

Global proteomics analysis of COVID + vs COVID - plasma samples to decipher the host response towards COVID pathogenesis

Submitted by-

Archita Dev Barman (PAWS20014)

Vellore Institute of Technology, Vellore

B. Tech in Biotechnology

School of Biosciences and Technology

OBJECTIVE

To analyze the given dataset of plasma samples of both COVID + and COVID – samples and determine the changes in the proteome to find proteins (if present) and their relation to biochemical pathways which corresponds to host response towards COVID pathogenesis.

INDEX

- Introduction
- GISAID Database
- Secondary Data Analysis using Metaboanalyst 5.0
- Pathway Analysis using
 - (i) Reactome
 - (ii) Metascape
 - (iii) Panther
- Biological Interpretation

INTRODUCTION

- The World Health Organization declared the outbreak of COVID-19 as a pandemic on 11th March, 2020 and till date, there has been approximately 29 crores of reported cases across the world.
- COVID-19 is caused by Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) which is an enveloped, single-stranded mRNA virus that belongs to a large family of viruses, coronaviridae.
- SARS-CoV-2 primarily infects the lower respiratory tract and lungs of human and is known to cause respiratory illness from mild to severe and sometimes even death.
- Seven human coronaviruses have been identified till date, viz. 229E (alpha coronavirus), NL63 (alpha coronavirus), OC43 (beta coronavirus), HKU1 (beta coronavirus), MERS-CoV (the beta coronavirus that causes Middle East Respiratory Syndrome, or MERS), SARS-CoV (the beta coronavirus that causes severe acute respiratory syndrome, or SARS), SARS-CoV-2 (the novel coronavirus that causes coronavirus disease 2019, or COVID-19)
- Out of these, MERS-CoV, SARS-CoV, and SARS-CoV-2 are way more pathogenic and are known to cause severe symptoms such as shortness of breath and eventually death.

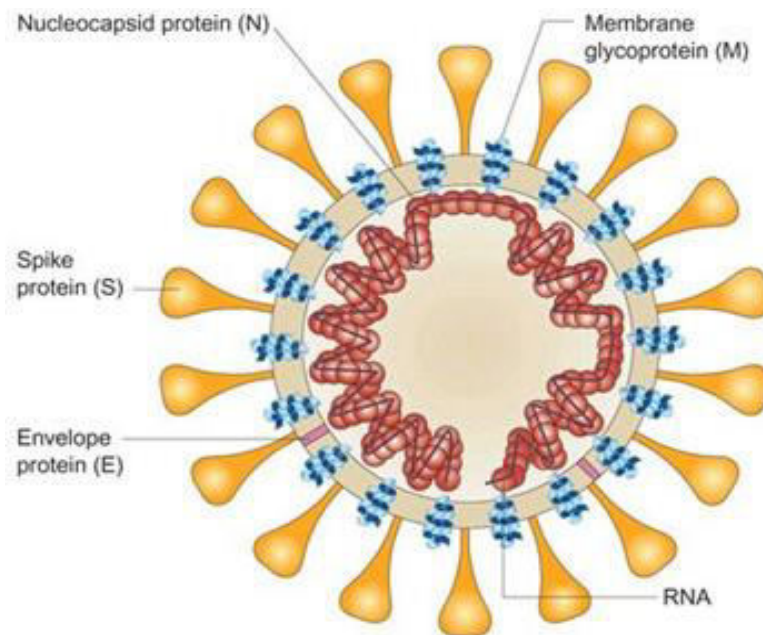


Figure 1. Coronavirus

- At the core of the coronavirus, the single-stranded RNA is present which is the genetic blueprint which enables the production of proteins.
- The nucleoproteins are attached to the RNA which aids in the structural formation and also helps the virus to replicate.
- The viral envelope which is made up of lipids protects the virus and it also anchors various structural proteins which are used by the virus during infection.
- The envelope protein is a membrane protein, which aids in the virion assembly and morphogenesis.
- The crown-like appearance in the coronavirus is due to the presence of the spike proteins. They allow the virus to penetrate into the host cell and cause infection.

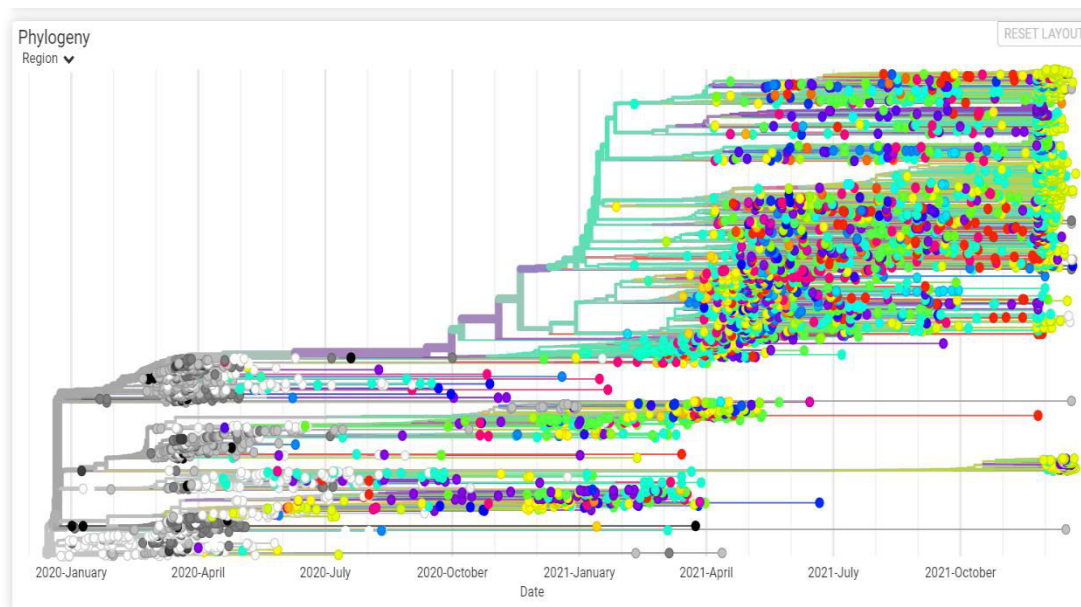


Figure 2. GISAID

GISAID (Global Initiative on Sharing All Influenza Data)

This diagram depicts the prevalence of various strains of coronavirus across India.

Description of the dataset

- COVID-19 Positive and Negative Plasma Samples
- Total samples: 71 (20 Negative & 51 Positive)
- 58.7 % missing values

Secondary Data Analysis using Metaboanalyst 5.0

Data processing information:

Checking data content ...passed.

Samples are in columns and features in rows.

The uploaded file is in comma separated values (.csv) format.

The uploaded data file contains 71 (samples) by 1206 (peaks(mz/rt)) data matrix.

Samples are not paired.

2 groups were detected in samples.

Only English letters, numbers, underscore, hyphen and forward slash (/) are allowed.

Other special characters or punctuations (if any) will be stripped off.

Non-numeric values were found and replaced by NA.

321 features with a constant or single value across samples were found and deleted.

A total of 36903 (58.7%) missing values were detected.

By default, missing values will be replaced by 1/5 of min positive values of their corresponding variables

Click the **Proceed** button if you accept the default practice;

Or click the **Missing Values** button to use other methods.

Missing value estimation:

Too many missing values will cause difficulties for downstream analysis. There are several different methods for this purpose. The default method replaces all the missing values with a small values (the half of the minimum positive values in the original data) assuming to be the detection limit. Click **next** if you want to use the default method. The assumption of this approach is that most missing values are caused by low abundance metabolites (i.e. below the detection limit).

MetaboAnalyst also offers other methods, such as replace by mean/median, k-nearest neighbours based on similar features - KNN (feature-wise), k-nearest neighbours based on similar samples - KNN (sample-wise), probabilistic PCA (PPCA), Bayesian PCA (BPCA) method, singular value decomposition (SVD) method to impute the missing values ([ref.](#)). Note for KNN, k is set to 10 (the default value). Please choose the one that is the most appropriate for your data.

Step 1. Remove features with too many missing values

☒ Remove features with > 50 % missing values

Step 2. Estimate the remaining missing values

☐ Replace by LoDs (1/5 of the minimum positive value of each variable)

☐ Exclude variables with missing values

☐ Replace by column (feature) mean

☒ Estimate missing values using KNN (feature-wise)

Missing value: proteins with > 50 % were removed

Estimation: k-Nearest Neighbor algorithm and Feature selection

Sample Normalization

- **Less than 250 variables:** 5% will be filtered;
- **Between 250 - 500 variables:** 10% will be filtered;
- **Between 500 - 1000 variables:** 25% will be filtered;
- **Over 1000 variables:** 40% will be filtered;

Please note, in order to reduce the computational burden to the server, the **None** option is only for less than 5000 features. The maximum allowed number of variables is 5000. For power analysis, the max number is **2500** to improve power and to control computing time. Over that, the IQR filter will still be applied to keep only top maximum features, even if you choose None option.

☐ Filtering features if their RSDs are > % in QC samples

☒ None (less than 5000 features)

☐ Interquantile range (IQR)

☐ Standard deviation (SD)

☐ Median absolute deviation (MAD)

☐ Relative standard deviation (RSD = SD/mean)

☐ Non-parametric relative standard deviation (MAD/median)

☐ Mean intensity value

☐ Median intensity value

Filtering was set to none as the data had less than 5000 features.

Sample normalization

☐ None

☐ Sample-specific normalization (i.e. weight, volume) [Specify](#)

☐ Normalization by sum

☒ Normalization by median

☐ Normalization by a reference sample (PQN) [Specify](#)

☐ Normalization by a pooled sample from group (group PQN) [Specify](#)

☐ Normalization by reference feature [Specify](#)

☐ Quantile normalization (suggested only for > 1000 features)

Data transformation

☐ None

☒ Log transformation (base 10)

☐ Square root transformation (square root of data values)

☐ Cube root transformation (cube root of data values)

Data scaling

☒ None

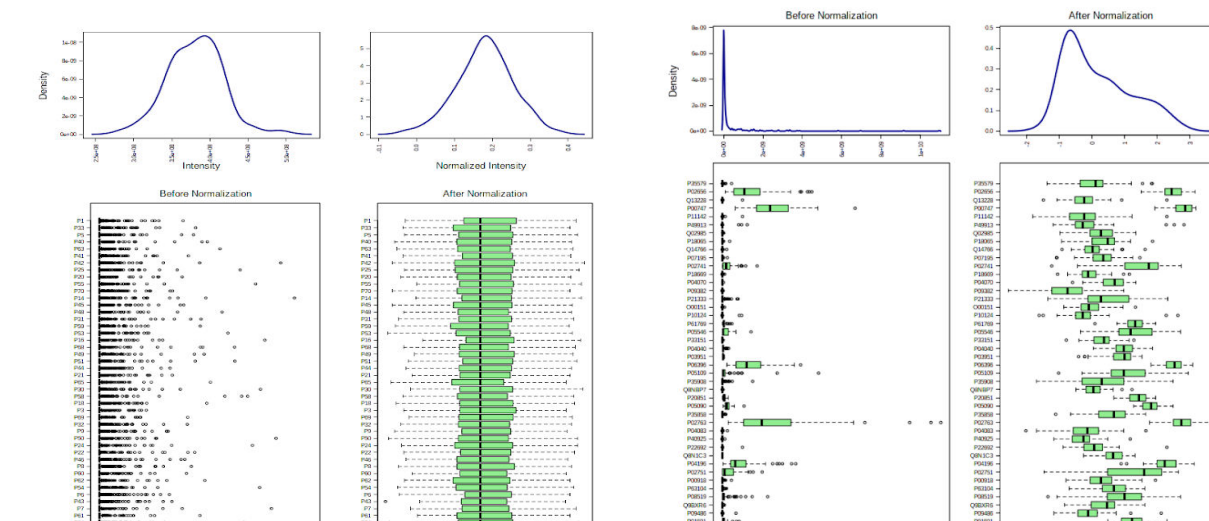
☐ Mean centering (mean-centered only)

☐ Auto scaling (mean-centered and divided by the standard deviation of each variable)

☐ Pareto scaling (mean-centered and divided by the square root of the standard deviation of each variable)

☐ Range scaling (mean-centered and divided by the range of each variable)

Based on the feature and sample view, the data was normalized by median, transformed by log transformation and data scaling was set to none.



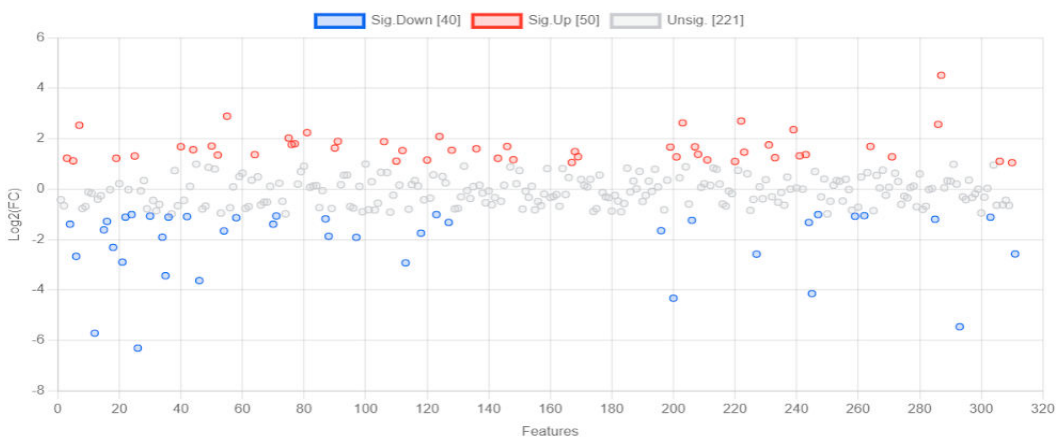
Sample View

Feature View

Statistical Analysis

1. Fold Change Analysis

Click a point to view; drag to zoom; reset zoom at bottom



FC analysis was performed to compare the absolute value of changes between two groups which basically depicts the change in protein expression.

Significant Upregulated proteins: 50

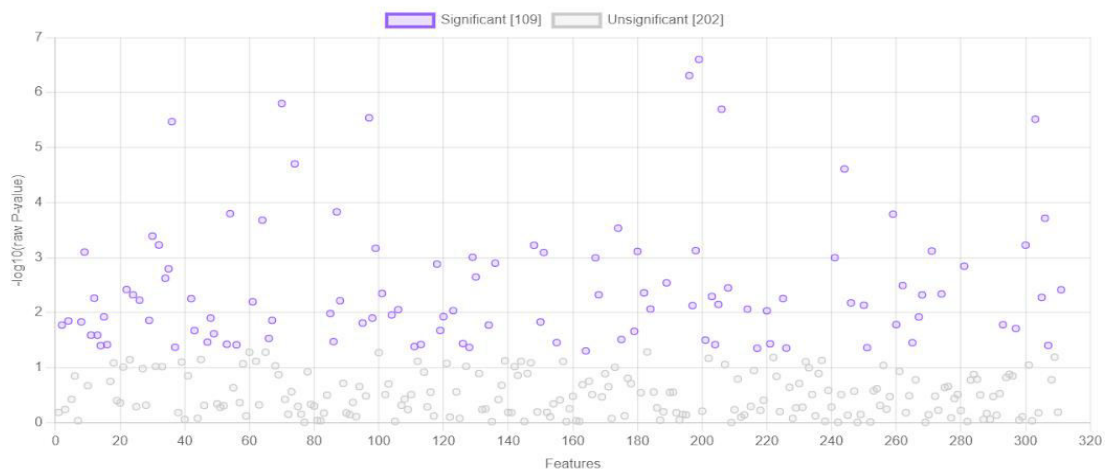
Significant Downregulated proteins: 40

Unsignificant proteins: 221

The positive half is the overexpression of proteins for positive samples and vice-versa.

2. T-test Analysis

Click a point to view; drag to zoom; reset zoom at bottom



T-test was performed to find the number of significant and insignificant proteins in the sample.

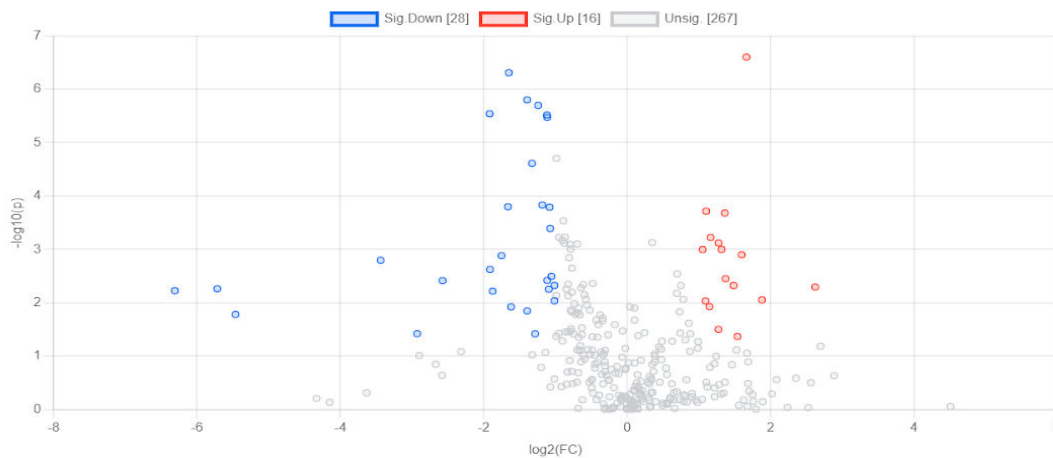
Significant proteins: 109

Unsignificant proteins: 202

The plot shows 109 significant proteins that pass the P-value threshold of 0.05 or 5 %.

3. Volcano Plot

Click a point to view; drag to zoom; reset zoom at bottom



Volcano Plot is a combination of Fold change and T-test.

Significantly upregulated: 16

Significantly downregulated: 28

Unsignificant: 267

The red dots imply the proteins that are upregulated in the positive samples while the blue dots denote the proteins upregulated in the negative samples.

Table 4: Important features identified by volcano plot

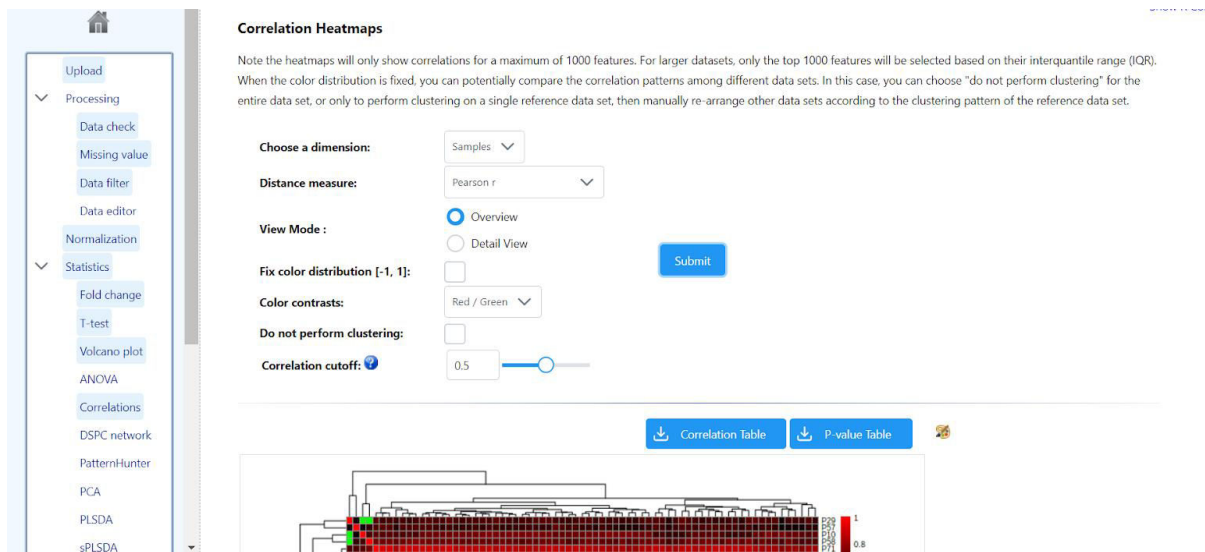
	Peaks(mz/rt)	FC	log2(FC)	raw.pval	-log10(p)
1	P04275	3.1699	1.6644	2.5028e-07	6.6016
2	P04196	0.31845	-1.6509	4.9184e-07	6.3082
3	P01344	0.38021	-1.3951	1.5834e-06	5.8004
4	P05452	0.42312	-1.2409	2.0116e-06	5.6965
5	P17936	0.26464	-1.9179	2.8862e-06	5.5397
6	Q16610	0.46085	-1.1176	3.0578e-06	5.5146
7	P13598	0.46205	-1.1139	3.3824e-06	5.4708
8	P14151	0.39877	-1.3264	2.454e-05	4.6101
9	Q13201	0.44047	-1.1829	0.00014782	3.8303
10	P08637	0.31594	-1.6623	0.00015935	3.7976
11	P25311	0.47248	-1.0817	0.00016258	3.7889
12	Q96IY4	2.1458	1.1015	0.00019292	3.7146
13	P22352	2.572	1.3629	0.00020917	3.6795
14	P03950	0.47513	-1.0736	0.00040689	3.3905
15	P01023	2.2414	1.1644	0.00059805	3.2233
16	P35908	2.4221	1.2762	0.00075947	3.1195
17	P13645	2.4917	1.3171	0.0010043	2.9981
18	P02675	2.0753	1.0533	0.0010107	2.9954
19	P00739	3.0255	1.5972	0.001264	2.8983
20	P05362	0.29694	-1.7518	0.0013132	2.8817
21	P49913	0.092249	-3.4383	0.0015997	2.796
22	Q8N6C8	0.26572	-1.912	0.0023764	2.6241
23	P27918	0.48116	-1.0554	0.0032215	2.4919
24	P05546	2.5908	1.3734	0.0035551	2.4491
25	P16070	0.46195	-1.1142	0.0038116	2.4189
26	Q9Y6R7	0.16784	-2.5749	0.003853	2.4142
27	P02679	2.8014	1.4862	0.0047541	2.3229
28	P33151	0.49503	-1.0144	0.0047546	2.3229
29	P05109	6.1635	2.6237	0.0050898	2.2933
30	Q86UD1	0.019011	-5.717	0.0054758	2.2615
31	Q9Y5Y7	0.46892	-1.0926	0.0055775	2.2536
32	Q02818	0.012614	-6.3088	0.005954	2.2252
33	P55072	0.27258	-1.8753	0.0060821	2.2159
34	P00918	3.6812	1.8802	0.0088487	2.0531
35	P01591	0.49516	-1.014	0.0092109	2.0357
36	P07737	2.1331	1.093	0.0092605	2.0334
37	P62937	2.2179	1.1492	0.011836	1.9268
38	P01714	0.32584	-1.6178	0.011947	1.9228
39	Q8WZ42	0.38031	-1.3948	0.014217	1.8472
40	Q02383	0.022674	-5.4628	0.016593	1.7801
41	P04406	2.4182	1.2739	0.031608	1.5002
42	P01619	0.1313	-2.9291	0.038048	1.4197
43	A0A075B6K4	0.41065	-1.284	0.038218	1.4177
44	O75636	2.9073	1.5397	0.042937	1.3672

44 significant proteins were found from volcano plot.

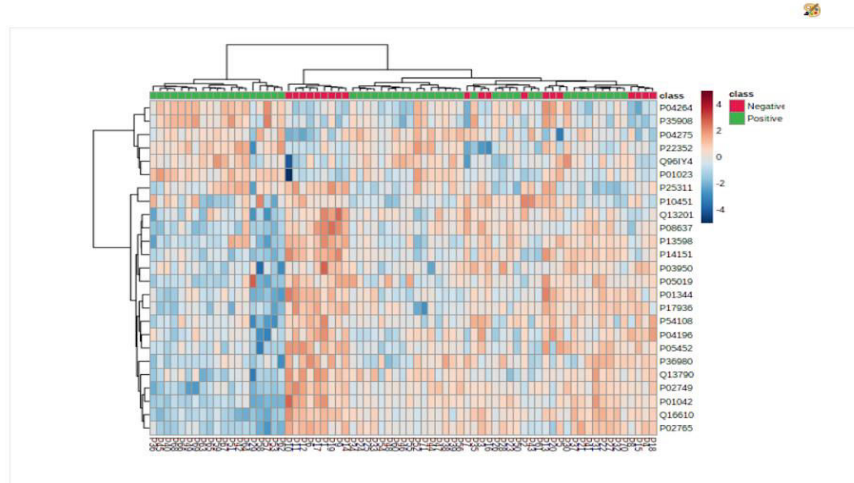
I have sorted out the upregulated and downregulated proteins in the volcano plot by filtering it from largest to smallest.

The upregulated proteins were taken for further pathway analysis using Reactome, Metascape and PANTHER.

4. Correlation heatmaps

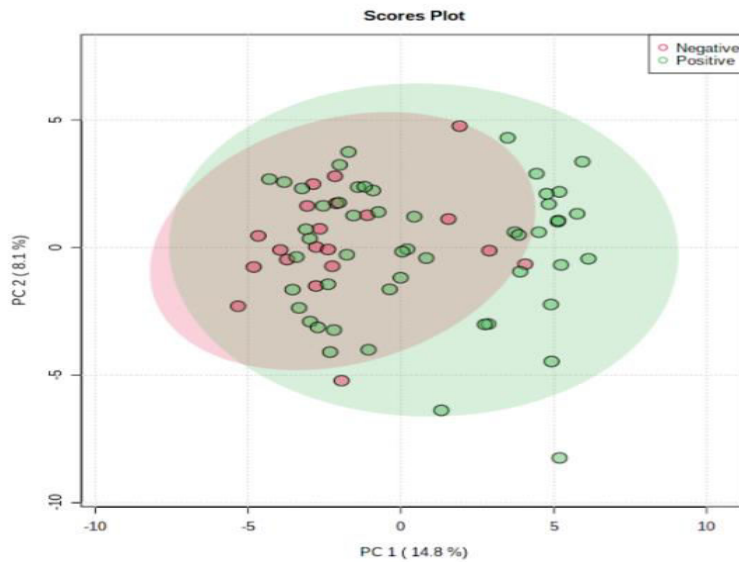


The correlation heatmaps is basically how one protein is related to the other. Here we are using Pearson r correlation which basically denotes if one protein is increasing how it affects the other.



In the hierarchical clustering heatmaps, the red ones denote the upregulated proteins while the blue ones represent downregulated proteins.

5. PCA plot



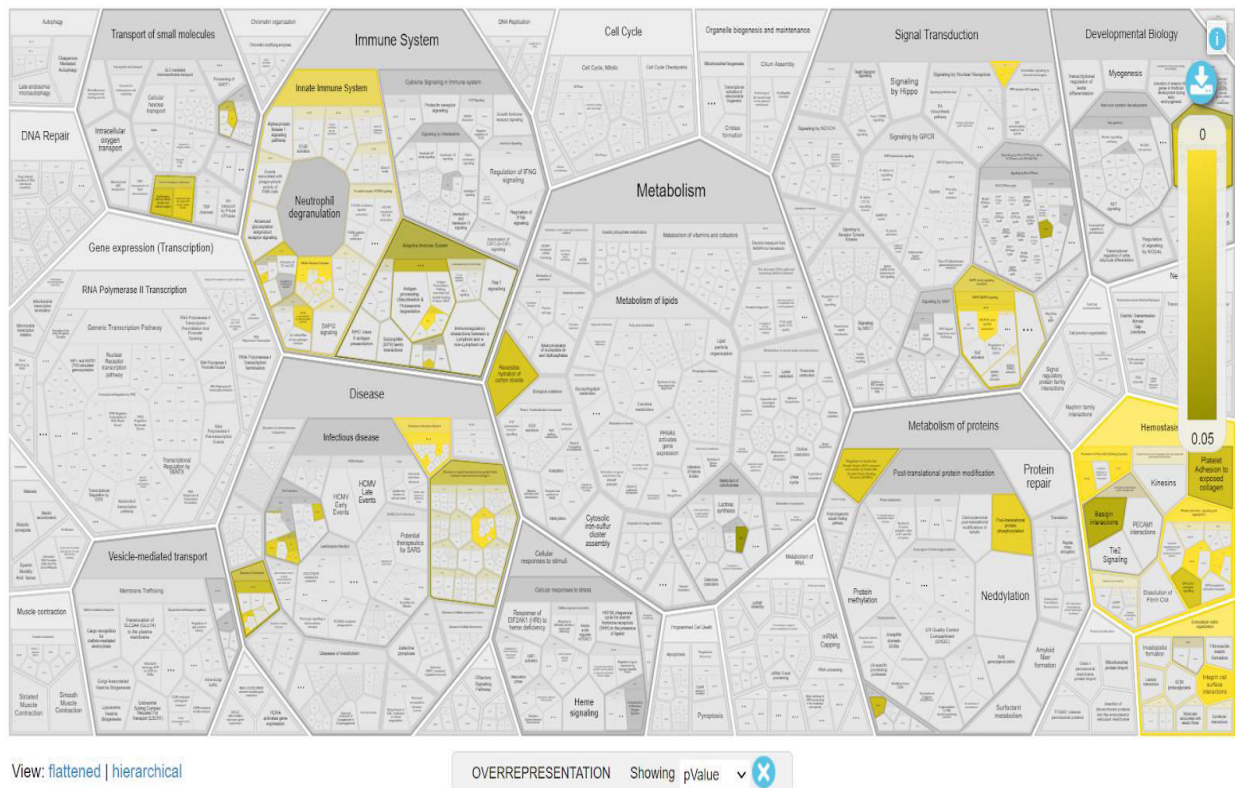
Proteins in PCA are overlapping because PCA is an unsupervised model, so no grouping is there to differentiate if the samples are severe or non-severe.

Here, we have two clusters of one positive and one cluster of both positive and negative. Possible reasons for overlapping would be that the data that has been shared might be bad due to any sort of reasons such as, even the negative samples might have some positive proteins, i.e, the people who have tested negative might not have been recovered properly.

Pathway Analysis

The 44 significant proteins that was found from volcano plot was sorted out as upregulated and downregulated proteins by filtering it from largest to smallest. A total of 16 upregulated proteins were found. 16 of the submitted entities were found, mapping to the 16 Reactome entities.

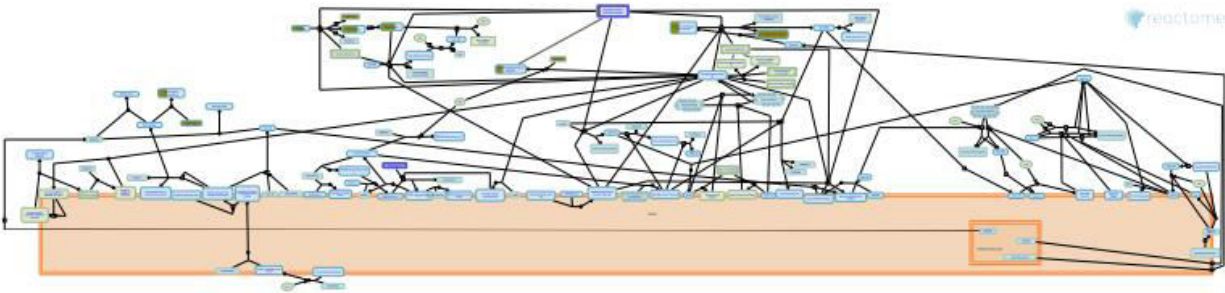
1. Reactome



Pathway name	Entities				Reactions	
	found	ratio	p-value	FDR*	found	ratio
Formation of Fibrin Clot (Clotting Cascade)	5 / 39	0.003	2.13e-09	2.81e-07	10 / 61	0.004
Platelet degranulation	6 / 128	0.011	1.59e-08	8.75e-07	2 / 11	8.10e-04
Response to elevated platelet cytosolic Ca2+	6 / 133	0.012	1.99e-08	8.75e-07	2 / 14	0.001
Platelet activation, signaling and aggregation	6 / 265	0.024	1.12e-06	2.10e-05	23 / 116	0.009
p130Cas linkage to MAPK signaling for integrins	3 / 15	0.001	1.31e-06	2.10e-05	3 / 3	2.21e-04
GRB2:SOS provides linkage to MAPK signaling for Integrins	3 / 15	0.001	1.31e-06	2.10e-05	2 / 2	1.47e-04
MyD88 deficiency (TLR2/4)	3 / 19	0.002	2.66e-06	3.73e-05	2 / 2	1.47e-04
IRAK4 deficiency (TLR2/4)	3 / 20	0.002	3.10e-06	4.03e-05	2 / 2	1.47e-04
Regulation of TLR by endogenous ligand	3 / 21	0.002	3.59e-06	4.31e-05	2 / 12	8.84e-04
Common Pathway of Fibrin Clot Formation	3 / 22	0.002	4.12e-06	4.54e-05	6 / 29	0.002
Intrinsic Pathway of Fibrin Clot Formation	3 / 23	0.002	4.71e-06	4.71e-05	4 / 24	0.002
Integrin signaling	3 / 28	0.002	8.46e-06	6.76e-05	15 / 24	0.002
Diseases associated with the TLR signaling cascade	3 / 34	0.003	1.51e-05	1.05e-04	4 / 15	0.001
Diseases of Immune System	3 / 34	0.003	1.51e-05	1.05e-04	4 / 15	0.001
Signaling by high-kinase activity BRAF mutants	3 / 37	0.003	1.94e-05	1.36e-04	4 / 6	4.42e-04
Platelet Aggregation (Plug Formation)	3 / 40	0.004	2.44e-05	1.46e-04	16 / 27	0.002
MAP2K and MAPK activation	3 / 41	0.004	2.63e-05	1.58e-04	8 / 12	8.84e-04
Signaling by RAF1 mutants	3 / 42	0.004	2.82e-05	1.69e-04	4 / 7	5.16e-04
Hemostasis	7 / 726	0.065	3.19e-05	1.91e-04	36 / 334	0.025
Signaling downstream of RAS mutants	3 / 47	0.004	3.93e-05	1.97e-04	4 / 7	5.16e-04
Signaling by moderate kinase activity BRAF mutants	3 / 47	0.004	3.93e-05	1.97e-04	4 / 7	5.16e-04
Paradoxical activation of RAF signaling by kinase inactive BRAF	3 / 47	0.004	3.93e-05	1.97e-04	4 / 7	5.16e-04
Signaling by RAS mutants	3 / 47	0.004	3.93e-05	1.97e-04	4 / 9	6.63e-04

Table. Top 25 pathways from Reactome

1. **Formation of Fibrin Clot (Clotting Cascade) (R-HSA-140877)**

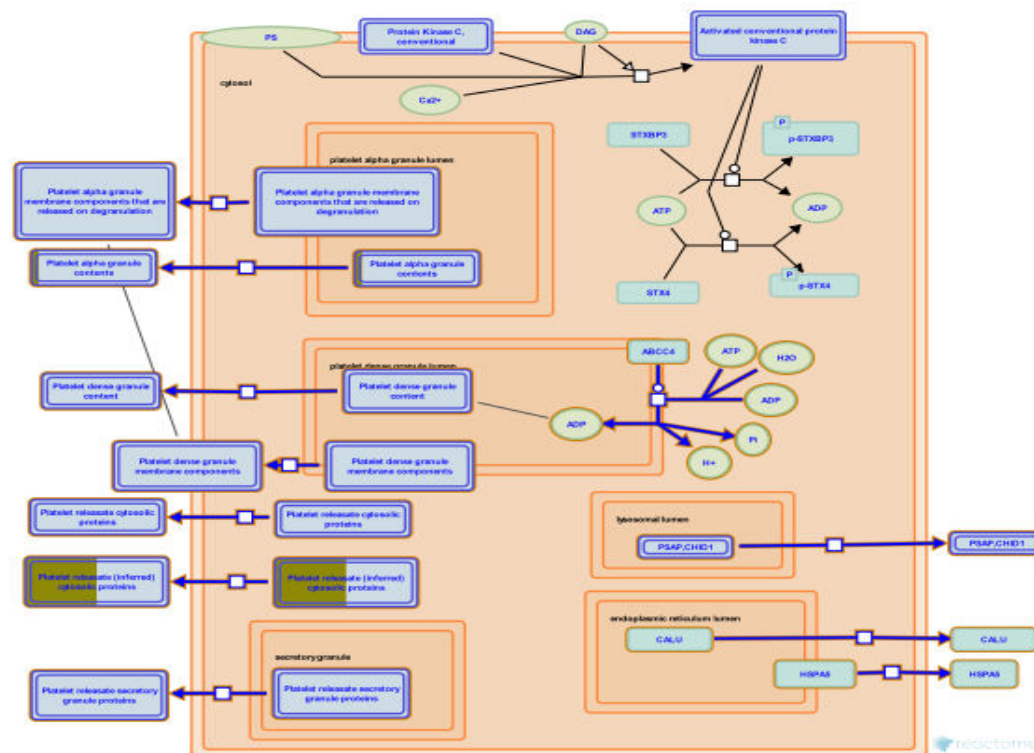


5 submitted entities found in this pathway, mapping to 5 Reactome entities

Input	UniProt Id	Input	UniProt Id	Input	UniProt Id
P01023	P01023	P02675	P02675	P02679	P02679
P04275	P04275	P05546	P05546		

The formation of mysterious blood clots in the various tissues and organs of COVID-19 patients is still a mystery. The high ACE2 expression in the endothelium of blood vessels facilitates the high-affinity binding of SARS-CoV-2 using spike protein, causing infection and internal injury inside the vascular wall of blood vessels. This viral associated injury may directly/indirectly initiate activation of coagulation and clotting cascades forming internal blood clots.

2. Platelet degranulation (R-HSA-114608)

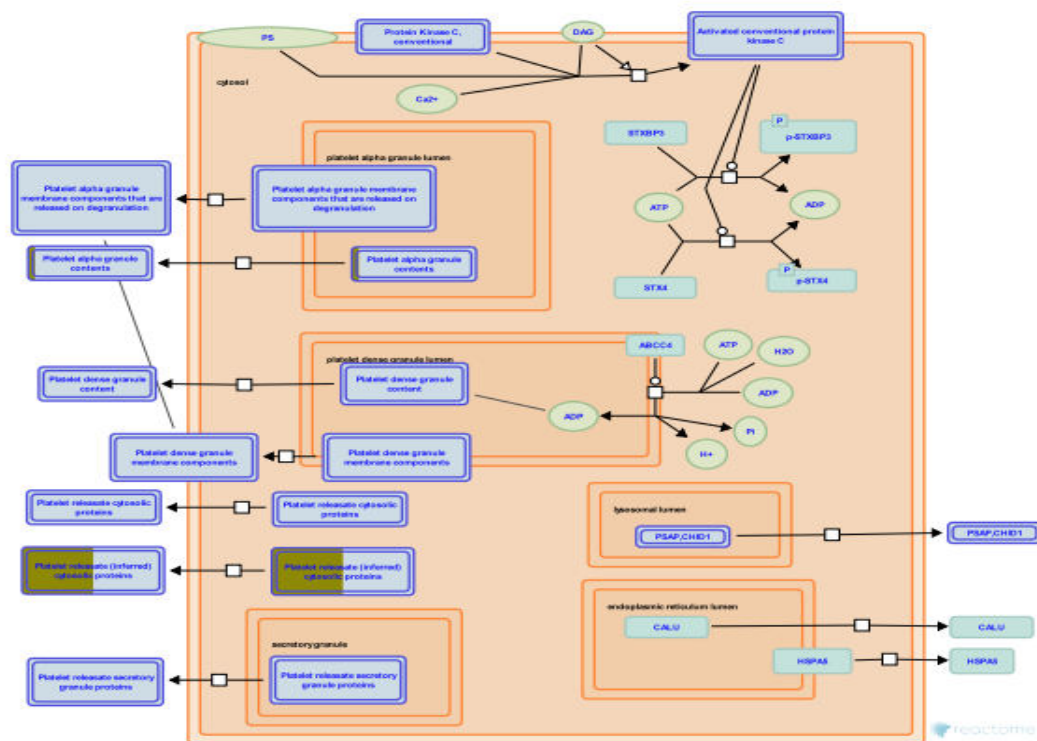


6 submitted entities found in this pathway, mapping to 6 Reactome entities

Input	UniProt Id	Input	UniProt Id	Input	UniProt Id
P01023	P01023	P02675	P02675	P02679	P02679
P04275	P04275	P07737	P07737	P62937	P62937

Platelets are at the frontline of COVID-19 pathogenesis, as they release various sets of molecules through the different stages of the disease. Platelets may thus have the potential to contribute to the overwhelming thrombo-inflammation in COVID-19, and the inhibition of pathways related to platelet activation may improve the outcomes during COVID-19.

3. Response to elevated platelet cytosolic Ca²⁺ (R-HSA-76005)



6 submitted entities found in this pathway, mapping to 6 Reactome entities

Input	UniProt Id	Input	UniProt Id	Input	UniProt Id
P01023	P01023	P02675	P02675	P02679	P02679
P04275	P04275	P07737	P07737	P62937	P62937

Activation of phospholipase C enzymes results in the generation of second messengers of the phosphatidylinositol pathway. The events resulting from this pathway are a rise in intracellular calcium and activation of Protein Kinase C.

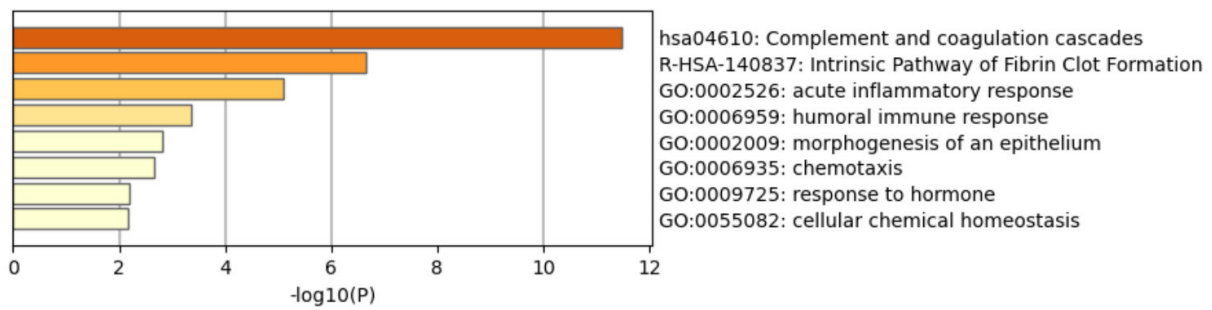
2. Metascape

Metascape Gene List Analysis Report

metascape.org¹

Bar Graph Summary

Figure 1. Bar graph of enriched terms across input gene lists, colored by p-values.



I have found almost similar results as Reactome.

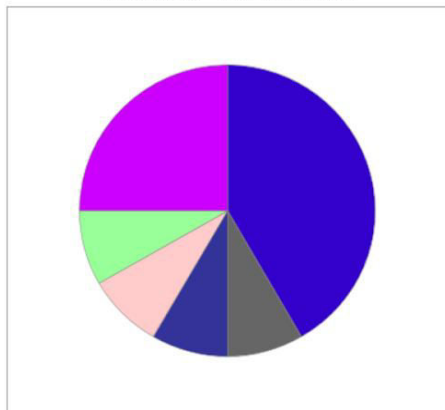
3. Panther

	Homo sapiens (REF)	Client Text Box Input					
PANTHER Pathways	#	▼ #	expected	Fold Enrichment	+/-	raw P value	FDR
Unclassified	17977	8	13.97	.57	-	3.41E-04	1.90E-02
Blood coagulation	48	5	.04	> 100	+	3.96E-10	6.61E-08
Plasminogen activating cascade	21	3	.02	> 100	+	7.69E-07	6.42E-05
Inflammation mediated by chemokine and cytokine signaling pathway	255	1	.20	5.05	+	1.81E-01	1.00E00
Huntington disease	148	1	.11	8.70	+	1.10E-01	1.00E00
Glycolysis	20	1	.02	64.36	+	1.62E-02	6.76E-01
Cytoskeletal regulation by Rho GTPase	82	1	.06	15.70	+	6.25E-02	1.00E00
Alzheimer disease-amyloid secretase pathway	67	0	.05	< 0.01	-	1.00E00	1.00E00
Alpha adrenergic receptor signaling pathway	25	0	.02	< 0.01	-	1.00E00	1.00E00
Adrenaline and noradrenaline biosynthesis	30	0	.02	< 0.01	-	1.00E00	1.00E00
Nicotine pharmacodynamics pathway	35	0	.03	< 0.01	-	1.00E00	1.00E00
Toll pathway-drosophila	2	0	.00	< 0.01	-	1.00E00	1.00E00
SCW signaling pathway	3	0	.00	< 0.01	-	1.00E00	1.00E00
MYO signaling pathway	1	0	.00	< 0.01	-	1.00E00	1.00E00
GBB signaling pathway	1	0	.00	< 0.01	-	1.00E00	1.00E00
DPP signaling pathway	3	0	.00	< 0.01	-	1.00E00	1.00E00
DPP-SCW signaling pathway	3	0	.00	< 0.01	-	1.00E00	1.00E00
BMP/activin signaling pathway-drosophila	4	0	.00	< 0.01	-	1.00E00	1.00E00
Xanthine and guanine salvage pathway	4	0	.00	< 0.01	-	1.00E00	1.00E00
Activin beta signaling pathway	2	0	.00	< 0.01	-	1.00E00	1.00E00
Vitamin B6 metabolism	3	0	.00	< 0.01	-	1.00E00	1.00E00

Select Ontology: View:

PANTHER Pathway

Total # Genes: 16 Total # pathway hits: 12



Click to get gene list for a category:

- [Blood coagulation \(P00011\)](#)
- [Cytoskeletal regulation by Rho GTPase \(P00016\)](#)
- [Glycolysis \(P00024\)](#)
- [Huntington disease \(P00029\)](#)
- [Inflammation mediated by chemokine and cytokine signaling pathway \(P00031\)](#)
- [Plasminogen activating cascade \(P00050\)](#)

Color picker powered by Web Colors by VisBone

**Chart tooltips are read as: Category name (Accession): # genes; Percent of gene hit against total # genes; Percent of gene hit against total # Pathway hits

Conclusion

- The dataset that was provided contained overlapped Covid 19 plasma samples containing 51 positive and 20 negative samples.
- Differential host response was seen during COVID-19 infections.
- Some common pathways from all the three pathway analysis tools which include fibrin clot formation, platelet degranulation, complement and coagulation cascades, blood coagulation, etc. showed significant alteration.
- Protein S100-A8, Carbonic anhydrase 2, von Willebrand factor, etc. are found to be differentially expressed proteins.
- The output obtained have the potential to facilitate possible therapeutic development and treatment. Thus, proteomics has indeed aided in monitoring

COVID-19 pathology and hence, it has become an emerging and promising approach for the study and treatment of COVID-19.

Bibliography

- Fard, M.B., Fard, S.B., Ramazi, S. *et al.* Thrombosis in COVID-19 infection: Role of platelet activation-mediated immunity. *Thrombosis J* **19**, 59 (2021). <https://doi.org/10.1186/s12959-021-00311-9>
- Zaid, Y., Puhm, F., Allaey, I., Naya, A., Oudghiri, M., Khalki, L., ... Boilard, E. (2020). *Platelets Can Associate With SARS-CoV-2 RNA and Are Hyperactivated in COVID-19.* *Circulation Research*, *127*(11), 1404–1418. doi:10.1161/circresaha.120.317703
- Biswas S, Thakur V, Kaur P, Khan A, Kulshrestha S, Kumar P. Blood clots in COVID-19 patients: Simplifying the curious mystery. *Med Hypotheses*. 2021 Jan;146:110371. doi: 10.1016/j.mehy.2020.110371. Epub 2020 Nov 6. PMID: 33223324; PMCID: PMC7644431.
- <https://www.wikipathways.org/index.php/Pathway:WP1903>
- <https://www.metaboanalyst.ca/MetaboAnalyst/home.xhtml>
- <https://reactome.org/>
- <https://metascape.org/gp/index.html#/main/step1>
- <http://www.pantherdb.org/>

Acknowledgement

I would like to wholeheartedly express my sincere gratitude to the entire organizing team of Proteomics Advanced Winter School (PAWS-2021) for giving me an opportunity to carry out this work. I am extremely grateful to convenor cum professor, Dr. Sanjeeva Srivastava for encouraging and guiding me throughout this training and workshop. Also, my thanks are due to my guides, mentors, peers, well-wishers, etc.

PLAGIARISM SCAN REPORT

Words 797 Date January 07, 2022

Characters 5823 Excluded URL

8%

Plagiarism

92%

Unique

3

Plagiarized
Sentences

34

Unique Sentences

Content Checked For Plagiarism

To analyze the given dataset of plasma samples of both COVID + and COVID – samples and determine the changes in the proteome to find proteins (if present) and their relation to biochemical pathways which corresponds to host response towards COVID pathogenesis.

INDEX

- Introduction
- GISAID Database
- Secondary Data Analysis using Metaboanalyst 5.0
- Pathway Analysis using
 - (i) Reactome
 - (ii) Metascape
 - (iii) Panther
- Biological Interpretation

INTRODUCTION

- The World Health Organization declared the outbreak of COVID-19 as a pandemic on 11th March, 2020 and till date, there has been approximately 29 crores of reported cases across the world.
- COVID-19 is caused by Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) which is an enveloped, single-stranded mRNA virus that belongs to a large family of viruses, coronaviridae.
- SARS-CoV-2 primarily infects the lower respiratory tract and lungs of human and is known to cause respiratory illness from mild to severe and sometimes even death.
- Seven human coronaviruses have been identified till date, viz. 229E (alpha coronavirus), NL63 (alpha coronavirus), OC43 (beta coronavirus), HKU1 (beta coronavirus), MERS-CoV (the beta coronavirus that causes Middle East Respiratory Syndrome, or MERS), SARS-CoV (the beta coronavirus that causes severe acute respiratory syndrome, or SARS), SARS-CoV-2 (the novel coronavirus that causes coronavirus disease 2019, or COVID-19)
- Out of these, MERS-CoV, SARS-CoV, and SARS-CoV-2 are way more pathogenic and are known to cause severe symptoms such as shortness of breath and eventually death.

Figure 1. Coronavirus

- At the core of the coronavirus, the single-stranded RNA is present which is the genetic blueprint which enables the production of proteins.
- The nucleoproteins are attached to the RNA which aids in the structural formation and also helps the virus to replicate.
- The viral envelope which is made up of lipids protects the virus and it also anchors various structural proteins which are used by the virus during infection.
- The envelope protein is a membrane protein, which aids in the virion assembly and morphogenesis.
- The crown-like appearance in the coronavirus is due to the presence of the spike proteins. They allow the virus to penetrate into the host cell and cause infection.

Figure 2. GISAID

GISAID (Global Initiative on Sharing All Influenza Data)

This diagram depicts the prevalence of various strains of coronavirus across India.

Description of the dataset

- COVID-19 Positive and Negative Plasma Samples
- Total samples: 71 (20 Negative & 51 Positive)
- 58.7 % missing values

Secondary Data Analysis using Metaboanalyst 5.0

Missing value: proteins with > 50 % were removed

Estimation: k-Nearest Neighbor algorithm and Feature selection

Sample Normalization

Filtering was set to none as the data had less than 5000 features.

Based on the feature and sample view, the data was normalized by median, transformed by log transformation and data scaling was set to none.

Statistical Analysis

1. Fold Change Analysis

FC analysis was performed to compare the absolute value of changes between two groups which basically depicts the change in protein expression.

Significant Upregulated proteins: 50

Significant Downregulated proteins: 40

Unsignificant proteins: 221

The positive half is the overexpression of proteins for positive samples and vice-versa.

2. T-test Analysis

T-test was performed to find the number of significant and insignificant proteins in the sample.

Significant proteins: 109

Unsignificant proteins: 202

The plot shows 109 significant proteins that pass the P-value threshold of 0.05 or 5 %.

3. Volcano Plot

Volcano Plot is a combination of Fold change and T-test.

Significantly upregulated: 16

Significantly downregulated: 28

Unsignificant: 267

The red dots imply the proteins that are upregulated in the positive samples while the blue dots denote the proteins upregulated in the negative samples.

44 significant proteins were found from volcano plot.

I have sorted out the upregulated and downregulated proteins in the volcano plot by filtering it from largest to smallest.

The upregulated proteins were taken for further pathway analysis using Reactome, Metascape and PANTHER.

4. Correlation heatmaps

The correlation heatmaps is basically how one protein is related to the other. Here we are using Pearson r correlation which basically denotes if one protein is increasing how it affects the other.

In the hierarchical clustering heatmaps, the red ones denote the upregulated proteins while the blue ones represent downregulated proteins.

5. PCA plot

Proteins in PCA are overlapping because PCA is an unsupervised model, so no grouping is there to differentiate if the samples are severe or non-severe.

Here, we have two clusters of one positive and one cluster of both positive and negative. Possible reasons for overlapping would be that the data that has been shared might be bad due to any sort of reasons such as, even the negative samples might have some positive proteins, i.e, the people who have tested negative might not have been recovered properly.

Sources	Similarity
www.the-sun.com › news › 2930758Coronavirus passed from dogs put kids in HOSPITAL, study finds May 21, 2021 · 229E (alpha coronavirus) NL63 (alpha coronavirus) OC43 (beta coronavirus) HKU1 (beta coronavirus) MERS-CoV (the beta coronavirus that causes Middle East Respiratory Syndrome, or MERS) SARS-CoV (the beta coronavirus that causes severe acute respiratory syndrome, or SARS) SARS-CoV-2 (which causes Covid-19) Possibly CCoV-HuPn-2018	13%

https://www.the-sun.com/news/2930758/coronavirus-passed-dogs-kids-hospital-study-finds//	
<p data-bbox="121 147 1209 203">B'more City Health on Twitter: "Other human coronaviruses - Twitter</p> <p data-bbox="121 226 1217 297">Another great question. Coronaviruses are not new, but COVID-19 is, at least for humans. Human coronaviruses were first identified in the mid-1960s. Another great question. Coronaviruses are not new, but COVID-19 is, at least for humans. Human coronaviruses were first identified in the mid-1960s.</p> <p data-bbox="121 315 834 342">https://mobile.twitter.com/bmore_healthy/status/1425254281276444674</p>	<p data-bbox="1361 226 1417 259">13%</p>
<p data-bbox="121 394 1043 421">pubmed.ncbi.nlm.nih.gov › 28382917 GISAID: Global initiative on sharing all influenza data ...</p> <p data-bbox="121 439 818 465">GISAID: Global initiative on sharing all influenza data - from vision to reality</p> <p data-bbox="121 483 558 510">https://pubmed.ncbi.nlm.nih.gov/28382917//</p>	<p data-bbox="1361 432 1417 465">5%</p>