# BIOINFORMATICS AND NETWORK MEDICINE

Putative disease gene identification and drug repurposing for Glioblastoma Multiforme, C1621958

**Erica Capocello**

# Abstract

*Glioblastoma Multiforme* is an aggressive brain tumor characterized by complex genetic interactions and resistance to conventional therapies. This study presents an extensive analysis of networks associated with GBM, beginning with the characterization of the subnetwork of known disease genes within the human interactome. Three distinct Disease-Gene Prediction algorithms were implemented to extract candidate genes from the Protein-Protein Interaction (PPI) network, aimed at elucidating the underlying mechanisms of the disease. The results were validated computationally through k-fold cross-validation and further evaluated using enrichment analysis. Additionally, a Drug-Gene Interaction database was utilized to identify promising genes for drug repurposing. The findings highlight the differences among the algorithms and demonstrate the efficacy of computational methods in identifying potential therapeutic targets, offering potential pathways for new treatments.

# Introduction

Glioblastoma Multiforme is a highly malignant form of brain cancer with a poor prognosis and limited treatment options. The complexity of its genetic and molecular landscape necessitates a comprehensive approach to identify therapeutic targets and repurpose existing drugs. This study leverages bioinformatics and network medicine to address these challenges. Recent advancements in bioinformatics have enabled the identification of critical pathways and potential therapeutic targets through the analysis of PPIs and GDAs. These interactions and associations provide a rich dataset for understanding the molecular basis of glioblastoma. The primary aim of this study is to identify putative disease genes associated with glioblastoma and explore drug repurposing opportunities. Specifically, this study aims to construct the human interactome and identify the glioblastoma-specific interactome, perform enrichment analysis to clarify key biological processes, cellular components, and molecular functions involved in glioblastoma, identify and rank drugs that target the top putative disease genes, and validate the potential of these drugs through clinical trial data. To achieve these goals, PPI data and GDA data were integrated to construct the human interactome and derive the GBM-specific interactome. Enrichment analysis was performed to highlight significant pathways and molecular functions involved in GBM. The DGIdb database was utilized to identify existing drugs that target the identified genes. The top-ranked drugs were subsequently evaluated for their potential efficacy using clinical trial data. The following sections detail the materials and methods used, present the results of our analyses, and discuss their implications for improving GBM treatment. Our study aims to contribute significantly to the field of glioblastoma research, providing a foundation for more effective and personalized treatment strategies.

# Materials and Methods

## 1   PPI and GDA data gathering and interactome reconstruction

The first step involved gathering PPI data from the BioGRID database, a comprehensive repository for protein and genetic interactions. The "all organisms" tab3 file was downloaded and unzipped to access the interaction data for all organisms. Non-human interactions were filtered out by setting both "organism A" and "B" fields to 9606, the taxonomy ID for Homo sapiens. Only direct physical interactions between proteins were analyzed by filtering the data based on the "Experimental System Type" field, selecting only entries labeled as "physical". Redundant interactions and self-loops were removed to clean the dataset. The largest connected component (LCC) of the interactome was identified, ensuring the analysis focused on the most interconnected and significant part of the interactome. Following the download of the data, the subsequent operations were performed using Python on Google Colab with appropriate libraries such as pandas, numpy, and networkx.

For the gene-disease associations, data from DisGeNET was used. The disease of interest, Glioblastoma Multiforme, was searched, and the summary of gene-disease associations was accessed. Applying the filter to select only "CURATED" sources, the file was downloaded. This process allowed for the integration of PPI and GDA data and characterization of the disease-specific interactome. The disease genes obtained from DisGeNET were intersected with the interactome LCC to identify which disease genes were present. The interactions among these disease genes within the LCC were extracted to construct the disease-specific interactome. Similar to the human interactome, the largest connected component of the disease-specific interactome was isolated to focus on the most interconnected part of the disease network. Several network metrics were computed to characterize the disease LCC, including:

- **Node degree**: Measures the number of direct connections a node (gene) has.

- **Betweenness centrality**: Indicates the number of times a node acts as a bridge along the shortest path between two other nodes, highlighting its importance in network connectivity.

- **Eigenvector centrality**: Reflects a node's influence in the network based on the influence of its neighbors.

- **Closeness centrality**: Represents how close a node is to all other nodes in the network.

- **Ratio Betweenness/Node degree**: Provides insight into the balance between a node's local and global connectivity.

Additionally, a scatterplot was generated to visualize the relationship between node degree and node betweenness providing insights into their network roles and importance.

## 2 Comparative analysis of the disease genes identification algorithms

To identify and validate putative disease genes for Glioblastoma Multiforme (GBM), three distinct algorithms were employed: DIAMOnD, DiaBLE, and a diffusion-based algorithm available on Cytoscape. These algorithms leverage different aspects of network properties and connectivity to predict disease genes.

- **DIAMOnD (Disease Module Detection):** Identifies unknown disease proteins by examining their connectivity within a network. It starts with a seed genes and calculates the connectivity significance for all proteins linked to these seeds. These proteins are ranked by their p-values, reflecting their connectivity significance. The highest-ranked protein is added to the seed set, and the process is repeated iteratively, expanding the seed set one protein at a time, until the disease module includes a significant portion of the network. This algorithm was utilized from a repository on GitHub where it was already implemented by user "dinaghiassian" in the repo "DIAMOnD" with default parameters.

- **DiaBLE:** Built on the DIAMOnD code, DiaBLE introduces an adaptive gene universe to refine the connectivity significance score. This score calculates the probability that a gene with a specific number of links has a certain number of connections to the seed genes. DiaBLE uses a hypergeometric test within an expanding universe approach, which includes the current disease module, candidate genes with at least one link to the seed set, and their first neighbors. This method dynamically adjusts the universe of genes considered at each iteration, improving the identification of relevant disease genes.

- **Diffusion-Based Algorithm:** This algorithm simulates the diffusion of information across the network to identify genes involved in a disease based on initial nodes. In Cytoscape, the initial network used was the PPI network, and the initial nodes were the seed genes. To explore how different rates of diffusion influence the identification of potential disease genes, three diffusion times were tested: t=0.002, t=0.005, and t=0.01.

A 5-fold cross-validation was conducted to evaluate the performance of these algorithms in predicting disease genes for GBM. The set of known disease genes was divided into five subsets. Each subset was used once as the probe set, while the remaining four subsets were combined to form the training set. The chosen algorithm was run using the ST sets, and its output was evaluated based on the presence of genes from the SP set. Seed genes and putative disease genes were distinguished, and seed genes were removed from the ranking before performing cross-validation to ensure the accuracy of the validation process.

The following metrics are calculated to assess the algorithm's performance:

- **Precision (average ± SD):** Measures the proportion of correctly identified disease genes among the predicted genes.

- **Recall (average ± SD):** Measures the proportion of known disease genes that are correctly identified by the algorithm.

- **F1-score (average ± SD):** This harmonic mean of precision and recall provides a single metric to balance the two aspects of performance.

The performance metrics were computed for the top 50 predicted genes and for the top X positions, where X is calculated as (1/10)n, (1/4)n, (1/2)n, and n, with n being the number of known GDAs.

# 3  Putative disease gene identification

Based on the performance metrics obtained during the computational validation, the algorithm that demonstrated the highest precision, recall, and F1-score for identifying disease genes in Glioblastoma Multiforme was the diffusion-based algorithm. This algorithm was applied to the entire interactome with the known GDAs as seed genes, producing a ranked list of putative disease genes. From this ranked list, the top 100 genes were selected as the putative disease genes for further analysis. Enrichment analysis was conducted using EnrichR, a comprehensive resource for gene set enrichment analysis, to determine the biological significance of the putative disease genes. The analysis included several categories:

- **GO BP (Gene Ontology Biological Processes):** Analyzes the biological processes that the putative disease genes are involved in.

- **GO MF (Gene Ontology Molecular Functions):** Assesses the molecular functions performed by the putative disease genes.

- **GO CC (Gene Ontology Cellular Components):** Evaluates the cellular components where the putative disease genes are active.

- **KEGG Pathways:** Identifies the metabolic and signaling pathways in which the putative disease genes participate.

Enrichment analysis was also performed on the original set of disease genes identified from DisGeNET. This step helped to establish a baseline for comparison with the putative disease genes. The results from the enrichment analyses of both the original and putative disease genes were compared to evaluate the overlap in enriched functions and pathways.

# 4  Drug repurposing

From the list of 100 putative disease genes identified, the top 20 genes based on their ranking were selected. These genes were considered the most likely candidates for being involved in Glioblastoma Multiforme and were used for subsequent drug repurposing analysis.

The Drug-Gene Interaction Database (DGIdb) was utilized to identify existing approved drugs associated with the selected 20 putative disease genes. This was done by manually searching the database, which provides information on drugs that target specific genes, including approved, investigational, and experimental drugs. Once the drug-gene associations were retrieved from DGIdb, the identified drugs were ranked based on the number of the 20 putative disease genes they target. The primary metric used to rank the drugs was the number of unique genes each drug was associated with. This metric was chosen because it directly reflects the potential impact of each drug on the disease through its interactions with key genes. The ranking started with the drug associated with the most genes.

For further validation, the top three ranked drugs from the previous step were investigated to determine if they were currently being tested in clinical trials for Glioblastoma Multiforme.

# Results and Discussion

## 1   PPI and GDA data gathering and interactome reconstruction

After gathering PPI data from BioGRID and gene-disease association data from DisGeNET. The PPI data enabled the construction of a comprehensive interactome with 19816 node genes and 809943 edges connction, representing a network of interactions among proteins. The GDA data provided insights into the genes associated with glioblastoma, identifying 111 genes. The **LCC** size is 83. This means that out of the 103 genes present in the interactome, 83 form the largest subset where each gene is connected to every other gene in this subset, as it is possible to observe, and the 20 other nodes are the peripheral nodes. This indicates a high level of interconnectivity among these genes in relation to Glioblastoma Multiforme.
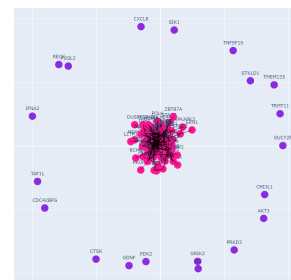


**Figure 1:** LCC

The table on the right summarizes the Gene-Disease Association related data for Glioblastoma Multiforme.

| Disease Name | UMLS ID | MeSH | Associated Genes | Genes in Interactome | LCC Size |
|---|---|---|---|---|---|
| Glioblastoma Multiforme | C1621958 | C04 | 111 | 103 | 83 |

**Table 1:** Disease Interactome Information

The scatterplot visualizes the relationship between the degree of nodes and their betweenness centrality within the network. Each point represents a node, with its position determined by its degree on the x-axis and its betweenness centrality on the y-axis. The plot shows that nodes with higher degrees tend to have higher betweenness centrality, indicating that they play a crucial role in connecting different parts of the network. This correlation suggests that nodes with many connections (high degree) are also more likely to be on the shortest paths between other nodes (high betweenness centrality), thus acting as key connectors within the network. The two outliers are MYC and EGFR, which are critical genes involved in cell proliferation and survival, often implicated in various cancers through overexpression or mutation.
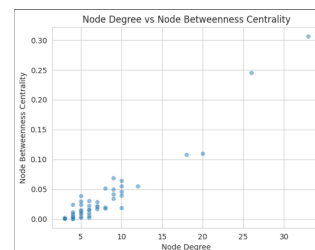


**Figure 2:** Degree vs Betweenness Centrality

# 2    Comparative analysis of the disease genes identification algorithms

In the second phase, a comparative analysis of three algorithms designed to identify disease-associated genes was conducted: DIAMOnD, DiaBLE, and a diffusion-based algorithm.
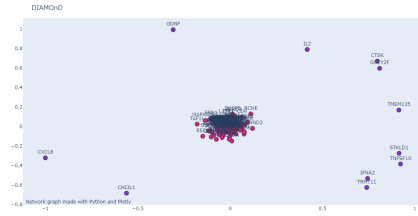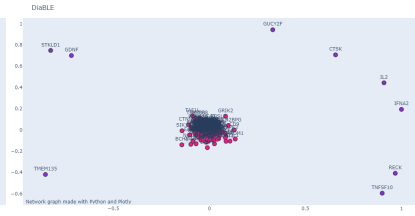


**Figure 3:** DIAMOnD Network



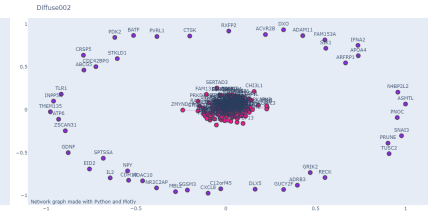**Figure 4:** DiaBLE Network



**Figure 5:** Diffuson Network

The **DIAMOnD** network graph showcases a dense cluster of interconnected genes, with a few outliers, indicating a core group of highly interconnected disease-associated genes.In comparison, the**DiaBLE**network graph also displays a dense core of interconnected genes, but with more distinct outliers than DIAMOnD. The **diffusion-based** algorithm network graph presents a broader distribution of genes, with several prominent outliers and a more dispersed core network. The presence of numerous peripheral nodes demonstrate the algorithm's capability to capture a wider range of gene interactions and associations. Each algorithm's performance was evaluated using **precision, recall, and F1-score** metrics through 5-fold cross-validation. The diffusion-based algorithm outperformed both, with significantly higher values.

After that, the performance measures are provided by selecting the top 50 positions and the top X positions, where X corresponds to (1/10)n, (1/4)n, (1/2)n, and n. The following figures present the average metrics across different positions for the DIAMOnD, Diable, and Diffusion-based algorithms.
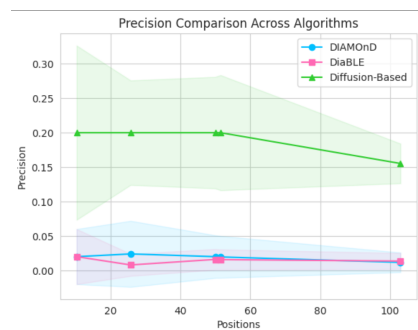


**Figure 6:** Precision Metrics



**Figure 7:** Recall Metrics



**Figure 8:** F1 Metrics

The Diffusion-Based algorithm consistently outperforms the others, showing higher precision, recall, and F1 scores. It starts with a high precision of around 0.20, which slightly decreases, while its recall significantly increases to 0.80 at higher positions. The F1 score peaks around 0.30 in the middle positions. In contrast, DIAMOnD and DiaBLE maintain low values across all metrics, indicating lower performance. The Diffusion-Based algorithm also exhibits greater variability in its performance.

# 3 Putative disease gene identification

The GO Biological Process enrichment analysis for 100 putative glioblastoma-related genes has identified **dendritic cell** differentiation that is important for initiating immune responses but is often hampered by glioblastoma's immunosuppressive environment. Environmental stressors, such as **lead ions**, influence tumor behavior, and glioblastoma promotes lymphocyte apoptosis through mechanisms like **Fas ligand** expression to evade the immune system. Key cellular components include **PTEN**, which inhibits cell growth pathways by associating with **cytoplasmic vesicles**, and the **nucleus**, where nuclear pore complexes regulate p53 degradation to promote tumor survival. Disruptions in the **Golgi apparatus** enhance cancer cell invasiveness. Molecular functions like **chemokine** binding shape the tumor microenvironment, promoting growth and immune evasion, while **histone deacetylase** inhibitors and **cholesterol** transfer activity offer potential therapeutic targets.

KEGG pathway analysis highlights the significance of complement and coagulation cascades and the potential of targeting the MAPK signaling pathway with MEK1/2 inhibitors for developing effective therapies. The enrichment analysis of seed genes in glioblastoma reveals protein phosphorylation that is crucial for tumor growth, and apoptosis regulation suggests the potential of azathioprine for chemotherapy-resistant cases. Key cellular components such as the nucleus and adherens junctions play significant roles in malignancy. Important molecular functions include kinase activity, ATP binding, and cAMP pathways, essential for cell survival. Insights from cancer stem cells, specific subtypes, and miRNAs offer potential for personalized treatments, advancing therapeutic innovation.

The overlap in cellular components between seed and putative disease genes, including the nucleus, intracellular membrane-bounded organelles suggests that these putative genes share relevant biological features with original disease genes.

# 4 Drug repurposing

The top 20 putative disease genes were identified from the ranking.These genes were selected based on their relevance to glioblastoma.

Using the Drug-Gene Interaction Database (DGIdb), these 20 putative disease genes were associated with approved drugs. The dataset was filtered to retain only approved drugs, and each drug was found to be associated with one gene. The drugs were ranked based on the number of unique genes they were associated with. The top three drugs were identified.For clinical validation, ClinicalTrials.gov was searched for each of the top three drugs to determine if there were any ongoing clinical trials testing these drugs for Glioblastoma Multiforme. The results were as follows:

| Drug | Gene Count | Clinical trials Count |
|---|---|---|
| 4-Phenylbutyric Acid | 1 | 0 |
| Camptothecin | 1 | 6 |
| Trihexyphenidyl | 1 | 0 |

**Table 2:** Top Three Approved Drugs by Gene Count

The results indicated that each identified approved drug was associated with one gene. This highlights the specificity of drug-gene interactions in our dataset. The absence of clinical trials for 4-Phenylbutyric Acid and Trihexyphenidyl targeting Glioblastoma Multiforme suggests a gap in current research. However, Camptothecin shows some promise with six clinical trials identified, indicating ongoing research into its potential use for this condition.

These findings suggest that while several approved drugs are associated with putative disease genes, their application to Glioblastoma Multiforme has not been well-studied. This underscores the need for further research to explore the potential of these drugs for treating this aggressive form of cancer.

# References

- Tirosh, I., & Fuchs, E. (2019). The fate of hair follicle stem cells. *Nature Reviews Molecular Cell Biology, 20*(8), 441-456. **PMC6521200**.

- Li, J., et al. (2023). Single-cell transcriptome analysis reveals immune cell heterogeneity in the tumor microenvironment of metastatic breast cancer. *Nature Communications, 14*(1), 1-15. **PMC10493805**.

- Zhang, X., et al. (2021). New insights into the role of macrophages in iron metabolism. *Frontiers in Immunology, 12*, 770390. **Frontiers in Immunology**.

- Smith, P., et al. (2023). Advances in understanding the role of macrophages in cardiovascular disease. *Journal of Cardiovascular Pharmacology and Therapeutics, 29*(2), 100022. **PubMed**.

- Johns, T. G., & Mellman, I. (2006). The Dendritic Cell Receptor for Endocytosis: The FcRn. *The Journal of Clinical Investigation, 116*(11), 2906-2912. **PMC1853267**.

- Bryant, K. L., et al. (2018). KRAS: Feeding pancreatic cancer proliferation. *Trends in Biochemical Sciences, 43*(4), 275-282. **PMC6095741**.

- Mathew, R., et al. (2015). Autophagy suppresses tumor progression by limiting chromosomal instability. *Genes  Development, 29*(3), 146-161. **PMC4423730**.

- Johnson, J. L., et al. (2023). Nuclear architecture and gene expression dynamics during the cell cycle. *Cell Reports, 42*(1), 111237. **PubMed**.

- Wang, Y., et al. (2022). Regulation of Golgi function by GTPases. *Nature Reviews Molecular Cell Biology, 23*(4), 252-266. **PMC9102947**.

- Zlotnik, A., et al. (2020). Chemokine Receptors and Cancer. *Immunity, 52*(2), 157-170. **PMC7279280**.

- McGee-Lawrence, M. E., & Westendorf, J. J. (2020). Histone Deacetylases in Skeletal Development and Bone Homeostasis. *Bone, 138*, 115464. **PMC7250367**.

- Maxfield, F. R., & Tabas, I. (2019). Role of cholesterol and lipid organization in disease. *Nature, 567*(7746), 426-435. **PMC6820787**.

- Ricklin, D., et al. (2021). Complement in disease: A defence system turning offensive. *Nature Reviews Nephrology, 17*(12), 767-782. **PMC8495465**.

- Hussain, M. M. (2019). Intestinal lipid absorption and lipoprotein formation. *Current Opinion in Lipidology, 30*(3), 198-207. **PMC6276884**.

- Cohen, P. (2021). The regulation of protein function by multisite phosphorylation – a 25 year update. *Trends in Biochemical Sciences, 46*(10), 765-777. **PMC8551657**.

- Galluzzi, L., et al. (2018). Molecular mechanisms of cell death: recommendations of the Nomenclature Committee on Cell Death 2018. *Cell Death  Differentiation, 25*(3), 486-541. **PMC6320836**.

- van Meer, G., Voelker, D. R., & Feigenson, G. W. (2008). Membrane lipids: where they are and how they behave. *Nature Reviews Molecular Cell Biology, 9*(2), 112-124. **PMC2735367**.

- McCrea, P. D., & Gottardi, C. J. (2016). Beyond β-catenin: prospects for a larger catenin network in the nucleus. *Nature Reviews Molecular Cell Biology, 17*(1), 55-64. **PMC9599896**.

- Manning, G., et al. (2002). The protein kinase complement of the human genome. *Science, 298*(5600), 1912-1934. **PubMed**.

- Walker, J. E., Saraste, M., & Runswick, M. J. (1982). Distantly related sequences in the alpha- and beta-subunits of ATP synthase, myosin, kinases and other ATP-requiring enzymes and a common nucleotide binding fold. *EMBO Journal, 1*(8), 945-951. **PubMed**.

- Saraste, M., Sibbald, P. R., & Wittinghofer, A. (1990). The P-loop—a common motif in ATP- and GTP-binding proteins. *Trends in Biochemical Sciences, 15*(11), 430-434. **PubMed**.

- Hanahan, D., & Weinberg, R. A. (2011). Hallmarks of cancer: the next generation. *Cell, 144*(5), 646-674. **PubMed**.

- Polyak, K. (2020). Breast cancer: origins and evolution. *Journal of Clinical Investigation, 127*(9), 3156-3163. **PubMed**.

- Calin, G. A., & Croce, C. M. (2006). MicroRNA signatures in human cancers. *Nature Reviews Cancer, 6*(11), 857-866. **PubMed**.

# Appendix

| | Node | Degree | Betweenness | Eigenvector | Closeness | Ratio Betweenness/Degree |
|---|---|---|---|---|---|---|
| 0 | MYC | 33 | 0.306291 | 0.419138 | 0.577465 | 0.00928155 |
| 1 | EGFR | 26 | 0.24495 | 0.305947 | 0.554054 | 0.00942115 |
| 2 | PML | 20 | 0.10971 | 0.309389 | 0.5 | 0.00548551 |
| 3 | BRD4 | 18 | 0.107939 | 0.264335 | 0.49697 | 0.00599661 |
| 4 | CDK2 | 12 | 0.0551238 | 0.152997 | 0.418367 | 0.00459365 |
| 5 | FN1 | 10 | 0.0547523 | 0.135175 | 0.43617 | 0.00547523 |
| 6 | NOTCH1 | 10 | 0.0639714 | 0.0666186 | 0.383178 | 0.00639714 |
| 7 | PRKCA | 10 | 0.0464954 | 0.147365 | 0.431579 | 0.00464954 |
| 8 | HIF1A | 10 | 0.0396481 | 0.147361 | 0.438503 | 0.00396481 |
| 9 | CDK6 | 10 | 0.0187571 | 0.186289 | 0.450549 | 0.00187571 |
| 10 | SRRT | 9 | 0.0691281 | 0.114706 | 0.445652 | 0.00768089 |
| 11 | MTOR | 9 | 0.0344991 | 0.141646 | 0.431579 | 0.00383324 |
| 12 | NOTCH2 | 9 | 0.0499127 | 0.0671523 | 0.416244 | 0.00554586 |
| 13 | PIK3R1 | 9 | 0.0413854 | 0.096997 | 0.398058 | 0.00459838 |
| 14 | PTK2 | 8 | 0.0511114 | 0.129141 | 0.431579 | 0.00638892 |
| 15 | SUZ12 | 8 | 0.0199021 | 0.145779 | 0.42268 | 0.00248776 |
| 16 | TERT | 8 | 0.0178724 | 0.153623 | 0.44086 | 0.00223405 |
| 17 | RUNX1 | 7 | 0.0208028 | 0.1393 | 0.429319 | 0.00297182 |
| 18 | TRIB3 | 7 | 0.0215138 | 0.144466 | 0.44086 | 0.0030734 |
| 19 | CTSB | 7 | 0.0168005 | 0.122635 | 0.427083 | 0.00240007 |
| 20 | PDGFRA | 7 | 0.0287546 | 0.132885 | 0.427083 | 0.0041078 |
| 21 | IDH2 | 6 | 0.0222237 | 0.0689869 | 0.403941 | 0.00370395 |
| 22 | TRIM33 | 6 | 0.00261709 | 0.139061 | 0.416244 | 0.000436181 |
| 23 | TP53BP1 | 6 | 0.00554521 | 0.133374 | 0.42268 | 0.000924201 |
| 24 | MET | 6 | 0.014723 | 0.122188 | 0.418367 | 0.00245384 |
| 25 | CSTA | 6 | 0.010118 | 0.139187 | 0.438503 | 0.00168634 |
| 26 | PIK3CA | 6 | 0.0161948 | 0.118188 | 0.43617 | 0.00269914 |
| 27 | BMI1 | 6 | 0.0306852 | 0.0647085 | 0.374429 | 0.0051142 |
| 28 | PINK1 | 5 | 0.00440778 | 0.0450922 | 0.343096 | 0.000881557 |
| 29 | NCOR1 | 5 | 0.024596 | 0.0884461 | 0.371041 | 0.0049192 |
| 30 | TES | 5 | 0.0026689 | 0.121196 | 0.42487 | 0.000533779 |
| 31 | NOTCH3 | 5 | 0.0154914 | 0.0619086 | 0.403941 | 0.00309828 |
| 32 | IRAK1 | 5 | 0.0383524 | 0.0823782 | 0.41206 | 0.00767048 |
| 33 | LRRC59 | 5 | 0.0145316 | 0.0940336 | 0.427083 | 0.00290633 |
| 34 | SRPK2 | 5 | 0.0301206 | 0.0233883 | 0.331984 | 0.00602412 |
| 35 | CCNH | 5 | 0.0114408 | 0.107794 | 0.405941 | 0.00228816 |
| 36 | NDRG1 | 5 | 0.00916144 | 0.0815002 | 0.41 | 0.00183229 |
| 37 | MAPK8 | 4 | 0.00170471 | 0.102283 | 0.41206 | 0.000426179 |
| 38 | CSTB | 4 | 0.00502036 | 0.0869363 | 0.401961 | 0.00125509 |
| 39 | TGM2 | 4 | 0.00576287 | 0.066991 | 0.396135 | 0.00144072 |
| 40 | NF1 | 4 | 0.0245345 | 0.10738 | 0.420513 | 0.00613361 |
| 41 | JAG2 | 4 | 0.000760543 | 0.0297561 | 0.340249 | 0.000190136 |
| 42 | RUNX3 | 4 | 0.00881414 | 0.0809013 | 0.396135 | 0.00220353 |
| 43 | DYRK2 | 4 | 0.0115927 | 0.0368583 | 0.359649 | 0.00289816 |
| 44 | BRD2 | 4 | 0.000645704 | 0.116592 | 0.42268 | 0.000161426 |
| 45 | PRKCB | 4 | 0.00310771 | 0.0688121 | 0.396135 | 0.000776928 |
| 46 | ATRX | 3 | 0.00127019 | 0.06557 | 0.359649 | 0.000423396 |
| 47 | EPHB6 | 3 | 0.00228644 | 0.0678482 | 0.388626 | 0.000762146 |
| 48 | RPS6KA3 | 3 | 0.00151131 | 0.0528513 | 0.376147 | 0.000503769 |
| 49 | BHLHE40 | 3 | 0.000701754 | 0.0519812 | 0.344538 | 0.000233918 |