

# Liver Cancer Key Genes Identification

Ashitha Ebrahim, Joby George

**Abstract**—Liver disease is perhaps the deadliest malignant growth on the planet. In momentum contemplation, the capabilities for being chosen as key qualities in illnesses is bit low, constraining the precision of the anticipated key qualities in infections. To distinguish the key qualities of liver malignant growth with high exactness, and coordinated different microarray quality articulation datasets identified with the liver disease utilized. At that point recognize their basic DEGs (Differentially Expressed Genes) which will bring about more exact than those from the individual dataset. The datasets are on the whole human microarray quality articulation information recovered from the GEO (Gene Expression Omnibus) database and need to discover differentially communicated qualities among wellbeing and liver malignancy conditions. In light of these qualities, a protein-protein association system can be built and dissected to recognize the qualities tests that are having a higher impact on the system. These quality examples are prepared by utilizing a neural system LSTM. From this prepared neural system, the key hubs can be recognized and they can be considered as the key qualities of liver malignant growth. In addition, the strategy can be applied to different sorts of informational collections to choose key qualities of other complex ailments.

**Keywords** : DEGs, GEO, GO, KEGG, LSTM

## I. INTRODUCTION

It is basic to find fruitful medications for liver dangerous development. Liver harmful development, hepatocellular carcinoma (HCC) explicitly, is maybe the deadliest ailment around the globe, and the pace of HCC is extending rapidly in the United States and other made countries [1]. In 2012, Liver harm is the fifth commonest infection generally [2], and its overall ailment inconvenience incited a liberal number of significant lots of life lost [3]. It has a poor perseverance rate given its unimaginably compelling nature [4]. In this manner, liver infection is up 'til now an overriding general clinical issue and gravely needs stunning medicines. To fix this dangerous development, it is basic to know its unmistakable instrument and pathways. In the part, a couple of characteristics may expect fundamental occupations can be portrayed as key characteristics of a disease. Recognizing these key characteristics can add to uncovering the arrangement of an ailment. These key characteristics have an ability of filling in as focal points of treatment against this sickness. Extending liver threatening development related educational files give a significant advantage for locate the key characteristics of this sickness. To find the particular

Instrument or pathways of liver threatening development, a goliath number of studies have been made, provoking enormous improvement of related datasets. For example, when the watchword "liver danger" was used to search for the educational lists of verbalization by group from GEO, 24788 human datasets and 2015 mouse datasets were returned, independently at the hour of creating. Among these educational lists, microarray quality explanation data is regularly used to think about the effects of explicit meds, sicknesses on quality verbalization. These enormous amounts of natural datasets have incited the progression of various computational techniques, for instance, AI. For example, some AI based instruments have been made regular progression assessment including DNA, RNA, and protein sequences [5]. These instruments are natural and can deliver features described by customers for downstream assessment, for instance, portrayal of genes [6]. Frameworks organization is a systematically and practical approach for mining the enormous regular data. Given its ability, it has been used in various examinations to locate the key characteristics of diseases reliant on at any rate one sorts of high dimensional natural data [7]. Long et al. separated typical differentially imparted characteristics between protein-protein affiliation and transcript factors-target orchestrates as focus point characteristics in coronary course disease [8]. A couple of researchers developed a protein protein orchestrate and picked proteins with high degrees as the potential biomarkers of infirmities subject to microarray data. We furthermore have made many related works including foreseeing focus characteristics of harmful developments by taking a gander at two differential co-explanation organizes under two unmistakable conditions (prosperity and disease) and perceiving key characteristics of schizophrenia through connection of least navigating trees removed from two various quality frameworks in two one of a kind states. In recurring pattern considers, the capacities for being picked as key characteristics in diseases is to some degree low, compelling the precision of the foreseen key characteristics in ailments. Likewise, the gauge of the key characteristics in diseases not withstanding everything ought to be improved

Cytoscape is an open source bioinformatics programming stage for envisioning sub-atomic cooperation systems and coordinating with quality articulation profiles and other state information. Extra highlights are accessible as modules. Modules are accessible for organize and atomic profiling investigations, new designs, extra record position backing and association with databases and looking in enormous systems. Modules might be created utilizing the Cytoscape open Java programming engineering by anybody and module network improvement is empowered.

Manuscript received on January 27, 2021.

Revised Manuscript received on February 02, 2021.

Manuscript published on February 28, 2021.

Ashitha Ebrahim, Department of Computer Science and Engineering, Mar Athanasius College of Engineering, Kothamangalam, India. Email: achuashitha5@gmail.com

Prof. Joby George, Department of Computer Science and Engineering, Mar Athanasius College of Engineering, Kothamangalam, India. Email: jobygeo@hotmail.com

## Liver Cancer Key Genes Identification

Cytoscape additionally has a Java Script-driven sister venture named Cytoscape.js that can be utilized to dissect and imagine diagrams in JavaScript conditions, similar to a browser. In advanced circuits and AI, one-hot is a gathering of bits among which the legitimate blends of qualities are just those with a solitary high (1) piece and all the others low (0). A comparative usage in which all bits are '1' aside from one '0' is now and again called one cold. In measurements, sham factors speak to a comparable strategy for speaking to downright information. One-hot encoding is mostly utilized for manifesting the ambiance of a temper gadget. While utilizing an equal or Gray cipher, a decrypt is a whirl on to harvest the context. These machine, needn't play with an etymologist as the temper machine is in the nth mind-set if and just if the nth piece is elevated. A ring counter with 15 brightly referenced states become an event of a state machine.

Long flitting memory (LSTM) is a fake irregular neural framework (RNN) plan used in the field of significant learning. As opposed to standard feed-forward neural frameworks, LSTM has analysis affiliations. It cannot simply technique single data centres, (for instance, pictures), yet also entire progressions of data, (for instance, talk or video). For example, LSTM is material to assignments, for instance, unsegment, related handwriting affirmation, talk affirmation, and peculiarity acknowledgment in orchestrating traffic or IDS's (interference disclosure structures). A regular LSTM unit is made out of a cell, a data portal, a yield entryway, and a disregard entryway. The cell recalls esteems over abstract time breaks and the three doors control the movement of information into and out of the phone. LSTM frameworks are proper for describing, getting ready, and making gauges subject to the time course of action data, since there can be slacks of the darkening term between huge events in a period plan. LSTMs were made to deal with the exploding and dissipating edge gives that can be experienced while getting ready ordinary RNNs. The relative absence of care toward opening length is an ideal situation of LSTM over RNNs, covered Markov models, and other gathering learning strategies in different applications. In this paper, to improve the desire precision of key characteristics of liver threat, we composed different microarray quality explanation datasets containing tests under the conventional condition and liver illness condition to build up a quality framework from various resources. Considering the framework, we perceived the characteristics with a high degree, and high balanced betweenness centrality, and these quality models are set up under the neural framework LSTM. It can perceive the key centres and that can be considered as the key characteristic of Liver Cancer.

## II. RELATED WORKS

### A. GLOBOCAN

Appraisals of the general event and mortality from 27 noteworthy illnesses and for all tumours joined for 2012 are at present open in the GLOBOCAN game plan of the International Office for Research on Cancer. We study the sources and systems used in aggregating the national danger recurrence and mortality checks, and rapidly depict the key results by dangerous development site and in 20 colossal "domains" of the world. All around, there were 14.1 million new cases and 8.2 million passing's in 2012. The most typically dissected dangerous developments were

lung (1.82 million), chest (1.67 million), and colorectal (1.36 million); the most broadly perceived reasons for harmful development destruction were lung illness (1.6 million passing's), liver sickness (745,000 passing's), and stomach threatening development (723,000 passing's).

### B. HCC Incidence in United States

Hepatocellular carcinoma (HCC) is the third driving purpose behind danger mortality around the globe. This danger happens more every now and again among men than women, with the most raised rate rates uncovered in East Asia. The event paces of HCC in the United States have really been lower than in various countries. Regardless, in progressing decades, HCC age-adjusted rate rates have duplicated and fundamental liver harmful development demise rates have extended speedier than death rates for some other driving explanation behind ailment. About 90 percent of fundamental liver infections in the United States are HCCs, while most of the remaining 10 percent are intrahepatic cholangiocarcinomas. The pathway provoking HCC, generally, begins with an exceptional hepatic insult which progresses throughout the decades. Fibrosis and cirrhosis are basic forerunners of HCC. Among patients with limited stage HCC, treatment options may consolidate.

### C. Pse-in-One

With the heavy slide of natural progressions made in the post-genomic age, one of the most testing issues in computational science is the way to enough arrangement the course of action of a character model, (for instance, DNA, RNA or protein) with a discrete model or a vector that can suitably reflect its gathering plan information or catch its key features concerned. But a couple of web servers and stay lone mechanical assemblies were made to address this issue, all of these instruments, in any case, can simply manage one kind of test. Additionally, the amount of their intrinsic properties is obliged, and consequently, it is consistently difficult for customers to calculate the normal progressions as showed by their optimal features or properties. With a much greater number of inalienable properties, we are to propose a significantly increasingly versatile webservice called Pse-in-One which can, through its 28 one of a kind modes, produce practically all the possible feature vectors for DNA, RNA and protein courses of action. Particularly, it can in like manner produce those segment vectors with the properties described by customers themselves. These part vectors can be viably gotten together with AI figuring's to make computational markers and examination systems for various endeavours in bioinformatics and structure science.

### D. BioSeq –Analysis

With the heavy slide of common groupings delivered in the post-genomic age, one of the most testing issues is the way to computationally analyse their structures and limits. Man-made intelligence systems are expecting key occupations at the present time. Customarily, markers subject to AI frameworks contain three central advances: feature extraction, pointer improvement and execution appraisal.

Though a couple of Web servers and stay singular instruments have been made to empower the natural course of action assessment, they simply focus on individual advance. At the present time, this examination a mind blowing Web server called BioSeq-Analysis has been proposed to normally complete the three key steps for building a marker. The customer simply needs to move the benchmark instructive assortment. BioSeq-Analysis can make the upgraded pointer subject to the benchmark educational file, and the show measures can be represented as well. Also, to grow customer's solace, preliminary outcomes showed that the pointers made by BioSeq-Analysis even defeated some front line procedures. It is predicted that BioSeq-Analysis will transform into a significant gadget for characteristic progression examination.

**E. t-LSE**

Protein-protein affiliation (PPI) frameworks give encounters into the perception of natural techniques, work, and the crucial complex transformative instruments of the cell. Showing PPI orchestrate is a critical and chief issue in structure science, where it is still of huge concern to find an unrivalled fitting model that requires less assistant assumptions likewise, is logically solid against the huge piece of uproarious PPIs. At the present time, propose another methodology called t-key semantic embedding (t-LSE) to exhibit PPI frameworks. t-LSE endeavour's to adaptively get comfortable with an estimation introducing under the direct geometric assumption of PPI frameworks, moreover, a non-angled cost work was gotten to deal with the disturbance in PPI frameworks. The preliminary outcomes show the transcendence of the assault of t-LSE over other framework models to PPI data. In addition, the solid mishap work got here prompts immense updates for dealing with the disturbance in PPI compose. The proposed model could thusly support further diagram based examinations of PPIs and may help assemble the disguised crucial natural data.

**III. PROPOSED SYSTEM**

**A. Microarray dataset**

A microarray is a lab instrument used to distinguish the outflow of thousands of qualities simultaneously. DNA microarrays are magnifying lens slides that are printed with a great many little spots in characterized positions, with each spot containing a realized DNA arrangement or quality. Frequently, these slides are alluded to as quality chips or DNA chips. The DNA atoms connected to each slide go about as tests to identify quality articulation, which is otherwise called the transcript me or the arrangement of delivery person RNA (mRNA) transcripts communicated by a gathering of qualities. The datasets were all human microarray quality articulation information recovered from the GEO database. We at first gather ten liver disease related datasets. Since we need to discover differentially communicated qualities among wellbeing and liver malignancy conditions, we sifted through some datasets lastly kept three datasets. Their GEO promotion numbers are GSE84402, GSE76427, GSE64041, separately. They are gotten from various sorts of tissues and have various examples. The fundamental data of every datum set is appeared in Table 4.1.

Accession	samples	Tissues	Organism
GSE84402	28	HCC tissues corresponding non-cancerous tissues	homo sapiens
GSE76427	167	Primary HCC tumor tissue Adjacent non-tumor tissue	homo sapiens
GSE64041	125	Tumor from HCC patients non-tumor liver from HCC patients normal liver	homo sapiens

Table 1. Micro array Datasets

**B. Differentially Expressed Genes (DEGs)**

GEO2R was applied to recognize qualities that are differentially communicated across tumour and non-tumour tissues. GEO2R is a pleasant web apparatus that can be utilized to look at least two gatherings of tests in a GEO Series to recognize DEGs between exploratory conditions in an intelligent manner. Balanced P-values were processed to lessen the bogus positive rate through the default technique for Benjamin and Hochberg bogus revelation rate. Our cut off standard of choosing DEGs is the balanced P-esteem  $\leq 0:01$  and  $j\log FC_j \geq 0:5$

**C. Common Differentially Expressed Genes**

For each GO articulation, we need to check the repetition (k) of mannerism in the examination set (n) that are identified with the term, and the repeat (K) of facet in the masses set (N) that are associated with a comparable term. By then we test how likely would it be to secure at any rate k characteristics identified with the term if n habit would be aimlessly inspected from the masses, given the repeat K and size N of the people. The best possible quantifiable test is the one-followed variety of Fisher's cautious test, in any case called the hyper geometric test for over depiction. Right when the one-followed variation is applied, this test will calculate the probability of finding in any occasion the model repeat, given the masses repeat. The hyper geometric dissemination evaluates effectively the probability of k triumphs in n draws, without replacement, from a restricted masses of size N that contains definitely K productive things

Pathway examination is an amazing asset for understanding the science fundamental the information contained in enormous arrangements of differentially-communicated qualities, metabolites, and proteins coming about because of present-day high-throughput profiling technologies. The focal thought of this methodology is to aggregate these not insignificant arrangements of individual highlights into littler arrangements of related organic highlights (qualities and metabolites), normally dependent on natural procedures or cell parts in which qualities, proteins, and metabolites are known to be included.

**D. PPI Network**

Protein-protein correspondences (PPIs) are the mortal association of high unequivocally created connecting at any rate two protein particles because of biochemical occasions guided by affiliations that solidify electrostatic exertion, hydrogen holding, and the hydrophobic impact. Many are obvious relationship with the atomic association between chains that happen in a cell or a living structure in a particular bio molecular setting. Proteins conflictly act alone as their capacities





# Liver Cancer Key Genes Identification

will, generally speaking, be controlled. Different sub-atomic philosophy inside a cell is done by sub-atomic machines that are worked from various protein segments filtered through by their PPIs. These planned endeavour's make up the alleged interatomic of the living being, while atypical PPIs are the explanation behind different storing up related contaminations, for example, Creutzfeldt–Jakob, Alzheimer's diseases. PPIs have been concentrated with different techniques and from substitute points of view: characteristic science, quantum science, atomic parts, and signal transduction, among others. This data empowers the advancement of giant protein joint exertion structures – like metabolic or natural/epigenetic systems – that step in the to and fro development information on biochemical falls and sub-atomic etiologic of ailment, also as the exposure of putative protein central purposes of mending interest.

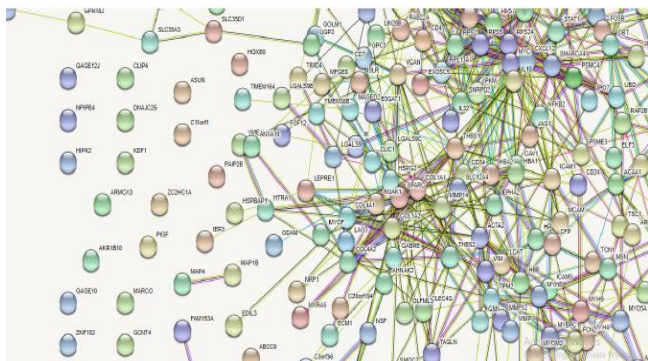


Fig. 1. PPI Network

## E. Analysed with Cytoscape

Cytoscape is an unfurled source bioinformatics programming stage for imagining sub-nuclear participation frameworks and fusing with quality enunciation account and other air data. Additional features are available as modules. Modules are open for organize and nuclear profiling examinations, new plans, additional report bunch support and relationship with databases and glancing in tremendous frameworks. Modules may be made using the Cytoscape open Java programming plan by anyone and module arrange headway is engaged. Cytoscape moreover has a JavaScript-driven sister adventure named Cytoscape.js that can be used to dismember and picture outlines in JavaScript conditions, like a portel. While Cytoscape is most regularly used for natural research use, it is realist to the extent use. Cytoscape can be used to envision and explore organize outlines of any kind including center points and edges (e.g., casual networks). A key piece of the inheritance think up Cytoscape is the usage of modules for specific features. Modules are made by focus architects and the more imperative customer organize.

## F. One-Hot Encoding

A major piece of the pre-processing is something encoding. This implies speaking to each bit of information such that the PC can see, thus the name encode, which truly signifies "convert to [computer] code". There's a wide range of methods for encoding, for example, Label Encoding, or as you may of speculated, One Hot Encoding. Mark encoding is instinctive and straightforward.

## G. Long Short Term Memory

Long transient memory (LSTM) units or squares are a bit of a monotonous neural framework structure. Discontinuous neural frameworks are made to utilize explicit sorts of fake memory shapes that can help these man-made mental ability

undertakings to even more effectively duplicate human thought. The dreary neural framework uses long transitory memory squares to offer setting to the way in which the program gets sources of info and makes yields. The long transient memory square is an incredible unit with various parts, for instance, weighted wellsprings of data, inception limits, commitments from past squares and potential yields.

## H. One Hot Decoder

In computerized circuits and AI, one-hot is a gathering of bits among which the legitimate mixes of qualities are just those with a solitary high (1) piece and all the others low (0). A comparative usage wherein all bits are '1' with the exception of one '0' is now and then called one-cold. In measurements, sham factors speak to a comparative system for speaking to straight out information

## IV. RESULTS

### A. Differentially Expressed Genes

ID	log <sub>2</sub> P.Val	P.Value	logFC	Gene.Symbol	Gene.Title	GO.Function
7892862	0.00261	2.95E-08	-1.91			
812971	0.00261	1.02E-07	-3.45E-01	CMBL	carboxymethyltransferase	carboxymethyltransferase activity
7991234	0.00438	3.95E-07	1.26	MFBFG8	milk fat g/l integrin binding	[phosphatidy]ethanolamine binding/[phosphatidylserine binding
807244	0.00593	1.14E-06	-1.42			
7942774	0.00593	1.40E-06	-0.61E-01	AQP11	aquaporin molecular function	water channel activity
7948944	0.00593	1.89E-06	-0.69E-01	GCAT	glycine-N-glycine-N-acyltransferase activity	[glycine-N-benzyloxycarbonyltransferase activity]/[protein binding]/[transferase activity]
8041179	0.00593	2.80E-06	1.12	CP19A	CP19A G-protein binding	
8126269	0.00593	3.00E-06	1.53	ALDOA	aldo-keto-1,5-bisphosphate reductase (NADP) activity	[peranrylseranylreductase activity]/[indanol dehydrogenase activity]/[protein binding]/[retinal dehydrogenase activity]
8126117	0.00593	3.12E-06	2.11	CG19A	gpg1 membrane binding	
8144866	0.00593	3.38E-06	-1.52	NAT2	N-acetyl-L-lysine N-acyltransferase activity	[L-tyrosyl transferase activity]
7927286	0.00593	3.94E-06	1.07	RASGEF4	Raf-1 protein binding	
8128707	0.00593	3.63E-06	1.31E-01	MICAL1	microtubule FAD binding/[Raf GTPase binding]/[Sh2 domain binding]/[actin binding]/[oxidoreductase activity, acting on paired donors, with incorporation or release of iron]	
7932985	0.00593	3.75E-06	1.38E-01	NRP1	neuropilin coreceptor activity	[cytokine binding]/[growth factor binding]/[growth factor binding]/[heparin binding]/[metal ion binding]/[protein binding]/[serine protease inhibitor activity]
8112668	0.00593	4.97E-06	1.53	GNAT4	glucosyl-N-acetylglucosaminidase activity	[beta-1,6-galactosyl-D-glucosyl-glycoprotein beta-1,6-N-acetylglucosaminidase activity]
7903301	0.00593	5.20E-06	1.66E-01	SLC35A3	solute carrier SLC35 family class III member 3	[L-tyrosyl transferase activity]/[protein binding]/[sugar protein symporter activity]
8064978	0.00593	5.33E-06	1.15	JAG1	jagged 1 Notch binding	[Notch binding]/[calcium ion binding]/[growth factor activity]/[protein binding]/[structural molecule activity]
7922402	0.00593	5.39E-06	1.15	GAS2L3	growth arrest specific 5 (non-protein coding)	[small nuclear RNA C/D box 47]
8142961	0.00593	5.81E-06	1.49	PCOLCE	podocalyxin protein binding	
8022442	0.00593	7.96E-06	2.29E-02	ZNF532	zinc finger metal ion binding	[sequence-specific DNA binding]/[transcription factor activity, sequence-specific DNA binding]/[transcription regulatory region DNA binding]
7939266	0.00593	5.51E-06	-1.36	BDKRL1	gamma-butyrolactone dihydrogenase activity	[gamma-butyrolactone dihydrogenase activity]/[iron ion binding]

Fig. 2. GEO2R

### B. Common Differentially Expressed Genes

```

jupyter GO Analysis Last Checkpoint: 03/07/2020 (autosaved)

writer.writerow(["Region"] +
                ["%s count" % c for c in conditions] + ["neglog p-value"])
out_info = []
for i, gene in enumerate(genes):
    counts = [int(work_count[c][gene]) for c in conditions]
    out_info.append([probabs[i], [gene] + counts])
out_info.sort()
writer.writerow(start + [prob] for prob, start in out_info)

In [11]: for row in csv_f:
         print(row)

[{"ID", "Name", "Size", "Expect", "Enrichment Ratio", "pValue", "Category"}]
[{"GO:0009184", "vasculature development", "478", "6.821885117", "4.10393211", "2.04E-10", "Biological Process"}]
[{"GO:0002939", "cardiovascular system development", "987", "6.51626221", "4.04023709", "2.70E-10", "Biological Process"}]
[{"GO:0001589", "blood vessel development", "109", "6.54399156", "4.12921307", "3.93E-10", "Biological Process"}]
[{"GO:0005235", "angiogenesis", "484", "4.87275747", "4.51489899", "3.78E-09", "Biological Process"}]
[{"GO:0048514", "blood vessel morphogenesis", "571", "5.74864365", "4.00994793", "1.57E-08", "Biological Process"}]
[{"GO:0005208", "extracellular matrix structural constituent", "137", "1.30862615", "7.99332059", "6.43E-08", "Molecular Function"}]
[{"GO:0032295", "tube development", "987", "0.93079736", "2.91844545", "1.50E-07", "Biological Process"}]
[{"GO:0042060", "sound healing", "525", "5.28935051", "3.78391244", "3.54E-07", "Biological Process"}]
[{"GO:0030198", "extracellular matrix organization", "346", "3.48345288", "4.59192036", "4.51E-07", "Biological Process"}]
[{"GO:0032239", "tube morphogenesis", "802", "0.87427785", "3.09625248", "5.00E-07", "Biological Process"}]
[{"GO:0042289", "RMC class II protein binding", "6", "0.05333749", "56.46120932", "1.34E-05", "Molecular Function"}]
    
```

Fig. 3. GO

```

jupyter KEGG Analysis Last Checkpoint: 03/07/2020 (autosaved)

[pathways, pathwaycount, pathwayrat] = kggpathwayrat(genes, data)
enrich = enrichment(GenesSet, pathwayMat)
result = out2html(GenesSet, pathwayMat, enrich, Genes, Pathways, geneLists, pathwaycount, ratio, Fileout, puidName)
return(result)

In [19]: for row in csv_f:
         print(row)

[{"ID", "Pathway Name", "Size", "Expect", "Enrichment Ratio", "P Value"}]
[{"hsa0512", "ECM-receptor interaction", "82", "1.13343604", "6.17929169", "1.27E-04"}]
[{"hsa0071", "Fatty acid degradation", "44", "0.68818393", "8.2219712", "3.23E-04"}]
[{"hsa00288", "Valine, leucine and isoleucine degradation", "48", "0.66573382", "7.53689736", "4.88E-04"}]
[{"hsa01180", "Metabolic pathways", "1288", "17.00282041", "1.74252973", "0.77E-04"}]
[{"hsa0474", "Protein digestion and absorption", "90", "1.24401591", "4.8210231", "0.00145597"}]
[{"hsa00268", "Glycine, serine and threonine metabolism", "40", "0.52884485", "7.23653465", "0.002145291"}]
[{"hsa04933", "AGE-RAGE signaling pathway in diabetic complications", "88", "1.354591488", "4.429379673", "0.002267734"}]
[{"hsa01212", "Fatty acid metabolism", "48", "0.65437392", "6.02877888", "0.004180548"}]
[{"hsa04392", "Hippo signaling pathway", "29", "0.40884801", "7.48412424", "0.007121301"}]
[{"hsa00380", "Tryptophan metabolism", "39", "0.539072123", "5.56511805", "0.016342897"}]
[{"hsa05143", "African trypanosomiasis", "34", "0.35148075", "11.38275504", "3.85E-04"}]
[{"hsa05144", "Malaria", "48", "0.48518647", "8.0278308", "0.001489513"}]
[{"hsa04418", "Regulation of actin cytoskeleton", "209", "2.13449153", "2.804447092", "0.019678561"}]
[{"hsa05213", "Endometrial cancer", "56", "0.578790882", "5.18321978", "0.01957993"}]
[{"hsa04725", "Cholinergic synapse", "111", "1.147246213", "3.48668988", "0.027451646"}]
    
```

Fig. 4. KEGG

### C. PPI Network

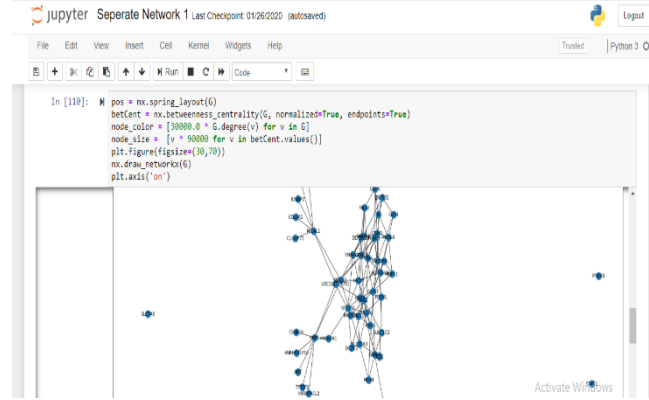


Fig. 5. PPI Network

### D. Analysed with Cytoscape

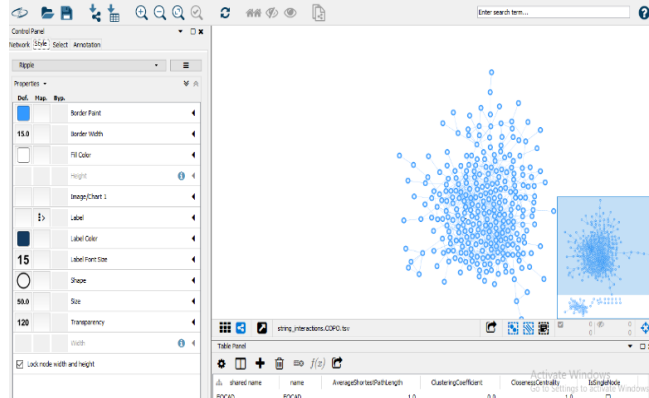


Fig. 6. Network Analysis

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	
SUID	Average	Betweenness	Closeness	Clustering	Degree	Excentricity	IsSingle	Name	Neighbor	Number	Number	Of	Radiality	selected	SelfLoops	shared	on	Stress	Topological	Coefficient
2	168	5.03862	0	1	9	FALSE	NMIRK1	5	1	0	0	0.88938	FALSE	0	NMIRK1	0	0	0	0	
3	239	4.24215	0	1	10	FALSE	IFNGR2	15	1	0	0	0.750599	FALSE	0	IFNGR2	0	0	0	0	
4	255	5.32699	0	1	12	FALSE	ETN2C	2	1	0	0	0.618715	FALSE	0	ETN2C	0	0	0	0	
5	308	6.33564	0	1	12	FALSE	CPY2J	5	1	0	0	0.589566	FALSE	0	CPY2J	0	0	0	0	
6	355	4.930796	0	1	10	FALSE	DAP1	5	1	0	0	0.697861	FALSE	0	DAP1	0	0	0	0	
7	394	4.608997	0	1	10	FALSE	ODAM	18	1	0	0	0.722385	FALSE	0	ODAM	0	0	0	0	
8	397	7.342961	0	1	11	FALSE	GOLGA3	3	1	0	0	0.512111	FALSE	0	GOLGA3	0	0	0	0	
9	411	4.961938	0	1	11	FALSE	CPY5A1	6	1	0	0	0.695236	FALSE	0	CPY5A1	0	0	0	0	
10	413	5.899564	0	1	12	FALSE	AQP11	2	1	0	0	0.623104	FALSE	0	AQP11	0	0	0	0	
11	426	6.33218	0	1	11	FALSE	SUC2A1	2	1	0	0	0.589822	FALSE	0	SUC2A1	0	0	0	0	
12	480	5	0	1	11	FALSE	CLEC4M	4	1	0	0	0.692808	FALSE	0	CLEC4M	0	0	0	0	
13	461	4.211873	0	1	9	FALSE	HNR1PA1	17	1	0	0	0.752994	FALSE	0	HNR1PA1	0	0	0	0	
14	473	6.892794	0	1	13	FALSE	CPY2A3	3	1	0	0	0.546713	FALSE	0	CPY2A3	0	0	0	0	
15	525	5.988616	0	1	13	FALSE	ROBO3	2	1	0	0	0.338588	FALSE	0	ROBO3	0	0	0	0	
16	529	5.854671	0	1	11	FALSE	CDXRC	2	1	0	0	0.628584	FALSE	0	CDXRC	0	0	0	0	
17	564	1	0	1	1	FALSE	MAP4	1	1	0	0	1	FALSE	0	MAP4	0	0	0	0	
18	565	1	0	1	1	FALSE	MAP1B	1	1	0	0	1	FALSE	0	MAP1B	0	0	0	0	
19	577	5.449827	0	1	11	FALSE	SERP1B	4	1	0	0	0.657706	FALSE	0	SERP1B	0	0	0	0	
20	579	5.449827	0	1	11	FALSE	SERP1A	4	1	0	0	0.657706	FALSE	0	SERP1A	0	0	0	0	
21	582	4.221453	0	1	10	FALSE	ZEB1	2	1	0	0	0.598035	FALSE	0	ZEB1	0	0	0	0	

Fig. 7. Features Extracted

```

trainnucleotidesquences - Notepad
File Edit Format View Help
CTGGCATCCCTTAACCCAG
GTTGAC TGAGGCGGAGGTT
ACAATCGGGCGGGCCGGG
GCAGAAGAGCAGGAGGAGCT
GCCATAGCAGGCTGTCTCT
AGATATATATACAGAAATG
GTCTGTGGGCACTTGTCTT
ACTAGGAGACCTGGGGCAAG
AGAGCCGGAGCCGCAACC
ACACTGGAGCAAACTGGA
ACACATCGGGAGCAGCGGA
AGATAGCTGTGACGCTTGC
ATTCTATCATCCCAATGAT
GTTCTGTGCAAAATTTGAT
GTTCTGTGCAAAATTTGAT
TTGGCGATGACCCGGGTT
CGAGTGAGAGAGGCGGGCG
CCCTGGGTAAGCCGAAGTC
GTTCTGTGCAAAATTTGAT
ATAAACAATCGGAAGTTTC
GGGCTCGGGCCCCGGGATG
AGACAGAAGAGATGGAGCT
CGGATCCCGGGGCGAGCCG
GTTCTGTGCAAAATTTGAT
GTTCTGTGCAAAATTTGAT
    
```

Fig. 8. Nucleotide Sequences

### E. One Hot Encoding

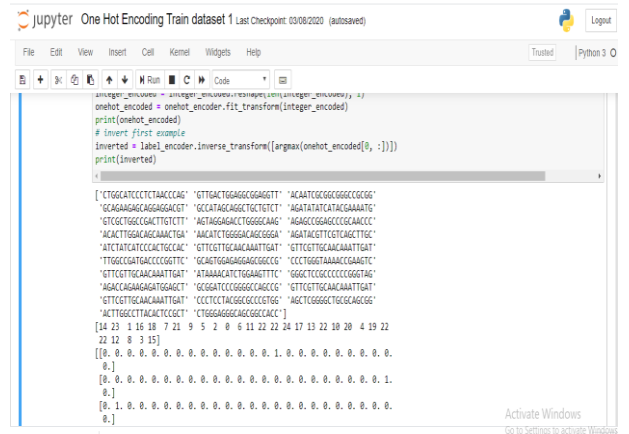


Fig. 9. One Hot Encoding Sequences

### F. LSTM

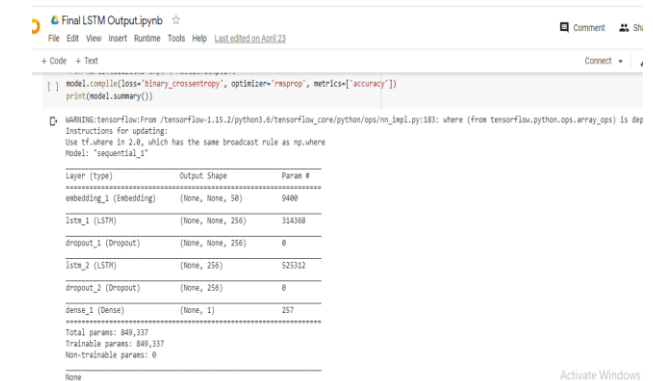


Fig. 10. LSTM Model

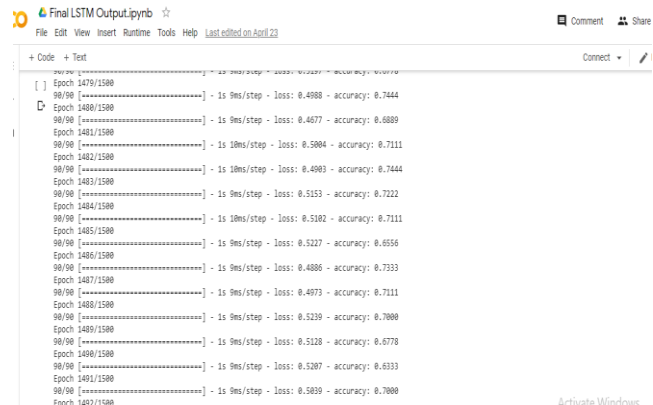


Fig. 11. Training Data

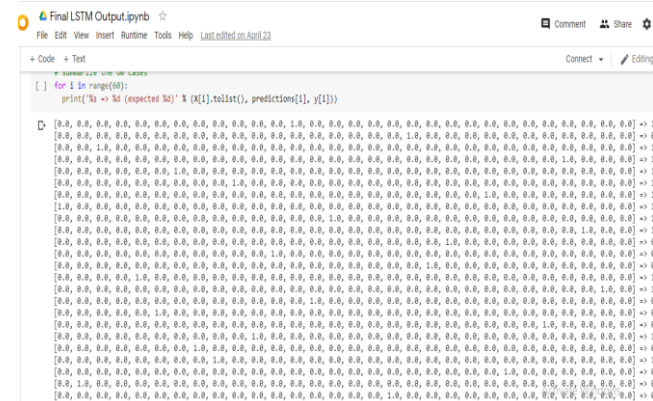


Fig. 12. Predictions

### G. One-hot decoding







**Prof. Joby George** is currently working as Associate Professor & HOD of Computer Science and Engineering department in Mar Athanasius College of Engineering, Kothamangalam, Kerala, India. He received his B-Tech Degree in Computer Science and Engineering in 1994 from MG university and M-Tech in Computer Science and Engineering from IIT Bombay in 2005. His research interests include Biometrics and Image processing.