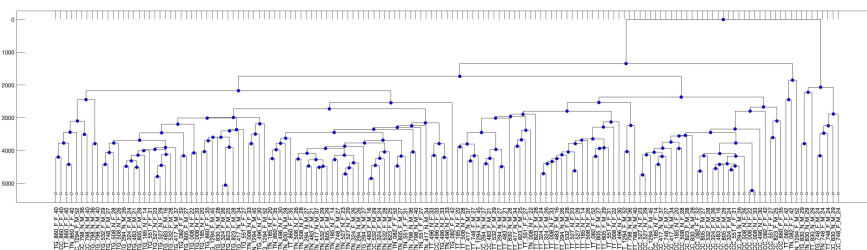
**Dendrogram with only protein coding genes** Perform hierarchical clustering on a subset of the gene expression data stored in the variable RBarretTNFATGFBTPM\_GMask with only protein coding genes.

```

thisrand = unique(randi([1 size(RBarretTNFATGFBTPM_GMask,1)],1,1000));
thisdist = pdist(RBarretTNFATGFBTPM_GMask(thisrand,:));
for i=1:9999
thisrand = unique(randi([1 size(RBarretTNFATGFBTPM_GMask,1)],1,1000));
thisdist = thisdist+pdist(RBarretTNFATGFBTPM_GMask(thisrand,:));
end
thisdistmat = squareform(thisdist/10000);
thistree = seqlinkage(thisdistmat,'average', RBarretsampleskeysTNFATGFB);
plot(thistree,'ORIENTATION','top')

```

Basically, the plot is clustered by their treatments.



**The percent of the top 100 genes** Calculates the top 100 expressed genes in each sample based on their transcript per million (TPM) values in the RBarretTNFATGFBTPM\_GMask matrix.

```

% iterates over 151 samples, it first sorts the TPM values of all genes in descending order
yall=[];
for i=1:151
[x y]=sort(RBarretTNFATGFBTPM_GMask(:,i),'descend');
yall=unique([y(1:100); yall]);
top100percent(i)=sum(RBarretTNFATGFBTPM_GMask(y(1:100),i))/1000000;
end

```

After calculating the sum of the percent of the top 100 genes, it is 58.25%.

In the end, we store the Gencode\_33\_Selected\_Geneid\_GMask.txt, Gencode\_33\_Selected\_Genename\_GMask.txt, Gencode\_33\_Selected\_MappSS\_GMask.txt, RBarretTNFATGFBTPM\_GMask.txt, and RBarretTNFATGFBcnt\_GMask.txt for our further analysis in R.

```
writetable(cell2table(Gencode_33_Selected_Geneid_GMask), 'Gencode_33_Selected_Geneid_GMask.txt')
writetable(cell2table(Gencode_33_Selected_Genename_GMask), 'Gencode_33_Selected_Genename_GMask.txt')
dmlwrite('Gencode_33_Selected_MappSS_GMask.txt', Gencode_33_Selected_MappSS_GMask, 'delimiter', '\t')
dmlwrite('RBarretTNFATGFBTPM_GMask.txt', RBarretTNFATGFBTPM_GMask, 'delimiter', '\t')
dmlwrite('RBarretTNFATGFBcnt_GMask.txt', RBarretTNFATGFBcnt_GMask, 'delimiter', '\t')
```

## In R

**Install packages** First, install packages BiocManager, BiocLite, IHW, DESeq2, and ggplot2. Then, read in RBarretTNFATGFBcnt\_GMask.txt, RBarretTNFATGFBsamples.txt, samplekeys\_Sam.txt, and Gencode\_33\_Selected\_Genename\_GMask.txt.

```
setwd("/Users/LuC/Desktop/Cedars-Sinai/PROJECTS/IBD_RNASeq/RBARRETTNFATGFB/")#setwd("/Users/LuC/Desktop/Cedars-Sinai/PROJECTS/IBD_RNASeq/RBARRETTNFATGFB/")
```

**Form samplekeys\_Sam.tab** Separate samplekeys\_Sam.txt by “\_” to get samplekeys\_Sam.tab before next step. Here are my code in terminal.

```
#In terminal
#Create an empty file to store the output
touch samplekeys_Sam.tab

#Loop over the sample names and split them by "_"
for sample in $(cat samplekeys_Sam.txt); do
    IFS=_ read -r col1 col2 col3 col4 col5 <<< "$sample"
    echo -e "$col1\t$col2\t$col3\t$col4\t$col5" >> samplekeys_Sam.tab
done
```

**Generate a sampleTableTNFATGFB** The sampleTableTNFATGFB contains Treatment, Line, Pheno, Sex, Pass, Factor, and Batch.

```
sampleTableTNFATGFB = read.table("samplekeys_Sam.tab")rownames(sampleTableTNFATGFB)<-samplekeys_Sam.txt
```

	Treatment	Line	Pheno	Sex	Pass	Factor	Batch
TG_008_N_F_22	TG	8	N	F	22	TG_N	8_22
TT_008_N_F_22	TT	8	N	F	22	TT_N	8_22
TN_008_N_F_22	TN	8	N	F	22	TN_N	8_22
CC_008_N_F_22	CC	8	N	F	22	CC_N	8_22
TG_008_N_F_33	TG	8	N	F	33	TG_N	8_33
TT_008_N_F_33	TT	8	N	F	33	TT_N	8_33
TN_008_N_F_33	TN	8	N	F	33	TN_N	8_33
CC_008_N_F_33	CC	8	N	F	33	CC_N	8_33
TG_082_F_F_52	TG	82	F	F	52	TG_F	82_52
TT_082_F_F_52	TT	82	F	F	52	TT_F	82_52
TN_082_F_F_52	TN	82	F	F	52	TN_F	82_52
CC_082_F_F_52	CC	82	F	F	52	CC_F	82_52

**DESeq2 package** The `RBarretTNFATGFBCntGMaskBatch` is created with the **Batch** information specified in the design formula, while the `RBarretTNFATGFBCntGMaskFactor` is created with the **treatment and phenotype** information specified in the design formula.

Next, the `DESeq` function is used to estimate size factors and dispersion values for the `DESeqDataSet` objects.

Finally, the `varianceStabilizingTransformation` function is used to perform variance stabilizing transformation on the `DESeqDataSet` objects. This transformation is important for reducing the effect of noise and heteroscedasticity in the data, making it more suitable for downstream analyses such as differential gene expression analysis.

```
RBarretTNFATGFBCntGMaskBatch <- DESeqDataSetFromMatrix(RBarretTNFATGFBCntGMask, colData= sam
```

## PCA

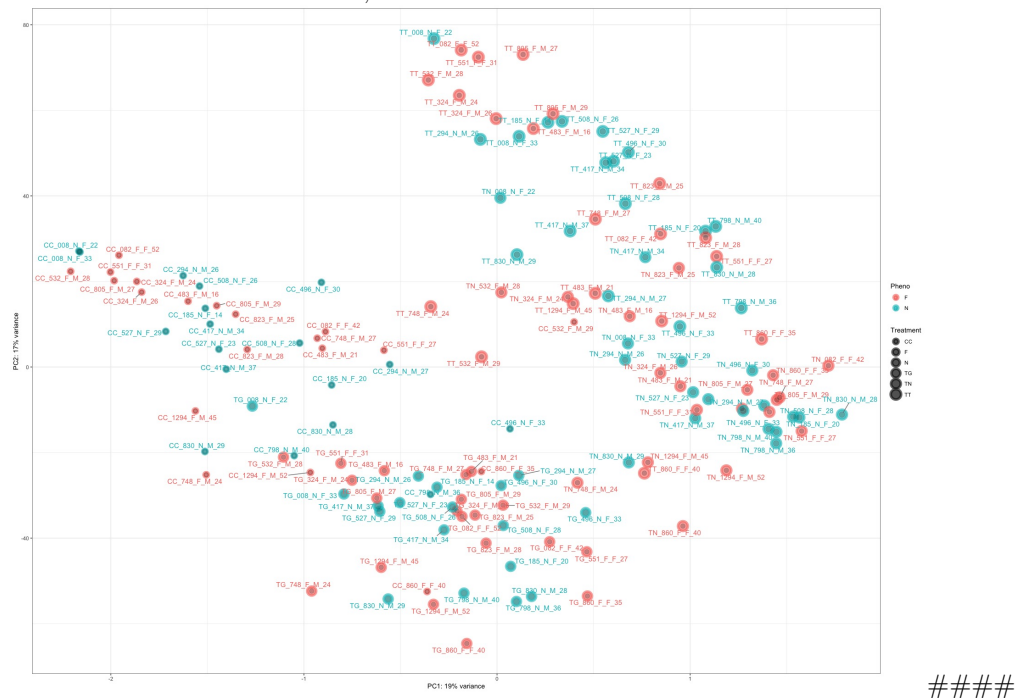
```
#perform principal component analysis (PCA) on the variance-stabilized counts data pcabatch
#give the percentage of variance explained by each principal component percentVarbatch <- rou
#pcabatch$rotation is a matrix containing the loadings of the principal components. aloadbat
#normalize the loadings in aloadbatch so that each column (i.e., PC) sums to 1. aloadrelativ
#pcabatch$x is a matrix containing each sample's coordinate on each principal component pcaba
#center the PC1 scores in pcabatchR to have a mean of 0. This is done so that the PC1 variab
```

## PCA plots

```
ggplot(pcabatchR, aes(PC1, PC2, color= Pheno)) + geom_point(aes(size= Treatment), alpha=0.6,
```

The plot shows the relationship between **PC1** and **PC2** colored by **Pheno** variable, with the point size indicating the **Treatment** variable.

Four distinct groups were formed based on their treatment: the CC group (untreated) is located in the right corner, the TG group (treated with TGF-b) is located at the bottom, the TN group (treated with TNF-a) is located in the right corner, and the TT group (treated with both TGF-b and TNF-a) is located at the top. These groups were differentiated based on PC1, which accounted for 19% of the variance.



A series of bar plots (one for each principal component)

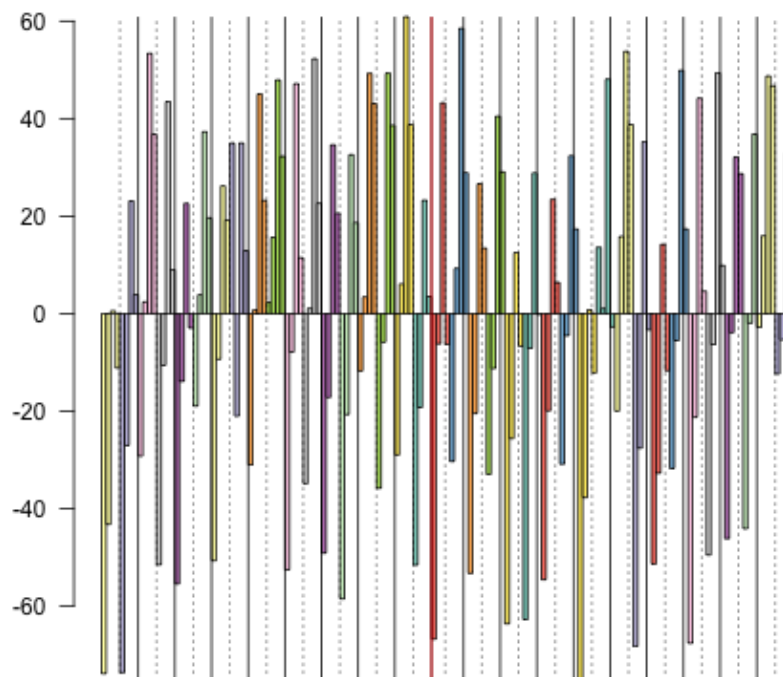
Each bar plot represents the loadings of all samples on a given principal component.

#the resulting vector could will contain 12 colors from the "Set3" palette.library(RColorBrewer)

A red line is drawn at position 72 in order to separate the non-fibrotic group and the fibrotic group. Within each patient, the treatment order would be CC, TG, TN, TT. The color of each bar represents the batch of the sample, with a unique color assigned to each batch. The vertical lines on the plot indicate the position of specific loadings, with thin and thick lines indicating different positions.



For example, this is PC\_1.png. From this plot, you can see that the highest expression is closely related with TNF-a. As for the first patient, compare to the control(untreated), the TNF-a group is much higher and the TT group (TGF-b+TNF-a) is not that high.



**PCA rank matrix** Take csv files and converts it to the txt files with the second column onwards. It does this by first removing the first row using awk, replacing multiple commas with tabs using tr, and removing the first column using cut.

#In terminal

```
cat aloadrelativeMask_batchmodel_filtered.csv | awk 'NR>1{print}' | tr -s "," "\t" | cut -f 2- > pca_batchmodel_x.txt
cat pca_batchmodel_x.csv | awk 'NR>1{print}' | tr -s "," "\t" | cut -f 2- > pca_batchmodel_y.txt
```

Read in the preprocessed data files created in the previous steps and store them in variables pcabatch\_samples and pca\_loadings, respectively.

```
#In Matlab
pcabatch_samples = textread('pca_batchmodel_x.clean.txt','');
pca_loadings = textread('aloadrelativeMask_batchmodel_filtered.clean.txt','');
```

Sort the three columns of `pca_loadings` in descending order and store the sorted values in variables `x1`, `x2`, and `x3`, and the corresponding indices in `y1`, `y2`, and `y3`.

```
%%%%%%%%%% PCA SUPP
```

```
%x = pca_loading number, y = its index
```

```
[x1 y1]=sort(pca_loadings(:,1),'descend');
[x2 y2]=sort(pca_loadings(:,2),'descend');
[x3 y3]=sort(pca_loadings(:,3),'descend');
```

Determine the rank of each row in the original order for the first three principal components and store the ranks in a matrix `pcarankmatrix`.

```
[x y z]=intersect(1:length(y1),y1);
pcarankmatrix(:,1)=z;
[x y z]=intersect(1:length(y2),y2);
pcarankmatrix(:,2)=z;
[x y z]=intersect(1:length(y3),y3);
pcarankmatrix(:,3)=z;
```

```
%contains the rank of each feature in the original order for the first three principal comp
```

```
dlmwrite('pcarankmatrix.txt', pcarankmatrix,'delimiter','\t')
```

To efficiently manage our data with a single glance, I have organized it into an Excel spreadsheet using a combination of command line, Excel, and R.

**Spreadsheet** The first sheet (patients) built on excel contains patients order, patients id, phenotypes, and sex. You can visit the sheet by clicking [here](#).

The second sheet (allbiotypes\_percents) includes the names and percentages of all biotypes, along with their respective minimum, maximum, and average values, providing us with a comprehensive overview. You can visit the sheet by [clicking here](#).

```
paste allbiotypes allbiotypescountspercents > combine_allbiotypes_percents.txt
```

```
#In R #allbiotypes_percentssheet2_1 <- list("Biotypes")sheet2_2 <- sampleKeyTNFATGFBcombined
#add their respective minimum, maximum, and average values on Excel
```

The third sheet is the main sheet that includes Genename, Geneid, Mapp, PC1, PC2, PC3, and patient's TPM values.

```
paste Gencode_33_Selected_Genename_GMask.txt Gencode_33_Selected_Genename_GMask.txt Gencode_
#In R
```

#spreadsheet

```
sheet3_1 <- list("Genename", "Genename", "Geneid", "Mapp", "PC1", "PC2", "PC3")sheet3_2<- sampleK
```

You can sort this sheet with PC1, PC2, and so on to see the correlation between the treatment and the expression level in each gene.

GeneName	Genome	Contest	Mapp	PC1	PC2	PC3	TG_OBR_N_F_22	TT_OBR_N_F_22	TN_OBR_N_F_22	GC_OBR_N_F_22	TG_OBR_N_F_33	TT_OBR_N_F_33	TN_OBR_N_F_33	GC_OBR_N_F_33	TG_OBR_F_F_TT_OBR_F_F	TT_OBR_F_F_TT_OBR_F_F	TN_OBR_F_F_TT_OBR_F_F	GC_OBR_F_F_TT_OBR_F_F				
CXCL6	FXS000000104875.10	104875.10	68.00	1.00	887.00	32.40	10887.00	5693.30	31.87	12.82	4951.60	870.40	31.70	13.20	4292.20	149.90	124.10	1.40	4126.20	843.90	7.07	
CXCL8	FXS000000104820.10	104820.10	1.00	2.00	380.00	0.30	4089.30	1729.20	0.40	0.30	1642.70	167.40	0.30	2.21	2033.80	804.61	7.38	0.67	4032.80	8072.00	0.46	
CXCL3	FXS000000103789.10	103789.10	111.00	8.00	1782.00	10.80	5487.70	2412.60	18.21	12.17	1873.10	911.30	19.33	17.31	2760.30	149.20	76.39	3.44	1397.00	1362.50	32.78	
C3	FXS000000102739.10	102739.10	581.00	21.00	4.00	531.00	9.10	948.90	238.21	0.64	0.09	330.99	23.94	0.71	0.10	310.76	2.27	0.09	1.99	117.54	11.13	18.66
ABHD8P	FXS000000104175.17	104175.17	8990.00	1104.00	5.00	1275.00	40.70	412.98	49.70	48.11	10.23	636.79	59.99	49.44	7.72	621.97	40.02	45.88	0.68	15.85	4.66	3.58
GNP	FXS000000105964.11	105964.11	2779.00	147.00	8.00	2037.00	167.00	3.04	62.18	0.06	34.12	0.19	302.38	0.08	54.69	0.00	38.93	0.08	18.30	0.11	12.97	0.12
IL1A	FXS000000104000.10	104000.10	5644.00	7.00	132.00	0.42	1.41	0.13	0.14	10.08	0.78	0.10	5.12	0.04	64.36	0.31	1.03	0.48	1032.30	15.43	20.23	11.67
CXCL3	FXS000000103734.4	103734.4	138.00	138.00	8.00	750.00	3.01	1397.70	799.40	1.76	3.72	185.10	0.89	0.10	514.44	11.46	1.84	1.84	1.18	3769.20	322.10	4.70
ISG15	FXS000000102096.18	102096.18	1879.00	56.00	9.00	117.00	12.40	1272.70	708.11	21.84	6.93	913.08	202.64	21.62	8.31	1386.80	84.31	10.23	0.63	706.20	230.10	41.28
CXCL5	FXS000000103793.7	103793.7	2337.00	81.00	10.00	530.00	0.34	880.40	221.78	0.64	0.17	151.01	14.67	0.17	0.00	95.85	2.82	0.05	0.31	1987.60	87.10	0.61
CCL5	FXS000000102084.10	102084.10	1373.00	1.00	11.00	21.00	0.16	206.41	212.72	0.00	0.29	549.11	381.46	0.07	0.14	494.50	133.94	0.00	0.04	96.00	833.98	0.00
PTG2	FXS000000103661.4	103661.4	1781.00	3156.00	12.00	4555.00	1346.10	526.50	3919.60	2158.20	184.73	6763.80	708.90	2271.00	55.27	3225.10	140.89	150.15	41.80	448.01	43.16	1158.50
CCL2	FXS000000100801.10	100801.10	4.00	13.00	172.00	1.70	1063.10	894.13	1.22	1.09	1280.90	806.48	1.90	6.20	1766.40	133.94	5.13	0.39	142.81	62.48	19.16	
CXCL2	FXS000000103413.1	103413.1	1038.00	393.00	14.00	1285.00	1.45	904.10	368.30	4.76	1.59	304.28	32.38	3.71	0.68	308.55	3.36	2.07	1.00	1236.00	140.94	9.98
VICAM1	FXS000000102021.12	102021.12	1284.00	592.00	15.00	2405.00	62.11	1263.00	462.26	284.87	16.17	2822.00	343.05	206.46	102.71	1191.00	765.00	290.26	46.82	537.60	479.14	19.14
MX1	FXS000000107001.14	107001.14	4535.00	413.00	16.00	495.00	1.44	278.78	62.79	2.31	2.09	116.67	6.18	2.21	0.61	487.83	9.96	15.53	1.48	37.89	22.66	19.08
IFP1	FXS000000102706.10	102706.10	921.00	125.00	17.00	540.00	107.01	1248.00	392.40	264.68	4056.10	307.41	265.16	102.71	1191.00	765.00	290.26	46.82	537.60	479.14	19.14	
TFEB	FXS000000103062.13	103062.13	6842.00	1074.00	18.00	844.00	0.13	74.20	3.17	0.37	0.54	76.34	8.95	0.22	0.10	139.07	0.62	0.40	0.09	0.14	0.08	0.02
IGFBP1	FXS000000103039.19	103039.19	2624.00	22.00	17.00	2180.00	5.56	1055.40	420.16	17.80	1.54	1145.40	421.10	29.86	27.86	1059.80	1064.70	129.72	7.60	482.29	905.00	50.26
BST2	FXS000000103083.13	103083.13	900.00	145.00	20.00	390.00	12.06	1170.20	201.40	14.94	0.55	151.86	71.96	12.86	0.55	1587.70	7.71	5.38	2.72	54.49	50.72	27.63
DNAI3	FXS000000102101.10	102101.10	7094.00	127.00	21.00	471.00	1.30	125.87	198.60	0.12	0.54	61.60	8.19	1.99	0.12	194.17	11.51	2.29	0.71	24.47	22.70	8.72
TRPA1	FXS000000103823.11	103823.11	5002.00	88.00	21.00	3508.00	0.94	17.40	183.22	0.04	0.03	98.37	83.83	0.06	0.10	35.11	0.33	0.02	0.03	355.15	55.77	3.51
MMP1	FXS000000103812.13	103812.13	2844.00	21.00	27.00	2816.00	3.10	137.17	938.10	0.12	0.34	141.11	1.24	0.05	105.16	2.07	0.09	0.26	2862.70	252.80	0.54	
DNAI1	FXS000000103817.13	103817.13	4894.00	315.00	24.00	484.00	0.25	83.85	111.32	1.32	0.13	88.14	2.10	1.44	0.10	149.43	3.93	2.78	1.11	14.71	15.95	14.70
HNRB2	FXS000000103898.4	103898.4	6476.00	5115.00	21.00	2145.00	3.46	289.40	107.19	35.96	1.17	169.10	146.70	15.36	1.84	550.72	12.12	119.07	9.99	297.11	165.18	64.92
TNFAIP2	FXS000000103113.11	103113.11	4873.00	47.00	26.00	786.00	2.96	480.20	114.66	7.49	2.45	324.42	10.14	7.21	8.45	493.27	60.84	0.80	0.72	108.12	66.00	2.27
PTG1S	FXS000000103484.11	103484.11	1468.00	1036.00	27.00	977.00	2.78	937.77	109.45	120.14	6.40	201.08	56.81	130.62	1.55	477.40	9.48	25.17	3.77	188.15	18.91	38.80
CMA2	FXS000000103101.12	103101.12	841.00	26.00	40.00	204.00	0.06	184.60	4.90	0.07	0.10	27.13	0.40	0.08	0.05	100.23	0.33	0.10	0.09	1.94	1.64	0.82
TNFAIP3	FXS000000103039.19	103039.19	5175.00	16.00	29.00	1440.00	11.01	588.75	471.41	12.45	12.07	527.70	324.42	12.11	5.71	743.53	304.41	0.40	7.88	186.80	359.82	18.54
COL1A1	FXS000000103039.19	103039.19	3719.00	426.00	30.00	1240.00	14.14	8.84	4.44	0.23	0.10	39.39	0.16	0.45	0.00	2.38	0.00	0.22	2.40	4.16	0.00	
CF	FXS000000103039.19	103039.19	1583.00	100.00	11.00	8610.00	15.94	325.09	80.13	57.83	26.44	99.34	57.94	77.74	1392.20	45.82	580.07	1.89	0.91	0.12	76.14	
EPST11	FXS000000103039.14	103039.14	1381.00	1297.00	12.00	1530.00	15.06	125.10	178.12	30.64	4.82	211.47	49.10	26.46	4.30	777.19	18.16	18.40	0.10	25.39	13.08	10.65
ANG01	FXS000000102101.12	102101.12	1812.00	979.00	13.00	1124.00	85.09	2.80	19.76	0.63	241.05	9.27	10.30	0.88	219.43	1.02	98.02	1.15	156.78	1.19	5.04	11.22
SUCRA2	FXS000000103039.17	103039.17	8019.00	1297.00	14.00	1464.00	0.19	201.51	78.47	60.84	1.73	101.47	56.86	1.11	94.47	279.42	50.21	1.51	405.95	792.75	57.41	
CN11	FXS000000103039.11	103039.11	1713.00	139.00	15.00	1053.00	10.47	0.13	228.89	0.63	52.76	0.46	230.10	0.62	103.95	0.19	138.63	0.74	776.71	8.49	777.11	1.70
SUCRA8	FXS000000103039.13	103039.13	1573.00	398.00	36.00	154.00	6.01	121.69	96.33	1.46	1.16	141.62	15.82	2.16	0.84	231.80	1.71	3.36	1.61	31.61	11.17	1.46
BIRC3	FXS000000102345.15	102345.15	7797.00	8.00	37.00	468.00	0.20	42.90	23.12	0.22	0.05	47.92	19.87	0.14	0.12	49.18	19.24	0.23	0.08	30.17	29.72	0.44
CDKN11	FXS000000102307.15	102307.15	4440.00	46.00	343.00	198.27	802.70	370.26	2146.20	0.00	107.42	18.86	0.00	0.00	361.14	1.06	0.20	0.00	12.87	6.02	0.00	
CXCL10	FXS000000102345.15	102345.15	3076.00	398.00	39.00	406.00	0.00	182.1	10.80	0.00	0.00	107.42	18.86	0.00	0.00	361.14	1.06	0.20	0.00	12.87	6.02	0.00
IFIT1	FXS000000102307.15	102307.15	2424.00	343.00	14.00	1424.00	24.62	816.70	12.64	59.05	6.90	380.20	25.11	5.13	2.82	208.67	14.63	107.10	18.20	47.48	45.39	18.13
IFAH	FXS000000102309.18	102309.18	4048.00	539.00	41.00	413.00	4.36	126.48	46.30	4.59	2.29	95.52	5.01	5.34	2.67	304.69	7.35	5.97	0.24	11.16	7.95	2.26
IL32	FXS000000100017.14	100017.14	2810.00	2.00	42.00	1031.00	0.20	194.80	142.13	0.87	1.08	412.84	329.95	0.27	3.48	412.81	483.17	4.46	4.01	114.71	305.85	3.70
LOC404	FXS000000100017.14	100017.14	480.00	41.00	41.00	91.00	0.61	0.02	0.10	0.02	0.10	0.02	0.10	0.12	15.34	1.17	14.01	0.46	14.29	0.20	0.20	
IFP1	FXS000000103484.11	103484.11	2208.00	1773.00	44.00	4440.00	4.82	862.12	113.00	146.01	12.20	111.61	32.33	143.41	0.30	89.01	1.96	4.22	0.42	1.83	0.46	1.14
CTSL	FXS000000103113.11	103113.11	3911.00	94.00	41.00	1064.00	0.17	180.80	184.1	0.90	0.10	84.66	16.63	0.10	0.10	129.91	10.59	2.78	0.09	79.83	96.88	1.18
CH2E2	FXS000000104884.14	104884.14	3551.00	224.00	46.00	1133.00	0.19	11.15	10.11	0.46	0.10	28.93	7.60	0.14	0.06	19.92	0.81	0.42	0.17	7.14	1.01	0.79
IL6	FXS000000103039.17	103039.17	1746.00	41.00	17.00	3060.00	0.16	403.09	10.78	0.00	0.13	61.80	1.11	0.12	0.67	117.20	33.02	0.12	0.27	303.18	47.11	0.11
TNFAIP3	FXS000000103039.19	103039.19	4819.00	48.00	48.00	608.00	3.01	88.89	23.19	0.97	0.03	88.14	5.92	7.44	0.39	88.43	0.84	3.78	1.00	27.89	2.10	1.18
PTG1S	FXS000000103039.19	103039.19	3993.00	48.00	18																	

8. Ingber, D.E., Human organs-on-chips for disease modelling, drug development and personalized medicine. *Nat Rev Genet*, 2022. 23(8): p. 467-491.
9. Brooks, I.R., et al., Functional genomics and the future of iPSCs in disease modeling. *Stem Cell Reports*, 2022. 17(5): p. 1033-1047.
10. D'Alessio, S., et al., Revisiting fibrosis in inflammatory bowel disease: the gut thickens. *Nat Rev Gastroenterol Hepatol*, 2022. 19(3): p. 169-184.
11. Carcamo-Orive, I., et al., Analysis of Transcriptional Variability in a Large Human iPSC Library Reveals Genetic and Non-genetic Determinants of Heterogeneity. *Cell Stem Cell*, 2017. 20(4): p. 518-532 e9.
12. Gleeson, J.P., et al., Development of Physiologically Responsive Human iPSC-Derived Intestinal Epithelium to Study Barrier Dysfunction in IBD. *Int J Mol Sci*, 2020. 21(4).
13. Workman, M.J., et al., Modeling Intestinal Epithelial Response to Interferon-gamma in Induced Pluripotent Stem Cell-Derived Human Intestinal Organoids. *Int J Mol Sci*, 2020. 22(1).