



## NETWORK-BASED ANALYSIS OF BREAST CARCINOMA DISEASE ASSOCIATED GENES USING WEB-BASED BIOINFORMATICS TOOLS

Sonal Kharivale\* and Palak Larai

Department of Biotechnology,

Pillai College of Arts, Commerce & Science (Empowered Autonomous), New Panvel, India

\*Corresponding author E-mail: [sonalkharivale@mes.ac.in](mailto:sonalkharivale@mes.ac.in)

Received: 03 February 2026

Revised: 28 February 2026

Accepted: 19 March 2026

Published: 30 March 2026

DOI: <https://doi.org/10.5281/zenodo.19479914>

### Abstract:

Breast cancer is a heterogeneous disease driven by complex genetic alterations and dysregulated signalling pathways, resulting in significant variability among patients and challenges in effective treatment. Understanding the molecular interactions and identifying key regulatory genes are critical for improving therapeutic strategies. In the present study, an integrative bioinformatics approach was employed to investigate the molecular basis of breast carcinoma. A curated list of disease-associated genes was obtained from a reliable database and analysed using protein-protein interaction networks and functional gene association tools. Further, enrichment analyses across biological processes and signalling pathways were conducted to identify recurring molecular patterns. Pathway mapping provided insights into the involvement of these genes in established signalling cascades. Genes consistently identified across multiple analytical platforms were considered potential hub genes. The analysis revealed a core set of genes, including *ESR1*, *ERBB2*, *EGFR*, *VEGFA*, *AR*, and *GATA3*, which are known to regulate hormone signalling, cell proliferation, and angiogenesis. Key pathways such as receptor tyrosine kinase signalling, *ERBB2*-mediated pathways, and *VEGF*-driven angiogenesis were prominently enriched. The consistency of these findings with established literature validates the robustness of the approach and highlights its potential to identify novel therapeutic targets in breast cancer research.

**Keywords:** Breast Cancer, Bioinformatics, Hub Genes, Signalling Pathways, Angiogenesis, Protein-Protein Interaction, Gene Enrichment.

### Introduction

Breast carcinoma is still one of the most common cancers diagnosed in women and one of the leading causes of cancer death globally (1). What makes it particularly difficult to manage is how much the disease varies between patients. Estrogen receptor-positive, HER2-positive, and triple-negative subtypes each behave differently at the molecular level, follow different clinical trajectories, and do not respond equally well to the same treatments (2).

That biological variation is a large part of why treatment fails in many cases, why tumours come back, and why drug resistance keeps emerging as a clinical problem (3).

Treatments available today include chemotherapy drugs like anthracyclines and taxanes, hormone-based therapies such as tamoxifen and aromatase inhibitors, and targeted agents like trastuzumab, pertuzumab, and CDK4/6 inhibitors (4,2). Survival outcomes have improved considerably because of these advances. But the same treatments come with side effects, including cardiotoxicity, bone marrow suppression, and general systemic toxicity, and they do not work equally well across all subtypes. Triple negative breast cancer, in particular, remains difficult to treat because it lacks the molecular handles that targeted therapies depend on, and multidrug resistance further narrows options as the disease progresses (5,6). There is a genuine need for better targets and for treatment strategies that are built around the specific molecular profile of each subtype.

Part of what makes breast cancer so hard to target pharmacologically is that no single gene or pathway is responsible for the disease. It arises from the combined effect of many dysregulated genes feeding into each other within interconnected networks (7). Designing drugs against that kind of system is genuinely difficult, and the traditional drug discovery pipeline — which is already expensive, slow, and prone to failure — is not well suited to handling that complexity (8). Computational approaches offer a more practical entry point: they can survey large gene sets, map the relationships between them, and flag the most promising targets before any laboratory work begins.

Bioinformatics has developed into one of the most useful toolkits available for this kind of problem. Public repositories like the Gene Expression Omnibus (GEO) and The Cancer Genome Atlas (TCGA) contain large-scale gene expression and genomic data that can be mined to identify which genes are altered in breast carcinoma and how (9,10). Disease–gene association resources such as DisGeNET go a step further by aggregating evidence from multiple sources to score how confidently a given gene can be linked to a specific disease (11). Once candidate genes are in hand, tools like STRING and GeneMANIA place them within interaction networks and reveal which ones are most heavily connected — the genes that tend to be hubs are often the ones with the greatest influence over the system (12,13). Enrichr then helps interpret what those genes are actually doing by testing whether they cluster within known biological processes or pathways (14). Molecular docking and virtual screening can be used downstream to model how small molecules might interact with candidate targets and narrow down which ones are worth testing experimentally (15).

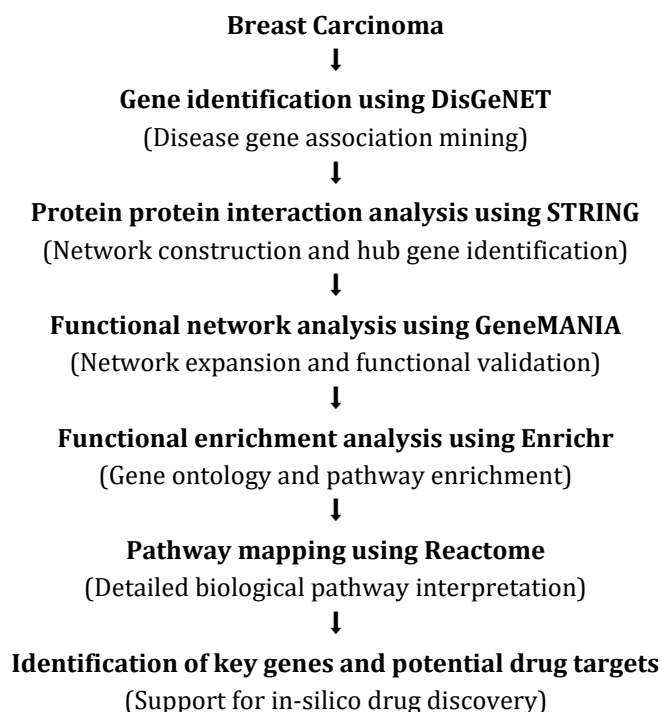
None of these tools is definitive on its own, but when they all point at the same genes, that convergence is harder to ignore (16,17). This study ran DisGeNET, STRING, GeneMANIA, Enrichr, and Reactome in sequence, using each to build on what the previous one found, with the end goal of identifying which breast carcinoma-associated genes are most consistently implicated and therefore most worth targeting in drug discovery.

### **Methodology**

No wet lab work was involved. Every analysis relied on publicly available, web-based tools and databases. The approach was staged: disease-associated genes were identified first, then examined in interaction networks, then tested for enrichment across biological pathways and gene-set libraries, and finally mapped onto curated molecular pathway databases.

**Bioinformatics tools used in the study****Table 1: Bioinformatics tools used in the study and their applications**

<b>Tool</b>	<b>Working Principle</b>	<b>Application in This Study</b>
<b>DisGeNET</b>	Integrates data from curated databases, GWAS, animal models, and scientific literature to establish disease gene associations with evidence scores.	Used to retrieve genes associated with breast carcinoma, which formed the primary dataset for further analysis.
<b>STRING</b>	Constructs protein protein interaction networks using experimental data, pathway knowledge, computational predictions, and text mining.	Used to build a protein protein interaction network and identify highly connected hub proteins involved in breast carcinoma.
<b>GeneMANIA</b>	Uses machine-learning to integrate co expression, genetic interactions, pathways, and protein interaction data to predict gene function and associations.	Used to expand and validate the interaction network and identify functionally related genes.
<b>Enrichr</b>	Performs statistical gene set enrichment analysis to identify overrepresented biological processes, molecular functions, and pathways.	Used to determine enriched biological functions and cancer related pathways associated with the selected genes.
<b>Reactome</b>	A manually curated, peer-reviewed database that maps genes and proteins to detailed biological pathways and molecular reactions.	Used to map prioritized genes onto curated pathways and understand molecular mechanisms involved in breast carcinoma.

**Work flow of the project**

**Results**  
**DisGenet**

Gene	Gene Full Name	N diseases	N variants	Scores	N PMIDs	N Chemicals	N PMIDs Chemicals	N variants	First Ref.
GATA3	GATA binding protein 3	233	128	0.2	5	0	0	0	2018
TRPS1	transcriptional repressor GATA class 1	550	338	0.2	2	0	0	0	2022
PIP	prolactin induced protein	352	3	0.2	2	0	0	0	2023
ESR1	estrogen receptor 1	3553	822	0.2	2	0	0	0	1998
MKI67	marker of proliferation Ki-67	1256	3	0.2	2	0	0	0	2008
ERBB2	erb-b2 receptor tyrosine kinase 2	550	105	0.2	2	0	0	0	2013
AR	androgen receptor	328	278	0.2	2	0	0	0	2012
MSLN	mesothelin	233	8	0.1	1	0	0	0	2008
FLT4	fnr3 related receptor tyrosine kin...	483	82	0.1	1	1	1	1	2018
PTHLH	parathyroid hormone like hormo...	256	17	0.1	1	0	0	0	1997

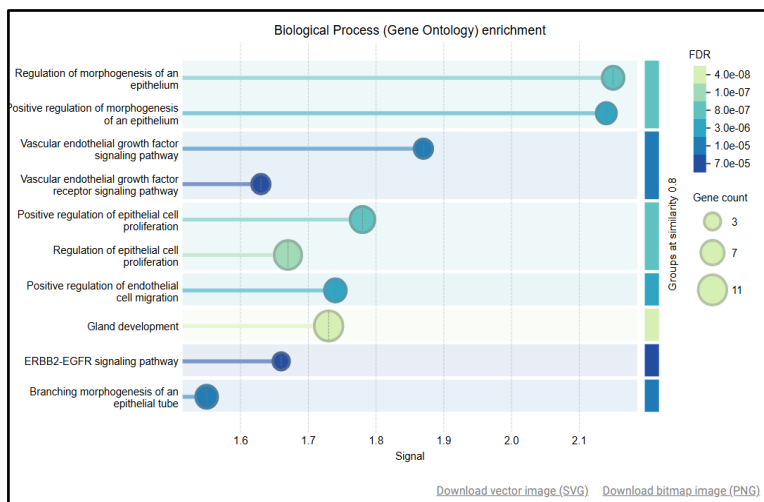
**Figure 1: Disease–gene association network of breast carcinoma-associated genes obtained using DisGeNET**

Querying DisGeNET for metastatic breast carcinoma produced a list that held together biologically rather than looking like a random pull from the database. The genes with the highest association scores were GATA3, TRPS1, PIP, ESR1, MKI67, ERBB2, AR, MSLN, FLT4, and PTHLH. Most of them are well-known enough that seeing them near the top of the list is not shocking — ESR1, ERBB2, and MKI67, for example, turn up in molecular profiling studies of nearly every breast cancer cohort. That familiarity is actually useful here, because it suggests the database query is returning clinically relevant genes rather than statistical noise. What the list also shows, when you look at it functionally, is a convergence on a few key processes: hormone signalling, cell growth control, transcriptional regulation, and angiogenesis. That overlap suggests these processes are not being independently disrupted in breast carcinoma but are failing together as part of a broader molecular breakdown. This gene list was carried forward into network and pathway analyses.

**STRING**



**Figure 2: Protein–protein interaction network of breast carcinoma-associated genes generated using the STRING database**



**Figure 3: Gene Ontology (Biological Process) enrichment analysis of breast carcinoma associated genes obtained from STRING**

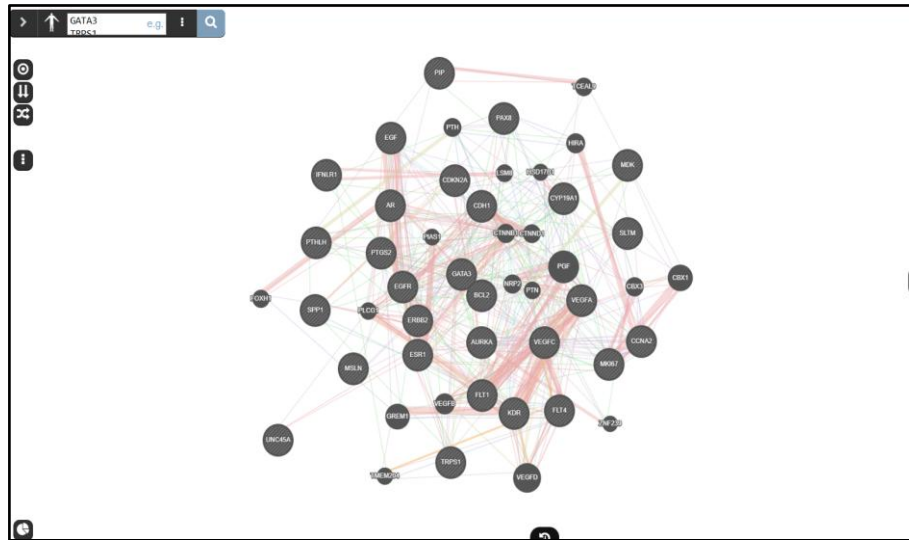
The GO enrichment output from the STRING network highlighted epithelial cell growth, tissue remodelling, VEGF signalling, and endothelial cell migration. If you know how breast cancer behaves, that list makes immediate sense: tumour cells need to proliferate, remodel their surroundings, attract new blood vessels, and eventually migrate to distant organs. The fact that ERBB2-EGFR signalling also appeared confirms that growth factor receptor activity is not incidental to this gene set — it is one of the core mechanisms at work. The network, in short, looks like breast cancer should look.

**Table 2: Summary of Gene Ontology (Biological Process) enrichment results from STRING analysis**

Theme	Representative GO Terms
Epithelial morphogenesis	Regulation of morphogenesis of an epithelium, positive regulation of epithelial morphogenesis
Tumor cell proliferation	Positive regulation of epithelial cell proliferation, regulation of epithelial cell proliferation
Angiogenesis	Vascular endothelial growth factor signaling pathway, VEGF receptor signaling
Cell migration and invasion	Positive regulation of endothelial cell migration
Growth factor signaling	ERBB2 EGFR signaling pathway
Tissue and gland development	Gland development, branching morphogenesis of epithelial tube

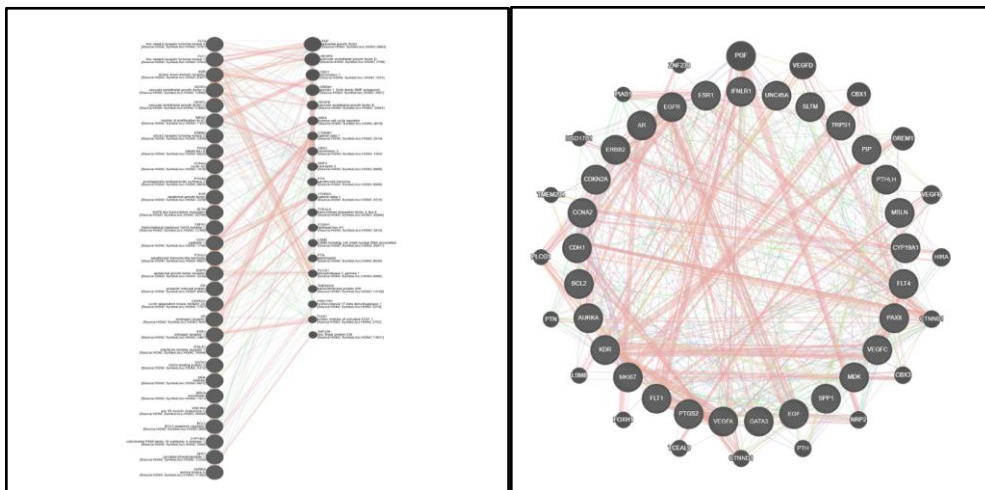
**GeneMANIA**

The GeneMANIA network was dense and well-supported, with connections backed by physical interaction data, co-expression patterns, pathway overlap, predicted functional links, and shared domain architecture. It was not a network built on one type of evidence with everything else missing — the same gene pairs tended to show up across multiple evidence categories, which strengthens confidence in the connections.



**Figure 4: Functional gene association network of breast carcinoma-associated genes constructed using GeneMANIA**

Sitting near the centre of all of this were ESR1, ERBB2, EGFR, AR, GATA3, MKI67, VEGFA, and FLT4, which are exactly the genes you would expect to find at the hub of a breast cancer network. Their position reflects genuine biological centrality, not just citation bias. The volume of connections converging on these genes suggests they are integration points in a system coordinating hormone signalling, cell cycle decisions, growth factor responses, and blood vessel formation simultaneously.



**Figure 5&6: Evidence-based interaction types supporting the GeneMANIA functional network**

The core of the network is built around ESR1, ERBB2, EGFR, AR, GATA3, and MKI67. Each of these genes has its own well-established role in breast cancer, but seeing them clustered this tightly together suggests something more than a coincidence of individual importance — they appear to be operating as part of the same regulatory machinery.

Enrichr

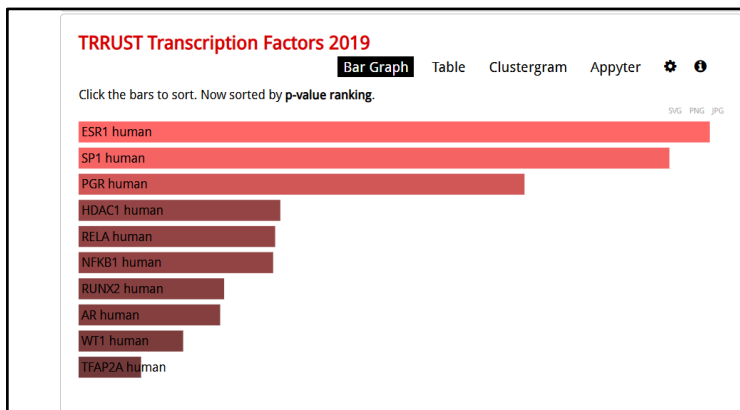


Figure 7: Bar graph showing enriched biological processes of breast carcinoma associated genes obtained using Enrichr

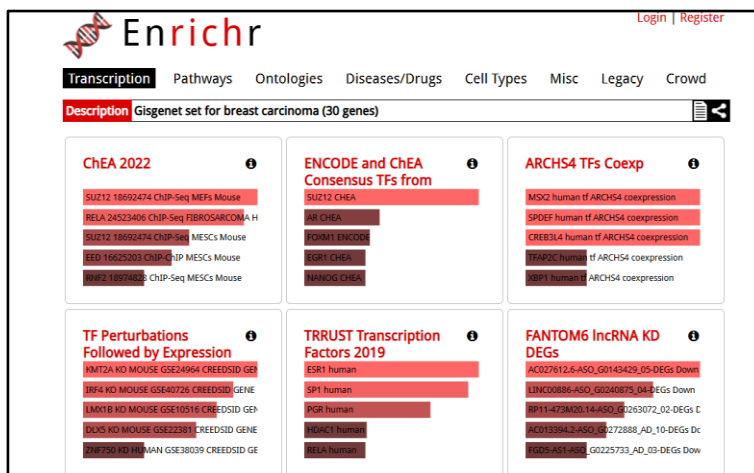


Figure 8: Enriched transcription factor regulatory networks of breast carcinoma associated genes identified using Enrichr

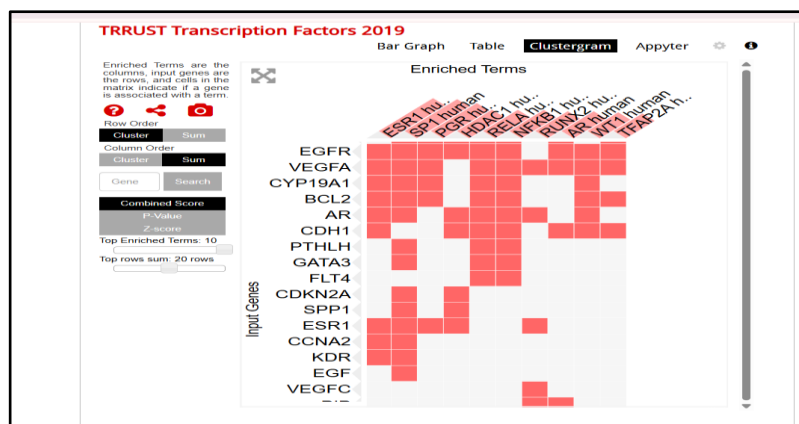


Figure 9: Disease and pathway enrichment analysis of breast carcinoma-associated genes obtained using Enrichr

**Table 3: Top breast carcinoma-associated hub genes supported across multiple enrichment databases**

Gene	KEGG / Reactome	ChEA (TF binding)	TRRUST (TF regulation)	Disease DB	Count
ESR1	✓	✓	✓	✓	4
EGFR	✓	✓	✓	✓	4
VEGFA	✓	✓	✓	✓	4
AR	✓	✓	✓	✓	4
ERBB2	✓	✓	✗	✓	3
GATA3	✗	✓	✓	✓	3

✓ = Significant association present; ✗ = Not significantly enriched

Count = Number of databases supporting each gene

To put the gene list in biological context, enrichment analysis was run through **Enrichr**. The analysis covered pathway libraries, transcription factor binding data, and curated disease gene sets. Any gene can score well in one library by chance; genes that score well across all three

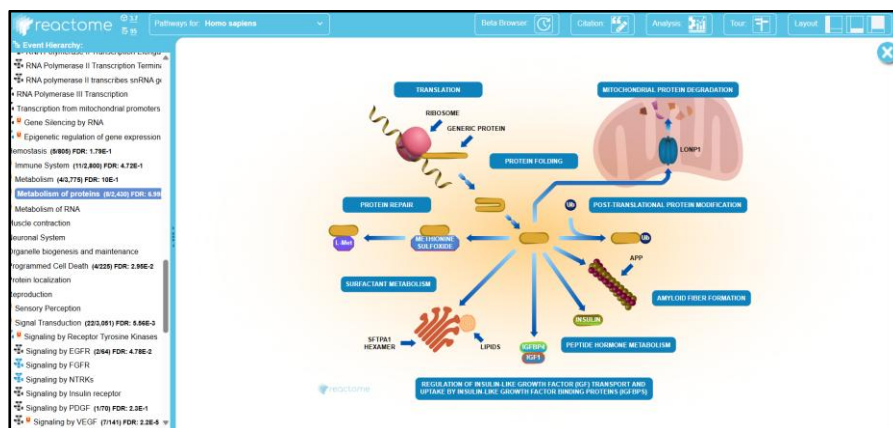
types of libraries are more likely to be genuinely important, so recurrence across library categories was used as the primary filter for identifying hub genes.

Six genes stood out by appearing across all three library types: ESR1, EGFR, VEGFA, AR, ERBB2, and GATA3. That level of cross-library consistency is not what you would expect from background noise, and it points to these genes being genuinely central to the regulatory landscape of breast carcinoma rather than incidentally enriched in one database.

**REACTOME**



**Figure 10: Reactome pathway enrichment analysis highlighting receptor tyrosine kinase signaling in breast carcinoma**



**Figure 11: Detailed Reactome pathway map showing ERBB2- and EGFR-mediated signaling pathways in breast carcinoma**

**Reactome pathway enrichment**

**Table 4: Reactome pathway enrichment analysis of breast carcinoma-associated genes**

Pathway Name	Number of Hits	p-value / FDR	Relevance to Breast Carcinoma
Signalling by Receptor Tyrosine Kinases	10	$2.82 \times 10^{-6}$	Key pathway driving tumor growth and survival
ERBB2 activates PTK6 signalling 3		$5.01 \times 10^{-4}$	Directly linked to HER2- positive breast cancer
VEGF binds to VEGFR leading to receptor dimerization	4	$2.82 \times 10^{-6}$	Central to angiogenesis and metastasis
PLCG1 events in ERBB2 signalling	3	$5.68 \times 10^{-5}$	Downstream growth factor signalling
TFAP2 (AP-2) family regulates transcription of growth factors	3	$4.9 \times 10^{-4}$	Transcriptional control of tumor growth
Signalling by EGFR	4	$6.88 \times 10^{-11}$	Growth factor mediated proliferation

Reactome confirmed the mechanistic picture that the earlier analyses had been building. The gene set enriched strongly in receptor tyrosine kinase signalling, ERBB2-associated cascades, and VEGF-driven angiogenic pathways — and not with marginal p-values. Receptor tyrosine kinase signalling sits at the intersection of nearly every major pro-tumour signal in breast cancer; ERBB2 pathway activation defines one of its most clinically recognised subtypes; and VEGF is the primary driver of the tumour vasculature without which a solid cancer cannot grow beyond a few millimetres. The statistics and the biology are pointing in the same direction.

**Conclusion**

Using integrated analyses of gene–disease networks, enrichment studies, and pathway mapping, a consistent molecular pattern in breast cancer was identified. A core set of genes—ESR1, ERBB2, EGFR, VEGFA, AR, and GATA3—repeatedly appeared across all tools, indicating their central role in disease-associated networks and

pathways such as receptor tyrosine kinase and VEGF signalling. Their recurrence confirms biological relevance rather than database bias. These genes are linked to cell proliferation, angiogenesis, and gene regulation. Future work should include experimental validation, functional studies, and clinical data integration (e.g., TCGA), alongside computational drug screening to support targeted therapeutic

## References

1. Siegel, R. L., Miller, K. D., Fuchs, H. E., & Jemal, A. (2024). Cancer statistics, 2024. *CA: A Cancer Journal for Clinicians*, 74(1), 17–48. <https://doi.org/10.3322/caac.21820>
2. Waks, A. G., & Winer, E. P. (2019). Breast cancer treatment: A review. *JAMA*, 321(3), 288–300. <https://doi.org/10.1001/jama.2018.19323>
3. Polyak, K. (2011). Heterogeneity in breast cancer. *The Journal of Clinical Investigation*, 121(10), 3786–3788. <https://doi.org/10.1172/JCI60534>
4. Harbeck, N., & Gnant, M. (2017). Breast cancer. *The Lancet*, 389(10074), 1134–1150. [https://doi.org/10.1016/S0140-6736\(16\)31891-8](https://doi.org/10.1016/S0140-6736(16)31891-8)
5. Holohan, C., Van Schaeybroeck, S., Longley, D. B., & Johnston, P. G. (2013). Cancer drug resistance: An evolving paradigm. *Nature Reviews Cancer*, 13(10), 714–726. <https://doi.org/10.1038/nrc3599>
6. O'Shaughnessy, J. (2016). Extending survival with chemotherapy in metastatic breast cancer. *The Oncologist*, 10(Suppl 3), 20–29. <https://doi.org/10.1634/theoncologist.10-90003-20>
7. Hanahan, D. (2022). Hallmarks of cancer: New dimensions. *Cancer Discovery*, 12(1), 31–46. <https://doi.org/10.1158/2159-8290.CD-21-1059>
8. Paul, S. M., et al. (2010). How to improve R&D productivity: The pharmaceutical industry's grand challenge. *Nature Reviews Drug Discovery*, 9(3), 203–214. <https://doi.org/10.1038/nrd3078>
9. Tomczak, K., Czerwińska, P., & Wiznerowicz, M. (2015). The Cancer Genome Atlas (TCGA): An immeasurable source of knowledge. *Contemporary Oncology*, 19(1A), A68–A77. <https://doi.org/10.5114/wo.2014.47136>
10. Clough, E., & Barrett, T. (2016). The Gene Expression Omnibus database. *Methods in Molecular Biology*, 1418, 93–110. [https://doi.org/10.1007/978-1-4939-3578-9\\_5](https://doi.org/10.1007/978-1-4939-3578-9_5)
11. Piñero, J., et al. (2020). The DisGeNET knowledge platform for disease genomics. *Nucleic Acids Research*, 48(D1), D845–D855. <https://doi.org/10.1093/nar/gkz1021>
12. Szklarczyk, D., et al. (2021). The STRING database in 2021: Customizable protein–protein networks. *Nucleic Acids Research*, 49(D1), D605–D612. <https://doi.org/10.1093/nar/gkaa1074>
13. Warde-Farley, D., et al. (2010). The GeneMANIA prediction server. *Nucleic Acids Research*, 38(Suppl 2), W214–W220. <https://doi.org/10.1093/nar/gkq537>
14. Kuleshov, M. V., et al. (2016). Enrichr: A comprehensive gene set enrichment analysis web server. *Nucleic Acids Research*, 44(W1), W90–W97. <https://doi.org/10.1093/nar/gkw377>
15. Pagadala, N. S., Syed, K., & Tuszynski, J. (2017). Software for molecular docking: A review. *Biophysical Reviews*, 9(2), 91–102. <https://doi.org/10.1007/s12551-016-0247-1>
16. Kuenzi, B. M., et al. (2020). Computational drug discovery and repurposing for cancer therapy. *Nature Reviews Drug Discovery*, 19(11), 673–690. <https://doi.org/10.1038/s41573-020-00107-4>
17. Li, X., et al. (2022). Bioinformatics analysis in breast cancer research and drug discovery. *Frontiers in Oncology*, 12, 843629. <https://doi.org/10.3389/fonc.2022.843629>