

Hierarchical Community Detection on Co-expression Networks for Functional Classification of Cell-Cycle Regulated Genes of *Saccharomyces cerevisiae*

Jaimielle Kyle Calderon¹ and Princess Angel Ventures²

¹ Algorithms and Complexity Lab,
University of the Philippines Diliman
jccalderon@up.edu.ph

² Algorithms and Complexity Lab,
University of the Philippines Diliman
pventures@up.edu.ph

Abstract. This study explored hierarchical community detection in gene co-expression networks to enhance the functional classification of cell-cycle regulated genes in *Saccharomyces cerevisiae*. Biological networks, which illustrate complex interactions among genes and proteins, required sophisticated methods to uncover their intricate structures. We evaluated three hierarchical community detection algorithms—Girvan-Newman (GN), Paris, and Local Optimization Function Model (LFM)—on a gene co-expression network of *Saccharomyces cerevisiae*. Our findings showed that the GN algorithm effectively identified distinct community structures at various hierarchical levels, demonstrating high modularity and closeness centrality. The Paris algorithm provided a comprehensive view of the network’s hierarchical structure, while the LFM algorithm revealed detailed sub-communities sensitive to the alpha parameter. These algorithms offered valuable insights into the hierarchical structures within biological networks, contributing uniquely to our understanding of gene co-expression patterns. The study’s results were validated using modularity, closeness centrality, and the Adjusted Rand Index (ARI), highlighting their potential applications in genomics and cellular biology. This research advanced the field by presenting new perspectives and methods for biological network analysis, paving the way for more targeted and effective approaches to understanding gene functions.

Keywords: Networks · Community Detection · Gene Co-expression · Hierarchical Structures · *Saccharomyces cerevisiae* · Biological Networks · Modularity · Closeness Centrality · Adjusted Rand Index (ARI)

1 Introduction

Prior research highlighted the urgent need for more sophisticated methods to unravel the complexities of biological networks. These networks, illustrating interactions among genes, proteins, and other biological entities, were incredibly

intricate. This complexity made it difficult to accurately interpret and analyze these networks, as noted by Besas and Redeakar. Their work underscored the importance of community detection in uncovering hidden structures and functional relationships within these networks.

Our study focuses on one specific challenge: exploring hierarchical community structures. This area holds great promise for uncovering the layered intricacies of biological networks, potentially revealing new insights into how different elements within these networks are interconnected and function. Building on foundational research by Rossetti and others, we aim to expand upon their methodologies and findings.

The goal of our study is to assess the capability of predicting hierarchical sub-clusters within the *Saccharomyces cerevisiae* gene co-expression dataset. By addressing this challenge, we hope to open new pathways for research and application in understanding biological networks.

We evaluated the performance of three hierarchical community detection algorithms—Girvan-Newman (GN), Paris, and Local Optimization Function Model (LFM)—on the gene co-expression network of *Saccharomyces cerevisiae*. The GN algorithm effectively identified distinct community structures at multiple hierarchical levels, demonstrating high modularity and closeness centrality. The Paris algorithm provided a comprehensive overview of the hierarchical structure, while the LFM algorithm revealed detailed sub-communities sensitive to the alpha parameter. Our findings indicate that these algorithms offer valuable insights into the intricate hierarchical structures within biological networks, each contributing uniquely to our understanding of gene co-expression patterns. Detailed evaluations using modularity, closeness centrality, and Adjusted Rand Index (ARI) further validated the efficacy of these methods, highlighting their potential for practical applications in genomics and cellular biology.

2 Background of the Study

2.1 Network Analysis

Network analysis served as our analytical beacon, guiding us through the maze of biological complexities. It's a versatile framework applied in diverse realms such as social networks, transportation pathways, and biological systems (Hanneman & Riddle, 2005; Hevey, 2018; Ma & Gao, 2012; X. Zhang et al., 2015).

In the context of biological networks, network analysis revealed patterns and relationships among biological entities. Our study focuses on the nuances of community detection in these networks, a critical step in understanding and interpreting the complex web of biological interactions.

2.1.1 Biological Networks

Important biological activities resulted from the coordinated effects of multiple molecules, cells, and genes interacting with each other (Ma & Gao, 2012). Initially, it was believed that biological processes operated in isolation and followed

linear sequences of events. However, as research progressed, it became clear that these processes form a complex, interconnected web of interactions (Bhalla & Iyengar, 1999). This revelation provided valuable insights into the integrated nature of biological functions, making network analysis a powerful tool in studying biology.

Biological networks can represent a wide range of factors such as protein-protein interactions, metabolic pathways, and gene co-expression networks (Yu et al., 2013). Several topological properties are usually analyzed in these networks, including degree and centrality for identifying hub genes or disease biomarkers, and motifs and modules/sub-networks for identifying functional units and recurring patterns (Ma & Gao, 2012; Rahiminejad et al., 2019).

Schwikowski et al. analyzed a global yeast protein network of 1,548 proteins with 2,358 interactions to predict functions. They successfully predicted functions for 72% of characterized proteins and assigned functions to 364 previously uncharacterized proteins. This network analysis revealed that proteins with established functions and shared locations tended to form clusters (Schwikowski et al., 2000).

Tang et al. used weighted gene co-expression analysis to identify potential prognostic biomarkers for breast cancer. Using gene expression data from the Gene Expression Omnibus (GEO) database, they found five distinct modules, with one module heavily correlated with the severity or grade of breast cancer. Within that module, they identified several hub genes, including CASC5 and RAD51, whose higher expression levels are linked to poorer survival outcomes for individuals with breast cancer (Tang, 2018).

2.1.2 Construction of Gene Co-expression Networks

The heart of our study lies in gene co-expression networks. These networks are intricate maps representing the functional relationships among genes, revealing how their expression patterns correlate under various conditions (Montenegro, 2022). This approach is crucial for understanding gene functions and plays a vital role in identifying novel metabolic pathways and regulatory networks.

Central to our understanding of these networks is the concept of communities non-overlapping groups of nodes within a network where nodes are more densely connected internally than with the rest of the network (Fortunato & Newman, 2022). Identifying these communities is instrumental in understanding the functional clustering of genes and can reveal previously unknown aspects of gene function and interaction.

However, community detection is fraught with challenges. It involves the automated discovery of groups of nodes that are either strongly connected or share similar features. The complexity of this task varies with the sparsity of the graph and the presence of generative noise. In some cases, polynomial-time algorithms suffice, while in others, the problem may require exponential time, making the hardness of community detection highly case-specific (Moore, 2017).

There are two methods for graph construction: rank-based and value-based. In graph representation, each vertex represents a gene, and an edge connects

two genes if they are co-expressed. With value-based construction, we compute the pairwise similarity of genes using Pearson’s correlation coefficient to form the gene co-expression network.

Gene expression data, reflecting gene activity levels in a sample, offers insights into the activation levels of genes and their co-expression patterns. Gene co-expression networks, showing relationships between different genes based on collective up-regulation or down-regulation, provide valuable information on functional relationships. Clustering gene expression data and analyzing it can offer possible annotations for genes with unknown functions based on other known genes in the same cluster (Oyelade et al., 2016; Yeung, 2001).

2.2 Community Detection in Networks

Exploring biological systems is challenging due to their inherent complexity and dynamic nature. Network science has revolutionized this field, especially with the introduction of community detection in gene co-expression networks. This methodology uncovers patterns and relationships within biological data, paving the way for meaningful discoveries.

Community detection in gene co-expression networks has evolved into a fundamental approach, with notable contributions from researchers such as Besas et al., who leveraged community detection algorithms to predict the biological classification of cell-cycle regulated genes in *Saccharomyces cerevisiae*.

Community detection involves identifying and classifying clusters where nodes (genes, in this context) exhibit higher connectivity within their group than with nodes outside it (Fortunato & Newman, 2022; Girvan & Newman, 2002). In biological networks, such as gene co-expression networks, community detection has been instrumental in unraveling hidden structures and functional associations.

A diverse array of community detection algorithms exists, each tailored to decipher specific types of networks. Fortunato classifies these algorithms into several categories:

- **Optimization approach:** This approach assigns scores to potential divisions of a network into communities, prioritizing divisions that yield high scores. Notable methods include Modularity, Spectral Clustering, and Hierarchical Clustering.
- **Statistical inference approach:** Here, communities are seen as primary drivers of network structure, with nodes connecting due to shared group affiliations. Methods such as Stochastic Block Model (SBM), Bayesian Inference, and Maximum Likelihood Estimation are popular in this category.
- **Dynamical process approach:** This method links community structure to dynamical processes occurring on networks. InfoMap and Random Walk Methods are exemplary techniques in this category.

2.2.1 Application of Community Detection in Biological Networks

Our challenge lies in choosing the right algorithm for community detection in biological networks. Selecting an appropriate community detection algorithm is

crucial as it significantly impacts the quality and validity of the insights derived from the analysis of biological networks. These networks, which can represent relationships between genes, proteins, or other biological entities, often exhibit complex structures that necessitate specialized approaches for accurate interpretation.

In recent years, the application of community detection to biological networks has been extensively explored. Redekar and Varma (2022) categorizes the commonly used algorithms for biological networks as shown in Table 1

Table 1. Comparison of Community Detection Methods (Redekar & Varma, 2022)

Category	Type	Methods
Type of community	Disjoint Community	Traditional Method (Graph Partitioning, Hierarchical Clustering, Spectral Clustering), Divisive Method (Girvan-Newman Method), Modularity-Based (Greedy Techniques, Simulated Annealing, External Optimization, Spectral Optimization)
	Overlapping Community	Clique Percolation Method, Line graph and Link Partitioning, Local Expansion and Optimization, Fuzzy Detection
Nature of network	Static Community	Graph Partitioning, Spectral Bisection, Hierarchical Clustering
	Dynamic Community	Random Walk, Spin Models, Synchronization

The Table 1 aids researchers in selecting the most appropriate method for specific network analysis needs, especially in the context of analyzing biological data like gene co-expression networks in *Saccharomyces cerevisiae*. Each method has its unique strengths and is suitable for different types of network structures and community detection requirements.

We identified three crucial open problems in existing literature: the identification of dynamic communities, an exploration of hierarchical community structures, and an analysis of overlapping communities. In response, we have curated a comprehensive list of algorithms tailored to address specific open problems in biological network analysis. This list, presented in Table 2, was compiled with insights from Vieira et al. (2020), Jaguzovic et al. (2022), Wagenseller III and Wang (2017), Rossetti et al. (2017), and Rossetti (2019). It provides a comparative overview of various community detection methods and approaches, facilitating the selection of the most suitable algorithm for each open problem.

The Table 2 includes a variety of community detection algorithms, each backed by significant research (Bonald et al., 2018; Coscia et al., 2014; Girvan & Newman, 2002; Gregory, 2010; Lancichinetti et al., 2009; Palla et al., 2005; Pons & Latapy, n.d.; Reichardt & Bornholdt, 2006; Rossetti et al., 2017; Rosvall & Bergstrom, 2008; Xie et al., 2011; Yang & Leskovec, n.d.). Each of these

Table 2. Comparison of Community Detection Methods for the Identified Open Problems

Open Problem	Algorithm	Approach	CDLib	Network	Communities
Identification of Dynamic Communities	Spinglass	Dynamical	spinglass	Undirected, Unweighted, Temporal	Crisp
	InfoMap	Dynamical	infomap	Directed, Weighted	Crisp
	Random Walk	Dynamical	walktrap	Undirected, Unweighted	Crisp
	TILES	Dynamical	tiles	Undirected, Unweighted, Temporal	Temporal
Examination of Hierarchical Community Structures	GN-Method	Optimization	girvan-newman	Undirected, Unweighted	Crisp, Hierarchical
	Local Optimization	Optimization	lfm	Undirected, Unweighted	Overlaps, Hierarchical
	Hierarchical Clustering	Optimization	paris	Undirected, Weighted	Crisp, Hierarchical
Analysis of Overlapping Communities	Clique Percolation Method	Optimization	kclique	Undirected, Unweighted	Overlaps
	Demon	Optimization	demon	Undirected, Unweighted	Overlaps
	SLPA	Optimization	slpa	Undirected, Unweighted	Overlaps
	Bigclam	Optimization	big_clam	Undirected, Unweighted	Crisp, Overlaps, Nested
	CFinder	Optimization	-	Undirected, Unweighted	Overlaps
	Copra	Optimization	-	Undirected, Unweighted	Overlaps

1 algorithms, with their unique approaches and applications, contributes signifi-
2 cantly to the field of community detection in biological networks, offering diverse
3 methodologies to address the identified open problems in the study.

4 2.2.2 Evaluation Methods of Community Detection Algorithms

5 Evaluating the success of community detection algorithms is typically done using
6 two types of metrics: internal criteria and external criteria.

7 Internal criteria, such as Modularity and Closeness Centrality, assess the
8 quality of detected communities without reference to any external structure,
9 providing a blind evaluation of the community structure.

10 External criteria, on the other hand, involve comparing the detected commu-
11 nity structure against a known reference, using metrics like Normalized Mutual
12 Information (NMI), Adjusted Rand Index (ARI), and Purity. These metrics of-
13 fer a benchmarked evaluation, essential for validating the detected communities
14 against a ground truth (Liu et al., 2019).

15 3 Related work

16 Having explored the complexities of biological networks and the nuances of com-
17 munity detection, we now delve deeper into the realm of gene expression data,
18 hierarchical community detection algorithms, and statistical analysis of these al-
19 gorithms. This section uncovers the gaps and opportunities in the current body
20 of knowledge. Specifically, we are interested in the sophisticated analysis of hi-
21 erarchical community detection algorithms, particularly in their application to
22 the *Saccharomyces cerevisiae* gene expression data.

23 3.1 Gene Expression Data

24 In the realm of publicly accessible resources, we identified a particularly pertinent
25 dataset comprising a subset of 384 genes, observed across 17 distinct time points.

1 This dataset, initially used in the study by Yeung et al. (2001) and originating
 2 from Cho et al. (1998), is a rich repository of information, accessible via the
 3 supplementary website for Yeung’s paper. The utilization of gene co-expression
 4 networks in our study was motivated by the accessibility and richness of such
 5 online datasets.

6 Cho et al. (1998) conducted a meticulous analysis of the *Saccharomyces cere-*
 7 *visiae* gene expression data derived from laboratory experiments. Their study
 8 categorized approximately 380 genes into five primary functional groups—early
 9 G1 phase, late G1 phase, synthesis phase, second growth phase, and mito-
 10 sis—adding further depth by identifying subgroups within each of these cate-
 11 gories. This classification provided granular insights into the functions associated
 12 with each gene ID.

13 Building upon this foundational work, Clemente (2022) conducted an anal-
 14 ysis of the *Saccharomyces cerevisiae* gene co-expression network. Using commu-
 15 nity detection techniques on the same gene set examined by Cho et al. (1998),
 16 Clemente’s study identified five primary clusters. These clusters broadly align
 17 with the five main functional groups identified by Cho et al. (1998), offering a
 18 macroscopic view of gene activity during various phases of the gene cycle. How-
 19 ever, this approach, while informative, leaves much to be desired in terms of
 20 specificity. The broad nature of these clusters obscures the nuanced functions of
 21 individual genes, presenting a significant limitation in our understanding of gene
 22 functionality within these groups.

23 This leads us to a conspicuous gap in the existing literature: a comprehensive
 24 analysis and comparison of results obtained from hierarchical community detec-
 25 tion algorithms, particularly applied to the *Saccharomyces cerevisiae* gene ex-
 26 pression data, remains unexplored. Such an analysis, especially one that focuses
 27 on identifying sub-clusters within the main functional groups, could provide a
 28 more detailed understanding of gene functions.

29 Recognizing the immense potential of sub-cluster analysis to yield finer,
 30 more detailed insights into gene functions—as suggested by H. Zhang et al.
 31 (2018)—our study aims to bridge this gap. By employing and comparing vari-
 32 ous hierarchical community detection algorithms, we aspire to uncover a more
 33 nuanced understanding of the *Saccharomyces cerevisiae* gene co-expression net-
 34 work, thus contributing a significant piece to the puzzle of gene functionality in
 35 biological research.

36 3.2 Hierarchical Community Detection Algorithmns

37 In the context of community detection within complex networks, hierarchical
 38 algorithms not only detect communities but also reveal the multi-level structure
 39 inherent within networks. In Table 2, we have highlighted three such hierarchi-
 40 cal algorithms available in CDLib: Girvan-Newman, Paris, and LFM. Each of
 41 these algorithms offers a unique approach to uncovering the layered structure of
 42 networks.

43 In the context of our study, these three algorithms – Girvan-Newman, Paris,
 44 and LFM – offer a diverse range of analytical lenses through which we can ex-

1 amine the gene co-expression network of *Saccharomyces cerevisiae*. Their unique
 2 characteristics and methodologies provide us with a comprehensive toolkit for
 3 dissecting and understanding the intricate web of relationships and structures
 4 within our dataset.

5 3.2.1 Girvan-Newman Algorithm

6 The Girvan–Newman algorithm stands out for its methodical approach in com-
 7 munity detection. It progressively dissects a network by systematically removing
 8 its edges. At each iteration, the algorithm identifies the ‘most valuable’ edge,
 9 typically the one with the highest betweenness centrality, and removes it. This
 10 process of edge removal continues until a clear community structure emerges
 11 within the network. The communities are revealed as the graph fragments into
 12 more tightly connected groups. The algorithm’s effectiveness lies in its ability
 13 to highlight the natural divisions within a network, which can be visualized
 14 through a dendrogram. The algorithm is particularly suited for undirected and
 15 weighted networks, though it does not accommodate directed networks. Its pro-
 16 cess, characterized by the systematic deconstruction and subsequent revelation
 17 of tightly-knit communities, provides a detailed and insightful view of the net-
 18 work’s architecture (Bonald et al., 2018; Clemente, 2022).

19 In terms of hardness, it has been known that Girvan-Newman is an NP-Hard
 20 problem (Brandes et al., 2008). Specifically, it is $(1 + \varepsilon)$ -inapproximability for
 21 dense graphs and logarithmic approximation for sparse graphs (DasGupta &
 22 Desai, 2013).

23 3.2.2 Paris Algorithm

24 Paris is an algorithm that combines the principles of modularity-based clustering
 25 with a hierarchical approach. It operates agglomeratively, basing its process on
 26 a unique distance metric between clusters. This metric is induced by the proba-
 27 bility of sampling node pairs within the network. Paris is particularly effective in
 28 its simplicity and its ability to reveal community structures at various levels of
 29 granularity. This agglomerative algorithm is adept at handling undirected net-
 30 works, although it does not support directed or weighted graphs. Its simplicity
 31 and the intuitive nature of its distance metric allow for a nuanced and layered
 32 understanding of a network’s community structure, making it an excellent tool
 33 for exploring gene co-expression networks where intricate relationships are often
 34 hidden within the data (Girvan & Newman, 2002).

35 3.2.3 LFM Algorithm

36 The LFM algorithm is a robust tool for detecting overlapping communities and
 37 understanding hierarchical structures within networks. It operates on the prin-
 38 ciple of local optimization of a fitness function. One of its key features is its
 39 adaptability, as it can uncover communities of varying sizes based on the alpha
 40 parameter. This parameter effectively controls the granularity of the community
 41 detection process, with larger values leading to smaller communities and vice

1 versa. This flexibility makes LFM suitable for undirected networks that require
2 a nuanced approach to community detection, allowing researchers to tailor the
3 algorithm to the specific needs of their data analysis (Lancichinetti et al., 2009).

4 **3.3 Statistical Analysis of Hierarchical Community Detection**

5 Having explored the capabilities of the Girvan-Newman, Paris, and LFM algo-
6 rithms in unveiling the hierarchical structures within the *Saccharomyces cere-*
7 *visiae* gene co-expression network, we now turn our focus to the critical aspect
8 of evaluating their effectiveness. The success of any community detection algo-
9 rithm lies not only in its ability to identify patterns but also in how accurately
10 these patterns reflect the true underlying structure of the data. To this end, we
11 employ both internal and external criteria in our statistical analysis, ensuring a
12 comprehensive and rigorous assessment of the algorithmic outputs.

13 **3.3.1 Internal Criteria**

14 Internal criteria evaluate the quality of the communities detected by the algo-
15 rithms based solely on the structure of the network itself, without reference to
16 external data. This self-contained analysis is important in determining how well
17 the algorithms have managed to capture the intrinsic clustering within the gene
18 co-expression network.

19 One of the primary metrics that we utilize in this study is Modularity. Modu-
20 larity measures the strength of division of a network into communities by assess-
21 ing the density of links inside communities compared to links between communi-
22 ties. A high modularity score indicates that the network has a strong community
23 structure, with more connections within communities than between them. This
24 metric is particularly useful in evaluating algorithms like Girvan-Newman, which
25 are designed to identify distinct and tightly-knit communities (Hajibabaei et al.,
26 2023).

27 Additionally, we consider the Closeness Centrality, which measures the de-
28 gree to which an individual node is near all other nodes in the network. In the
29 context of community detection, this metric helps in understanding how central
30 or peripheral each node is within its assigned community, thus providing insights
31 into the community structure’s coherence and the algorithm’s ability to allocate
32 nodes appropriately (Jarukasemratana et al., 2014).

33 By employing these internal criteria, we can critically assess the intrinsic
34 effectiveness of our chosen algorithms. This assessment not only validates the
35 detected communities in terms of the network’s own topology but also ensures
36 that the algorithms are effectively deciphering the complex, layered relationships
37 inherent in the *Saccharomyces cerevisiae* gene co-expression network. Such an
38 analysis forms the bedrock of our evaluation process, ensuring that the commu-
39 nities we identify are not only statistically significant but also meaningful in the
40 context of the biological data.

1 3.3.2 External Criteria

2 Central to our external validation process is the the Adjusted Rand Index (ARI).
 3 ARI provides a quantitative measure of the similarity between two data cluster-
 4 ings – one generated by our hierarchical community detection algorithms and the
 5 other representing the true clustering, as established by prior biological knowl-
 6 edge.

7 The Rand Index (RI) measures the similarity between two data clusterings
 8 by considering all pairs of samples and counting pairs that are assigned in the
 9 same or different clusters in the predicted and true clusterings. The Adjusted
 10 Rand Index (ARI), an enhanced version of RI, adjusts for the chance grouping
 11 of elements. It provides a more accurate measure of the similarity between the
 12 predicted and true clusterings. The ARI is calculated as $ARI = \frac{RI - E(RI)}{\max(RI) - E(RI)}$.
 13 A maximum ARI score of 1 indicates perfect alignment between the predicted
 14 and actual clusterings, while a negative ARI score implies an agreement level
 15 below random chance (Clemente, 2022; H. Zhang et al., 2018).

16 By incorporating ARI as our external criteria, we can gain a comprehensive
 17 understanding of how effectively the community detection algorithms categorize
 18 gene sets based on their biological functions. These metrics enable us to evaluate
 19 and compare the algorithms.

20 4 Research Aims

21 The comprehensive analysis we’ve undertaken thus far has not only shed light
 22 on the strides made in understanding gene networks but also brought into focus
 23 the critical gaps and uncharted territories in this field. These areas, brimming
 24 with potential for discovery and innovation, form the bedrock of our study’s ob-
 25 jectives. In the next sections, we articulate the specific challenges we addressed,
 26 alongside the objectives and scope of the study.

27 4.1 Problem Statement

28 In this study, we confront several critical gaps in the field of gene co-expression
 29 network analysis, particularly focusing on the *Saccharomyces cerevisiae* genome.
 30 Our objectives are threefold:

- 31 **1. Outdated Gene Function Classification:** To ensure the robustness and
 32 relevance of our findings, an updated and more extensive dataset was em-
 33 ployed by us, building upon the foundations of past research while incorpo-
 34 rating the latest advancements in genomic data. Since the classifications and
 35 sub-classifications from Cho et al. (1998) were outdated, we utilized existing
 36 databases for Yeast genome from sites like Kyoto Encyclopedia of Genes and
 37 Genomes (KEGG) and Saccharomyces Genome Database to update the gene
 38 function classifications.

- 1 **2. Limited Precision in Gene Function Identification:** Given that the five
2 major classifications from clusters found in the study by Clemente (2022)
3 were too general compared to the groups and subgroups seen in Cho et al.
4 (1998), this research pinpointed the roles and contributions of individual
5 genes with greater precision, addressing the limitations of existing classifica-
6 tions and paving the way for a better understanding of gene functions within
7 the *Saccharomyces cerevisiae* gene co-expression network.
- 8 **3. Absence of Network Analysis on Hierarchical Sub-clusters within**
9 **our Chosen Dataset:** A previous study by Salido (2016) has investigated
10 hierarchical clusters within our chosen dataset, but it limited its analysis to
11 traditional cluster methods, particularly Non-metric Multidimensional Scal-
12 ing (nMDS). Our study aimed for a network-based analysis approach to
13 identify the hierarchical clusters. The performance of each algorithm on the
14 gene co-expression network of Yeast cell-cycle regulated genes was systemat-
15 ically compared, contributing to a nuanced analysis of different hierarchical
16 community detection approaches in the context of gene co-expression net-
17 works.

18 4.2 Significance of the Study

19 The significance of this research lies in its potential to transform our understand-
20 ing of hierarchical structures in biological networks, with far-reaching implica-
21 tions:

- 22 1. This research aims to contribute significantly to the theoretical underpin-
23 nings of biological network analysis. By focusing on hierarchical structures,
24 it addresses a critical aspect of network science that is often overlooked but
25 is crucial for understanding complex biological systems.
- 26 2. Understanding hierarchical structures in biological networks has direct im-
27 plications for practical applications such as drug development, disease mod-
28 eling, and personalized medicine. Insights gained from this research could
29 lead to more targeted therapeutic strategies and a deeper understanding of
30 disease mechanisms.
- 31 3. The findings from this study could have interdisciplinary impacts, bridging
32 gaps between network science, biology, computational methods, and even
33 extending into fields such as ecology and sociology, where hierarchical struc-
34 tures also play a crucial role.
- 35 4. This research can also enrich academic discourse, providing valuable material
36 for educational purposes. It can inspire new curricula and teaching methods
37 in the fields of biology and network science.

38 4.3 Objectives of the Study

39 The study seeks to provide a comprehensive analysis of hierarchical community
40 structures within biological networks. The main objective is:

- 1 – To assess the capability of predicting hierarchical sub-clusters within the
2 *Saccharomyces cerevisiae* gene co-expression dataset, which can be accessed
3 through this link, using hierarchical community detection algorithms.

4 Alongside the main objectives, we have pinpointed specific objectives crucial
5 to the overall methodology of the study:

- 6 1. To update the gene expression dataset to reflect the latest gene classifica-
7 tion from Kyoto Encyclopedia of Genes and Genomes (KEGG) and Saccha-
8 romyces Genome Database.
- 9 2. To implement a Python script, integrating the CDLIB library and other
10 necessary imports, to automate the detection of sub-clusters in the *Saccha-*
11 *romyces cerevisiae* gene co-expression network. This tool will also facilitate
12 the comparative analysis of different algorithms using statistical methods
13 such as the Adjusted Rand Index (ARI) and the Rand Index (RI).
- 14 3. To identify which existing hierarchical community detection algorithm is the
15 most effective in detecting sub-clusters on the *Saccharomyces cerevisiae* gene
16 co-expression network. The current list of candidate algorithms can be found
17 in Table 2.

18 4.4 Scope and Limitations of the Study

19 The study’s primary focus was on identifying sub-clusters within the gene co-
20 expression network of *Saccharomyces cerevisiae*, with a specific emphasis on
21 cell-cycle regulated genes. The subset of 384 genes with available expression
22 data was the center of this analysis. Our goal was to uncover the underlying
23 structure and potential functional groupings within this network by employing
24 various hierarchical community detection methods. The effectiveness of these
25 methods in revealing biologically relevant sub-clusters was a key aspect of this
26 study.

27 In specifying the scope of this research, it is important to note that we did
28 not address overlapping nodes or communities within the gene co-expression
29 network. Our analysis focused on distinct sub-clusters, and the methodologies
30 we employed were tailored to identify non-overlapping hierarchical structures.
31 This decision was made to maintain clarity and specificity in the interpretation
32 of results.

33 In terms of methodology, we streamlined the process by adopting the most
34 effective gene co-expression network generation method identified in the study
35 by Clemente (2022). This method had been proven effective in constructing gene
36 co-expression networks for various genomes, including the human genome (Lee
37 et al., 2004). Notably, Clemente (2022) demonstrated its efficacy, showing that
38 this method was the best algorithm for constructing graphs in the specific yeast
39 genome dataset that we also utilized.

40 Furthermore, we did not venture into the development of new hierarchical
41 community detection algorithms. Instead, we leveraged existing algorithms, rec-
42 ognized for their efficiency and accuracy in various contexts, including biological

networks. This approach ensured that we utilized proven methodologies while concentrating on the unique aspects of gene co-expression in yeast.

To evaluate the performance of these community detection algorithms, we relied on existing scoring metrics. The choice of these metrics was crucial, as they needed to be apt for assessing hierarchical community detection in a biological context. We carefully selected metrics that not only provided quantitative assessments of the algorithms' performances but were also relevant and interpretable in the realm of gene expression and cellular biology.

By focusing on established methods and tools, we aimed to efficiently utilize our resources and expertise to delve into the complex interactions and organization within the *Saccharomyces cerevisiae* gene network. This focused approach enabled a more in-depth understanding of the network's structure and the biological significance of its sub-clusters, potentially leading to valuable insights in the field of genomics and cellular biology.

5 Methodology

As we move forward from defining the precise aims and objectives of our research, we now step into the methodology. In this section, we outline the comprehensive theoretical framework, which harmonizes computational analysis with experimental validation alongside the flow of our research methodology.

5.1 Theoretical Framework

A comprehensive approach to gene function identification involves integrating in silico network analysis with wet lab experiments. This blend of computational and experimental methods allows for the development of robust hypotheses and testing strategies, ultimately leading to the updating of gene annotations in biological databases. By synthesizing these two approaches, we can achieve a more holistic understanding of gene functions and interactions, paving the way for significant advancements in the field of gene co-expression network research.

To visually represent the intricate interplay between wet lab experiments and computational analyses in our research, we present a theoretical framework diagram in Figure 1. This diagram serves as a visual roadmap, illustrating the connections between the outcomes of wet lab experiments and the subsequent application of community detection algorithms on gene co-expression networks. The integration of in silico network analysis (*e.g. process done by Clemente (2022)) with wet lab experiments is foundational to the comprehensive approach for gene function identification. This showcases how the collaborative efforts of wet lab experiments and computational analyses converge to advance our understanding of gene functions within co-expression networks.

Now, let's take a look at the specific steps in our methodology, which are outlined in Figure 2. We provide an overview of the processes involved, from data collection and pre-processing to network construction and community detection script implementation, highlighting the systematic approach we used in this study. Each step will be discussed in detail in the following sections.

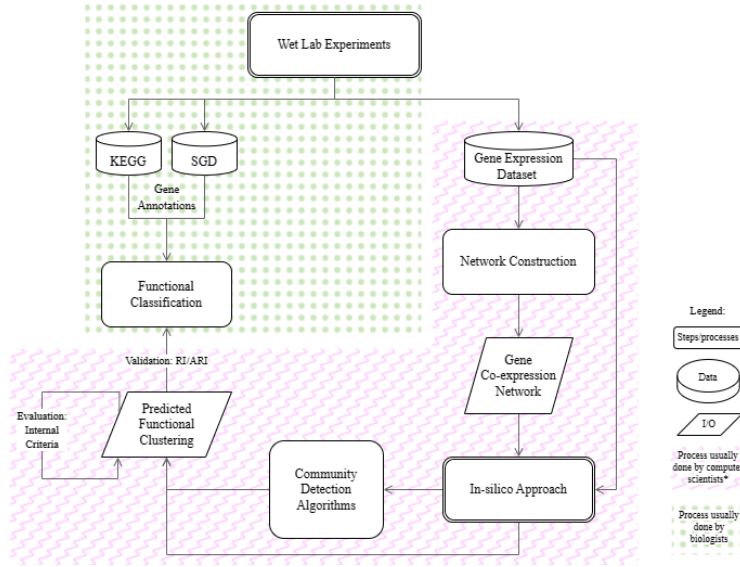


Fig. 1. Theoretical Framework Diagram

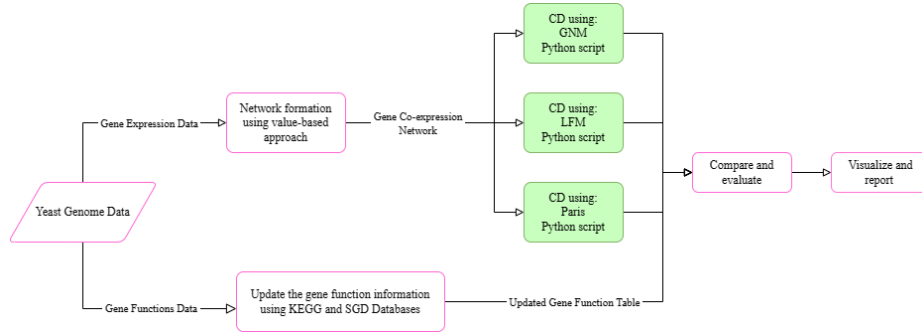
5.2 Data Collection and Organization

The first step in our methodology is to collect data from the Kyoto Encyclopedia of Genes and Genomes (KEGG) for each gene present in the dataset used by Cho et al. (1998). Our objective is to achieve a more precise grouping and classification for each gene, facilitating a comprehensive understanding of the dataset's hierarchical structure, essential for hierarchical community detection. This phase involves scraping relevant details for each gene. We have observed that there are three consistent levels of classification for every gene data that was scraped, each with a unique identification number. We will then further refine our scraped data by extracting only the specific ID for each classification at each level.

At this point, we have acquired the classifications for 242 genes (out of the 384 genes from the dataset used by Cho et al. (1998)). However, the information for the remaining 142 genes was incomplete in the KEGG database, lacking essential classifications. Consequently, these genes were removed from the dataset.

To gain a clearer understanding of the groupings, we chose to generate a graph visualization. This visualization was represented as a forest, with each tree root corresponding to the level 1 classifications. The subsequent tree levels depicted the level 2 and level 3 classifications, while the leaves represented the genes themselves.

Figure 3 displays the generated graph for the 242 genes and their respective classifications at each level. The hierarchical structure of the classifications within KEGG is clearly depicted in the graph, as certain classifications are positioned "under" the classification at the level above them, and each classification

**Fig. 2.** Methodology Flowchart

in the lower levels belongs to a classification in the upper levels. Additionally, we observed from this graph that there are no overlapping nodes; therefore, we did not need to filter out anything, as overlapping nodes are not within the scope of this study.

The next step is to investigate potential relationships or correlations between the groupings at any level from the KEGG data and the two levels classified by Cho et al. (1998) in their study. To accomplish this, we utilized the Adjusted Rand Index (ARI). To proceed, we must filter the KEGG data to include only the genes present in Cho et al. (1998) more detailed classifications. This constitutes a smaller number (149 genes), as Cho et al. (1998) did not include all 384 genes in the more specific classifications.

We generated a graph for the groupings from Cho et al. (1998) which is shown in Figure 4. ARI was performed for each combination of levels, with the highest score being around 0.16. The scores for each combination are shown in Table 3.

Table 3. Relationships between KEGG and CHO classifications with ARI

KEGG Level	CHO Level	ARI
1	1	0.022699
1	2	0.083272
2	1	0.057357
2	2	0.118445
3	1	0.043423
3	2	0.156492

Collecting and organizing the data provided a foundational understanding of the hierarchical structures within the dataset and identified potential relationships between classifications from KEGG and classifications from Cho et al.

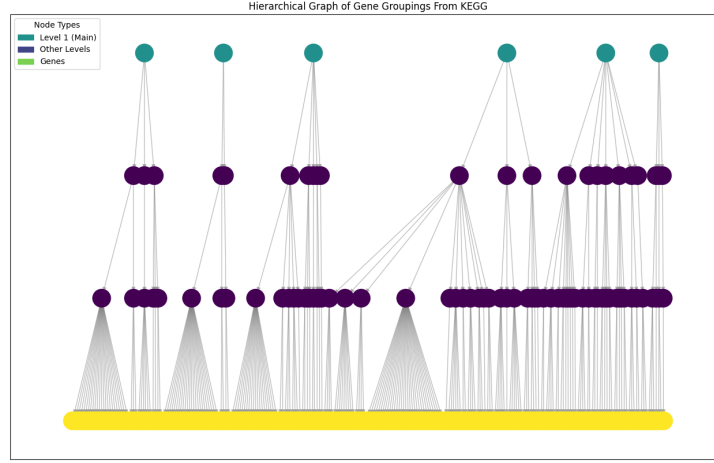


Fig. 3. Graph of Gene Groupings From KEGG

(1998). We will leverage this groundwork to construct gene co-expression networks.

5.3 Gene Co-expression Network Formation

Building upon the insights gained from the previous section, our next step is to delve into the construction of gene co-expression networks. Using the gene co-expression data provided by Yeung (2001), we employ a value-based algorithm for network construction, a method validated for its superior performance, as demonstrated in Clemente (2022).

We've organized our coding process in a Jupyter notebook, ensuring clarity and facilitating navigation across project stages. Utilizing the NetworkX library in Python, we aimed to construct a gene co-expression network by calculating pairwise gene similarities via Pearson's correlation coefficient. This method enables precise assessment of gene expression relationships.

Introducing the δ parameter, we established a crucial threshold for determining co-expression between gene pairs. This parameter's flexibility allowed us to fine-tune our network based on varying degrees of gene expression similarity.

In our analysis, we employed value-based graph construction, exploring thresholds from 0.70 to 0.95. This broad spectrum allowed us to understand how network sensitivity varies with different similarity thresholds, providing comprehensive insights into network dynamics.

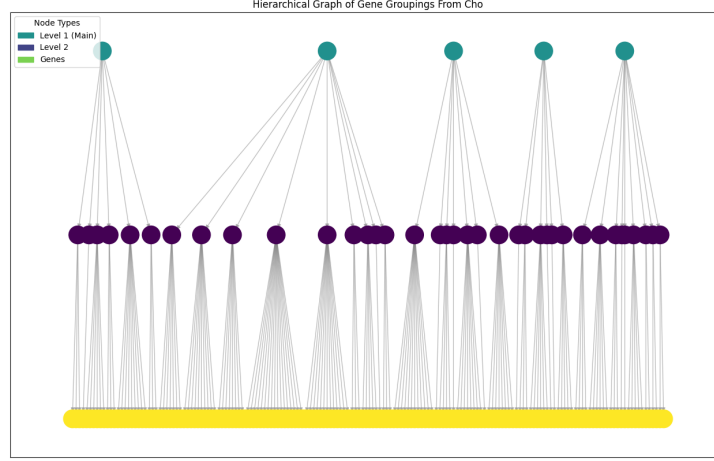


Fig. 4. Graph of Gene Groupings From Cho et al. (1998)

Our examination revealed that a small portion of nodes constitutes the largest connected component, with few singletons. Around 63% of nodes are integrated into non-trivial connected components, underscoring the network’s complex and interconnected nature.

We produced a graph depicting the gene co-expression network, as depicted in Figure 5. This visualization facilitated interpretation of our findings and conveyed the complex relationships within the yeast gene expression data.

Our graph illustrates that setting δ to 0.85 results in a compelling structure, marked by the highest number of non-trivial connected components.

Remarkably, this mirrors the findings of the study by Clemente (2022) which we sought to build upon. This reveals analogous patterns of connectivity and functional grouping, affirming the credibility and significance of these co-expression relationships.

5.4 Script Implementation for Community Detection

For the script implementation of the automated community detection, the process began with a thorough review of available algorithms across different libraries, particularly NetworkX and CDLib. Notably, while NetworkX supports various algorithms, LFM and Paris were absent, emphasizing CDLib’s unique capabilities.

We tested CDLib’s community detection algorithm GN, which successfully identified meaningful communities within our network data. A pivotal parameter

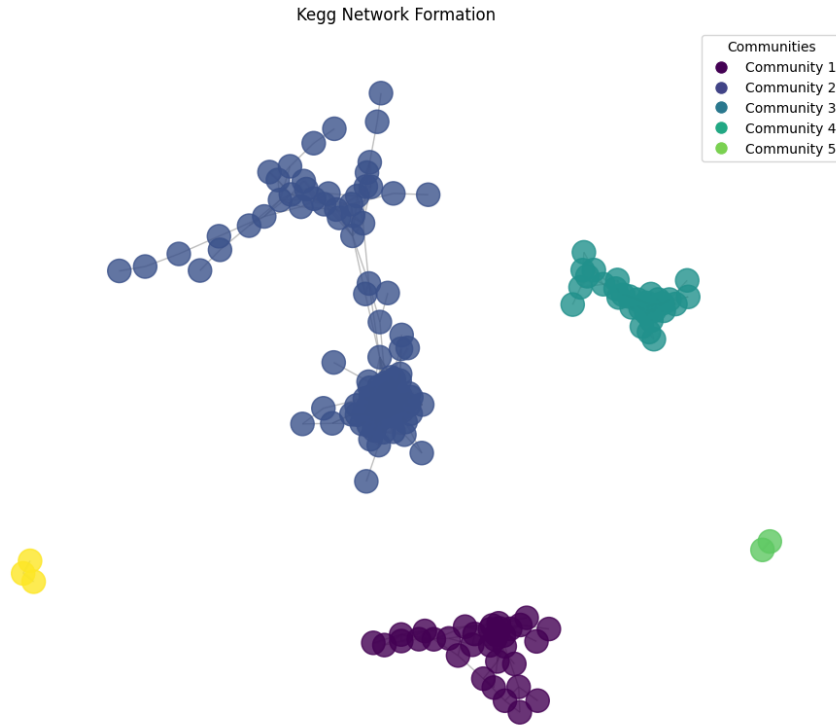


Fig. 5. Gene Co-expression Network

1 in this process is the "level," determining the depth of analysis and indicating
 2 the stage at which we opt to cut the dendrogram to observe the community
 3 structure.

4 As the level increases, we transition from a macro-view of the network's
 5 community structure (comprising broad, large communities) to a more micro-
 6 view (highlighting smaller, more specific communities). These insights will be
 7 elaborated upon in the following section, accompanied by detailed discussions of
 8 the generated graphs.

9 Similarly, CDLib is employed for the LFM algorithm, which also unveils a
 10 hierarchical structure of communities. The alpha parameter regulates the res-
 11 olution of community detection, influencing both the size and the number of
 12 communities identified.

13 We experimented with varying alpha values from 0.5 to 1.0. The results
 14 indicate a significant impact of alpha on the outcome of community detection,
 15 with more communities emerging at higher alpha values. This underscores LFM's
 16 sensitivity to the alpha parameter and its capability to adapt to different scales of

1 community structure. Graphs illustrating different alpha values will be presented
2 in the subsequent section.

3 Lastly, CDLib was used for the Paris algorithm. We have found that this
4 specific method did not involve changing any parameters.

5 Now that we have finished implementing the scripts for the three algorithms,
6 our next step is to evaluate each algorithm by comparing its performance (groups
7 formed) against our dataset. This is discussed in the next section.

8 5.5 Visualization, Evaluation, and Analysis

9 To understand how the three community detection algorithms compare to one
10 another, we employed several evaluation methods on the communities formed
11 by each algorithm. Specifically, we used modularity, closeness centrality, and
12 the Adjusted Rand Index (ARI) to assess the quality and effectiveness of the
13 detected communities.

14 In addition to these quantitative measures, we created visual representations
15 of the networks using graphs. These visualizations help provide a clearer and
16 more intuitive understanding of the community structures and the overall net-
17 work.

18 5.5.1 Girvan-Newman

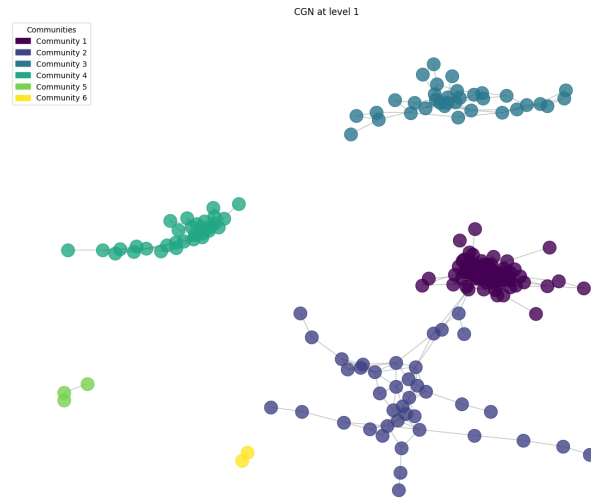


Fig. 6. Communities Detected Using GN at Level 1

19 At Level 1, the Girvan-Newman (GN) algorithm successfully identified six
20 distinct communities within the gene co-expression network, as illustrated in

1 Figure 6. Each community is represented by a different color, highlighting the
 2 clear separation and clustering of nodes.

3 The visualization reveals that the communities are not only well-defined but
 4 also vary significantly in size. Community 1 and Community 2 are the largest,
 5 containing the majority of the nodes, while Community 6 is the smallest. This
 6 distribution of community sizes may reflect the inherent hierarchical nature of
 7 gene interactions, where certain groups of genes are more densely interconnected
 8 due to similar functions or regulatory mechanisms.

9 The identified communities are spread across the network, showing distinct
 10 clusters that are interconnected by fewer edges. This clear demarcation between
 11 communities indicates that the GN algorithm effectively captures the modular
 12 structure of the gene co-expression network at this level.

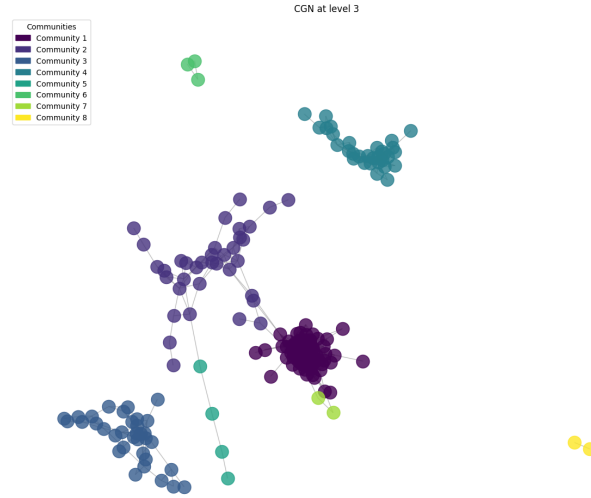


Fig. 7. Communities Detected Using GN at Level 3

13 At Level 3, the Girvan-Newman (GN) algorithm further refines the commu-
 14 nity structure within the gene co-expression network, as depicted in Figure 7. The
 15 algorithm successfully identifies eight distinct communities, each represented by
 16 a unique color. This level of analysis reveals more granularity compared to Level
 17 1, allowing for a deeper exploration of the network's hierarchical structure.

18 The visualization shows a notable increase in the number of smaller commu-
 19 nities, indicating that the GN algorithm is capable of detecting finer subdivisions
 20 within the larger clusters identified at earlier levels. Community 1 remains promi-
 21 nent, but it has split into several smaller groups, reflecting more detailed inter-
 22 actions among the genes. Similarly, Community 2 has fragmented into smaller,

1 more specific clusters, providing insights into sub-communities within the gene
2 network.

3 The spatial distribution of the communities in the graph suggests that the
4 GN algorithm can capture not only the broad divisions but also the subtle, intricate
5 relationships among genes. The presence of smaller, isolated communities,
6 such as Community 8, highlights the algorithm’s ability to detect outliers or
7 specialized clusters that may have unique functional roles.

8 This increased granularity at Level 3 offers a more comprehensive understanding
9 of the gene co-expression network, revealing sub-structures that might
10 be obscured at higher levels of abstraction. These detailed clusters can provide
11 valuable insights into specific gene functions and interactions, facilitating targeted
12 biological investigations and hypotheses generation.

13 5.5.2 Paris



Fig. 8. Communities Detected Using Paris

14 The Paris algorithm, as illustrated in Figure 8, identifies five distinct com-
15 munities within the gene co-expression network. Each community is marked by
16 a different color, providing a clear visual representation of the clustering.

17 The communities identified by the Paris algorithm exhibit a variety of sizes
18 and structures. Community 1 is the largest, forming a dense cluster at the center
19 of the network. This central cluster indicates a highly interconnected group of
20 genes, possibly reflecting core functional relationships or regulatory modules.

21 Community 2 and Community 3 also form substantial clusters, albeit smaller
22 than Community 1. These communities are relatively compact and appear to be

1 well-separated from each other, suggesting distinct functional groups within the
 2 network.

3 Community 4 and Community 5 are the smallest, with Community 5 con-
 4 sisting of only a few nodes. These smaller communities may represent more spe-
 5 cialized or less connected genes within the network, highlighting the algorithm's
 6 ability to detect both large and small functional groupings.

7 The spatial distribution of the communities in the graph shows a clear separa-
 8 tion between different clusters, with some nodes and edges connecting these clus-
 9 ters, indicating potential interactions or overlap between different gene groups.
 10 The Paris algorithm effectively captures these nuanced relationships, providing
 11 a detailed view of the gene co-expression network's modular structure.

12 5.5.3 LFM

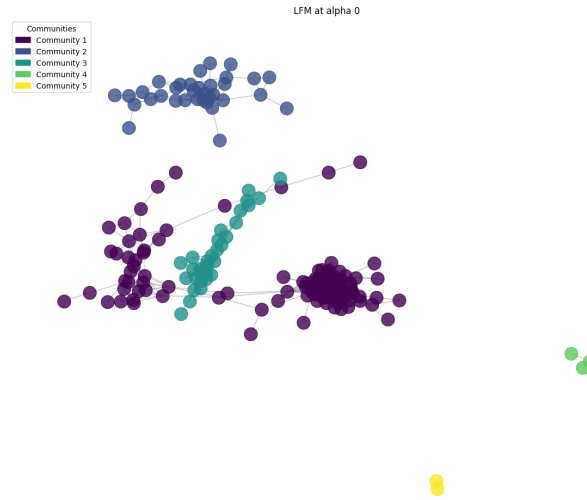


Fig. 9. Communities Detected Using LFM with Alpha=0

13 The visualization of the LFM algorithm at alpha 0, shown in Figure 9, identi-
 14 fies five distinct communities within the gene co-expression network. Each com-
 15 munity is distinguished by a different color, illustrating the algorithm's ability
 16 to cluster nodes based on their connectivity.

17 Community 1 is the largest and most densely connected cluster, positioned
 18 centrally in the network. This suggests that Community 1 encompasses a sig-
 19 nificant portion of the network's core interactions, potentially representing a
 20 primary functional group or regulatory module.

21 Community 2 is another substantial cluster located towards the upper part
 22 of the network. Its size and density indicate a secondary but still prominent

1 grouping of genes, possibly related to a distinct set of functions or processes
2 compared to Community 1.

3 Community 3 and Community 4 are smaller, more compact clusters. These
4 communities are relatively isolated from the larger groups, highlighting their
5 specialized roles or interactions within the gene network. Community 3, in par-
6 ticular, bridges between Community 1 and other parts of the network, suggesting
7 a possible intermediary role in gene interactions.

8 Community 5 is the smallest cluster, consisting of only a few nodes. Its
9 position on the periphery of the network indicates that it may represent a unique
10 or less connected set of genes, potentially involved in specialized or less central
11 functions.

12 The overall structure of the network under the LFM algorithm at alpha 0
13 shows a clear separation between the communities, with well-defined bound-
14 aries and varying sizes. This visualization underscores the algorithm's capability
15 to identify both large, central clusters and smaller, peripheral groups, provid-
16 ing a comprehensive view of the hierarchical organization within the gene co-
17 expression network.

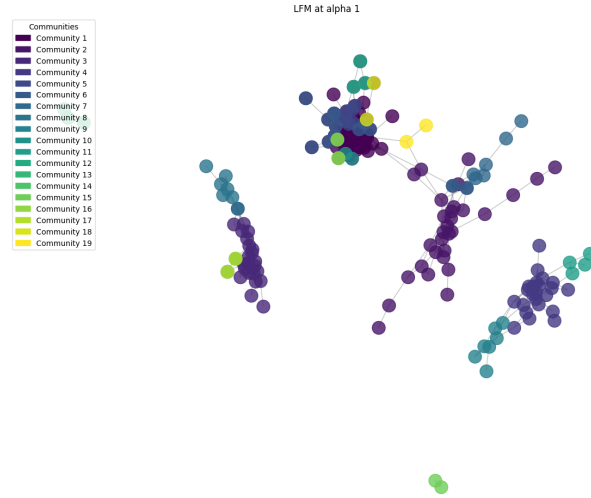


Fig. 10. Communities Detected Using LFM with Alpha=1

18 The visualization of the LFM algorithm at alpha 1, shown in Figure 10,
19 identifies nineteen distinct communities within the gene co-expression network.
20 Each community is represented by a different color, showcasing the algorithm's
21 ability to detect a high level of granularity in the clustering.

22 At this alpha level, the network is divided into many smaller communities
23 compared to the clustering observed at alpha 0. This increase in the number of

communities indicates the algorithm’s sensitivity to the alpha parameter, which controls the resolution of the community detection process. As alpha increases, the LFM algorithm identifies more fine-grained structures within the network.

The largest clusters, such as Community 1 and Community 2, remain central in the network, but they have now split into several smaller sub-communities. This splitting reflects the detection of more specific functional groups within the larger clusters, highlighting intricate relationships among genes.

Other communities, such as Community 8 and Community 11, are relatively small and compact, indicating specialized groups of genes with strong internal connections but fewer links to other communities. The presence of numerous small communities, like Community 18 and Community 19, scattered throughout the network further emphasizes the detailed segmentation achieved at this alpha level.

The overall structure shows a dense network with many inter-community edges, illustrating the complex and interconnected nature of the gene co-expression network. This high level of detail provides valuable insights into the modular organization of the network, potentially uncovering subtle functional relationships that are not visible at lower alpha levels.

The collected results from both the quantitative evaluations and visual analyses were then thoroughly examined. The insights gained from this analysis are discussed in detail in Section 6.

6 Results

Level	Number of Communities	Nodes in Largest Community	Singletons	Modularity	Closeness Centrality	ARI (Level 1)	ARI (Level 2)
GN	1	4	64	0	0.331094126	0.4292961301	0.5947369851
	2	5	64	0	0.3305609729	0.5090427336	0.5999552912

	14	17	58	6	0.3218033538	0.7451175912	0.574702
	49	52	38	37	0.244454	0.9030717181	0.303390071
Alpha							
LFM	0	3	85	0	0.2577457511	0.1565445117	0.4436266995
	0.5	4	65	0	0.331094126	0.1565445117	0.5947369851
	0.6	5	65	0	0.330336056	0.1565445117	0.5963115436
	0.7	11	64	0	0.327002455	0.1565445117	0.6045236205
	0.8	8	64	0	0.3298426235	0.1565445117	0.5940667658
	0.9	9	64	0	0.3269775116	0.1565445117	0.6062407526
	1	10	64	0	0.3237162358	0.1565445117	0.6038008742
PARIS							
		3	85	0	0.1142872613	0.2106443044	0.71500005
Highest Values:							
		52	85	37	0.331094126	0.9030717181	0.71500005

Fig. 11. Evaluation Results on Cho Dataset

6.1 Cho Dataset Results

From our analysis of the Cho dataset evaluation results, which can be seen at Figure 11, we observed that the Local optimization Function Model (LFM) and

1 Girvan-Newman (GN) algorithms demonstrated a slightly higher average mod-
 2 ularity compared to the Paris algorithm. This suggests that LFM and GN are
 3 more effective in forming tightly-knit communities within the gene co-expression
 4 network. Specifically, the LFM algorithm showed a significantly higher Adjusted
 5 Rand Index (ARI) at Level 1, indicating a better alignment with the ground
 6 truth clustering results. This highlights LFM’s effectiveness in capturing the
 7 underlying biological structure in its initial clustering.

8 Interestingly, the Paris algorithm achieved the highest ARI at Level 1 among
 9 all algorithms, suggesting that Paris might also be highly effective at correctly
 10 identifying the true community structure at this level. However, GN excelled
 11 in closeness centrality, particularly at Level 49, where it achieved the highest
 12 value (0.9030717181). This implies that the clusters identified by GN are more
 13 central within the network, potentially indicating more meaningful or compact
 14 groupings of genes.

15 Moreover, at Level 2, GN’s ARI slightly surpassed that of LFM and Paris.
 16 This indicates that GN maintains better consistency in its clustering across dif-
 17 ferent hierarchical levels, making it a robust choice for multi-level community
 18 detection in gene co-expression networks. GN’s balance between high modular-
 19 ity, closeness centrality, and consistent ARI scores across levels underscores its
 20 potential as a reliable method for uncovering hierarchical community structures
 21 in biological networks.

22 Overall, while LFM and Paris show strengths at specific levels, GN’s perfor-
 23 mance across multiple evaluation metrics and levels makes it a strong contender
 24 for community detection in the Cho dataset. This comprehensive analysis em-
 25 phasizes the importance of considering multiple metrics and hierarchical levels
 26 when evaluating community detection algorithms for biological applications.

27 6.2 KEGG Dataset Results

28 Moving on to the KEGG dataset evaluation results, which can be seen at Figure
 29 11, our analysis revealed several key insights about the performance of the com-
 30 munity detection algorithms. The Girvan-Newman (GN) algorithm consistently
 31 displayed higher modularity values compared to the Local optimization Func-
 32 tion Model (LFM) and Paris algorithms, with the highest modularity observed
 33 at Level 3 (0.3438994678). This indicates that GN is more effective at identifying
 34 strong community structures within the dataset.

35 In terms of closeness centrality, GN again showed a significant advantage.
 36 At Level 49, GN achieved the highest closeness centrality value (0.8191581058),
 37 suggesting that the communities it identifies are more centrally positioned within
 38 the network, indicating tightly knit and meaningful clusters.

39 Additionally, GN outperformed both LFM and Paris in terms of average
 40 Adjusted Rand Index (ARI) values across various levels. At Level 1, GN and
 41 LFM shared the highest ARI (0.03375774012), but GN showed better consistency
 42 in ARI values at higher levels, with the highest ARI at Level 2 (0.081454). This
 43 indicates that GN provides more accurate clustering results that closely match
 44 the ground truth.

	Level	Number of Communities	Nodes in Largest Community	Singletons	Modularity	Closeness Centrality	ARI (Level 1)	ARI (Level 2)	ARI (Level 3)
GN	1	6	86	0	0.3434830981	0.5612222908	0.03375740132	0.063101	0.036326
	2	7	86	0	0.3430306206	0.5739693694	0.031901	0.061727	0.036023
	3	8	84	0	0.3438994678	0.6287934906	0.036905	0.065596	0.035593

	42	47	64	27	0.3084383113	0.803585	0.04781246601	0.081454	0.055733
	43	48	64	27	0.3049668155	0.8050001952	0.04798033944	0.080478	0.056259
	44	49	63	28	0.3038680595	0.8052317503	0.049356	0.08014	0.052442
	45	50	62	29	0.3023369761	0.8054931566	0.051025	0.080189	0.04908497428

	49	54	60	32	0.2941402547	0.8191581058	0.041532	0.078276	0.046109
Alpha									
LFM	0	5	126	0	0.2560283436	0.1295114562	0.04922694425	0.069688	0.03168
	0.5	6	90	0	0.3434830981	0.1295114562	0.03375740132	0.063101	0.036326
	0.6	6	90	0	0.3434830981	0.1295114562	0.03375740132	0.063101	0.036326
	0.7	9	86	0	0.3436331397	0.1295114562	0.038579	0.065264	0.037519
	0.8	8	86	0	0.341914754	0.1295114562	0.03775841077	0.066443	0.03625
	0.9	15	86	0	0.3290094215	0.1295114562	0.035189	0.076349	0.042517
	1	19	86	0	0.3339341458	0.1295114562	0.02676239691	0.038042	0.034276
PARIS									
		5	126	0	0.1253258451	0.1749535266	0.030291	0.049484	0.02355150125
Highest Values:									
		54	126	32	0.3438994678	0.8191581058	0.051025	0.081454	0.056259

Fig. 12. Evaluation Results on KEGG Dataset

The LFM algorithm, while slightly behind GN in modularity and ARI, showed significant adaptability with different alpha values, highlighting its flexibility in detecting community structures of varying granularity. Despite this, LFM did not reach the highest values in any metric compared to GN.

The Paris algorithm, although it provided a baseline comparison, consistently had lower performance across the metrics, with its highest ARI at Level 1 being 0.0492484, which is still lower compared to the ARI values achieved by GN and LFM.

Overall, the results from the KEGG dataset strongly support the superior performance of the GN algorithm in detecting hierarchical community structures within the gene co-expression network. Its high modularity, significant closeness centrality, and consistent ARI values make it a robust choice for uncovering complex biological interactions and functional groupings in the KEGG dataset.

7 Recommendations

Based on our findings and the insights gained from this study, we propose the following recommendations to further advance the field of gene co-expression network analysis and community detection:

1. Future research could investigate the application of dynamic community detection algorithms. These algorithms can handle time-varying data, capturing dynamic changes in gene co-expression networks over various time points. This approach could provide deeper insights into gene interactions.
2. There is a need to develop standardized benchmarking frameworks to systematically compare the performance of different community detection algorithms. Such frameworks would help in assessing the robustness and accuracy of different methods, facilitating the identification of the most suitable algorithms for specific biological datasets.

- 1 3. The methodologies developed in this study can be utilized for gene co-
2 expression networks of other organisms. This approach would assist in as-
3 sessing the generalizability of our findings and confirm the effectiveness of
4 the community detection algorithms across various species and datasets.

- 5 4. Collaboration with domain experts, including biologists and geneticists, is
6 essential for the successful interpretation of results. The experts can provide
7 valuable biological insights that can guide the computational approaches and
8 help in the accurate interpretation of the detected gene groups.

9 With these recommendations, future research can build upon this study to
10 achieve a more comprehensive understanding of gene co-expression networks with
11 community detection algorithms and their applications in biological mechanisms.

1 References

- 2 Bhalla, U. S., & Iyengar, R. (1999). Emergent properties of networks of biological
3 signaling pathways. *SCIENCE*, 283, 381–387. [https://doi.org/10.1126/](https://doi.org/10.1126/science.283.5400.381)
4 [science.283.5400.381](https://doi.org/10.1126/science.283.5400.381)
- 5 Bonald, T., Charpentier, B., Galland, A., & Hollocou, A. (2018, June 22). Hier-
6 archical graph clustering using node pair sampling. Retrieved November
7 27, 2023, from <http://arxiv.org/abs/1806.01664>
- 8 Brandes, U., Dellinger, D., Gaertler, M., Gorke, R., Hoefer, M., Nikoloski, Z., &
9 Wagner, D. (2008). On modularity clustering. *IEEE Transactions on*
10 *Knowledge and Data Engineering*, 20(2), 172–188. [https://doi.org/10.](https://doi.org/10.1109/TKDE.2007.190689)
11 [1109/TKDE.2007.190689](https://doi.org/10.1109/TKDE.2007.190689)
- 12 Cho, R. J., Campbell, M. J., Winzeler, E. A., Steinmetz, L., Conway, A., Wod-
13 icka, L., Wolfsberg, T. G., Gabrielian, A. E., Landsman, D., Lockhart,
14 D. J., & Davis, R. W. (1998). A genome-wide transcriptional analysis
15 of the mitotic cell cycle. *Molecular Cell*, 2(1), 65–73. [https://doi.org/](https://doi.org/10.1016/S1097-2765(00)80114-8)
16 [10.1016/S1097-2765\(00\)80114-8](https://doi.org/10.1016/S1097-2765(00)80114-8)
- 17 Clemente, J. (2022). Predicting the biological classification of cell-cycle regu-
18 lated genes of *saccharomyces cerevisiae* using community detection algo-
19 rithms on gene co-expression networks. *Philippine Computing Journal*,
20 16. <https://doi.org/10.48550/arXiv.2208.10119>
- 21 Coscia, M., Rossetti, G., Giannotti, F., & Pedreschi, D. (2014). Uncovering hier-
22 archical and overlapping communities with a local-first approach. *ACM*
23 *Transactions on Knowledge Discovery from Data*, 9(1), 1–27. [https:](https://doi.org/10.1145/2629511)
24 [//doi.org/10.1145/2629511](https://doi.org/10.1145/2629511)
- 25 DasGupta, B., & Desai, D. (2013). On the complexity of newmans community
26 finding approach for biological and social networks. *Journal of Computer*
27 *and System Sciences*, 79(1), 50–67. [https://doi.org/10.1016/j.jcss.2012.](https://doi.org/10.1016/j.jcss.2012.04.003)
28 [04.003](https://doi.org/10.1016/j.jcss.2012.04.003)
- 29 Fortunato, S., & Newman, M. E. J. (2022). 20 years of network community
30 detection. *Nature Physics*, 18(8), 848–850. [https://doi.org/10.1038/](https://doi.org/10.1038/s41567-022-01716-7)
31 [s41567-022-01716-7](https://doi.org/10.1038/s41567-022-01716-7)
- 32 Girvan, M., & Newman, M. E. J. (2002). Community structure in social and
33 biological networks. *Proceedings of the National Academy of Sciences*,
34 99(12), 7821–7826. <https://doi.org/10.1073/pnas.122653799>
- 35 Gregory, S. (2010). Finding overlapping communities in networks by label prop-
36 agation. *New Journal of Physics*, 12(10), 103018. [https://doi.org/10.](https://doi.org/10.1088/1367-2630/12/10/103018)
37 [1088/1367-2630/12/10/103018](https://doi.org/10.1088/1367-2630/12/10/103018)
- 38 Hajibabaei, H., Seydi, V., & Koochari, A. (2023). Community detection in
39 weighted networks using probabilistic generative model. *Journal of In-*
40 *elligent Information Systems*, 60(1), 119–136. [https://doi.org/10.1007/](https://doi.org/10.1007/s10844-022-00740-6)
41 [s10844-022-00740-6](https://doi.org/10.1007/s10844-022-00740-6)
- 42 Hanneman, R. A., & Riddle, M. (2005). Introduction to social network methods.
- 43 Hevey, D. (2018). Network analysis: A brief overview and tutorial. *Health Psy-*
44 *chology and Behavioral Medicine*, 6(1), 301–328. [https://doi.org/10.](https://doi.org/10.1080/21642850.2018.1521283)
45 [1080/21642850.2018.1521283](https://doi.org/10.1080/21642850.2018.1521283)

- 1 Jaguzovic, M., Grbic, M., Dukanovic, M., & Matic, D. (2022). Identifica-
2 tion of protein complexes by overlapping community detection al-
3 gorithms: A comparative study. *2022 21st International Symposium*
4 *INFOTEH-JAHORINA (INFOTEH)*, 1–6. [https://doi.org/10.1109/](https://doi.org/10.1109/INFOTEH53737.2022.9751314)
5 [INFOTEH53737.2022.9751314](https://doi.org/10.1109/INFOTEH53737.2022.9751314)
- 6 Jarukasemratana, S., Murata, T., & Liu, X. (2014). Community detection al-
7 gorithm based on centrality and node closeness in scale-free networks.
8 *Transactions of the Japanese Society for Artificial Intelligence*, 29(2),
9 234–244. <https://doi.org/10.1527/tjsai.29.234>
- 10 Lancichinetti, A., Fortunato, S., & Kertesz, J. (2009). Detecting the overlapping
11 and hierarchical community structure of complex networks. *New Journal*
12 *of Physics*, 11(3), 033015. [https://doi.org/10.1088/1367-2630/11/3/](https://doi.org/10.1088/1367-2630/11/3/033015)
13 [033015](https://doi.org/10.1088/1367-2630/11/3/033015)
- 14 Lee, H. K., Hsu, A. K., Sajdak, J., Qin, J., & Pavlidis, P. (2004). Coexpression
15 analysis of human genes across many microarray data sets. *Genome*
16 *Research*, 14(6), 1085–1094. <https://doi.org/10.1101/gr.1910904>
- 17 Liu, X., Cheng, H.-M., & Zhang, Z.-Y. (2019, February 18). Evaluation of com-
18 munity detection methods. Retrieved November 27, 2023, from [http:](http://arxiv.org/abs/1807.01130)
19 [//arxiv.org/abs/1807.01130](http://arxiv.org/abs/1807.01130)
- 20 Ma, X., & Gao, L. (2012). Biological network analysis: Insights into structure
21 and functions. *Briefings in Functional Genomics*, 11(6), 434–442. [https:](https://doi.org/10.1093/bfgp/els045)
22 [//doi.org/10.1093/bfgp/els045](https://doi.org/10.1093/bfgp/els045)
- 23 Montenegro, J. D. (2022, January 16). Gene co-expression network analysis. In
24 *Plant bioinformatics* (Vol. 2443). [https://doi.org/10.1007/978-1-0716-](https://doi.org/10.1007/978-1-0716-2067-0_19)
25 [2067-0_19](https://doi.org/10.1007/978-1-0716-2067-0_19)
- 26 Moore, C. (2017, August 23). The computer science and physics of commu-
27 nity detection: Landscapes, phase transitions, and hardness. Retrieved
28 November 27, 2023, from <http://arxiv.org/abs/1702.00467>
- 29 Oyelade, J., Isewon, I., Oladipupo, F., Aromolaran, O., Uwoghiren, E., Ameh,
30 F., Achas, M., & Adebisi, E. (2016). Clustering algorithms: Their ap-
31 plication to gene expression data. *Bioinformatics and Biology Insights*,
32 10, BBL.S38316. <https://doi.org/10.4137/BBL.S38316>
- 33 Palla, G., Derenyi, I., Farkas, I., & Vicsek, T. (2005). Uncovering the overlapping
34 community structure of complex networks in nature and society. *Nature*,
35 435(7043), 814–818. <https://doi.org/10.1038/nature03607>
- 36 Pons, P., & Latapy, M. (n.d.). Computing communities in large networks using
37 random walks.
- 38 Rahiminejad, S., Maurya, M. R., & Subramaniam, S. (2019). Topological and
39 functional comparison of community detection algorithms in biological
40 networks. *BMC Bioinformatics*, 20(1), 212. [https://doi.org/10.1186/](https://doi.org/10.1186/s12859-019-2746-0)
41 [s12859-019-2746-0](https://doi.org/10.1186/s12859-019-2746-0)
- 42 Redekar, S. S., & Varma, S. L. (2022). A survey on community detection methods
43 and its application in biological network. *2022 International Conference*
44 *on Applied Artificial Intelligence and Computing (ICAAIC)*, 1030–1037.
45 <https://doi.org/10.1109/ICAAIC53929.2022.9792913>

- 1 Reichardt, J., & Bornholdt, S. (2006, July 18). Statistical mechanics of commu-
2 nity detection. In *PHYSICAL REVIEW E* (Vol. 74). [https://doi.org/](https://doi.org/10.1103/PhysRevE.74.016110)
3 [10.1103/PhysRevE.74.016110](https://doi.org/10.1103/PhysRevE.74.016110)
- 4 Rossetti, G. (2019). *Algorithms' table*. [https://cdlib.readthedocs.io/en/latest/](https://cdlib.readthedocs.io/en/latest/reference/cd_algorithms/algorithms.html)
5 [reference/cd_algorithms/algorithms.html](https://cdlib.readthedocs.io/en/latest/reference/cd_algorithms/algorithms.html)
- 6 Rossetti, G., Pappalardo, L., Pedreschi, D., & Giannotti, F. (2017). Tiles: An
7 online algorithm for community discovery in dynamic social networks.
8 *Machine Learning*, 106(8), 1213–1241. [https://doi.org/10.1007/s10994-](https://doi.org/10.1007/s10994-016-5582-8)
9 [016-5582-8](https://doi.org/10.1007/s10994-016-5582-8)
- 10 Rosvall, M., & Bergstrom, C. T. (2008). Maps of random walks on complex net-
11 works reveal community structure. *Proceedings of the National Academy*
12 *of Sciences*, 105(4), 1118–1123. [https://doi.org/10.1073/pnas.](https://doi.org/10.1073/pnas.0706851105)
13 [0706851105](https://doi.org/10.1073/pnas.0706851105)
- 14 Salido, J. A. (2016). Identification of candidate gene function group of yeast cell
15 cycle genes using gene expression data.
- 16 Schwikowski, B., Uetz, P., & Fields, S. (2000). A network of protein–protein
17 interactions in yeast. *Nature Biotechnology*, 18(12), 1257–1261. <https://doi.org/10.1038/82360>
18
- 19 Tang, J. (2018). Prognostic genes of breast cancer identified by gene co-
20 expression network analysis. *Frontiers in oncology*, 8(374). <https://doi.org/10.3389/fonc.2018.00374>
21
- 22 Vieira, V. D. F., Xavier, C. R., & Evsukoff, A. G. (2020). A comparative study
23 of overlapping community detection methods from the perspective of
24 the structural properties. *Applied Network Science*, 5(1), 51. <https://doi.org/10.1007/s41109-020-00289-9>
25
- 26 Wagenseller III, P., & Wang, F. (2017, December 2). Size matters: A compar-
27 ative analysis of community detection algorithms. Retrieved November
28 3, 2023, from <http://arxiv.org/abs/1712.01690>
- 29 Xie, J., Szymanski, B. K., & Liu, X. (2011, November 10). SLPA: Uncovering
30 overlapping communities in social networks via a speaker-listener in-
31 teraction dynamic process. Retrieved November 27, 2023, from <http://arxiv.org/abs/1109.5720>
32
- 33 Yang, J., & Leskovec, J. (n.d.). Overlapping community detection at scale: A
34 nonnegative matrix factorization approach.
- 35 Yeung, K. Y., Haynor, D. R., & Ruzzo, W. L. (2001). Validating clustering for
36 gene expression data. *Bioinformatics*, 17(4), 309–318. [https://doi.org/](https://doi.org/10.1093/bioinformatics/17.4.309)
37 [10.1093/bioinformatics/17.4.309](https://doi.org/10.1093/bioinformatics/17.4.309)
- 38 Yeung, K. Y. (2001). Cluster analysis of gene expression data. [https://api.](https://api.semanticscholar.org/CorpusID:61503158)
39 [semanticscholar.org/CorpusID:61503158](https://api.semanticscholar.org/CorpusID:61503158)
- 40 Yu, D., Kim, M., Xiao, G., & Hwang, T. H. (2013). Review of biological network
41 data and its applications. *Genomics & Informatics*, 11(4), 200. <https://doi.org/10.5808/GI.2013.11.4.200>
42
- 43 Zhang, H., Zeng, J., Tan, Y., Lu, L., Sun, C., Liang, Y., Zou, H., Yang, X., &
44 Tan, Y. (2018). Subgroup analysis reveals molecular heterogeneity and
45 provides potential precise treatment for pancreatic cancers. *OncoTargets*

- 1 *and Therapy, Volume 11*, 5811–5819. [https://doi.org/10.2147/OTT.](https://doi.org/10.2147/OTT.S163139)
2 S163139
3 Zhang, X., Miller-Hooks, E., & Denny, K. (2015). Assessing the role of network
4 topology in transportation network resilience. *Journal of Transport Ge-*
5 *ography*, 46, 35–45. <https://doi.org/10.1016/j.jtrangeo.2015.05.006>